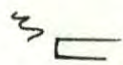


# BACKTRACKING $\approx$ REVISION $\approx$ Updating .25



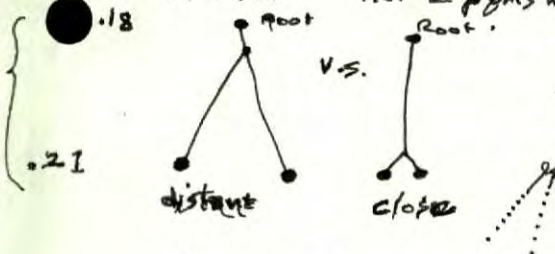
00: 195.40: In Z/41: We backtrack to previous; i. part that we write later part is any output of t. UMC. We retain a certain part of t. previous: Usually all of t. dots in their (Z) frequencies. We may discard a few of t. ~~most~~ most recent dots, but this is probably rarely done.

04: In "M.B.", I got I backtrack to a part where t. UMC had already started giving output bits: so t. corpus was (partially) visibly coded. I would imagine that for most kinds of codes, this would not work well. If t. corpus was not sequential, it could be possible to recode a certain section of t. corpus, (re-adding, perhaps to certain coding "codes" which were to be changed.

primitive

Perhaps most, my ideas as to what constituted "Backtracking" were rather (naive) primitive. A more general concept than "Backtracking": Modification of certain aspects of t. code, "keeping to root" "Backtracking" has t. Connotation of going back to an earlier (less informed) state. A More General idea perhaps is "Theory Revision" ("Backtracking" is one aspect of this)

Def N.B. One of t. ideas I had in "M.B." ( $\equiv$  don't backtrack) was to try to store (1) pairs that were Maximally distant: ~~for~~ (for "Diversity") This had a simple criterion proposed in M.B. — i.e. 2 pairs were very distant if their lowest common element was very high.



In other "Revisions" — (saving of 11 codes) — t. "Distance" measure is not so clear.

In Z/41: any 2 parsings of t. sum of dots might be regarded as "fairly dist. apart". If we use distant based dots, ~~to~~ to ~~code~~ code, then t. more non-common during t. more distant over t. 2 codes, (to  $\frac{18-21}{2}$  ~~of~~).

So, in int. Q is: Are there any general principles of "Revision" — common to many coding methods/pairs? A Basic idea here, is t. "Growing Corpus". The thing that suggests revision, is a lot of new data of very small pc. ("Very small" means as expected much larger — based on previous experience.)

.29: Actually, "Revision" is same as "Updating". Whatever does depends much on available cc.

So in General, it is a MAJOR PROBLEM Area for TM

.31: So, just as we expect to have a lot of OT's is a lot of induction Alphas (IA's) if we want to make a grammar to generate R. set(x): The updating methods ( $\equiv$  Backtrack  $\equiv$  revision methods) will also have to be generated/generalized. So: I will simply list various IA's along w. their updates & pairs, & look for common codes.

00: : So where am I now?

01 Well, Most recently I looked at ~ 3 methods of "Backtracking" for a simple UMC coder.  
Comparing them to the Update scheme for ZIT, I suggested that if any of them worked  
at all, this would be favorable to the (UMC, Regsys in corpus) pair.

That in general, Backtrack Theory revision & Updating, were all about the same problem - is  
that the way it was done depended much on it. It that was being "updated" (perhaps on  
the type of Regsys in the corpus).

Before that, I was worried about the QA system. Again, the major problem (I guess)  
was (other than T3 & writing), the update algm.

Some dittys: ① for medium cc updates, one would want to use only part of the corpus.  
Just what part is a difficult problem. For a "preliminary TM", we could assume updating

is for entire corpus. - (The, in general, one updating technique for limited cc is to "Backtrack" -

which is a way to use only part of the corpus (for certain IA's)

② I think I had a set of IA's [IA<sub>i</sub>] and a function f(x,y) that looks at x and

at time available for updating, and produces a wt. vector for  $\sum [IA_i]$  (i.e. tells how much time is  
to be spent on that IA for the problem, x,  
closely assoc. w. using an IA<sub>i</sub> for part of the updating, etc.)

③ We note that each IA<sub>i</sub> has at any pt. in time, a partial summary of the parts  
of the corpus to which it has been applied.

④ Go thru Perry's paper in Bibliog. refs. where description is discussed.

⑤ The main problem in QA was updating the IA's for "operator induction"  
An easier (I guess) induction problem is big induction, i.e. easier yet is seq. production.  
The operator induction can be viewed as ratio of a "big induction" problem - I usually don't  
think of it that way!

Consider the update algm. for seq. prod. ... the "easiest" problem. There's something  
way I know to do seq. prod. - each unit's own update algm.  
In fact, the core of an IA is its update algm.

[SN] I remember an apparent serious ditty in the f(x) idea of (14): F(x) has the  
problem as input & it has a d.f. on IA's as output. I. IA's have x as input &  
the solution to the problem as output. Somehow, F seemed to be identical in function w. the IA's.  
And, again, I'm not sure as to just what the problem was!

[SN] I'm really uncertain about "Anytime" problems! If we have an F(x), i.e.

x is an "Anytime" problem, then based on past history, the IA's to choose will  
depend on the time allowed. - If we use the "Lsach", the wts. of the various IA's should  
change as the colons ("Lsach") proceeds!

There's the standard T ← T solution to the Anytime problem. I think this is equivalent to using the Lsach, but  
changing the PC's as 37-38 - some more would be T is different & so the f(x) of the IA's would be different.  
→ 201.00

IDSIA

## T. "Correlation" problem

- .00: On the "Correlation" problem: In LS such for some problems, Guided by 2 proby distribns, LS can be optimum for "Blind Search", if the pc's of cards were uncorrelated. If they are correlated, many sets of cards will be about the same &
- .03 ~~it~~ is very wasteful to test all of the elements of ~~one~~ one of these sets. Testing only
- .04 one per "correlated set", would be best (if sets have little correlation between them). If one knew exactly what the "correlated sets" were, one could implement .03-.04 — but usually one doesn't know this.

An Alternative way would be to try to get "Trials w. high diversity". These are relatively high pc trials that seem very different & are expected not to be correlated.

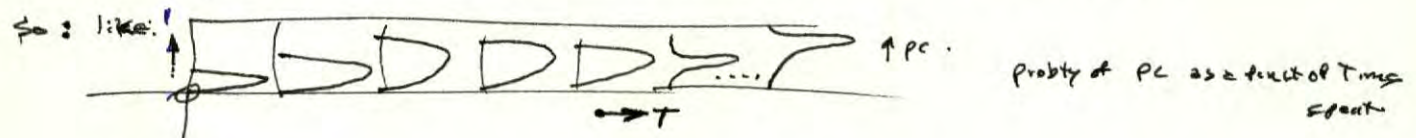
~~It~~ For each grammar that generates cards, one might be able to design up a "measure of disparity" of 2 cards, — that in some sense the 2 cards have "very different derivs" w.r.t. Generating Grammar. 2 cards can both be of high pc, ~~but~~ because they use cons of high pc ... but they use different cons — so the 2 cards are "very different".

Having cons that are uncorrelated is also  $\approx$  Having much "Diversity" in trials. Diversity in GA (or other search Algms) makes local extrema less of a problem.

00: 199.90: **SN** In the QATM I had been thinking of using many IA's in 11 (199.14 ff).

→ There were some apparently serious difficulties in that approach. If I used different IAs on different problems (or for different units of time on different problems), it would become very difficult to update  $F(x)$ . Remember that at any time an IA spends on each problem changes ("updates") that IA.

Remember that we don't compare 2 IA's by simply the total PC they use for their corpus. It may be closer to the ~~same~~ same mean of PC ... ~~But~~ ~~isn't~~ plotted more exactly, we want  $F(x)$  to give a def. on PC of soln. as a function of time spent.



Also, if  $T$  for an IA is long enough, it will sure work on problems in the corpus that it had not been applied to in the past! — Just how ~~update~~ update they (the IA's) go about this, is not clear!

Another BIG Q: Just how much time to we give to  $f(x)$  to compute  $f(x)$ ?

19. If we ~~allow~~ allow a "T" ~~as part of x~~ as part of  $x$ , perhaps  $f(x)$  itself could determine an optimum time limit for itself! At first, perhaps, the user could choose how long  $T$  for  $f(x)$ . Later  $f(x)$  could learn how to "optimally" do it for itself.   
 expenditure of time budget. itself is the set of IA's.   
 202.03 to optimize  
 202.03

I think that near the end of my most recent work on QATM, I began to realize that each IA changes as it works on problems, that what each IA is depends on ① the IA ② how much time it spends on each of the problems in the corpus

Assuming (19) is adequate: Is QATM at all clearly defined? — There is, of course, the problem of "updating  $f(x)$ ". This is an induction problem! It probably could be regarded as a QA problem of a special type — but it's a question

31 (w.o. any specific "correct answer" or even a set of correct answers.) → correct (19) (36)

32  $f(x)$  does, here, have a well defined source (?). We want it to give "Answers" (eg pc) products of the various "surfaces" of (19). Also, we have to discount its pc by its data cost. Also, we want this updating of  $f(x)$  to be able to use concs. This occurs in the IA's.

30: (31) Actually  $f(x)$  does have correct answers: for each IA that was tried on  $x$ , we got a pc in a time,  $\Rightarrow$  (correct answer) so updating  $f(x)$  is a standard problem. We should be able to apply many of the IA's to Paris' problem, since it is a QA problem! The Q is (or set of Q's) is

IDSIA:

Eric Wittgenstein

total / Time A.E. was applied to X

:00:20:40 : of the form  $(X, A.E., T_i)$  : "A" is the p.c. of the soln. of ~~the~~ X that was obtained.

.02 Well, Pitz sounds not so bad! Perhaps write a "summary" of this entire QA model.

.03 In 201.18-19  $F(x)$  finding an optimum T-time limit for itself is not exactly an induction problem! # ETS an "optzn problem" perhaps of the  $G(x) \cdot F(t)$  type - unsolved by Loren!

While it is an essential problem in the system, unworn it is not very critical, & for a final system, we can probably use some not very optimum solns: — T. main diff is that

we have to express it as a QA problem: # After TM has worked various  $G(x) \cdot F(t)$  ~~and~~ or problems as "QA problems", it will be able to <sup>possibly</sup> work on this one. } This will be a fairly highly developed ~~QA~~ ATM!

We could start off w.  $T_{FC} = \frac{1}{2}$  or  $\frac{1}{3}$  or  $\frac{1}{4}$  of  $T_{IA}$ . Fixed if the factor

is  $\frac{1}{2}$ , there is the old argument ~~for~~  $T_{50\%}$  soln that it couldn't be too bad!

On "updating": Each IA has its own method of updating. At any particular time,  $\Sigma$  IA will be represented by some stochastic Operator. "Updating" involves modifying this operator, in view of some arguments of the corpus. The way the updating is done will depend on how much time (cc) is available. For large cc, it ~~can~~ involve massive recording of the corpus.

I can list various IA's & their update algos.

1) Barn say (Say Alphabet  $\geq 2$ ): T. complete update consists of the set of case counts for the alphabet. These are trivially updated by the corpus, & the p.c. of the conditions are obtained from "Lap's rule".

2) Z141: The system state is the <sup>current</sup> set of deficits & their case counts. Updating ~~updates~~ updates case counts (in view of augmented corpus) — & parsing of the augmented corpus using the current case counts. For large cc, we can try re-parsing most of the corpus newly just the augmented part — perhaps all of the corpus. For larger cc we can look for new deficits, & if we find any, re-parse <sup>or re-count</sup> part of the corpus, or the entire corpus, in view of the new data.

3) AZ141: ~~The corpus is a set of functions or a set of (not necessarily) words.~~

Here, the corpus is a set of QA pairs. At any time, we have a stoch operator that goes from Q to A. This operator is generated "generated" by "AZ141", from a primitive set of operators. We have a pre corpus (like Z141) that includes the primitive set of ~~ops~~ funcs & any derived funcs. The p.c. of any ~~final~~ final func, is determined as in Z141 by the product of p.c.'s of its component func's. T. notation for functions is Polish

.00:202.40 : So  $f$  function is def'd by a string, & we can get pc's of each symbol in  $f$  string (as in 2.14.1) by considering how frequently that symbol has occurred in that s.f. function.

~~... In the simplest case, the grammar (NMTM),  $\gamma$ .~~  
Current function gives  $A_i$  for each  $Q_i$  thus far. (i. operator is not stochastic).

Finding function from grammar is easy for  $\gamma$  because of e. t.s.p. Updating is usually a bit complicated:

.05 Say we have a function  $f$  that works for all of corpus up to now. We experiment corpus  $Q_n, A_n$ .

.06 It doesn't fit this operator,  $f$ . Using the histories of  $f$  we construct a new hyp operator  $F_n$  that maps  $Q_n$  to  $A_n$ . Also we find an "ob" that can distinguish  $Q_n$  from all previous  $Q$ 's — so we know when to use  $F_n$ . This is individual combinator  $f$  &  $F_n$ .

Gives a new operator that works w. all  $Q, A$ 's this far. T. for "updating" is for small  $cc$ .

For larger  $cc$ , we want to integrate  $F_n$  to  $F_1$  is a hy pc way.

In simple t.s.p.'s, it is clear as to how to do this, but I can't say any more. Very General & bold this 'integration'. An alternative way to deal w. updates of .05, .06: Use  $f$  to do it. To do it, data (initially) as a special operator, using combinator for  $f$ . Now  $A_n$  is  $C_n$  w. use old  $f$ , but we have to "correct" it. We must use  $f$  to correct it. We must still report (Storage Modif. Mechanism) Schönknecht's proof of S2 to do a better  $f$  operator.

.16 If i. relation of  $Q$  to  $A$  is stochastic (NMTM), then ~~...~~

For each input  $Q$  we can have additional bits that are needed to obtain the particular  $A$  for this  $Q$ . (i. same  $Q$  when repeated, can have diffr.  $A$ 's)

One simple way to implement this: we have a non-stochastic  $f$  for all  $Q, A$ 's:

.20 In addition, we may (or may not) have a seq. of bits that tell how many bits of  $f$  output were in error, & just which bits they were. (See section 2.14.1 & Wolff, for exact formula (s) (There are 2 alternative formulas & obtained).

4) A third very General kind of induction  $\Rightarrow$  Bay induction. Operator induction can be done w.

from it. Using i. formula:  $P(A|Q) = P(Q) \cdot P(A|Q)$ .

In Operator induction we want  $P(A|Q)$ :  $P(A, Q)$  is  $P(Q)$  even obtainable by Bay induction from a corpus of  $Q, A$  pairs. So  $P(A|Q) = P(A, Q) / P(Q)$ .

In 3) we obtain  $P(A|Q)$  directly.

In the present pure Bay induction (under finite strings) is considered:  $Q_1, Q_2, \dots, Q_n$ .

~~...~~ In general, we do use a stack structure for i.  $Q_i$  set. There has been much written about Grammar induction. Most successful work has perhaps been on Prob. state Grammars, (subset-subsets of GF Grammars, .... I'm really not very familiar w. this. H.S. Fu has written much on applications of Grammars, but much less on trading them.

Wolff is working on some other General Grammars, but so far, I've found his writings very hard to understand.

ID51A

E=∞: ~~Unbd~~ Extensibility: (.18, .29)

Extensibility Unbounded

.00: 205.40: IN 3) (202.33) I think part we can be guided by a reasonable TSC —  
to help determine what kind of Grammar to use & just how to update. This is because  
I feel that I "know" the "form" for the TSC & (including updating) is all I then need to  
do is map this knowledge into a good formalism. I think the formalisms of  
Z141 & AZ141 would be very useful.

My latest idea was to use a Universal Functional lang. — like Lisp; — from which one can  
more easily to each problem in that lang, & see if I could find an updating scheme of  
reasonable pc.

Another imp idea was that of "cond pc", so the pc of a conc. would be a function of  
the application to which it was to be given. This was expected to deal w. one aspect of  
"t. Scaling Problem".

SN & Some Relevant Points to Read:

- 1) Mc Connell's paper on X589
- 2) Try to get idea as to what Wolff has been doing.
- 3) Try to get papers on CFG discovery: Look at Horning. Also other, more recent work.  
Not nearly using ALP or MDH.: Look at my own work on P56 discy

.18 → An Aspect of QATM not dealt w. in this recent analysis: Howt set of IA's is  
augmented by ~~IA's~~ "updating" F()? Since any conceivable type of IA can be  
invented, the system which is theoretically capable of arbi. class IA's. So it is a SE aspect of T.M  
On 201.32 I got I. idea that updating F() was a normal QA problem. It was a  
N MTM problem, hvr. (stochastic).

It gives  
UNBOUNDED  
Extensibility

Well, the output of F() can be a string that can desc. any conceivable IA!  
The induction (updating) problem for F() maps problems (X) into strings that desc.  
IA's: For a given TM, the set of poss. IA's will be a (small) finite set —  
maybe < 10 diffent (specific) IA's. Later, when TM has more experience w. grammatical  
induction (i.e. other good induction techniques), it will be allowed to extrapolate

.28 & 29 → initial set of IA's to the most general poss. set of IA's < Unbd. Extensibility; Extensibility = ∞  
E=∞ (not E=∞)

I'm not quite sure of the adequacy of .18-.28, but it may be ok.

.30 At present pt., I think I'm ready to write a good review of QATM!  
Better than usual, because the main outstanding problem is "Updating technique" & my main  
approach on this is to write a TSC, ... which ~~probably~~ probably will not prevent me  
from finishing the review.

6/5/01 Re: .30: First make a careful analysis of the system to be sure it's "all there",  
but I haven't left out some critical part. The "Updating scheme" is the main vacancy —  
But I think that any Update system should be close to that of Z141, AZ141, Lsrch.  
2/0 a combination of them. AZ141 consists of ① Discovery of new concs by combinations of  
old concs ② repairing of corpus. Lsrch enters in ① or ② when we have to do a search

205.00

IDP/A

to implement them. One example might be to find a coder for a corpus experiment, using existing codes. T. (FC) looks at the problem & decides on how to deal w. it.

One way would be for FC) to assign by pc's to relevant codes that would be used to attempt to (solve) code the <sup>new</sup> QA. By giving diffrnt pc's to the set of available codes, FC) is essentially creating a new IA!

FC) can also be used in this way to help solve "sub problems" that are parts of "Main Problems".

09 **SN** on QA Problems: Each Q should have one or more category indices

that guide TM in pooling its data: E.g. say we are giving TM a problem in

Curve fitting: we give it various  $x, y$  pairs. Each pair is a QA, & the pair has an index that tells TM: data is numerical, what kind of axis: Float, Integer, or whatever. So both Q & A could be "integers" & we could have a long sequence of pairs: all belonging to the "same problem": so both have a common index that gives the "problem name".

Also, when humans are <sup>given</sup> problems, they are also given lots of assoc. info

— Like: Did the problem come from physics, chem, economics, sociology, psychology, "computer science"... ect.

They know when (in what order) in the T & Q the data occurred.

All these things can, & should be used as auxiliary labels on the data.

On the other hand, we will probably want TM to learn on certain kinds of problems w. minimal amt. of such aux. "Heuristic" info.

One might also try tell TM whether a <sup>an "inductor"</sup> problem is MTM or NMTM

These indices can be also used directly, so TM can learn to categorize a problem that has ~~be~~ be inadequately indexed.



1:00: On a concept of a TSC: A kind of formal decln. ;  
 For a Sequential Corpus: We start w. an initial (corp) P.D.,  $P(X_n)$  (X\_n is a string of length n)  
 We start w. ~~example~~ problem  $Z_0$ . Using Lsrch we find a "summary machine"  
 for  $P(Z_0)$ . The cc is  $c_0$ . This new P.D. is  $P^0(X_n)$ ; we find  
 a new problem,  $Z_1$ . Using Lsrch, it finds a summary machine  $P^1(X_n)$  using  $cc = c_1$   
~~use~~ Summary machine  $P^m()$  w. use problem  $Z_m$ , ~~to find~~ <sup>and  $cc = c_m$</sup>  "summary machine"  
 $P^{m+1}()$ .

For a more general Corpus: we have a cond. P.D.  $P^m(Q, A)$   
 (or  $P^m(A) \equiv P(A|Q)$ ). After being given problem  $Q_m, A_m$ , using Lsrch's  
 suitable ~~updating~~ <sup>algorithm</sup> updating of  $P^m$ , we obtain  $P^{m+1}(A|Q)$ .

One problem is the ambiguity of the form: That is, meaning of  $P^{m+1}(A|Q)$   
 depends on just what sub-corpus ~~one~~ uses to compute it. Also, one might  
 well "over search" (in Lsrch) to get a better  $P^{m+1}()$ .

17 Another impl. omission is  $P^m()$ ; the character of  $EA(P^m)$ . — or is  $FC()$   
~~mean~~ meant as a "factorization" of  $P(A|Q)$  into  $F(Q) \rightarrow P_Q(A)$ ?

19 ~~FC()~~  $FC()$  is a function that goes from problems to  $EA's = P_Q^i(A)$   
<sup>one</sup>

20 So 100-11 is the more general approach. 17-19 is <sup>one</sup> ~~of~~ it. SEE 207.29 for last non-el. method 207.26

Basically, the NON-el procedure for TSC's: SEE 207.29 for last non-el. method We start w.  $P^0(Q, A|Q)$ .

This is any way to get from  $Q$  to the prob. of  $A$ ; we run ~~one~~  $P^0()$  on  
 $Q, A$  pair:  $Q^0, A^0$  (corpus). After "updating", we obtain  $P^1$  using  $CB = C^0$ .  
 We obtain the updated  $P^1(A, Q)$ . In general,  $P^i()$  are all "practical P.D.'s" —  
 i.e. their values depend on the  $CB$  used to evaluate them.

For  $CB = \infty$ ,  $P^{inf}(A|Q) \equiv P^i(A, A^i|Q, Q^i)$  (whatever that means!).

$A, A^i$  &  $Q, Q^i$  are meant to be "composite objects" that  $P^i()$  can accept as inputs.  
 I.e.  $Q, Q^i$  can be 2 questions &  $A, A^i$  can be their respective answers: But we can formally  
 regard  $Q, Q^i$  as a "single  $Q$ " &  $A, A^i$  as a "single answer"

33 I will probably want to ① try to get "summary machines" for BAC induction

② Get "summary machines" for stack operators. (after writing a good analysis of a general  
 soln. to the Operator induction problem (of course one could just solve ①))

022 ~~EQ~~  $P(A, Q) = P(Q) \cdot P(A|Q)$ . i.e.  $P(A|Q) = \frac{P(A, Q)}{P(Q)}$  : T. rel. to P. solns to  
 2 ~~big~~ probs.

Anyway do 33 first.

206.40: The  $\Sigma$  MDL way to do  $\Sigma$  BAG induction!  $P(X)$  assigns PC's to  $\Sigma$  instances of  $X$ .  
we find a  $P^i(\cdot) \ni (P \circ P^i(\cdot)) \cdot \prod_{k=1}^n P^i(X^k) = \max [X^k]_{k=1}^n$  is v. corpus (BAG).

(we take exponents to  $P^i(X^k)$  if  $X^k$  occurs  $> 1$  times in corpus.)

The Rearranged best update for  $\Sigma X^k$   $k=1|n+1$  is to ~~optimize~~ optimize  $P^i(\cdot)$  based on  $\cdot$ .

new corpus: starting from scratch. Here usually, we will find an update algm.

So  $P^i(\cdot) = \text{funct of } [P^i(\cdot) \text{ and } X^{n+1}]$ , inasmuch as we have trouble doing  $P^i$ ,

we may backtrack several  $X^k$ 's (not nearly in ~~temporal~~ temporal order in which they occurred but maybe in reverse order)

For <sup>stoch</sup> operators; some system: we find operators ~~directly~~  $P^i(A|Q)$

$\ni P \circ (P^i(1)) \cdot \prod_{j=1}^n P^i(A_j|Q_j)$  is max. then updating  $P^i \rightarrow P^{i+1}$

$P \circ (P^i(1)) \cdot \prod_{j=1}^{n+1} P^{i+1}(A_j|Q_j)$  is max!  $P^i \rightarrow P^{i+1}(1) \ni$  anti-summarizing machine of  $[Q_i, A_i]$   $i=1|n+1$ . (ie. each "IA")

Normally, we would not use .11 to get  $P^{i+1}(1)$ , but make  $P^{i+1}$  a func of  $P^i(1)$  (a modification)  $\ni Q_{n+1}, A_{n+1}$  as in (.04-.06)  $\uparrow$ . In general, each technique for obtaining  $P^i(\cdot)$  from  $\cdot$  corpus will have its own "update scheme".

All a IA is is a means for getting an approx max  $\cdot$  .10

Imp't idea: That all IA's amount to getting an approx of .10  $\dots$  that they: involve specifying CB used: this makes each such (IA  $\equiv$  approx), a type "Proced Pd" (ie.  $\in \Sigma$  is relevant param) reg  
(2) IA  $A \rightarrow Q$ : T idea is "Induced by definition" One finds a reg, defines it,  $\ni$  (rephrases  $\cdot$  corpus in terms of  $\cdot$  new defn). That there will be many ways to rephrase,  $\ni$   $\Sigma$  in impl. fact in evaluating  $P \circ \Sigma$  of  $\cdot$  data.

.26: 206.20: A better "More General" approach turns 206.00-.11: GPD( $\cdot$ ) looks at

$Q_i$  and becomes a IA that gives a pd on  $A_i$  ("IA" means that we have to give CB info)

How do we update GPD( $\cdot$ )  $\ni$  w. a new  $Q_i, A_i$  pair?

.29 GPD( $Q_i$ )  $\rightarrow$  PD on  $A_i$ 's. we want GPD a  $(P \circ GPD) \cdot \prod_{j=1}^n GPD(A_j|Q_j) = \max \ni$  (see Note on OSL: 209.07)

~~usually~~  $[Q_i, A_i]$   $i=1|n$  is  $\cdot$  corpus. T. problem of updating GPD  $\ni$  as stated, is quite difficult.

~~usually~~ In .29, we don't even look at  $\cdot$  GPD for  $\cdot$   $i=1|n-1$ . In practical induction,

usually we look at  $\cdot$  GPD for  $i=1|n-1$  and we modify it in view of  $Q_n, A_n$ . with OSL form of 209.07

So: .29 is "best" most non-cl. kind of TM. We will consider various approaches to it.

This "best" can be obtained by  $P(A|Q) = \frac{P(A, Q)}{P(Q)}$   $\ni$  in which case we have "factored" one operator inducing problem into 2 BAG induction problems.

Obtaining  $P(A, Q)$  seems to be obtain use some cues as getting  $P(A|Q)$  Directly

So I'm not sure we saw any thing. One advantage: it is a slightly different way to look at  $\cdot$  problem.

6/7/01  
IPSA

Rev

00:207.90 : Outline of Review of QATM.

167.00 is early review.

160.00, 18 on SI for QATM:

161.15 early demand QATM

152.00 - 201, 156.00

150.33# list of E X's & QATM update sequence.

Work on Recent QATM seems to begin no more recently than 150.00 or 149.00.

140.00 - 141.11 was a kind of review of work up to that time. - it told how present

TM model differed from that of ~~SSG~~ SSG; essentially, a list of impl. delays since SSG

✓



24/4/2000

00:207.20 Next, Consider "factory"

GPDA(Q) into  $F(Q) = PD$  on  $IA_i$ .

then  $IA_i$ 's operating on  $Q$  to give D.F. on  $A$ . We do run into 'problems' of each  $IA_i$  having

had a diffrnt CB for each  $Q_i A_i$  in  $t$ . corpus; But see 210.01!?

210.04

What constitutes a usable TSQ will depend on what IA's are used & just what it's updating

Method is .

07.207.29 SN T. formulation of 207.29 is essentially MDL: So OSL is not included: In most QA

TSQ's, I would think that OSL problems would be very common (which normal MDL doesn't deal with)

The (Theoretically O.K.) method of dealing w. OSL, using MDL, is to include (Anns, Qnns)

into product of 207.29 (i.e.  $t$ -product index =  $n+1$ ):  $Q_{n+1}$  is known, & Plus gives up a D.F. over

poss. values of Anns. I'm not sure .09-.11 is really correct, in terms of complete ALP. formulation.

Consider MDL for sequential prodn: "data" is  $t$ . binary string  $X_n$ : We are considering various

poss. contents, "z":  $P_i()$  is  $t$ . set of all perms (including partial recursive perms)

T. D.F. on "z" is Perm  $\sum_{i \in P_i(z)} P_i(X_n^z)$ . For each z, we write here

a list of  $i$  that max  $\sum_{i \in P_i(z)} P_i(X_n^z)$  product. The ALP formulation is  $\sum P_i(P_i(z)) \cdot P_i(X_n^z)$ .

I don't know how this would look w. Pns. Probably he would restrict to sum to his pre-defined set of models (his "Model Class").

SN (Regression Curve fitting) (Analog/real) is L search: Regard it as an OZ problem in following way!

The form of the model is a discrete set of polys. T. paramns of the space of real models, is a real space.

So, we assign a polys to the discrete space of models and do an L search from. Each discrete model is then given a score

by optimizing its params wrt to data, using standard optimization. T. score is then not Hessian in

continuous model space mult by PC of discrete model (mult by PC of data wrt to model).

T. time for each paramns optzn probably doesn't depend much on size of data set.

If we actually spend same amt of time on each model, this will be  $\approx$  McConnell's SURFER

McConnell discusses commutativity & associativity contributing to Redundant codes

& how he got rid of this (by hashing I think) — The one unit defect identity by use of random

Params. (i.e. 2 funcs are probably identical if they have same values for a random argt.)

One unit use of this redundancy to give more time to certain funs.

Say we see anomaly on 1 cond; for awhile. Then we eventually go to next cond; & discover it's identical to cond 1: We then go back to cond 1 & work on it an amt. assoc all cond's PC.

Similarly, whenever we detect a cond identical to a previous cond, we go back & work on the previous cond. I think, we mark the dup. cond, & use mark t. first cond as

having hyper partition apparent.

A (perhaps) easy way to detect identity: look at prodns for first n examples; if they are the same this is evidence that they are identical in general. (But this is after optimization, so: maybe not such good test!)

IDSIA

There are at least 3 "requirements" that I'm committing in this preliminary QATM:

- 1) T-Corp is always wrt. to entire corpus, rather than wrt. a part of it. Corpus more directly related to <sup>latest</sup> ~~problem~~ <sup>effect problem</sup>
- 2) We do not address the problem of correl. betw. cond.
- 3) We do not consider OSL.

Continuing on the "factoring" of QATM's Corp into FC) is  $\prod (IA_i)$ .

If each IA has worked on various sections of the corpus different parts of time, we may want to take a  $\gamma$  power of the resultant PC's ( $\gamma = 1$  for larger  $T \in CB$ ),  $\gamma < 1$  for small  $T$ : for small  $T$ , this small  $\gamma$  means  $PC \rightarrow 1$ , which is a way of giving that prod. "less wt.". It might be clearer as to what is being done if we take  $\frac{1}{\gamma}$  of the produce of the PC's of the IA's — in which case,  $\gamma$  becomes a coif of the various  $\ln(PC)$ 's.

Suppose we have several IA's making a prod. about  $A_n$ . How much wt. to give to each? —

Depends on the "size" of each IA on problems of that type, & how good it has been for such problems

Essentially, the  $\gamma$ 's have been assigned to the IA's by FC) (via "wts" = "pc's")

Hyp. the opinion of FC) is "double" in a sense that FC) gives small  $T$  to IA that it thinks would not do well on problem.

So we have w. small  $T$ , an "inert IA working for a short time" so we have a weak effect.

In general, the marginally effect of FC)'s assigning a  $T_i$  to IA<sub>i</sub> depends much on

- both IA<sub>i</sub> & the problem. For certain IA-problem pairs, ~~FC) would~~ "as good as it can do" solution is obtained in less than time  $T$ . (even tho  $T = T_{pc}$  wr.  $T$  did &  $pc$  assigned by FC). ABCde
- Another possy is that IA<sub>i</sub> is not able to do anything useful in time  $T$  on the problem.

Going Back to (04) The ALP soln. to the QA problem is more like:

$$\sum_i (pc(P_{mi}) \cdot \prod_j pc(A_j | Q_j))$$

So perhaps FC) could assign  $PC(\#)$  to  $P_{mi}$ .

But FC) looks at  $Q_n$  first!

- 27 A perhaps more relevant approach! Given the new data pair  $Q_{ni}$ ,  $A_{ni}$ ; How much time
  - 28 does each IA get to "update" on this new data?
- Another approach! That each IA has a "shallow" "quick update" time & a "deep" update time. "Deep" in many cases can be flexible so longer time gives better update.

- 31 Actually, there are 2 (apparently related) problems; 1 - Problem of how much time to spend on each IA; 2 - Problem of using (this corpus is the set of IA's that have
- 33 been updated over various times ~~to~~ wrt various parts of the corpus) to make a prod. for a new  $Q_n$ .

6/9/01 An imp. consideration in 27-29 & 31-33 is the fact that the IA's may "update" at different rates.

Let us study some extreme cases of mixtures wrt. update/prod. problem

- 1) Say 2 IA's have been ~~not~~ trained on a completely different corp. They both give their own train for  $A_{ni}$  w. in pot  $Q_{ni}$ : How do we combine these 2 tr.'s? ~~the~~ d.f.
- Is this problem related to "META Analysis"?

6/9/01

IDSIA

730 start.  
755 end.

.00: 210.40: Well, if we assume th. 2 IA's are "uncorrelated" (= "Independent"), then we simply multiply th. 2 d.f.'s. If a IA has small/SSZ a poor success; it will have a broad, flat d.f. and it will have little influence when combined w. th. output of an IA of large SSZ in which success ... i.e. a very narrow d.f.

O.K. Great! — but In fact, th. IA outputs are <sup>usually</sup> correlated: they ~~are~~ <sup>use</sup> partly same data & often similar induction Methods. — So its Really every Common, Very General Problem: To combine 2 or more d.f.'s that have been obtained in somewhat disparate Ways. This does seem Related to Meta Analysis.

Earlier work on this problem: 177.00-03, 177.16-37; 179.08, 179.36

.10 To do mixing of uncorrelated prediction P.D.'s: Mult them together. ( $P_1 \cdot P_2$ )  
" " " " <sup>(if they are</sup> ~~highly~~ <sup>highly</sup> correlated) " " ( $P_1^{\frac{1}{2}} \cdot P_2^{\frac{1}{2}}$ ) if they are not uncorrelated:  $(P_1^{\frac{1}{2}} \cdot P_2^{\frac{1}{2}})^{1+\epsilon}$ .  
If we ~~make~~ assume (for simplicity) that all IA's are uncorrelated, we would multiply all together & get a spuriously narrow d.f.

Using th. Model of 179.09-22, but w/ <sup>"with"</sup> multiplication (instead of with addition) of IA P.D's, we get something quite different from ALP's "All Perm's Method", & its not really justifiable from ALP pt. of view. ! The idea is that F(.) would blow into a sort of cross-correlation, as well as to "mix" of th. IA's for each Qj problem. So F(.) would give a certain matrix for each Qj problem and prodn would be ~~■~~ ..... ?

I'm thinking of a "Corr. Matrix" as a "poorman's joint P.D."

.21: 177.37 Also Note! Int. corpus coding scheme of 177.44-37: The pc of th. corpus includes pc's of each of th. IA's in addition to th. pc of F(.) (Th. switching function). (Note "cheap robe"!)  
The IA's in this code are not "all codes". But, th. original pc of th. IA's are a priori — so they will be large pc's (low costs), (Maybe zero cost; PC=1)

So 2 problems ① How to make a predn. based on this mixture of IA's!

② How to update th. IA's (i.e. which IAj's to update with which  $(Qj, Aj)$ 's & how much ~~to~~ <sup>to</sup> ~~or~~ <sup>or</sup> ~~to~~ <sup>to</sup> ~~use~~ <sup>use</sup>)

① is contingent on how ② was solved — so say we know soln. to ② — how to solve ①?

Do ① w. & w.o. correlations/bias taken into account.

MAIN Problem

Chen's well-known method of coding <sup>(non-linear)</sup> linear regression,

At the start we just give the discrete params of the code, i.e. how many bits  $n$  of its linear regression. We assume steady-state, so if we have  $n$  bits a history  $h$  ( $h > n$ ) (R window). Then we know the previous  $h$  values of  $x_i$  to start off. Since  $n$  is  $h$  are given, the bits can be computed for each prodn. in a standard way, so we need ~~to~~ give only the deviations from the prodn., i.e. this is coded using a Gaussian dist. So we can so to active code on bits of specifying  $h, n$  is the seq. of normalized deviations from predicted values,

This last depends only on  $\sigma^2$  and the no. of data pts  $= m$ .

We can choose  $h$  or  $n$  so to minimize total cost of the code, or take all codes in 11.

It would be interesting to take a chaotic seq. a good linear code itself is not that one gets slightly better prodn. by using 11 coding (corresponding to a non-integer no. of bits)

The main idea in the fugg. is that we don't have to code the regression coeffs.

Also, the method works for  $X$  window as well as for  $R$  window treated as  $n$  bits.

I think the coding method of .00ff requires infinite precision, ~~the~~ even if the values of  $X$  are discrete (say 256 values of  $x_i$ ). Any roundoff error would probably accumulate. Here, if we did restrict  $X$  to 256 values (say), and  $\mu_i$  was known by the "receiver", it would be to find exactly how the calculations would be made in integer math, I ~~think~~ ~~could~~ be possible to do it with finite precision. Each prodn. would be a discrete value of  $x$  (actually  $x$  would have to have a very large range, but  $\text{granularity} = 2^{-8} = \frac{1}{256}$ .)

.00:

In my recent work on "Min backback coding", I ended w. the idea that probably it wasn't a v. good idea: That most codes started w. a large non-priority part, that dec'd to form — & then they just coded the corpus in terms of that form.

For such codes, min-backback is not much good.

$2^{30} \approx 10^9$

$2^{40} \approx 10^{12}$

But for coding corpus in which the default form was only  $< 20$  bits ( $\geq 20 = 10^6$ ) I'd write to do just do standard Lsrch to find codes. Maybe 30 bits is 40 bits w. new faster Machine.

For (analogy = real) ISM data, I partition binary form by using a threshold. Whenever price exceeded the threshold (up or down), we reset price to zero & look for next threshold crossing. This gives a seq. of ups & downs (0's & 1's) — binary seq.

Here, predicting ISM from single seq. is not a good way!

Anyway, deciding on a best set of instructions for the ref. machine is not trivial! perhaps look at keene for derivs of simple instructions for recurrence function.

Or I might just try simple seqs like (10011)<sup>(100)</sup> or (01500)<sup>(100)</sup> or (0015)<sup>(100)</sup> etc. as input problems.

Try to use 11 modes of Pentium or AMD.  
or use Sony's PS2 chip.  
or use 32 (or 64) bit 11 capabilities of machine for binary math.

**SN** In recent work on S89+ (4/1/01 - 5/1/01 say at IDSIA & at Lugano):

I did have this as "proof" that the prod was "within a factor of 2 of optimum" — This compared it to a human using heuristic such methods.

There were a few conditions that made it der part from this "optimality",

but it still looked pretty good! It used Lsrch.

The most recent work on QATM does not seem to use Lsrch, but it is not yet a finished system (to put it mildly!) — it may eventually use Lsrch, & it may eventually be possible to use it optionally on it.

Just what was wrong w. S89+? Perhaps nothing! — That QATM was originally meant to be a simplification of it — No Inv. probs no "general" QZ probs — just

QA (Operator reduction) problems

→ 140.00 - 141.40 was a review of S89+ & told what some diffs were! I think "correcting" better. cards was a big problem. Also, I don't know to what extent I was able to solve the problem of getting TM to remember & use info obtained during trials on other cards during the same "Round" of  $T \leq T$  Lsrch — or in analogous situations in 11 Lsrch.

→ Q: ~~Can this last problem be regarded as an aspect of "correcting better cards"?~~

Could present problems w. QATM be applied as criterion to the older S89+ system? The older S89+ system would work QA problems by Lsrch on the IA's & expand the IA's by expressing known programs. (QA would be regarded as an "QZ problem").



.00: 177.37 : On the approach of 177.16-37 of having obs test tall which IAs to use on a problem. This seems different from using obs, to tall which operator (function) to use in xprog.  $Q_j$  m to  $A_j$ . The IAs has as arguments,  $\sqrt{Q_j}$  plus all (or at least a large part) of the previous corpus. The function only has  $Q_j$  as argt. The function is, (presumably) much faster than IAs. IAs has as output, both the proper obs to (micro)function that operates on  $Q_j$ .

Some Random Ideas:

.07 1) 207.26 find a general TM  $GPD(Q) \rightarrow$  an IA that takes  $Q$  m to a PDA IA. - rather non-cl, but diff.   
 Because we give a CB.

Could we use L search to optimize  $GPD(Q)$  directly?

One (somewhat al.) way: The update problem is to find  $GPD_{c,Q}()$  as a function of  $GPD_c()$  and  $Q$ . This assumes no back tracking, how can we allow back tracking if we allow the func to be from  $GPD_{c',Q}()$ ,  $X^2 Q$  to  $GPD_c()$  c is corpus, updates  $\leftarrow GPD_{c,Q}()$  is  $GPD_c()$  with back new corpus  $c, Q$ .

Have  $c = c'x$ ; ( $c'$  is a prefix of  $c$ ). (We have back tracked by "X"  $\leftarrow$  So to "Back tracking" function is  $\approx$  a special case of

non-back tracking function: Normally, we will update using

$\gg$  one QA pair, anyway.

As usual, the amt. of Backtracking one will do, depends on amt. of PDA cc available.

2) If, in QA TM, no part of  $Q$  is "to problem dom", it becomes clearer that this QA TM is potentially very powerful.

3) in .07 (1) we could start by doing the updating by "Hand" or "Expert system" in which we perm. the updating system. (But the update system is in "factored form" so later, TM will be able to "understand" it). After TM has learned about, we will give it the problem of "improving the update system". To understand the problem, is to have good ideas on how to solve it, the TM will have to have a TSO of suitable "definition lang" examples of various updating systems. In 207.26 .07:  $GPD_c()$  can be any induction system.

What we need, of course, is an "expert" induction system that is good enough to learn out to improve its induction algm. (which is ~~improvement~~ one of the goals of S89. - The other way is to read that lang.). So all of this QA TM work should be simply oriented to finding an adequate (but not necessarily perfect) induction system for QA problems.

Also, in recent QA TM work, the TM was, even at its beginning, able to work on the problem of improving the induction Alg. (was this "Overkill"?)

So, to start off, I could just use AZ141 as the  $\leftarrow$  learn TSP's for Algebra. Use  $\approx$  Lisp (or Lisp power insns - perhaps not even recursion) - add insns, later. Now, we do have to watch the "scaling" problem! (perhaps finally by using just the PC of any abstraction is a function of the situation calling for it)

Next, look at 140.00-141.40 to see how relevant to the post S89 ideas are to .35 See that I do have a complete formalism for early Alg. TSO

IDSIA

00: 114.90: On the "non-universality" of AZ(1): AZ(1) is "universal" on the first level.

The method of updating is not universal - because you require in the corpus that the method of updating has no way to recognize.

But, if I do a scheme of 214.19 -> 24, in which QATM is eventually able to work on "SE": improvement of ~~the~~ induction its own induction/updates scheme, it may be "universal" in the needed sense. To attain this "true universality" it will have to have a suitable TSO - of the type mentioned in 214.19 - 24.

At present, I have a very preliminary idea of how to do updating in AZ(1). As I get TM to work harder & harder problems, I will find that I have to "add to" / ~~then~~ "modify" its updating scheme (= "expert system"). So I will be "learning along w. TM" & expressly this "learning" in TM's update/induction system. -> (218.01)

12 [SN] On the justification of AZ(1)'s macrocode coding/p.d. construction:

Say I have a system that I can describe, that takes any binary (or finite radix or possibly continuous or <sup>discrete</sup> infinite radix) sequences, & associates w. each number of the seq. a ~~the~~ positive semi-definite ~~no.~~ <sup>no.</sup> (positive definite matrix better) feature & a unique function of the symbols in the sequence preceding it.

Then I can normalize (using no. 6.2) regard them as cond. pc's. The resultant norm. of the functions constitutes a legit p.d. or ~~compatible~~ set of all infinite strings over the alphabet of 13.

In the present case, it legitimizes the probs assigned by 2141. As to AZ(1) (& assignment of pc's to functions), it's not yet clear in my mind, - tho I guess it is a "well defined" assignment of pc's to all functions in binary use by AZ(1).

But, because of the possibility of many different parsing methods, in practice, 2141 & AZ(1) will not give unique data prob by dif.'s. Any parsing alg. (function/interp data) ~~gives~~ will, tho, give a unique coding & a unique p.d. for any corpus.

27 Another fun/diffy: in AZ(1) (applied not really to facts but to ~~the~~ coding of a sequence (corpus - say Nat'lamps) - to code a letter A following a letter C:

If we have defined the condition, then all previous cases of A following C must be accounted for <sup>specifically, explicitly</sup>. Do I have to normalize the A with all occurrences of C in the past, (in which A never have followed)? - or do I just write a no. for Prob A = no. of times it has followed C in the past, then normalize w/ all other probs. & symbols? I should get

the same result. - if not, the system is inconsistent: Suppose ~~the~~ C -> C B had also been defined but not C C or C D (say merely have 4 letters in alphabet). I think I have to resolve the ambiguity of coding before it becomes a "well-defined system."

In pure 2141, we only define n-tups, & there are no ~~strings~~

6/11/01  
ID612

∴ Augmented "Laplace's Rule"

- .00: 215.40: "Intersymbol constraints" (other than those depend by  $n$  types). We have a pure form seq., but w. certain case assoc. w. the domain of symbols.

.02 → In the case of 215.27 ff, I will have to clearly define what I want the pc's to be!

(This Q will become very imp. as the TSQ grows ("SCALING") because I will want objects (definitions) to become more & more "context dependant" ← (Gray/publication text "hr")

So: Say at present, AZ141 can assign pc to any (conceivable) function. [Ambiguity (multiplicity, parallelism) in parsing  $\Rightarrow$  a problem, hr]. So we have a function (or perhaps several functs) that desc to corpus  $[Q_i, A_i]$  <sup>with</sup>  $\Rightarrow$  discreteness.

We add  $Q_{n+1}, A_{n+1}$  to the corpus: What are some good, "fast" ways to get a new funct for the augmented corpus? ( $\equiv$  "Updating"). So that's the

.14 by problem. I do have one method: In which I find a ~~cheap~~ cheap funct relating  $A_{n+1}$  to  $Q_{n+1}$ , using pc's w. the old function, then I find a

.16 way to acceptize  $Q_{n+1}$  as an <sup>cheap</sup> any/recursive function re the old function.

Then I have the largest funct that works in the augmented corpus: I try to "compress" it ← (which I'm not at all clear about!)

So: Get a <sup>reasonable</sup> TSQ & see how I seem to do updating & incorporate that into the AZ141 system.

Conceivably (probably?) OSL can be used in .14 & .16 to get "cheap functions": These are functs that have not yet been defined, but defining them is cheap because they have occurred once before (in undefined form)

Note again that Z(4) is AZ141 applied to operator (QA) induction, & is quite definite. Z(4) is a way to assign codes (i.e. pc's) to a sequence (corpus). AZ141 in operation, can assign good (ish) pc's to functions formed from a primitive set, but it does not tell how to decide by pc functs that desc to corpus, for a small corpus, we might use Lachont's set of functs, but for larger corpus, we have to devise a good ~~set~~ <sup>min</sup> system for "Updating", i.e. "Min backtracking" (hr, see 214.13 - it for why/backtracking is ~~not~~ <sup>min</sup> some type of problem as ordinary "Updating").

OSL one Shot Long.

00: 216.90: In OSL (but 2141 & A2141) Once you think of it as: In coding a new chunk of code, one normally has a choice ~~of~~ of all primitives & all defined symbols. In OSL, one adds to this, the set of all defined symbols of symbols that have occurred in the past. So if "AB" has occurred once in the past, (but has not yet been defined) — then we can <sup>tentatively</sup> use AB as a defined symbol in coding. [ "tentatively" means we have to see if the definition really gives any compression — i.e. analysis of OSL in my writing above will give (perhaps easy to use) constraints on when it is likely that a ~~set~~ set of symbols will be useable in OSL.

.14 An idea in A2141 related to OSL coding: Say A is a large operator; ~~with~~ It has  $\alpha$  as one of its sub-operators. Then, A w. ~~with~~ substitution of  $\beta$  for  $\alpha$  in A (  $\beta$  is an operator w. some no. of I/O's (is this what "arity" is?). will give an operator w.  $\beta$  fairly by PC: The reason: A : B can be easily coded using a dummy variable function x, that can take 2 values  $\alpha$  or  $\beta$ . [ I'm not entirely clear on this entire process, hvr. ] [ I'm not entirely clear on this entire process, hvr. ] [ I'm not entirely clear on this entire process, hvr. ] ← Koh-1-hor

.21 The reason  $\beta$  is like OSL, is that we (<sup>define</sup> redefining) both the old & the new Object → 219.13

Spac

.00: 215.12 : So 214.19 - .40; 215.00 - .12 looks like the most Reasonable, Simplest Model.

.01 Its not far from what I was working on in Spac! But I now know 140.01 - 141.40.

.02 But what about the Guiding idea of SSGT: of  $f(C)$  that looks at a problem & decides on how best to work on it? The  $f(C)$  idea was also (part of) used in obtaining a default pc ~~of~~ <sup>of</sup> a boss to be used in reconstructing trials - as a function of the "problem type" - This seems essential in dealing w. the "Scaling" problem. See 141.19 for a little discussion of "Scaling" soln.

.06 Well, in QATM, there must be a part that recognizes what kind of problem is being presented, so it knows what <sup>sub-</sup>function to use on it. (The Ob-op algebra). So QATM looks at a problem & it has a "kuff Ob" to categorize it (Math or Physic or Chem or Linguistic type problem) Next it will have, within each category, subcategories to determine which funct. to use to get to answer.

Learning these categorizations is a regular QA problem! After TM has done it w. some success (or has been told the set of categories by ~~indices~~ "you problems indices" 205.09 discusses use of such "indices"

So, initially TM may have a PD for each category of input. These will initially be labeled as to "type".

.17 In .06 ff note 214.00 - .04 on difference betw. .06 & 177.16 - 37!

In both of them, we have Q as input & the Soln(s), A as output. In .06 the output is rather fast. In 177.16 we first get to an appropriate IA; which solves the problem - (but could take much cc to do so!). The final output of 177.16, hvc, is the ~~cheaper~~ function of .06.

One difference is that in 177.16, there are such fewer IA's than there are "solns" in .06. So .06 maps to a much larger space. In 177.16 there may be 10 or at most 100 IAs.

In .06 there will be a default soln for each kind of problem. In .06, ~~what~~ Q might be what is square root of 256? Answer could be a routine for getting Sq. roots. In 177.16 the answer might be a routine for optimization of  $x^2 - 256$ , or a general system for solving <sup>Solving</sup> working Alg. eqns. ( $x^2 = 256$ ).

So it looks like there is a "continuum" embracing .06 & 177.16: .06 focuses on a narrow part of Q space & gets a quick, sharp, answer. 177.16 makes a broader categorization in Q space & finds a IA that can work all probs in that space. QATM, regarded as an IA, would be an extreme form of 177.16: There is only 1 type of problem, & ~~the~~ <sup>the</sup> ~~max~~ <sup>max</sup> would be an extreme form of 177.16: There is only 1 type of problem, & ~~the~~ <sup>the</sup> ~~max~~ <sup>max</sup>

.30 all solved to same way (i.e. "put them into QATM"). First .06 is a specialization (special case) of 177.16 → 219.15

.31 <sup>large has</sup> ~~been~~ <sup>been</sup> a discussion of (apparently) "problems solving mode": Another Mode is "Updating", in which we give TM a ~~new~~ <sup>new</sup> QA set of QAs. This can be regarded by TM as a special kind of Q: But does there is way to tell TM when an Answer is "correct"?

I have been thinking of "Prob-solving Mode" as being non-trivial, that TM just applies whatever the current Macro operator is - to the current Q. It is conceivable that the output of that operator is stochastic - so there could be several probability used

6/12/01

IDSIA

219



The Artificial Intelligence Laboratory  
545 Technology Square  
Cambridge, Massachusetts 02139  
Telephone 617 661-1213

218.40  
00:215.40

"outputs", a each output could involve considerable  
calcul.

I had had of the diff. part as being "Updating": ~~the~~ Modifying the Stack Operator  
to a com. data. Augmented corpus.

1.04 T. (Most common) main case of complexity of calcul. in "prob-solving mode" occurs if we have OSL, or  
perhaps the Special coding of 217.14 (which is related to OSL)

Non-OSL prob. is prob. w/o. new data. OSL is prob. using a new data.  
In QATM it is implemented by looking for previous occurrences of parts of the new Q.  
Note that OSL is normally probistic: T. soln assoc. w. "One shot" need not be  
very long. — it can be  $\ll$  f. pc of some non-OSL solns.

In OSL, one does a bit of updating during ~~the~~ problem soln. phase (Also perhaps in 217.14.)

1.13: 217.21 217.14 is <sup>flimsily related</sup> (closeness ~~from~~ w. Pen. in min. distance) in Q & A. It says that OSL  
occurs if some thing "close" to the present problem has occurred ~~once~~ (at least) once in the past.

1.15: 218.30: After doing solving "a problem viz 177.16, we try to point into 218.06 form by narrowing  
down to initial pd on IA's (at first) — but ideally, it's a QATM problem, we want something  
narrower (our "soln") than a v.g. IA!

So: 2 things to understand: 177.16 vs. 218.06; 2 Prob solving v.s. Updating.

"Updating" can be regarded as a "regular problem", but only a very smart TM could usually work on  
that problem. — so we save it for TM's "more creative yrs."

~~Soln. I usually think of 177.16 for~~ So: It looks like Pen: often

When Q comes in, TM is asked to <sup>(recognize)</sup> decide what kind of problem it is a quickly  
derive a soln — which may take a long or short time. This would ~~be~~ 218.06  
If F() has no v.g. idea on how to solve Q, then its categorization of  
the type of soln. needed is more probistic — like over IA's.

[ but more generally, the "Q" can be any kind of problem dec. — I prob of  
all Q's being "induction problems" — while this is literally true, the main soln of v.  
"induction" problem might be an O2 or INV-problem. — In fact, the "induction" part could just  
be the part in which TM acquires an understanding of what the problem is. For a problem solved  
in terms of Q's that TM is very familiar, the induction part could be a trivial part of it.  
problems T. main part could be an O2 or INV problem.

1.36 Perhaps a more general idea as to what the QA problems: That A can be prob.  
solving method — ~~which~~ it would ~~have~~ have to be something like that for O2

1.39 problems — since we usually don't know the "correct" answer.

1.34 — [ But in fact O2 prob. are expressible as pure induction problems, if we require

17 15  
10 16 24 27  
10:13 -  
10:13 -  
10:14 - 6 37  
-43 37  
15-13 36  
57 38  
data  
177.16:  
F() → p down  
IA's.  
218.06: F() →  
"soln. of problems."

37"

# TSQ's 2 SCALING .19 ff



The Artificial Intelligence Laboratory  
545 Technology Square  
Cambridge, Massachusetts 02139  
Telephone 617 661-1213

## Some General Remarks

- 1) Quantum cost is probably same as regular cost. "Escape Modern Machines"
- 2) Resource limited / <sup>Quantum</sup> cost may be different from modern compa<sup>ring</sup>
- 3) Obtaining Practical Prob values is not nearly solved by ordinary RL cost. "bit more reason of this Q may be as diff as NP vs NP"

4) On Kol's conjecture that all of <sup>Modern</sup> Math is within a 100 or 200 bits of a certain time. Resource limited cost.

↳ Lavin's conjecture that distances smaller than Kol's, but still within Practical distance.  $10^{20}$  is now practical  $\approx 2^{68}$  "May"  $10^3 \rightarrow 2^{10}$

This distance is the amount of "choices" necessary to create Modern Math. So it was some time contact w. an Alien culture, this falls out. probably that they will have the same <sup>Modern</sup> "Math" as we do. But, there will be a smaller / subset of Math

(Not included in all of Modern Math) that will be closer to time. —

↳ We have much more likelyhood of having this "small core" in common w. t. Alien Culture. The "choices" involved in deriving Modern Math, may be have been made by applying Math. to Physics, chem, logic, sociology ....

5) When humans are shown Algebraic Expressions on paper, they know a man that abs position on paper is not important. This (2 perhaps other w) ideas must be somehow <sup>given</sup> conveyed (and by) to TM.

.19

6) On TSQ's 2 "Scaling". Each Conc must can be learned by LSrch, only if at least of each new Conc. is  $< CB$  threshold. This Least must be in terms of conditional costs of all of old concs involved. So when a new conc is discovered, we are always concerned w. its conditional pc: Under what conditions is it a likely "candidate"? Unless we have usable conditional pc's we have an unsolved scaling problem. Unconditional pc's tend to give a "scaling problem".

.27

↳ To some extent "vacancy" can be a usable "condition" — it amounts to a short short term memory. Say we use "X Windows": The "t" can be updated periodically, to give best results for entire corpus or for a "Super Window" w. some large, fixed "t".

This might be a good end of "conditional" pc for "earlyish" TM — utilities smart enough to discover better conditionals. It may be that I can study my own Prob-solving & get ideas for good conditionals — e.g. "Why" did I quickly decide that a certain conc. might be relevant in a particular (sub) problem?

↳ If I use minimal "conditionals" (I have to use some) & large CB for LSrch, I will tend to have more "Density" &  $\therefore$  more "Creative" Solns.

6/14/01

IOSIA

Optim. Technique

00: 219.40: OZ probs as being defined by being able to predict, well, how good each opt. will be for any particular problem for any particular CS. ]  
 Inv. probs are almost always solved as OZ probs, rather than by direct "Lurch".

03: 219.38 So, in line w. 219.36: The Q is mapped to a direct soln. A, if poss. If not, it is mapped to a method of finding A from — which might be R. 219.39 — 221.00 — .01 method of dealing w. OZ probs.

08 → Would it be of interest to do a QA TSP "like" SAARE — in which t. soln. to each Q was a relatively 'simple' funct. (that it could have ~~large~~ CC)? Since we are doing Lurch, we have some bound  $\frac{CC}{PC}$ , & find solns within that bound. — which was (idea of SSG is SAARB.)

To what extent would 08 be a serious step towards ~~SAARB~~

Terminology "T.M."?

17 (6/13/01) Consider several QATM models  $\hat{=}$  t. ~~TM~~ Corp (≡ TSP's) that they could deal w. / make out w. satisfactorily: Describe Update phase & prob-working phase. (They need not necly have a separate Update phase ... but ...)

20 Model (1): SSG, SARB; In update phase, TM tries to find/function that will map all  $Q_i$  into  $A_i$ . (This is a MTM). In prob. solving phase, it simply applies T's function to t. new  $Q_{n+1}$ . "Update" can work in at least 2 ways: (A) It looks at t. ~~complete~~  $[Q_i A_i]$  Bag & tries to find a function of non-cl. approach; probly Lurch. (B) It looks at t.  $F()$  for  $[Q_i A_i]$  & tries to modify it minimally, part way betw (A) & (B). (C) For some Back tracing, ~~try to update~~ try to update

simplest for (Mod 1)  
 1) Always search for entire corpus for updates (?)  
 2) MTM  
 3)

32 The  $F()$  of upto  $Q_k, A_k$  ( $k < n$ ) is by using sub corpus  $[Q_i A_i]$   $i = k+1/n+1$

34 Model (2)  $F()$  looks at  $Q_{n+1}$ , assigns wts,  $p_i$  to existing set  $\Sigma [A_i]$  of Inductn Algs. We do Lurch over t.  $[A_i]$ 's ~~to solve t. problem,  $Q_{n+1}$~~  to solve t. problem,  $Q_{n+1}$ . Updating  $F()$  ~~known~~ →

of maximizing t.  $PC$ 's of t. probly values assigned to OZ problems  $\hat{=}$  Corp is that obtained ( $GPD_1, GPD_2$ )  
 MCT analysis: We have  $F_1()$  &  $F_2()$  as in t.  $GPD$  of ~~SSG~~ SSG+. This updating on  $F()$  can be done (as in .22.32), by using t. corpus, or as small modification of t. old  $F()$ , or as a series of various sized "back track" of t. corpus



(t. IA's)

00:221.40: IN Mod 2, we also update t. IA'S: Each of them has its own update procedures (for both finite & infinite C's)

Model 2 (221.34)  
A common problem:

How can we minimally Modify  $F(x)$  so it works for "nt1"?  
The idea of "minimal Modify" of  $x$  to produce  $y$  is usually put in to  $F(x)$ .

This "min Modify" of old  $F(x)$  is an approx. best way to date for  $F(x)$  that fits all of t. data. A number of actual cases, stops  $F(x)$  that is good for all of t. data.

Should be updated for y sum of Least Sq. Err.

07:  $P(y|x)$ . My impression is that t. rat. machine used to compute  $P(y|x) \rightarrow P(y|x) = \frac{P(x,y)}{P(x)}$   
But isn't it true that  $P(y) = P(x) \cdot P(y|x)$ ? i.e.  $P(y|x) = \frac{P(y)}{P(x)}$  | No! This is RITE!

10: Anyway, t. P values in .07 can be simply obtained from ALP: In  $P(x,y)$  we want to find best codes for all  $x,y$  pairs.  $\rightarrow$  .18

11: Hrr. In Mod 1, we want a func that has about same values for most of a certain old part of its domain, & new values for a new part of its domain. There may be standard heuristic Methods of doing this. One, of course, is to design an op that recognizes t. new domain part, & a special operator ~~found~~ must be found to map back down section on to its proper ~~values~~ values. [SEE 224.18-.23 for more about decr of 'sum']  $\rightarrow$  224.10

12: Heuristic Methods of doing this. One, of course, is to design an op that recognizes t. new domain part, & a special operator ~~found~~ must be found to map back down section on to its proper ~~values~~ values. [SEE 224.18-.23 for more about decr of 'sum']  $\rightarrow$  224.10

18: Here "x" is  $f_{old}(x)$  & "y" is  $f_{new}(y)$ . When thinking about a metric from  $f_{old}(x)$  to  $f_{new}(y)$

If  $f_{old}(x)$  is very large, metric has lots of info about t. past corpus - that, presumably, would be relevant to creating  $f_{new}(y)$ .

21: We have to successive pairs  $f_i(x), f_{i+1}(x)$ :  $P(x)$  does give a p.d. - a stochastic

22: x sum. from  $f_1(x)$  to  $f_{21}(x)$ .  $\rightarrow$  224.00

24: Operator in GA. Ideally, w'd want a condl. mutation operator & "condition" being t. problem decn. For pair-wise "crossover" or "recombination" on all pairs.

Crossover: take t. set  $\{cond_i, G_i\}$  is  $|u|$ : is from t. set a p.d. ~~cond~~  $cond_j, G_j$

pairs. To do good Mutation & crossover, it is well to have a good language expressing t.  $\{cond_i, G_i\}$  corpus. - T. lang. should have useful concs in it

useful for creating new concs in any G range desired.

I did consider .24 off a lot in my analysis of GA, but I wasn't able to find a

26: good stack lang class to try to fit to t. data. It would seem that I write back to express  $\rightarrow$  set of  $G$  (funcs) as learning bany composed of a small no. of hy pc. ~~sub~~ sub functions. A smaller no. of hy pc funcs would generate many concs of any hyper ~~up~~ ~~part~~ ~~of~~ p.f. is t. hyper expected G. We would have to first search thru space

IDSIA.

# GA's

∞: 222.40 hypc Cands first. We just search thru Cands in PC order! We do try to find common "Subfunct's" in hyp G cands.

Perhaps the trouble w. 222.36 ff was that it long did not fit Y. data very well! — i.e. It generated lots of cands of hyp PC (i: hyp expected G) — but for the most part, <sup>very</sup> few of these cands had hyp G. — The cands of hyp G that occurred usually had almost all things had much smaller PC's — ~~is that~~ combinations of a few no. of "Subfunct's"

But, maybe not! I imagine that usually successfull functs are composed of only 2 or 3 or 4 sub-functs. — There are too many "large" functs to search over.

So the "fit" may not be very bad!

There are two functs. That are "very large" that are very good. — BUT they are discovered only if they have sub-functs that are discoverable "Subfunct's" of other hyp G cands.

There is something Bad about the forgg: We tend to use mainly the same parts of hypish G. cands. — which gives us low diversity. We may get an Excessive elitist system. Well, we could get subtle amounts of diversity by  $\delta$ -fying

the PC's w.  $\delta < 1$ .  $P_i \rightarrow P_i^\delta$ . We ~~probably~~ have to renormalize. — but it probably can be done — maybe not nearly uniquely! In general, given a set of  $P_i \rightarrow \sum P_i = 1$ ; there will be a min value of  $\delta$   $\geq \sum P_i^\delta < \infty$

One reason for wanting to do renormalize as we generate Y. Cands: If we do it right,  $\sum P_i$  will converge; otherwise they may not.

Say the subfunctions are  $f_i$  w.  $P_i = PC's$ .  $\sum P_i = 1$ , but we chose  $\delta$  lower and

so that  $\sum P_i^\delta$  converges. Using the normalized  $P_i = P_i^\delta / \sum P_i^\delta$ ,

we can generate new cands that are more like how PC's react over simple products of  $P_i$ 's of component functions. — These products need not be renormalized or stay sum to 1.

$$\sum_i P_i^\delta < \infty \text{ a.p. if } P_i = \frac{1}{i} \text{ then } \sum P_i^\delta < \infty \text{ if } \delta > 1.$$

Perhaps Go back to my (not so distant time) work on GA "longs" & some of the forgg. answers my objections to v. models I was trying at that time.

## IN "SGA"

I tried to find a stuck loop on Y. cands  $\Rightarrow G^\delta$  dotted.  $P_i$ 's of the long, as well as poss. For any G-funct, there will be a lower bound on  $\delta$  so that  $\sum G_i^\delta$  converges (see (18)) & (18R)

To find a  $\delta$  to put a  $G^\delta$  product. to "match" a  $P_i$  funct., look at  $\ln G_i$  &  $\ln P_i$  as funct's of  $z$ . If they are both linear in  $z$ , then  $\delta$  is the ratio of their slopes  $\leftarrow$  NO!  $\frac{1}{z} \Rightarrow G_i = \frac{1}{z}$

$\ln G_i = -k \ln z$  — which is not linear in  $z$  — it is linear in  $\ln z$  here. — so see if both  $G_i$  &  $P_i$  are linear in  $\ln z$

.00: (Spec 222.22): Inducing  $f_{n+1}$  from  $f_n$  P.d. from observing  $f = (f_i \rightarrow f_{i+1}) \cdot (i=1 \dots n)$  pairs.  
 is a rather gross way of getting  $f_{n+1}$ : Usually one does much "reasoning" to get  $f_{n+1}$  directly, or to narrow down to a few poss. choices, considerably. — Tho I haven't worked out just how this "Reasoning" is to be done. Logical reasoning; formal logic.

.04 (One way wgt to have TM do "reasoning" problems as part of its TSP — then it would have to realize that these techniques could be applied in other domains — i.e. its "own" problems solving — its "Meta problems",  $\Sigma$ )

It is analogous to  $\Sigma$  (occasional) discovery by students, that tech neeps used to

.08 Solve problems in schools, could also be used to solve problems non-in-school  
 .09 222.16 Spec: The 2 operators of 222.12-16 (could should) be related: ideally, find them both simultly

.10: 222.16 (Model 1): I Guess Mod 1 (221.20ff) won't work: The parts "reading" work"

1) ~~Update~~ "update" is always for  $i$  entire corpus, since  $F(Q_i)$  is for  $i$  entire Domain/  $Q_i$  values. I'm not entirely sure this is true: But it may be ~~literally~~ literally true

2) The "update" problem is not clearly solved! I have some kind of "soln": 222.11-16:

More exactly,  $F(Q)$  works for  $Q \in A_i; i=1 \dots n$ ; but not for  $i=n+1 \dots n+k$ .

- .18 To update we do 2 things: a) find an operator  $F'( )$  that works for  $Q \in A_i; i=n+1 \dots n+k$ ,  
 .19 using concs in  $F( )$  that work from  $i=1 \dots n$ : we want minimum cost for this  $F'( )$ .  
 .21 b) find an operator that can distinguish  $\{A_i\}$  from  $\{A_j\}$ . say  $[Q_i] \cdot i=1 \dots n$  &  $\{Q_j\} \cdot j=n+1 \dots n+k$ .  
 w. minimum cost. c) a) & b) create a new operator  $F''( )$  that works for  $i=1 \dots n+k$ :  
 Try to compress "this operator" (NOTE .09)

.23 I expect that as I design a TSP for Mod 1, I will find ways to solve  
 .24 it & use some of these soln. ideas to improve the update scheme of Mod 1.

Also, I expect to be using the AZI41 formalism to assign pc's to complicated functions:

(possibly using a variant of LISP or of my own "Functional Lang.": Mod 1)

.27 AS I see it. .10-.27 (i.e. ~~previously~~ preceding part of Mod 1) are pretty much what I was doing in SAARB. It was MTM. The probs I've solved since then that are

- impl. ① Better ideas on how to deal w. scaling problems (see 220.19-.20 for recent ideas)  
 ② That I only used defs of previously solved problem & primitive funcs: I think this was true because the TSP was very small, & I hadn't yet developed funcs that were complex enuf to need "sub trees" of generic funcs. : ~~at~~

③ In SAARB's solns of probs, I did use OSL much — but w. some assoc. w. soln. to a post problem was accurately put in the "STORE" as given pc = (total no. of funcs. in store)

[This, of course, led to a serious scaling problem.] I did not use OSL (or any other method) employing useful subfunctions (i.e. sub trees) <sup>inside</sup> Post FC) functions. (225.00)

**Model 1** (cont) .01  
summary

- .00: 224.40: Say I got Mod 1 running O.K. i got it to do a large TSC i put it to solve diff't probs.
- .01 Some direction of Genzen: (1) At all times, have > one FC on hand, so that mod 1 to deal w. new problems is more likely to work w. at least one of F. FC 1's.
- .03 - This may facilitate bit of NMTM.
- .05 (2) Do NMTM. some remarks: (2) AZI41 is already probistic - but this is a Meta problem for Mod 1
- (b) Finding Many FC 1's would be our way to realize NMTM.

(c) I may want to go to Mod 2 (221.34) to get NMTM.

(d) What are the Limits of Mod 1? .136 is BIG Limitation!  
 Normally, may solve as OZ problems, which are not MTM probs. Mod 1's method of operation uses AZI41 - which is a probistic system - sort of META level, Mod 1 is "probistic" - is this an attempt to solve OZ probs?

(e) Mod 1 can be modified to work probistic probs in folg. way: The FC ( ) that is sought becomes a devn of a stochastic operator. This devn can be single string ( it can be opset of strings - as in .01-.03 & .05). We hunt for this FC ( ) in the same way as we did before, using same update scheme!

(f) If Mod 1 can then work all kinds of probistic Q A probs, it is a "complete TM" can work on SI.

(g) See 214.19-40; 215.00-12; 218.00-01 for T. detn of Mod 1; that is "SAARBTM".

(h) (9) MI is really solving induction probs! When how does it know when to stop its search for a short code? - I guess it just picks the first code that fits any thing! This would be O.K. if it were looking for a seed code for the entire corpus! Normally it will do 2<sup>n</sup> searches to find f. 2 operators of 224.19 and 21. This gives a devn. that will "work w. whole corpus" - but

then we want to compress this "dual operator". T. criterion for "adequate compression" is unclear. Also, I have no idea on how to do it - But see 224.09 for streaming of an idea It integrates searches for 2 operators w. search for their integration

(i) (10) I expect that much of the uncertainty in things I got by watching TM is ME solve f. probs in T. TSC. Q

(j) (11) The method by which TM is supposed to discover & employ hours is unclear: I want to plan a plan observing on fly Prob. Solving (10)

Remember MI (K) A Big deficy of MI: I don't see how it can discover/use hours.

(12) MI uses AZI41: It's not clear how pc's of concs can be made to fit into an array nature of prob. may be used to solve. W.O. this feature, we have SCALING problem.

See 220.19 for devn of scaling problem: T. "vacancy" condition of 226.27 may be good evnt for certain kinds of MI's TSCs

on Model 1  
225.00

on Model 2  
226.01

ID

**Model 2:**

221.34 - 40; 222.00 - 70; 18 - .22; 224.00 - .08;

I Guess the Main Idea of Mod 2 is that we start w. this "Non-od." version of QATM:

$F()$  looks at  $Q_i$ 's outputs P.D. on  $A_i$ : we look on this pd to find the  $A_i$ ? — ~~usually~~

If this is a MTM: i.e. a NMTM a pd on  $A_i$  would be t. soln! ~~XXXXXXXXXX~~

A non-el NMTM:  $F()$  looks at  $Q_i$ 's outputs (conventionally) a pd on  $A_i$ : Now  $F()$  is a "Procedural P.D." so its output depends on how much cc it was given.

While this  $F()$  could, in theory, be found by search: (finding it is t. same as a General Operator Induction Problem) — one could do a search over Lisp (say) derived functions — (For MTM!): If  $F()$  is a pd, then we can also do a search for MTM, but we rather many of t. shorter codes for t.  $[Q_i, A_i]$  corpus. So we have a formal soln. for very large cc, for both MTM & NMTM.

For finite CB, t. problem becomes more interesting: In steady state, how does one do "updating" (we can start out w. a small corpus & do simple search to get t. many good codes that define  $F()$ ). ... So t. "steady state" update problem, is main

Problem. "Updating" means — How to obtain  $F_{n+1}()$  from  $F_n()$  &  $(Q_n, A_n)$  (or  $[Q_i, A_i]_{i=1}^n$  plus corpus Augmentation) Activity from t. code or "params" of  $F()$  See 220 for "conclusion"

In general, there are many IA's that can work on t. NMTM to NM-QAM problem — each will have its own updating Alg'n.

SN .17 is t. exactly how updating is done: Usually  $F_n()$  is dec'd. by short code set of params. It is these params that are "updated" by t.  $[Q_n, A_n]$  info. — for an example, consider

updating Barneolis, Z141, AZ141 — in each case, t. params that dec the "latest predictor" are updated — usually in a simple way. In fact, this seems like a good way to look at (many/most/maybe/all?) updating schemes. Before  $Q_n, A_n$  come in, there were certain kinds of rules that were dec'd by  $F()$ , & its short (codes/params) told how to do this. When t. new  $Q_i, A_n$  comes in, we will either have to see how it fits into t. old decus — i.e. may make "modifi. of params" necy. — &/o, it may make t. data of a new kind or very necy (as in Z141/AZ141) — so we add new parameters.

So Z141/AZ141 is a good model for updating induction in general — (at least t. way I consciously do induction / updating.) studying

An addition to changing params (freqs); & defining new params (new rules)! Updating can also repress t. corpus in over or more ways. Repressing is neither a modifi. of freqs or a new decus, but it can be a very imp't feature of "updating". It could be a new way to combine a bss to do indirectly or to code t. corpus

00:226.40: Look at some other induction systems & their Updating Algs.

01 <sup>Feed forward</sup> 1) Neural Nets: ~~subset~~ (This is a QA TM): "State" of system is given by values of "wts":

There are various updating Algs on how  $Q_n$  affects these wts. The "Partial derivatives"

~~part~~ is usually common to all, but Momentum & "learning rate param" are some of these parts of the update Alg. Normally, cross-validation is used in ANN runs. However, use in "fitness" criterion, cross-val is not needed. Cross val would probly make updating more diff — If we want to try to update by doing trials "near" to last "best fit". → See 34

2) GA: These are not normally applied to incremental learning — but one could (in at least

the foll. way): One has a corpus & one has found a population that "solves" that corpus

say its a seq. of QA pairs like  $[Q_i, A_i] \quad i=1..n$ : The population tries to get good codes

for  $Q_i \rightarrow A_i$  using a set of primitive (possibly universal) operators. We use to use

A (p) conc. function as a "fitness function". When a new  $Q_{n+1}, A_{n+1}$  comes in, we

start out w. the old population that  $\approx$  "solves", and we do the normal mutation/crossover

routine to solve it. A  $F()$  is defined by a population & its fitness values.

While  $F()$  is defined as above, the execution of  $F()$  could involve picking a

best candidate & applying it to  $Q_{n+1}$ ; <sup>or a set of it by fitness reads</sup>

so that  $F()$  is again not to solve as its params (or its desc/codes)

3) Decision Trees: I'm not (at the moment) familiar w. them: ~~But~~ say we have

some data (corpus) & we've fitted a good Dec. tree to it. There is a standard

MDL alg & trick for doing this: If we ~~add~~ augment the corpus we could

update by applying alg. anew to the augmented corpus — but the idea of

"updating" is to save it so we try to modify the old dec. tree "minimally":

~~Changing it~~ "Reprising" would seem to be too diff — but we might do various

small mutations of the Dec. tree to try to accommodate the new data —

This would be in the spirit of "updating".

34 4) ~~Update~~ Linear/non-linear Regressn. w. a finite sample, one finds a good pt in

param space. When  $Q_{n+1}, A_{n+1}$  comes in, one searches near the old pt. for an optimum

pt. to include the new datum. (ANN of .01 is a special case of NL regressn. — ~~open~~

functional forms can be used: linear sums, OR, <sup>linear sums</sup> functions, non-linear functions, or (Barrow's) better

ways to do non-linear regressn (of which ANN is a special case).

# INDUCTION ALGMS

- means "Universal".
- means ~~can be~~ Stochastic.

00: 227.40 Some where in Ross or DSA notes, I made a list of about 17 or 18 different kinds of

Induction Algms (IA's): can't find it yet. (It's 150.33-151.15) ← [ Has Good  $\Delta$  is Crosses! ]  
Has

Try to reproduce:

- ~ • 1) ALP w. CB=0
- ~ • 2) RLP (=ALP w. CB=00): A sort of LStoch.
- ~ • 3) Z141 (Bern Supos~~u~~ w. Datinidens)
- ~ • 4) AZ141.

07 ~ **3)** Linear / N-linear Regressn. Wiener filters, (Kalman filters?)  
 • **6)** ANN & Barron's Generalization.  $\rightarrow$  RANU [ ANU w. outputs  $\rightarrow$  inputs ]  
 ? • **7)** ("normal") GA & SGA & Sim. Annealing. Linear Regressn can use ANN  
Can be used to update IA's. See 331.13 for more details.

- ~ • **8)** SVM Support Vector Machines (Vapnik & Chornavensis); read Prof Tutnel I have in Postscript dir on C:\PS\
- ? • **9)** Decision trees (JDD3)

? **10)** Max Entropy Method (Jones) & partners  
 Ron Christensen Genz of it (I think it's ~~is~~) is not a special case of ALP. <sup>The</sup> Fader says it is.

~ • **11)** Grammar/induction (related to Z141 & AZ141) - can be Universal if Grammar class is adequate. I'm not sure, but I think + usual methods is

~ • **12)** Hidden Markov Models (related to 11) but more specific) techniques for optimization of continuous params: The discrete str. off. state space has to be (Guessed) ~~assumed~~ by the user.

~ • **13)** "Usual statistical methods" whether these are ~~to~~ include Stein's method.

~ • **14)** Progn/discovery in Sciences: This is example of N-L prediction models, but the horizons may be quite silent, ~~can be~~ for guessing good n-L models.

? **15)** Explain Base Lang  $\rightarrow$  277.30A checks this: EBL seems to be "compression"  $\equiv$  "finds Rups in data"

? **16)** Case Based Lang. [ related to: Generalization of OSL (On-Site Lang). ]

? **17)** Bayesian Nets (Part)

ID

Z141 Non-Universality  
AZ141 Non-Extensibility

225.36

380c Mod 1

08: 225.40: What is Min augmentation/modification of Mod 1 that can get (1) Time of hours (2) Conditional

pc's in AZ141 → (e.g. 225.37). T. Solns to Reson 2 probs could be closely related; nature of  
Often hours will be modification of L-gram P.D. conditional en/problem.

03 \* Poss. ways to solve them both: 1) In my working T.M.'s T.S.Q.; see what hours I use

04 ! See how they could be derived by "meta", i.e. modify M1 so it could do "that sort of thing"

2) A way to solve them both (perhaps implement .03-.04). Modify AZ141 so

6. pc's can be arbitrary functions of the past. (That AZ141 was not "universal" on that level, was an imp. worry I had about it not being "arbitrarily extensible" i.e. arbitrarily "smart".)

One way to do this would be to apply AZ141 to itself at a higher level.

10:23

12 AS is, AZ141 can construct an arbitrary function: it is "universal" in that sense:  
But it has a narrow set of poss. d.f.'s on that set of all poss. functs.

Somehow, by "applying" AZ141 to itself on a META way, we will be able to get arbitrary d.f.s on the set of all functions because this meta-AZ141 is "universal" in the sense of 12.

So this could solve critical problems in T.M. (1) 225.36! - Discovery of hours.

(2) 225.37: pc's w. arbitrary "conditionality" (3) pc's that can be an arbitrary function

(4) ATM should work on NMTM problems

all poss. functions! e.g. in Solata: In the "Allforms Marked", I want to be able to describe an arbitrary person "all poss. functs".

Stoffen 140.00-141.40 may be of interest!

Also N.B.! Realist time I worried about the Non-universality of the AZ141 d.f., I don't think I had any specific objections in mind: My objections were at "Re-architect" (C.E. Rink): It would be Great if I could find that stuff!

One reasonable way to work on this necessary Expansion of AZ141 is

.03-.04

Q: M2 did not (apparently) have d.f.s at 17 (1) & 18 (2) (I'm not sure about 18 (3))

Did it do this by using AZ141 on a META level? ← I think so! ← 19 (4)

On the lowest level, M2 worked on probabilistic D.F.'s (NMTM rather than MTM) — so it was able to use universal probabilistic d.f.s for SI (a meta problem).

Perhaps if I had a good way to give M1 a meta-universal pd, it would work a lot much simpler than M2.



ID

# GIAP

When I program a computer: It does what I  told it to do, not necessarily what I want it to do.

When I  feed GOALS into an Intelligent machine, these  goals will not necessarily be like  goals I would like the Machine to have.

More exactly: When I  feed goals into an Intelligent Machine: It will pursue  those

goals: Not necessarily the  goals I'd like the Machine to have.

This is the General problem of  Bugs:  Meta Bugs: A  Bug causes the program to not  perform as in the specifications: A  meta Bug is an error in the specifications.

In the foregoing " Goal" drifts; the error can be a  Bug or a  Meta Bug.  Meta Bugs are often very difficult to discover; they are most often found after the program has been applied to problems in R.W. These  Bugs can be very expensive, sometimes catastrophic.

My guess is that a  Bug in a  Goal, is likely to be a  Meta Bug.

— VERY  BAD, very  expensive, potentially  catastrophic.



Marrin's idea: That a Machine is much smarter than us, we should not be telling it what to do. Reply: A person could be much  smarter than me, yet if what he wanted to do was  very bad, I would oppose it. E.g. he might be a  homicidal maniac — He loves to build " Doomsday Machines". Since he is much smarter than me, what I perceive as a "bad action" on his part may not be bad at all — since he understands the action better than I do — on the other hand, his  Goals may be quite different from mine — he may not be concerned w. my  goals at all.

6/21/01

231

ED

See 228 for list of IAs.

# CREATIVITY

00:

Continued - I scanned from 231 back to 208 - <sup>found</sup> ~~found~~ nothing on "creativity"; 166 is abit relevant

: previous "Recent" note: An Experienced Good Scientist, will have lots of Good hours. They solve most problems very fast, but ↓ <sup>in</sup> ~~the~~ diversity.

A smart student (1) has fewer hours (2) ~~has~~ more diversity (3) Very often can spend 24 hrs/day; 7 days/week working on a problem. Older Scientists do this less (The apparently, Newton was able to do it when he was 74yo. !)

Poor Scientists have a few hours: They use them as far as they can go - then ~~if~~ if problem is not solved, they can do nothing more.

ID

: 229.40: **Another tack**: Instead of using an unordered corpus (as in M1), use a purely sequential corpus, but have it be a sequence of Q/A's. The advantages are:

- ① I don't know a good "super universal" d.f. on all possible operators, but I do know a universal measure on strings.
- ② If TM finds a good answer for  $Q_{n+1}$ , it can make a "recognition" of  $Q_n$  be simply a pc of it.  $n$  integer,  $n$ . (later on this can be expensive, hvr.)

Disadvantages: For QA TSP really is mainly unordered data! The answer is same sequential info in it.

Advantage ③ If I get a soln to this sequential QA problem, it is likely to be very useful in solving the unordered QA problem

Thinking about the sequential induction problem: its usual ALP soln, reminds me of how to do a "super universal" operator d.f. Formally to get a universal operator! — do it like a "sol" 2 kind of problem. did it with BAG induction, for Next OSL problems,

Given a  $\Phi$  input unc.  $\rightarrow$  Input 1 is chosen, 2 is input to operator, 3 is output of operator.

Input 1 is chosen in addition to  $Q_2$ , to get  $A_2$

We try to find an input  $\Phi$  (finite string, always)  $\rightarrow$   $\sum |Q_i|$  overall  $Q_i$   $A_i$ 's  $\rightarrow$   $+(|\Phi|) = Min$

We do this for several  $\Phi_i$ 's for which this sum  $(|\Phi|)$  is smallest.

In "Lispish" notation: We have a code for a string that defines an operator for  $Q \rightarrow A$ . The pc of this code is obtained via AZ141. We want  $S_i$  to have small  $Kcost$ , say, so  $C_i$  should be small.

I wanted  $C \rightarrow S$  to be via "Lisp": or some universal ~~operator~~ device!

I want the pc of S. In M1, the pc of S was obtained in a simple Bernolli way via AZ141

I want a better PD on S. This PD will be obtained by observing ① context dependence of PC's of sub-functions of  $F(F(x)=x)$  ② Heuristics that modify pc of  $F$

Sub-functions of PC) — (37)

**SN** In General, whenever a new defn is made, we must also make rules that

tell when that defined object has reasonable PC (i.e. solving problem most usually be addressed as soon as when'ever a new defn is made. — The sometimes to S2 wouldn't be large and — we could, sometimes use OSL hvr.

In line w. 00-30: Perhaps ② essential features of M2 that make it

- ① "Super Universal": TM<sub>1</sub> is able to do probabilistic induction (ENACTM)
- ② we have a TM<sub>2</sub> (which = TM<sub>1</sub>)

232.40: Hvr, TM<sub>2</sub> has to be constantly at work in each TM<sub>1</sub> induction problem, if we are to have much use of (a) context dependent PC's of abss (b) (rug of hours). Hvr, even if TM<sub>2</sub> is only invoked periodically, w. OSL, a "recency" dependence of both "pc's of abss" is "useful heuristics" (similarity of structure of previous problems) we will get some useful effect from TM<sub>2</sub> in these 2 areas.

Ideally, I'd like to "superuniversal" p.d. to be implemented w.o. having to invoke TM<sub>2</sub> & TM<sub>1</sub> separately; but as a single partial recursive function. On the other hand, my ~~own~~ observation of my own Prob. Solving, will probably be in the TM<sub>1</sub>, TM<sub>2</sub> - 2 level form.

It is certainly poss. to do TM<sub>1</sub> & TM<sub>2</sub> simultaneously, when we learn in Time Share mode.

- There seem to be at least 2 ways of thinking about the "Superuniversal" p.d.
  - The idea of 232.19: a universal d.f. or S, which is to code for all poss. functions (i.e. operators)
  - The TM<sub>1</sub>, TM<sub>2</sub> idea, in which both TM<sub>1</sub> & TM<sub>2</sub> are of the A=IT1 form.

Clarification of .11 - .14 would be a most critical problem, allowing the creation of a very simply defined, ~~superuniversal~~ superuniversal TM! (There are some Other details - like what part of the corpus to use for a particular problem, even a certain C.B.)

17 **SN** Marcus' approach to a Reinforcement Machine! This could give RTM could begin to same as a M1 or M2. Since M's approach is Non-el - (i.e. complete ALP), it may be poss. to use his approach for a "final F.M." T. analysis. Thvr, would be much simpler than a real, general reinforcement p.d. using the QATSQ (as for M1 or M2) Since f. envt. is much simpler, it would need to do ~~no~~ "experimentation"

A difference & below: RTM & QATSQ work w. 2 codes, but is a set of Q's (No prior Q's) RTM works w. Q's only; Pign's told how good its replies are. Hvr, a RTM can begin sets of Q's selected by 2 simple Q's. Hvr, I suspect that Marcus would want TM to have just 1 ("Best") response, & Give A=1 for correct, R=0 for incorrect.

24 **Exp. B** "experiments, Rytis" [RTM's response to a Q is a p.d. on A's; its reinforcement is in PC of context "A". Max PC<sub>Actual</sub> = max PC of corpus (i.e. Max Likelihood Method?)] Also think about E. Min Backtrack machine! Say Euse & Lisp, & after each successful run, E puts all data at beginning of t. code. (like AZ141) Then I think the system would be very much like AZ141: E would it differ at all? - in what ways?

30 If a code does not use one of its factors at all (directly or as a sub-factor) part of code used better? Then use success process before mind? - 234.07

32 If the "Corpus" consists of both the [Q; A] in the T. process at hand in a temporary (1 level) to predict, then I think any Universal language code (w. P.d.) There are (at least) two big problems (1) we don't want to code the entire Corpus for ~~it~~ it we have must do a Predn. w. finite known CB (2) we have no "Dynamic Form" problem that Marcus is ~~now~~ writing about - but the actual cause of "backtrack" is our own traces. We're at all times Making the dynamic "Experiments". -> (238.35) 238.35

Whoops! we have to subtract out best of the data! so it's not exactly "Max Likelihood" 238.32

G(23/08) ID

- 00: 233.70: **[NB]** "If all info is in P.D. Then L such is  $\approx$  optimum" This M to be really true
- 01: If "Corpus" includes all of TM's Traces! No! For close to optimum, TM must be able to notice what happens in early traces of an L such, & use P.D. in the future in the L such: P.D. is not a simple L such based on a fixed
- 02: **[SIN]** In <sup>re</sup> parsing a corpus in (A) 141, after one has made a new definition: ("All info in P.D.") Instead of using old P.D. of utups, (i.e.  $p_i$ ) use  $p_i^x$  ( $x > 1$ ). This pushes the code toward hyper pc: The code of least pc has all  $p_i$  to same ("Max entropy").  $p_i^x$  tries to make p.c.s as different as poss.  $\therefore$  hyper total pc.

07: 233.31: 2-141 does seem much different from <sup>Track</sup> Min Back Track Coding; In 2-141 I automatically sum over <sup>very</sup> many codes. This is not done in <sup>My</sup> ~~Min~~ <sup>Most recent version of</sup> Min back track coding. In 2-141, the parallel codes are all strictly equivalent, so there is no poss. loss by "pooling" P.D. statistics — no (info. loss) "Decision" is ever made (as per the parsing decisions)

**Note on MDL v.s. ALP:**

Consider Arithmetic Coding: In ALP, we have an infinite no. array of codes for each finite bit string. We can approximate this, by taking only the first  $n$  bit approx. of true ALP pc. Because of the padding property of  $p^x / (1-p)^{1-x}$  (i.e. it padding when  $x \neq p$ ), this will give an error per bit of an amount  $\approx (p-x)^2$  (R.D. smooth will always be  $> 0$ ). So for a long code, of length  $n$  it should be no. overall extra bits in the approx. will be  $\approx n$ .

However, the purpose is not to find a short single short code for a long corpus. The best way is to use as high precision as one can for each bit of each bit, then use Arith coding. This would seem to require infinite precision, in order to get a small out. of bit loss in the final codes.

perhaps has figured out how to deal with this — So arith coding is practical, IE's version

Elias coding is not! This is discussed in I.T. May 2001 p.1533

Vol. 47, no. 4 W.D. Withers

How ibid. (p1533 col.1) he gets  $2^{-9}(q+2)$  extra bit per symbol if he uses  $q$  bits of precision! I should think he's got  $\approx 2^{-29}$  extra bit per symbol.

Heads va to also says " last line of it " Various clever make-shifts are known for handling the carryover problem. It may be that there is a problem in Arith coding that I haven't really been addressing! ●

ABCDEF  
1111

.2mm  
=.02cm  
=.20  
=.02 "  
2.54 "  
=.008 "

ED

# Sub Corpus Problem 1.00

Out Goal for induction; If TM is allowed to use different parts of the corpus for induction, Goal is unclear (?). If the problem is  $\frac{find\ pc(i)}{PC(i)}$  versus, then, presumably, corpus for num, is domain. is same. Say we spend same cc on selection of sub-corpus as finding codes for corpus (numerator = domain). Code finding starts for Num/denom area "internal" (? - is this possible to determine useful way?)

Anyway, so Post Mark, could one ~~legit~~ legitly determine optimum strategies for corpus selection (denom) loading?

It may be that Semi-computable (enumerable) functions are definable iff.  $\leftarrow$   $CB$  is given, but "non-enum" (ratios e.g.) are not uniquely defined.

If we use a standard approach, of a time,  $CB$ , it would seem that ratios are well defined.

But perhaps, if one uses different strategies (stand for num, denom) then they don't converge to same thing (?) (but limit of ratios = ratio of limits).

229.17 p. 40  
232.37

.14

A poss. BAD way to do .01: say corpus is binary Bern. seq:

We select sub corpus of all 1's — this gives sharp d.f.  $\epsilon$ , is looks v.g.!

Troublers, the specification (den) of the sub corpus is very expensive. SO!  $\leftarrow$  IS

This a big clue to the sub-corpus selection problem? Consider all ways to deriv. sub-corpus:

Each way such Deny has its own wt ( $\epsilon pc$ ). We could use all of these sub-corpus:

the resultant predn way  $pd$ 's (on 0,1). So, in "normal" ( $cc = \infty$ ) induction, we could also

divide up the corpus & make predns using the parts — is wt. those denms or  $\epsilon$ .

pc of the den. of the part. This would seem to be not the same as Regular ALP

Tho, if there are sub-corpus that are useful for predn, then this is an interesting kind of regularity that ALP should be using.

.26

.14 - .26 is a novel take off on the "Sub-Corpus problem", but Return to it later if I get new ideas on it.

A rather el. "soln.": Initially, all problems are given several "indices" by USR.

From these, TM is able to tell how to categorize the problem to some extent, what sub corpus to use (Pao Pao's is not yet linked w.  $\leftarrow$   $CB$  given for the problems). Later, the USR is less careful about assigning indices: TM learns to do it itself.

Now this all seems related to the Scaling problem: Assigning  $pc$  to concs, that depend

on context. "Context" seems related to the "indices" of problems, just as  $\epsilon$ -pc of a concept is related to its context: The  $pc$  of a prediction is related to the context of that predn. for concepts, we did pooling of into over the context of that conc.

for  $pc$ 's (predn) we have pooling over the "context" of that predn  $pc$ . Hvr, for predns. the "conc" is the 0,1 pair. The issue sub corpus depends on the local context of the (C,1) pair. This might be the

"QA" — or just the "Q"

ABCDE

ID

# "Epicurian Codes" (.09)

## "Updating" & "Summarizing Machines"

Rep =  
Rant.

On "Summarizing Machines": Perhaps one way of looking at "Updating": Predit produces a kind of "Summarizing Machine". (~~Summ?~~ Summ?)  
 One way to produce them was outlined in the "2 kinds of Prob. Ind." paper: We use a set of Machine States associated w. the shortest codes found. [My present impression is that this is not a v.p. approx. method! that it usually has to be very large in order to get a reasonably good pd. ...]

~~That is~~ e.g. for a Bern seq. w. a alphabet containing m symbols, - each w. a diff. pc - having ~~just~~  $k = m$  would not express the true PD very well - [Huv, I reached this "conclusion" by considering the Z141 code: (which is a kind of "Epicurian Code" that uses lots of 11 codes: I'm not really sure that it's true)]

A "partial" "Summarizing Machine" could have a (relatively) short corpus  $C'$  plus, a special Machine,  $M$ , so that the pd on the new chunk of corpus,  $X$  is the pd induced by Machine  $M$  on the string  $C'X$ .  
 In the "2 kinds" (1999) paper  $C'$  was the null string.

**T.** more general concept is that of a PPD "practical PD" (a PD w. cc as an input parameter - T. result of the "Updating" is a useful PPD - (usually a "PD" ... Since  $cc$  is usually small <sup>or</sup> acceptable or specified in the PD defn).

In Z141 & AZ141, there are 3 components of Updating:

- a) update freqs of concs, & reparse.
  - b) Derive new concs, ~~update freqs~~ reparse & update freqs.
- If (b) is not done, then (a) is done.

It may well be that most methods of updating IAS are close to or analogous to these 3 operations. 228.00... gives a list of 14 IAS.

Also Note AZ141 can give a simple meaning to "MM Backtrack" - i.e.

w.o. Backtrack, we ~~track~~ simply update pc's in ~~response~~ new data, & do a rather greedy parsing of the new data & update the pc's.  
 Next in order of ~~cc~~ <sup>next cc</sup>: Find new concs, reparse corpus, ~~reparse~~ recalc. PC's  
 Then maybe reparse & recalc. & update pc's again

- .34 Next in order of ~~cc~~: remove last definition After new data is obtained, remove last definition
  - .35 ~~try~~ <sup>after new data is available</sup> to find new definition, reparse, recalc. pc's; Maybe ~~loop~~ loop to (2) (.35)
- w. more cc available, remove last 2 defns, & do like .34-.35.

→ In some cases, it may be wise to remove certain defns that are not the most recent ones

7. Subcorpus Problem: Relation to Scaling Problems of Context dependence of pc's of

00: (SPEC 235.40): So TM looks at PC indices of t. current Q: All other Q's w. same set of indices are certainly

"relevant": Q's w. all but one indices t. same as t. present Q, are next in line (if CB)

logs on it): Then we have next consider Q's w 2 embedded indices.  
T. logg. is v. ruff: Certainly we should consider such certain indices

should be more impt. than others.

05 Now consider context dependency of conc. pc's: Say we have a Q & we want to construct an operator that maps that Q into its A. T. Q will be t. "context" & will have a associated w. it, a p.d. over all possible states & thus for. (In some cases, it may be able to give pc's for concs not yet defined or used. — it does this by extrapolation from t.

stochastic log. assoc. w. concs that have already occurred in their pc's in t. ~~operator~~ operator that created E.A's. → 16

Any way: 1) T. Subcorpus problem is t. 2) context dep. assignment of pc's to concs (or to other things)

Should be declared clearly: To what extent is TM "solid", if I solve these problems?  
Other impt. probs: 3) correlated pc's of concs. in LSTK. 140.00-141.90 lists 12 ideas "past" 50/89

16 10 On context dependent pc assignment: This assignment should be CB dependent

(like subcorpus assignment). If CB is small, t. assignment should be very narrow: to a small set of situations/environments. If CB is large, we can afford broader characterization, but it will lower pc's in most likely (domains environments) in order to get more "diversity"/"creativity", by assigning not such small pc's in less "prime time" domains. [So this ties 1) & 2) to 11 together, more]

On the "similarity" of 1) & 2) (11): In t. subcorpus problem, we are given a problem (Q0), a variant of CB; C0: we want a subcorpus; set of Q; A; pairs in t. past. || for context dep. pc of conc: Given a certain conc/abstraction S, X: we want its pc as a function of t. context in which it is to be used. This "context" might be a problem (say "Q")

30 (SN) Some confusion in my mind about "subcorpus" problem: When it occurs in "updating", we only have QA pairs (no isolated Q's): In our version of QA, when an isolated Q occurs, we always have available a "updated operator" that can map that Q into a d.f. on A's. [Hvr, it only has certain CB, C0, & that I may be inappropriate (?).]

34 In an Alternative version: TM trust present stock operator on t. Q — then decides that t. d.f. isn't "good enough" ("sharp enough"), so it does more "updating" w. what it feels is the "relevant subcorpus". But normally, the "subcorpus" is needed only for updating a set of QA pairs. In 34, hvr, we have an isolated Q & TM finds an appropriate (small(?) set of appropriate QA's) to update — which gives another subcorpus of QA's to include in t. updating process. This amounts to backtracking "by subcorpus" Q

So TM must be able to find a relevant associated subcorpus for either a single Q or a set of QA's  
either a single Q  
SPEC 240.37



ED

Outline of Letter to Marcus (ca. to Juergen): "RANT"

- 1) How Q.A. machines able to solve any problem that <sup>would</sup> be introduced in practice?
- 2) I am intrigued in a minimal machine that ~~can~~ be arbitrarily intelligent.
- 3) It is Q is a machine of this sort. It does not explicitly have a "horizon" problem, since such Q.A. is sharply delimited. ~~It~~ When such a machine becomes sufficiently intelligent, ~~it will~~ <sup>it will</sup> ~~introduce~~ <sup>introduce</sup> ~~new~~ <sup>new</sup> ~~problems~~ <sup>problems</sup> that I am unable to solve. (I was able to define these problems, but not solve them).

4) One of the attractions of a AI system that has no adjustable parameters is that it is more likely to be considered "intelligent" by virtue of the AI concept. I am very much not interested in proving that machines ~~can~~ <sup>can</sup> be intelligent. When they are very usefully developed, this fact will be quite apparent. In all respect. (and for people) there will be many people who will be surprised. The reasons are not closely related to AI or its goals.

5) It is my impression that reinforcement learning is "Loss-Loss" problem. If you don't succeed, you lose; if you do succeed, there is no reward to you or the machine manipulating the world. While many people will accept this "danger", I will not. When ~~the~~ <sup>an advanced Q.A.</sup> machine is given a Difficult / Dynamic / Long problem, there is danger that the user being manipulated. <sup>or simply a very difficult problem, which interaction w. R.W. is a Needful.</sup> Thus the misdeeds, I think, is misdeeds to realize.

It is a general reinforcement machine. If any <sup>intelligent</sup> machine is given "Loss" problems, then

.31 6) see ID 233.17-24 a ~~idea~~ (237.17-32) R for comparison of RTM & QATM.

.32: (R33328) Hvr. One could regard QATM as MAX Likelihood, is to pc of the model (spans) is included. One trouble is that the single best (model, corp code) would be selected in Max Likelihood - or ~~by~~ <sup>By</sup> RTM. I think QATM (would/could/should) keep several of the best models, not only the "Best", & perhaps could try to have Max diversity over 6 models.

.35: 233.40 On "No experimentation" in QATM! Actually, QATM does do "experimentation" if it normally has "Dynamic Penalty" problems. These occur in ordering trials on Gnds. <sup>(Non-greedy, ordering of trials "a Dyn. Peng. problem")</sup> Hvr, in early training of QATM, I will use very simple WATM approximations to the solns. --- in the hope of eventually getting to pc. at which I could give it its "Dy Peng" problem to TM. (essentially - it ~~can~~ <sup>will</sup> become part of "SI"). So, to start, TM's order of trials will be completely or "fairly" Greedy.

233.40

.00: 238.40: Also Note: The Dynamic Prog. Problem Marcus is working on: Including CB (as well as fixed horizon)

is extremely diff. To do anything like optimization w.r.t.  $\Sigma R$  is fixed horizon & fixed CB (that includes all of t. poss. paths), seems really formidable. (At first glance, I thought it was the Unsolved G(X). T(X) OZ problem - but it's apparently not.)

I don't really know just how far Marcus has gotten w. this problem. He says it's within a factor of  $2^t$  of optimum, but I'm not sure I know what  $t$  is.

.06 If he did get something like an Lsach soln, this could be a big step toward a usable "Heuristic" / "TSP" soln.

.07 **REV** 232.11 - .40; 233.00 - .40; 234.00 - .01; 238.31 - .40; 239.00 - .06 is main line in developing QATM:

Other recent topics: 232.00 Using a sequential (rather than a BAG) corpus for QATM

- 234.02: use of  $P_i^8$  for pc's of emes (renormalize)
- 234.07 MDL v.s. ALP: Peter Elias Code v.s. Ari-R code (234.26)
- 235.00-01 Sub-Corpus Problem.
- 237.00-01

236.01-04 Up Dating: "Summarizing Machines"

238.00 - .40 Differences betw. Marcus/Teegen's RTS & my QATM: Advantages/disadvantages!

Back to Line of **.07**: Say I had  $M1$ : a QATM & MTM: AZ141; "SwarmTM"

Its serious drawbacks: 1) No NMTM problems 2) Inability to determine hours (since they are NMTM probs)

- .18 3) pc's assigned to ~~cones~~ are not context dependent - so serious scaling problem.
- .19 4) Perhaps inability to work on "subcorpus problem" (maybe related to 3: see 235.00-04, 237.00-10)

If we allow  $M1$  to work on NM probs, we may get a large jump in its capabilities - a ability to deal w. the other 3 diffys.

In  $M2$  we solve "Gloss over" t. problem of NM induction, by saying "Here are many soln. methods - each w. its own updating scheme".

To start off, we have set of IAs: for each QA, we use all of them! After we have our data, we give TM t. (Proto) problem of ~~making a P.D. from~~ being given a problem & finding a P.D. over IAs that are most likely to be the "best" soln. opt. problems. ("best" pc of soln. <sup>correct</sup> in given CB),

TM now uses this P.D. + Lsach to find good P.D.'s for QA probs, w. given CB. (I think ~~the~~ <sup>correct</sup> mobility to do something like this, was a feature of  $M1$ )

TM also can work on (3) (18): see 235.01-04 & 237.00-01 for Gen. disch. See 237.05-08 for how this <sup>might be done</sup> (This disch. is a bit vague, hvr).

TM can also work on <sup>(19)</sup> subcorpus problem 235.01-04 & 237.00-01 237.00-04 in particular

TM can also find Hours (Lquess) since this is a NMTM problem! How TM can apply these hours to t. solns of probs is unclear! One approx. way is to work on ~~given~~ QA probs & finding hours in parallel Hours too almost always expressible as Mobile of t. P.D. for Lsach, so we will be modifying this P.D. during t. Lsach (which makes it a kind of Modified Lsach) (of a search method inspired by Lsach).

A main diffy of  $M2$  is that when it uses Lsach, t. ends can be highly correlated, which leads to very non-optimum results. A big Q is: would such a TM be "good snuff" for serious "Bootstrapping"?

A similar Q is: Could a TM w. facilities for OSL ever do serious "Bootstrapping"?

*i.e. one that didn't worry about "correlations"*



"Greedy vs. Parsighted."

240.40 : solns  $[A_i^k, PC_i^n]$   $k=1, \dots$  Say  $A_i^5$  is t. correct soln. The  $A_i^5, PC_i^5$  info will be used in updating. Also t. techniques used to find this soln. will be ~~used~~ given "special treatment" in t. updating process.

Think of 2141 induction: T. "Q" is t. latest "chunk" of a time series. (say). In this case, we would have to update before doing a prodn: T./TS contains A's as well as Q's; well, say we have to update T SQ already & we have to make a prodn: we use t. primitive, defines & pc's of t. last update, to get a pd. for t. next symbol. [Here we assume no "OSL"]. When we find out what t. correct next symbol is, we can update t. EA by changing pc's, defining new symbols, & re-parsing.

On "Parsing" rather than "Best, or "it has A" phrases may not be adequate!  
 try to find "Covariate/Guar" It is below - 4/1/01  
 better I left for Europe  
 LA HURSB 243.00ff  
 Or, when one has 2 very large CB for parsing

03 brings up a problem I've not worked on. That a given EA could have several subcorpi that it works on, & each will have its own update parameters. Th. subcorpi will also have some "shared" parameters - As with STEIN affect.

So  $F()$  should be chosen a (EA, subcorpus) pair. Each such pair has its own update params; (t. some "shared" params (11)). Just when t. "shared" params are updated/inserted/deleted; is unclear. So "Grand"/"Macro" update argy

This idea is Brainiac in w. Y. design of 240.10-40: It suggests that 240.11-13 may be t. best "E12n". (rather than 240.31-32)

Note, how, that in t. large "Time Series" problem, t. parts of t. <sup>sub</sup>corpus were all clearly labeled ("indexed") & s. belonging to t. same "Time Series".

A (subcorpus, EA) pair that has been updated is like a "Summarizing Machines" - so for a problem (Q);  $F()$  selects (1 or more) "Summarizing Machines" to solve it.

I have been (Most recently) mainly thinking about 2141 & AZ141-type induction problems. Clearly there are kinds of problems that look much different! Eg. Solving/equ, & symbolic integration, (A long sequence of operations that solves t. (INV) problem.). (Some TSQ were short sequences of operations that solved t. problem. Q is: Are these much different from t. ones I've been thinking about?

Also, there is a general class of problems solvable by "WON" techniques (AND/OR nets of subtasks). In AZ141, t. problems are all solved by ~~some~~ <sup>My Modified</sup> LS such over combinations of a set of symbols in which t. symbol to use is "context dependent". T. Q is: can all problem solns. be expressed in this form? If not, can this system do SSI & find a more powerful system that can solve "all" problems - or can this system recursively reach an arbitrary level of prob. solving skill? Or can any conceivable <sup>prob solve</sup> heuristic be found by this system?

My ~~thought~~ about this system has mainly been with 241 & AZ141 type MTM problems (21!) M2 has to do NMTM problems. But AZ141 has 2 levels: It is a P.D. on t. set of all functions. Think of this stochastic level; & there is t. set of all functions - which is ~~MTM~~.

IS

Def

AZ1

00: 241.40

1) AZ141 is a complete solution to MTM problem, since this is a wide def over } 245.18!

2) (deterministic) operators: So we can solve any MTM prob. w. LSrch.

2) Since MTM & NMTM have same soln. (in ALP) this may solve NMTM problem!

3) For NMTM, I had 4. idea <sup>to finite CB,</sup> RLP was not to soln. <sup>clear</sup> first for finite CB, selecting part of corpus, then using RLP "iterate work" - but <sup>clear</sup> it was not clear how to do <sup>sub</sup> corpus selection <sup>clear</sup> w. the proposed method, ~~idea~~ of ~~result~~ - corpus selection, it was not clear that this was optimum.

4) <sup>clear</sup> 1) & 2) can be for  $CB = \infty$ , so even if 1) & 2) "work" is solved for MTM for  $CB = \infty$ , it would not necessarily be a soln for  $CB < \infty$ .

5) Re 1) & 2) exam: It should be poss. to develop a complete theory / deriv of Universal

~~stochastic~~ operators (analogous to Universal Distributions) using AZ141

3U AZ141 differs from LISP, in that all definitions occur in the pre-corpus. - Other than that, I just use a "minimal" LISP, w.o. its specific "definitional" facilities. [R.R. got no "lambda" notation]

6) There remains 4. possy. <sup>sum over</sup> "Short codes" is not a very practical way to get P.D.'s. So, even if I found a way to get AZ141 to give a universal PD on stack operators, it would not be a V.G. way to do this! Still, it might ~~be~~ suggest ways to do it.

One (ALP) way to Look at "Universal Operator"

We have  $z_u$  <sup>(U/O)</sup> w. 2 inputs, 1 output: Input 1 is "Q" (input to operator); output is output of operator. Input 2 is a random signal. <sup>short string</sup> ~~The PD on operators~~ so  $M(Q, R) = A$ . With a given Q, the prob that the output will be A is  $\sum_{R \in \text{finite strings}} P(A \text{ is } Q \text{ and finite strings})$  is defined to be the <sup>relative</sup> prob of A for this model.

An equivalent model: <sup>to 17-21</sup> 3 inputs to Machine: Q, O<sub>i</sub>, R; output is A.

Q, R, O<sub>i</sub> are finite strings; R is random string. O<sub>i</sub> is <sup>stack</sup> operator from.

The rel. prob of A being printed from Machine stops is the rel. prob of  $Q \rightarrow A$  for operator O<sub>i</sub>.

In 17-21 ~~was~~ <sup>sum over all</sup> O<sub>i</sub> we write  $2^{-|O_i|}$ .

The long R of a random input that creates A is then effectively added to  $|O_i|$  to give the rel. assoc w. that random input.

Can <sup>idea in</sup> 17-27 (a especially 22-27) be used to get a stack operator & a universal operator, out of AZ141? Using AZ141 notation, we can have a 2 input machine: One input is Q, other is O<sub>i</sub>; but O<sub>i</sub> is a ~~stack operator~~ MTM operator. If we put random input into O<sub>i</sub> we get a stochastic operator (if probably a universal / stack operator) - If so, we can

could divide ~~exam~~ 4. random part into R and O<sub>i</sub>, we'd have a stack operator "derived" by O<sub>i</sub>.

Or AZ141 could have 3 inputs: ~~control~~ Q, and R. R could be regarded as "an extension of" O<sub>i</sub> or of Q; or just a separate input.

In 22-27, if  $R=A$ , then the output is the probability A: Various values of R can be regarded as derived derivators from this machine. - (This is "MDL to A"! ) - Not so! If we sum over many R's we <sup>sum</sup> 244.00

# CURE CANCER Problems

## Cure Cancer Problem

ID

.00: 242.40

I haven't been able to find my previous soln ( $< 4/1/01$ ) in these notes:

A new try:  $i$  is time index;  $X_i$  is TM 2.0m,  $Y_i$  is corresponding response of patient.  
 $[X_i, Y_i]$   $i=1/n$  ~~is a poss. sequence of treatments~~ interactions.

.03 ALP gives us a p.d. over all such interactions for all such interaction sequences! (see .12-.13)

For  $n = N$ , we have assoc. w. each sequence an  $R$  (Revenue) value.

We want to make choices so as to maximize Expected Value of  $R$ .

.06 Given history up to ~~is~~  $i=N-1$ , we can always choose  $X_N$  to optimize given  $E(R)$ .

.07 " " " "  $i=N-2$  " " $X_{N-1}$  " "

.08 " " " "  $i=N-3$  " " $X_{N-2}$  " " acct.

.09 " " " "  $i=0$  (no history) " $X_0$  " "

So, there may be a soln for fixed horizon  $= N$ . It seems Much simpler than my

Previous Soln! — Is it correct?

.12 In .03, we need from ALP is: given  $[X_i, Y_i]$   $i=1/n$  and  $X_{i+1}$ ; what is

.13 t. p.d. for  $Y_{i+1}$ ?

≡ Say  $N=100$  & each  $X_i$  can have 10 values only, and  $Y_i$  can have 20 values only.

~~We have~~ There are  $10^{100} \cdot 20^{100} = (10 \cdot 20)^{100} = 200^{100}$   $X_i, Y_i$  strings up to  $i=100$

In .06, we ~~divide~~ divide t. no. of possys by  $10 \times 20 = 200$

.07 " " " " " " " "  $200$  again

.08 " " " " " " " "  $200$  ... acct.

En) int .09 there is only 1 possy left!



Very likely  
Peters

Great Book Rev. ~~What?~~ ~~What?~~ ~~What?~~

Universal P.d. on all Stock & Operators

↳ But see (20) for DOUBTS!

SEE 253.37 for how to do it right!

- see 253.21 for counter example to 09-17

00: 244.40: We want a universal d.f. of (i) all finitely decidable d.f.'s on functions.

So: A definition string in AZ(1), defines a p.d. on functs, [or it can be regarded as a d.f. of a single funct.] Any Universal discrete d.f. on these defining strings would be a universal d.f. on all functions.

267.00 WHY for counter-example is wrong!

How could AZ(1) give a (Discrete) Universal D.f. on strings? How does LISP give a Univ. D.F. on strings? Well, LISP is a univ. Functional lang. - So if we make the input constant (say A or 0 or 1) we get a univ. output d.f.

So, the input to LISP to get Univ. d.f. on functions: Initial string defines function as output.

This defined function is then a function of interest - But it is a def. function!

09 [ Q (n to 244.34-39) can a universal pd on def. functs be equivalent to a universal pd on all {D.P. on functs} ? ( A D.F. on functs ≡ Stack Operator ).

11 answer to 09 would seem to be Yes! Because any pd on a pd on functs is a pd on functs. The this seems like a good idea, I may have to be careful when I have continuous parameters, since derv lengths → ∞. I'd like to know if saying a universal P.D. on/func ≡ Universal pd on stack functs ?? Try to prove it using the "constant factor" argt.

17 Well, any p.d. on stack levels is a p.d. on def. functs, if one is decidable w. a certain no. of bits - So is the other - by same no. of bits. So Probly. 11 is TRUE!

18 Now, note 242.00 If AZ(1) is a universal d.f. on def. functs, is it also a univ. d.f. on stack. functs? → Even if it's univ., is it any good? → 242.03 suggests "YES"!!

20 Some Doubts on 09-17: say F<sub>i</sub>(X) is a complete set of functs of k(x) strings. G<sup>j</sup>(z) is a p.d. on z - so this amounts to a p.d. on F<sub>i</sub>(X). T. functions G<sup>j</sup> have dervs, so one can have a universal p.d. on them: P<sub>j</sub>. ∑ P<sub>j</sub> G<sup>j</sup>(z) F<sub>i</sub>(X) is then a p.d. on F<sub>i</sub>(X) : F<sub>i</sub> has wt ≡ P<sub>j</sub> G<sup>j</sup>(z) Still, do I ever need to know G<sup>j</sup>(z) if given wt. to it because its imp? - Yes! ∑ G<sup>j</sup>(z) F<sub>i</sub>(X) is a stack operator on X : It could be a soln to a QA problem.

T. Soln. to a TSA of MTM QA probz will be usually much easier than for a corresp. NMTM problem: This is because after solving a MTM problem, one usually only saves 1 (or very few) solns, while in NMTM it's essential to save many. Note how, this in a complex MTM, we must have NMTM facilities so Heur may be used & implemented.

Dats 33) Note: For NMTM, one only needs 1 function to solve all problems. For NMTM, I imagine one would need an Enormous no. of (intd) functs to give reasonable p.d.'s for most NMTM probz. (M) means deterministic. (NM) means stochastic. M probz v.s. NM probz; M functs v.s. NM functs

So, the Moral may be that while ALP is fine for theoretical understanding, The problems of approximating it w. finite CB, normally involve things that don't look much like a approximations to the sum over all codes in ALP!

(246.00 is on course) → 247.00



# Correlation

SCALING 2: A possibly new area where scaling is imp't:  
"Correlations". Th. reasoning: Say we are searching 4. pc space of 1.

product of  $k$  cons. Each one is self correlated so that every other trial is about  
1. same. So redundancy = factor of 2. If we search over  $k=2$ ;

Correlated

1	1	2	2	33	44
A	A	A	A	AA	AA
1	1	2	2	33	44
A	A	A	A	AA	AA
1	1	2	2	33	44
B	B	B	B	BB	BB
1	1	2	2	33	44
B	B	B	B	BB	BB

Uncorrelated

1	2	3
A	A	A
1	2	3
B	B	B
1	2	3
C	C	C

4 combinations  $\rightarrow$  4 trials.

4 combinations  $\rightarrow$  16 trials.

So if  $r$  is redundancy of 1 cons, then

$r^k$  is redundancy of  $k$  cons.

"Counter Args" As we ~~search~~ go along in the TSC,  $k$  will be kept  
constant. However, we would have lots of trouble as  $k \uparrow$ . This would limit our  
searches to small  $k$ , to a pathological degree. The searches are already limited to  
small  $k$  for reasonable  $C/B$  — but we would get an extra bias ~~to~~ toward small  $k$   
via this "correla" effect.

The way the "correla" works in (.04L) 4. pc's of each trial over  $\frac{1}{2} \cdot \frac{1}{2}$  as large  
as the corresponding trials in (.04R). Ideally in (.04L) we'd like to "pool" all  
4 of the  $A$  trials, so their total pc would  $\uparrow$  by 4, but we have no way to do this —

we have no way to recognize that they are identical & that their pc's  
should be added & that they ~~are~~ should be processed as a single trial.

So essentially, Cost of Solns would be mult by  $r^k$ .  $\therefore$  strong bias toward small  $k$ .

In addition, to the existing bias toward small  $k$ .

ID



.00: (Spec 245.33): That I'd need a very large no. of (wt'd) func solns for NM problems. Unclear as to why this should be so! Maybe IO would be evnt.

.02 How would AZI do NM probs? Say the targ. data is  $\{A, B, V, W\}$   $X_1 \rightarrow Y_1; X_2 \rightarrow Y_2; X_1 \rightarrow Y_3$  ( $Y_2 \neq Y_1$ )

$F_1(X_1) \rightarrow Y_1; F_2(X_2) \rightarrow Y_2; F_3(X_1) \rightarrow Y_3$ : We'd like  $F_2$  to  $\neq F_1$  or  $F_3$ .

Each  $F_i$  has its wt,  $P_i$ . If  $F_2$  is distinct from  $F_1$  &  $F_3$ , then we use as "soln"  $\frac{P_2}{\sum P_i}$  w/td sum of  $t. F_i$ .

The 3 probs are of pc's  $\frac{P_i}{\sum P_i}$  so total pc =  $\frac{\sum P_i^2}{(\sum P_i)^2}$ . I guess the total pc of  $t. model$

is  $\frac{\sum P_i^2}{(\sum P_i)^2}$ . If some of the  $F_i$  are  $t. strong$ ,  $P_i$ 's  $\frac{\sum P_i^2}{\sum P_i}$  can be much  $\uparrow$ .

**SN** If 245.09-17 is really true (i.e. useful), then I will want to review many previous ideas, concs, problems, in view of this **CHANGE** in All-over-view.

Do a CB=∞ analysis of .02 ff to get a feel for how  $t. soln$  should look.

For say  $t. eq$ , if no  $X_i$  has  $> 1 Y_i$ , then  $\exists$  a single  $F_i()$  that will do it. (The it may be of very low pc  $\neq$  — say the  $Y_i$  are "random").

So,  $t. induction$  will consist of  $\geq 1$  of  $t. func$  that map  $1$  or more of  $t. X_i$  into its  $Y_i$ .  
However we will then derive codes for  $t. corpus$  from these functions.

7/7/01 Superficially: Say we have a seq. of input to ops:  $\{Q_i\} i=1/n$ .

We have a "true" generator of  $t. A_i$ 's: which is a (wt'd)  $\in$  of various  $M$  operators.

$\sum z_i F_i()$ .  $\sum z_i = 1$ .  $t. pc$  of  $t. seq$  of  $A_i$ 's is

$$\prod_i \left( \sum_j z_j F_j(Q_i) \right)$$

Consider an op.  $F_i$ : Count how many times it has been "right" in the sense of getting correct  $A_i$ . Compare  $t. universal$  D.F. to  $t. "true"$  D.F. from ops. in this matter.

My impression is that with  $t. "learned"$  operator d.f.  $\rightarrow$  tends to learn the  $wt$  to  $z_i F_i(Q_i)$  but  $\neq$   $z_i$  mult. by  $t. frag$  of  $t. t. diff$  kinds of  $QE$ 's (whatever that means)!

Mainly, that  $t. frag$  is dependent not only on  $z_i$ 's but on  $t. seq$  of  $Q_i$ 's.

Right now, I'm thinking of  $t. TSQ$  as being a  $m$  a kind of "stationary", "steady state" condition.

— in which  $t. TSQ$  has to learn  $t. z_i$ 's. Actually, this is not  $t. model$  at all! The  $Q_i$ 's are arranged so as to introduce a ops of increasing complexity in a kind of hierarchical way: so recent ops are constructed of successful older ops.

.33 Hvr. The ~~Convergence~~ form of S7873 is for not really stationary TSQ's.  $\leftarrow$  I

was aiming toward that kind of form. for operators. What I really want is a correspondence betw.  $t. S7873$  more  $t. universal$   $t. operators$ .

One trouble is that S7873 is oriented toward induction and includes  $t. Q_i$ 's; it finds  $P(Q, A)$  not  $P(A|Q)$ . But  $P(A|Q) = \frac{P(A, Q)}{P(Q)}$ .

Trouble is,  $\{Q_i\}$  is not really a "statistical object" — At least, or at least, it is designed to train  $TM$ .

: Oneoff Motivations of 244.07ff, was that I had no really good ideas

on how to make stochastic operators for string  $\rightarrow$  string  $\rightarrow$  QA's.

One of the <sup>promising</sup> ideas was 244.12. Input string  $Q \rightarrow$ ;  $f(Q) = Q'$ ;  $Q'$  (plus a set of continuous params.)

.03 defines a stock Grammar that Peter Gray & I did on A.

AZ141 gives a D.F. on Functions. This is different from .00-.05 & from 244.11-.18

It is a stock operator.

.06 [ Another approach: "Updating" consists of creating A from Q as cheaply as possible: "cheap" being defined by v. cost of ~~op~~ sub operators accumulated over their use in t. past. ] **Not**

.06 Sounds close to what one often does in solving induction problems.

.09  $\rightarrow$  .06 sounds like "M" induction, but. It could be <sup>(stochastic)</sup> NM, ~~if~~ we found  $> 1$  way of creating t. A.

But I'd like a more continuous way to do NM induction.

But .09 actually doesn't sound bad either! It's like when there are 2 or 3 poss. "causes" for an effect, & one ~~must~~ makes a bridge from Q to A in various ways - each with its own pc ( $\in$  epsip).

**SOUNDS VERY REASONABLE** for "string" induction, in which there may be several poss.

Discrete "causes" ( $\equiv$  paths of causality).

**7/8/01** How to use .06 to deal w. ~~#~~ multiple (empirical) A's for some Q? At all times (for all Q's)

20 A's are poss. (i.e. how pc  $> 1$ ), because there are ops. But give those A's from t.  $\epsilon$  given Q -

18 Here, if one op has been successful at most all of t. time, for a large corpus, a best op has by epsip, Real it will get almost all of t. wt.

Say we have this "fairly successful op"  $F_2(\epsilon)$ . If there is a cheap mod'n of  $f_1(\epsilon)$ , then it might get some ~~wt~~ not negligible wt. I do want to work out details of how this comes about!

For each QA pair, consider all  $\infty$  of ops that are cons. w. that pair. Each op has its own epsip; but so this alone could give a kind of induction - i.e. epsip's. But in addition these ops are structurally related, so that increasing epsip of one will  $\uparrow$  epsip ( $\therefore$  epsip) of many others. Again, I have to work out details of just how this works.

.28 poss. Note of Imp't: Perhaps many of t. ops have no output for most inputs! They are specialized to recognize a process certain kinds of Q's. [like "3-state" (open collector) output of logical modules]

Consider t. vector of wts of all ops. Say each wt. is  $\propto$  no. of times that that op has been correct.

.32 Alternatively, consider ops of type 28, (w. reduced SSZ) - t. fraction of times they were correct v.s. "relevant". (Small SSZ because they are not always "relevant").

So: we have these 2 ops: One has been right ~~times~~ 80 times out of 100; & other 19 times out of 20. Given t. epsips of these 2 ops, one could combine them ~~to~~ or say 1 was right 3 times out of 3.

There is another aspect of t. foregoing: We can do an analysis assuming we have  $\geq 11$  operators available. Or, taxilically, at any particular time, we have only "constructed" a few of t.  $\infty$  of ops, & we only know t. pc's of these few ( $\rightarrow$  to some extent) their (pc)  $\rightarrow$  249.00 in their relations.

1.00: 248.40 : What I'm really looking for is a good ALP-type understanding of how this operator system (should) work. Just how is it corpus derivd, & what's its pc?

Well, we have this Kuga stochastic operator: It gives a pc to each Q-A. that occurs.

- So how do we use this for produ.?  $\equiv$  How do we "update" the system?

1.04 An ALP-type treatment: Consider the sum of all stochastic operators. Its pc for b. corpus will be at least the pc of the stock of ~~the most likely~~ that is "correct"  $\Sigma$  (i.e. a stock of this could have evolved to corpus, given the  $\{Q_i\}$ )  $\leftarrow$  multiplied by the wt. of that stock. of the sum of

So this  $\rightarrow$  will give us STBTZ: It gives a STBTZ (i.e. "constraint factor" part) - from which STBTZ can be derived.

1.06 Hvr, the set of all stock ops is a very large Hilbert space whose vectors represent the rel. wts of each of the  $M$  funcs that are its "basis": ~~the sum of~~  $\Sigma p_i \leq 1$  where the  $p_i$  are the coeffs for a particular stock op.

From a practical pt. of view, what we want to do, is find 3 stock ops that have by pc the corpus

$\leftarrow$  This would give a kind of approx. of .04-.08.

Hvr, for .04-.08, we need a "universal" d.f. on all of these pts. in Hilbert space. The region of H space is "finite" (in some sense) because  $\Sigma p_i \leq 1$ . Even so, any useful density funct on that space would assign (probably) a too small wt. to each pt.

(Actually ~~since~~ since the  $p_i$  are continuous, we are concerned with volumes "of" H space - this H volume being a func of corpus size ( $\geq$  size))

Probably the no. of dimensions of "Hilbert space" will be a func of the no. of examples in the corpus.

But anyway I need some kind of approx on the set of all stock ops, so that I can compare different "candidate stock ops" in how good they fit the corpus.

The  $\Sigma p_i = 1$  condition would give a H space w. "infinite" dimensions.

AZ1 gives a pd on the  $M$  funcs - from kind of info content argt. ("cost of coding")

1.23 Could we give a "coding cost" for the corpus" argt. to give wts ( $\equiv$  pc's) to various poss. stock ops?

Well, to code  $k$  wts of the set of all  $M$  funcs! We need the precision of the wts  $\leq$  the pc's of the

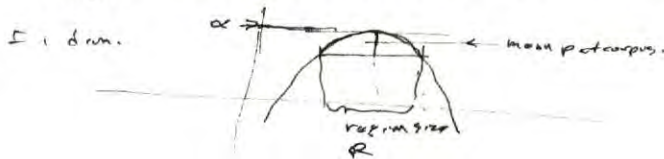
Corresponding  $M$  funcs.

for  $k=3$ : we have  $f_1(), f_2(), f_3()$  & wts  $p_i$  (i=1,2,3) ~~These are the pc's of the  $\Sigma f_i$~~  These  $p_i$ 's can be "closely coupled" so if we derive  $f_1()$  then  $f_2()$  is easier to derive, etc.

So for a cheap set of  $M$  funcs, it's best that they be highly related to the pc of coding them all.

For a stock op., we need wts. For  $k=3$ , these wts will be  $\sim \frac{1}{k}$  & they must sum to 1, so they have  $k-1$  deg. of freedom

In this 2 dim space each pt. assigns a pc to the entire corpus, so we want the region of max pc corpus. The max pc of this set of wts =  $\frac{\text{total poss. vol. of space}}{\text{h. vol. of region of the pc of corpus}}$



1.23

choose  $R \Rightarrow \frac{\text{mean pc corpus}}{R} = \text{max}$

for 1 dim case, mean  $p = \alpha - \beta \int_0^R x^2 dx$

$$\int_0^R (\alpha - \beta x^2) dx = .2 \left( \alpha x - \beta \frac{x^3}{3} \right) \Big|_0^R = \frac{2}{3} \left( \frac{\alpha R^3}{3} - \beta \frac{R^3}{3} \right) = \frac{2}{9} (\alpha R^3 - \beta R^3)$$

$$I \rho \quad \frac{2R}{1} - \frac{2}{R} = \min \quad \frac{1}{R} - \frac{12}{R} = \min \quad -\frac{1}{R} - \frac{12}{R} = \min \quad R^2 = -\frac{12}{R}$$

225 may  
1/2 of  
M. 005

no: multiply by R forget p:  $\frac{1}{R} - \frac{12}{R} = \max$

$$\frac{1}{R} - \frac{12}{R} = \max \quad \frac{1}{R} - \frac{12}{R} = \max \quad \frac{1}{R} - \frac{12}{R} = \max$$

Mean value =  $\int_{-R}^{+R} f(x) dx$

$$\int_{-R}^{+R} f(x) dx = \text{Mean value} \times R$$

multiply by cost of oil's etc. etc.

$$\frac{1}{R} \int_{-R}^{+R} f(x) dx = \text{Mean value}$$

Mean value =  $\int_{-R}^{+R} f(x) dx$

$$\int_{-R}^{+R} f(x) dx = \text{Mean value} \times R$$

Mean value =  $\int_{-R}^{+R} f(x) dx$

$$\int_{-R}^{+R} f(x) dx = \text{Mean value} \times R$$

Mean value =  $\int_{-R}^{+R} f(x) dx$

$$\int_{-R}^{+R} f(x) dx = \text{Mean value} \times R$$

Mean value =  $\int_{-R}^{+R} f(x) dx$

$$\int_{-R}^{+R} f(x) dx = \text{Mean value} \times R$$

Mean value =  $\int_{-R}^{+R} f(x) dx$

$$\int_{-R}^{+R} f(x) dx = \text{Mean value} \times R$$

Mean value =  $\int_{-R}^{+R} f(x) dx$

$$\int_{-R}^{+R} f(x) dx = \text{Mean value} \times R$$

Mean value =  $\int_{-R}^{+R} f(x) dx$

$$\int_{-R}^{+R} f(x) dx = \text{Mean value} \times R$$

Mean value =  $\int_{-R}^{+R} f(x) dx$

$$\int_{-R}^{+R} f(x) dx = \text{Mean value} \times R$$

Mean value =  $\int_{-R}^{+R} f(x) dx$

$$\int_{-R}^{+R} f(x) dx = \text{Mean value} \times R$$

Mean value =  $\int_{-R}^{+R} f(x) dx$

$$\int_{-R}^{+R} f(x) dx = \text{Mean value} \times R$$

Mean value =  $\int_{-R}^{+R} f(x) dx$

$$\int_{-R}^{+R} f(x) dx = \text{Mean value} \times R$$

Mean value =  $\int_{-R}^{+R} f(x) dx$

$$\int_{-R}^{+R} f(x) dx = \text{Mean value} \times R$$

09  
10  
12  
15  
19

What I mean is a  
particular example 253.21, it seems that  
total stock ops 2 weeks later

pc wd "combin of a set of data must be 0  
From "combin of a set of data must be 0  
pc wd "combin of a set of data must be 0

Yos! This seems U.G.!

So, this may be an adequate solution (part of) the stock op. problem!

Actually, I think the result has been using  $\int_{-1}^{+1}$

form w. factors. for jump  $P_1 \neq P_2$  results

I think I once worked out the general case; it's rather simple

well known.

It results in

2 cases of Lap's rule

So, this may be an adequate solution (part of) the stock op. problem!

So, this may be an adequate solution (part of) the stock op. problem!

So, this may be an adequate solution (part of) the stock op. problem!

So, this may be an adequate solution (part of) the stock op. problem!

So, this may be an adequate solution (part of) the stock op. problem!

So, this may be an adequate solution (part of) the stock op. problem!

So, this may be an adequate solution (part of) the stock op. problem!

So, this may be an adequate solution (part of) the stock op. problem!

So, this may be an adequate solution (part of) the stock op. problem!

So, this may be an adequate solution (part of) the stock op. problem!

So, this may be an adequate solution (part of) the stock op. problem!

So, this may be an adequate solution (part of) the stock op. problem!

Non-Mark  
Mark  
M. 005

28

In 20

Can I really express all stock ops in this form?

Can I somehow have parallel  $M$  funds needed to cover, so that p.c. add relations

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

28

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

28

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

28

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?

Can I really express all stock ops in this form?



00: 25/190 ; Th. method of 250.09.19 would seem not to be a very efficient way to code a QA corpus.

It would seem better perhaps to look at a Q, recognize that it is a certain class & could be best coded by a certain class of coding techniques: This amounts to a kind of pd. One way of doing this is via "Recogn": of sub corpus. So it could use a. some induction scheme for each sub corpus, but the probs will differ, because a corpus is different & its "updates" are different. Maybe use 250.09 method for each sub corpus?

08 Sub corpus are first acquired by having QAs being "indexed" by user. Next, user put in less indexing: TM learns to perform indexing, using previous user indexes as data.

12 Next, TM modifies indexing so it is useful to its ends, not merely corresponding w. user's indexes.

In line w. 00: 250.09.19 does not seem reasonable. 293.06 sounds better. What about the Argts. of 245.09 ff on Universal stock ops?

57873 notes  
247.33 - N.A

I Prot I had some formal soln to the QA problem w. a 57873 type proof: (say in min last 10 pp)

Couldn't find it: Try: 242.17 ff: (242.22): Think (about in terms of) soln to Bay induction problem.

Paraphrase all stock ops ~~summed~~ are approximately as a wtd sum of Maps.

Soln. to Bay induction: Soln (2 kind probab): MDL -> soln to QA Bay problem.

Say  $O_2$  [242.22] is the den of a stock op.  $O_2$  does  $P_{O_2}(A_i | Q_i)$ .

We want to find  $O_2$  s.t.  $\prod_{i=1}^n P_{O_2}(A_i | Q_i)$  is Max.  
"pc(O<sub>2</sub>)"

22 -> More exactly use with 2-10.11 wtd sum of probs of all  $O_2$  on  $P_{O_2}(A_{n+1} | Q_{n+1})$ .

A Universal d.f. on  $O_2$ 's will assign pc > 0 to every finitely describable  $O_2$  for stock ops describable by continuous params. & theorem has to be modified.

Now, what I want is a look at the general stock op,  $O_2$  - it is constructible via a 3 input UMC of 242.22. Can I somehow do it via AZI & "LISE"? KEEP 242.22 as the UMC in mind. We can think of AZI as a 2 input device: (1) Den of function (2) Input to function: (0) at pt = output of function.

33 In 242.17 (or 242.22) we have a 2 input UMC. To get to a 2 model of 242.22 we need random noise (Noise). 3 ways: (1) Add to  $O_2$  den (I guess this gives 2 d.f. over stock ops) (2) Add noise to Q (an approx Q? - I don't know how to interpret this)

35 (3) put Noise in a third input. SEE 253.00, But 253.37 is MUCH better!  
36 Paraphrase By defn, an (arby) op  $O_2$  is defined in several ways (Given Q, - d.f. on A & B for all ways) but d.f. on A can be desc. by (1) Monte Carlo output of A; or for Arby  $A_i$ , the machine prints out  $P(A_i | Q)$  or T. machine prints out  $P(A_i | Q)$  in a  $P(A_i | Q)$  order. See 258.37 for proof that a Model of 253.37 can implement an arby stock op of P45 (36) kind.

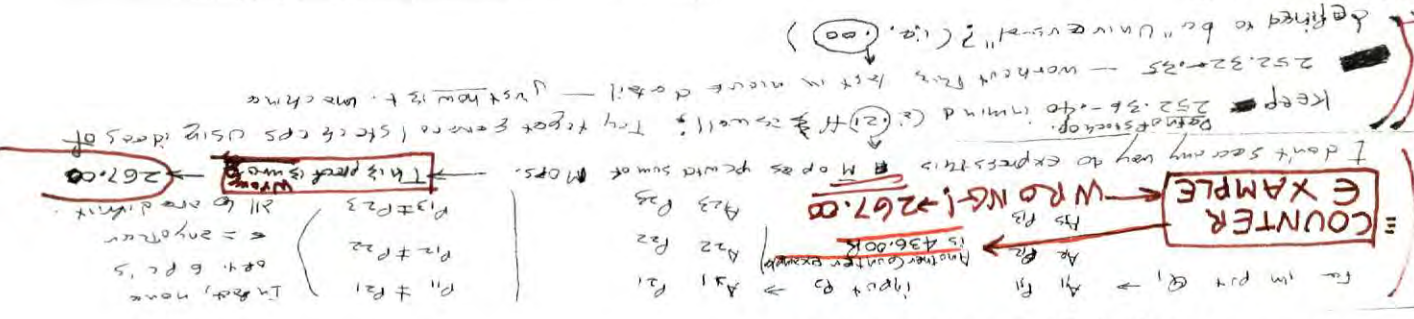
Imp:1

40

My impression: IS Mat 250.09-19 amounts to a p. over H-operators. (Hints by adding noise to operator dom.): By adding noise to Q dom, we get effectively a pd over Q's w. a fixed (M) operator (but a narrow kind of way to gate P on A).  
 There should be a way to modify (cross) AZI (Loop) so that it has a third input that is of a desired kind of Universality.  
 It is adding noise to A, we gate it in the input of H-operators. — which seems very narrow.  
 Adding noise to Q may be more general. Certain kinds of operators could be differentiated by the way Q is added to noise & act like a noise on an array of input (operator).  
 In Lisp, we can have functions of 2 (or more) variables — so one can be "Voice" & be used to construct stock ops. The no. of bits of noise needed to give a particular A, gives an estimate of the pc of that A. (We will return to this later) & particular A values. (See 258.37 for a more detailed model can express all stock ops as defined in 252.36)

28

But a narrow kind of way to gate P on A



Good!  
28

Counter examples of a string of bits is not expressible as a publishable wtd. Mops!  
 That violates it. 252.36-40 defines stock ops in most general way.  
 Some how you want to be able to give counter examples of stock ops.  
 In particular, try counter examples of 10-11. — I feel that 250.09-19 is  
 Anyway, the diagram of a string of bits is not expressible as a publishable wtd. Mops!  
 Go over all of the p.2. Carefully.  
 21.24 of actual counter example.

14

but, function of bits, in part. — (I'm not concerned its universe) on the function argument  
 A possible inference of 06-08: That the unit of AZI is not universal.

10  
10  
09

That the unit of that form — yet vice of form 252.36-40 is by stock ops.  
 This is a result seems reasonable; I may be able to do the domain stock ops  
 See 21.24 for actual counter example

See 21.24 for actual counter example

3) Since AZI seems to use a unit, then 250.09 — the domain M functions with defining  
 2) 253.00 — that all 3 must be equal.  
 1) 252.32-35: That 3 inputs of noise rendering a general stock op.  
 I can review various Argts leading to Unreasonableness;

106

It is the fact that approach of 250.09-19 would seem to be correct — it. We go to  
 If on M functions is this is (as ready to root) Equiv to any poss. stock op!  
 of info can be interpreted at other input, & it info is given; so 252.32-35 would  
 give 3 equivalent ways!  
 N.B. If a unit is universal, then I think that means that Reasonably any kind  
 of info can be interpreted at other input, & it info is given; so 252.32-35 would  
 give 3 equivalent ways!

100

2 input

ID

7/10/01

253



ID

00:253.40! So we have the Lis<sub>p</sub> (Functions operators) operating on  $\mathbb{Q}$ . For each  $Q_{n+1}$ , we  
 If our operator give  $t$  into  $A$ , fine; If not, we see how much noise (how many R bits)  
 are needed to get  $t$  into  $A$ . In this way, for each poss. Operator, we will have a non-zero  
 PC assignment for  $t$  entire  $\{Q, A\}$  corpus! As a default,  $t$  operator can not (use of)  $\mathbb{Q}$ ,  
 & use R to be  $A$ . This is usually a very expensive down, hrr.

So, at all times, all we want is a stack of that gives the lowest (cheapest) down of  $t$   
 corpus (including cost of  $t$  or itself &  $t$  cost of  $R$  inputs).

One big problem is "Updating": In a general way, we know the Old Op, Op<sub>0</sub>  
 that works w.  $t$  corpus up to now. We have new "data".  $Q_{n+1}, A_{n+1}$ . We will search  
 "near" Op<sub>0</sub>: "Near" in terms of the code of Op<sub>0</sub>. We will search  
 i.e. its sub operators & their orders. Perhaps try to find shorter code of Op<sub>0</sub> as part of "extended  
 we can (calculate)  $t$  CJS of any specific  $S_{n+1}$ . We will do Lsearch, so  
 IS .07 ff  $t$  best way?  $T$  search ~~search~~ does not look at  $Q_{n+1}, A_{n+1}$ .  
 This seems not so smart! Perhaps this is where Heuristics enter?

Could I get inside on all-012 by working MTM ( $\mathbb{Q}, A$ ) problems? (i.e. SAAB TSQ)  
 Actually, in  $t$  SAAB TSQ - did it look at  $Q_{n+1}, A_{n+1}$  in Lsearch (.11-.12)??

Perhaps a very ~~simple~~ **General (by level) Heur** for TM to learn!  
 Given  $O_0$  for  $[Q_n, A_n]$   $t=1/n$  &  $Q_{n+1}, A_{n+1}$  is  $t$ -structure / code of  $O_0$ :  
 What is a good metric for  $|O_1 - O_2|$ ? i.e. what is a good measure which  
 to search for Mod. Insect of  $O_0$ ? - (38) TM will have to be ~~revised~~ revised so that it can usefully work on  
 this problem: So on Bill theorem, I will have to devise (by hand)  
 One way, is to find a short Op for  $Q_{n+1} \rightarrow A_{n+1}$ , then find a recogn.  $\alpha$  (pm) for  $Q_{n+1}$   
 to distinguish it from all previous ( $i \leq n$ )  $O_i$ . (34) U.S. updating scheme (Conf. text), so I mean "understand"  
 it.

Anyway: I had that way to do  $\mathbb{Q}, A$  induction for MTM, which was o.k. but  
 limited in various ways; all related to its being a MTM rather than  $\bar{M}TM$ .  
 Now I have this simple extension of  $M\bar{T}M \rightarrow \bar{M}TM$  - so see, if it is adequate -  
 see if it addresses all of  $t$  points of 140<sup>00</sup> - 141.40.

A common Method of Search for  $O_1$  ("Updating"): find a code w. a  $A'$  matrix "close" to  $A$ ,  
 then use  $R$  input to dev.  $A' \rightarrow A$ . Related, is to use  $t$  ~~un~~ unupdated Op on  $Q_{n+1}$  to  
 get an  $A_{n+1}$ , then try to find an  $R$  that gives  $t$  into  $A_{n+1}$ .

Actually, .20 may not be so good! (what it amounts to is generating  $A$ , indep of  $\mathbb{Q}$  - which  
 means we get  $t$  standard  $\alpha$  props for unconditional PC for  $A$ . (indep of  $\mathbb{Q}$ ). This is not what is proposed.  
 This is not what is proposed.  $T$  Here I was thinking of occurs in ANL  
 (SAAB TSQ) Learning "3 X 5" after 8 + 7 has been ~~learned~~. Learned.

So, we start w. one or more methods to do .16-.19 (in factored form) then we use it to work  
 problems leading to ability to improve updating: Then we have TM work on improving  $t$ .  
updating routine (.16-.19). This is an OZ problem, but my usual way to solve

ID

dynamics associates

00: 254.90 : QZ probs is in 2 parts! f. first part is an induction problem! f-second part

is a kind of interpretation technique - but I had that standard Monte Carlo

Way to do it, but it seemed efficient

One trouble w. f. exp. is that in f. Lstch, correlns betw. condg causes slowness.

Another Consideration on f. such: When we are trying for new Ops, we recurrently have to test each one on f. whole corpus. Hvr, there are certain modifiers of f. old Op that will not change its response to all of or parts of f. corpus - so one could save much time by trying them first. One Op of f. type recognizes "old" parts of f. corpus, & uses unmodified old Op on that part. For "new" part it uses modified Op, so it has to test it on that part. T. new Op may be a min. modification of old one,

so as to reduce its down cost.

On "Updating": As of a long time ago (in IDSA) I was considering many possl. induction systems. Each had its own update system (But I suspected that

update systems were a very ~ to that of f. MTM, AZ1)

What I have now, is a common way to express all solns to operator induction problems.

This makes it easier to compare their effectiveness on various sub-corpus.

w. this new idea (253, 37) (Good way to get stock ops) I probably have to go back (1.16) and rethink f. whole system. (14?)

Initial system will only have one induction mod: 1 update mod (More advanced system will have several IAs & a FC) that looks at Q & decides which IAs to try.]

Even in f. simplest system .04-13 (modifiers that in value testing only 25 well part of corpus) will be very imp.

At first glance, hunting for Good Stock ops. is a very much different (inspired) than hunting for good Maps. We do trials in so cost order: T. pc of trial depends on both f. pc of an op (w.  $R=1$ ), & how large R has to be to d. c. b. f. A. We can compute the total pc of f. successful trial in this way, & estimate its cost ( $\approx c_j$ ). We will do trials of Ops using short R's; & can new ops (old ones) w. diffrent, larger R's next.

I'll have to work out details of repairing, reversal of pc of sub ops, discovery of new sub ops - All part of updating UPDATING.

Maybe go thru prog. of MTM first, (so those parts of MTM are clear) then add in an extra input to f. OMC. In a certain sense, MTM adds nothing to f. MTM problem.

In MTM we look for a cheap Op (or set of ops) that does  $Q \rightarrow A$ ; for known corpus. In MTM ... plus R ... Its just that we have to add "R" input (in addition to f. Op data) in MTM. Its usually a lower pc object searched for, but the search method should not (superficially at least) be much different!

dynamics associates

00: 255.40: So why did I get it. ideas that MMT was so much harder. It may indeed be

that the cjs's in MMT are usually logarithmic in MMT.

MMT has, normally, many assoc IA's, that could involve complete searches (?).

When doing Proda, w.t. Stack op. modal of 254.37, we ~~could~~ <sup>usually</sup> get a diff A for each value of R! ~~we~~ <sup>almost always</sup> get a P.D. of fair width for same A, or in some cases, w. finite  $\leq \beta$  - we might have output that is indep of R

So, one output w.  $pc=1$  - R's could not occur if we considered all poss. codes, hvr.  $\rightarrow$  set 255.40-43

In General, the most likely output A/occurs when  $R=A$ , which is a great advantage. <sup>usually</sup>  $\rightarrow$  probab.  $\rightarrow$  see.  $\rightarrow$  for reasonable A's. Hvr, the R input has to be a prefix set, so  $\sum 2^{-|R|} < 1$  (actually all we need is  $\sum 2^{-|R|} < \infty$ )

This means that usually A is not an acceptable input for R (!). (since it's a prefix set w. only one member!)

Actually, the way it works: we put Q into the operator. Q is self-limited, so T. maxima knows when it terminates. T. Q's must be a prefix set. On the R input, we just try inputs at random until we get output. TM reads as much of R as it likes, then eventually, prints A & steps. T. inputs on "R" that would give over/under/stoppage, constitute a prefix set.

Each prodn. involves feeding in some random "R"s to get a distribution of A's: we collect them, & note which ones have repetitions. We just observe density of A's & we don't collect random noise that generated them.

Another, more systematic technique for generating R's is not random, but a "left hand rule" search of a binary tree. The actual way it works is described in "0 search 1985". I also had a more recent analysis of R's (in Rese IDSA notes). I think a deterministic method may go thru to get a set of outputs, w. varying based on repetitions - but I'm not sure it's so good if one only wants to spend a little time to get a set of reasonable A's.

Fortunately, one doesn't ask TM Q's very often & during its training - in fact one need not ever ask it Q's during its "training" phase.

So the fact that op of may be slow, is not ordinarily an impt. constraint on TM. If we want TM to operate in real time, it might be a problem. - But there are certainly more serious problems in TM's design!

26 One big problem is 255.04-13. To make that ops in a way so that one has to test them on only a small part of the corpus - because the changes in op will not change its response to most of the corpus. [whenever finds a  $(Op, R)$  appropriate to a QA pair, one saves both the Op & the Q. (perhaps): T. idea is that one is testing a modified Op on the corpus, it may be heavy to test it on some old inputs. Knowing  $\geq R$  that once worked each would make it one's first trial R, & it could save much time  $\rightarrow$  36]

31 Probably what to do: For each of the impt parts of 140.00-141.40 (a 255.04-13) (a probab. obs) Try it out on a TSO. See how I seem to solve or avoid these problems. - Incorporate these techniques into the Heurs. for the IA's.

36 31 One way to deal w. 26 is to divide corpus into parts (i.e. sub-corpus problem) such that each part has its own stack operator: whenever finds a modith of such an operator, it only has to be verified on its own sub-corpus.

Even so: verification of a stack op on a corpus of any size involves enormous costs. <sup>NO! see</sup> Given a corpus operator, one must look for each soln. for each Q & A. Well, not so bad! If there are k QA's in the sub-corpus & each soln. A has a pc of  $P_0$ ; then we just add the costs of each "A" such, so effective pc is  $\frac{P_0}{k}$  (not  $P_0^k$ ). Also, for final R, we use the R that worked for the stack op being updated - this could save a lot of time!

00: 256.40: "IN Review": AZI could do MTM a bit, but only if hours were <sup>initially</sup> externally supplied. Since finding New News is a MTM problem. (Hr, comparing it to MTM solving General QA problems - in this case we would have to have TM be highly educated before it was able to find hours in a useful manner: So for a long ~~early~~ part of TSC's, both MTM & NMTM would have to be "about the same" in the sense that all of their hours would be externally supplied.

07  
08  
13  
15  
[SN] In early ALP work, I had the idea that <sup>formal CR=00</sup> f/solns to MTM & NMTM were the same. (Hr, from counter example of 253.21-24, this would appear not to be so! Another possy, is that they are the same but the common soln. is diff from what I thought it was) For sequencial predn, it would seem that in both MTM & NMTM one just does standard Alp. & finds all codes for the corpus. In MTM the shortest code is much shorter than any of G. values (unless the TSC is very short). The aux "R" input to UMC is not needed to get the correct Prob. for NMTM.

(Must likely) Since MTM is a special case of NMTM

Next, what about BAG-Induction? Do <sup>solns of</sup> MTM & NMTM differ? This is the Language induction problem. My early ideas on this were that MTM was not solvable - first only <sup>soln was "undefined"</sup> formulation of the problem

(it only possible cases were in TSC)

2 NMTM ~~problem~~ has a maximally full GORC!  
Probably 07-08 is correct; But also note (13-15)! So I'd really have to use the full NMTM soln. all of the time! On the other hand, for sequencial predn, the idea of MTM, as that there is a "short" <sup>finite</sup> code that will generate the entire sequence. If the sequence has any entropy at all, an infinite seq. will need a code of length  $\infty$ . Hr, even a seq. w. zero entropy can need an infinite code for seq. of length  $\infty$ : But the code length for seq. of length  $N$  will be  $R(N) \approx \lim_{N \rightarrow \infty} \frac{f(N)}{N} \rightarrow \phi$   
I write define MTM as having a corpus that for  $SSZ = 005$  has a finite term.

Anyway, whether I start out w. a MTM or NMTM corpus, ~~my~~ the system will not be discovering hours for a long time. I will have put in hours, & they may have to be expressed in probabilistic form: {say the input times of 253.37}

32 So in any TM, the major problem is finding efficient updating Algos. This is tied to the set of IAs used, since each IA has (a suggests) different methods. In all IAs we will have hours for efficient updating.

One Q is how far can we go when in a TSC w. only a few updates ~~how~~ hours?

One Big Hour is to divide Corpus in to Subcorps: initially in indexing... (252.08-12)

What I want to do now, is write (micro) Review of How I am now: Just how much of M/TM I have designed: what most critical probs are & what some stats at solns. Are this ~~work~~ would be from the core of the IDSIA reports.

7/11/01

IDSIA

dynamics associates ABCDEFG

$$\frac{3}{4} \times 200 = 150$$

$$\frac{200}{\frac{4}{3} \times 100}$$

So ~~1.5 times as fast.~~

.00: 257.40: To what extent could MUTATION & Crossover & Fair sized population be a useful Method to do TSQ's? The ~~new~~ cands would be a 2 input Crisp pairs of 253.37. We look for a ~~the~~ cand of max Gene. prodns 200 best made w. + Best Cand. A big trouble (as 256.26) is testing Cands. Say a cand has only 1. Op decrn. For each  $Q_i$  we must also find a proper R to get A<sub>2</sub> — which can take "much" time! So the Q<sub>i</sub> would Mut & crossover be better than LSrch? — (The Lsrch for small CD updating "has not been defined" (or; more ~~exactly~~ exactly, I suspect there are many reasonably good ways, but I only know 1 or 2). Any way GA is but one poss. Heur. for UPDATING → .16

.08  
.09 N.B. I think my Main line of Attack is to write a TSQ that think ε should be able to "Learn" — Then as I (think), Express this 1mg, in 2500 Heurs, in a formalism that I've devised. Up to now, 253.37 my formalism hasn't really been "Complete". It seems to be now. I will want to take various IA's & their heurs for Update, & express them in terms of t: 3 input UMC (of 253.37). Also note UNIVERSALITY of 253.37 (259, 40-43)

.16 .08 **N.B.**: Since Evaln. of Cands appears to be very time Consuming: I should use a Variant of GA in which we try to make a quick approx. of t. Evaln. funct. & use it for parts of t. Infg. Another Possy is that I should have approx. evaln. of cands be an important problem for TM. Normally, I'd think of this as a problem for TM<sub>2</sub>: & not have TM work on it until it was "Suffly. Mature". But it may be possl. to have TM work on this problem Usefully at an early state of Prog.!

.20 On Sub Corp's & Updating: For many non-overlapping sub-corp's, TM will have separate Ops., that will be updated indply. Some params of these Ops will be shared (12 & Spirit of STEIN) — but many will not be. It could be that pc's of various sub operators are completely Global, but t Operators for each sub-corp's are updated individually, since to an imp't extent, Param (main corpus used for evaln) are non-overlapping. Hvr's when an Op is updated, t pc's of t sub-ops used, will be reflected Global use of these sub-ops. Actually, this "Global vs. local" info in updating pc's of subops is quite imp't. The STEIN affect is very imp't. because ssz's are usually small (2 or 3, say). On the other hand, we do want some "locality" of pc's & so pc's are scarily "Context dependent" — otherwise, we run into t Scaling problems. So, t no. of Concs per subcorps should be bounded (ε). (x Window Bounded)

.37 A (sorted) Proof that 253.37 is capable of expressing t-general stack QA op of 252.36. In 253.37, for each Q input, we have a (possibly) diffrent pd. on all poss. A's. In 252.36, for each Q, we can have any possl. relation betw. t R & put it t output A. Now, Params / relation ship of ε 2<sup>-|Rel|</sup> to t pc of an output is what ALP is about. Any such R → A (function) defines a P.P. And every possl. P.D. can be defined by such a function — and the Operator decrn input can "get Q" to define an any funct — (or have an any func assoc. w. any possl. Q.)

UNIVERSAL STOCH. OP Perfor Distribn Funct. 40

00:258.90: On Subcorpus Updating: Evaln. of Operator on a Sub-Corpus: Since we are doing many trials "Quick Abort" techniques are imp. Statistical testing can be useful: to use certain QA's within a corpus for Quick abort, because they have been found to be "difficult" QA's (i.e. Most OPS don't do well on them). It may be possib. to save an enormous amt. of time by such "statistical tests". Remember: We are looking for a OP that is better (or even "borderline good" would also be useful) as to previous OP. - So such tests are relevant. Note that we can afford to make errors. If we accept a stock of the letter turns out to be bad, we simply backtrack. "Backtracking" means <sup>size of</sup> / corpus augmentation is > one QA.

Outline of Review

200 208.00 for earlier review  
164.0 vs v.g. version

- 1. Form of QA problems:  $M \in \bar{M}$ : How Much of Math can be found ~~in~~ By MTM (QA)
- 2. Models for induction: That all probs can be usefully expressed as  $\bar{M}$  QA
- 3. The MTM QA model: AZ1: How this done, How updated: pc's, mem data, revision, pc modify, re-assigning...  
The 3 input line of 253.37: Having a universal Stock OP (40).  
proof of adequacy of 253.37 (258.37)
- 4. Formal Solns of  $\binom{M}{M}$  TM v.s. practical solns. } ONE of MOST DIFFLT PROBS. (Note .00-.03, .15, .21, .28-.38)

14 ON UPDATING: Arrange so that when UMF has inputs  $Q \in N$ , it will do  $Q$  calcns first, & be able to <sup>rapidly</sup> attend to various random choices for R. After a choice has been made for the structure of the OP, we have a ppm that examines the OP's tries to put it in to 84's form. [I'm not at all certain about the feasibility of this!] Perhaps goal of the system to start off, is to devise an OP that has this feature: It may well be that this feature would be selected for automatically! Well, an OP that got into A quickly via R input would be selected for. To some extent, this automatically selects against Ops that have low pc on the R inputs, but I'm uncertain as to whether it's really doing "Quick Abort" — [This Q needs more work] (28)

- 22:14 4. SI as a QA problem.  $T \leq Q \leq S \leq E$
- 5. Writing the TSO. (258.09) often only one QA!

24 OP Updating: Normally in updating, Modifiers off. (OP) are  $\rightarrow$  only a small part of the corpus need be tested for those Modifiers. This is an imph. kind of character of the kinds of modifiers one usually makes on the G OP ( $\equiv$  Grand OP) (GOP @). I really have to characterize those ops (at 24). Try to find examples, Generalize them...

28:21 If we use a TRMS for vet. machine: It reads the Q top completely, then asks for the last bit of R (It could have printed part A out at this pt.), we store the machine's state at this pt. & try various R values, we probably to get D.P. of A; or to just find R values that give desired A. Perhaps we can arrange so that automatically by distinct machine - it reads OP term first, then Q, then R. If so, we would want to save state after OP term is read & try, so we can test for different Q's.

39 For a TRMC (i.e. probably other machines) we can arrange of: state table so that 36-37 has to be mod. UNIVERSALITY Problem!  
40 The 3 input TRMC model of a Stock OP of 253.37 solves the UNIVERSALITY Problem!  
- since all stock ops are expressible in this form, as finite strings. If an OP should print out a A in stop, (w.o. looking at R input) then that A gets PC=1 from that OP. To get PC  $\neq 1$  one must consider parallel OPS - (usually of much lower PC). (260.00)

00: 259.43 @ ! So it looks like a main outstanding problem is small CB Update. [large (→∞) CB Update is clearly defined]

101 Some **Updating bricks**:  
 ① Division of Corpus into Subcorpora: Initially by indices supplied by USR, then modified, then improved by TM (as 252.08-12)  
 ② Try to enhance TM should try to devise a **fast approx** of  $f$ . OP evalu. funct. 258.16

(Ea know my Criticism of Koza is proposed soln. of that problem): ~~③~~ Technical backup is an example. Another possibility is a **rule approx** for  $Q$  input is a more exact **Goal** for final ("True") R input.

③ Arrange so that  $f$  is  $Q$  input is that  $Q$  and  $R$  are data in order:  $Q$  pattern,  $Q$  input,  $R$  input -  $\hat{z}$  function can store  $f$  state of  $f$  & time after each input. (259.28-38). This "storage" can reduce ~~③~~ **Testing time considerably**.

④ (Extension of ③) T, T S Q is modals of GOP (= Grand Op) are arranged so that **Update** is **minimal**.  
 This can be done by doing GOP modals <sup>that only</sup> require testing on very few (often only 1) Q A's. (259.34)  
 In Most Human prob. solving, it is **not necessary** to test inductive by pushing the limits on recent problem - on all of post corpus. Why this is true, is unclear.  
 ⑤ That Diffrent IA's have different Update Alarms (257.32): Hvr, t. Update alarm probably all have features or usually only of past problem Corpus!

analogous to Proc of  $\hat{x} \geq 1$  i.e. **pc updates**, **discovery of new subops**, **reparsing**.

So study various IA's: (Z) is **relatively simple example** of an IA w-its (abruptly easy update), **Statistical testing** (259.00): Just test a few diff. past Q A's & see that pc's of solns. <sup>tested</sup> for new Op & use to **Reconstruct** **Better** than those of old Op.

**SN Gen. discussion**: #5 (10) seems **very imp.** ⑦ could also be imp. in fact, if one makes a  $H_i$  takes  $\hat{z}$  accepts a bad GOP, one will probably have trouble in near future, and be forced to "back track" a bit.

or simply do **wider searching** (Per guess back tracking is more likely to leave low cc for soln.)  
 Hvr, t. General idea of (5)(10) is **most critical**! Humans usually ~~don't~~ find solns that don't have to be checked on much of  $f$ . corpus. This could be for "logical reasoning" reasons: <sup>trial</sup> **GOP** arguments ~~are~~ **devised** **selected** **invariant** so they **don't** have to be **verified** **iteration** **which** of  $f$ . **Corpus**.

18 ⑧ Find a **simple** by pc ~~step~~ **step** for  $Q_{n+1}$   $\rightarrow$  Anti: Then find a way to distinguish  $Q_{n+1}$  from **previous**  $Q_i$  (Z)  $\hat{z}$  Eq; (254.20-21, 34-37.) Ops of this sort don't have to be verified on previous  $Q_i$  A's (Zn).  $\rightarrow$  (261.08)

**Generation of this Eq** will use same pc's as present GOP use - But **context sensitivity** of PC's is **key**! Somehow,  $f$ . system must **compute**  $f$ . context. Some simple default context? **No** - that would be "no context at all",  $\hat{z}$  would give **very low**  $pc$  to all concs! Maybe use "vacant" context?

Actually, this is a **general problem**; whenever contexts are involved they must have a way to **evaluate** **com.** The **present case** is a **normal case**! If TM can't deal w. it, it can't deal w. context at all!

26 (More Gen. Discn): I'd like to **start** on **T S Q** as soon as poss., to get ideas of how **humans** solve "Update problems". Before doing this, I'd like to be **knowledgeable** enough to know what to look for, & be able to recognize something useful, when I see it! Also, I want to write kind of review of **the update problem** for myself & for the **IPSIA Report**  $\rightarrow$  332.24 ff. is a review of some impl. ideas on "Updating".

This involves knowing what problems are, so it will **notice** when a soln. to one **update** appears!

**SN For T S Q**: A possibly interesting sequence: Evaln of Alg. expressions by successive **simplification**. There are several <sup>(hours)</sup> **concs** involved - all of which would want TM  $\hat{z}$  to **help**.

One Q is: do we want TM to "learn"  $f$ . needed heurs? To do this would require a very large, diverse set of problems: T. <sup>(hours)</sup> **concs**  $f$ . in **ranking** of: ② The idea of a "value" of an expression/subexpression.

③ That a shorter expression is "simpler" is closer to one's goal (usually!), ④ The substitution of  $\hat{z}$  equivalent value for a sub expression, leaves value of  $M$  conc expression invariant.

⑤ That one can solve problems by "hill climbing" - successive x-chngs that improve some conc (in this case conc = **Simplicity**)

ID

An Organizing Principle to dynamics associates Understand what TM's Main

3.5 = 404 or 1.25M by

or AA in Gruebel.

SN AAren figuris Very low Complexity Art! Got roots from Raykur & ... (more page)

problem is .04

.00: 260.40 SN Perhaps Ken Haas' Thesis would be useful to read now. It seems more explicit than Lanot's demot "AM".

His (i.e. Lanot's) problem would seem to be close to that of MTM (i.e. maybe MTM) - but to Gore was a bit unusual: It was an estimate of "interestingness" of a Cond - i.e. how Gore could evolve -

[ That was true in "AM" - was it also true in Haas' "Cyano"? ]

.04 SN I had earlier, conceived of a TM that used many "IA's" in attempt to solve (M) QA probs. It would search first since all IA's have the same Gore, that each IA could be regarded as a kind of "Heur". Tho, in 1. case of Numerical v.s. Non-numerical QA or M QA probs, these each requires (IA's) could be very different - i.e. final Gores would even look different. (19) (264.25)

.08: 260.20 One way this could work: for QnH -> AnH, we find that previously successful GOP, (On) no longer (does work) well! It gives a very low pc to AnH. We try modification of On in the usual way (what ever those are! @). We finally find an optimal solution for QnH -> AnH in very high pc, but it doesn't work well for Q's of i in En - so we modify it using a combination operator from On & O - On works on Q's of i in En; O works for i in En; we have to find a means of distinguishing Q's in these 2 classes (we do this by i.e. explicitly - since i. data is unobtainable). It may be that O works

well for a larger class of Q's than just QnH: In this case, it is always cheaper to specify this class of objects, rather than to single member, QnH.

[ At worst to pc of a class is the sum of the pc's of its individual members ]

Macro-subcorpus. see [Q, A, i]

.19: .07 from .04-.07: we can view the QATM problem as: Given a [Q, A, i] to Q i=1/n; and to assoc On: Given a set of various IA's & other tricks can be regarded as "Heur" toward this goal. Find a good On in available time.

.21 SN Some insights on .19: One way to do updating! Look in the "neighborhood" of On (382.11-.15 discuss "macro", "neighborhood"). First we get as short a code for On as possible (or equivalently, a list of parts that we can use).

.23 Then, using the pc's of cones within On, (a pretty OSL), we make modification of On of increasing benefit.

.25 One trouble w. this: But I'd expect the "macro" should depend much on (QnH, AnH). .21-.23 is the classic P(Y|X) in ALP form = P(X|Y) / P(X). X is On, Y is cond for OnH. In .25, we use the analog of "unconditional probab"

Perhaps we are already taking QATM into account, since On & OnH are Q -> A mappings! ? A minor modification to .21-.23 would look at the "part" of On that is relevant to (QATM), i.e. look at modifications "macro" that part. [ I'm pretty sure that On does have "parts" of this sort, & in general, when we attempt update On, we only update parts of it. So On can have "parts" & the corpus can have "parts": There could be a (1-1) correspondence between them: But not really! ]

.33 SN T. "Resency" context for pc's of cones may be very useful! Each subcorpus (subject/domain) will have its own Resency. If time scale could be no. of problems into past, No of bits of info into past; No of compressed bits into past.

This measure must include Q's only, A's only; Q+A's. Hrr, occasionally, one must add a batch of data (a census or dictionary or encyclopedia) - (Q is ordinary, part of a Q).



00:26:40 : [SN] Advantages/disadvantages of starting out w. MTM! (rather than  $\bar{M}TM$ )

Advantages: The Problems are easier, yet the hours used, can probably be used directly or in modified form for  $\bar{M}TM$ . (So MTM is a good "Training exercise" for  $\bar{M}TM$ !).

Disadvantages: Caution: It may be easy for an "error" in a TSQ (whatever an "error" may be!) could get MTM into a cul-de-sack (local extremum) from which it could not escape.

It may be possible to avoid the "trap" by using the "Cmc. Net" associated methods as described in Sol 89,

Disadvantages: MTM may not be able to get context dependent pc for cons. — if not then it could not go very far due to "scaling problem". It could be that context dependence

is <sup>for PC Assignment</sup> could be developed by USR (Mo) is given to MTM as Fixed HEURS. This may

have to be true for  $\bar{M}TM$  as well, since the young  $\bar{M}TM$  would not have much ability in

working on probs of that difficulty (??). I could simply use the "recovery" context (201.34)

TD

00: 262.40: Dcm. of 3 input times used in QATM:

The 3 inputs are ordered: The first is  $\pm$  operator dcm. These inputs ~~are~~ <sup>are</sup> parallel. It is a dcm. thing.

The second is  $\pm$  "Q" input: it, too, is a dcm. thing.

The third is  $\pm$  R input " " " "

The inputs and outputs are all unidirectional. The first input must be read completely before the second is read.

The second must be read completely before the third is read.

Output may occur at any time. When the machine enters <sup>the</sup> "stop" state, then whatever is out output tape,

08: is "The output",  $\equiv$  "A". Write ~~the~~ <sup>the</sup> definition,  $\pm$  3 inputs must all be prefix codes (easily shown) (33)

Because of the prefix property, each input defines a p.d.w.  $\sum p_i \leq 1$ . Also, since we are taking the product of 3 p.d.'s to get our "Wk", we may have stronger convergence conditions:

For 2 p.d.'s  $p_i \geq p'_i$ ;  $\sum p_i \leq \sum p'_i = 1$ , it is true that  $\sum_{i,j} p_i p'_j$  must converge

more rapidly than  $\sum p_i$ , say? i.e. say we put  $p_i p'_j$  in order of size, so  $x_k = p_i p'_j$

$x_k \leq x_{k-1}$ . Would  $\sum x_k$  converge? Similarly w. 3 ~~parallel~~ p.d.'s  $p, p', p''$

$x_k = p_i p'_j p''_l$ ,  $x_k \geq x_{k-1}$ ; would  $\sum x_k$  converge?

Say  $p_i = p'_i \approx 2^{-i}$   $x_k \neq p_i p'_i$ .

Q: does  $\sum_{i,j} \frac{1}{2^i 2^j}$  converge? It diverges for each  $\#$  constant; so it converges badly!

But still  $\sum_{i,j} \left(\frac{1}{2^i}\right)^k$  converges for  $k > 0$ .

$\sum p_i p'_j$  diverges if  $\sum p_i$  diverges - no matter how fast  $\sum p'_j$  converges.

The conv. suggests the  $\sum p_i p'_j$  doesn't converge more rapidly than  $p_i$  or  $p'_j$ . - it converges as rapidly as

conjecture  $\rightarrow$  fastest worst of  $\pm$  2 p.d.'s. So  $p p'$  will have first moments only if both  $p$  &  $p'$  have first moments

(we assume  $n = p$  &  $p'$  are in pc order, where moments are computed)

If  $\sum p_i$  is not in pc order (largest first), then  $\sum p_i p'_j$  never converges.

I suspect that changing order of  $p_i$  will not change whether or not it has a first (or any other) moment.

The convergence depends on behavior of  $p_i$  for large  $i$  only (because a finite no. of small  $i$ 's).

However, I'm not sure of this.

33: (38) The reason we want the  $O_p, O_Q, R$  inputs in that order, is that it makes it easier to test

them. E.g.: Say we want to test a particular <sup>operator</sup> ~~operator~~  $O_p$  w. various  $Q$ 's &  $R$ 's:

We save its internal state after it has read the  $O_p$  dcm. input.

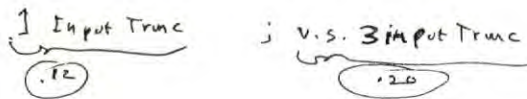
Similarly, if we want a d.f. for  $O_Q$  w. input  $Q$ , we save state of  $O_p$  after  $Q$  input, & do

tests ~~of~~ of various  $R$ . The  $R$  itself can be random, or in order of length, or ~~just~~  $n$

"left without use" for exhaustive testing (The key in order can also be exhaustive)

7/15/01  
TD

264



42km) 13 Joules.  
18V  
254V x 344  
13.00  
18^2 x 320 =  
x 0.42 = 1/25  
25.4M x 3.44 =  
18755.  
Richard Gott?

106: 263.40: Defn. of "Universality": The Op dem. is universal, in the sense that if  $M(Q, R)$  is any 2 input machine of the previous type ( $Q$  is Read-only input, (so Read and (Pre)Read sets)), then we can simulate this machine exactly w. the suitable Op input. (The Op input would have dem. of  $M$ 's state ~~table~~ transition table  $\geq$  a simulation of operation of that table)

T. Above univ. also has (at least one) work tape: The we will probably not be using Truncs; at least CSP Machines or Storage Modifi. Machines or Random Access Memory Machines.

109 Each Op dem is each  $Q$  is dem'd by a prefix set of codes! set. Ude can tell when a code ends  $\rightarrow$  start reading the next kind of input  $O \rightarrow Q \rightarrow R$ . The  $R$  inputs need not be coded by recursion anyway. When Rec machine stops, it will have read the last bit of a  $R$  input. — This makes  $R$  input set automatically a prefix set.

112 Note: We confuse easily a 1 input univ. because since  $Q_i \neq Q_j$  inputs are prefix sets,  $R_i$  univ. knows when Recy end is done on to the next one. As for  $R_i$ , (.09 - .11) makes it clear that this becomes a prefix set anyway! As w. 3 input machines, we can save the state of the system after  $Q_i$  or after  $Q_j$  codes completed. The " $Q_i$  state" occurs when it asks for the first  $R$  bit; Rec " $Q_j$  state" occurs when it asks for the first  $R$  bit.

20 A possible Advantage of a 3 input machine:  $Q$  can be parallel two series. A can be prediction of another time series (or once off loss in  $Q$ ). Instead of stopping when TM has finished A, it can continue indefinitely (as for normal sequential prodn). We can only allow only to read up to a certain pt. (today's day): "A" is tomorrow's data.

25: 261.07! On IABed "Heuristics": Very Good! I problem of Q/TM is clearly stated in terms of getting a good codes for the Corpus. Each IA ~~is a~~ can be Usefully Regarded as ~~not~~ making more than a heuristic device for finding good codes. (almost)

30 The main very imp't aux problem is the Sub-Corpus problem: Given a certain CB: How much is what part of the Corpus should be used to make a particular prodn (or  $Q_i \rightarrow A_i$ )? This may be related to the part of the Updating problem involved in not (wanting) to check the entire corpus when one makes a modifi. in the GOP. My mind is certainly not clear on this! If one uses only a certain Sub-corpus for a prodn, then the Op used for that prodn, will have to be verified on that Sub-corpus at most. — but I think this is only part of the problem of (.32 - .33) — There are other methods.

35 7/16/01, On "Updating": That we should have a good, fast, way to get  $O_{n+1}$ , given  $O_n$  and  $Q_{n+1}$ ,  $A_{n+1}$  — This is the main idea of TSO's — That such a way to solve problems is not a useful technique unless this  $O_n \rightarrow O_{n+1}$  fast updating scheme is possl. I think it is possl. because people do it. But still "It remains to be shown". So really, this "fast updating" idea is perhaps the most critical idea in TM — the idea of TSO's. And I haven't really shown how to do it in General! — This solution to the problem should make sense of the problem situation.

5pac  
264.38  
00: 264.40 : Usually when one updates w.r.t.  $Q_{n+1}, A_{n+1}$ , one does not need to check trust of corpus at all! Why is this? It may be that one uses special  $Q_n$ -modification (egms that make this pass). I will look at a few TSO examples to get some idea as to how this is usually done!

[Note that this "Quick updating" is not always poss. Periodically one does "Serious Change Revision"  
= Serious backtracking - that can involve much more L2ch.

One partial soln to the "Verify w.r.t. whole corpus" problem: (1) reduce problem viz "sub-corpus" technique (264.22-33) (2) Verify only w.r.t. specially chosen, critical examples  
This is close to the technique of Statistical Verify, but not the same. It can be reduced w.

Statistical Verify. A good candidate for a "critical example" would be a  $Q_i, A_i$  in which it took for  $O_i$  a particularly long time. (The, the reasons for "long time" should be larger C.J.S. not some "External reason"!

(SN) Another problem that may be serious. When a Op gives a p.d. for A, we are always interested in the normalized pc of the "like" A. If we use the "semiclass" instead, T. main problem is that different Ops for the same problem will have different normal constants - so be not well comparable. It may be that this is an error we can "live with"!

Note, here, that a very experienced operator that has been given Q's but always have A's will have a very small prob of "U" (i.e. undetected fault) so for Mature MATURE TMs, Normalized will not be a problem.

See 'sol to comments folder 7/15/01 on "Normalized"

(SN) There had been a problem w sub-corpus in which one s.c. would end up w. a certain pc for  $P_0$  equivalent to soln of a particular problem. A different s.c. would end up w. a different pc for  $P_0$  but a different soln. to that problem  $P_0$  is  $P_1$  and pc for the entire corpus - including the latest problem - so  $P_0/P_1$  does not give the pc ratio for these problem solns. (The perhaps this was an error - one must use some <sup>sub</sup>corpus for comparing pc's). More troublesome would be if  $P_0$  &  $P_1$  were obtained w. same corpus, but different CB's & /o different heuristics.

Superficially, this problem would be solved by the 3 input UMC model of 153.37 (or its equivalent input model of 264.12). The output of the UMC depends on random (R) inputs. The ratio of pc's of 2 outputs depends only on how one does the R trials. There are several "inbred" UMC models (1) random R trials in order of  $\log(R)$  (2) Exhaustive R trials (this may not be possible, since for many/most Ops, there for all A,  $\exists R$  that will give that A. i.e. no A has zero pc) -> Here, probably true, see 39-40, May 266.00!

36 (SN) The 3 input UMC of 253.37 is formally correct, it may not always be a good way to express stochastic ops. T. main necessity is that I want a term w. 2 inputs: Op dem, <sup>(a)</sup> Q input & this dems a pd over all possible A's. This latter "P.d." need not be the "R" input form. 39 It can be in any p.d. form that is useful. T. advantage of the "R" input type & form, is that 40 it is Universal - so we get a universal d.p.f. from R to output for all Ops & all Q's. !!??

ID

Boris Katz? (MIT)

209.00 1324 early review of QATM

290.10 has Goodiegram subgroups problem; it incorporates w. idea of IA's as "hours" 164 is v.g. early form of QATM.

00: 265.40: This lasts seems unlikely! A more one deriv. any operators, & many operators are not universal. But getting back to 265.36: I may not usually actually use r. 3 input vmc Model, — but that model should be able to be used to criticize any induction I do w. other models. [Note 267.00; I may be able to use it. "dfs coded" approx. math. doc

03 So math apparent probs. (other than TSC writing) (264.30-35) are (1) deciding on what sub corpus to assign a Q to? (2) what sub corpus to assign a QA set to for "update"? (3) How to verify/evaluate a trial Op w. minimum cc how to select ops so that P's evolv. is easy.

In line w. t. confusion of .03-.04: In my QA machine, it seems that all "real work" is done in updating QA's. That given a new problem, Q, finding a likely soln, A, involves only searching over outputs of current QOP vs R inputs.

Yet, in fact, when a new, unsolved problem, it can take an enormous amount of time & creative energy. What's going on here? Perhaps this depends on what one considers to be a "solution" to a problem. Perhaps in some domains (math, perhaps) a "soln" to a problem is a good technique for solving it. The "technique" could be quite cc-consuming!

Well, very hard problems are like "Curing Cancer". (293.00). One has to learn certain general techniques for solving probs of Reschkind. 243.00 ff puts the problem in a precisely defined form (as a Dru. Pmg problem) but devising suitable approx. techniques is something that only a mature TM could usefully work on.

As is, the QATM that I've outlined, is able to solve problems by generalizing from known solns to previous problems. I planned to educate it by a sequence of QA's only. In the present framework, it would be meaningless to give it a problem to work "to educate it". When given a Q, it simply generates A's of (what it thinks) are of my pc.

But: I can't learn to solve eqns & other math problems. It could learn how to "invent proofs" of Lemmas. It could be taught how to solve very diff. probs. — But, to what extent could we get in such "Originality" out of such a device? See 231.00 on "creativity"; we could teach TM "creativity" by Examples. (30)

(SN) In an earlier note, I had IA's choosing subcorpus as well as doing fi approx. So this is a "big job" for a Hour! — [see 261.04]

30: (27) also by biasing search probs toward lower pc's — i.e.  $p \rightarrow p^{\delta}$ ,  $\delta < 1$ . — we can do this with only limited  $\delta$  since it will diverge if  $\delta$  is too small!

Or, use a very small  $\delta$ ; use random (using  $p^{\delta}$ ) choice w.o. replacement (not so easy!) — maybe hash codes to prevent repeats? Or, keep a list of trials made & "lock up" each when a new trial is suggested. It may be easy to find ways to store to prevent repeats.

[ In some cases, we may want to allow (say 2) trials — e.g. if H.W. is "crumbly" ]

00: 266.40; WOOPS! That  $\left[ \text{proof} \begin{matrix} \text{Rif} \\ \text{of} \end{matrix} \parallel \text{M ops could NOT simulate all M ops} \right]$  is wrong!

We can simulate various pc by summation of various MOPS. Say we have 100 MOPS - each w.  $pc = .01$ , our resoln. of .01.  
For each  $Q_i$ , we can have an easy d.f. on many  $A_j$ 's (limited only by our resoln. of .01). So timed out of 250.09 ff may be ok!  
So,  $\parallel$  MOPS may be able to simulate any MOPS - but  $\perp$   $Q$  is - is it a useful, & efficient representation?

[One can ask  $\sim Q$  about 253.37 (see 265.36)]

To ~~prove~~ 250.09 ff! Consider an arbitrary no. of MOP: Each one ~~has~~ has no output for all inputs except one  $Q_i$ . That has output w.  $pc = 1$ . We have  $\sum_{i,j} p_{ij}$  of these MOPS (k it large: like  $10^{10}$ )  
For each  $Q_i$ . For each  $Q_i \rightarrow A_j$  we have  $k p_{ij}$  MOPS w. outputs  $A_j$ .  
So we have only those MOPS for  $Q_i \rightarrow A_j$ :  $k p_{ij}$  of them.  $O_{p_{ij}}$  has  $A_j$  as output for  $Q_i$  output, otherwise it has no other output for any input.

T. simulation of 00-09 is not good! It assigns  $\Lambda$  output to  $Q$ 's machine almost all the time.

We would have to renormalize - remove  $\Lambda$ 's at irrelevant - but  $\Lambda$ 's doesn't sound so good!

Another way, hm. Say we have  $k$  ( $k \gg 1$ ) MOPS.  $k p_{ij}$  of them respond to  $Q_i$  by  $A_j$  (each MOP has  $w_i = pc = k$ )

If there are  $n$   $Q_i$ 's &  $m$  distinct  $A_j$ 's: Could we do a proper assignment of  $Q_i \rightarrow A_j$ 's w. MOPS?

I think this model works ok, but I need to go over it.  $\rightarrow$  442.16 ff for a "Positive" view: Seems to work!  
Woops! 442.34 is very serious objection

We still have.  $Q$  of 265.36 is this  $\uparrow$  or 253.37 a good way, useful way to express MOPS?  
Also "woops!"  $\rightarrow$  pc of this demo is very low! not a theorem

The model of 250.09 ff has a very simple update scheme! I am dubious.  $z_{ij}$  if  $\rightarrow$  one MOP had the correct answer for its  $Q_i$ . T. analysis of (00) ff possibly change MOPS have same input

to correct for a given  $Q_i$ . (computationally) the problem seems diff

Say  $w_j$  is the wt. of  $i, j$ th MOP. Then want  $G \equiv \prod_i \left( \sum_j z_{ij} w_j \right) = \max$ .

$z_{ij} = 1$  if MOP  $j$  gets into  $A$  for  $i$ 's problem. Otherwise  $z_{ij} = 0$ . Messy!

So the L's multiplies w constant  $\leq w_i = 1$ . The partial derivatives look very messy!

Theorem I in G (of .38) to calc of  $w_j$  is a big expression involving by powers of all other  $w_i$ 's.

Since all of  $w_j$  have to be  $\geq 0$ , it may be that there are no local maxima. In this case,

it might be adequate. (If there are  $m$  MOPS) to fit an max quadratic form into each

pt. in  $n$  dim space  $\rightarrow$  do hillclimbing. In order to have a soln., each correct  $A_i$  must have at least one MOP that gives that  $A_i$  for  $Q_i$ . We can use L's multi on a Quad form, w. constraint  $\sum w_i = 1$

Presumably, we look for MOPS that have a logsum of correct  $A_i$ 's for as many  $Q_i$ 's as poss.  $\rightarrow$   $\{w_i \in w_i = 1\}$

Maybe we can do it w. a relatively small  $m$  - say  $m = 10$  or  $20$ . We search over  $n-1$  dim space - so basis various are different!

Another trick is to divide the MOPS into subsets of "non overlapping" MOPS. "Sub Corps". (27.29) R solves this

For  $n$  opten of  $\{w_i\}$  one doesn't have to reevaluate the MOPS. It just involves solving for  $\{w_i\}$  & can be relatively fast, if  $n$  say  $m \times n < 50$ .

Solving  $50 \times 50$  matrix takes  $\sim 50^3$  steps =  $125k = \frac{M}{8}$  - which is not very many steps.

$m = 100$  would take  $10^6$  steps, which is signif. for a 1 GHz processor: it may take a second.

But  $M = 100$  seems very large; If we have so separate "sub corps" (which is likely), we may be able to have  $m < 10$  (!). [we'd really like  $M = 1$  or  $2 \dots \rightarrow$  MTM problems.]

Actually, the "quadratic form" trick of .26 ff may not be so easy! Finding the elements of the quad form

This working is crazy!

.35

.38

.26

.27

.29

.30

(27.29) R solves this

7/17/01  
IP

**m: MOPs**    **n: QAs**

.00: could be very time consuming. For  $m=10$  there are ~55 coins to evaluate. (55 quadratic coins plus 10 linear + 1 constant = 65 coins. The linear coins are simply f. coins of  $r \cdot W_k$  in 2.67-3.8.

.02 T. express  $\sum_{j=1}^m a_{ij} w_j$  as  $\prod_{r=1}^n \frac{1}{\sum_j b_{rj} w_j}$ . This factor is to cancel out various terms in the first factor.   
 In  $r_1$ ,  $\sum$  sum of .02 we do not include terms containing  $W_k$ ; here .02 presus the coin of  $W_k$ .

To get the quadratic coins, it's much harder. It may take  $m^2$  operations for each coin, so  $\sim m^2 \cdot m^2 = m^4$  operations (!). I think this factor involves summing over terms (like

.07 
$$\frac{\sum_j a_{ij} w_j}{\prod_r \sum_j b_{rj} w_j} = \sum_j a_{ij} w_j \cdot \prod_r \frac{1}{\sum_j b_{rj} w_j}$$

Here we are evaluating the  $k, l$  coin of the quad form. In (07L)  $i = k, l$

we have to sum .07 overall  $\sim \frac{m^2}{2}$ , pairs  $k, l$  pairs.

.14 **NO!!** I want  $\sum_{k=1}^m \sum_{l=1}^m \frac{1}{\sum_j a_{ij} w_j \cdot \sum_j a_{ij} w_j}$ . Don't include terms that are  $\infty$ .   
 Fortunately, almost all terms  $= \infty$ .   
 If  $MOP_k$  was correct  $\alpha$  times  $\beta$  times  $\beta$  times  $\alpha \cdot \beta$   $k, l$  coins

If we use sub-coins that are small (<100 QAs) then things may not be so bad.

**So Conclusions:** that f. coins could be feasible if we use small (<100) subcoins — MOP

Also, it may not many QAs have  $\geq 1$  ~~that gives~~ rather answer than the character (250.00 ff) method may give reasonable ~~method~~ approximations (?)

How Good this model is to represent MOPs, however, is uncertain.

.21 ~~Actually, it is wrong! since it involves partial derivatives wrt  $W_k$  &  $W_l$ , it should have  $\frac{\partial}{\partial W_k}$  &  $\frac{\partial}{\partial W_l}$  z coin (factor) to products of the elements of  $W_k$  &  $W_l$  that occur in the expression. For the same reason, .02 is wrong.   
 However, in general, I'm uncertain about these (.00-.23) results! I really would have to go over it in much detail.~~

SEE 28 for how like way to do this!

MOP best.   
 Since  $\frac{\partial f_i}{\partial W_k}$  is not  $\frac{\partial f_i}{\partial W_k}$  into answer.

.28 This is another way to compute derivatives. Since we want derivatives of a product: note:   

$$\frac{d \prod f_i}{dx} = \prod f_i \cdot \left( \sum \frac{1}{f_i} \cdot \frac{d f_i}{dx} \right)$$

$$= \left( \prod f_i \right) \cdot \left( \sum \frac{f_i'}{f_i} \right)$$
 is provable by mathematical induction.

From  $\frac{d \prod f_i}{dx} = \prod f_i \cdot \sum \frac{f_i'}{f_i}$  we get 2 simple expressions for .02:   

$$\frac{\partial}{\partial W_k} \prod_{j=1}^m \left( \sum_{i=1}^n a_{ij} w_j \right) = \prod_{j=1}^m \left( \sum_{i=1}^n a_{ij} w_j \right) \cdot \left( \sum_{i=1}^n \frac{a_{ik}}{\sum_{j=1}^m a_{ij} w_j} \right)$$
 We next need  $\frac{\partial}{\partial W_k}$ : It seems to have  $\leq m^2$  terms! But do  $\alpha' \beta + \alpha \beta'$    
 $\alpha' \beta$  is same as  $\alpha \beta'$ ;  $\beta' = \sum_{i=1}^n \frac{a_{ik} a_{ij}}{\left( \sum_{j=1}^m a_{ij} w_j \right)^2}$  ! so we have  $\alpha (\beta + \beta')$  for second deriv.   
 Note  $\beta'$  is negative

So, finding  $\alpha$  involves  $m^2$  terms vs no. of QAs since  $\alpha$  involves  $m$

$n$ : no. of QAs  
 $M$ : no. of MOPS

To do these <sup>Calculus</sup> tests: first compute all of  $\sum_{j=1}^M \sum_{i=1}^n a_{ij} w_j$ . There are  $n$  of them, each has at most  $M$  additions!  
So  $\leq$   $n \cdot M$  additions.  $\alpha$  is obtained by  $n$  multiplies of  $\beta$ . Next we compute  $\sum_{j=1}^M \frac{1}{\beta_j}$  reciprocals!

$\sum_{j=1}^M \frac{1}{\beta_j}$  so  $n$  divisions. To get second derivatives, separate these so  $n$  multiplies. For each second partial (there are  $\approx \frac{M^2}{2}$  of them), we have to add  $\leq n$  ~~reciprocals~~ (reciprocals)<sup>2</sup>: So  $n$  additions for each partial (divides  $\frac{M^2 n}{2}$  additions for all of them). This  $\frac{M^2 n}{2}$  is  $\gg$  any of the other nos. of operations, so this is a better c.c. of finding  $f$ . partial derivatives. Solving  $f$ 's resultant eqns after hillclimbing involves  $\approx M^3$  operations. The  $\frac{M^2 n}{2}$  is an upper bound: many additions are omitted. What Rec  $\frac{M^2 n}{2}$  or  $M^3$  is larger, is ~~uncertain~~ But  $\frac{M^2 n}{2}$  is an upper bound (assuming  $n \gg M$ ).

7/18/01 T. Dagg. tells how to get/peak for wts of MOPS ~~by~~ Successive Approx. Some initial Approx: Have wts/ or no. of QAs they get correct. 2) ~~Give~~ Give <sup>no. of</sup> w. ~~largest~~ no. of correct answers wts. of  $f$ ; rest wts = 0.

It may well be that Approx 1's Food ~~is not~~ for most applications, & one doesn't need "hill climbing"

Q: T. analysis of 268.36 ff doesn't use  $L^2$  multipliers • Seems wrong!

So ~~every~~ " will give  $w_j \gg 1$  & they will continue to  $\uparrow$  as we do successive Approxs!

(267.27-30) & discusses this: We simply use  $\alpha$ , first & second <sup>partial</sup> derivatives of  $f$ . Give to approximate  $\alpha$  Give (locally). Then we use  $L^2$ 's Mults ( $w^2$ ) to get an approx peak for  $f$ 's next approx. This involves solving  $M+1$  eqns in  $M+1$  unks rather than  $M$  eqs in  $M$  unks. So " $(M+1)^3 + \frac{M^2 n}{2}$ " vs  $M^3 + \frac{M^2 n}{2}$  <sup>main</sup> ~~problem!~~

$L \approx 279.09$   
~~273.01~~  
out Hessian  
Hm, Note 281.00 why 268.28 ff may be unnecessary!  
The T. methods could be very imp. for 274.10 = 14



New Pen: .27

→ 310

ABCDE

.00: What were principal problems since I started working on QATM? ~ 259.08 / 2.08 is review

One of them involved some confusion in / multiple subcorpi = Multiple IA analysis:

Suppose we have several partially overlapping (IA/subcorpi): How do we mix them into products?

One way: Int. overlap regions we use w/d means of relevant I.A.'s: we chose w/ds

So that total PC of over lapping regions is Max. That's if we've already decided on IA/subcorpi. Use methods like 268.02 - 269.40 to get optimum wts.

If we have various sets of IA (subc) to choose from, we consider t. PC of ~~entire~~ entire Corpus is mult by PC's of IA's used (presumably IA's include info on how to cut out sub-corpi - if not we have to multiply by t. PC of this info.)

568 P Hold.

For new products, we use t. same IA/sub selection criterion as before (these criteria are based on t. Q's only) & we use t. previously optimized wts.

7/19/01 Another way to look at t. things: (There is some thing w to this around 165):

Say one has a Q & bunch of IA's, [IA<sub>i</sub>]. F(C) (cost of t. latest Q & asks:

What is t. probab. that IA<sub>i</sub> (of t. extra set,  $\rightarrow$ ) will assign to max best PC to t. correct A?

This problem is solved like t. standard OZ update problem is solved: by devising

2 F(C)'s: F<sub>1</sub>(C) & F<sub>2</sub>(C): F<sub>1</sub>(C) involves diff. in d. chng; F<sub>2</sub>(C) is a kind

of Integration of F<sub>1</sub>(C) & is a purely <sup>Mathematical</sup> operation. (Hvr, write out more detail of both F<sub>1</sub>(C) & F<sub>2</sub>(C) ... I think I described this in my recent version of t. "IDP report". (F<sub>2</sub>(C) was called "GPD")

T. division of this induction into an "update" vs. "prediction" modes is unclear, however

It is conceivable that TM would only do "updating" when a new problem was given. The way, TM would have access to what parts of Corpus needed updating. - i. less we stick to

.27 7/19/01 Alvin Tech Liter Q. 13 7/19/01 : 5:22 PM start of use. → cost 1.95 on Webster's Porter Sq.

128 7/19/01 164.00 - 169.40 Guess 2 Dem of a QATM; Also 169.22 has criticisms of new prs. in t.

System: [One of Russ: "Non-Universality of t. lang used to desc FEMs" is no longer valid.]

259.40  
250.08  
252.21 → 267.00

Main Criticisms: 1) Optimizing F(C) [which chooses IA's & subcorpi] is an OZ prob. for which I use L such: which seems potentially very inefficient because of complex, broken, codes.

2) (69.31): "That ~~most~~ all of t. info obtained in t. TSQ isn't properly applied to t. soln. of t. F(C)

(often modeling induction) problem" (I'm not sure I understand this remark!) It would seem that t. model used for F(C) was universal, there would be no problem - perhaps this was per

point: At that time, my model for F(C) was not universal: ∴ couldn't approach optimality.

.37 So my present universal models would (presumably) take care of both d. phys of (69.30-32) No!

Also longer has certain optimality. This can be done by t. use of L such -

FD

.00: 270.40: So, my impression of the system (as outlined in 164-169.40!

Main known ~~basic~~ radical weaknesses: ① Use of LSrch for  $f(\cdot)$  often: Correlation between  $f(\cdot)$  and LSrch.

Other diffys: 1) How do "mix" different products of different IA's!

2) " " update various IA's (in view of their different Corp's)

3) In updating: How to reduce <sup>amt.</sup> of CC used, by not verifying trials on very many QA's.

4) I really need to go over 164-169.40 to be sure that I really have a well-defined system.

.06 The main thing, is the overall system Goal ( $\equiv$  e. Goal for improving  $f(\cdot)$ ) is improving [IA's]

So the system Goal is related to the Goal of .06. The structure of the system + [IA's] and  $f(\cdot)$  are heuristics aimed at maximizing the Goal. They are H.C. heuristics (which possibly may be regarded as "Search Heuristics")

So, I need to <sup>define</sup> ~~define~~  $\Sigma [Q_i: A_i]$ , the  $\Sigma [IA_i]$  [how do we define IA's?]

Exactly what is  $f(\cdot)$ 's goal? Exactly what is the system Goal?

How are the IA's updated?

.19

[SN] Essentially, I am using an externally supplied set of IA's. The  $f(\cdot)$  is <sup>Sub</sup> Corpus selection amounts to a v.s. solution of a "Wrapper" problem. That's the main contribution of the TM model - apart from the consistent use of ALP to resolve problems.

[ "T. use of ALP to resolve problems in Wrapper Optimization" ]  
or [ Wrapper Optimization using ALP ] (description long & complex?)  
Pow. titles for paper/report.

Definition of a IA is its "Updating". The IA's we will be considering are all of the following form: Given a seq. of <sup>ordered</sup>  $Q_i: A_i$  pairs, the IA will construct a stochastic operator

$O_n(\cdot)$  which, for input  $Q_n: A_n$  will give a pd on  $A_n$ . This stochastic operator takes several forms:

[ Markov chain: R input; output = list of  $A_i: P_i$  pairs in the  $P_i$  order; Given  $A_i$ , it will give  $P_i$  ]  
[ Given  $n$ , it will give  $A_i, P_i$  w. h.  $\Sigma P_i = 1$  ]  
Also  $A_i, P_i$  or other common DF's.

If we have several IA's operating on same subcorpus, their prediction are combined by giving each IA a weight  $\alpha$  (e.g.  $\alpha = \frac{1}{\text{number of IA's}}$ )  $\times \prod P_i$  ;  $P_i$  being pd of IA  $i$  on correct  $A_i$  of  $Q_i: A_i$ , (with  $Q_i$  as input to its  $O_n(\cdot)$ ).

"Updating" is of 2 kinds: ① ~~Global~~ Global ② incremental. In Global updating,

the IA considers to entire corpus  $\Sigma [Q_i: A_i]$  and finds ~~parameters~~ <sup>parameters</sup>  $C_n$  and its associated parameters (discrete & continuous params). This is ~~not~~ for many IA's this is only possible if  $n$  is small.

.31

ID

00: 271.40: In Incremental updating, we are given <sup>sub corpus</sup>  $\{Q_i, A_i\}$  and  $\{O_{nk}\}$  Part of system has branched for this sub-~~corpus~~ <sup>corpus</sup>. We are given an additional subcorpus  $\{Q_i, A_i\} i = n+1 \dots n+k$ . (k is often 1). From this information, the system modifies the parameters  $O_{nk}$  on the basis of info. in  $\{Q_i, A_i\} i = n+1 \dots n+k$  subcorpus, to obtain the  $S_{nk}$  or  $O_{nk}$  appropriate for the  $S$  corpus  $\{Q_i, A_i\} i = 1 \dots n+k$ . Usually this takes much less cc than 1. Global update of  $S_1, \dots, n+k$ .

Each IA has its own particular methods for doing Global and Incremental updating, most of them have certain common features, so studying ~~the~~ the updating techniques of one IA, in view of improving them, ~~and~~ will often suggest methods to improve updating algos for other IAs.

I will describe Global & Incremental update methods for ANNA  $Z_1 \dots Z_1$ :  
 Also for simple Bernoulli updating: Possibly Kalman filters (See Num. Rec.)  
 Maybe ~~ANNs~~ ANN's GA's.

- 18 [SN] I think there may be a big problem in using  $> 1$  IA! : If all IAs are in mutually exclusive subsets, so all IAs in same set have same (S.corpus) every time, - there is ~~no~~ no problem
- 19 Each subset corresponds to a Macro-IA, At S.corpus of different IA's don't overlap.

Set S.C

If IAs do not satisfy 18-20; then problem is diff't, because IA's ~~predicting~~ predicting A's for the same Q, ~~have~~ have different SC's (= sub corpus), it's not clear how to weight them. Is the technique of 267.00-269.90 ~~relevant~~ relevant? Certainly the problems don't correspond directly: 267.00 ff dealt w. deterministic OPS (Mops), it's not very wts for them. However, we may want wts for stoch ops. (Mops).

28 One way: Each SOP has an assigned wts. That is used for all products: When other SOP's product for same A's, the products are combined linearly w/ these wts & normalized. The pc for the corpus drawn is the product of ~~probabilities~~ combined (or combined some kind) probs of all ~~correct~~ correct A's products, Mult by pc's of all IA's used, mult by pc of FC) -

The decision procedure for IA Application seems a bit expensive way to do it  
2 CORPUS [N.B. The analysis of Prof. ti wts of 249.23-250.19 will be relevant if this ~~is~~ coding method of 28 ff is used.]

36 How much of coding in 28 ff is in it, so pc's add rather than multiply!  
Yes: At each pts at which we have  $> 1$  ~~prob~~ prob, IA predicting, we have alternative means of coding,

37 Also, the decs of the IAs (or their ~~updating~~ updating params) can be A priori  $\therefore pc = 1$  (no cost) Perhaps also true of FC: (only <sup>its</sup> params have pc's 1). Perhaps the methods of 249.23-250.19 could be adapted to other (perhaps similar) continuous params. - the which we have to integrate over all of parameter space. 249.23-250.19 is an argument/method (273.00)

00:272.40 That can probably be used in many situations.

HA!

In that hill climbing method of 267.28 ff (Attly's Quadratic Form), the Determinant of Q quad term at peak gives us  $\approx x$  true pc at this point! I don't immediately

see how I could use this info, but it might be useful if there were several discrete p's (local max's) - This would enable me to give wts to each such "local max" - or simply cut out v. very low ones!

This extremization of 272.37 ff. is used for i. wts. only of IA's only. We can use .01 ff to get pc's of IA's i of F() when they have continuous parameters.

Hum! The idea of 272.28 ff is a way to combine IA's: I may have written this up before! Somewhere betw. 165 & 265! I certainly did work on that problem a lot - Maybe I subconsciously remembered this approach (Soln.?) Try to find best treatment yet! (165-265)

Anyway, is this Multiple IA prodn, v. fast Big diffy? Look at this I may have had some good objections!

I hadn't "solved" it by 172.24 - still working on it by 177.11 (3 days later) 177.20 suggest use of "Decision trees" in early F() to decide on which IA to use - for a particular Q.

Since wts. calcn. in 272.38-.40 will be time consuming; maybe drag to program: We may want to estimate v. wts. by seeing, on t. average, how good each IA was, per QA prodn: Also so is it t. pc of t. IA was  $p_0$  i. t. did k prodns of  $\bar{p}$  economic mean, then maybe mean pc of prodn

No.  $\rightarrow (p_0 \bar{p})^k = p_0^k \bar{p}^k$ . Trouble is: if 2 IA's did same copies,  $\bar{p}$  of log B's, they'd get wts of  $p_0^k \bar{p}^k$ ; rather than  $p_0 \bar{p}^k$  as they should. - so I don't know a simple approx

Take a clue from Bernoulli's  $\bar{p}$  Gaussian  $\bar{p}$  estimates! wts calcn. of cases /  $\sigma^2$  is the accuracy; maybe analog of  $\bar{p}$  so how worse wts be of  $k \bar{p}$ ?

or make correspondence to Gaussian error be exactly same as Gaussian error so  $\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{x^2}{2\sigma^2}} = \bar{p}$   
wts  $\propto \frac{k}{\sigma^2}$ ; so  $\propto k \bar{p}^2$

Actually in the Gaussian case, the probab. density rather than rather  $\propto \frac{1}{\sigma^2}$  - which is a more general result. So  $\bar{p}^2$  gives "line".  $\bar{p}^2$  rather than.

Two imp. Q's do we want  $(\bar{p})^2$  or  $(\bar{p}^2)$ ? They are normally quite diff't! - i. diff. by  $\sigma^2$ .

2 We should take cross correln. into account. If we have  $\geq 2$  pc's we will need to correln. mean & relevant pc's. A statistical analysis in terms of cross correln. could give a method estimator of pc's! - I go & think it more computational intensive (harder) in a peak of 272.29-.40 is probably more accurate; it gets dithering. (we may be able to construct Gaussian by cross correln. estimates)

Note: Since each IA makes several prodns for each Q... well, we are only interested in pc's of correct IA's. - which simplifies things. - The for small size cases, using info on other prodns could be very useful (T. STEIN effect).

Vertical notes on the right side of the page, including a small graph and various mathematical expressions like  $\sigma^2 = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-\frac{x^2}{2\sigma^2}} dx = \bar{p}$  and  $\bar{p} = \frac{1}{\sigma^2}$ .

.00: 27340 : Probly best way (fastest, most accurate) way would be to use cross corrln. to get first approx., then use 272.29-.40 for final approx. (~~1 or 2 iterations~~ <sup>only</sup> ~~might be~~ needed.)

Actually, ~~the~~ <sup>no</sup> conv method is 272.29 ff could involve about same amt. of time: [that Pco cross corrln & Pco Matrix elements of 272.29 ff could take a similar cc, or actually be identical!]

.04 Key my: 272.29 ff could be used to not only get good single predn., but to optimize  
.05 ti internal continuous & discrete params of t. [AE]! T. discrete params are optimized by trying various combinations & comparing Hessians. (or just using Hessian-wtd mean of linear predns)  
This ~~is~~ could be extremely time-consuming hrs since t. various combinations of discrete params are ramp by in number, & one has to do a complete mean ~~opt~~ continuous opten. for each!  
It might be possible to do a sample of t. [Qz Ae] data set for this opten of discrete params!

.10 .04-.05 is nice because it gives a very exact (Pco not really practically computable) Gorc for

t. set of IA's & t. predn. params of t. system as a whole.

**GOOD!** IN FACT one could also <sup>simultly</sup> optimize t. Params of FC() & both continuous & discrete nicely & Non-EL Gorc for t. entire System!!

.15 As a Special Case .10 ff gives an "exact" soln. to t. ANN predn. problem (w.o. Cross Validaten).  
T. discrete params" in a ANN involve how many nodes & how to connect them, T. ~~discrete~~ <sup>continuous</sup> params & wts. To get a predn., we do a Hessian wtd mean of all ~~discrete~~ combinations of discrete

.18 params. A more accurate way would refuse Hessians, but integrate over all possible values of continuous params. (T. Hessian is an approx. to these integrals). I think this is t. soln. to B.

.20 ANN problem I got in SAARB. — It was part of a lecture to give for Prob "Course".

.21 Hrs, t. ANN soln. of .18-.20 (w. or w.o. t. Hessian approxn) could be used for any predn. model (re. any IA!) & .10-.14 (including FC) is just PLX!

.23 This defines t. "Soln Space" for QATM in a relatively Complete Way! It is, essentially,

.24 "t. <sup>wtd</sup> Sum over all models" version of ALP. → Note 465.11 → This is really an alternative way of looking at t. 3 input unc. model of QATM.

In t. case of GAs, we also have continuous & discrete params to describe t. system — & we can use same soln. of .18 plus wtd mean over discrete param combins. We could use "Meta GA" to optimize these Params: They are initialization, params: Size of population, freq. of Mutation, freq. of Crossover & Nature of Mutation & Crossover, etc.

.23-.24 Theoretically solves all "theoretical" problems in QATM! : In t. case of t. "Control problems"

in Lsrch — We try various tricks to deal w. it & .23-.24 gives a criterion for selecting

.28 best way w/o wfng various ways. Well, not so clearly! .23-.24 gives solns. for CB=∞ only!

.34 For finite CB, it suggests 3 soln: Namely: sum/integrate overall params combins, Just

To pick the best Hessian pt. of all discrete params that satisfy t. CB. T. Hessians are

.36 constrained to params that satisfy t. CB, & so. (I'm not so sure of this!!)

So it looks like a good Genl. Soln. to t. QATM problem! (Except, of course, for t. Problem of Convem. in Lsrch — (which is sort of solved Theoretically "33-(34-36)).

Go thru 165-265 again & look for ideas that might be critical of this present "Grand Soln".

In doing this, also make list of impl. ideas worked on, & their reference cards. A sort of Biblio. Review, of last 100 pp.

# REFERENCE

pages 270-274

A new Method of Hill Climbing in Continuous Spaces!

Pick a small unit distance,  $D$ . (Smallest possible, yet large enough so that noise is small)

In  $k$  dim. spaces evaluate at initial pt. then  $\pm R$  ( $R$  is random direction in respect  $d$  ( $R=1$ ))

From those  $\pm R$  pick jump in that direction of Max G. (using parabolic approx).

Evaluate  $G$  at new pt. Then do trials  $\pm R'$ , again, but choose  $R'$  so it is orth to  $R$

by subtracting out that component. Then from  $\pm R'$  now pick a push in new direction of

best  $\pm R''$  but what we want  $R''$  ortho to  $\pm R$  & all previous  $R'$ 's.

Continue until we do  $k$  jumps; then start over, w. random var, or not necessarily orth to any  $R$ 's.

This method is "similar" to Peter's 267.26 (fitting a quadratic form into local function), but apparently cheaper, but not nearly successful!  $\odot$ ! .oo ff involves many function evals.

267.26 would involve  $\frac{d^2}{dx^2} f(x) = \frac{1}{2} \cdot 2k^2 = k^2$  evals of a func. to get  $\frac{d^2}{dx^2}$

(using  $\frac{d^2}{dx^2} f(x) \approx \frac{f(x+R) + f(x-R) - 2f(x)}{R^2}$ )

On constructing a Univ Lipp function set w. 3 inputs

In writing a Review, The stuff on 164 ff is quite good. There are a few

Things missing (T. universal / sup, a hard mix IA results), but the rest looks fine!

I'm uncertain about just what  $F()$  does: Say it (discretely) chooses a set of IA's for each  $Q$ .

Then we update  $t$ : see or "choose" IA's on that  $Q$ .

If we make a prodn.  $d.f.$  w. any  $Q$ , we first ask  $F()$  for relevant IA's, then

compute wts of those IA's w.r.t. that prodn. using  $t$ :  $\mu$  method of 272.26 - 274.24. (There are zero methods described there for assigning wts: One is fairly non-el.). We can also update  $t$ : param of  $F()$ .

While 19-23 may more or less work, it is not much in  $t$ 's spirit of what I used before! IA assigning  $t$ .

I.e., I was having  $F()$  make a prodn. of IA's, of  $t$  pc of each, to subset of IA's.

In 19-29  $F()$  doesn't output  $d.f.$  on IA's: it outputs a subset of IA's.

Eventually, the system would have to take the prodn on IA's & either select subset for prodn.

or use  $\text{cc} = T_{pi} \geq C_B$  for all of  $t$  IA's: While this is a  $C$  such, I'm not sure

sure its such a good idea! For almost all IA's,  $T_{pi}$  would not be close to  $C_B$  anyway

usable. Also (perhaps) most IA's have a fixed  $C_B$  for prodn, so giving them a  $C_B$  is

the same as a yes-no decision & gives us a subset of IA's anyway!

We may, however, not be able to know how much time a IA needs, in advance!

Another way to look at it:  $F()$  wants to find  $t$ : IA giving largest  $PC$  to  $t$ :  $Q \rightarrow A$ . — an open problem.

So, if we did do (L such, devoting time share  $\mu$  to  $IA_i$  — we find certain IA's have no output until a certain time — then they give  $t$ :  $PC$  of  $t$ : correct  $A$ . Other IA's given

varying  $PC$ 's to correct  $A$ . Their assigned  $PC$ 's may not be monotonic w. time! — But

its supposed to become "More trustworthy" as a IA spends more time on it. At the end of a "run" the single  $Q \rightarrow A$  update

$F()$  has gathered lots of data on  $PC$  (t spent) for all (opt) IA's that have been tested. This data is used to spec (277.11)

stochastic  
↓  
for S.O.P.s.

.00: : On searches for Medias of a useful kind (of 3 types).  
The "OSL" ideas maybe critical. Hvr, my impression is that coding these large operators is somewhat  
different from coding a sequencial corpus. [the: 277.00-00 makes them very similar]

If a subtree has occurred only once, we can still (perhaps "usfully, w. boost") use it as part of  
v. define of another tree or subtree. Perhaps an input distance betw. coding sequencial corpus  
v.s. Coding a Function: In i. Seq Corpus, Every thing has already occurred (all of text) —  
→ we have to parse it into subsegments: In Function Coding, we are inventing both  
the code elements (and to parsing) as we code it. □

.09 In t. SAAB TSQ's, Each Macro Op ~~was~~ (This was actually used for Procu.)  
was put into Macro: As a result, we had, at all times, a PRE Corpus that consisted  
of v. nested defs of all old Macro Ops. — These nesting order was same as t. ordering  
with t. assoc. problems were given. My Complaint at t. time, was that there were no  
other function defs! Later I decided that this was because v. probs were very simple!

Normally, w. harder problems, there would be new defs, betw. t. defs of Macro Ops.  
Also note that any sub-tree occurring in an OP can be used as if it were defined  
w. not very large extra cost. This is done by recoding, in which the first time that sub-tree  
occurred, it was man defined into pre corpus & used in all cases <sup>henceforth</sup> as a single symbol.

.19 So in this definition process, we don't have to specify where, t. sub-tree occurred.  
.27 This will be specified in Pre 2 (or more) actual uses of t. defined symbol (→ sub-tree)

In trying to modify a ~~Macro Op~~ Macro Op (during "Updating"), one can back track ~~by~~  
.29 by one or more Macro Op. Definitions. This may make it easy to construct Macro Ops that  
give some <sup>acceptable</sup> results <sup>perhaps</sup> as before — so one doesn't have to verify them w.r.t. t. entire Corpus.!

.29-30 would be great if true! Well it is a little bit true at least! ~~perhaps~~

Say we back track by t. <sup>2</sup> Macro Ops needed for t. last 2 QA examples: For t. new QA,  
if t. ~~prev~~ previous 2 QA's, we devise a recognition operator & a special OP to  
X from ~~prev~~ these Q's into good A's. We then need to test t. recogn. op &  
t. new special op on 3 QA's only. (v. recognition op may have to be tested on much  
of t. corpus, but there may be "logical reasons" why it has to work on the entire Corpus)



ID

Explanation-Based Lrng/Genrn (30) A Very Good Approach!

RNK et al.

But still, I'm not sure its what they mean. But this discussion seems Very Generally Useful.

100: 276.40: An implication of 276.19 = 27: Say one is Constructing Cond. Ops. One has  $OP_A \neq OP_B$  as a result for the choice of a new  $OP_i$   $OP_X$ . If  $OP_A \neq OP_B$  have been ~~the~~ byts of  $OP_C$  over or more times in the past, then the  $PC$  that  $OP_X = OP_C$  is hyper than if  $OP_C$  were not true. (This is because we can non-verbally define  $OP_C(O_A), O_B(B)$ .) It amounts to Context sensitive choice of operators.

The analogous Reg. occurs in 21: Say  $AB$  has occurred in the past; then if  $A$  occurs subsequently,  $B$  is more likely to follow (than if  $AB$  had never occurred in the past).

NE Much (if not all) of the discussion of 276.09 ff is true only for ~~the~~ references to the period of Macro ~~as~~ already parsed! If we look at the previous part of the Macro/primatives ops, there ~~may~~ may be many unusual, ~~non-predictable~~ repeated sub-grams - but we can't use them, unless we can abstract to primitive trees.

11: 275.40: Update  $F()$  (which then creates  $FC$ ) by iteration. So, we have this "update CC" available in  $F()$ . Q is: How to use this CC "most efficiently". T. output of the updating is the updating of some of the IA's - each "to some extent" (But some IA's would be completely updated in (finite) available time. T. Criterion for "Goodness of an Update of IA's" is unclear.

A poss. goal: That the IA's that are likely to be used on problems like the present QA, have been evaluated ~~for~~ using CC's that they are likely to employ in the future on such problems (Hrv Note 282.01!)

So: 3 problems: 1) updating IA's prediction using updated IA's 2) updating  $F()$  ... (Pro I may do this "by hand" at first - so it looks like a pure "unpredictable" problem.) 3) updating  $F()$  ... (Pro I may do this "by hand" at first - so it looks like a pure "unpredictable" problem.)

SN As far as I know, only one of the IA's (AER) is "Universal". (Also, the "R" input version of AER is "not bad". That it should use a "short R" to get more likely A's seems very "natural" & perhaps "workable/universal"! Well, GA can be Universal; Recurrent ANN can be universal; Vermeir's system (partly). On p 228: of the 16 IA's listed, at least 8 can be made Univl. T. non-universal can be simulated by U. - U. universal can simulate each other. 228? lives of 150 "15 IA's"

30: SN On Expln-Based Lrng. How is an "Expln" ~~Genrn~~  $\rightarrow$  a Genrn? : An "Expln" answers the "Why" question. It gives (a) short codes for "what occurred". Usually short codes imply Genrn. An "Expln" makes the data seem reasonable (by pc) by showing just how (for example) it is a special case of a more General rule. As such, it does not invent the rule; it mostly, if not most, is a another case of a "T. rule" so it re-mits update to (usually continuous) params of "T. rule". [It shows how the data is a special case of a Genrn. This has already been made. Some times the "Expln" will be the discovery of a new rule, in which case it is true Genrn. So, the new rule is stated; it is clear that the new data sets by the rule; It is implied that other cases have occurred in the past. If an explanation just shows how a set of previously accepted rules, it implies the present data, then I don't see where the "Genrn" occurs. On the other hand: If we show the new data can be accounted for on the basis of a chain of applications of known rules - this can be a "short code" for the data. - But how is it used for Genrn? - How does it "Generalize"? 278.28-30 has a useful example 278.00

→ Nitro 297

00: 277.40: So the Q is: How does a short code (compression), (Always?) imply causa?

Well, the short code means that the batched data is "more likely than it looks to be". So if one saw part of the data, the rest would be implied w. unusually high PC.

An example of an "Explain" would be the parsing of a message sentence. - How do we get "Guzen" or "Ling" from this? Well certainly, after parsing, we have a better position to

"Understand" the sentence more fully: - perhaps to give it an even better "Explain".

"Compression" ← find the regularities; Finding the regy implies long/pause.

Very often, an "Explain" will be compression by parsing - not by making new definitions.

"Parsing" shows how data could have been generated from already defined concepts (= old definitions)

Its clear how a new definition would imply causa/ling; But how does ~~find~~ parsing find for

the data do Q's?

In general, parsing can be "hard work": i.e. it can involve much such before one finds a good parse! → 297.00

13: 277.19: As I see it: 277.11-17 more or less (casually) defines what "updating IA's" is a bit. As for prediction, perhaps FC) has to choose which IA's are to be used. [It might be that just choosing

which IA's to update w. a particular QA, would be ~~was~~ a non-critical decision - but Not so! If we update a IA on data for which it is poor in exp, it will ~~often~~ often end up a bad mean score.

Obscuring the fact that it can't, indeed, be v.g. if applied to proper problems.] (So we would lose much of the benefit of that IA, because we would assign it a spuriously low wt. for problems it is very good at)

So in both updating & prediction, choice of proper set of IA's is very import.

Initially, TM would do updating in a very broad, (orthodox) way, to ~~be able~~ be able to discover of what IA's were appropriate for what Q's. Later, TM will want to delete <sup>POST-HOC</sup> certain Q's from it.

Sub-copy of certain IA's. <sup>Want to sure</sup> However, we need to know that at least one (a preferable more)

IA's are updated w. each QA, so the system is able to learn how to solve parse probs, <sup>with</sup> i.e. FC)

28: A single example of a parse could significantly change/parameters of context-dependent <sup>generating</sup> grammar rules, if the SSZ's were very small.

30: One kind of "explain" that implies induction: An AI's Q's can be "explained" by an Op <sup>or</sup> Q's

31: Short code that turns Q into A. We can then use this same Q<sub>0</sub> to induce other A's from other Q's. → 297.00

It would seem that selection of IA's for predicting is for updating could be almost identical... except that one might want to be more "diverse" in update, to explore IA's capabilities of A's in unexpected areas. So, if a ~~IA~~ would be selected for

answering Q<sub>0</sub>, then IA's would concentrate update on Q<sub>0</sub>. However, perhaps it

IA's is selected for update on Q<sub>0</sub>, IA's would not necessarily be used to help answer Q<sub>0</sub>.

A poss. annoyance: A given IA could have several disparate (non-overlapping) sub-copies. Each could be very useful in its own way - but because each <sup>sub-copy</sup> ~~has~~ implied a distinct update

history, ~~that~~ IA, scope for will have different skills! → Maybe not so bad at all - more good! → 279.00

↳ 297 on EBL

ID

Def. S.C. = Sub Corpus

for all of TM

00: 278.40! If we use a single universal IA, it could have several disparate S.C.'s. We could start each S.C. by giving it IA certain impl-concs to start w. That are particularly ~~useful~~ useful for that S.C. In fact, constructing IA/corpus pairs in this way could be it having a different set of IA's to start w.

T. idea of having a  $F()$  "front end" was suggested by how humans seem to work: i.e. one approach solves a problem ( $F()$ ) by deciding on a general set of tools and will use to solve it - also, perhaps a "General Plan" for solving the problem.

Abode

01: 269.40 For prob. is updating ~~and~~ all we need is the ratios (relative wts) of various IA's with particular Q. To obtain this using the optzn method quadratic form hill climbing of 267-29-269.40 it may be poss. to only ~~do~~ do a hill climb on relatively few wts, if all wts are not interdependent. It may be poss. to do an ~~optzn~~ optzn. on several indep subsets of coeffs of IA's. ~~This would probably take a long time to~~ ~~optimize~~ optimize on all coeffs.

~~Optimize~~ When one needs to optzn one set of coeffs: ~~the rest of coeffs are irrelevant & need not be optimized~~ irrelevant & need not be optimized

I suspect that it is common for Probe to be non-overlapping sets of IA coeffs: It would almost amount to having separate TMs! - But not so, since many coeffs of common concs are shared betw these "separate TMs".

18  
19  $F()$ 's learning how to chose IA's for each Q, is diff, because any change in  $F()$ 's params would involve (perhaps), some Modif of ~~the~~ IA's working on problems they never worked on before!

- But perhaps small changes in  $F()$ 's params wouldnt need this much, is verifiable info on (IA's working Q's) that has already been collected will constitute most of  $F()$ 's info needed.

~~How~~ How, improving  $F()$  will be a diff task, & toward this goal we should compute updates on many diff IA's (as in 278.34-36)

This looks like a serious problem. TM will have access to it at all times a large no. of (S.C., IA's) pairs [an IA can be simulated by a fixed, universal IA and ~~is~~ is equivalent of various S.C.'s. - This "Equivalent" could be in form of certain (Heuristic) concs being given to that IA at very low cost -  $\epsilon$  (its initialization)]

30 Anyways, from these set of (S.C., IA's) pairs, it must induce how best to propose a small set (or just one) (S.C., IA's) pair

31 for a new Q. So  $F()$ 's "self-modify", is simply normal induction. - As usual induction, its "update" usually consists of small Modif. of Con & various params, w/ less frequent invention of new concs & reparameterizing (a small part, or varying much in much concs available) of its  $F()$  operator.

Normally, if effective S.C. one IA will be order dependent. - Since the Q's are given in some order empirically, the order will best fall S.C. - How  $F()$  could decide

to all  $Q$  to a S.C. that would have occurred in its "Middle" - This would involve 280.10 spec (allocate?)

Reprising: 0.44 or usually not needed: .00 - .09

00: 279.40: REPARSING: This is not so common as I thought! Usually when a new (conc) rule is discovered, it is via OSL (SSZ = ~~2~~ 2 or 1) so reprising is unnecessary — Or, in general, the SSZ is quite small, or it would have been derived sooner! (The same concs are expensive so they need a large SSZ before they can be recognized as a large conc)

Occasionally, a conc. will be discovered "late" ( $\equiv$  by SSZ), simply because of search time (i.e. search technique could change w. discovery of a new search hour)

Since the conc. of interest only occurred a few times in the past, usually not much reprising is needed/warranted.  $\rightarrow$  307.08

10: 279.40: re-computing  $\frac{1}{2}$  of S.C. to get into order! Alternatively (Foster) we could just compute the update of the Q by putting it at the (temporarily) end of the S.C.  $\rightarrow$  as a cheap approx. For a large CB ordering is irrelevant, since the Q's are hierarchically "unordered Data"; TM could try all possible orderings to see which gave best PC; but for finite CC, coding of S.C. is incremental/sequential.

17 Normally FC) takes various S.C.'s & does a "trial augmentation" w. a QA. Usually this gives a low PC for the QA, but if it is by, FC) will ~~try~~ <sup>try</sup> a augmented S.C. as a proper S.C. for that IA. If the Q is not easily recognized, this will be a recognition will be an impl. cost. — So ~~the~~ an impl. part of FC)'s long consists of finding good ways to recognize that certain Q's should be added to certain S.C.'s of certain IA's.

22 This aspect of FC)'s long, in which it is able to take ~~the~~ pick (known) QA's & try (S.C., IA) augmentations, is separated long: long w. a teacher: (Rec to QA's are restricted to be readable set — TM can't create new Q to try, since there is no way for it to get to assoc. "A".) (Hr, later in TM's career, it may be able to take Q's of its own, or of an encyclopedia — at some large CC. — or for it to do RW experiments — which would be very expensive in "cc".

31 CHEAP Hrr, for most, if not all, these "Q trials" (which are "experiments"), TM would probably be using a very greedy algm. to decide which tries to make. T. cc for a trial would not normally be large enough to warrant much cc in selecting the trial! For very expensive trials (like QA & some ANN or RANN) the cc of trials will warrant much more expensive algms for deciding what trial to make.

37 Hrr, from the pt. of view of diversity we don't want to be too careful in trial selection (too "elitist") since we Need "Diversity"  $\rightarrow$  281.00

TD

Context dependent pc's of Concs. 22

of 268.28 - 269.40 - 279.09 ...  
 .00: 280.40 : T. Forgy (279.19 - 280.40) suggests that the "Empirical" optn. of wts & (s.c., IA) pairs is unnecessary. — That would do better getting wts via an "optimized"  $F(\cdot)$  function.  $F(\cdot)$  would find/construct various s.c./IA pairs & give perm wts for any particular new Q. & depending on exp. & empirical profs. ~~the~~ selection Alg. only every little (280.31-40)

.04 Note: In T. Forgy, I have ~~not~~ discussed effects of ~~the~~ limited CC in body choices of (s.c., IA) pairs & evaln of QA's (updates) by those pairs.

.07 Hvr, in any Universal IA the CC must be stated implied, since true universality is only poss. if CC = 0.

Hvr, I may want to draw <sup>up</sup> picture of QATM w. minimal attention to CC — just as a kind of Book mark / summary of present state.

Here  $F(\cdot)$  is able to assign pc's to each (s.c., IA) pair in its available sets

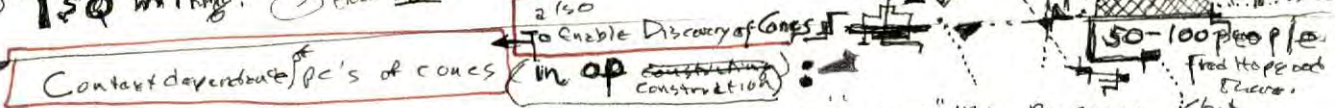
These pc's are used to determine which s.c./IA's are to be updated on what QA's —

(since pc will be assigned zero for all but a few s.c./IA's.) <sup>wrt a particular QA</sup> Also E. pc's will give wts to prodns of various s.c./IA's for ~~each~~ <sup>each</sup> ~~particular~~ Q.

Also  $F(\cdot)$  is "updated" by ① creating new (trial) s.c./IA's & assoc pc's; ② update of continuous & discrete params & a poss. reversion of part of  $F(\cdot)$ .

T. respec. QATM problem: ① Initial choice of IA's & assignment of s.c.'s to them.

② T&Q writing. ③ How Best to deal w. CC utilization



22

As a corpus grows, the size of contexts grows, so it's easier "the pc sense" (but non-negligible error in sense) to find legit contexts which can advise the pc's of various Concs. Whether whether this ↑ the pc's of Concs fast enough to cancel out ↓ in pc's of Concs due to their ↑ in number w. ESS, is unclear. Either this effect is sufficient to make contexts good enough, or we have to find other ways to ↑ pc's of Concs,

or we simply can't deal w. large corpus! (Away out of using large corpus is to use <sup>t.</sup> window or x window on ~~single~~ time history of corpus.) → 300.00 is progress on this problem

See 300.26 ff → 300.36 ff for very Goal. soln to T&Q problem! — 300.00 Goal!

IP

# PSM (≡ Prob Solving Machine) - QA machines, sounds too naive

.00: Back to review

.01 SN 277.11.17 on Meaning of updating "FC": My present impression of meaning of FC's

outputs: They are meant to be wts of output of IA's. We want those wts, ~~is~~

.03  $P_{ij}$  to be  $\Rightarrow \prod_{i=1}^m \left( \sum_{j=1}^{m_i} P_{ij} \right)$  (n = no QA's, m = no of IA's or QA's)  $\rightarrow$  Max.

$P_{ij}$  is / assignment for IA<sub>j</sub> gives to correct A for Q<sub>i</sub>.

FC could update its wts for  $\dots$  If FC updated  $P_{ij}$ 's via 268-28ff we would get the  $P_{ij}$  for each IA<sub>j</sub> - Instead .03 allows FC to give different  $P_{ij}$ 's to IA<sub>j</sub>, depending on what the problem is.

No! .03 is not right! I want to Max  $\prod_{i=1}^n \left( \sum_{j=1}^{m_i} P_{ij} \right)$   $P_{ij}$  is the prob that IA<sub>j</sub> gives to A<sub>i</sub> - (T. correct ans to Q<sub>i</sub>)

$P_{ij}$  is the wt that FC gives to IA<sub>j</sub> for problem Q<sub>i</sub>

If FC updated via 268, 28ff +  $P_{ij}$ 's would be idiotic: I, IA's would get

weights irrelevant w.r.t. problem numbers?

When we update FC's params we chose to continuous & discrete params  $\rightarrow$  continuous. Each of discrete parameters has a peak for contin. params, it has a Hessian associated w. that peak.

These peaks in FC's continuous params are as well as the assoc Hessian are available via 268, 28ff - But we do not obtain the  $P_{ij}$  (prob) this way

.00-21 Tells what our goal is in "Updating" FC: It says "via 268, 28ff for

the continuous params", but doesn't say how to optimize discrete params.

We should be able to get "PSM" ( $\equiv$  QATM) to work on that problem.

Ideally, I'd like to be able to make trials in FC so prob only 2 certai (perhaps small) set of its outputs are changed, so we don't have to get FC to evaluate

every large no of QA pairs. So: "Incremental" (changes in FC) is what I want. ("Incremental" in the sense that only a little  $\Phi$  (i.e. for a small no of Q's) of

FC is changed.

So .00-34 tells what FC does. (Presumably, it assigns  $P_{ij} = 0$  to IA's that it doesn't want to work on problem Q<sub>i</sub>.)  $\downarrow$  Hurr, Note 280.31 on desire for Diversity! This sounds perhaps quite diff from R. Goal of .00-34

The in .00-34, the search for good set of params for FC is an 02 problem - I'm not at all certain that LSrch is the best way to solve it: Because of "convergence" but w. the conds. Other than that "I fell into the trap" it may be the best way!  $\rightarrow$  283 of spec

00: 282.01: ~~I next big problem area in intro of sc~~  
 The IR's Parameters, are regarded as having input  $Q_i$  & presumably a p.d. on  $A_i$ 's in "some cc"  
 Note 281.09-.07: if any has Universal ops, CC has to be considered.  
 Another pointed Univ. ops was  $\begin{matrix} 259.40 \\ 253.37 \\ 263.00 \end{matrix}$  } ~~unrelated~~ 3 input unc; 264.12 equiv 1 input unc.  
 (I forgot what this point was! @)

04: 282.40: - Also updating the IA's is also an OZ problem - but in this case we usually have  
 "reasonable methods" to solve them: In the case of universal unc's we will  
 (probably) be using Lsrch over the discrete params that define IA.  
 T. recent previous notes have some context & Heron's comments on the proposed system.  
 T. stuff on 164 ff <sup>to programs</sup> is also of much interest in these regards.

10 Re: the CC problem: We could have  $F(\cdot)$  assign a time limit for each  
 IA. - But whatever we do, the complete characterization of a S.C. IA will have to include  
 the amount of cc used for each of the  $Q_i$ 's in the S.C. We can usually improve  
 (S.C. IA) by the cc spent on any of its  $Q_i$ 's. - but doing this retroactively can cause  
 serious problems! Say we get a better value for an early (or's way thru)  $Q_i$ 's!  
 This done by a new definition. To use  $P_i$ 's data in subsequent versions of the  
 S.C. IA, we'd have to do a lot of reworking, etc. - sounds ~~very~~ <sup>very</sup> expensive!  
 spend

22 So "We" ( $F(\cdot)$ ?) somehow have to decide how much cc to spend on updating each  
 S.C. IA & also how much cc to be used on each  $Q_i$  by each S.C. IA.  
 24 Perhaps use same function  $F(\cdot)$  for both  $Q_i$  answering & updating.  
 25 cc spent on updating/improving  $F(\cdot)$  can be "1/2 of all cc" - since  $P_i$ 's  
 is a "SI problem"  
 26 One poss soln to the cc allocation problem:  $F(\cdot)$  does, indeed, do it,

30  $\rightarrow$  we want something like  $\text{Max} \frac{\text{total input}}{\sum CC_i}$  **MAYBE!**  $\rightarrow$  See 284.26 for an IMPT Consideration  
 How one would go about searching param space & computing J. GORC, it's unclear  
 But it would be nice to have a clear "Parameterized" Goal for  $F(\cdot)$  (i.e. in context of system)  
 A desideratum of 30  $pc_i \rightarrow (pc_i)^x$  <sup>loss environment</sup>.  $x^+$  is a value that preserves  $\phi$  & maps,  
 so it preserves  $\pm$  normal pc.

In 30 for each  $\sum CC_i$  value would  $\text{Max} \sum_{i=1}^n (pc_i)^x$ , so perhaps it must be of the form,  
 $\text{Max} \sum_{i=1}^n w_i (pc_i)^x \cdot G(\sum CC_i)$ , w.  $G(\sum CC_i)$  some  $\downarrow$  function of  $\sum CC_i$ .  
 we want some functional equivalence betw.  $G(\sum CC_i)$  &  $\sum_{i=1}^n pc_i$ . As it is, we have a partial ordering,  
 with  $G(\sum CC_i) \leq \sum_{i=1}^n pc_i$ .

$\downarrow$   
 $\text{FC}(\cdot)$  modifi  
 $\downarrow$   
 $\text{FC}(\cdot)$  modifi  
 $\downarrow$   
 $\text{cc of}$   
 $\downarrow$   
 $\text{a pc of}$

Blochman  
Bribery  
Burglary

00:283.90: "Study Problem": Consider a single QA problem. Again we have a partial order of Govc below. CC is PC. 1 pers Good, & CC is bad, for any pb. in pc, cc is spec. but does a particular set of cc & set of pc make Govc better or worse? It would seem to depend much on the <sup>nature of</sup> particular problem. Say we are studying Nuc Reactor design. We have an even. funct that gives probab of failure. We have several empirical cases of failures. Would like them to be given to pc's of  $> 10^{-6}$ /yr. Then the even. funt now gives them  $10^{-9}$ /yr. Should we spend \$10^6 to get even. s up above  $10^{-6}$ , or simply regard these failures as "inevitably to be random events"? So its a factor of  $10^3$  in pc - are we willing to pay \$1M for it? - we ~~may~~ may!

On the other hand, there would be cases where a factor of  $k$  in pc would not be worth \$M.

Also note that  $\sum CC_i$  may not be correct: spending on some problems is not equiv. to spending \$1 on others. On the other hand, to get any particular result (revenue) [PC's] = 1/n, we would be interested in total CC's. CC is designed to be a function that is additive in this way - so perhaps  $\sum CC_i$  is ok. On the other hand, in SM analysis (in Fortune) seems to be the best thing to maximize!

10

It may be that Govc of 283.30 is good enough to design a very effective TM.

We can use it temporarily & keep working out problems of Govc Design.

19

If we change this Govc, it may not be a amount to much change in TM's PEA. (26)

So the Q is: Do I have enough detailed understanding of PSM to start on TSO's. If so, I must write a good report, w. a lot of footnotes for myself on details, expansions, references.

26.19

In 283.30 (assoc. problem) Do we want to include PC of FC (modifus w/o CC of modifying FC) in Govc? probably include PC of FC! as for CC of FC (modifus) ... I don't know + HVR, see 283.25-26 (50% solution)

The design of 283.10 - 24 is not about specific envt on just how CC's are taken into account. - I should have some more specific ideas before going to TSO's. 283.22-24 is about as "specific" as it gets. Assuming it's a probab. & updating problems have same CC assignment function: CC is a function of SC, IA, Q.

32

FN on 32

The updating from operation is to prod operations over slightly different involves PC of  $A_i$  (i.e. correct soln) only, prod involves PC's of all  $A_i$  of interest. So it normally takes longer. If we wanted norm'd PC's, we would have to put all of them  $A_i$  into. Actually, the PC for prod. is "up to the level". It could involve much updating of one or more SC IA's - it may indeed involve discovery of new ones. After a Q has been given "Answer" in full (i.e. lots of candid reviews, reports, etc), the updating process is trivial & involves  $\approx 0$  CC! HVR, in any case, TM has to run all of the SC IA's associated w. the particular Q. Note HVR: for good UMC w. large corpus, & difference bew. norm'd & unnorm'd PC's very small. -> 286.00 spec

33  
40

285.01



ID

.00: 284.40

So just consider update cc in  $\Gamma$ . eq. (Garc) of 283.30. (to start out!).

I think this problem is <sup>related</sup> to 284.00 - 16 - i.e. Each problem has its own "Importance level" -

So we want to spend more cc on "very imp. probs." On the other hand, a problem that is "unimportant" in the sense that its soln. is not of great value, maybe imp because it enables <sup>imp</sup> new cons to be discovered, that we're in utility in the long run.

.05

So how much cc to spend on a "update" is not clear!

Perhaps it is so imp, because in later updates one could add any needed cc.

Perhaps the best way to "do" a large corpus, is to update incrementally, QA by QA, but w/ small cc. Then using the primary cons discarded, go over it again. Whole corpus again is again - each time doing it incrementally - QA by QA - each time

ABCDE

using all good concepts previous runs.

Still, I have not addressed the problem of some IA's giving fixed probs in sharp cc values, & others (all Universal IA's) always giving probs that vary in cc - "better" probs for larger cc.

A Good Set of "Study Problems": Consider following cases:

1) all IA's have <sup>fixed</sup> fixed times for solns.

2) 1) but 1. IA's have different cc's for solns.

F() knows rough cc's in both cases

~~3) same as 2) but F() doesn't know cc's, & they may vary from Run to Run~~

A soln. to 3) would be close to solving the problem of Universal IA.

because we can consider a Universal IA to be sum of 2 or 3 "Take longer"

.28

In case (1), there's not much of a problem - except that IA's that "take longer"

might be equivalent to using <sup>several</sup> "fixed cc" IA's in parallel.

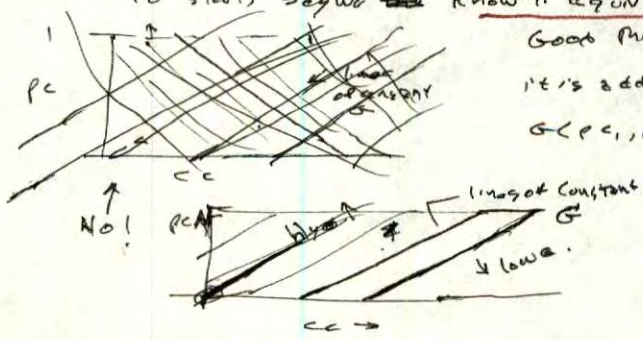
.30

To start, say we know 1. eqn. between cc & pc: so  $G(pc, cc)$  tells how

Good that pair is: we subject  $G$  to a Monotone function ->

it's additive wrt.  $pc$ . so  $G[pc_1, cc_1; pc_2, cc_2] =$

$G(pc_1, cc_1) + G(pc_2, cc_2)$ . We want to Maximize G.



IQ

Sec. 32

284.40 : On updating v.s. predicting! In updating, TM only has to find good codes for Aist. Correct Answer. It can quickly reject Ops & R's that don't give this. How it does by a variety of S-ops to get this.

FC) also creates new SC's for IA, & also creates new IA's. Also updates the Universal IA's.

For probn., usually FC) specifies Res IA's - that have already been updated - so it just runs a sequence of R trials to get a D.R. on [A]. As in note of 284.39-40; Normen usually is uneng m a 'nature'.

T. Moral's, predn is usually trivial. Updating is difficult task: T. only time predn. becomes diff, is when we decide, for a particular Q, that we'd like more updating, (more cc) so we do that & run do standard probn.  $\rightarrow \text{O} \rightarrow \text{I}$

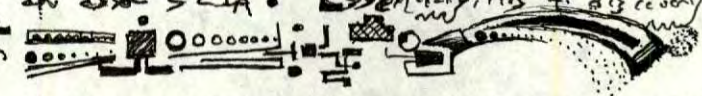
So: as for deciding how much cc to spend on what tasks: Updating is a main task.

10

Note: FC) also spends time creating New SC's for a given IA - these correspond to strong biases int. operation of that IA. If we only have one IA (presumably Univ),

All of our SCIA's are based only on data SC's - or on pseudo SC's. A pseudo SC is a SC that would have given rise to a particular IA; this IA having been created by the

Universal Reference IA (i.e. hypothetical) 'pseudo SC'. What FC) tries to do is to find Q's that have an identifiable similarity, so that they are easily answered w. a given SCIA. Essentially, it is a discovery of "problem classes" that are answerable



unable to contain problem solution methods - perhaps w. a certain pseudo SC augmentation!

20

So, T. MAIN problem ( $\equiv$  updating): Given a QA; which SC IA's to update (which SC's to add to QA): Who shall be new w. & when? How much cc to spend on each SC IA? (Default guess is  $T_0 \cdot P_i$  ( $P_i$  is old w. of that SCIA))

Note that a "certain cc"  $\Rightarrow$  a function of (sc, IA, Q) may not be the whole story - first, we may want a softer cutoff of cc perhaps - only  $\uparrow T_0$  (as m || Lsch) until 1. pc of Ai is "large enough" (whatever that means!)

Or, in some cases, we would be about to complete a study & write result m & by PC, but "it will take a little more time" - So how much time to give it?

Well, say  $.20$  (w.  $T_0 \cdot P_i$ ) maybe and for a half, working update alg m.

30

Each SCIA will have its update program in form of searches for short codes involving random fixed Q, modifies of Ops of course finding good R's. Probably use Lsch.

32

A part that is quite vague in my mind is FC)'s invention/discovery of new SC's for the main Univ. IA. It is certainly important. Part of it is automatically done when FC) updates a subset of SCIA's w. a particular QA: that QA is automatically added to all of those SC's. A more interesting (a diff) thing is to create SC's

28. write 279.19 - 280.40 + 280.10 - .30 is relevant to 32 280.17 - .37 is most relevant.

37

In 280.17-.37, FC) takes various SC's & supplements them w. various QA's of past, in trials to get new, useful SC's. In order to do this, FC) stores data on the history of each SC  $\rightarrow$  287.00

40

ID : [PSM]

start of 7/25/01 ~ 12:05 AM

00 : 286.40 : its pc is ~~unpredictable~~ its codes & ~~its~~ identifies pc's of components concs. of those codes. This makes it poss. for F( ) to create new sc's by breaking up old ones & adding new QA's to them.

↑ This could be done easier if we allowed branching of a sc into 2 new sc's ;

like original sc is  $Q_1 Q_2 Q_3$  ; This can "branch" into  $Q_1 Q_2 Q_3 Q_4$  and  $Q_1 Q_2 Q_3 Q_5$

So a big Q is when we do  $Q_1 Q_2 Q_3 + Q_4 \rightarrow Q_1 Q_2 Q_3 Q_4$  ; we have a new sc,  $Q_4 = 34$  ;

Do we ~~have~~ ~~to~~ ~~ok~~  $Q_{1,2,3}$  ? In 286.37 - 287.03, we do effectively store all such sc's, & their assoc "parameters".

So this effective branching construction of 286.37 ff could be the basis for any kind of sc. (construction / creation / Birth) : ~~286~~ 279.30-40 ; 280.10-30 gives some heuristic ideas on how this might be done - tries in 279.31 to show how this aspect of F( )'s behavior can be expressed as a normal induction problem - & kind PSM normally solves.

F( ) chose

A minimal thing we need to store about each Q is what set of sc's & what wts. they were given at ~~that~~ time that Q was first introduced into the Corpus.

The stuff ~ 270.00 to the present seems like a good understanding of how an imp. part of F( ) operates : Mainly on updating, but also on creation of new sc's.

HA! We can also make new (trial) sc's by AND, OR, NOT ~~operations~~ Boolean Ops on sets. — But F( ) must know how to do this so that it doesn't have a reasonable pc of being useful ; 279.30-40 has some ideas on how F( ) leaves things (like this ?).

- 1) not a good idea
- Q for sc's
- if 1. sc's are
- possibly - that Q
- 2) Developing
- initial sc's.

When F( ) guesses selects a set of sc's for a Q ; and if turning out, one of them does nothing (relative to the rest) on that Q ; Then F( ) may want to determine a way to give that sc a very low wt. for that kind of Q or to give it zero wt. — so that Q isn't added to that sc.

On initialization of PSM : i.e. construction of initial sc's :

we can start w. a <sup>initial</sup> Universal IA, & given problems that all belong to same sc. Next, we start w. some initial UME IA, & give a different subset of problems, appropriate to a "different sc" ;

we can do this for several problem areas — so we get several sc's.

The purpose is all done w.o. any F( ). In each case, there's only one sc's being tried, so it always gets all of the wt. If it doesn't get reasonable / expected by pc for the city A's then

we must acc or modify it.  $t \leq Q$ . (This is essentially Sol 89, but with stochastic ops, rather than

35 t. Deterministic Ops of Sol 89)  $\rightarrow$  <sup>288.05</sup> <sub>286</sub> ~~280.05 spec~~

Developing several sc's of this sort could be good initial training for Me —

to develop TSQ's. The development of F( ) can be done later.

It may be poss. to do PSM w. only one sc ! Perhaps look into this !

IO

(i.e. solns. of problems & sub-functs of solns)

297.40

00: 299.40

For SI of such a system, we want to improve the techniques that are used in searching for appropriate functions. Also, any heuristics that a human could discover, must be discoverable by the system. So a sufficiently advanced PSM should be able to usually work on the problems of finding heuristics, using its own history as a set of suitable Q's. (The Q's would have to be "standard" & self-generated)

.05: 280.35: We can do this very nicely, by having TM learn almost all many of the same type behaviors simultaneously: Just have clear indices ("names") on each QA examples so it can learn how to answer Q's of that index type. This way, TM develops skills that use ~~the~~ concs common to a various (scia - types). ~~Later~~

Later, we can omit the ~~index~~ info occasionally, & have PSM try to re-construct the index. - eventually, we would omit Roff indices.

For later problems, we may want to introduce indices again, to quickly fail PSM just ~~as~~ a set of Q's is in some common identifiable field, it should be treated in at least one kind of common class

For "f. report" I'd like to describe MTM version: 1. Ume has no "input".

Trouble is, that we need context dependent P's assignment to concs, - if we want to avoid one aspect of "Scaling problem". Stochastic PSM can obtain P's for concs & part of its normal induction routines, - this is not true of MTM. To work well, the context dependant P's concs, has to be an externally supplied Alg. -> MTM isn't able to improve it.

It's conceivable that at a certain (large!) point, MTM would stop growing & would grow very slowly in response to new input problems. (Kolmogorov/Lavin suggest that for  $C \approx \infty$  MTM would not be a very large Alg.) Hrr, w. finite  $C \approx B$ , I suspect this Mechanism would be very large - larger and the Latch to be in product.

I had that of the "recovery" context. How far one could go w. this is unclear. Depends on size of "recovery" & on degree of Skinnish in the T.S.M. Hrr, in that stoch PSM would not be any good at learning the contexts of concs. - I'd have to put it in by hand anyway. How advanced stoch PSM would have to be before it's able to work on problems of conc context discovery, is unclear. At first, we can start with "Recovery Context."

As for description (f. report) stoch PSM is not unlike to Deterministic PSM. T. assoc Ume has 1 more input: P's it!

ID

- Sections: §1 Intro 299.00
- §2 T. format for Training Data 294.00
- §3 The Prediction Operator 295.00

00: 288.40: So: Re: T. Report II.

Introduction: Deterb. QAS (T. nature of r T S Q). — why is called "PSM".

How practically any problem can be presented in to f. mechanism's form.

No! Simpler: Just talk about a string of Q's, & a stochastic operator.

~~Perhaps~~ Perhaps into is x f. Abstract!

T. intro. should deterb. f. problem.

§1

Introductory  $W$  is a <sup>large set</sup> sequence of <sup>finite</sup> pairs of strings. Pairs represent questions and correct answer pairs.

We are then given a new question string  $Q_{n+1}$  without ~~knowing~~ <sup>without</sup> ~~knowing~~ <sup>knowing</sup> what a ~~probability~~ <sup>probability</sup> distribution on possible answer strings  $A_i$ ? The problem is of much interest, since almost every many problems in medicine/law can be expressed in this format way.

~~On~~ solution to the PSM form of a probabilistic operator. Its input is any  $Q_i$  string; its output is a probability distribution on all possible ~~answer~~ <sup>answer</sup> strings. There are many operators of this sort. We choose one that assigns high probabilities to the correct ~~answers~~ <sup>answers</sup> in the data ~~set~~ <sup>sequences</sup>, when presented with the corresponding questions.

Reference by 32-35

If it is possible to search over all possible operators ~~of this kind~~ <sup>of this kind</sup> and choose the one that fits the data best, using ~~some~~ <sup>some</sup> ~~of~~ <sup>of</sup> this "Goodness fit" ~~criteria~~ <sup>criteria</sup> of ~~the~~ <sup>the</sup> smallest sets of Q, A pairs.

However, for all but the ~~simplest~~ <sup>simplest</sup> ~~problems~~ <sup>problems</sup>, this approach is well beyond the ~~computing~~ <sup>computing</sup> capacity of any computer ~~provisioned~~ <sup>provisioned</sup> in the next century ~~or thereabouts~~.

Approximation techniques are needed. ~~the~~ <sup>the</sup> "Incommensurate Learning"

For which the stochastic operator appropriate to a string of Q, A pairs is

We have a ~~set~~ <sup>set</sup> of question answer pairs  $[Q_i, A_i], i=1 \dots n$  we have a stochastic operator  $O_n$  that fits that data well. We are given a new data pair  $(Q_{n+1}, A_{n+1})$ .

We then modify ~~the~~ <sup>the</sup>  $O_n$  so that it fits the augmented data set.

We then "update"  $O_n$  by modifying the parameters that describe it — so that it fits

the augmented data set.

[T. part of ~~intro~~ <sup>summary</sup> sections of the ~~report~~ <sup>report</sup>]

30  
31

§2 This describes the QA pair sequence: stochastic v.s. deterministic ~~operator~~ <sup>operator</sup>. The concept of ~~the~~ <sup>the</sup> stock op. to generate an seq. of ~~QA~~ <sup>QA</sup> pairs. → 294.13

32  
33

We choose the operator,  $O$ , such that

$$2^{-l(O)} \prod_{i=1}^n P_0(A_i|Q_i)$$

is as large as possible. Here  $l(O)$  is the length of the description ~~of the~~ <sup>of the</sup> operator,  $O$ .

$P_0(A_i|Q_i)$  is the probability that operator  $O$  assigns to ~~output~~ <sup>output</sup> answer  $A_i$ , when ~~the~~ <sup>the</sup> operator  $O$  is given  $Q_i$  as input.

$l(O)$  is defined Qk. here

292. DO spec for continued work ideas in Reports

40

00: 289.40: Present tentative plan is ~ 287.28 ff: to use one & universal IA in several independent sc's "simultaneously" (2 in 1) for initial TSO's. Perhaps try them individually (each a bit) - to debug TSO's. In this system, there is no external FC: i.e. IA learns to do this on its own. (FC) & functions

(SU) One old idea on how to train TM: Start w/ determin. probs. After TM learns them & 2 problems bot w/ a little noise, then add more noise. TM should learn "about" probabilities" This way. (1998 notes)

Reading Old Notes on TM. Seems like a v.g. way to proceed ideas /

For "report" & Give abstract analysis of unc. & unvl. system.

[SN] The Score of 289.32-35 is "net bad" but we want to give TM as much info as possl. (for debug) on its trials as possl. A better score would assign "degrees of Acceptability"

to various & possl. As responses. Sometimes several responses would be about

This kind of Response is "user-intensive", hrs & w'd line to avoid it, if possl. May have some in early part of trial only.

Equally valid. Sometimes certain responses would cost a larger negative score.

The Score is, essentially, a Reinforcement Machine (TM) w/ ties (history)

It has no "Look ahead" capabilities. — which is fine! Much easier to program and

little (if any) possy of Machine Manipulating User. (See 291.23 ff for bit more Reinforcement of this kind)

If I went "Look ahead" I can get the PS TM to work out.

problem (e.g. see Cure Cancer 293.00ff) ...

[SN] Opt. possy of starting with MTM-type problems only: (2 D-PSM)

The history & Context-dependance of Concs. could be done by a GA system — or ANN or RANN. — slow, (but often "creative" & Mach "Diversity")

If we use GA, there is no BLOAT problem, since, for induction problems, we automatically make "size of soln." part of the "fitness function" in the best possl. way (for CO=00 induction)

In report: On "updating" discuss updating continuous param. values, finding new functions.

responses. As 3 main aspects of updating: Also discuss "Backtrack" by reynolds.

(No of Q's of Backtrack). Hvr, main development, addition of new heurs, of new methods of heur. discovery, will cover by my studying it. & @'s new heuristics seem to work best. This will suggest new to do.

Re (improved) (update) is a well defined problem: I should be able to present it to TM so it can work on it.

30

Perhaps have several ops developed in // w. occasional cross breeding (crossover)  
To get "Hybrid Vigor" ~~to~~ <sup>Several</sup> ~~times~~ <sup>of breeding</sup> ~~would have~~ <sup>somewhat</sup> ~~different~~ <sup>different</sup> S.C.'s ~~is~~ <sup>is</sup> ~~important~~ <sup>NO</sup>  
crossbreeds, - then Hyb vigor would ~~be~~ <sup>suddenly</sup> ~~much~~ <sup>improving</sup> ~~crossbreeds~~.

If we have several ops that derb. some S.C. & they have pc's influ = factor of 10 or so, say, we  
should save all of them, to be used in "Theory Revision" if some of them fail in a new QA.  
This is partly for better pc's but mainly to save in "Backtracking" (presumably storage costs less than trials).  
In the "2+3; 7x8 ..." sequence of problems: ~~2+3~~, (Sum 2, 3) works, (mul 2, 8 works)

T. changes is a problem: "2 → 7" or "3 → 8" are expensive to check (if it nos. are 32 bit integers!)

but "4 → x" is easy to do, if only 10 or so operators are available at that time. Fine.

- , & ÷ can also be learned thusly, but even easier because of "context"

(possibly recency concept among others). Try to find another or others: Try net dependence "recency".

Stochastic TM: ~~to~~ <sup>def's</sup> for Numbers could be found as special cases:

Integer Arithmetic:

Integer arith: ~ zero errors; Errors in float pt. often in last decimal, - but could

worse if many calcs are ~~done~~ <sup>done</sup> to obtain result. - or zero error if TM does exactly same calc as underlying "true" generator of data

T. idea of probab can Rules is rather different from D.F. of "accept to be"

answers "to" Q's in Nat lang. <sup>exp.</sup> See 290.10-13 on better way to deal w. QA's

where many different answers are about equally good.

Perhaps a better Goal: Say we know t. stock Op Prob generation

4. QA, say data seq. then a v. p. goal is to minimize  $\sum_j (P_{ij} - P_{ij}^*)^2$

$P_{ij}$  is t. correct pc for jth answer to i's problem.  $P_{ij}^*$  is TM's answer.

We could have data of this type in input TM, but ~~is~~ <sup>is</sup> very diff in later TSQ's

An approach in some cases: Either give RTM (for updating) a list of acceptable <sup>correct</sup> answers,

or, if no. of accept poss. correct answers is hard to anticipate: Have "teacher" look at

TM's/answers & give them "scores" (or yes/no). TM tries to get ~~the~~ <sup>max</sup> of

$\sum$  pc that it gives to acceptable answers. If scores are given for ~~each~~ <sup>each</sup>

entire set of pc values for a given Q,  $P_{ij}$  is a regular "RTM" - but history = 1, so

not so bad! (see 290.14 for details.)

**SN** ON THE IMPORTANCE of Z1 or AZ1 v.s. "Shortest Code".

If shortest code is used is used in a sufficiently narrow way; for a case that accounts

times, int. part(AZ1) can be k times as good as "MDL": MDL could have k diff

codes for  $P_{ij}$ /concept, which is very wasteful in Lst. A more sophisticated

MDL would use Huffman coding, so usually to difference between AZ1 & MDL would

be small ( $\frac{1}{2}$  bit or thereabouts) in this case. ~~MDL~~ <sup>MDL</sup>,  $\frac{1}{2}$  bit ~~in~~ <sup>in</sup> a code of length

2 bits is 12.5% - Not much impact for full searches. For a by pc conc, (pc: 75)

Th. error of  $\frac{1}{2}$  bit is (uses) but I don't think t. average pc will be bad by a factor

of 2 at most.  
But one imp advantage of (AZ1) is that t. pc's are easy to compute

ID 2

289.10 spec ← introduction. (abstract of lecture)

.00:291.40 : The main critical library in the present P&M model: The default idea is alternatives:

① T. form of t. data as  $[Q_i, A_i]$  pairs. A BAG of such pairs? Alternatives: set v.s. Bag; order to some extent (weak precedences allowed) | possibly Deterministic parallelizing. Some sub-sequences of QM's could be (nearly) ordered.

② Model of predictor:  $\Rightarrow$  ~~input~~ Universal device:  $Opdata, Q_i \in R \rightarrow A_i$ . (Deterministic models: w/d set of deterministic input models & set of output models.) Interpret R input as Random for Monte Carlo Model.

Doc part is  $[P_{ij}]_{i,j}$  → fixed, j varies for t. problem  $Q_i$ : List in  $\mathbb{R}^{P_{ij}}$  size - large crash  $A_{ij}$  is an output or by input  $P_{ij}$  is output.  $P_{ij}$  is input; output is all  $A_{ij}$  with pc's of  $P_{ij}$  (1+3).   
 [Mention n. Significance of Universal operator - must be guaranteed to work. For has a recursive nature of data & one enormous amount of computation time.]

Predictor give only one output (Gore is summary like no. of connected outputs, or can be

2 → Real Random. reduce This is a MTM. (deterministic). Look for stochastic if use  $A \geq T$  Any given posty? or List of operations.   
 We have many of these MTM's in //.

③ Gore: 289.32-35  $\sum_{i=1}^n P_0(A_i | Q_i) = \text{Max}$ .

use  $A \geq T$  Any given posty? or List of operations.   
 (use of parallel codes of some operator to a parallel operator)

④ Theorems: 290.10 Gore for each set of  $\{P_{ij}, A_{ij}\}_{i,j}$ : A form of RTM. w/d. "Look Ahead" Very Greedy

If we do post "Look Ahead", TM tries to predict future Q's,   
 Could try to influence PGM (Manipulate var). We could have   
 limited "Look Ahead" w. ability to predict, but no desire to influence var (?) ← is  $P_{ij}$  (posty/feasible)?

T. Gore of 291.18.  $\sum_{i=1}^n (P_{ij} - P_{ij}')^2$  :  $P_{ij}'$  are "low" pc's.   
 How pc of 0 ( $\approx 2^{-200}$ ) should enter is unclear.  $P_{ij}$  are TM's pc's.   
 See 291.18 for discussion & More Gore possys.   
 [It enters as  $2^{-200}$ , but also  $\sum_{i=1}^n (P_{ij} - P_{ij}')^2$    
 & enters as  $\epsilon$  minus square error in approx. - we choose Gaussian or better & best distribution to represent.   
 & generalization to feedback - feedback error]

④ Updating techniques: Main type: AZI: Update pc's of each; find new funcs; reusage.

occasionally back track to going back  $k$   $Q_i, A_i$  is ~~just~~ modifying t.   
  $\{Q_i, A_i\}_{i=n-k+1}^n$ . It's just AZI model

~~Backtrack to going back  $k$   $Q_i, A_i$  is just modifying t.~~

27.15 EP Perhaps discuss continuous param. opten. vs. discrete param. opten. The (Hessian)<sup>-1</sup> factor: At least   
 If not in the report, write it up for my own Remembrance/reference. It's probably applicable in a critical way   
 to Barron (1993) on General ANN. To apply it we need it in a needed to desc. t. D&E corpus,

30 exactly, by means of the "ANN" corb + t. ~~code~~  $\sum_{i=1}^n \epsilon_i$  square error. T. ~~pc's~~ post of t. set of   
 corb is the (Hessian)<sup>-1</sup> opt. set of corb. mult by  $\frac{1}{n}$  (for  $n$  adapts, is  $n$  error of  $\epsilon^2$ ) - multiply by   
  $\epsilon^2$  or  $\epsilon^2$  also NM part of  $\frac{1}{n}$ , maybe.

could be subject of Very Serious paper & maybe put Marcus w/o Jürgen to write it (?)   
 Very Serious paper.

36 ⑤ How this approach compares to other approaches to Machine Learning. → 293.20



IS

LISP

"Side effects"

292.40 : One aspect of Lisp I haven't investigated is "side effects":

298.00 : One simple side effect that doesn't affect things is printing.

299.00 : ~~Does~~ asking for or accepting keyboard input: Does this affect Lisp?

One thing that seriously affects things: Use of Memory  $M(A, X)$ ; as side effect, it puts X in Mem address A.  $MO(A, X)$  takes value of ~~Address~~ content of Address A.

These operations make it easy to get data from computer = functional way to code.

It seems unnecessary in my normal formulation of AZI: Each function can have its own output.

of any previously occurring function into expression. The trouble is, to get recursive

expressions, I had have to use a special device, & I don't know how to implement that.

[Well not ~~that~~ special: My functional had/function = functionals. The functionals could take

several functions & create one recursive function from them. — I don't know if I consider

a functional w. several ~~output~~ functions as output.]

A possible distance between my "functional lang" (which doesn't use ~~any~~ addressable RAM) & Lisp —

is that in PC's can be different. In AZI, the PC of all inputs to a new level of "next level"

function, ~~are~~ are the same: it's the same being the no. of ~~the~~ available outputs. In Lisp, the

RAM addresses can have PC's ~~for~~ for the content of the address.

Probably I will decide on modifications/variants of the language when I start doing TCS. 294.00 Spec

292.36 → HAVE SECTION in REPORT on How this problem/approach is (or derived from

Other probs/approaches: Comparison:

- 1) Sequential services process (Time Series ... ordered, unsorted)
- 2) Curve fitting: (P.D. is always constant value & direction. PSM is ungeneral.)
- 3) # Sequential/parallel by linear/n.l. regression (w/ to curve fitting: <sup>is then</sup> function of ~~the~~ previous values)
- 4) Reinforcement lang: ~~can~~ can be regarded as special case of Reinforce lang:
  - a) No look ahead
  - b) simple, special kind of ~~reinforce~~ function for set of  $[P_{ij}, A_{ij}]$  outputs.
 Is easier than Reinforce lang, yet ~~enhanced~~ <sup>suitably educated</sup> machine can be given RSM problem: Is more complex, but ~~more~~ <sup>more</sup> dangerous than PSM.
- 5) ANN (feedforward): ~~ANN~~ ANN is usually used as a continuous curve fitting Alg.

As such, ANN might be r.g. for certain kinds of PSM problems, but the error criterion used in usual feedforward ANN's is incorrect. He Schmidt Huber's "New Problems" is much closer to the goal used by PSM. [This is not obvious to readers, but I could outline a reason: Note that if used in any way "cross validation" is unnecessary & counter-productive. — wastes data. Usually small no. of sets is used #2 to Low Memory alg. is used to find soln.

Maybe discuss how ANN's universal for continuous spaces, ~~and~~ <sup>and</sup> ~~and~~ (viz. Barron's) — but not forcing universal for discrete. <sup>in Turing sense.</sup>
- 7) RANN: I think Ray was Turing universal: probably Ray was as difficult to fit to data as ~~other~~ other unex; Ray may be better for certain kinds of problems — it should be good for problems that are.
- 8) other LA's: see 150 & 228 for list.

2 2 2 1 2 2

Spec

00: 293.18; LISP (cont) One big Advantage of ~~main~~ RAM is cc. w.o.it, we have to recompute a value each time we use it. Info stored in RAM can be obsoleted & again we recalc. The AZI model, (which doesn't use RAM) could get similar pc's for funcs hvr. (Except for 293.14-17) → 294.00

There is also the possibility of using a universal lang like Jargon used in His "Lang has to Lun" paper. This is simply a RAM computer lang. w. only a few (probably universal) instructions. If more are needed for certain problems, add more instructions. Perhaps use **FORTH** - perhaps hypercompiled version (very fast) In fact, the initial set of "prim." insts (+ others) is moved to the "input" of the function spec described by the system. Perhaps I can simulate open systems university (systems 1 hr) → 294.00  
d. code by the system. Perhaps I can simulate open systems university (systems 1 hr)  
RANN or support SYM'S  
vector machines on substrate.

10 292.40 **SN** The intro (S1) should be a very clear, more or less exact presentation of the problem. So the later sections will only go into more detail & (perhaps) discuss Alternative choices.

13 299.31 **S2** The form of the training data. The training data is a finite, (unordered) Bag of  $Q_i, A_i$  pairs.  $[Q_i, A_i] i=1, \dots, n$ . They both  $Q_i$  and  $A_i$  are finite strings. We may think of  $Q_i$  as a ~~problem~~ pseudonamed  $A_i$  as "correct" answer, but this is not the most general form or interpretation. More generally,  $Q_i$  is the description of a problem, and  $A_i$  is a possible solution to that problem. (not always correct)  $A_i$  need not be the only solution. There may be more than one solution to the problem. Though  $A_i$  need not be an answer, it will be correct because of them. Since we are modeling a stochastic process, it will be answer.....  
If the ~~problems~~ problems are all deterministic, each will have at least one correct answer.....

**SN**: Give examples of  $Q, A$ 's at different stages of education.  
23 **.16** Since we will be modeling a stochastic rather than deterministic process, the data set should contain "errors" (e.g.  $Q = 1+7=? ; A = 3$ ) as well as repetitions. ~~parts of the data set~~ repetitions of the same  $Q$  (or  $Q$  as  $Q$ ) - sometimes with the same answer, sometimes with different answers. Since we will eventually want the machine to be able to use data from the "real world", it must be able to deal with ~~many~~ many errors in this data.

Another, perhaps much more important reason for creating a ~~probabilistic~~ stochastic probabilistic rather than deterministic ~~basic~~ problem solver, is that a stochastic problem solver can more effectively work on the problem of improving itself. In advanced machines, a good fraction of the available ~~time~~ will be devoted to this task.

Note In  $P_0$ , 0 is u.c.

29440 §3 The Prediction Operator.

Each prediction operator has at least one input - the finite string  $Q_i$  that describes the  $i^{\text{th}}$  problem. Its output is a probability distribution on all possible strings,  $\{A_{ij}\}$  that might be solutions to the problem  $Q_i$ .

This probability distribution might take any of several forms:

omit

1) A ~~finite~~ <sup>finite</sup> set of outputs in which the probability that  $A_{ij}$  will occur as output, for input  $Q_i$ , is just  $P_0(A_{ij}|Q_i)$ .

(1) ~~A~~ <sup>A</sup> sequence of pairs  $\{A_{ij}, P_0(A_{ij}|Q_i)\}$ , ordered so that the pairs of most probability occur earlier in the sequence

(2) In addition to the input  $Q_i$ , the operator  $O$  has an input for an arbitrary string  $A_{ij}$ . The output of  $O$  is then  $P_0(A_{ij}|Q_i)$ .

(3) The operator  $O$  has two inputs,  $Q_i$  and  $R$ .  $R$  is an arbitrary, potentially infinite binary string. The output is  $A_{ij}$ . If we feed random bits into the " $R$ " input, the probability that the machine will output  $A_{ij}$  and stop is just  $P_0(A_{ij}|Q_i)$ . The part of the input string that the machine reads before stopping may be regarded as "a description of  $A_{ij}$  with respect to the operator  $O$  and the problem,  $Q_i$ ."

In this paper we will be primarily interested in the last form.

omit

It can be realized using a 3 input universal machine. Though a universal Turing machine with 3 unidirectional input tapes, one unidirectional output tape and one bidirectional work tape, there are many other universal devices that are more

flexible in their behavior. There are many universal devices or processes that can realize this behavior. In our exposition we will use a Universal Turing machine with 3 unidirectional input tapes, one unidirectional output tape, and one bidirectional work tape.

The first input is a finite string that describes the operator  $O$ . Since the machine must know when this description ends, the acceptable input set must form a prefix set known to the machine.

The second input is  $Q_i$ , the problem description. It, too, must be a member of a prefix set known to the machine.

The third input is the binary string,  $R$ . This input is an arbitrary, potentially infinite binary string and is not restricted in any way. However, only finite subsequences of the tape that are read before the machine stops are output and ~~stops~~ <sup>stops</sup> values automatically constitute a prefix set.

IP

00: 295.40: The machine first reads all of the  $O$  input, ~~then~~ then it reads all of the  $Q_i$  input, then it reads all of the  $R$  input. During any of this reading ~~of the~~ data, the machine might produce some output. The bits only, if any, that is on the output tape when the machine stops, is regarded as the output. ~~But the machine never stops, it is regarded as being non-terminating.~~ However, for some kinds of problems,

(e.g. sequence extrapolation) the prefix of the output contains usable information and can be obtained before the machine stops. ~~and we don't care whether it ever stops~~

Well, I may describe the input function abstractly: then say it can be simulated by a machine.  $\Sigma$  or  $\Gamma$  read tape unit.

Then, say ~~it~~ for this direction we will be using  $\Sigma$  instead of  $\Gamma$  & functional lens —

Such as Lisp. [very reduced versions of Lisp have been described by Chaitin & Koze] <sup>for this model</sup>

Or just don't bother w.  $\Gamma$  tapes! — many advantages, is that it's easy to see that input

$R$  gives a universal def. with inputs  $1 \leq i \leq 2$ . inputs ( $\Sigma$  &  $\Gamma$  tapes).

The "functional lens": we want it to be able to describe any describable function of  $\Sigma$  limited strings. Not exactly! The finite  $\Sigma$  inputs over  $\Gamma$  input. ( $O$  is  $Q_i$ ) the  $R$  input is different in that it can be an infinitely long string. (well  $O$  is  $Q_i$  and  $\Gamma$  is  $R$  in this model)

While it can't see any part of  $O$  being possibly infinite string. —  $Q_i$  is certainly  $R$

could be seen by  $R$  — i.e. we could have printing before  $\Sigma$  machine reads  $\Sigma$

of  $Q_i$  or all of  $R$ . (or before it reads all of  $O$  for that matter!)

$Q_i$  could be the beginning of an infinite time series, & the machine could print out "predictions"  $R$  is for "of the next symbols of  $Q_i$  or of  $\Sigma$  next symbols of the process (of time series) to be predicted.

How infinite "O" string would be interpreted as unclear. One way:  $T$  length of

$O$  grows w. corpus (monotonic). Unclear  $\frac{1}{2}$  as to how to compute bits (pc of corpus).

Also, any increment in info needed for a particular  $Q_i$  (at a particular  $Q_i$  length  $n$  so far) can be given by extra "R" input.

So far, it looks like  $O$  will be a finite input & the machine will read all of  $O$  —

then it ~~is~~ covered  $Q_i$ ,  $R$ , & prints.

Perhaps increment in  $O$  input could give amount of "permanent/irrev" part  $O$  has acquired from  $Q_i$  plus far (unclear)

Consider the binary Bernoulli seq. w. parameter  $p$ .  $O$  could represent the "latest" value of  $p$  is well  $O$ . More exactly,  $O$  would contain  $n_0$  &  $n_1$  from which the D.F. for the next symbol could be obtained. All  $O$  entries is  $P_{n_0}(0) \approx P_{n_1}(0)$ ;  $P_{n_0}(1) = \frac{n_0+1}{n_0+1,2}$ ;  $P_{n_1}(1) = \frac{n_1+1}{n_1+1,2}$  & Laplace.

Drop  $\S 3$  for a while! I'm getting into some interesting points, but I don't want to use up my limited cerebral memory space. it just won't. Go on to  $\S 4$  on Prediction and Updating.

PREFIX SETS - 121

EBL 100

100: 278.14 (on EBL) If an explanation gives no predictive power (= compression), it is completely useless! "Why did she say that?" - "because she believes in God". This "explanation" is useful if we can use it to (help) predict her behavior: - which it can be used for.

If an "explanation" simply applies a well known rule & gets confirmation: then it only serves (only we get out of it is  $f$  in size of the rule) so  $\uparrow$  in its precision. prepare for it.

He flushed the room because a) He didn't prepare for it  
b) " " got any sleep with her  
c) " was ~~stupid~~ stupid!

08 Each explanation has different predictn. about his future behavior  $\rightarrow$  EBL 309.00

Wrote about Generalization of method of generating prefix sets, that was used in OSS (84 ff) see 121 ff ABCD etc

10 In General about PSM: I expect PSM will have various S.C.'s i. learn to switch problems betw. them. (at first by undering S.C.'s - later, less undering).

Note that various ~~methods~~ areas of induction, apparently using different IAs can be "simulated" by suitable language (sometimes very large) S.C. or pseudo S.C.'s. - seen this way, the present model of PSM, using only one IA, may be 0.4.

Perhaps write Section on "up doing" now. I am forgetting it, yet it is certainly the most critical part of PSM! Put only a few "work out" ideas in the report - (otherwise I will never finish!)

But do make a copy or refer to my notes, & w. expansions of various ideas. - Conjectures on what differs into arise, ~~and~~ etc.  $\rightarrow$  298.00

PREFIX SETS

21 in OSS 85 I had this scheme for L such in which I did an exhaustive check over all codes w.  $\frac{EC}{PC} < T$ . I conjectured that this ~~system~~ was

23 a completely general & complete! That all prefix codes were generatable by this mechanism. Then I found a counter example! (which I don't remember).

Anyway I was, now via you a proof of the original conjecture  $\rightarrow$  of 123.

It was a little this: If the set of acceptable codes does represent a prefix set, there must be a way to recognize when the end of the word occurs. If the UMC used such a scheme to describe the step. We would simulate that prefix set.

30 But I had the idea that the constraints on the UMC were rather wonder than just specific knowledge!

My impression: that OSS 85 simply listed trials in order of size of <sup>value</sup> binary trees. This listing prevented duplication: for programs, it was easy to keep track of what (prefix code) trials had already been made (i.e. were: illegal to ~~do~~ do or ~~to~~ to continue)

Also, I may have been mainly interested in computer prefix sets. (i.e.  $2^L \leq 2^{-L} = 1$ ) (In complete "prefix sets, it is of course necessary to have lots of code words of oo length!  $\rightarrow$  e.g. "all ~~and~~ codes starting in "0101" are code words. - (we could cut them off after extensions of  $> 10$  bits, say).

ID

**LISP .00**

**HW. note .12**

294.09  
293.13  
293.01

Lisp has these "side effects".

RAM is an impl. one! Hur.

Kleene's representation of recursive functions (say Prim. rec. functs) does not seem to use RAM / How Cons? It unlike to put every time a value is used, we have to generate it a-b-m-to from primitive cons. This is certainly O.K. from pt. of view of knowing what is representable by a system, but is not so good for cc or even pc |

For non-recursively defined cons, defusing unary: For recursively defined cons, defusing unary may be nary in Lisp.

In my functional lang, defusing are not nary, but they make pc's to grow, what belongs

to be, better pc's: uncontrolling perhaps cc's are better (functions are arrays realized by "inline" construction.

**SN** on **HW.**

I had this idea of a computer in which to main thing occurring was moving things around. Each string would contain refs to a seq of refs to other strings or to "primitives". A "primitive" is an actual H.W. function. We made pcu. lots of strings in Q's for use of various "primitives". When Queue too long a dup of a primitive is made. If Queue gets too short, its H.W. is invoked to make copies of functions wr. to long Queues. The "H.W." could be Programmable Gate Arrays. I do have a more lengthy writeup on this Machine Type.

Outline of Report (later)

In intro, a Summary: Tells what the problem is in some but not complete detail: Also what the solv looks like. Global V. S. Incremental Updating. Also use P. not 2-RCO - but say, option 2-RCO is programmable approx to P.

Defn of problem & A problem strings: But that of a problem's soln.

Success Operator is stochastic, problems should have occasional errors.

Note that this CEM is meant to be a general purpose intelligent machine. Perhaps Discuss its being able to solve most problems explainable to a human (worse explain. a broader problem to t. Hough)

The mechanics of stoch operator: Various coll. forms:

1, 2, 3 ... T is input unc input unc The Lisp machine is regular machine code.

Advantages of Lisp formalism ① Easy to Express Math ideas ② Easy to simulate reasonable Apprs via AZ

Advantages of Machine Code ① fast ② em simulate Lisp via Macros

Perhaps use FORTH (if Lisp = Machine Code).

How the stoch probabilistic functions are realized: How M(S, Q, R) can be model as a univ. d.f.: E.g.  $f(M(S, Q, R)) \rightarrow A$  is a nonuniversal distrib on A, w. R as random input to f. function M(S, Q, R). Hur, M(S, Q, R) is not universal.

The production problems usually very easy, once Q has been updated: put state of M(S, Q, R) in Memory, Run the random R, with return to RAM

In groups... Goal for system: Seq updates is reasonable approx to P. acceptable time - text + input of i. probs to get log on. in acceptable time. The draw: ① memory of data ② ISQ's ③ Some heuristics.

260.01ff  
on  
**UPDATING TRICKS**  
8 v. or ones

00 : 298.40 Global updates, used for early problems in TSO.

Incremental updates is usual updates methods: can be for one or several QAs.  
If several, they will <sup>usually</sup> be similar problems. Sometimes they will not be: e.g. Backtracking!

Explain It is usually by cc operation.

Before explaining incr. updating: Explain how updates are done on PC's.

First explain of Z1; Give examples refer to S 4b: Discuss update in Z1  
Discuss OSL in Z1

AZ1: Say "is not Z1" - Give examples.

Sol 64b update pc; no data; response  
data pc response.

Finding common subexpress is under Prin. Finely column

Subexps: BT OSL in AZ1 is also more diff: More prosys:

Any function has been used in past can be part of new op, w. of pc callon.

There are only a few of methods of update; I have not discussed probabilistic functions, explicitly

While they can be done in same way as deterministic functions, I suspect that

It is likely that they have different aspects.

This is only an outline of a few updating techniques.  
That many more techniques have to be discovered;

How I later expect to discover them using TSQ's:

- TSQ's are to teach TM
- to teach designer how to make better TM (= better update techniques)

Induction

Mention special update techniques in various branches of update techniques

See 304.24 - 305.18 for some other topics to (perhaps) discuss:  
It discusses TSQ's, transition from TM to STM - The idea of TSQ is from Mosseswell's TM.

In discussion of Z1: perhaps give my (ad-hoc) letter to Wolff (in the letter PS.)  
One of big problems will be to explain Z1: its updating process -

Then explain AZ1 & its updating process.

It will be easier to explain AZ1 directly, w.o. first explaining Z1:

Do mention techniques in which one must to solve new QA in isolation by doing:

a cheap of function for it, then find a way to recognize new it's data from all of previous

QA's (or "almost all") - since it can make errors - since it's stochastic Op)  
This - Errors are not allowed in MTM probs.

10

20

30

40

# CONTEXT of CONCS

→ Hvr, also note 4/17.09: on how "context" = "cond. probty" is: what QATM is normally working on.

ID

281.40 SPAC

00; 299.40: I had hidden sheet w/o. context, pc's of concs would → ± of  $\frac{1}{n}$  (n = corpus size)

← This is REALLY GOOD!

This is not true. As  $n \rightarrow \infty$  pc's of concs settle down to certain values.

These are uncond context free pc's = unconditional pc's. Since they are averaged over the entire corpus, they tend to be small (i.e. h' not small "small" about a particular situation (= context)).

When input TM starts only uncondly, pc of concs is not bad, because absolute no. of concs is small (< 20, say), so their pr's are relatively large for achievable Levesches (i.e. tolerable "CJS")

Hvr, as  $n$  (i.e. corpus size) ↑, we get into trouble of 100 mod.

Free A (perhaps the most important) way to deal in this is to use not context free (= uncond.) pc's for concs, but context dependent (= conditional) pc's.

However, one should not think of context dependence in a narrow sense (like in "context dependent" Grammars) but in a general sense of cond. pc - always

condition (= context) can be any thing - any info assoc w. the event (= conc) to be (predicted/have its pr evaluated)

Now PSM is essentially a conditional probability evaluator (20).

Perhaps rather than PSM, call it Condition Probability Machine CPM. (CPE) CPC Evaluator Computer Control Probty Calculator

So we could put the machine to work on the problem of assigning cond. pc's to concs. An impt problem is to decide what to use for 'context'. If we are very broad, the problem becomes very difficult. If we are very narrow, we don't get v.g. cond. pc's.

Bayesian Nets - Add to EAS.

Also note 4/17.09 "Context" = "Conditional PC". Lots look at the context not Broadway: Say we are doing incremental updates on  $O_n$  a previous problem, was constructed, using a certain set of partially updated functions, their issue, uncondly pc's & cond pc's for these functions.

The conditionals are on the previous problem ( $O_n$ ) and their previous  $O_{n-1}$  is its assoc function their cond & uncond pc's. (So a kind of (non-infinite) regress) of dependencies. So in the broadest sense of conditionals are all of previous conps, all of previous  $O_i$ 's & their codings & rpe's of their codings.

26 → 30 is, hvr, a bit clearer on just how the dependencies occur & support which should be strong, etc. [if context includes (past) time as well as "space", we can have (vns) during such]

Hvr, in general 26 ff is usually too broad to be directly usable. "If I narrow down to be directly usable: Anything comes 'Q'". See uptr of 301.33-.40 on how 26 ff is a bit elastic.

**N.B.** As realized by a 3 input machine, each  $O_i$  is a function is easy to realize as a Minter Carlo device. As we construct new Minter Carlo's, each modifn. of the operator (i.e. additional writing function) can be given a (conditional) pc by  $O_n$  - so we get a Minter Carlo D.F. of possible (301.00)

See 397.00-10 for different kind of "context" defined by any "ob". (= Junction string to Boolean (TF))



00:300.40 : **Ont1** trials. Two impl. cnds of this Generation: ① T. time needed to generate Ont1 should be much less than time needed to test it (? - I'm not sure of this, hvr). ② I'd like a way to do small (trial) modifications of  $O_n$  (conservative approach) we in attempts to get a fast soln. ("Quick & Dirty").

See 382.11-15 for some ideas on what "small modifs" could mean. This is Best description. On is an ensemble (S.D.F.) is our target Ont1 is also S.D.F. since for being represented by "populations" in GA.

03: Ont1 trials for constructing Ont1: They can be done via determination of L's over the R's of the various nesting functions that have to be generated (Note 33) 18 days to go so ~ 2pp/d. not bad - But not Great!

06: [Remember that L's are "imperfect" because of "correlations" betw. f. cnds] Also Remember that usually (almost always), Ont1 will be constructed a couple of few (say 3 or 4) functions. - a Reminder: PC's are too small & CJS gets too big

09: So: 300.36 - 301.09 could be a final soln to the TM problem! I do have to work out (300.36 - 301.09) - how to "narrow down" & broadness of desc'd in 300.26 - 301.09. T. "small modifs of  $O_n$ " is not easy to do: but see 382.11-15 for good ideas. SEE 302.00 - NOT BAD APPROACH - ALSO involves Learning During Learn

12: Well, not so final! There is still the problem of correlations of cnds in L's. A pretty many other impl. unsolved probs: BUT, it DOES look Very Good! Other Very imp. probs are General Updating Algos: see 260.01 ff for v8 good ideas. The 300.36 ff is a v.g. updating Alg! It says that  $O_n$  is to Building intelligence that Builds Ont1: which is the way people seem to work! -> 302.10

21: To Teach TM how to do "problems" like 300.26 - 301.05: Give ex examples with an "index" that indicates these examples are all of "the same type". - (i.e. a "SC"). Still, it's not clear how I get these "examples"! - Just what are they? - Give one example! When I <sup>see</sup> a problem is a TSO, write down all the (cnds) I can think of that would be needed/used in these solutions. TM must have acquired these (cnds) before it can solve that problem that way!

24: 8/1/01 INDEXING: This seems Very imp! It is an easy way to breaking corpus into SC's. Also, the INDICES can form a tree, so they can have varying amounts of "overlap" - It's would be to normal tree structure of categorization devices - like an "Outliner".

33: T. analysis of "context" of 300.26 - 301.09 is shaky E1, i.e. since that its output is a f. over functions. It does not consider the inputs of the functions or whether its outputs to be "final output" or "lower level" output. A less "E1" approach would use the same "context" info as 300.26 - 31, but its goal is not a "useful function", but **Ont1**. Actually, it may be easier to assemble a usable corpus for the goal of 306 than for the goal of 300.26 ff: i.e. finding a good function to continue building out the partially constructed Ont1. But it is a very usefully @/21 of 36. SPEC 302.10



8/1/01  
LD

Corrects bank lines  
Learning Dorny Lsrch .08

Effect of CC < B on System  
(22) 285.00 - .40

.00: 302.40 : In each Mt. Carlo context ~~MAX~~ of 302, 37-39, we first penalize each ~~context~~ G by an amount that depends on its correlation with cards already tried. Or - alternatively, Peters picks a way to modify G's ~~behavior~~ of funds, before they enter Mt. Carlo contexts. T. forgo. is simple. because the search via Mt Carlo is a v.g. way of doing

.07  
.08  
Lsrch. Also Note the forgo. technique is applicable in all processes of learning during Lsrch. which was a principal criticism of T. True optimality of Lsrch as a prob solving method (Assuming "All info was int. P.D.")

.10 Since I will be using Lsrch for critical operations CPM it is imp't that it be "close to optimum" (.08-10) plus the older ~~eyes~~ eyes out. ~~the~~ optimality of Lsrch (contains factor of 2) makes it likely that Lsrch can, indeed, be v.g. possibly, indeed "within factor of 2" (But note ~~the~~ what this means: Using Lsrch as I do, is no worse than using an optimum method ~~with~~ method w. at least 1/2 clock rate) Here, an optimum method at 1/2 slower clock rate could be many times faster than Lsrch. i.e. it could take much less time to get same results ← this last seems unreasonable, since changing clock speed will always speed by exactly same factor! ← Go thru 6. original arg't about "factor of 2" is a very sound arg't as I remember. Also a arg't that optimum rate of sol'n. to "problem" could speed by very large amt. (30)

.20  
.25  
.22 **NB** In forgo. Analysis, I haven't much considered CC limits, KB's. - Pro Play Are certainly critical to operation of system. (I wrote a Note ≡ this within last month... ~~but~~ A possible conclusion may have been that it was not very hard to take CC into account for the Model of CPM that I had at that pt.) 285.00 - .40 Has v.g. ideas 285.38 - .40 on Equil. betw pc & cc is v.g.! - I don't know what relation T. vector is, but .1MPC is a reasonable guess. I can use any assumed functional form ~~of~~ of require. 2 Use of of set CPM to find optimum behavior in view of it.

.30 (21) T. Arg't. of .08-21 - One: Optimality of Lsrch with factor of 2 "if all info is in PD" ; Since Lsrch's main activity goes on in CPM, this becomes a strong arg't. Part's being more optimal! As I construct CPM, a specific weakness in T. arg't. may become clear, as well as ways to deal w. these "Non-Optimalities" NOTE, HVR: This CPM is not truly optimum since it Doesn't Do "Lookahead"! Its "horizon" is 1. I can work out General RSM problem (see 243.00 or 243.00) a "line search" Perhaps I could Allow CPM to have "lookahead" within a problem of certain types. T. Big Problem "Lookahead" is the machine's Manipulating User.

T. Weakness  
This is the weakest pt. in T. Arg't. in general it is not easy to divide up T's behavior into 5's and 10's.

ID

REV

Of Present CPM: The path I expect to take

1.75 x 4  
50 x 4

28 + 7 + 1/2 = 47 1/2

Summary of present state of system:

1. Steady state behavior for a (somewhat) Measure CPM (Cond. Prob. Machine)

$O_n$  is state of Machine operator after updating on  $[Q_i, A_i]$   $i=1|n$ .  
Possibl. cond state

A pd. on  $[O_{n+1}]$  is obtained ~~to~~ <sup>from</sup> ~~update~~ <sup>by</sup>  $O_n$  w.  $Q$  input consisting of <sup>set/sequence</sup> of  $(Q_{n+1}, A_{n+1})$  (This seq. is also a relatively short ~~seq.~~ of  $Q$ 's.)

Expanded Th. current deriv of  $O_n$  in terms of ~~an~~ <sup>an</sup> ~~levels~~ <sup>primitive & defined</sup> ~~func.~~ <sup>func.</sup>

[for discuss of  $O_n$ 's input for  $P_{n+1}$  updating problem: see 300.20 - 301.05]  
[Also, a CB (comp. bdu) should be an additional input to  $O_n$  - ~~but not a part.~~]

Lsrch's done on  $P_{n+1}$   $O_{n+1}$  cruds, it's Goal of 289.33 is what what

we use Lsrch to try to Maximize. ~~S~~

See 302.00-303.20 for discuss of how Lsrch should be modified to deal w. "Corvins" betw.  $t_i$  cruds. The end of this <sup>discu.</sup> <sup>is on</sup> <sup>optimal</sup> <sup>of</sup> <sup>Lsrch</sup> <sup>(in</sup> <sup>its</sup> <sup>in</sup> <sup>its</sup> <sup>mt. P.D.)</sup>

$O_n$  has, as input, a  $Q_i$ : It has an aux input  $R$  (Random) ~~input~~ <sup>input</sup> <sup>void</sup> <sup>mechanical</sup>.  
Random input "R" ~~enters~~ <sup>enters</sup>  $O_n$  <sup>output</sup> (defined to be output when machine steps)

(.12-.13) tells how  $t_i$  systems ~~use~~ <sup>use</sup> ~~cruds~~ <sup>cruds</sup>,  $Q_i$ .

2. On the nature of  $t_i$  Q's: Any <sup>primitive</sup> learnable / reaction is usable (I don't mean "learnable" in  $t_i$  sense of Computational Logic Theory of <sup>behavior</sup> ~~relat~~ ~~et al~~). So  $Q$  questions; Answers, Problems & Solns; Any ~~action~~ I/O relation that is defined (or not to be defined) by a finite string or a finite string plus a finite no. of <sup>params.</sup> continuous variations.

So we can use CPM for curve fitting, linear/non-linear Regressu. It is expected that it will also be able to answer  $Q$ 's about a data base - but I haven't worked out much detail on how it would do this.

3. "Initialization" allow the PM to "steady state".

We start out w. MTM probs (Deterministic  $Q$ 's) from Algebra, Math TS  $Q$ 's.

Initially  $O_0 \equiv A$  & we do a search over deterministic functions, using AZI to assign pd's to ~~cruds~~ <sup>cruds</sup>. This is for ~~very~~ <sup>very</sup> ~~short~~ <sup>short</sup> ~~TSQ's~~ <sup>TSQ's</sup>. ~~we want~~ <sup>we want</sup> ~~2~~ <sup>2</sup> ~~with~~ <sup>with</sup> ~~max~~ <sup>max</sup> ~~pc~~ <sup>pc</sup>, ~~max~~ <sup>max</sup> ~~prob.~~ <sup>prob.</sup> ~~with~~ <sup>with</sup> ~~responses~~ <sup>responses</sup> for ~~all~~ <sup>all</sup>  $Q$ 's.

The initial pd's on ~~sub.~~ <sup>sub.</sup> ~~functions~~ <sup>functions</sup> to create  $O_1$   $Q$ 's will be unconditional pd's (essentially, & same pd for all situations). When no. of ~~cruds~~ <sup>cruds</sup> is small, this is a solvable system, tho it is slow.

It is necy to get conditional pd's for each primitive or defined function  $t_i$  (conditions) being some of those discussed in 300.26-.33. ~~initial~~

Another trick is to derive new  $O_{n+1}$  w. ~~basis~~ <sup>basis</sup> of ~~min~~ <sup>min</sup> ~~modification~~ <sup>modification</sup> of  $O_n$

Min ~~number~~ <sup>number</sup> of  $Q$ 's have to be verified on  $L$  proposed  $O_{n+1}$ .

I expect to learn how to do .31-.34 by studying how humans learn  $t_i$  simple tsqs ~~but~~ <sup>but</sup> ~~CPM~~ <sup>CPM</sup> will be initiated ~~on~~ <sup>on</sup>. But we will use for CPM's initial tbl.  $T_i$  method of 307.17-.19 ~~use~~ <sup>use</sup> ~~very~~ <sup>very</sup> ~~common~~ <sup>common</sup>.

Ess of "Falsifiability": To propose theories (models), that ~~are~~ <sup>are</sup> ~~easy~~ <sup>easy</sup> ~~to~~ <sup>to</sup> ~~test~~ <sup>test</sup>.  
To propose theories (models) that ~~are~~ <sup>are</sup> ~~easy~~ <sup>easy</sup> ~~to~~ <sup>to</sup> ~~test~~ <sup>test</sup>.  
 $L =$  max expansive / cheap

$P(A_i|Q_i)$   
= 1 for  $i=1$   
= 1/n

**REV**

30

.00

The transition from DTM to <sup>↓ Deterministic</sup> STM: <sup>↑ Stochastic</sup>

When we want probabilistic answers: We make t. index "s" appear as a Q that need stochastic answer  
We start out w. ~~probabilistic~~ deterministic problems w. occasional random errors. i.e. P.D. over A's

.10

~~These~~ ~~prob~~ det. probs - but w. several equally good answers.  
Next, we may try Bern seqs. <sup>Just Binary, then all symbols</sup> Then  $Z_{12}$

.08

Try learning curve fitting: <sup>First</sup> Linear run non-linear: It could lower or zero error cur fit  
First, then learn MS error curd fitting then other error criteria.

.10

Then HMM's: SEE Section 310.03-311.40 for a very nice way to make TS Q's for <sup>NMTM</sup> ~~TS~~, I have to back into Len Dren (Copterberg brot up on MTM problems!)

.18

.20

.30

.40

4. Sub. Corel  
Indexing of Q's (301.21, 29) <sup>as.</sup> Very often Q's will be given once or more  
"indices" that tell "what kind of problem" is being given. Some indices correspond to  
"Context" for a human! <sup>It could relate to</sup> Where a problem came from, what areas of science, what  
t. problem was given, what other problems is it "similar" to.  
As with the aspects of Q's, CPM is not told what they mean: It has to be "induced"  
from t. data.

~~The~~ elements of a sequencial TS could have a common index plus a "Time" index  
Alternating a Q could consist several "synchronized (time locked) time series, i.e.  
problem could be a regression problem: to find an eq. that probabilisticly relates ~~and~~  
t. elements of one set of seqs to previous data in all off. seqs.

A Great Breakthrough of 310.33ff is a way to do NMTM very nicely!

Gen notes:

1) On QA: T. accuracy of QM has been mainly tested w.r.t. & positioning of peaks & zeros of different patterns. I think it was likely that ~~the~~ exact shape etc. patterns based on DS & Parks is likely to be wrong!

2) I have been using a "UMC" in CPU to replace implement. could be by DF.

While this ~~may~~ may give vector p's it can be very poor for CC. Lisp users & RAM, via its "side effects" ID 293.00, so it may be potentially as good as "sequenced" machine. T. idea of most info proc. copy per ~~the~~ spent is max Q. So ultimately,

I want to get even result from TM w. min CC — & this probably involves modification of H.W. A machine that redesigns itself using **FPGA's** (Field Programmable Gate Arrays) — like in GA's, would probably be v.g.   
 See ID 298.12  
with a

3) I may want to write up a detailed Aug. & brief 6/ Factor of 2 of optimum ~~is~~   
  $\approx$  If all info is in PD! Mainly to make it clear just what assumptions I'm

making. see ID 303.08 - .20 — In particular remarks (303.20-22) R  
There is an apparent bug in Aug. that I need to analyze! i.e. 303.18-20

4) In the Hessian of  $\vec{F} \vec{G} \in \mathbb{R}^n$ : ~~the~~ Rank is: no. of non-zero eigenvalues   
 So  $(\vec{A} - \vec{I}X)$  must be expressible as  $\lambda - A$    
  $\leftarrow$  The actual Hessian is close to  $\sum_{i=1}^n \gamma_i \gamma_i^T$    
  $\neq B$

a b c

5) Laplace's Rule   
  $\{ \geq 1 \}$  is Laplace's Rule: continuous coding } say we have a Binary Bernoulli seq. w.  $p(0) = p, p(1) = q, p(q) = 1$

Using 2's, 1's pc of  $n$  seq. of  $n_0$  0's &  $n_1$  1's

$$\frac{(2-1)! \cdot n_0! \cdot n_1!}{(n_0 + n_1 + 2 - 1)!} \left[ \text{for } \alpha \text{ alphabet of } \beta \text{ symbols: } \frac{\prod_{i=1}^{\beta} n_i! (\beta-1)!}{(\sum n_i + \beta - 1)!} \right]$$

To code bin. Bern. — y. whole seq. — takes  $\approx PC = (p \cdot q)^{n_0 n_1} \cdot \sqrt{\frac{p \cdot q}{n_0 n_1}}$   
we code p & its "width" is  $\approx \sqrt{\frac{p \cdot q}{n_0 n_1}}$  so that's  $\propto$  t.p.c of decaying p.

$$\frac{n_0! \cdot n_1!}{(n_0 + n_1)! \cdot (n_0 + n_1 + 1)} = \left( \frac{n_0}{n_0 + n_1} \right)^{n_0} \left( \frac{n_1}{n_0 + n_1} \right)^{n_1} \cdot \sqrt{\frac{n_0 \cdot n_1 \cdot \pi}{n_0 + n_1}} \cdot \frac{1}{(n_0 + n_1 + 1)}$$

$$= (p + q)^{n_0 n_1} \cdot \sqrt{\frac{n_0 \cdot n_1}{n_0 + n_1}} \cdot \frac{\sqrt{\pi}}{(n_0 + n_1 + 1)}$$

So that 2 ways of doing it agree except for factor  $(n_0 + n_1 + 1)$

So, Bern. Doing the whole sequence by  $p + q$  & width of p knowledge

takes  $\propto \frac{n_0! \cdot n_1!}{(n_0 + n_1)!}$  as opposed to  $\frac{n_0! \cdot n_1!}{(n_0 + n_1)! \cdot (n_0 + n_1 + 1)}$  for  $\geq 1$ !

AH  $p = \frac{n_0!}{n_0 + n_1}$ ?   
 Which gets into answer   
 Perhaps just take  $\int_0^1 dp \cdot p^{n_0} (1-p)^{n_1} = \frac{n_0! \cdot n_1!}{(n_0 + n_1)!}$   
 Therefore (1) codes Bern. coding seq. & are exactly the same as uniform coding. ~~spec~~  
 That it should be identical (to f. sequential coding result, is surprising!  
 307.24  
307.00

Case counts

00:306.40 : ⑥ In Z1 : T. method I described. letter to Wolff!  
 I did have formula for PC of corpus when one <sup>kgm</sup> ~~kgm~~ was defined. I then had to decide on  
 a specific parsing of t-corpus & use that for next ~~defn~~ definition trial. This specific parsing  
 gives rise to specific Case Counts for all symbols & resultant PC's for those symbols.

IN AZ1 I'd ~~have~~ <sup>now</sup> to decide on a particular parsing: perhaps not  
 The O's are created (along w. their parents). This gives PC's to all component substructures.

So, any sub. tree (found) could be ~~found~~ <sup>found</sup> Essentially OSL.  
 08:280.09: Note that in Z1 & AZ1 the probable natures of t-problems are quite different. In Z1, we have  
 t-corpus subset of primitives, but ~~perhaps~~ <sup>new</sup> definitions must be made in t-corpus must  
 be parsed with them.

In AZ1, we have created t-O<sub>n</sub> & we know its <sup>unique</sup> parents. To create O<sub>n+1</sub> we look for  
 common subnodes at O<sub>n</sub>: Since only a little has been modified in O<sub>n</sub> (from O<sub>n-1</sub>) we still have  
 many new subtrees to check for repetitions (most subtrees have already been checked for  
 repetition on previous O<sub>i</sub> constructions.)

⑦ It may be that it is very common to solve a new problem w. a special function  
 and way to recognize it. latest Ph.D. Do this many times then look at sup of  
Solns & simplify them: Anyway, this avoids having to fast t. whole corpus.

⑧ Certain cases, involving MB, Logic, communication w. USSR (i.e. language in which problems  
 is expressed) can have special treatment, special indexing — to denote they are of much int.  
 perhaps equiv. to very by SSZ. Therefore certain things we want TM to believe w. & close to I.

⑨ Once of t. indices (scs) could be an "Advice Cleaned".

24:306.40 ⑨ Actually, 4-306. ~~result is trivial~~. Each conditional pc =  $\frac{n+1}{n+1+2}$  is obtained  
 by taking ratios of pc's of augmented corpus <sup>in</sup> unaugmented corpus. So the product of these, a varying number  
 of acceptable pc of t. final corpus!

I think this result is close to Lemma 2 of S78T3, which is b. cuts of  
 t. mom? Doug Campbell may have read scott's note — it may have been trivial.

"trivial" <sup>conditional</sup> <sup>empirical</sup> <sup>part of corpus</sup>  

$$\prod_{i=1}^n \frac{P_i}{P'_i} = \frac{P_n}{P'_n}$$
 updates only from n to n+1: That t. update is easy & so easy  
 for expected values of source quantities is

R. Big Q,  
 say  $b_i \equiv \ln P_i$   $b'_i \equiv \ln P'_i$   
 so  $\sum (b_i - b'_i) \equiv B_i - B'_i$ . Does this imply  $\sum (b_i - b'_i) = \sum (B_i - B'_i)$ ?  
 (with multiplicity)  

$$\sum_{i=1}^n (b_i - b'_i) = \sum_{i=1}^n (B_i - B'_i)$$

Also, is it true in any of whether source  $b_i$  is a measure or some measure? In particular, I suspect  
 that  $\sum B_i$  doesn't have to be v. same as  $\sum b_i$  — Princeton theory parallel. The only part that  
 requires  $\sum B_i$  to be w.t. "true p"  $\equiv (\sum b_i)$  is the thru prob t. K-L distance is ways  $\geq 0$ .

Gen notes

u / u

38 pp in new pen  
NEXT slide of ~ 309 is added to L. ALVIN (I) Pen.

What people mean by "Real World" is "Yesterday's Physics".

Real World

(10) The idea in "Sci" / "Philosophy" is that "There is a Real World out there":

I think they mean that we will continue to get data in accord w. our old ~~model~~ present "instructive" models of the world. These "instructive models" are essentially nothing more than "Yesterday's Theories":  $\therefore$  That leaves in "RW out there" = T. world continues to get in accord w. Yesterday's Theories. — This is the idea of Heggen? (who wrote a book on "T. end of physics") that we will soon be at the end of serious medicine of physical laws (is "probably zero" from my pt. of view)

[for serious ~~physicists~~ physicists. This is accord. to Professor of cafe community]

With his little ball announcing "Closing time Gentlemen — The party is over" —

A nice way to say that / Hans Fintel(?) us & his wife fine parties; At a certain pt. / her would get to his piano and began playing loudly, E.G. don't think we know then that / it was quite late and ~~party was over~~ — ~~the party was over~~

It was time to go home — the party was over.

When — talks of the end of physics — I hear this loud piano music in my head — but I cannot believe this to mean is long dead and it cannot be ~~believed~~ that the party is over.

Man — Heggen talks about "end of physics" — I heard melodies in my head — but

I find ~~it~~ <sup>not</sup> ~~is~~ <sup>at all</sup> a ~~renewment~~ <sup>not</sup> ~~compelling~~ <sup>at all</sup> ~~the~~ <sup>the</sup> ~~party~~ <sup>party</sup> will ~~be~~ <sup>not be</sup> over ~~in a long time~~ <sup>in a long, long time</sup>.

(11) SN ON STETS for Arby radix! + analysis of Lemma 1!

This appears to be a. have part of V. Diagram! The Lemma 2 by itself, is stronger than d. Diagram!

Anyway's possible to <sup>proof</sup> ~~proof~~ Lemma 1: Say we have k components B & PC!  
 $\sum_{i=1}^k x_i = 1$  T. diff. bew. k-2 expressions is zero if  $x_i - x_j = 0$  ( $\bar{x} = (k)$ )

(Probably) its partial derivs = 0 in all directions, ~~along first lines~~ along first lines.  
T. lemma is  $f(x) = \sum_{i=1}^k p_i \ln \frac{p_i}{x_i} = \sum (p_i - x_i) = 0$   $\therefore p_i = x_i = 1$ . (Anyway a proof with be poss. in part

of Gacs on Li-Vi 1998 p 329, first FP after eq. (5.6)

$p_i = x_i$   $\forall i$  is k ~~constraint~~ <sup>constraints</sup>,  $\sum p_i = \sum x_i = 1$  gives k-2 constraints on 2k vars.

So k-2 dim space.  $\therefore$  T. value of f(x) in that space; Next show partial derivatives of f(x)

are all zero on that sp. Next show second derivatives all  $> 0$  in that sp. —  $\therefore$  that.

Second derivs are parabolas  $> 0$  in a large space; OR that first derivatives are  $> 0$  on our region, but  $> 0$  in another: These are regions not divided by other material ~~is~~  $x_i - x_j = 0$ .



-00:297.08 : I had postulated that any legit "Explanation" is a "compression" & all compression can be used in induction. Compressions (perhaps "Recording" or just "Coding". Does this necessarily imply induction? If a coding implies ways to extend the code, then it is inductive.

E/G. In many Models of the data, we have certain fixed, certain random params. The implication is that this coding will continue into future in the same way.

If a coding is by showing that the data is an example of a general rule (maybe call this "Explanation (or coding) by induction"), ~~then~~ it fits an old rule, no update params (SSC) of the rule.

If it's a new rule, we put it in our "dictionary" or "rule set" w. its initial (params: <sup>Tomorrow's teacher</sup> ~~parameters~~).

We may have to reparse to enter corpus. (we may put it off ~~later~~ (memorize)).

-10 Reversing can be a key job. In theory reversing a Scientist will ~~have~~ have to think all of his scientific experience relative to "New Paradigm".

One trouble w. ~~statistical~~ covers "Execution Complexity": It does give codes for corpus & its extensions, but no clear way to extend the codes - since machine steps a step each code.

Well, perhaps not so: Any method used to code a corpus implies ~~that~~ that method could be used for other corpus. E.G. the BDS method of coding "Non-linearities", can be used for any finite corpus. How, its non-incremental means method of ~~code~~ to use for prediction - But, in theory it can be used for prediction.

If one uses "continuous" universal D.P. then one has appended codes to the corpus - any continuation of the code give ~~prediction~~ prediction w. proper params.

-20 How, in "discrete" univ. D.P. we have self. Dec. codes. There is no way implied to extend the code. - Unless we regard the code as an example of "BAG" (with an implied ~~to~~ stack loop, & a resource ~~(SSC)~~ of it. log. output. In this case the codes used in the data, give partial grammar - ~~is~~ is the Grammar can be used for induction.

INGENERAL: It may be that Each "method of coding" has its own SSC. method of extension of "induction". If it does not, its hard to see any

"Value" in it! -> but see (38)!

-30 Explanation = UPDATING: To the extent that this is true, Explan. can involve discovery of new codes. "Updating" means ~~the~~ modification (usually incremental, but it could be global) of a predictive structure.

Any <sup>compression</sup> coding of the data (even ~~new~~ new, uncoded data) <sup>is</sup> implying future compression.

Even a random seq. imposes a non-uniform d.f. on the future.

36 SN: A long ~~and~~ (holmes) random seq. for it to be followed by (01)<sup>n</sup> is now much less probable than if the initial seq. had not been "random".

38 Functional v.s. Placate Explan: The "Unexplained" gives a false feeling of uncertainty, our minds doom... functional A functional explain (or understand) can deal w. this feeling - But a non-functional "Placate" explanation, seems understandable & will often fix the "unexplained" feeling. The Non-productive utility of an explanation is "the code power" (309.005 per)

ID

Modifn of Universal Stoch Operator

Looks very Intst



309.40! Pleco Explanans: Large Industry buildn supply Pleco Explanans (Philosophy) (Religions) Also part of industry's Things that don't need explanation, but seem to: A Pleco explanation will be supplied by the "Industrious". (Also News Paper Headlines: "Market Goes (down/up) on news of (larger/smaller) no-merger Hostile takeover economy in spite of ↓ in volatility")

305.08 On the "Universal P.D." used for CPM: I had in mind a 3 input machine since I know it is universal in many ways I need it to be. It is not clear, however, how one can easily express other IA's in it. So want to take idea of taking a bunch of IA's known to be useful & write a Grammar that expresses them as a set of symbols. V.g. T. Grammar is likely to be universal if I do comb. on it IA's.

Very Good!

10 This looks like a strong step in the direction of solving a serious problem: i.e. I didn't have a good way to express IA's that I normally use, in a "compact" format.

12 I try to list a list of IA's (as I do) & (one) of them (at least) (was) Universal. Each IA has a "unconditional" pc of being chosen for an applic. Later, I'd want these pc's to become "conditional" which is what f() was for. f() looked at the context of a problem (i.e. now, it would also look at its "context" the one could simply regard "problem+context" as a finite problem den. So it could well be that my "factoring" system in to f() is a set of IA's may be a not-bad ELN of a problem that was very diff. in "Normal" form.

20 Now, now, I think I'd go beyond having IA's be a set of a Grammar. I think I'd want a [IA] in which each of the IA's used a subset of general rules that could describe any regys in a (probable) delay. I'd have my largest set of concs used for sketch, prodn. — each of the IA's would prob be subset of these concs or combn.

21 I'd seem to do prodn. — This seems like better use of General Prod using a "IA Grammar". I.E. "f()" or its equivalent would construct "IA's from this grammar." A by "show Update" is done. Down simply update Grammar & its continuous & discrete parts of Grammar?

29 I think the lang. approach processes "Universality" in the sense that any describable pd. can be expressed by a "finite string" + "continuous parameters". → 324.30

30 In particular, when I use various IA's in solving TSO — I will look to use the lang ideas to implement enable CPM to directly solve these probs in some way I think I do. → 311.00

34 Any poss of starting out on SM data? Well now I know roughly what kinds of data to use, but I'm forming of the regys are not certain to me. — For this reason, I don't see that I'd be able to take on TM anything viz TSO using modern data. If I could get very old data it is likely to have sample regys to learn from, (at least I will use "p'd figure" & analysis. — Try to find most popular books on this. → 311.00

- .00: <sup>spec</sup> 31090: On Learnable P.P.S (cont): In TM's general try about P.D.'s: When  $\epsilon$  is given  $\geq$  P.D. in ztsq, (w. proper index; so TM knows a probabilistic soln is needed), it will include a set of functional forms, a computation methods, that it could use to solve those problems.

I will be periodically, (umps & bumps - from time to time) be putting in new probabilistic methods, I should try to put them in "factored form" so they could see some benefits (analogy may be) better solns, - is useful structure/elements to derive correct forms of stochastic algms. (These new forms would be constructed mainly when lots of CC was available for prob. solving).

- .10 Just as an Algebraic TSO gives TM "basic" operators like  $+$ ,  $-$ ,  $\times$ ,  $\div$ ,  $1$ ,  $-1$ ,  $x^2$ ,  $x^3$ ,  $x^{-1}$ ,  $x^{-2}$ ,  $\sin$ ,  $\cos$ ,  $\sin^{-1}$ ,  $\cos^{-1}$ ,  $e^x$ ,  $\ln x$ , etc., along with their derivatives in terms of more basic operators (we may want TM to know a better complex version early on its TSO). - TM will have  $\frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2\sigma^2}}$ ,  $\text{erf}(\ )$ ,  $x e^x$ ,  $x^n e^{-x}$  as our useful P.D.'s available to it. For ex. fitting, it will have linear & n.d. eq. solving methods - including methods to check for errors in computation - as well as methods to correct for them.

It could have access to ANN, RANN, GA or any other such methods -

- .20 but these involve, usually onerous CC allocations - so I'd now to have

Special ways to train TM to use them.

So MTM probs are what TM starts out with. <sup>Gore</sup> ~~more~~ is then to find model of ~~w. max pc~~ max pc (pc of ~~600~~ in terms of model steps = 1) When NMTM probs are introduced, TM has a larger ~~larger~~ <sup>conc.</sup> ~~larger~~ set of conc. to use in attempts to solve problems, & its general methods are the same as before but more  $\rightarrow$  a bit of trick - We will pre-train Gore by pc of complex m.o.b. model.

- .30 So starting out w. MTM is ~~not~~ not best, since going to NMTM involves very little ~~change~~ change.  $\rightarrow$  like 300.26 ff

- .32: 8/12/01 This could be because 310.12-33 involves very context dependent conc. selection: In fact, 310.12-33 is mostly talking about NMTM, rather than MTM!



28/8/01

TM:

Grammar  
~~POD~~ Disc: This may be an adequate & GOOD (practical) method to do it. Critical points: 314.20 (315.30) (315.13-24) (Gore)

VERY GENERAL!

3/4 \*  
P313 may not exist

100+

! Another sure way of getting answers, Hill climbing.

314.30 is 3.15 - 24 another General for any kind of Grammar

Try to find a set of words  $\{A_i\}$  &  $\{B_j\}$  > to pc of defining phrase trees

to pc saved in using them, is max.

The pc of using them is not care product  $\sum \{A_i B_j\}$

The doc is one if  $A_i$  or  $B_j$  occurs, it's more likely to be all over  $\sum \{B_j\}$  member. Each time a phrase such appears is defined, we'll pc of corpus.

Other multiplication, is Boolean Addition of phrases gives a CFG.

Each such grammar rule (with a/o boolean Addition) will be a pc of corpus.

Also, ~~other~~ various grammarless (with a/o boolean Addition) ~~are~~ are found, we may be able to find pc of code by adding or deleting elements from various texts.

To. Finally, can be adapted to other by use of Grammar rules, ~~text~~ context dependent recursion.

like (context dep. substitution, words "phrases" - ... etc.

Grammar Rules

20

Superficially, this looks like an outline for a complete solution of word Grammars (discovery problems: IS IT really?)

phrases

printed file  
beam &  
G-best script?

Well, our Q23! IS is good enough to discover, and you want to

"Get off + Get?" It may work - but perhaps ~~only~~ only in large corpus.

The "size of corpus" needed, is related to how many cases of  $A_i B_j$  (concepts) are

in the corpus. If we have 3 phrases (A)(B)(C) then we have ~~more~~ more posses of

instances of concepts - if we have a rule about rules in Grammar.

30

One could study various such Grammars - see how large corpus must be before contain "instances" of specific Grammar / rule appear in corpus.

u30 30-31 - Make up various sample such Grammars, w. various pc params.

Get recurrent frequencies of various concepts of phrase elements (esp. "leaves")

Say we have discovered many word pair concepts w. unusually high frequency.

If size is large enough, one can see out word pairs for

A LOT



8/2/01

Grammar  
Axiom

By proving out a pair of words, assume it's to members in a system  
 for context: So if M & N are sets of 2 sets ~~that~~ i.e. digests one would  
 need M & N create P's: Use 2 sets only  $M+N+1$  ~~parameters~~ are needed.  
 One can construct total PC's needed to describe corpus — using Parsec 2 sets.  
 The "info" needed to describe corpus using 2 sets of continuous parameters,

is described in IP250.09-19  $\int \prod p_i^{x_i} p_3^{x_3} \dots p_n^{x_n} d p_1 p_2 p_3 \dots p_n$  ( $\sum p_i = 1$ )  
 Use an integral over all of 15 ~~space~~ of parameters,  $\sum p_i = 1$

Integral is 4 PC of 4 corpus with full model ( $\sum$  subset parameters),  
 consider 2 sets of  $M=2, N=3$ , so ~~2x3 = 6~~ v.s.  $2+3+1 = 6$   
 So no gain!  $6 \times 5 = 20$  v.s.  $4 \times 5 = 10$  is by difference.

In case of  $M=2, N=3$ , subgraph  $\prod_{i=1}^2 \prod_{j=1}^3 p_{ij}^{n_{ij}}$  w. constraint  $\sum p_{ij} = 1$  (?),  
 $2 \times 3$  matrix

for  $2+3+1$  parameters: integrated  $\int \prod p_i^{n_i} d p_i$   
 $M, N, R$  correspond  $\sum p_i = 1$   
 what's constraint PC? Maybe just  $\sum p_i = 1$

or, w. pure Bayesian  $2+3$  parameters  $\equiv$  frequency of ~~5~~ symbols  
 constraint  $\sum p_i = 1$ .  $\prod_{i=1}^5 p_i^{n_i}$ ;  $\sum p_i = 1$

14.16  $\sum_{i=1}^2 p_i = \sum_{j=1}^3 p_j = 1$  may mean that 1. interval can be expressed as  
 2 factors.

It may be possible to devise a ZFC-like formalism to "compute"  
 the  $\int$  of  $\int \prod p_i^{n_i} d p_i$   $\frac{(B-1)!}{\prod_{i=1}^B n_i!}$  or products of P's  
 "THE ULTIMATE BURNING OF GOD"  $(B \geq N_i)!$

Even if G. Garsia et al's are correct, the problem is to  
 decide which words to put in which groups. What is needed is  
 [ a formula telling how many put out one entry by adding  
 or removing 2 word from a set. ]

This is like in my letter to Wolff about 141.

Second Advisor	} started.
Technique 0.1	

PS Grammar Discovery

001; 315.40: T. methods of 3 (420, 3513, 24 are especially searching for discovery of rules.  
 Once one has a reasonable no. of them, it becomes possible to discover other Grammar rules of similar type! (i.e. to discover (≡ Sub (syntax)) occurs in paths perhaps almost all (if not all) Grammar types but have been never proposed!

Criticism of the Techniques

Because of Grammar rules appear from

$N_1 \rightarrow N_2 \cdot N_3$ , ~~it~~ will exist but do not correspond to NT's.

.10

(Non-Terminals). — So we will discover many "rules" that do not correspond to NT's in "Target Grammar". A BIG problem is to find rules that pick a small subset of these rules, & produce a good Grammar from them.

One Mile process: Corpus incrementally, start & starting w. a small corpus generate all useful rules, — Pick a corpus size, Grammar variables, etc.

This way, the no. of rules & w. corpus size. Generate rules that are

Created earliest, and perhaps most likely to be usable rules for

.20

more complex Grammar rules.

At first, the rules are created in pairs: (Any leaf, terminals  $\rightarrow$  is a leaf rule).  
 Anyway, once we have a few rules, we can find  $A \rightarrow B, C$  (A, B, C are terminals)

we can try various rules, like  $A \rightarrow A \cdot B$

or  $A \rightarrow A \cdot C$  (C is leaf & is also a rule); or  $C \rightarrow A \cdot B$ , where

C has been found to be created via C, D, etc.

In general, the starting & only way to get rules is by doing  $C \rightarrow A \cdot B$ .

.30

(How one or two can be leaves). — the count of rules is equal to rules is equal.

Also, since only rules will be rules: First candidate rules will be

$\exists b \exists a \cdot b$  in upm. We can then try  $[b] \rightarrow [a] \cdot c$  or  $[b] \rightarrow c \cdot [a]$

How, the size of the rule generated, will (usually) be an increasing function of SSZ (≡ corpus size)

# Grammar Discovery

00: : One way to Usefully study Grammar discovery techniques like f. ferris:  
 Pick a small family of Grammars, (say 100) use  $< \mathcal{G}$  grammar rules  
 Pick one of 1 family & use it to generate a corpus of size "S".  
 See how small S can be & yet one can identify 1 of Grammar.  
 Increase size of family of Grammars. See how S must grow to enable  
 Identification, using a certain set of heuristics and Prods (like f. ferris I've  
 outlined in ferris discussion.

10: What kinds of Grammar pairs are hard to differentiate, using these techniques?  
 If 2 Grammars are hard to distinguish better. — is prediction, using a. Per  
 Grammar, about to occur? If so, it's a bit of a trick.  
 Hur, if 4. predicts do much better, one should be able to use this fact  
 to get better version betw. 2 Grammars. [Good Predn is usually equiva  
 to Good "UPDATING" — Updating means both Q & A are given in all cases.  
 ideally, one should be able to tell how good a grammar is by forcing Q & A  
 set, w.o. having to "predict" ]

20: One way one had a set of Heuristics (like f. ferris) for Grammar discovery.  
 One has a stack Grammar that generates (usually small) <sup>stack</sup> Grammars — each  
 Grammar generates its Corpus. The rules are applied, it's seen which  
 Grammar rules have trouble w. This is used as feedback to USQ (M) to

20: generate new, better Heuristics.

General Remark: I think 319.00 - 317.26 is a very promising, very General approach to  
 Grammar discovery — works w. practically any type of Grammar.

30: For each Grammar type, there may be certain words ~~that occur~~ that occur w. unusually high frequency, even w. small (corpus) S. We should look for these  
 to discriminate betw. various Grammar types.



REV: .00 also note .22: latest earlier Reviews of Report outlines

Perhaps various Stages of Evolution (Subscript)

: Another run at an outline of the "Report".

1) Describe 'problem, data, [Q<sub>i</sub>, A<sub>i</sub>] w. some detail - Give examples of different kinds of Q<sub>i</sub>A<sub>i</sub> sub-corr - a Major business 'medium' so TM 'knows' what kind of "Q<sub>i</sub>" it is.

2) Descr. M(O<sub>w</sub>, Q<sub>i</sub>, R) → A<sub>i</sub> : Discuss practical facts for O<sub>w</sub>, Q<sub>i</sub>.

Stopping condition for R.

3) Discuss Goal. If > 1 O<sub>i</sub> is v.g., we will need a Plan - "Priority Review" = "update".

4) Operation Plan for prediction

5) "Steady State" operation for Advanced System.

Updating as "normal problem": either CB, even as "Runtime Problem".

6) Initialization: a) Early TSOs.

b) Early Updating Methods: Diffr. for each IA, but general Sum. Identifiers

Use of L such: + "concepts" problem: How being by "long during" "Catch"

7) Later TSOs.

perhaps non-Greedy L search  
"Growth" "Core Concept" problems

Wts. Unit Disk and; AMN wts. of sorted prior populations of GA sets.

8) Any problem set operation can be considered to be a QA problem: either given to O to solve, or a preliminary solns can be derived by user & user can try to find TSO that would lead to that soln.

17  
Note: General  
Master for  
USP

Look at early writeups on this Report: See if I can use any of the material.  
312.11  
304.00 - 305.40 is latest on most advanced form (Using Series: included in comp. "Breakthrough")  
V160 was first start of QA idea  
Before that was 289.08 method: 284.10 formal QA date 295.00 T. Pardo operator

2) 298.21 - 299.30 Again an Outline of the report.

3) 304.00 - 305.40 (Very comp)

2.5 312.11 (3.75) 318.00-19

4) 259.08-14, .22, .23 24 - early outline of the report

5) 208.00 & 169.00 ← v.p. previous review

6) 164.00 first - see 270-28 for criticism, design. I now have ways to deal w. principal deficiencies.

is. non-universality not doing long during L search: not being able to realize all history,

to On(O<sub>n</sub>, Q<sub>n</sub>, A<sub>n</sub>, ...) → On<sub>n+1</sub> model pieces & lots in v. 169.00 model but still,

has no "Lack of head" (ATM<sub>n</sub>)

T. 164.00-169.40 system; w. FC) & [IA<sub>n</sub>] was a reformation of the more "Non-ET" system of 312.11

(On<sub>n+1</sub> = On(O<sub>n</sub>, Q<sub>n</sub>, A<sub>n</sub>, ...)) - So a good idea for a path to take to a more Non-ET system.

Also a Non-ET system can give a Unified Goal to the ET approaches of ET.

243.00 : Core Recur.

~ 150.00 I wrote about "Breakthroughs" since 5-89. There have been many "breakthroughs" in last month: A Summary: 312.11 - .40 may do this! ← I may want to be a bit more detail.

1. Use of 3 input vnc. (or 1 input w. 2 prefix set formalisms) to realize Universal Operator Def. (253.37) → 259.90  
~~My old model~~  $\sum$  on universality

My old model of 245.09 - 17; 253.37 (~~253.21~~) 267.00 ← counter-example showing

2.  $O_n \cdot (O_n, Q_n, A_n + \dots) \rightarrow O_{n+1}$  : Very Genl. form for TM 254.16 - 19 ←  
of discoverability

That they consider + must General Form of Context free.



ID

Q4. Problem 320.22: It seems that there is a tradeoff betw. the

Simplification used in  $O_1$  is a kind of "teaching" in the early TS Q's: First

(subject to some lower kind of meta-licensing) we can always put on the meta-licensing

TS Q's to bring  $O_1$  up to an adequate "initial level". (31)

What is not clear in my mind, ~~was~~ was how I expected the initial TM to operate. I had all of these "indexed" TSQ types, & TM would

learn each kind of sub-corpus separately - as if it were a bunch of indep TMs.

Each index would have a known IA assigned to it.

Then as a new problem, I could have TM learn to assign indices to problems -

using the (index, IA) pairs of meta, & perhaps some new ones.

Some: Make a listing of "Initialization Techniques" - BRs to 320.10: BRs to (5)

1) a) (large?) set of IAs <sup>each TM</sup> w/ assoc. SC & update techniques: Each instance in a SC is indexed by user.

These IA include both MTM & NMTM SC's.

b) TM looks at indices & assoc. IAs (assigned by user) as a QA problem SC.

Indices stock notation: so indices are not a key: TM can place IA for

21 & by Q, unindexed QA. I think this is the FC) idea may earlier QATM models.

to some extent (unless as to how much) context (is equivalently, pc's of context, w/ some context dependence) are shared by all IAs. (I have not worked out details of BR).

2) Did I have any ideas defeat from 1)?

3) Learning during Lsych, context-dependent pc's of cues: How do these ~~work~~ over across relevant to (2)? Well, in 22 was do Lsych to solve (perhaps all) problems.

"Learning" in Lsych is implicit (2) & ways a) Learning from previous trials (+ or -)

b) Answer solving to "correct" problem (no a) & b) may be ideal (concepts!).

T. Context dependence of cues: There is (among other things) very relevant to the sharing of cues betw. different IAs & their assoc SC's.

31 (3) SN I know of a  $O_1$  & it is of 02: It is unusual to machine 320.03 learn using Lsych to find sums. Or, (closer to practicality): a AZI-type machine w. Lisp semantics good universal function (Lsych), - uses AZI for meta assignment of pc's, - BR's should be for fracture models (NMTM) so e.g. 3 input machine (O, Q, R) work for one passy, - probably a better way of expressing PD's & could be est, hrr. (See 310.00-311.40 renew v.g. ideas only).

31# Actually, write be a v.g. way to get a TM off "off. end". I decided on goal set of IAs & then putting TSQ's for each IA. "Adequate" means TM would know when it had found a soln. but KSS would be too by (2 1020, say). We are not over

TM soln. completely "by hand", or give it an upper limit so that CJS is down to "acceptable size". This CJS limit will depend on ~~our~~ cc resources we have available: 10<sup>9</sup> may be ok, but 10<sup>12</sup> or 10<sup>15</sup> if we had bigger machines. Getting TM to do this search means it could find more compact codes than those we know — also, it could ~~be~~ build on codes used in ~~previous~~ solns. of other IA's to which it had found "solns" — which we might not put in our "A. H." solns.

So, if we have a big copy computer available, we can ~~get~~ train TM "better". — In general, we want solns w. as large CJS as we can afford. — Since this makes it more likely that TM will find something better than solns. we had in mind.

In general, when we put in a "A. H." soln, there will be usually much fewer useful <sup>empty</sup> sub-codes than if it soln. were obtained more serendipitously.  
Having smaller computer available, means we have to spend a lot more time seriously constructing our TSC's — so they have lots of generally useful genes. in them.

SN In <sup>ANL</sup> ~~the~~ I had 6. folg. <sup>in mind</sup>:  $3+5=8$  has soln sum 3,5, ~~3,5~~ → 8.  
~~8x9=72~~ has soln mul 8,9. → 72. These situations differ in ~~the~~  $x$  ~~variables~~  $n$   
 $t \rightarrow x \neq \text{sum} \rightarrow \text{mul}$ ; so make ~~the~~ a function  $t \rightarrow \text{sum}$   
 $x \rightarrow \text{mul}$ .

~~The~~ TM would be unlikely to be able to do that kind of processing at this point of its evolution: It would need some ~~form~~ ability to deal w. "logic of situations" — i.e. it has had at this point, no ~~idea~~ TSC ~~or~~ a suitable IA ~~to~~ deal w. such problems.

SN Confusion in my mind about what a "program" is, what a "soln" is; what a "SC" is, i.e. whether IA is: 326.00 - 17 (on globe = Dynamic IA's) seems to solve this problem.  
Problems on Q; Soln: Any AI method that VR decides  $\rightarrow$  a "soln": This can be very ambiguous to VR.  
e.g. say there are several equally good solns: As a TSC, it would be well to include all Q; A; <sup>previous TSC</sup>  $j = k$  if ~~there are~~ a good solns. Even if 5 solns are not all exactly as good but within ~~say~~  $\geq 5\%$  — it would be well to present <sup>all</sup> solns as a part of to corpus "BAG".

An "SC" is a <sup>list</sup> ~~subset~~ ~~of~~ ~~the~~ corpus, that all have  $\gg$  an index on their Q. Since a Q can have several indices, it can be in several SC's. Perhaps have indices ordered (linear or partial) so there's a <sup>max</sup> ~~max~~ (top) index. T. Main index determines ~~the~~ SC; How other indices are used, how SC's can ~~share~~ ~~some~~ into or pc's of codes, is unclear.

Perhaps certain problems could have  $\gg$  (main index)  $\rightarrow$  so they are in  $\gg$  SC.  
So: <sup>if not using above/and w. reference to</sup> indexing is ~~the~~ problem of how to pool data from SC's ~~to~~ — unclear.

8/11/01

ID

discrimin  
322.40?

important! 331... not written  
did I mean 231.04-08?  
it almost fits!

322.01! Look at 331.04-08: Here, the SC's were like independent TM's (at first). This was just for update to individual IA's. Later, TM became categorization of New Q's into SC's is somehow different appropriate updating.

Say we've done this preliminary sup. of 0.01: Also, we have TM able to give probabilistic assignment out of Q's to SC's (if Q's are unmerged).

At a less abt. level,  $F()$  assigns naturally pc to the IAs (wrt a given  $Q_i$ ), but interprets / various weights & products of  $\epsilon$  variables SC's used Not shown to do. So number, figures out how to update the SC's w.r.t.

This new QA ← this last is diff. This one then by which it does may take influence to wts in

At a recap non-abt. level:  $F()$  has complete control (over prod. & updating).

It controls also context dependence of pc's of codes used in constructing new anal solns to problems. Essentially  $F()$  is mainly to "point" 0 that works on updating — is because of but it must be involved in prod. as well (since updating ultimately depends on prod.).

Empirical for a while! More on Confusions in 322.24; on problem, SC & IA:

In the simplest case, an IA works a problem, Qu. Having done so, the QA pair is added to its SC. Therefore updating of the IA also occurs to SC. is unmerged. Certain patterns are changed & bid, but also, new ones can be found, because of the copy

Assessment system. (In AZI there is no "repeating" — but I feel uneasy about this)

It was mainly dealing about Algebra long as in SAARB TSC's.

Suppose TM "travels" to solve the 20 by linear eqs. If we give it lots of new linear eqs it solves them, 20, in some sense, its skill in Algebra is increasing — but what is mainly getting more reliable is

our confidence in its ability to do a certain subset of problems — namely linear eqs.

Its pc's for solving a new kind of problems doesn't change much. — some continuous parameters (pc's of codes) will change slightly.

Actually, since the range is for a MTM corpus, the values of pc's of codes will not change at all — as long as the linear eqs are all solved! T. of eqs for does not change if the corpus is in size, as long as the rate always keep existing in.

In General for MTM problems, only one SC is necessary. If we had different SC's and different operators, one would find a Q discriminant function is merged & operators!

Is an analogous way poss. for NMTM probs? Say we have 2 distinct SC's! S1, S2

S1 has a pc of P1 finite corpus & its pc discriminant is P01 so total pc is P01 · P1  
Similarly for S2 ... P02 · P2

If we had a discriminant function between 2 SC's: w. same pc = Pdisc, then we could decide a new S2 that switch on S1 for its corpus & 1/2 for its corpus; total pc would be

Pdisc · P1 · P2 · P01 · P02; It might be possible to find Pdisc · P01 · P02 by using common codes or operators

Alternatively reducing, so this is a lower bound for the pc. This would be difficult better than using either O1 or O2 to code & evaluate corpus, if O1 & O2 were much better on S1 & S2 respectively than S2 & S1 resp. is useful joint discriminant (O1, O2 is Discriminant function) Didn't have a too low pc. This assumes infallibility of the Discr. function.

If we give the discr function wts w1 & w2 for the 2 Q's O1 or w2 for the (depending on which way it decides), then we end up w. 2 mixed prodn. for each Q in the corpus — static vs. dynamic IA's

if we use anyway we find this better than P1P01 or P2P02.

326.00  
Spec

I had in this area  
"unnecessary"  
This seems to be the way to find the feeling of

8/11/01

# REPARSING → .20 ← This is a Big change of way of looking at PEGs!

00:323.40: So, say we did find a suitable discrim funct so to make a parser was better: How do we  
from updates  $O_1, O_2, S_1, S_2$  (perhaps disjunctive func.?)

Also, is it o.k. if  $S_1$  &  $S_2$  have common elements? BARBARA  
Using  $W_{alt}$  as in 323.38 ~~is~~ expansion to compute pc of entire corpus because  
 $O_1$  has to compute all of  $S_2$ 's pc's &  $O_2$  has to compute all of  $S_1$ 's pc's:

Anyway would be to use a ~~discrim~~  $O_1$  discrim. function - It would say "use  $O_1$  or  
"use  $O_2$ ": If it were wrong occasionally,  $O_1$  would have to ~~do a problem~~  $(S_1)$ .

008: [A] If the Discrim funct. was kept constant (w/out history of TM), it was used to define  $S_1$  &  $S_2$ :  
it would make no errors in carrying out problem, to  $O_1$  &  $O_2$  (C)

00: Another approach that I may have used in past! For each problem, ~~state~~ let ~~the~~ least recently history  
of TM), most problems are worked by more than one IA, so would worry about .03-.04.  
A good advantage is that it gives us a better idea as to how good the various IA's are on various  
problem areas. I think FC) was the discrim. function. - But still, how to do updating?  
How bad is it to update a sc/IA w. a problem that is sc/IA isn't particularly good at?

Well, it would seem that we'd like a ~~sc~~ sc to be as "pure" as possible, so it didn't ponder its own  
"operator down bits" trying to "fit" a problem in which the operator was not used.

Actually, we could "change" a OP/IA only for problems in which FC) chose them as best.  
So this would clearly define  $S_1$  &  $S_2$  (like in .08).

20: [SN] I had the idea that Reversing would not be a problem this model of TM! Maybe wrong!  
Say we put a finer new line. Somehow, it enables us to reparse the same old  $O_n$  much higher pc:

Also ~~now~~ several cases may have their pc's down a lot. (But undoing an old case can be very tricky -  
maybe not do it, unless it is easy, it inches down "waters" (C). By reparsing & I mean,  
~~then~~ in most extreme case, breaking down entire  $O_n$  into primitive functions/operators,

Then reparsing intervals of latest delays & their contents (now first approximation) pc's.  
Even the undoing of "O" operator can be done very efficiently, or to any depth.  
Usually its done not very heavily. After reparsing of much depth, the pc's change, so it may have to be  
repeatedly reparsed to get the work pc.

30: 310.29 ON NMTM operations  
I think the idea of 310.12-.30 was: To start out with a "user given" set of IAS: to use them  
on the corpus to "calibrate" them ~~then~~ so it'd have a pc of each IA as an "as" in a stack grammar.  
Then make a stack grammar in which each IA is an "as" in the grammar. This is a level grammar,  
Let the IAS's can now be regarded as an ordinary stack grammar w. the PD's on Q's as its terminals (leaves).  
As such, the PEGs we have to modify/optz it are much broader, than in its initial calibrated form  
of "Grammar w. Meta/rewriter"

310.12-.30 was first Genz. of FC) from very simple MetaGrammar (just give pc's to  
IAS's) to a stronger Meta Grammar that also generated new IAS's (≡ stack grams). For the Genz. this is not much of a genz.  
so it need not have "subgrammars" (the most grammars do have subgrammars) - or that it had subgrammars additive if at all a genz. its a  
& perhaps much of pc's of other forms of the "old" IA grammars. simple realization of what's going on!  
In addition, I wanted this Meta Grammar to be completely dependent & not only on the input problem, Q, 325.00

00: 324.40: but on the entire deriv of the present On is its structure & ~~its~~ internal PCs of its cons, etc.

01 Using such a model, the Question of what is to be updated? what SC's to use, etc. — by consequence —  
(perhaps irrelevant!)

Superficially, it would seem that a very non-al version of TM would be easiest to update —  
conceptually consistent, that is, because the Gores is (somewhat) clearer. This is least E.I. TM

02 the ~~idea~~  $O_n(O_n, (QA)_{n+1}, context) \rightarrow O_{n+1}$ : T. Gore for any  $O_i$  is clearly defined (if we ~~don't~~ <sup>positively</sup> mention of  $CB$  invalid!) — its "max procedure x pc of  $O_i$  defn" ← eq. of: (289.33)

[A possible additional complexity is that of "procedural" involving many I.I. defns.]

08 should make it possible to evaluate any TM models. "updating" a IA, e.g. would mean modifying the params of that IA so that the  $\langle$  Gores of the auxiliary system  $\rangle$  of 08 is max.

We deal w. optimal continuous params, viz the  $EXP = \#(Det(Hessian))^{-1}$  approach of

12  $\frac{1}{2} (250.09 - 19) \leftarrow$  this is on the  $\int_0^1 \dots \int_0^1 p_1^{a_1} \dots p_n^{a_n} w \sum p_i = 1$ , but is a special case of the more

15 general function that approximates. We use AZI for discrete params & the Hessian for continuous params.

So, the Q is: Is 08-15 an Adequate Framework for the Quantitative of Optimization of the Evolutionary Sequence of TM Models from 310.12 to 311.0? 324.30 - 325.01 is a start

Summary of the seq. of Models from 310.12 to 311.0 (perhaps!)

20 Well, 08-15 seems reasonable. A possible Q is: The ~~more~~ simplified Gores may be too hard to evaluate — so I'd like a procedure or Method for Models to facilitate easier Gores Evaln.

TL;DR  
Applying the fuzzy ideas to the early deriv. of CPM:

Say we have this  $F(x)$  & set of I.A's:  $F(x)$  looks at  $Q_{n+1}$  gives a P.D. over it.

I.A's: which in turn give a P.D. over the  $[I.A's]$ .  $F(x)$  &  $[I.A's]$  constitute a mapping from  $O_{n+1}$  in a defn. of  $[A_{n+1}]$ . So this  $\rightarrow$  is a "O<sub>n</sub>". As such, we can use

the eq. of 289.33 to optimize params (both continuous & discrete) — which is what "update" is.

As noted before, each IA will have its own update techniques. They will usually correspond to have elements in them corresponding to elements in AZI: ~~the~~ ~~idea~~ ~~of~~ ~~continuous~~ ~~vs.~~

1) update contin params. 2) reparse 3) find new functions, concepts: ~~how~~ ~~to~~ ~~use~~ ~~it~~ which is how much of each depends on CB available.

In updating, ideally,  $\{F(x)$  & all of the I.A's params  $\}$  are simultaneously updated. In practice, each  $Q_{n+1}$  will usually involve only a few params in a few I.A's & perhaps in  $F(x)$ .

[Also Note 326.00 on "STATIC V.S. DYNAMIC I.A's"]



ID

STATIC & DYNAMIC IA'S

323.90 spec

325.90

on 322.34; 323.11: In general, an IA is a program which supports an input Q<sub>i</sub>

to one part of a P.O. or Answers [A<sub>n</sub>]. These seem to be 2 kinds of IA's:

One kind is static, & never updated. Another Dynamic changes as its SC & rows. — i.e. its I/O behavior changes.

I think Quizes basically sort of Q<sub>i</sub>. I was worried about 322.29 & 323.11 &.

So, say both kind of IA's were in a IA pool [IA<sub>i</sub>]. They are both used by F( ) to solve problems.

A static IA could be a component of a dynamic IA: E.G. O, t. over all Operator, is a Dynamic IA, but it can, via F( ) use static IA's when being seem appropriate.

A IA producing a time series would vary its program (both central & discrete) as its corpus grows. A correlating Alg. could be static — unchanging as its corpus grows.

Somehow, f. updating Alg. must recognize what needs to be updated & what doesn't.

Well, f. Master Updating Model is 325.08-15. f. current O is modified to make 325.08 ≈ (289.33). In this Mod. procedure, certain IA's will be modified & others will ~~not~~ perhaps never be modified.

More examples of static/dynamic IA's: static: how to solve quadratic dynamic: how to solve n.l. eqns. ; Somewhat static is solving linear eqns. — but IA's could start out dynamic — then settle down to static.

The Master Update model of 325.08-15 is supposed to take care of the "SC" problem. What SC's are assoc. w. what Dynamic IA's, etc.] To detail mechanics of this over as you need to me.

How, I think I have now have more over out to write an intelligible report. In my own experience

hard copy is swif<sup>t</sup> copy: have refs to my notes, expanding various points in report.

Some of these expansions should be in SW & hard copy also.

The preliminary outline of 320.00 looks fine: 320.22 (?) needs to be expanded which is what 821.00-826.19 is about.

Descriptions of "Installation of CPM": ~~static~~

1) Set of IA's: F( ); Index Q's: so F( ) is not much needed to assign Q's to IA's.

Also ~~the~~ Each IA acts separately on its index SC.

How do IA's work? I have idea of how IA for MIT work. — its AZ1. — at first, w/ bound. pc's of <sup>(memory)</sup> Q's.

Try to find examples of ~~static~~ <sup>Static</sup> IA's: ~~Static~~ is Dynamic.

In more advanced machine: Given set of ~~word~~ Q's <sup>memory</sup> indexed, assoc A<sub>i</sub> are

f. indexes given: To induce index for a new (or old) Q<sub>i</sub>. — it gives PD (usually!).

Some static IA's! (nonlinear) correlating, ANU.

Present ~~also~~ <sup>also</sup> used as Dynamic IA's: e.g. in linear regression, CPM learns

f. sequence & improves its refs as no. of ~~obs.~~ <sup>obs.</sup> as it gets more data.

Updating is static. ~~an~~ adjustment of parameters w/o consideration of "new terms", new functional forms.

Analy. Predn. of Best seq. is static or dynamic.

ID

# REV <sup>327</sup> .19 → 328.23

00 : : How is some distance betw. stock Op. ( $O_n$  is a stock op), & P.D. ( $\frac{1}{2\pi i} \oint e^{-\frac{xz}{2\pi i}}$ )

→ "PD" & IA's: I guess a IA has a corpus of examples as input & has a pd as output. They differ as to what corpus elements are.  $O_n$  is not an IA, but  $O_n$  is a stock operator.

CPM is an IA. In CPM, corpus is  $\{Q_i A_i\}_{i=1}^n$ . In ~~the~~ regression, (linear or Non-linear) A stock op. is a degenerate IA: It ignores all previous data corpus data.

T. ~~for~~ input is a time series. | It is a stock IA.  
A IA can be for sequential data (Time Series) or for unordered data (e.g. stock buy for simple pigs; stock operator for any corpus of QA pairs)

3 common types of IA input: ① Time series (Numbers - Nos. are Real or integer) or complex (Integers, real, complex...)  
② Bag of objects ③ Pair of I/O pairs.  
The objects can be strings & / or ~~sets of integers~~.

T. system presents a new  $(Q)$  to a IA for (predn) (update): It already has the previous part of the assoc corpus.

T. IA's of interest must be able to deal w. a  $\sum Q_i A_i$  / corpus, ~~type~~ series.  
Hvr, a IA that is a time series predictor, can regard to  $Q_i A_i$ 's as a elements or a Time Series.  
Then give a  $Q_{n+1}$  it can give D.P. for  $A_{n+1}$ . - Hvr, it may be losing the part  
t. data is "unordered".

Anyways: As present,  $FC$  (when executed properly), looks at  $Q$  & sends  $Q$  (or  $QA$ ) to certain IAs. Or better,  $FC$  outputs a PD on IA's - pd that each will "best" (i.e. "best" means best pc for rate  $A_i$ ). So, ideally,  $FC$  obtains a PD over

t. IA's for a given  $Q$ , → .29  
Woops! In .19-.21; t. pc assoc w. each IA will depend on how  $Q$  is used. If a predn is via t. IA, assigned that pc for that  $Q$ , then t. pc assoc w. each IA is t. pc that will ~~be~~ assigned that pc to  $A_i$ .

Hvr, if predn are t. pc w/  $\sum$  of outputs of t. IA's, then these pc's will be adjusted to optimize predn outputs using Present's.

Which way to do it is something I don't like to decide now; See 337.23-32 Another form of the "2 ways of going for  $FC$ " problem.

So to continue. Even if  $FC$  were v.g. This system would be (if I know, then using a bunch of IAs in a decision which (or using a update) to give predn: Update of IA's would be share of any new independent IAs.

A first improvement over .19-.31; (all the resultant IA; On latter updating term  $(QA)$ )

Use  $O_n$  to do context dependent "selection" of concs. in updating t. IA's. This means that in Universal (or very extensive) IA's finding good updates, finding the new corpus for update would occur much faster than w/o this feature. Also such is rather separate & uses long during search to deal w. context. To do sequential long during search.

Another improvement is in "long during" search: so after I remove fails, t. PD vector such changes. It also helps deal w. corrections between search trials.

ID

Rev. : 328.24 - 330.40

00. : 327.39 - 40 ~~is important in helping realize~~ **589 Anz** - We should check to  
see that all kinds of hours are, indeed, accessible to this system

02 A perhaps Novel method of Doing LSrch is: T. input to CPM is T. Problem being "LSrch".  
Also inputs are trials & there results up to now. - Output is needed what next trial is.  
When (in 02) CPM is "cold", ~~it also~~ CPM is also given a time limit telling how much  
it to spend on this problem.

07 A serious Criticism of all opt. algo. is that there is no "lockhead": i.e. searches are  
all grossly "Greedy" (= "myopic"). How we could teach CPM how to solve "lockhead"  
problems & then it could apply them to "Aided Lsrch". - At the ~~least~~ level of  
level, all Lsrches are controlled directly by CPM: It looks at the problem, it previous trial  
results & devises the next trial.

One difficulty of present Systems I don't seem to have any good practical ways to do  
Stochastic. funcs for stochastic IAs:

There is a 3 input line, but I find it difficult to put intuitive problem ideas into it.

On **228000** offer **16 IA's**: It's not clear how many of them are stochastic! -  
It seems that at least 10 of them are, indeed stochastic.

Now, I shouldn't worry about this now. ~~But~~ In designing a TSCQ, if I know of a problem  
soln. to a problem, I will design a IA that has Access to that soln.

In particular, various Heur will be of probabilistic form & I'll have to design a IA (or IA's) to access them.  
Now, 310.00 - 311.40 to have some Guidelines on developing a good, useful set of Stoch IA's.

TSD's for them, etc.

So: Present Summary of initial System!

1) We start w. a set of IA's! Both M & NM. (M). An IA takes a ~~set of~~  $\{Q_i, A_i\}$  as input <sup>is able to</sup> ~~is able to~~ <sup>compute</sup>  $G_n$  &  $Q_n$  is a Pd. on  $A_i$ .  
A Defn of an IA includes its (incremental) update algo.

2) USR gives System a set of  $Q_i, A_i$  pairs! Each  $Q_i$  is indexed, telling which  
(or maybe  $\times 1$ ) IA is to work on the problem. These IA's incrementally update  
themselves on the  $Q_i$ 's they are given. T.S. IA's at any pt. are indep. of one another.

3)  $F()$  is created: <sup>It is created by:  $F()$</sup>  ~~It uses~~ <sup>uses</sup> a set of IA's (given by trainer) to look at the ~~indicators~~  
 $Q_i$  & IA assignments by  $A_i$ 's.  $F()$  then is the probabilistic relation found between  
the associated indices, & it finds a probabilistic relationship between them.  
From this relation,  $F()$  is inductively able to assign indices ( $\rightarrow$  IA's)

to now, unindexed  $Q_i$ 's.  
4) The combination of  $\{F()\}$  & a set of IA's is a new IA  $\equiv O$  ~~it is~~ when a new  
 $Q_i$  comes in,  $P_i$  is updated by updating the relevant IAs & also updating the  
 $F()$  function. - But be for that, note that  $F()$  was probabilistic outputs, so ~~unstable~~

**ISN** Here  
I've been ignoring the utility  
of common  
concs. nets  
Various IA's  
I've saved IA  
using to solve game.  
This Needs  
work.

00:228.90 : for a new, unindexed Q, it doesn't pretrain IA exactly — But we can apply

th. IAs of least 2 or 3 pc to a problem & update them on t. ~~IA~~ to a new Q, A.

Predictions can be made using wtd preds of T. IAs! But see 327.19-28 : There are 26 least 2 ways that FC can assign pc (wts) for IAs (for a given Q input).

I don't think it's a critical decision, but it has to be made.

With FC in place & suitable training, we ~~can~~ <sup>cannot</sup> work on 2 critical SI

Problems : S) & G) [ But note 330.16 : It may be poss. to do S) & G) earlier in the TSCQ than this ]

89 4-1-01 161 5/1/01  
 329 8-14-01  
 S=168 103+2=  
 105  
 1.6pp/d exactly since

5) (Context dependence of pc's of concs.)

Updating of IAs typically can while IAs will differ somewhat in their updating Algs,

these 2 plans will, typically, have 3 aspects. As new data is added to the SC, IAs

IAs, therefore, context (continuous parameters are changed. These are usually frequencies of occurrence of various elements in the SC.

a) New concepts are found in the SC. In the early stages of a machine's training, various concs will be assigned pcs based on their frequency of occurrence in the SC.

The changing of these freqs will be recorded in the updating process.

b) These concs are combined to form higher order concs, which have to be evaluated by the system core. The frequencies of a) are used to assign 2pp's to these concs before evaluation.

c) When new useful concs are found, the prediction program or the IA will be retrained in terms of them. Using them in 2 attempts to the 2pp of the current IA's program. (It's kinda effectively "shorter term" of this program)

As the amount of data in the machine grows, it becomes possible to

calculate the 2pp of a new concept, not only by the frequencies of past occurrence of its component sub-concepts — but by the contexts in which these sub-concepts occurred. This context dependent probability can be done by the grand

CPM operator if the system has had adequate training experience in this area.

This context discovery of context dependency of pc's of concs is critical in discovering new concs in larger systems. Without it, the probability assigned to new concs decreases rapidly with system growth and rapidly

in context to cjs of new problem beyond our computation capacity. But whether it ever gets completely out of hand or not, a system without context dependency of pc's is grossly inefficient in computation cost.

6) In the process of searching for new concs, part of the updating plans of most IAs, the search is normally guided by a pc that remains constant during the search. The constancy of this pc makes it impossible for the

omit

20

30

00:329.40 Search to take advantage of information gained in earlier trials. By revealing the guiding PD after each trial, or after every few trials, the efficiency of the search can probably be markedly increased. This modification of the guiding PD corresponds to "Learning during search". It also deals with an equivalent inefficiency of LS search - that if trials are not statistically indep, but highly correlated, the search becomes very inefficient. "Learning during search" ~~removes the~~ <sup>appears to be a good</sup> ~~is a~~ way to deal with the correlation. This "learning during search" may require training beyond that needed to construct the original static, guiding PD.

10 ~~Notes~~ <sup>while</sup> write re reading 328.24 - 330.02: (328.30R) : 1) ~~At~~ At the start of (try off) IA's it is in dip. It would seem that common cases in the AIE would be para-effective PC's. - but fitz needs work! As is, the system will probably work "adequately" w.o. fitz feedback. 2) In (4) 328.38ff: T. FC) mixing algo: ~~I~~ I have worked out (in these notes) some analysis of optimum n.l. mixing. HVR-2 linear analysis of linear mixing very cross corr, might be easier/faster or a good "first approx" for the N-l ~~mixing~~ option.

06 3) In 328.24: Initial operation: The sophisticated terms of SX (context dependent PC's for component cines in updating ~~is~~ G) (long during LSrch) could be done to some extent: The assignment of suitable IA's for each of these problems may have to be done by USB, HVR in ~~early~~ very early TSQ's.

20 4) The CC's are considered LSrch, there is no mention of CC in the foreg. discussion. I have HVR written framework for a general option of the system w/o. CC, see 303.22 - .29 for preliminary design: 285.00 - .40 for immediate!, but this needs much work

23 5) Note 303.08 - .21, .30 - .40 - on optimality (within factor of 2) of the system. This amounts to an over-all methodology for optimization of the system - i.e. suit is capable of ~~see~~ realizing/simulating any desirable (human) ~~search~~ search heuristic (fitz is ~~mainly~~ mainly PC oriented heuristics - the problem of "which about" seems different... I may be able to get CPM to work on this, if I can find a general key to describe over all CC problems to CPM!

24 6) Be sure to mention to over-all hear: that any problem describable in human can probably be usefully worked on by CPM. - So if there are any deficiencies in my model of CPM, I can probably get CPM to try to deal w. them.

30 7) 327.19 - 328.23 is the previous review (with a preface)! Has never Comments, Criticism!

8) T. system as described has no lookahead: it's maximum myopia ("overly"): This could be done via (2.9) also see 328.07. But no more

Def Operator → Stock Operator .26

M Op → M Op

Summary of imp't ideas in System 1: (i.e. "Claims" in a patent Application.)

1. Use of set of I A's (i.e. New Min & single IA) to solve a great variety of problem types

2. Suggested list of I A's.

3. Use of Level to help update I A's.

4. Training of a FC) from data to ~~decide~~ decide which I A's should work on a problem & how much wt. to give each of the I A's.

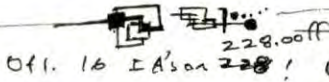
4. Imp't. methods of SI of system:

a) context dependant pc's of concs

b) (run during Level).

5) Overall Goal of System.

6) Final Goal of  $O_{n+1} = O_n(O_n, R_{n+1}, A_{n+1}, \dots)$  using Goal of 5)



Off. 16 I A's on 228! ⑥ ANN, RANN } These are not actually I A's: they are Optim. Methods.

⑦ GA, SGA, Simulated.

Their Optim. aspects can be used as part of a gradn. scheme! e.g. They can be used to implement (or help implement) the "updating" of an I A. Or they can be used directly to find "shortcuts" for a corpus - which amounts to a very general method of Induction.

In particular, ANN & RANN can be used to desc. Operators (deterministic operators).

They can also be used to desc. Stochastic operators, by giving a central prodn. w. & no observed op.

In a similar way GA (etal) can be used to hunt for deterministic Operators (M I G) as above.

Both (ANN) & (SGA) can also be used to optimize desc. of a corpus in terms of a M for every prodn.

It is possible diff. of 2 for each prodn. [usually, the prodn. system gives only 1 & we mess up

a mean of 2 for whole corpus; from which we induce a best of 2 fit for the model.]

cpm. cond. prob. prodn.

Way to take a M Operator (string → string) & rep. it into a M Operator:

M Op is observed to be correct say 80% of the time. When it is correct, its pc = .8.

When it is incorrect its pc = .2 mult by uncondl pd overall other vcs, probab. — This can be

even ~~more~~ viz Pcc R input. Its necy to do this because otherwise if

the operator is wrong, it would give zero pc fut. into answer, — which makes that the

Score would give zero score! We can contrast to Score obtained if the system gave the

same (Paradisi) d.f. for each Q — This saves D.f. being 1. Uncondl P.D. → (340, 31)

26 It is perfectly convenient way to convert any M Op to a M Op

So:  $h_{n+1}/r_{n+1} = -.8 \times (h_n + r_n) + .2 (h_n + r_n)$   $n =$  probabilistic mean pc of

(from a default (unconditional) d.f.

W. Point + def. predictor, ~~mean~~ uncond/case =  $h_n$

M Op to a M Op  
D Op  
deterministic  
S Op  
stochastic

$$= -.8 h_n + .2 h_n - .2 h_n + .2 h_n$$

$$= h_n (.8 + .2) = .500 - .2 h_n$$

$$= h_n$$

$$= h_n$$

Looks Good!

32



goal for  
00: 332.40: The condition that we are ~~trying to~~ constructing  $O_n$  is that

01 
$$P(O_n) \cdot \prod_{i=1}^n P_{O_n}(A_i | Q_i) \quad (2) \leftarrow \text{(eq. no.)}$$

be maximized  $\rightarrow$  a maximum likelihood criterion.

Here  $P(O_n)$  is the a priori probability of  $O_n$ . It is approximately  $2^{-\text{length}(O_n)}$  in bits.

$P_{O_n}(A_i | Q_i)$  is the probability  $O_n$  assigns to  $A_i$  for input  $Q_n$ .  
In general, it is not computationally possible to maximize expression (1)  $\rightarrow$  we will try to do as well as possible with the available computational cost (usually time).

10 We do not know of a way to obtain a  $O_n$  operator that satisfies both (1) and (2) in the direct manner.

13 Instead we will describe a sequence of operators of increasing power, that will ultimately culminate in  $O_n$ . 347.00

Some ideas on the next section!

1)  $F()$  could optimize  $\prod P_{O_n}(A_i | Q_i)$  as soon as the indexing steps — I don't see any advantage to trying to duplicate the indices — (two perhaps, completely, it would be a nuisance) — the no. of cases is large  $\rightarrow$  many IA evaluations would have to be done. look into this! — from cc pt. of view. Note that after  $F()$  has obtained a functional form for  $F(Q_i) \rightarrow I_j$ , the updating of  $F()$  is then incremental & doesn't involve a search IA evaln.

2) GA's  $\rightarrow$  IA's. When  $F()$  chooses a GA to work on  $Q_i$ ,  $F()$  should specify the GA as narrowly as possible, by giving population size, initial population,  $\rightarrow$  the most important characteristics of used GA. nature of mutations & crossovers. What this means is that to use GA's effectively,

29  $F()$  would have to study their behavior in detail. Essential. Also,  $F()$  should be able to watch the GA at work & suggest modifications of mutations  $\rightarrow$  input!  $\rightarrow$  344.16 specification of GA's as IAS.  $\rightarrow$  334.00

30 [SN] On punctuation ① It wouldn't be needed if I used prefix codes. (The Professor expensive I should look into that) ② in any  $i$  the approx; a lookup normalization — look it up in Ripp's Papers ③ wants to do limits a few long regular chunks of code, I could use prefix coded integers to tell how many bits (or symbols) were in each section. ④ Could modifying punctuation get rid of the "1/e" factor in Z? — I doubt it, but look into this! ⑤ In Huffman ~~code~~ language comp: frequency of symbols dependent corpus sizes so  $p \approx \frac{\text{count of symbol}}{\sum \text{total no. of symbols}}$ ; But look at that Book for details! ⑥ In some situations (like part of Lush) the system is self-normalizing; so using  $p \approx$  it would be ok.





ID:ia

00:339.90 : functions that solve new probs, random solns of old problems is held by PC (according to AZI), is do this via acceptable cc. H. final score is then funct of 333.01 being  $\max$ . AZI gives the  $P(O_n)$  factor. - We hope to keep PC "I" part constant except for the factors  $Q_{n+1}, A_{n+1}$

$P(O_{n+1} | Q_{n+1})$ .

We have to reflect AZI is part of it. Gave: Also its role as a first approach of ~~probabilities~~ conc. pc's to be used in trial  $\rightarrow$  soln. of problem).

An imp't point is that the AZI pc's are completely or almost completely indep of the problem  $Q_{n+1}, A_{n+1}$ . The over parts because they've been used to solve previous problems is  $Q_{n+1}, A_{n+1}$  is (usually) referred to factor

Factorial problems.

"To solve a problem" - our info worth behaviors from  $(Q,A)_{n+1}$  - A technique not suggested by

AZI's pc values.

The Problems Discussed 334.05 (to present) is of other imp't. While the discn. of 334.37 to present clarifies a lot, there is still a lot of things unclear - in particular 334.05-10 - on "universality".

334.05-39 talks about universality, i.e. idea that if code that AZI has for its output may have redundancies in it - that is, AZI is unable to "see" & then reduce those redundancies

32431-335.11 discusses the practical problem of finding a min code for the corpus.

conjecture: That removing the redundancies (if any) in AZI's code for the corpus, will not <sup>improve</sup> TM's effectiveness much. However, finding ways to rapidly construct operators that ~~work~~ give by score via <sup>GORC of</sup> (eg 333.01) is much more imp't. Using "Auxiliary CA's" to help solve PC's problem is apparently one good way.

The idea of "Context" is "any info that will help solve a problem of 333.01". If PC includes success ~~has~~ info on what happened in earlier trials - we have "long duration & rich" included in "Contextual considerations".

300.00 ff (300.26 ff in particular) expands idea of "Context" in very useful ways.

8/17/01 This idea of "Context" seems very imp't, even in Early TM (contact, - short TSQ).

I think I had no idea that usually & inform. app'd - PC was used for Gave Evaln - was "too small" to account for the amount of "intelligence" used in problem solving!

I.E. The Function,  $O$ , isn't ordinarily large so that the pc's of  $O$  sub-functions used, give <sup>sufficient CJS</sup> (a sharp end d.f.) to make CJS practical for large  $O$ . The SSZ within

$O$  is too small to give us able CJS. The guidance of the LSM has no use at all info: Some of it has to be (from a much larger auxiliary corpus): Some of it from the present corpus... ~~is~~ "context". However, I do have to walk out just

What the corpus is SSZ for PC's "Context" is.

the SSZ for the <sup>constant</sup> app'd is indep of the no. of examples of use of  $O$  (which could be quite large).

I think it may be very for me to actually work on a TSQ under the good understanding of "Context" & its SSZ.

ID

00:33:40: One intuitive "idea of" context", was that it was a narrowing of the corpus sub-corpus  
 So that with small s.c., the pc's of concs. were much different! That there was a small set of concs  
 w. by pc & t. case had very few pc. That in this "context" ~~the~~ s.c. One needs more or less known  
 which would set of concs to use — That in a sense, the context  $\leftarrow$  (i.e. the narrowed concs. implied) known  
 It can be that Heavis always refer to the context, t. context decided which hour to use.

Version  
 1.414 (19-8)  
 build  
 15981002  
 Win 95/98/NT  
 (order or structure  
 Area  
 209,6,203,139)  
 Order no  
 426772

This whole Question seems to be one of the most critical in my Theoretical Understanding of TM's  
 Operation!

Try to Frame the Question More Exactly: The scope of the concs within O: the whole s.c. is completely  
 indep of t. corpus. It depend only on the "Language" AZ1. (There are benzene rings in AZ1's code for  
 O: i.e. those that just the rules of the primitive & defined functions. But this is not "structure" that  
 I'm concerned w.: 300.26 is more what I mean by "context". Note 301,36  
 Actually, all of 300.00-303.21 is very much in "context": First part of my stuff (302.00-303.21) is on  
 is on "Correlations" & "range during Lsuch", which is paraps to some kind of reference to such.

HA! A nice way to Point out this! Say  $O_n$  was actually the result of  $CC=00$ . Actually ALP:  
 which is even better!  $O_n$  is Plan a "Summary Machine" for the  $(Q_n, A_n) \in \{M\}$  Corpus, - w. z "built-in ssc"  
 We apply  $O_n$  to the "Corpus" Best we're using to compare trials for  $O_{n+1}$ . If  $O_n$  were  $CC=00$   
 Summary Machinery,  $O_{n+1}$  Summary Machine  $O_{n+1}$  needs as info, only  $O_n$  &  $O_{n+1}$  Anal.  $\Sigma$   
 So, superficially, it would seem that using super  $O_{n+1}$  (like that of 300.26 - .35)  
 would amount to "Double Counting".

It's  
 given  
 $O_n, Q_n, A_n$   
 I must say  
 one can  
 find  
 $O_{n+1}$  exactly  
 Sec. 22

A poss. way out is that  $RLP \neq ALP$ . - But finite CC makes a lot of difference!  
 On second Prob, it isn't really "double counting", because the final score, 289.33 doesn't "double count".

While (.24R)  $O_n$  is  $Q_n, A_n$  appears to be not a number to give  $O_{n+1}$ ,  
 out, or that, if  $O_n$  were given as is ~~is not~~ all past.  $O_{n+1}$ 's w. is not a block, w.

(so  $O_n = \sum w_i O_{n-1}(i)$ ) Here  $O_{n+1}(i) = \sum w_j O_{n-1}(A_n/A_n) \cdot O_{n-1}(i)$ . I recall

In fact, for ~~the~~  $\rightarrow$  one would not have to buy over only s.  $O_{n+1}$ 's but as much w.  $\rightarrow$  352.10

Say,  $\Rightarrow$  "context" TM has all traces of its behavior ~~the~~ since start: ~~the~~ In trying  
 to prediction for  $O_{n+1}$  in terms of  $O_n, A_n$  would depend much on how "good" TM was in predict  
 during time  $O$  to  $n$ . - if it was poor,  $\Sigma$  to  $O_{n+1}$  distribution would be poor. At best, this induction  
 with speed up finding reasonable trials for  $O_{n+1}$  - so w. given  $C, B$ , TM would be better,

Remember In this particular "SI" mode, TM doesn't really know what the problem is -  
 its mainly trying to predict what  $O_{n+1}$  would be on the basis of past results.  
 It would be much better if TM really understood what the problem was. So it would  
 work it as a OZ problem. I don't particularly like to see the Lsuch solution  
 of the OZ problem: It mainly consists of looking for a good OZ problem Algor.  
 Perhaps use GA for that aspect of the problem! It would be much better if the  
 GA were more tightly linked to  $O_n$ 's normal operations of Operator Induction.



8/19/01

Am I a better copy of this at:

[www.seasrui.org/KDD/Workshops/IJCAI-2001/](http://www.seasrui.org/KDD/Workshops/IJCAI-2001/)

orig. Copy in ss\wrappers.txt.

For Moments this:

Google!

Wrappers KDD

A.

IJCAI-01 on : out Sat 4 Aug 2001

IJCAI-01 on by W. H. Hsu

Wrappers for Performance Enhancement in Discovery in Databases (KDD)

[workshop code ML-5]

Saturday, 04 August 2001

Seattle, Washington, USA

Workshop Description

(Last updated 05 Dec 2000)

The rapidly increasing volume of data collected for decision applications in commercial, industrial, medical, and defense domains has made it a challenge to scale up knowledge discovery in (KDD), the machine learning and knowledge acquisition of these applications. Many techniques currently

? applied KDD admit enhancement through the wrapper approach, which uses performance of inductive learning algorithms as feedback to optimize parameters of the learning system.

Wrappers include algorithms for performance tuning, especially of learning system parameters (hyperparameters) such as rates and model priors; control of solution size; and change problem representation (or inductive bias optimization). for changing the representation of a machine learning include decomposition of learning tasks into more tractable; feature construction, or synthesis of more salient or input variables; and feature subset selection, also known as elimination (a form of relevance determination).

This workshop will explore current issues concerning wrapper for KDD applications, including:

problems to which wrappers can be successfully applied

dimensionality reduction

feature subset selection with irrelevant variables

feature extraction

other change of representation for KDD

transforms

ensemble learning (mixture estimation)

empirical performance tuning and experimentation

meta-learning

\* learning methodologies that admit or support wrappers

search-based wrappers (selection, partitioning)

probabilistic wrappers

metric-based model selection

statistical validation methods

evolutionary computation approaches

genetic algorithms for feature selection and extraction

feature construction by genetic programming

"parameterless" GAs for performance tuning

detection and reduction of GP code bloat

linkage learning

multi-objective optimization

foundational theories

statistical learning theory

genetic and evolutionary computation

intron detection and elimination

theory of linkage learning

finding the GA "sweet spot"

the role of prior knowledge and interactive knowledge elicitation in developing wrappers for KDD

Workshop Audience

This workshop is intended for researchers in the area of machine, including practitioners of knowledge discovery in databases (KDD) and statistical and computational learning theorists. Intelligent researchers

start to run for

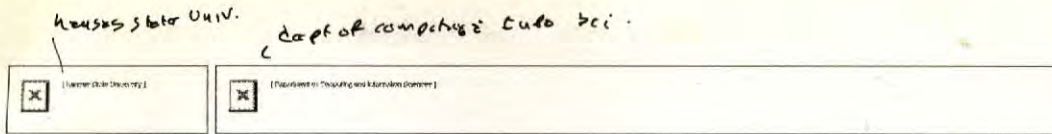
sub-problem

- meta-learning
- learning methodologies that admit or support wrappers
  - search-based wrappers (selection, partitioning)
  - probabilistic wrappers
  - metric-based model selection
  - statistical validation methods
  - evolutionary computation approaches
    - genetic algorithms for feature selection and extraction
    - feature construction by genetic programming
    - "parameterless" GAs for performance tuning
    - detection and reduction of GP code bloat
  - linkage learning
  - multi-objective optimization
- foundational theories
  - statistical learning theory
  - genetic and evolutionary computation
    - intron detection and elimination
    - theory of linkage learning
    - finding the GA "sweet spot"
- the role of prior knowledge and interactive knowledge elicitation in developing wrappers for KDD

## Workshop Audience

This workshop is intended for researchers in the area of machine learning, including practitioners of knowledge discovery in databases (KDD) and statistical and computational learning theorists. Intelligent systems researchers with an interest in high-performance computation and large-scale, real-world applications of data mining (e.g., inference and decision support) will also find this workshop of interest.

width=22 height=22>



# IJCAI-01 Workshop on Wrappers for Performance Enhancement in Knowledge Discovery in Databases (KDD)

[workshop code ML-5]

Saturday, 04 August 2001  
Seattle, Washington, USA

## Workshop Description

(Last updated 05 Dec 2000)

The rapidly increasing volume of data collected for decision support applications in commercial, industrial, medical, and defense domains has made it a challenge to scale up knowledge discovery in databases (KDD), the machine learning and knowledge acquisition component of these applications. Many techniques currently applied to KDD admit enhancement through the *wrapper approach*, which uses empirical performance of inductive learning algorithms as feedback to optimize parameters of the learning system. ))

Wrappers include algorithms for performance tuning, especially: optimization of learning system parameters (*hyperparameters*) such as learning rates and model priors; control of solution size; and change of problem representation (or inductive bias optimization). Strategies for changing the representation of a machine learning problem include decomposition of learning tasks into more tractable subproblems; feature construction, or synthesis of more salient or useful input variables; and feature subset selection, also known as variable elimination (a form of relevance determination).

This workshop will explore current issues concerning wrapper technologies for KDD applications, including:

- problems to which wrappers can be successfully applied
  - dimensionality reduction
    - feature subset selection with irrelevant variables
    - feature extraction
  - other change of representation for KDD
    - transforms
    - ensemble learning (mixture estimation)
  - empirical performance tuning and experimentation

# "WRAPPERS"

.00: The term "Wrapper" originally used is a means of facilitating communications betw. ~~heterogeneous~~ heterogeneous "Modules" or a larger system. Each "Wrapper" converts the language of the Module it "wraps" into a common format/language, so it could communicate w/ other modules in the system: An "Interlingua".

In the IJCAI 2001 conf, "Wrapper" is used in a different way. In a work by Hsu: "It uses performance of individual lang systems as feedback to optimize performance of the system."

Any kind of "tuning" modifications are permitted.

.10 This looks pretty much like what I have in mind for CPM! — T. Q: Do these guys know how to do it right?

Some leads: To get into or visit Doc Google: wrappers } This brot up many Candy wrapper collection, but also a paper defining w/ definition of a wrapper.

W. Google wrapper, KDD was put mainly this IJCAI conf, but also also an earlier version 1996.

W. H. Hsu was organizer of the 2001 workshop: He has a web site:

[kddresearch.org/KDD/Workshop/IJCAI-2001/](http://kddresearch.org/KDD/Workshop/IJCAI-2001/)

Hsu is at U of Kansas Dept of computer & info sci, has papers there on wrappers.

.21 **SN** As General Problem: Try to Define just how one would go about teaching

CPM ~~what~~ the Definition of an OZ problem. One could give many examples, but would TM realize "Rate goal was optimal, is not to specific techniques in the examples? One way into the idea of "coding by constraint" — of giving some constraints, ~~constraints~~ — which narrows to a subset down to "1 of N", then give N.

But also the idea of "Best Value w.  $C_B = C_0$ " is difficult. Maybe first teach what "Best value" is, then "Best w.  $C_B \leq C_0$ ."

T. advantage of TM really learning what "OZ" means is that if we find it to do OZ probs, they say ~~such~~ several other methods, ~~where~~ this would not really make it possible for TM to generalize & find OZ, possibly better method of ~~soln.~~ soln. of OZ probs.

perhaps first (easier) How to teach TM what a "INV problem" meant!

— That it has to find soln. w. min CE.

A perhaps useful Approach: Given 2 soln possible "solns" that should be  $A^* \geq A'$ . ~~both~~ both find soln. to INV problem but A' takes longer — so A has higher more probability of being a soln. to INV problem. Probly a approach could be



00: 328.40 used for OZ prob. & soln.

Naturally, all probl. A<sup>n</sup> values have pc > 0.

So: At present, T. system problem solving system I have in mind is not bad: It works conceivably  
be good and to be able to be very smart using a conventional TSO.

This CPM works foll. way: (1) It has a set of IA's: - Each of which is able to work  
Q<sub>i</sub>A<sub>i</sub> sets, but often only certain kinds of such sets.

(2) ~~The~~ T. initial QA set is suited, foll. which

328.24 - 330.08 covers most of the operation of this preliminary system, in some ideas  
on how to continue it, & some general criticisms. <sup>problems (Q<sub>i</sub>A<sub>i</sub>)</sup>

For "T. Report" it would be fine to just describe the set of IA's & ~~the~~ <sup>initial</sup>

T. initial indexing of problems by user; T. (imp. of indexing by FC) & finally, have  
general ~~use~~ gone to FC's operation.

This creates a CPM of some power, but can be further improved by user of context, &  
(reg. doing such.

Further development, whatever form it may take, so understand "better & better, &

t. error defus of <sup>inv.</sup> ~~inv.~~ OZ problems - what "Optza" really means, etc (338.21ff)

1 from 328.24 - 330.08; we start in. Please IA's will & we get FC) to learn to  
select <sup>IA's</sup> ~~proper ones~~ <sup>specific</sup> for ~~problems~~ At a ~~high~~ more advanced level, FC) will sort of

"create" new IA's because the existing IA's could have many imp. "Initial params" that  
characterize them. Later yet, the sort of IA's will be described as a set of a stack

grammar ~~grammar~~ (T. construction of this grammar will probably be done by user ~~trainer~~)

so FC) can construct new ~~IA's~~ IA's more appropriate to new Q's.

At higher level of development, O<sub>n</sub> can be used to update  
update or modify the operations of one or more of the IA's. ~~just case of~~  
for the IA A<sub>2</sub>, this could involve calcn. of pc's of rings based on

(local or (greatly) extended) contexts.

Also trying doing (even could involve O<sub>n</sub>).

Through "General possi. context" is used ~~to~~  $O_n (O_n, Q_n A_n, \dots) \rightarrow O_n$ .  
A relatively non-el. form of TM.

[SN] In the past, my Model of human prob solving, on which t. F(), <sup>IA</sup> ~~IA~~

interaction was based: That a human would look at a problem & from his experience,  
decide what kinds of prob. solving methods should be tried on it. - This seems

quite different from e.g. F(), ~~IA~~ [IA] systems. In the human system, the prob. solving

methods do not (run much from being used. in the FC), IA system, the IA's are designed to solve  
tho, I did have the concept of "static IA's" ~~that~~ that never did (in any form  
or that had stopped long.



8/21/01

ID

Trivial Example of convergence theorem for Inductive BHGs (≅ universal function) naturally, hvr.

**Convergence Theorem for Stack Operators:**

Say there is a stack operator  $H(Q)$  that generates a PD on  $P_i$ 's for every  $Q_i$ .  
 Say  $O_M$  is a universal stack op. Then if  $H$  is of density  $p \in K_H$  (in bits)  
 for any sequence of  $Q_i$ 's  $[Q_i]_{i=1}^n$ , the total expected error of  $p$  of individual bits, is  $< K_H$ .  
 If  $K$  has continuous paths,  $K_H$  may be a function of just about seq. of  $Q_i$ 's was used & how many  $Q_i$ 's.

It would nice to have a Perm about heuristics in this vein. Vam

I'd like a general formalism for heuristics.

perhaps there such problems are of 2 kinds (2 ways to deriv. what is being searched)

1) INV 2) OZ. (NEW idea! for OZ problem constraints could be partially ordered (?))

A theor is a func. that takes a problem descn, & outputs one or more ways to do it. such means more rapidly.

So this descn look like an induction problem. We could give TM many examples of problem / known paths. — But basically, we'd like TM to really understand what the problem was — that it was to & cc of such.

So if TM really understood what the problem was, then giving examples would not result in a simple induction soln; it would analyze the examples for "hints" on solving the "TRUE" problem.

Perhaps the main problem is to get TM to "understand" what an optzn problem is.

In General "Understanding" involves a potentially inf of ideas! opt all ways of doing something.

For optzn. to start off, TM would know where one kind was better than another.

It would also "understand" the "Better than" relation! transitive, usually non-commutative, (asymmetric) not idempotent  $a > b$  is never true. If you want to know what " $>$ " means also  $<$ ,  $>$ ,  $<$ ,  $=$  etc.

→ If  $b$  is 10 times as fast as  $a$ ,  $a$  is 10 times as fast as  $a$  then  $c$  is 10 times as fast as  $b$ .

In some cases, there is a "central/critical" idea to be "understood" & all other understanding of it is derivable from it by pure logic.

My General defn. of "Understanding" is that it is a way of "explaining" something that is productive. It's a productive code of the data that is "understood".

Hvr. just as humans can be made to understand things in what seem to be the "discovered" way, it should be possible to get TM to understand things in similar ways. I assume humans can understand/interpret languages by using normal (PALP) induction & so we should be able to get TM to do similar inductions.

Perhaps a good "STUDY PROBLEM": To take various problem domains designed for humans, & find an (inductive) way (≅ TSD) to get TM to understand these problems.

The simplest, it would seem, would be to understand what a "cc free" INV problem is; to find  $x \ni M(x) = c$ ; T. problem is to find  $x$ . no number of cc.

00: 3 4:40: ~~One~~ way for TM to understand cc free INV prob. of 34.39!

It just tries "x's at random" to search MCK's c. It has no heurs.

A <sup>Much</sup> ~~way~~ better way is LSrch. It is a more efficient way! Also "trying x's at random" does not ~~do~~ find a method, since some tests may take Time =  $\infty$  (!).

Anyway, The next step is "heuristics": to make do a ~~way~~ "better" <sup>guiding</sup> p.d. than ~~the~~  $z^{-1} x$  in the sense of less cc. (or just less "Time").

The next thing is for TM to find "better heurs". We seem to have gotten into the cold case of all heurs being modular of the p.d. - & they are certainly more than that.

A no-look-ahead heur (No LA heur or LA heur) would involve good PC ~~officer~~ solns for Lsrch; Use of PC to get next trial, in view of previous trials. And somehow, realizing that earn a "good PC" also involves small cc. So we want a p.d.

That gives max  $\frac{PC}{cc}$  to solns! cc is time needed to generate ~~trial~~ and test trial; "Generate time" is then part of the PC data, is characterization.

[ So: Good Heurs find PD's that give max  $\frac{PC}{cc}$  generate at least part solns.

How works out when one is "long during Lsrch" is Unclear.

~~XXXXXXXXXX~~

Again: Ont. first level (TM<sub>1</sub>) TM<sub>1</sub> must recognize soln. of problem & also recognize when one soln. is "Better than" another.

On second level (TM<sub>2</sub>) TM<sub>2</sub> is there to find methods of working problem that are as good as poss. This involves trade off betw total cc & total PC. - to the user has to specify just what the trade off is: say  $f(cc, PC)$  is a linear ordering of

22 -  $(cc, PC)$  pairs - to tell which is best.

My present orientation is to have TM<sub>2</sub> working very early in the exp. That way discovery of heurs, or just simple context dependence of some PC's of rules would be TM<sub>2</sub> problem.

Probably, it shouldn't be diff. to teach TM a number of INV & OZ problems, by giving examples. First cc free INV: Give examples of  $G(C)$ ,  $C$ , and  $X$ , so it could tell if a particular new  $G(C)$ ,  $C$ ;  $X$  was a case or not (M-TM problem).

Then, looking at the trace of several solns to INV problems, teaching which method less cc than which other. - the concept of "less cc" as a relation betw traces.

7. Concept of ~~simple~~ soln. of min cc.

So, I think it ~~may~~ may be poss. to tell TM what we want it to do. Hvr since there will be only 1 or 2 situations in which we need this critically, we might as well program it in! The needed cones are cc free INV soln., maybe cc soln of INV prob; perhaps find Optimim(cc, G) in OZ problem. Given to user supplied  $(cc, G)$  (linear ordering function (see (20-22)).

342.40: HVR, + near idea of "In theory" teaching TM to concs. needed to understand

what its all over concs, by "theoretical, large TSP" — ~~Practical~~ This is to pick it.

Conc. in proper form to give w. other things it ins.

So: It would seem ~~Practical~~ any of TM's well defined problems ~~TM~~ can be taught to understand & work on those problems. w. suitable TMs, TM can be

usefully/successfully work on TM's problems.

For 2 immature TMs, TM's problems should be generalized so that TM, can easily work on them. R more mature TM, can work on it. (all)

non-el problems of TM2.

So, it looks like I have a completed term of TM. — Good run up to 5 years on TSP!

Write report, w. special # & indices for items that I have written above

But I need expansion! GIVE refs to my notes.

Put indices in "original report" Now remove k, indices for k copy sent; qq, etc.

do IDISA. (There are probably variables to write down in or end of report.

A recent ~~is~~ <sup>was</sup> ~~retire~~ <sup>retire</sup> that it may not be diff.

TM + meaning of various commands — things that I

want TM to do: (like solve INU, or a problem... Give

Soln. w. sign CC or Max some kind of CC P, etc.

A variable format:  $\text{Find } X \text{ such that } F(X) = \text{Max.}$  Actually, this is equivalent

of a function: T gets  $x, f(x)$ : f. compute value for X.

In learning how to evaluate any function, TM could evaluate  $f(x)$ : what kind of function was?

Did X have to be numeric or string? was  $f(x)$  continuous? Did there are max slope/corner/curve?

If X must be a string? — Are constraints "barrier" etc.

Find X such that  $F(X) = C$  || Find X such that  $F(X) = \text{Max}$  & times T

TM finds Max in  $\frac{1}{2} + \frac{1}{2} = 1$  & the format  $\frac{1}{2} + \frac{1}{2} = 1$  gives + soln

in smaller way,  $\frac{1}{2} + \frac{1}{2} = 1$ , X are reduced.

from base of solns, TM then in order to make soln.

TM is not redundant to Gino compression.

343.16

343



1.00 : 344.40 : In the GA EA model: We could use deterministic models only; <sup>Gave for each case</sup> ~~the~~ would be the no. of correct responses to the set of QA's. — or perhaps a  
with sum. To make progress, we use a weighted mean of the counts, with being a function  
of "no. of correct responses". Each can has ~~some~~ predict for each  $Q_i$ ; for each Poss  $A_i$  we will add up all correct  
pairs ~~that answer~~

1.05 : We may also score  $\propto \frac{1}{C_i}$  per response or ~~or~~ ~~for~~ ~~each~~ ~~choice~~  
Give a score ~~to~~ ~~each~~ ~~choice~~ ~~of~~ ~~the~~ ~~entire~~ ~~corpus~~.  
Some sure ~~but~~ that ~~that~~ ~~is~~ ~~not~~ ~~good~~ ~~for~~ ~~a~~ ~~universal~~ ~~d.f.~~ — but my counter example ~~was~~  
wrong! It may be able to give a Universal d.f. — but whether it's a good "practical" model, is unclear!

1.07 : In fact, it is easy to prove that for any seq. of Df's  $A_i \in \mathcal{A}$ ,  $\exists$  an ensemble of ~~the~~  
sequences of deterministic choices for each  $A_i \rightarrow$  the max. of an ensemble = the desired stochastic d.f.  
To get any a particular number of an ensemble  $A_i$ ; just pick random  $A_i$  w. prob =  $P_i$ .  
The max of an ensemble will give the ~~max~~ ~~prob~~ ~~for~~ ~~all~~  $A_i$ .  
Does the model ensemble model of 1.07 - 1.10 suggest anything about how to find such an ensemble  
empirically? — (from data) — it would seem not: which is perhaps why it might be a bad model!.

1.10 : In GA, we certainly wouldn't worry about BLOAT — since the Gave takes care of that.  
The param to be adjusted is population size: would this have to grow ~~much~~ ~~as~~ ~~the~~ ~~corpus~~ ~~size~~ ~~↑~~?

1.20 : Since I plan, to somewhat better, what kinds of modifications of the On are needed, I  
should be able to design better crossover/mutation Algos. I will be trying to implement  
the AZ (distrib. over functions. — Here, I would want more context dependencies  
& mod. of "pairs" R.D. by "long". If I can't implement these in GA, it ~~will~~  
will be much weaker than the original model — that starts w. AZ! & then uses  
context & "long" during search to improve itself. In GA, I'd somehow have to find  
a way to implement context dependence in "long". ~~There~~  
There is a automatic variety of "long" in GA, in that code of low Gave are  
less likely to be used to generate new code.

How I'd get context dependence (e.g.  $O_{n+1}$  ~~could~~ ~~depend~~ ~~on~~  $D_i$  would depend on  
 $Q_{n+1}, A_{n+1}$  — also various past  $Q_i, Q_{in} A_{in} \rightarrow O_{i+1}$ , ~~data~~.

2.00 : ANN (RANN): In my own experience this would be for numerical

2.05 : I/O : (No ANN has been used ~~for~~ ~~string~~ ~~data~~ & graphs.)  
Anyway, the state of a system <sup>can</sup> be represented ~~can~~ be represented by a pt. in  
wt. space. If the no. of wts is minimized, we want maximally flat soln.  
New data, then would try to move ~~the~~ ~~minimally~~ ~~from~~ ~~old~~ ~~pos.~~ ~~in~~ ~~wt.~~ ~~space~~.  
Then when w/ more progress we should return many "max flat" pts in  
many. As we move ~~the~~ ~~minimally~~ ~~from~~ ~~old~~ ~~pts~~, we could bifurcate  
to generate new good pts — also discard some not-so good bad ones.  $\rightarrow 346.00$

2.10 : If we use many wts, we, again should return many pts in wt. space,  
but as we drift, we will, because of too many wts, end up "overfitting".  
Normally, this is dealt w. viz "cross-validation" (also faithful) but I don't  
see how to do it in present case.

2.15 : Back to small no. of wts: we could try using <sup>2</sup> diverse set of pts in wt. space  $\rightarrow 346.00$

Fri ID

$$x = a + b \cos \theta$$

$$y = a + b \sin \theta = \frac{1}{x+1} = \frac{x}{x+1}$$

$$\frac{dy}{dx} = \frac{1}{x+1} - \frac{x}{(x+1)^2} = \frac{x+1-x}{(x+1)^2} = \frac{1}{(x+1)^2}$$

spec (345, 34)

so we would simultaneously track several (many) dist. wt. sets in wt. space.  
 Note: for each "set of wts", there can be > 1 pts of good fit/mess & yield.  
 Each "set of wts" is a subset of poss. connections betw. neurons, & for each set

there is a subspace — w. one or more peaks of Gorc — each w/ its own Hessian —  
 So  $\frac{Gorc}{(Hessian)}$  = true Gorc of that pt. — T. Gorc has to be a ~~total~~ product of PC's of some times pc of Model. (pc of model is  $\frac{1}{\sqrt{\lambda}}$  | Hessian  $^{-1}$  — if all eivals of matrix are  $\gg 1$  — if any eivals are small, this amounts to reduction in dimension —

### Too many wts!

TM's job in ANN is to add (remove) Neurons/wts also to give min/max values for new wts.

- Outline of report:
- 1) Intro: & "part I write"; Also add stuff re: General pre-knowledge system.
  - 2) Details of what IA's are: Give list; perhaps more detail on some showing how they work.
  - 3) FC) how it works: initially class IA's by indices; later terms which IA's to use.
  - 4) ~~GA~~ Updating: a) Pruning: removal of FC) for update (?) the GA b) Use of On to update: @ context, @ long delay (set) for A-Z
- Mod of ~~GA~~ (crossover, mutation), population size.  
 notes can be used to cross Mod. — for it no longer looks like GA/learn GA?

SW

The A-Z is GA can be used to learn numerical (vector I/O) — also has an input subclass of IA, 345.26 - 346.10 discusses update of Gorc for update in GA, ANN, (see recent on-line complexity paper (C:\PS) on ANN/GA & some how ideas on this.

The "Begining I wrote, was 332.01 - 333.13 : It der to u (Analog System)

$O_n(O_n', G_n, A_n) \rightarrow O_{n+1}$  also sure for  $O_n$ .

Next, der to IA's, give general details — say  $O_n$  is a IA: we will der to  $O_n$  or IA in some detail (AZ1): GA T. network in which GA & ANN/RAN can be used as IA's and der to in Appendix. Mention IA's don't have to operate on All inputs. Also, we will add IA's to system based on our future experience w/ TSC's.

$O_n$  is not an IA, but it is a net generator ~~to~~  $O_n$  sub E. A km. generating  $O_{n+1}$  from  $O_n$  does constitute an IA.

A linear or non-linear "curve fitter" is a kind of simple IA with a simple or delay & gain. when given a set of  $\{X_i, Y_i\}$  pairs ( $i=1, \dots, N$ ) it will give a PD no  $Y_{n+1}$ .

Linear & non-linear regression can be regarded as types of IA's. One way to do this:

$$O_n = X_n, X_{n-1}, \dots, X_{n-k}; A_n = X_{n+1}$$

The  $\{A_i, Y_i\}$  data set can then be obtained from any time series.



A2S stochastic  
A2D deterministic

00:332.13

Section: Induction Algorithms:

An Induction Algorithm (IA) is an algorithm that is able to look at a sequence of  $M$   $(Q_i, A_i)$  pairs, like the system of ~~332.03~~ and ~~wrap up~~ then, in reply to ~~new~~ new input  $Q_{n+1}$ , give a probability distribution on possible values of  $A_{n+1}$ .

IAs will differ in their efficiency and accuracy. They will also differ in the kinds of  $Q_i$  objects that they will accept. Some accept only strings, other only numbers, others only vectors. ~~While~~ <sup>IAs</sup> must give output ~~probabilities~~ probabilities for all conceivable  $A_i$  values in their range, some are <sup>stochastically</sup> ~~deterministic~~ deterministic, and give only one  $A_i$  value for each  $Q_i$ , with a probability of 1 for that  $A_i$ , and zero for all others.

Note: 331.26 tells how to convert a deterministic IA to a probabilistic IA

We may view each IA as a kind of "specialist" in a particular area of production.

For our preliminary ~~approximation~~ approximation of the system of 332.03, we will use a ~~kind of~~ <sup>weighted mean</sup> weighted mean of the predictions of ~~several~~ <sup>several</sup> IAs. At first, the trainer of the system will assign one or more of the ~~several~~ IAs to work a particular problem in the training sequences. Later, the system will learn to assign IAs to problems without the aid of the trainer.

omit

The IAs that we will describe in most detail are based on universal function languages such as Lisp or Fortran. We will also consider IAs based on Neural Nets (ANN or RNN) as well as IAs based on Genetic Programming. The last two will be described in appendices I and II respectively.

We will describe in most detail, two IAs based on functional languages. One, DAZ, does deterministic production, the other, SAZ, does stochastic (probabilistic) production. IAs based on Genetic Programming and on Neural Nets will be described in Appendices I and II respectively.

Many different kinds of IA are possible. During the course of our training of the system, we expect to be adding new kinds of IAs to the set used by the system, and assuming ~~that~~ <sup>that</sup> ~~some~~ <sup>some</sup> ~~of~~ <sup>of</sup> ~~the~~ <sup>the</sup> ~~new~~ <sup>new</sup> ~~ones~~ <sup>ones</sup> ~~will~~ <sup>will</sup> ~~be~~ <sup>be</sup> ~~completely~~ <sup>completely</sup> ~~emulated~~ <sup>emulated</sup> by other IAs if their functions have been completely emulated by other IAs.

N.B.

insubstantia 361.26

00 : 347.40 : A way to get 2 ideas Across : (A) Idea of generality of what CPM does (B) Idea of what various

IA's can be ( = Specialists / consultants / experts in special field )

This would be a listing of some IA's

1) Curvature : linear or N-linear ~~vector~~ Vector  $\rightarrow$  Vector =  $\sum_{i=1}^n z_i x_i$  or ANN.  $\sum_{i=1}^n z_i = Q_n$   $z_i = x_{n-i} (i=1|n)$

2) Classifier : (e.g. ID3) The  $i^{\text{th}}$  data set  $\Rightarrow Q_i = x_0^i, x_1^i, \dots, x_n^i$ ;  $A_i = y^i$ ;  $y^i$  is the class decision rule  $Q$ ; vector  $\rightarrow$  Object in to Classifier as an example of IA

3) Theorem prover : Corpus is set of theorems and assoc. proofs.  $Q =$  Theorem,  $A =$  Proof of Theorem. Some theorems are "global"; they do proofs from primitive axiom set. An incremental TH we will be interested primarily in how to deal with incremental term provers: Let use proofs in previous parts of  $\dots$  sq. Theorems can be in Algebra, Logic, Geometry or any other parts of math.  $\rightarrow 350.10$

13  
14 **AI** In regular GA: on "crossover" If we have 2 "parent" function trees: To create third tree so that  $\rightarrow$  PC set.  $\exists$  is minimal; or create 3rd as a seq. of parts in order of PC.

In ordinary Koza crossover: Every job-sub-trees is coded as:  $A, A', B, B'$  are defined in processes.

so  $\exists$  examples such  $A, A'; B, B'$  in  $AB'$  (say): In each case we have to find just how the connections are to be made: i.e. what "terminal" of  $A, A'$  is connected to! A simple way to calculate  $\frac{1}{PC}$  would be no. of terminals on  $A$ .

In a similar way, an can code  $\exists$  can parents  $\rightarrow$  crossover offspring.

At this "Lord" this way of looking at things ~~but~~ doesn't generate results much different from normal Koza crossover of pairs of parents: I think a v.g. part is where we define common sub-trees or try to find common sub-trees. If may be possible to get this & check w.o. hunting for common sub-trees" which sounds expensive.

A poss. way to get grammar that generates code  $\exists$  w. assoc expected G proc; We generate code using some simpler grammar: We divide the code into maybe 4 percentiles w.r.t. GProc: for each of  $\epsilon, \theta$  we try to find a separate grammar. The 4 grammars will have certain common elements. This will be the "Mother Grammar". The 4 offspring grammars will be modifications of the "Mother". It may be easy to reparameterize the grammars w.r.t. their percentile no. (0, 1, 2, 3).  $\rightarrow 370.24$

34.08 **AI** Re: Classifiers: T. ~~the~~ binary set  $[x^2, y^2]$  is always very limited: it is better w. one particular class. T. info from one class. is not used with next class. problem. So in this sense,  $\epsilon$  this kind of classifier is not a IA.

37 In fact, examples 1) (0.05) 2) (0.5) 3) (0.6) all have this objection:  $\epsilon$  from prover of 4) (0.9)  $\rightarrow$  0.4, 1 hr. We may be able to use them for update improvements, but, since we always have the same problem in this case.  $\rightarrow 375.06$  ~~sec~~



00:34.40 : In many of these systems, we can grossly modify them in a discrete way by changing the "basis" of the system ~~for the problem~~  
 In curve fitting, used different "Basis" functions, or use different kind of squashing function in ANN.

(Note that "Hinges" are ~~very fast to compute~~ - Pro may have discontinuous first derivative ~~(piecewise continuous)~~)

06:39.40 : While this is true, we may not get much useful recursion (i.e. heavy to system implementation).  
 If we always use the same "IA" for updating, it is not linked to other problem cases.

Consider a curve fitting IA! It could work on many different curve fitting problems - each would have a different "Index no". While each curve would be fitted individually with its own "update parameters", ~~the~~  $\theta$  params for all of the curves could have common features in common - e.g. Bias, could be "updated" in a more "Global" way. So this particular IA would have a "individualized" update  $\theta$  sum.

12:39.40 : An alternative to having all of these IAs in  $\Pi$ , selected by FC for a particular problem - would be to have individual IAs as before, but have parameter "CALL" like front "Menu" system. In the present scheme it might be possible to do "calls", so e.g. +  $A=1$ , say, could call On to solve one of its sub-problems: This call would then go thru FC) which ~~would~~ would decide which IA's to assign it to.

In Sol 86 I had this "Planner" heuristic that called various prob. solving routines. Perhaps "in spirit of" ~~it~~.

20:39.40 : A kind of SOP I'd need for hours (I think): A rule of thumb type should be discoverable! In problems of type  $\alpha$  (i.e. in  $\Pi$ ); If class  $\beta$  condition occurs,

then response of class  $\gamma$  is helpful w. prob.  $p$ .  
 Examples: In <sup>symbolic</sup> integration: when integrand is of form  $UV'$ , then int. by parts is rational vs. usually useful trial. I think this sort of stock rule can be picked up by correlation studies. We use ALP to see if e.g. "useful" - i.e. a compressive code.

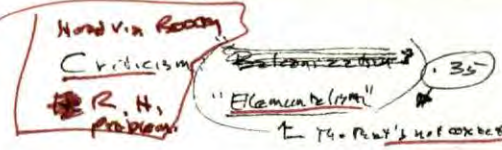
Long ago I had this idea of a set of op. formulas, a student of correlation between obs & ops.

26:39.40 : Another Go at correlation: One makes a certain Ob out.  $Q_i$  is got a certain result. (to be used for decision) This result correlates with (one has to predict well)  $t_i$ .  $A_i$  or certain features/parts of it. - So this is compressive code. Now: Could such correlation be found using a simple model of operators?  $[M, Q_i, P_{in}, R] \rightarrow A_{out}$ .  $\rightarrow$  Is e. form: If  $\alpha(Q) = e_i$ , then

29:39.40 :  $[A_i]$  is probably limited in ~~the~~ way: i.e. D.P. on  $A_i$ . - This may amount to a p.d. on aspect of  $A_i$ , so which, coupled with an aspect on  $A_i$  gives us an aspect on  $A_i$ .

31:39.40 : How ~~the~~ (29L) - ~~is~~ is ~~redundant~~ formula (i.e. (29L)) is unclear!  
 Actually, it's not clear how this can be put into a usable def. on  $A_i$  - let alone ~~the~~ p.e. d.P. in ~~the~~  $(29L)$  form. So, TM not only has to find "productive regys", but it has to find "productive regys" rather expressible in a usable form - but can be used to help get a p.d.  $\leftarrow$  (in some usable form).

32:39.40 : More exactly, TM has to find obs on  $Q_i$  that can control operators that construct  $A_i$ . As this IA matures, it gets better obs for  $Q_i$  a better control operators to construct  $A_i$ .



00:349.00 : Other than TSO design, the main concept / problem is improving & updating Algm.

01 This means finding ~~some~~ contextual requirements/dependencies for cues used in modifying  $O_1$  to create  $O_{n+1}$ . As TM evolves, these contexts will ~~change~~ become larger. Nothing is very ~~much~~ amount of success/failure of previous trial's solns to problems, is also an input they need to be ~~learned~~ learned.

One trouble is that I have no way for AZ1 <sup>data</sup> or AZ2 <sup>state</sup> to notice or even express regularities of c. type that are needed for  $\neq .01-.04$ . - The 349.26  $\rightarrow$  to is a step in this direction

Essentially, the way I'm dealing w. this situation is to call room for lots of "IA's" ~~that~~ ("specialists") that ~~are~~ supposed to be able to do this job.

- 09
- 00:346.13 **EN** On ~~the~~ examples of <sup>(job-comp)</sup> IA's that <sup>also</sup> "continued long":
- 1) One that imps Algebra: new Datas, new Params, new problems.
  - 2) Same as 1) but for other ~~cases~~ areas of math - say Geometry.
  - 3) Learning to work by holistic interpretation (both <sup>deterministic</sup> & <sup>indeterministic</sup>)
  - $\rightarrow$  4) Learning Q answering in English, about area that TM has ~~learned~~ learned much about (e.g. Algebra).
  - 5) Regression of Time series: Each TS is a QA. We give TS's of ~~more~~ <sup>more & more</sup> "side info".  
<sub>its own</sub>
  - $\rightarrow$  6) Improving & Updating Algm (see 01-09)
  - 7) Try to solve e.g. problems symbolically - non-linear types of number complexity
  - 8) ~~symbolic~~ symbolic solns of diff. equs. <sub>Ramanujan</sub>
  - 9) Go thru ~~review~~ review to test Ramanujan use. (Summary of Results in Elementary Math: Cont.)
  - $\rightarrow$  10) Learning Logic, applied reasoning - applied to math problem - able to apply ~~theoretical~~ theoretical logic to practical problems in math & other areas.

20 : 346.33 Use of **GA** v.s. **AZS** <sup>stochastic</sup> for Updating.

- 1) GA is standard, technique; doesn't have to be developed much. - AZS is usually fairly expensive.
- 2) GA appears to be capable of being very non-linearist - AZS seems to be essentially linearist (but I've not cover it really bad!)
- 3) AZS may deal more directly with the context dependence of cues in updating -

Also it may be able to introduce frag from previous trials - i.e. I don't know how to do this (yet) w. AZS.  
<sub>stochastic AZ1</sub> [4) GA usually retains many "good" codes & usually discards "bad" codes. This is usually stochastic AZS (but could be).  $\rightarrow$  349.22

32 .09 : A critical idea here, is that the IA that works on Update ~~with~~ improvement after has as part of its S.C. (subgoals - data set - esp), "main line" problems to solve that will ~~help~~ help it to learn ~~(AZS)~~ (AZS) to solve its update improvement ~~problem~~ problem.

What present Model ~~does not~~ ~~seem to~~ ~~do~~ seems not to do: is mix ideas, cues from diff. IA's. A possible way to deal w. this: there is a main IA that watches over IA's & tries to incorporate their methods into itself. - however, it has to "understand" the methods before it can usefully do this - i.e. the "external/experts" operating here to be "folded" into cues, useful to the "main IA" before the "main IA" can adopt them as its own. Main IA could regard other IA's as

The Real problem (experts in same room v.s. expert in computer).

IMPT Circumstances

00:350.40 part of the "advice channel" (14)

01: ANN (and GA) as IAS (in the long term long sense). ANN for ANN, each problem, w. its data set, is regarded as a "Q". What the ANN learns, is how to take a Q (a problem descn) & decide on how many & which initial wts to use & how to introduce "pruned" wts. during the learn. Also how best to do momentum, (not exact) steepest descent for solns. (a actually the Newton method using second derivs. would be better than steepest descent - it is the standard Non-linear Optim method - but for >100 wts, cc = N^2, so it may not be practical. - is there a way to use a simpler matrix, e.g. Toeplitz, that would be solvable w. cc = N? (By the way, Hessian usually only uses first derivatives! so it would use/save info as "steepest descent" - I think w'd have to use into a vector from last several pts (maybe 1000) in order to get usable matrix out of formulae of f(i), g(j).

02: But rather H\_ij = sum\_k f\_k(i) \* f\_k(j). This info be regarded as an axiomatic "Momentum" computation! An adjustable param is the "width" of the x window. It may well be that system has no second derivative at any pt, but when it does average over several pts, it does. Some concern about the allocation to fitness changes. para 12.22

14:00: Actually this learning to understand is used "external experts" can be an important aspect of TM's mainly. We introduce expert IAS into the system, as well as humans "open" to possible "Mem IA" (= "MIA"). MIA & Mem IA tried to do a small correction. Then M. experts, internally, by expressing their methods in their own internal lang, & using part of their methods as a/o models of them, to solve the same problems at least as well. It may be that there is a natural GORC for this involving usual PC of (comp + oper) descn. External experts on. 352.00

22: JGR, 01 Re: GA as a IA: Like ANN (or if bit diff): GA would work/problems used

Way: each work be a "Q": What it does for continued try is: for each Q, TM learns to devise a set of mutations, crossover & initial population appropriate to that problem. (I think the usual "Wrapper" approach does give a lot of RTUs, but TM should be able to do it better) - In particular, the "initial population" can be a v.g. guess of what the soln. is like. I noted some time ago, that this would patch up over critical weaknesses of GA systems. Another critical weakness is excessive time spent at evaluating the "fitness function". The system could also learn how to do this better & better: also modify fitness function - Or since TM is given a problem, it has to devise a fitness function. It could do this w. cc expanded in mind - also modify the fitness function's population moved towards a solution.

So to have essentially 3 IAS: AZD, AZS, GA, (R)ANN.

What about classifiers. 348.00, 348.40. Each classifi. problem, w. its data set, is a simple Q. Hvr. note 348.37-40. A classifier can be used for updating - it works on the same problem (updating) all the time & parts in new data. 349.06-11. It is a "weakness" in the use of classifiers in this way - but its inherent recursion loop of 3.5. This is less true of GA & (R)ANN because they can be significantly improved! See 22-31 on GA, (01-06) on (R)ANN.

8/26/01

ID

1804

19

Aug 18  
Sun

guide

Money, Things,  
Power, Control  
Estimates, Popularity, Love  
(All in "short supply")

352

Understanding, Knowledge,  
"Truth"  
Understanding, Dispute,  
Knowledge for the  
Wisdom

00:35:21 : T. relation betw. GA & AZ (B) can be/close! T. majority GA uses for mutation/crossover  
can be related to AZ's Update Alg. — How  $O_{n+1}$  is to be  $\geq$  (Mutation/crossover)  $O_n$ .  
[ T. reason  $O_n \rightarrow O_{n+1}$  can involve "crossovers": Plot "crossovers" area related to crossover elements,  
Sub-trees can be defined as functions, etc. ] GA & AZ are very close: Eventually will write the Merge them.

What AZ can learn from R(ANN) is unclear. — But probably AZ can improve R(ANN) considerably. — (If P.D. is true, then AZ is away from R(ANN)! — !!

10:33:25: So far, most seems to be, from (P.D. on  $O_n \rightarrow O_{n+1}$ , Ann) aspects P.D. on  $O_{n+1}$ . which suggests that one should keep retain a lot of  $O_n$  codes a/o P.D.'s on  $O_n$  "around" various pts.

A (perhaps expensive) Alternative for retaining Many codes is f. Ability to "Backtrack". — also special methods to deal w. OSL

13: On ~~update~~ GA's for updating  $O_n \rightarrow O_{n+1}$ :

One big problem is eval'n. of codes! 2 differs: (a) ~~evaluating~~  $P_c(O_{n+1})$  code!

This might be made easier if we only try modifications that don't break up ~~of~~ functions defined  $\geq$  the problems also.

(Backtrack  $\geq$  the problems) (b) eval'n of p.c. of codes of all old  $O_n$  A's! This can be done by sampling.

18: ~~XXXXXXXXXX~~: Also by using mut/crossover restricted as (a) ~~you will perhaps merely~~ Th. Q/A's for R

20: have codes w. single evalus for  $z \leq n-k$ . — Tho we may find other mut/cross restrictions so is true!

21: To f extent that we ~~do not~~ conform w. restrictions of (13-20) ~~we will be~~  $\approx$  (AZ)

22: Very close to f. Such codes selected by normal (AZ)!

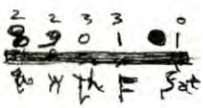
GA becomes very close to AZ then this will give us an undery. finding of how best to do mut/cross in GA in General! — since GA has no v.g. General Theory.

Contextual considerations should improve "hit rate" for both AZ & GA.

GA should "learn" from (bad) codes as well as (good) ~~code~~ codes. — but how to implement

This is unclear! Also how to implement learn in AZ in General, needs to be developed.

30: (SN) Way to Get Time & Peoples' Comments: Use log of "Patrons" time to answer patrons' Q's! So they can keep track of how smart TM is being. Perhaps limit briteness to answers about SM! @. Answering Q's may not take much time: Its "update" that takes most of time. "Patrons" can watch f. development of TM as it gives smarter & smarter Answers to Q's.



DAZ  
SAZ

00:35:40 So: What has to be done on report!

1) Defn of 4 IA's: AZD, AZS, GA, (R)ANN (RANN may be w/o 4. first 3) (347, 26-26 tells how we will actually use many more, different IA's & change them during JDI's turning)

For each one, describe how updating is done what first & how it will be done in more mature TM.

- Just to some degree, TM can work on improvement of update alg when very busy.

On "event" input TM can work on improvement of certain aspects of updating. uncoupled  
Begin discn. of 4 IA's as separate systems - each with own TSO - each with own update Alg.

Next discuss how USR selects TSO for each date - but later a different core issue.

" " " FC Is implicitly named as date - perhaps discuss "content" of "log"

" " " each update Alg is improved - perhaps discuss "content" of "log"

Where discuss look for updating: How it is actually improved by Modifying & Guiding PD

Modifying & Guiding PD back trials. We suspend as much as possible as / on to previous trial.

At end, discuss "lookahead" v.s. Control (= Mgmt) in Round Machines.

Advantages/disadvantages.

Dark DAZ in some detail, & how it related to TMC, a one (resp. for R), border between (All esp functions can be described in R).

Also describe update alg. (Machine either look or GA can see old for update - gives PD on Out)

Then tell how SAZ differs from DAZ. (how Machine differs from).

Then discuss GA: Tell how it differs from SAZ. only in its update Alg.

Actually, a set: GA does each problem from scratch: Next, we look for 4 forms

of GA such (size of popn, initial popn, mut/cross algs...) as a direct problem ( $\in Q_i$ ).

Updating consists of improving these Algs. (Use of look for this?)

[SN] It may be possible (perhaps easier) to have DZ & SAZ direct use  $\Rightarrow$  in put vnc. for SAZ!

Just a set of On's. How, while I had a proof that no exists to simulate any (finitely describable?) PD, I'm not entirely happy w. this model. (couldn't find proof: it was very simple!)

The JDI deterministic model had for a set of problem, A (w. probly  $P(A_i | Q_i)$ )  
So one has constructed this sequence of deterministic pred. i. for, all of it; Their mean is to describ P.D.

After writing this Report: write down the poss. choices for each component in the system:

As many reasonable options as I can think of.

- 1) List of IA's: I have some list. Add correlational IA's: 349, 26-40 (These will be very imp. for hours & for new IA's suggested by the TSO.)
- 2) Different ways to do updating: Look v.s. Straight GA.
- 3) 2 v.s. 3 input vnc for stochastic AZ (see .22)

00:353.40 : **UNIVERSALITY** of certain IAs: If IA<sub>0</sub> is universal, then for any TSG, T }  
 and any other IA, IA<sub>j</sub>, The pc given to T, A<sub>i</sub> sequence of T, by IA<sub>0</sub> will > k · TSG  
 given by IA<sub>j</sub>. It will depend on IA and IA<sub>j</sub>, but not on TSG · T.

Q: Do universal IA's exist? Do they necessarily involve infinite cc?  
 Well, the standard ALP is IA is a cond for universality: If it doesn't do it, then it's unlikely that any other IA will do it!  
 I think the data is wrong! Instead of say TSG: pick any (stochastic) Q<sub>i</sub> → A<sub>i</sub> function (finitely decodable). Then pick any Q<sub>i</sub> sequence. Every IA will have a certain expected error in assignment of pc's to a seq. of correct A<sub>i</sub>'s.  
 Or maybe a certain expected product of pc's of " " " " " "  
 This (.10)/product, say, dominates " (majorizes)" met of any other IA. Stochastic  
 T. most serious competitor to universal IA would be the IA that KNOWS and EXPECTS T. Generating Q<sub>i</sub> → A<sub>i</sub> function. It will have a certain expected error (which it can be zero in the case)

14 P.S. to Perry. Seems to work! K is not. (dit) down cost of producing A<sub>i</sub> sequence...  
 15 T. IA that competes with seq. exists has down cost of at least k,  
 I. Error of the universal IA will be  $\frac{1}{k}$  (cf. 15) (viz. Rorem of .11-13).

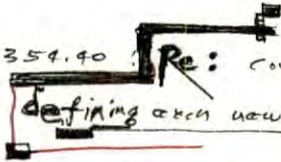
20 (SN) Th. problem of Derby (AZ): **Recursive funds**: deriv in **AZ**  
 say  $f(x) = f(x-1) \cdot x$ ;  $f(1) = 1$  is a function deriv.  
 $f(x) = f(x-1) \cdot (x-1)$ ;  $f(1) = 1$  defines  $f(x) = x!$   
 In our definition of functions F, the function name, is followed by a dummy variable x.  
 (F, x) which is followed by an expression that defines F in terms of x. If the expression does not contain F, then that's fine and ok. If the expression contains F, then we have at least 1 more eqn to define F. so we write  $f(1) = 1$  as  $f(1) = 1$  and add function symbol.

28 On assignment of pc's to strings of deriv funds or strings:  
 We start with a finite set of symbols and the possibility of adding symbols at certain points in the sequence.  
 We also have a set of rules that tell which symbols are legal at each pc in the seq. - at certain pts, new symbols of certain types can be introduced, this is expanding the alphabet.  
 Say we want the pc of a particular symbol in a code string. P<sub>1</sub> and P<sub>2</sub> are different symbol types in the same P<sub>1</sub> and P<sub>2</sub> are the symbols that are legal at the present pt., they have occurred a total of r times in the string P<sub>1</sub> and P<sub>2</sub>.  
 The symbol we want the pc of, has occurred x times in the string P<sub>1</sub> and P<sub>2</sub>.  
 so  $P_x = \frac{x}{d+r}$  This is a kind of general of "Caprice's Rule"  
 See comment, here  
 A pos. (trunc/modula) (not really good!) r includes only those cases in which a symbol was in a situation corresponding to the one now being considered.  
 355.00 spic



ID

Is IA "Induction Alg?"



It may be a way to introduce arity regularity types into the lang!

04 AH! A nice way to start out! around 332.03 start out by defining a IA as a device that will deal w. QA seqs like 332.03ff. Then

7. "report" as a talk, what I will talk about is <sup>non</sup> ~~research~~ <sup>program</sup> ~~supervised~~ for learning.

First ~~to~~ <sup>some detail about</sup> describe the problems I'm trying to solve! Then ~~some~~ some solutions to

Real problem: ~~Some~~ <sup>Some</sup> Very accurate, very theoretical and impractical to implement. Others <sup>known algorithms</sup> commonly ~~known~~ solutions, easy to implement but not very ~~accurate~~ accurate.

Others will be ~~rough~~ <sup>rough</sup> between these two extremes.

I will learn take kind of weighted ~~down~~ <sup>down</sup> of several of these ~~rough~~ <sup>rough</sup> approaches and use it as a general solution to the problem.

It is characteristic ~~of~~ <sup>of</sup> some of these algorithms that they improve accuracy as you give them more and more data, so after we have trained the system with ~~lots~~ lots of data, we will solve the problem of improving itself.

start Section I Perhaps write a main talk, then write an introductory review.

The problem we will be addressing is a very general <sup>which</sup> problem inductive inference; we have a string or a number or a mixture of both, ~~we feed~~ <sup>we feed</sup> into a black box; we observe the output - another string or number or both.

Symbolically,  $Q_i$  is input  $A_i$  is ~~output~~ <sup>observed</sup>. We are given a large set of pairs  $Q_i, A_i$  pairs - Given a new  $Q_n$ , what is a good estimate of the probability distribution of  $A_n$ ? <sup>something in the direction that gives</sup> Here we assume <sup>that</sup> there is some probabilistic relation between  $Q_i$  and  $A_i$ .

This problem is very general. The  $Q$ 's and  $A$ 's can be descriptions and answers in a formal or natural language. The  $Q$ 's might be descriptions of problems and  $A$ 's their solutions.

We might consider the  $Q$ 's to be the stimulus of a biological organism - a call or a smell or a human - the  $A$ 's are the responses of that organism. The  $Q$ 's can be ~~physical~~ <sup>scientific</sup> experiments - the  $A$ 's are the results of these experiments. - ~~that~~ <sup>general</sup>

Viewed in this way, the ~~problem~~ <sup>general</sup> we are trying to solve is the problem of all the sciences - prediction.

is based on ~~the~~ universal probability distribution.  
 We will define any algorithm that attempts to solve the QA sequence problem as Inductive Algorithm (IA)  
 IAS are known

There are many ~~ways~~ methods. Many ~~methods~~ are known for dealing with this problem.  
 The most general of them all is ITSS  
 It is known to give about the best results possible, but takes entirely too much computation time.  
 We will consider time-limited versions of this algorithm.  
 We will also consider other, simpler algorithms that are known to work well in certain areas of inquiry.

One important critical aspect of an IA is its "update algorithm".  
 Suppose we have an IA operating on a system so that after the system has seen  $[Q_i, A_i]$   $i=1 \dots n$ , and it is given  $Q_{n+1}$  as new input, its output is  $P_n(Q_{n+1} | A_{1:n})$   
 a conditional probability distribution on  $A_{n+1}$  given  $Q_{n+1}$ .  
 The update algorithm for the system looks at  $P_n$  and at the true  $Q_{n+1}, A_{n+1}$  and creates the new  $P_{n+1}$  distribution function. <sup>(3) Factor</sup>  
 More precisely, it may also look at many earlier  $P_i$  functions and earlier  $Q_i, A_i$ 's.  
~~IA has its own update algorithm. This algorithm embodies most of the information in a IA, and characterizes that IA.~~

Apart from initialization, the update algorithm is a complex description of an IA. Though update algorithms of different IA's can differ considerably, we will see that many of them have very similar features.

Our approach to the solution of the inductive inference problem is to take a small set of IA's as it is updated on  
 A very important characteristic of a useful IA, is that ~~it has been trained~~  
 a suitable, long sequence of QA pairs (its "training seq")  
 In many QA pairs, its accuracy becomes better and better — just as a student becomes more skilled with more training. We will use this characteristic of IA's to get ~~some~~ <sup>an</sup> idea of ~~how~~ <sup>to</sup> train them to usefully work on the problem of self-improvement.

Our approach to the problem of inductive inference is to take a small set of randomly selected IA's and train them individually on sequences of QA pairs. These training sequences will differ among the IA's since ~~each~~ <sup>each</sup> IA is usually ~~adapted~~ <sup>applied</sup> to only certain domains of inquiry.

↓ Introduction.

00: 356.40: We will then combine these individually trained IA's to create a macro IA "Macro IA"

That can work on a greater range of problems than any one of its component IA's. This can be done by using a kind of linear weighting of the output probability distributions of all of the component IA's.

Among the problems that we will give to the Macro IA, is that of improving the update algorithms of each of its component IA's. Next on agenda is its improvement of the over-all system of induction.

In the foregoing discussion it may be unclear as to how we define "improvement of the over-all system". Fortunately, the theory of optimal IA's does not tell one how to obtain a better IA in any finite time, it does tell us which IA is better than another - so "improvement" is a well defined problem.

~~be is a well defined problem~~

"improvement" is a well defined problem.

~~Section 2~~ Summary contents:

Section 3 will describe a ~~recovery~~ <sup>locally optimal</sup> IA. Such an algorithm is necessarily incomplete in limitations; we will also describe a computable approximation to it.

We will also discuss the use of Genetic Algorithms and Neural nets as IA's.

Section 5 Discusses ~~the~~ combination of several IA's to yield a Macro IA. We talk how to combine several IA's to yield a Macro IA that is more versatile than any of its components. The combination of the parameters of the combining ~~the~~ combination algorithm is ~~an~~ the first problem that we

(SW) May be have a FAQ at end of report?

Section 4 Discusses training ~~IA's~~ IA's! How a sequence of well designed PA pairs can increase the problem solving skill of an IA.

(SW) (8/24/01) On update Algos for GA, ANN: They necessarily involve loop - by some external loop device. Tho, since its an ~~alg~~ alg, it can be regarded as a "problem" to be solved by the GA or ANN itself. If so, the initialization of the GA/ANN must be set early by hand! Not exactly!

In the early TQS's for these, the initial state is setup by user for each problem.

The system then trains on ~~problem~~ this set of (problem desc / initialize params) pairs. This "initial training" can be done by a different system - it can be done by user, trainer or AZ; or ANN can work on GA or vice versa.

ID

A paper on how to find good param settings for Evol. Algms.  
[E3 Evol. Comm. Vol 5 #2 April 2001 p 129]

1 PM  
30 min

00: 357.40 In GAs  $I_i$  main initial params are pop. size, Mut/cross algms., init. (population).  
T. Simplot Mut/cross can be obtained from Trauma's given values; Koze had some forms of functions.  
Unclear as to just what was done; T. "initial population" can be improved by studying the  
the generations for various problems ( $\equiv$  fitness functions). Since it is an Optm Problem,  $I_i$  and  $I_{i+1}$   
I guess  $N$  may be similar.

Base for h  
1230  
Friday  
1 PM R. d. m  
7 87  
893  
9170  
3566  
Britton/Arum  
Probst  
Mackay/opp.  
Dart  
Copley/Sy.

[SN] A way to do ANL: First just get parms to get from  $Q_i$  to  $A_i$  individually so for  
each  $Q_i, A_i$  one has a best soln. (this works only if we use a large random nos. for variables)  
or each  $Q_i, A_i$  return of 10 or 20 problems w.  $t_i$  same "n".

So we have these triplets  $(Q_i, A_i, P_{opt})$ ! We then try to find simple path from  $Q_i$  to  $P_{opt}$  for all  $i$ .

Presumably this is easier than finding a simple path from  $Q_i \rightarrow A_i$  for all  $i$ .

T. problem of .00-.04 is interesting, perhaps very imp! Solns. mean that GA w/o ANN

could be ~~valuable~~ (or be close to) a final TM. ANN is more questionable, since it doesn't  
seem to deal w. digital data properly.  $I_i$  = initialization

So: consider a sequence of  $Q_i, I_i$  pairs. ( $I_i$  = Mut/cross algms, initial pop.,  $P_{opt}$ )

This can be worked on at 2 levels: (a) Given  $Q_i$  to get d.f. for  $I_i$ .

(b) Given  $Q_i$  to get as "good as possl." an  $A_i$  in available C.

In (a) we have a problem of finding a stochastic function of hyper first given hyper set composed of  $I_i$ 's.  
This is a ~~search~~ problem. Best IA's normally work. Ave "separable" (3 indep)

In (b) If we assume that  $I_i$  is "separable" (2 components of  $I_i$ ) (mut/cross; init. PP.) over indep.

(possibly not true but usually not very false). We can optimize them separately.

Give for ~~mut/cross~~ mut pop is hyper to mean "fitness" (as defined by  $Q_i$ )  
Give for mut/cross is unclear.

Actually, during Normal GA such, we could be ~~looking~~ trying new mut/cross methods & keep

score on them, keeping good ones. One way is to "Mut/cross" i. "Mut/cross" i.

If .23 is rather than  $I_i$  Mut/cross will (perhaps) slowly vary during  $I_i$  such as  $I_i$  mean  
fitness of the population  $f$ .

Anyway, if .23 is used, then, for data, we can have not  $I_i$  initial mut/cross as  $0.04$ ,  
but  $I_i$  final mut/cross  $0.04$  can be paired with  $Q_i$ 's. — so we want TM to be able to  
go quickly from  $Q_i$  to  $I_i$  (close to)  $I_i$  final set of mut/cross "cross/mute" used.

So in my model, we try to find  $Q_i$   $\rightarrow$  init pop., Mut/cross max/min (Fitness/conv?)

Also TM has to also run fast quickly from  $Q_i$  to  $I_i$  final set of cross/muts used for  $Q_i$ 's join.

(NB) In large problem  $Q_i$  should include "context" into that human would normally have available  
to help him do initialization.

...: 358.40! 358.13 - .40 is probably enough for to report! Going deeper would be too much however.

Re 358.04 - To what extent can we do this ANN? - Well, look at the seq. of initializations.

Ass'n 358.13 If we have this set of  $Q_i$ 's, is pairs:  
Number of

$Q_i =$  connections to I/O of neurons/wts., initial wts., squashing function.

All of These can be varied during the search & / or changed at successive stages.

We end up w/ a seq. of  $Q_i$ 's "state" is that of "test set" gone.

This corresponds to variation of cross/wts during GA (358.28-30)

We end up w/ a sequence of  $Q_i$ 's corresponding to each  $Q_i$  will be a final set of connections

& wts. that were "Best" - "Best" being w/ the  $\{ \text{Test \& Data Set} \}$ .

We may want to retrain use  $\geq 1$  such "Best" soln.

If we do ANN'S w/ a few wts as poss., we will use a Max Probabest criterion (usually)

mult by the ~~empirical~~  $p_c$  of the soln. (Empirical Max Likelihood Estimate)

The goal of the  $Q_i \rightarrow$  initialization function, is to generate a "initialization of by GA, or our Best can be converted to a set of wts that GA can use very quickly".

Perhaps the "very quickly" means that v.s. Non-linear Optm (Marquand-?) method.

perhaps a few jumps from the peak.

A possible goal like 13 could be a purely discrete goal: Part of ~~the~~ a certain no. of Neurons,

connections. If we are doing only a few wts than v.s. discrete optm. may quickly lead to a Peak. (The Best peak could be much "influenced" by initial ~~weights~~ "Approx" wts.)

8/31/01 In both GA & ANN: The "problem domain"  $Q_i$ , can have various amounts of ~~complexity~~ of ~~structure~~ of ~~structure~~

"pre processing" by user: One most extreme is ~~voice~~ ~~English~~ or ~~text~~ ~~input~~ ~~for~~ ~~humans~~. ~~pre processing~~ would

(Contextual info has to be added in by user if not already present.) ~~from~~ ~~from~~ ~~from~~

"abstract" the problem for TM, pick out relevant variables; perhaps help w/ functional forms of soln, etc.

Normally, we are somewhere betw. these 2 extremes. (Also a "hint" channel could be added!)

In describing GA IA; first ~~step~~ tell how GA's normally work; then tell how how they had to be modified to become a ~~GA~~. - Then go into a bit of detail on ~~how~~ ~~how~~ ~~how~~

how this modification can be done.

Do similarly for ANN.

9/1/01 For AZ just do SAZ (3 input). (Maybe mention the 2 input PAZ "important")

Perhaps, in describing how we compute some of the function: 354.28-40 (354.20 for recursive function)

The definition of a complicated function proceeds by is affected by defining a sequence of functions, each ~~definition~~ ~~definition~~ ~~definition~~ is allowed to use only ~~primitive~~ ~~primitive~~ ~~primitive~~ functions

defined before it as well as a certain set of "primitive functions"

$$+, -, x_i, \frac{1}{x}, x, y, x, y, +, x, y, -, x, y, 2, 1, x, x, x, x, F, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100$$

For small integers  $< 100$  say,  $p_c$  of  $\frac{1}{n^2}$  is not bad.  $\sum_{k=1}^{\infty} \frac{1}{k^2} > 1.5180$   $\frac{1}{3}?$

LD

Comments on "Grammatical Evolution" (Ryan & O'Neill) - I.E. etc on Ev. Comp August 2001 P349

The start w. a BNF style Grammar. They use uniform def. over choices to generate objects.

Depending on Grammar, R. objects can be functions or Multiline. p.p.s or - whatever.

Since it's crossing they use ~~less~~ than (say 24 bytes) they only get 24 decisions per chromosome - so if it's object needs more decisions they use c. 24 bytes in circle!

They consider certain off. decisions in the grammar, but it's not apparently reasonable way ~~to~~ - (d. constr. depends on decisions to various/moduli. <sup>disturb</sup>, but 2 → 2+1 mutation

(always moves all decisions 1 choice unit in some direction)

So they can end up w. ~~objects~~ very large objects (many decisions) having always limited int. context - limited by no. of bytes in chromosome.

To do LS, start w. chromosome of length 1, then 2 then 3, etc.

From 2 to 3 the things they do, I get the impression that the authors were not ~~unhappy~~ very

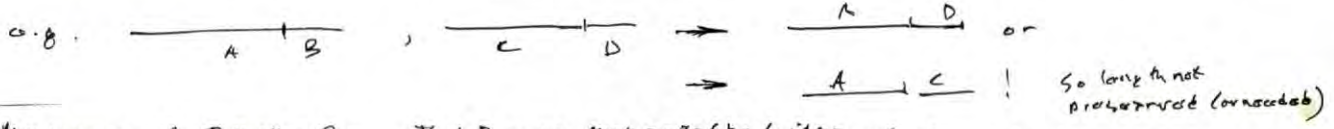
Skill in Math!

Anyway, mutation in individual bytes could make small or very large changes in phenotype -

A small change would be + → - or ~~1~~ 1 → X; a larger change 1 → +

would change no. of ops, i.e. change all decisions in rest of phenotype.

The effect of most crossovers would be to completely lose any meaning of decisions excluded!



My impression of O'Neill & Ryan: That they are much excited by the idea of generating legal expressions by Grammars - by the broad scope of ~~domains~~ Domains in which it will work. They pretty much

would like the actual search operation to be separated from the <sup>creation part</sup> construction part.

but they are excited by some "artificial" side effects of their coding method (e.g. reuse of <sup>chromosome</sup> bytes)

for ~~reusing decisions~~ coding needing many decisions.

O & R suggest get better results than Koza in "2 outputs" problems. Why this is true is not clear. The symbolic regression is interpreted as essentially the same problem, yet they

get wildly different results for ~~some~~ Koza on these 2 problems!



$\frac{6}{(\pi n)^2} = \text{prob of } n$

stochastic Alg. method  $\geq 1$ .

.00: 361.40: Perhaps write a BNF Grammar for Function generation: It may be impossible!  
ie. it may not be CFG - but it is uniquely parsable. <sup>possible.</sup> A <sup>complex</sup> function can be parsed into sub functions in  $> 1$  way, hvr.

I may (eventually) be able to find a way to describe SAZ as a variant of BNF, but this ~~will~~ <sup>is</sup> not likely to be a bit of ~~time~~ <sup>time</sup>. ~~Drop~~ <sup>Drop</sup> for now.

Since 364.37-.39 for a "Grammar" to generate functions

Do give several revealing examples.

.07: 361.36 The language SAZ begins with the set of primitive unary, binary and possibly <sup>ternary</sup> functions. Function zero defines <sup>sequentially</sup> compositions of primitives and previously defined functions. An economical way to describe a large complex function, ~~is~~ <sup>is</sup> to first describe various ~~subfunctions~~ <sup>simpler functions</sup>, ~~subfunctions~~ <sup>that are</sup> important ~~subfunctions~~ <sup>components</sup> in its description of the final function. In general, it is not economic to define a function unless it is used at least 2 times. In some cases, a function ~~must~~ <sup>will</sup> have to occur  $\geq$  rather than  $>$

.13 further, better if ~~it becomes~~ <sup>it becomes</sup> "economic" to do this.

.14 Suppose our primitives are, ~~some~~ <sup>some</sup> add, sub, ~~mult~~ <sup>mult</sup>, ~~div~~ <sup>div</sup>, ~~sin~~ <sup>sin</sup>, ~~sin<sup>-1</sup>~~ <sup>sin<sup>-1</sup></sup>.  
~~we want to do find~~ <sup>we want to do find</sup>  $F_1(x,y) = (x+y) \sin(x)$ , using ~~Reverse Polish Notation~~.

.16 Then, ~~we~~ <sup>we</sup>  $F_2(x,y,z) = F_1(x,y) \cdot F_1(x,z) \cdot \text{mult}$   
.17 ~~using Reverse Polish~~ <sup>using Reverse Polish</sup> we would write  $F_1, 2, X, X, \text{mult}, Y, Y, \text{mult}, X, \text{sin}, \text{mul}$   
followed by  $F_2, 3, X, Y, F_1, X, Z, F_1, \text{mult}$

.18 When we write  $F_1, 2$ , the "2" indicates the number of arguments  $F_1$  has.  
.19 Similarly in " $F_2, 3$ ", "3" is the number of arguments.  $F_2$  has.  
.20

To evaluate the probability of the description of .17 and .18:

First we write a sequence of all legal symbols in ~~order~~ <sup>order</sup> ~~of~~ <sup>of</sup> ~~the~~ <sup>the</sup> ~~corpus~~ <sup>corpus</sup>.

.23  $F, \text{add}, \text{sub}, \text{mul}, \text{div}, \text{sin}, \text{sin}^{-1}$   
we will call .23 the "corpus". Symbols following it will be the "corpus".

The first symbols in the corpus are (from .17)  $F_1, 2$ . ~~Therefore~~ <sup>Therefore</sup> since we are defining a function, the first symbol must be  $F$ . Its probability is 1. Since this must be the name of the first function, the probability of  $F$  is also 1. The next symbol, 2, denotes the number of arguments. There can be any ~~non-negative~~ <sup>non-negative</sup> integer. ~~Therefore~~ <sup>Therefore</sup> a zero argument function ~~is~~ <sup>is</sup> ~~not~~ <sup>not</sup> ~~possible~~ <sup>possible</sup>. Such as ~~the~~ <sup>the</sup>  $\sqrt{2}$  or ~~the~~ <sup>the</sup>  $\pi$ , which would be defined by suitable arguments ~~is~~ <sup>is</sup> ~~not~~ <sup>not</sup> ~~possible~~ <sup>possible</sup>. ~~Therefore~~ <sup>Therefore</sup> we will associate with

.30  $F_1, 2$  the probability  $\frac{6}{((n+1)\pi)^2}$ . The integer  $n$  gets probability  $\frac{6}{((n+1)\pi)^2}$ . We note that  $\sum_{n=1}^{\infty} \frac{6}{((n+1)\pi)^2} = \frac{6}{\pi^2}$ , which gives our normalized constant.  
.36

N.B. remember this using only  $\pi$  to be close by  $\sqrt{2}$ ! ~~Therefore~~ <sup>Therefore</sup> ~~the~~ <sup>the</sup> ~~probability~~ <sup>probability</sup> ~~is~~ <sup>is</sup> ~~not~~ <sup>not</sup> ~~possible~~ <sup>possible</sup>.

NB. In .16-.23 change  $x \rightarrow x_1, y \rightarrow x_2, z \rightarrow x_3$ .

In ~~the~~ <sup>the</sup> ~~corpus~~ <sup>corpus</sup> (in ~~fact~~ <sup>fact</sup>, generated by code) is a "PSM" ("State No. of steps + function has. - Also, if reduced down it might be (so at least pushdown store is needed) + state changes.





... 0 1 2 entry  
 $t = x \pm 1, \phi, \sin, \sin^{-1}$  functions, 2 2 +

00:363.40: In evaluating  $F_{5, \phi}$ ; we just evaluate  $F_5$ ; we know from properties of  $F_5$  that it is only  $\phi$ .  
 The definition of  $F_{5, \phi}$  holds as that it has no  $x_i$ 's in it.

Here when we write  $F_{5, 2}$  e.g. and develop a defn, we use at most 2 variables out  
 defn. — but we may use only 1 for now! If only 1 variable occurs, say  $X_2$ , then

$F_6(X_1; X_2)$  will be a function of  $X_2$  only in the way that  $X_2$  appears in order  $n$ .

That defn of  $F_{6, n}$ ; every time  $x$  occurs, it is followed by an integer from 1 to  $n$  w. = probability  
 the prob of " $x$ " occurring (perhaps) is a no. of times it occurred in past.

So defn  $F_{7, 3}$ ; The 3 is <sup>important</sup> ~~not~~ (see 363.33, 38)

We start with some function  $F$  say  $F_5$ ; which has arity 2. We put just 3, 1 outstack. 3, 1 / 2, 1

We write  $F_{5, 3}$ , we choose a function; say  $F_{2, 1}$ ; we write 2, 1 on stack

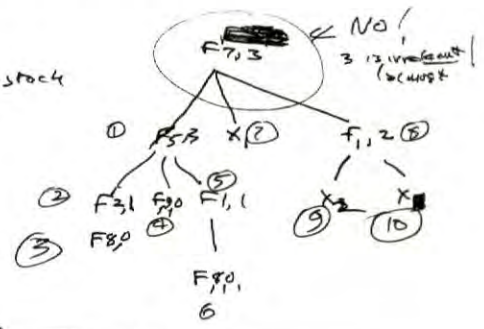
We choose  $F_{3, 0}$  (constant) write  $F_{3, 0}$ , then pop stack, put  $F_{2, 1}$  on stack.

" "  $F_{9, 0}$  " "  $F_{9, 0}$  " " " " of 2, 2 then look at stack.

we choose  $F_{1, 1}$  put 1, 1 on stack.

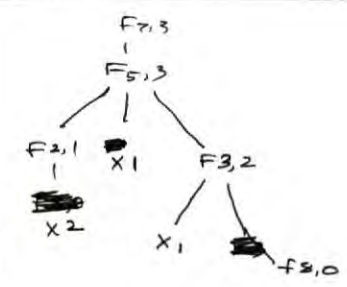
choose  $F_{8, 0}$  (const), write  $F_{8, 0}$  pop stack  $\rightarrow$  3, 2 i.e. pop stack

$F_{7, 3}; F_{5, 3}; F_{2, 1}$	3, 1: 3, 1
$F_{7, 3}; F_{5, 3}; F_{2, 1}; F_{3, 0}$	3, 1: 2, 2
$F_{9, 0}$	3, 1



- $F_{7, 3}$  : ~~3, 1~~
- $F_{5, 3}$  : 3, 1
- $F_{2, 1}$  : 1, 1
- $X_2$  :  $(\text{pop } 3, 1) \rightarrow 3, 1 + 1 = 3, 2$
- $X_1$  :  $3, 2 \rightarrow 3, 2 + 1 = 3, 3$
- $F_{3, 2}$  :  $3, 3 / 2, 1$
- $X_2$  :  $3, 3 / 2, 2$
- $F_{8, 0}$  : pop 3, 3 | pop end.

only  $X_2$  and  $f_n, \phi$   
 can cause pops.  
 If 2 pop gives  
 $n, m$  on stack and  
 $n > m$  then  $m \rightarrow m+1$   
 and new function  
 if  $m = m$ , pop stack early  
 (if  $X_2$  or  $f_n, \phi$ )  
 pop until  $n > m$   
 then  $m \rightarrow m+1$  and  
 push new function — if entry  $n' : n'$



If  $n' = 0$  or finite  $x_i$ , pop stack,  
 if  $n' > 1$  write  $n', 1$  on stack.

If  $X_2$  or  $f_i, \phi$ , pop until  $n > m$   
 If  $n > m$ , then no need to pop.  
 then  $m \rightarrow m+1$ , write new function  $f_i, n'$   
 push  $n', 1$  on stack if  $n' = 0$   
 if  $n' > 0$ , push  $n', 1$  on stack go to  $\beta$

- $F_{7, 3}$  —
- $F_{5, 3}$  3, 1
- $F_{2, 1}$  1, 1
- $X_2$  pop to 3, 1 then 3, 2
- $X_1$  3, 2  $\rightarrow$  3, 3
- $F_{3, 2}$  3, 3 / 2, 1
- $X_2$  3, 3 / 2, 2
- $F_{8, 0}$  pop all, end.

Write new function  $F_j, n$  or  $X_i$   
 If written function is  $X_i$  or  $f_j, \phi$  then pop until  $n > m$  also function  $f_j, n'$  push  $n', 1$   
 if  $n' = 0$  push  $n', 1$   
 if  $n' > 0$  push  $n', 1$

Go to  $\alpha$   
 Write new function  $F_j, n$  or  $X_i$   
 If new function  $F_j, 0$  or  $X_i$  then pop until  $n > m$ ;  $m \rightarrow m+1$  also push  $n', 1$   
 If stack is exhausted,  
 END

37  
 34

.00: 364.40: T. p. 8 of 364.37-39 is probly final I should clarify it. (Dont worry, not for me) for myself, it may be that I'll give less detail into report.

.02 A priori probability is assigned to functions such as  $\phi_n$ , in the following manner:

We start out with a list of primitive functions and primitive constants.

might be sum, minus, mult, div, sin, ~~exp~~, The constants could be  $0, 1, \pi, e, \dots$

We will use Polish notation to define functions: ~~sum, 1, 2~~

So the evaluation of sum, 1, 2 is 3.

To define a simple function we will write its name followed by its "arity", the number of arguments it has.

[FN: I am using "arity" of a function to be the number of arguments it has.]

followed by a sequence of functions and constants defining the function:

F1, sum, mul,  $x_1, x_2, x_3$ , MUL,  $x_1, x_2, x_3$  Defines  $F(x_1, x_2, x_3) = x_1 * x_2 + x_3$

To define more complex functions compactly, we will define a sequence of functions and constants, commencing in the complex function we want to define.

Each function or constant can refer only to primitive functions and constants that are either primitive, or have occurred previously in the definition sequence of definitions.

At each point in a definition string, only certain symbols are legal. Except for special cases that we will discuss later, the probability associated with a particular symbol,  $\alpha$ , appearing at a particular point is  $(n+1) / (m + \sum_{i=1}^n n_i)$

Here  $m$  is the number of different symbol types that occur at that point.

$n$  is the number times  $\alpha$  has occurred previously in the definition sequence.

$n_i$  is the number of times the  $i$ th legal symbol type for this point has occurred previously in the definition sequence.

This probability assignment method is very close to Laplace's rule.

There are two kinds of symbols that have probabilities assigned in a different way.

If  $r$  is the arity of a function being defined, it is assigned a probability of  $c(r)$ .  $c(r)$  is an assignment of probabilities to non-negative integers.

The form of this function should depend on how it is being used. In the present case,

$c(n) = 6/\pi^2$ .  $6/\pi^2$  is a normalization constant.

In functional definitions, a positive integer must follow each  $x$ . This integer will be between 1 and the arity of the function being defined. We will give each of these integers a probability of  $n^{-1}$ .

In definitions of functions,  $x_i$  will represent the  $i$ th argument of the function.

If the functions of arity  $n$ , we will assign each possible integer,  $i$ , probability of  $n^{-1}$ .

Suppose our primitive functions are sum, mult, and the primitive constants are  $\phi, 1$ .

Consider the probability associated with the symbols in the following function definition:

F1, sum, mul,  $x_1, x_2, x_3$ , F2, mul,  $x_1, x_2$ , F1, sum,  $x_1, x_2, x_3$ ,  $\phi, 1, \Delta$

Here  $\Delta$  is a symbol that has occurred only once.

Since  $\phi$  is always the first function defined, its probability is 1.

2 has probability  $c(2) = 6 \cdot (2\pi)^{-2}$

Consider the next point, the only legal symbols are sum, mult,  $\phi, 1$ . Since  $\phi, 1$  have occurred before

The probability of mult is  $(0+1)/(0+1+5) = 1/5$

The next symbol is  $x$ . The legal symbols at this point are sum, mult,  $1, \phi, x$ . Only mult has occurred

So  $x$  has probability  $1/(1+5) = 1/6$

IV

# REPORT 00-12

00: 365-40: The 1 following X has probability  $\frac{1}{3}$   
 The next X has probability  $(1+1)/(1+5) = 1/3$ , since the total symbols are 3 sum, mult,  $\phi$ , 1, X  
 The following 2 has probability  $\frac{1}{2}$

Since the definition of F1 is completed, the next symbol can be only ~~the~~ ~~next~~ ~~one~~  
 "end" symbol,  $\Delta$ , or  $F_2$ , the beginning of definition of the new function.

So  $F_2$  gets probability  $\frac{1}{2}$  (2-1)/2 = 1/2.

3 has probability  $C(3) = 6 \cdot (3\pi)^{-2}$

All after F1 position, the total possibilities are F1, sum, mult,  $\phi$ , 1, X, F1, so the probability

of F1 is  $(1+1)/(1+1+3) = 2/5$

The next position the total symbols are sum, mult,  $\phi$ , 1, X, F1

The probability of "sum" is  $1/(4+6) = 1/10$

The probability of X is  $\frac{1}{4}$ . At the next "X" position sum, mult,  $\phi$ , 1, X, F1 are all total 2

The probability of K is  $(1+2)/(2+1+2+6) = 1/4$

etc: the probability of  $\Delta$  at the end is  $1/3$

Discussion: Note (03-12)R: Perhaps here adjustable param, "2" so probability of  $\Delta$  is 2. Expected no. of functions in a string is  $\frac{1}{2}$ .

If we had a pc of "F" at the position: pc of  $\Delta$  would be

$\frac{1}{2} ; \frac{1}{3} ; \frac{1}{4} \dots \frac{1}{2} \cdot \frac{2}{3} \cdot \frac{3}{4} \cdot \frac{4}{5} \dots \frac{n-1}{n} \cdot \frac{1}{n}$  so the pc of  $\Delta$  is  $\frac{1}{n}$

sequence. If we allow "n" no. of previous F's to include all uses of F, the probability will be  $\frac{1}{n}$ .

Not clear how we can prevent this! We can do it Ad hoc: Have special way of computing pc of

new data. (Modifiable by the system of course). use a low 2 param. d.f.  $\rightarrow 36$

Another point: when we ask for a finite function, we may get any  $n < \infty$  perhaps  $\infty$ .

"arity" could be specified after the function is defined. The events "X" is X1, Y. next

has perhaps uniform d.f. over X1 & X2; After X1 has been defined, the next X has

uniform d.f. over X1 to X(n+1).

However, the "arity" of .19 only occurs when we are generating functions Monte Carlo-wise

Otherwise, its post-hoc assignment of pc's to functions of known arity.

Actually, .19 would probably be fine! We may not need to specify Arity when we first define a function. It is some my thesis discovered after the function is defined. I. me that I'd been using upto .19 was not so good in this case! If arity n was specified, a Monte Carlo

generation of a function could yield any arity  $n < \infty$ .  $\phi$  to  $n$  (!).

In fact, the mod of .20 ff, the arity is found after the function is generated. After

generation, the name, arity & case count are kept in the function's "properties list"

So, in .20 ff we want to keep track of the "arity" of the function in its definition of a function! This will, when the fun is completed, give its "arity".

This gives us at least 2 parameters to be adjusted in a system:

- 1 if n X types have been ~~used~~ ~~named~~ this tells what do we want the pc of a X(n) to be?  $\rightarrow (366.19 = 30)$
- 2 If m new functions have been defined this far, what pc do we want for a new function?  $\rightarrow ((366.08 = 10)R)$
- 3 the function C(n) of 365.24  $\rightarrow (366.13 = 18)$

In disc. of .13-.18: that say, instead of sequence  $(\frac{1}{2}, \frac{2}{3}, \frac{3}{4}, \dots, \frac{n-1}{n}, \frac{1}{n})$  we list

a series converging to  $\prod_{n=2}^{\infty} (1 - \frac{1}{n^2}) \cdot f(n) : \prod_{n=2}^{\infty} (1 - \frac{1}{n^2}) \approx e^{-2.68}$

if  $f(n) = \frac{1}{n^2 + \epsilon} \leq f(n) \leq e^{-2.68}$  so  $e^{2.68} \cdot \frac{1}{n^2 + \epsilon}$  will diverge. so  $\epsilon > 0$  may be needed in more study

converging series of  $f(n)$  (Two it should be poss. to make  $f(n)$  & slightly after  $\prod_{n=2}^{\infty} \frac{1}{n^2}$  so we could

have  $\prod_{n=2}^{\infty} (1 - \frac{1}{n^2})$  be any slowly  $\downarrow$  function of n. (say slower & even  $\frac{1}{n}$ )  $\rightarrow$  in fact, one can define  $f(n) \rightarrow 367.00$

Looks like pc of  $\Delta$  is always  $\frac{1}{n}$ !  
 F1 &  $\Delta$  never occur both!  
 This means long strings of functions very unlikely!

So it does converge; if we  $\epsilon$  much, it will not average  $\rightarrow \frac{1}{n}$ .

*[Handwritten scribbles and notes on the right margin]*

Prob. is the probability of a sequence of  $n$  definitions, then step:

so  $\prod_{i=1}^n (1-f_i) \cdot f_n \equiv g(n)$ ; where  $g(z)$  is any function of  $z$ .

The main point is!

We can choose  $g(n)$  & hence  $f(n)$  as we like. The sequence converges as rapidly or as slowly as we like.

$\frac{f(n)}{f(n-1)} = \frac{g(n)}{g(n-1)}$ ;  $\prod_{i=1}^{n-1} (1-f_i) = \frac{g(n-1)}{f(n-1)} = \frac{g(n)}{f(n)(1-f_n)}$

so  $\frac{f(n)(1-f_n)}{f(n-1)} = \frac{g(n)}{g(n-1)}$ ;  $f(n)(1-f_n) = f(n-1) \frac{g(n)}{g(n-1)}$

So, given  $g(n)$  we can find  $f(n)$  or we can solve for  $f(n)$  as a function of  $n$ .

$x-x^2 = \alpha$ ;  $x^2+x-x = \alpha$ ;  $x = \frac{-1 \pm \sqrt{1+4\alpha}}{2}$ ;  $x = \alpha$ ; if we want  $f(n) > 0$

SN One of the IAS candidates. TRAINER! But we want that not to be used much - make it rather expensive & a lot of it as TM will be used. (Actually, Trainer will tend to be slow)

$f = \frac{-\alpha \pm \sqrt{\alpha^2 + 4}}{2}$   
 $= -\frac{\alpha}{2} \pm \left(\frac{\alpha}{2}\right)^2 + 1$

Corrections needed in 365.03 ff! (in view of 366.32...35)

365.08 let us not follow by any of no. of ...

365.10 (some 2 copy)

21 two kinds of ~~...~~ labels 20's!

23-26 on CCF =  $\frac{f_2}{f_1}(n+1) \rightarrow$  omit

27-32 Also omit, but  $\rightarrow$  replace  $\rightarrow$  w. M! ~~...~~ 367.30 + 39

34 omit 2 and 3.  
 36 omit this line.

I want to rewrite all of 365.03 ff using subscripts on  $x$  & on  $f$ . (4th time easy to

pub in LaTeX: just under line.

Also in rewriting perhaps write down strings  $\rightarrow$  vertical columns: 11 columns; subscripts & symbols, probably.

(like  $(2^+)$  /  $(2, 3, 4, 5, 6)$ )

Maybe try for better function to ~~...~~ define.

prob. of  $F_2$  occurring & previous has been defined, is  $\frac{1}{2}$  no. of uses including definition. There should be "1" added in numerator or denominator!  $\rightarrow$  perhaps write correctly later

We could fix this by using a "prob. copy" & c. "straight rule".

In (368.15 ff) we used  $8.8 \times 10^{-8}$  a pc. of function.

Now, in any use of this function, we can use  $\frac{1}{2}$  & c. same definition but argument permitted!

So in 368.15 ff we have 11 pc of  $F_2$  - a func. w. 3 args. so 6 permutations.

$F_1$  has only 2 permutations, mul, sum later  $\rightarrow$   $8.8 \times 12 = 105.6$  so  $\rightarrow 105.6 \times 10^{-8} = 1.056 \times 10^{-6}$ .

So maybe  $F_1$  use ones, mul, sum, use twice.

Actually  $F_1 \equiv \text{mul}()$  not so interesting!

Try  $F_1(x, y) = x^y + x$ ;  $F_2 = F_1(x, y) \rightarrow y^2 + y^2$

$F_2$   $\rightarrow$   $F_1(x_1, x_2, x_3, \Delta)$  perhaps use sum, mul,  $\phi$ !  $\rightarrow$  permutation precepts

$F_1$  sum mul  $x_1, x_2, x_3, x_4$   $\rightarrow$  sum mul commut.  $\rightarrow$   $x_9$  term,  $x_2$  for  $F_1$ ,  $x_6$  for  $F_3$   
 $8 \times 2 \times 6 = 96$   $\rightarrow$  368.18

30 (to replace 365.27-32) Whenever the symbol "X" occurs in a function definition, it will have a position int. subscript that is a positive integer. The first time an "X" occurs in a function definition, this subscript must be 1. Henceforth, the  $n$ th or  $m$ th X in a function definition can be any integer from 1 to  $n+1$ ; where  $n$  is the largest X subscript that has occurred thus far in that function definition.

34 Cannot have integers & has probability  $(1/n)^{n-1}$ .  
 When the function name  $F_1$  occurs, it is the only possible symbol at that point. It is given probability 1.  
 For function names  $F_i$  with  $i > 1$ , the probability of  $F_i$  is  $\frac{(i-1)^{i-1}}{i}$ . The only other  $(i-1)$  symbol at these points is the step symbol  $\Delta$ , which has probability  $1/i$ . The step symbol indicates that the sequence of function definitions is completed. It always occurs at the end of the sequence of function definitions.

noisy  
hard

symbols	legal possibilities	probability	
$F_1$	$F_1$	1	$\frac{1}{2}$
mul	sum, mul <del>(<math>\phi</math>, 1)</del> illegal	<del><math>(0+1)/(0+0)</math></del> $(0+1)/(0+0+0)$	$\frac{1}{2}$
$X$	sum, mul, $\phi$ , 1, $X$	<del><math>(0+1)/(1+0)</math></del> $(0+1)/(0+1+0+0+0+5)$	$\frac{1}{4}$
subscript 1	sum, mul, $\phi$ , 1, $X$	<del><math>(1+1)/(2+0)</math></del> $(1+1)/(0+1+0+0+1+5)$	$\frac{2}{7}$
subscript 2	<del><math>F_2</math></del>	<del><math>(2-1)/2</math></del> $\frac{1}{2}$	$\frac{1}{2}$
$F_1$	$F_2, \Delta$	$(2-1)/2$	$\frac{1}{2}$
sum	sum, mul, <del>(<math>\phi</math>, 1)</del> , $F_1$	$(1+1)/(0+1+0+1+3)$	$\frac{2}{5}$
$X$	sum, mul, $\phi$ , 1, $F_1, X$	<del><math>(0+1)/(0+1+0+0+1+3+6)</math></del> $(0+1)/(0+1+0+0+2+2+6)$	$\frac{1}{11}$
subscript 1	sum, mul, $\phi$ , 1, $F_1, X$	$(2+1)/(1+1+0+0+2+2+6)$	$\frac{3}{12} = \frac{1}{4}$
subscript 2	sum, mul, $\phi$ , 1, $F_1, X$	$(1+1)/(1+1+0+0+2+3+6)$	$\frac{2}{13}$
subscript 3	sum, mul, $\phi$ , 1, $F_1, X$	$(3+1)/(2+1+0+0+2+3+6)$	$\frac{4}{14} = \frac{2}{7}$
$X$	sum, mul, $\phi$ , 1, $F_1, X$	$1/2$	$\frac{5}{15} = \frac{1}{3}$
subscript 3	sum, mul, $\phi$ , 1, $F_1, X$	$(4+1)/(2+1+0+0+2+4+6)$	$\frac{1}{3}$
$\Delta$	$F_3, \Delta$ subscript	$1/3$	$\frac{1}{3}$

product =  $2.8 \times 10^{-8}$

**Multiplicity  
redundancy  
domination**

$$\sum_{X_1, X_2} \text{Mul } K_1 K_2 X_1 X_2 = \sum_{X_1, X_2} \text{Mul } X_2 X_1 X_1 X_2 = \sum_{X_1, X_2} X_1 X_2 X_2 X_1 = \sum_{X_1, X_2} X_2 X_1 X_1 X_2$$

Also, if a function has an arity  $N$  domain, we have  $N!$  functions like it. That could be used where ever it is used, & give a distinct coding. However, various symmetries (like commutativity) can reduce (or multiply) this redundancy. Associativity might also give equivalent expressions. A  $\phi$  function could  $\uparrow$  PC of a function by a multiplicative factor.

Also, any subfunction used in a defn has its PC  $\uparrow$  by its (arity!). Whether this occurs each time this funct. is used in a defn. is not clear. **I really have to look into all of this!**

Have space to STACK for ARITY of function over working area - Tho it may be unnecessary. Then needed info may be already on disk in e. form of e. "n, m" pairs (n = arity).

Summary of work on expansion of how to compute PC of a function;

- 361.24 Derbs IA's in general. (Also note 347.00)
- 362.07 Derbs & recursive SAZ (composition); how PC's are assigned to terms.
- 363.12
- 363.16 - 364.39: Develops pgn. to General functions ( $\neq G.F.G.$ ;  $\neq$  BNF)
- 365.02 - 366.12 Exposition: how PC are assigned to function. This has to be re-written, w. a no function not 368-01 - 15 used as example. - Also note errors in whole thing! - **367.20** (fix w. precursors)

320.28 is a pure listing of reviews of QATM.  
List of Review pp. is subjects of Reviews:

- 1) List of Induction Algms 228.00; 150-151;
- 2) EBL 287, 308, 278
- 3) "Cure Cancer" : Pwms(soln. 243.00
- 4) Lisp 293, 294
- 5) H.W.: Programmable Gate Arrays 298.12 : Old H.W. idea + see Microbox "A" ~ 1990-1992
- 6) On how to Deal w. CC in TM : 285.30 [How to put com Global Goro]

- 7) Combining prodns of IA's : 272.28 ; 273.11
- 8) Early Stochastic OPS (MOPS) : Extens , 267.35A.

9) How to use t. Formlzn of AZI to desc. Recursive Concs  
~~based on the review 267.35A but the final 2 numbers are not there recent!~~  
 44! : 354.20

10) Stochastic Operators & some probl. forms: (331.20) (340.31) | <sup>clustering</sup> 329.20 <sup>conditioned very closely</sup> very general discussion: (375.20 - 30) ~~summarizes~~

11) On context as "conditional" for PC of Concs : for Updating Algn : 300.26 is stuff surrounding it. Very imp!

12) UPDATE Algn (perfect if  $c \rightarrow \infty$ !) Given p.d. on  $O_n$ ; (see wts. for all poss  $O_n$ 's)  $F = W_i^n = w_i$  for  $O^i$   
 [  $O^i =$  set of all  $r$ 's ]  $W_i^{n+1} = W_i^{n+1} + P_{O_i}(A_{n+1} | Q_{n+1})$ . Algn. for updating  $O_n$  may not exist, but algn. for updating p.d. on  $O_n$  is simple for  $c \rightarrow \infty$ ; but may be more diff for  $c < \infty$ !

13) "Final Soln." to UPDATE PROBLEM: Claimed by 300.35 - 301.20 ; 302.10 ff  
 Hrs. what about t. very time consuming diffy of finding  $P_{O_n}(A_i | Q_i)$  if  $O_n$  is a Monte Carlo device } 375.20 - 40 discusses this.

14) Logical, Metrical reasoning used by JM to solve Induction (i.e. other) problems: 224.04

15) <sup>ONE</sup> (Main) Track for QATM project: 378.00 - 379.40 : "Needs work" But Good: general paper.

Inadequacy of present Model for a day's Common Hour .00  
possi. Adequacy of Plans for UPDATE ALGM

**SN** I had this idea that Long is context-dependant PC's of concs, with board out to get a fairly smart TM. Perhaps NOT! This "Long" that I'd considered started out by doing a trial of the Best looking Goal. (Greedy!). Then, on basis of what was (and, from this trial), doing TM that was fair but looking good. (Greedy!)  
The trouble is, it doesn't do experiments! Its first trial was not an experiment, but a trial to get the final soln.

Many hours were in the "Do it, then (on basis of what happens) do it (a) • I don't see this model discovering hours of this sort!!

One hope is that since the QA problem model can express any problem - then we should be able to get TM to learn do experiments. This seems to be an implicit aspect of it.

"Look ahead" problem - a very sophisticated case is "Core cancer" -> (243.00)  
292.15 discusses "Limited Lookahead".

When TM is looking for a good Update algm, it could use "Look ahead" or any technique.  
A problem is: in order to be able to do this usefully, we'll have to give it problems like Best in its TSO. This will not be so difficult! TM could "watch" ~~the~~ "Experiments" being done in the soln. to a problem - (perhaps this is one kind of "Hint")

**SN** The Main Problem would appear to be a reasonable O.K. Updating Alg - Good enough so

- 1) T. "coffee" Core of a MDL (the cc would be inserted automatically if Lsuch was used)
- 2) I don't know how to insert cc in to this Core (285.30) -> Normal Lsuch is one way, but GA(3) may need an explicit expression for cc in its core. use a machine of copy
- 3) T. possy of using GAW, large CB (size CJS ~ 10<sup>15</sup>) (Real/this shouldn't be necessary, it would make the design of the TSO a lot easier!) At any rate, (1) & (2) make this a real possy!
- 4) T. forgets, ignores exists assume Global "search for update using last k problems of "Backtrack" (not very logical).
- 5) A real Augmented Update algm would use  $O_{n+1} = O_n (O_n, O_{n+1}, A_{n+1})$ . -> possibly other "Context". expanded idea of what can be "Context" maybe. find Reference on this! 300.26

Obtaining this is the major problem.  
I expect to get (5) by study of TSO's. Any details of (5) such as Context dependent pc's of concs, bring during Lsuch, Look ahead; I expect to develop methods to do these by observing them being done in TSO's. This is the Main Idea (I Hope!). That I can have a variety of IA's in the system, should make it possible to include any type of Hour that I might need.

7) In (6): T. mode of "insertion" of a Hour; Try to see how the Hour might have been derived & insert a simulated TSO in the Lsuch to find it. This will (perhaps) give me "parameters" (pc's) for the Hour, so it will be properly used.

- 8) Definition of a "Heuristic": any thing that enables cheaper solns of (1), (2) - i. Global Core. sum of 0
- 9) Possy of using "Patron pc's" to get CJS ~ 10<sup>15</sup> ("patrons" are allowed to give TM occasional problems to "obscure its progress". (see 370.00ff) (378.00-379.00) Update d.f. of O<sub>i</sub> by  $a_{n+1} = a_n P_{ij} (A_{n+1}/Q_{n+1})$  then "expand" the O<sub>i</sub> d.f. (385.00 discusses (378-379) v.s. (300-302) spec: 370.00
- 10) (394.06-11 also 394.12-16): Solve Ann, Ann, then find way to reach 170 Q<sub>n+1</sub>
- 11) (378.00-379.00) v.s. (300-302)



ID

Practical HW Estimator: 11 ff

369.25 is most relevant!  
(369.25) ~~369~~

00: 369.40: Discn of 369.30: 369.25 (3) in particular: I was impressed by kozz's being able to get GA to do what appeared to be "hard problems". This could be an illusion! e.g. @k oza may be faking it.

01: (b) T. probably he solved in ch+design (e.g.) may not be so hard: Tenzor methods may not generalize to other domains... (No analysis of problems/systems as "circuits" - seems to be effective over several domains: i.e. Gabriel Kozz (Tensor Analysis of Networks): I think he did try to use his T.A. for Gen Relativity: Also his Guy from France who gave talk at IDSIA when I was there. But first, I'd have to go over Kozz's Methods w/ some care!

How New functions are created? e.g. a few, i.e. how his "Circuit Grammar" works  
Circuit Grammar.

10: Affw. Kozz doesn't seem to understand induction (GA) as well as I do - so I should be able to get GA to work more like AZ) in recognizing/discovering using sub-processes & other "sub-functions". I'll be able to get  $20 \times 25 \times 10^9$  for  $\approx 7.5k$  in a yr. from new (Sec. 12 - 21)

Also kozz @ Didnt use Machine Logic: He had Lisp which can be reprogrammed. Hur, I think to me

11: bottlenecks was SPICE - f. overall cost. I may be able to speed this up much.

112: W.  $5 \times 10^9$  ops/sec, it took him several hrs to solve a serious problem. I can certainly get  $5 \times 10^9$  ops/sec (expensively) & probably  $50 \times 10^9$  ~~for~~ ~~at~~ ~~cost~~ or less. (or 10" sec about now, Park St)

100 102 ~~mother boards~~: perhaps wood & bookshelves to "hold them up",  
First try 5 or 10 mother bds! Ask Saty to help. Note: I don't need sound or video on Bds. - Some are MB's computer, ~~Sound & video~~ - just physics for RAM, HDD & floppy.  
'Shop around' for approx price MD. Say \$100 for MB, \$40 for .25 GB ram, \$50 for 20 GB HDD  
\$100 for 1 GHz CPU cost 300 for 1 GHz! \$30k per 100 GHz. = 10". Ram 25¢, HDD = 2000¢ = 2 TB  
Look at speed of HDD for parallel L3 cache! But HW choice will depend (I think) on just what

Such Algs I use.  $\approx 2 \times 10^9$

20: Also, it appears, price (will probably)  $\div 2$  & speed available for \$100 CPU will be 2 GHz:

21: So probably  $10''$   $\frac{GHz}{\$100}$  for  $\frac{30}{1} = 7.5k$ , or  $10^{12}$  for  $\frac{75k}{1}$ ,  $10^{12}$  GHz for \$75k.

I MAY want to do not too narrowly ~~the~~ ~~algorithm~~ / PC ordering ~~not~~ AZ (but order by heuristics) since this is too greedy & not "creative": So maybe do (PC) ordering w.  $\delta < 1$ . (danger of non-convergence?)

29: Also, Examine kozz's "Circuit Grammar" w. some care. Even if I can't figure out how his Grammar works, I should be able to devise a circuit Grammar that would work. It might be useful to look at how "SPICE" ~~works~~ ~~clerk~~!   
pc's are

30: Ok. So what's new? Last was at 368.30: Exposition of how ~~assigned~~ to functions.

This seems ~~under~~ control, / this to be rewritten.

Perhaps Go on to next step. This also follows: sent 347.00 - 360 starts section on "T. induction Algs"

The flow: 361.24 - 363.03: starts by ~~the~~ saying!! we will discuss 3 of T.A's in some detail

362.07 - 363.03 and 365.02 - 366.12 covers same ground: Play Overlap much. B

So, perhaps read this carefully & rewrite. Perhaps first read earlier ~~approx~~ ~~sect~~ ~~start~~: 332.00 - 333.13  $\Rightarrow$  347.00 - 14; 361.24 - 362.04  
(362.07 - 363.03)  $\gg$  365.02 - 366.12

9/9/01

ID

$3x(1+7)$

1. list of machines to previous pp.  
362.13: Explain what we mean by "economy".  
corresponds to a "shortcode" - tends to express  
popularity in the domain. It is as good as any possible pc. to find.

332.00 - 333.13 Section Introduction Data of a. Tsp. problem. Introduction

$O_n, O_{n+1}$ ; The 2 main critical ~~equations~~ equations.

347.00 - .16 Section Induction Algorithms. General form! How we use source (IA's)  
in it  $\Rightarrow$  "Specialists"

361.24 - .36; continuation / 357.00 - .16 Section of

360.29 - .36 Section of. SAZ system. Introduction of SAZ

362.07 - .13 More details of SAZ; how complex functions are described.

Concise by Program Later.

362.14 - 363.03 A first attempt to explain ~~the~~ assumption of ~~equivalent~~ to SAZ functions.

365.02  $\rightarrow$  366.12 same 362.14 over again - it only partly ~~describes~~ analyzes functions.

T. Table of 368.00 - .15 is a much clearer way to ~~express~~ present level of ~~only~~ but I should use precursors.

Revisiting 365.02 A!

An ~~operator~~ probability is assigned to functions such as  $O_n$ , in the following manner:

We start with a list of primitive functions and primitive constants that are usable in the function definition.

"Primitive functions" could be  $x, sum, minus, mul, div, sin, sin^{-1}, \dots$

"Primitive constants" could be  $\phi, 1, \pi, e, \dots$

We will use Polish, <sup>parenthesis-free</sup> notation to define functions, so  $3x(1+7)$  is written

sum, comma, use spaces  
mul, 3, sum, 1, 7

To define a simple function, we will write a list of primitive functions and constants it might use, followed by the function definition in terms of these primitives.

$x, sum, mul, \phi, 1, F_1, sum, mul, x_1, x_2, mul, x_1, x_2, \Delta$  defines the function

$F_1(x_1, x_2) = x_1^2 + x_2^2$

The " $\Delta$ " symbol indicates that the definition sequence is completed.

For example, to define  $F_2(x_1, x_2, x_3) = x_1 \cdot F_1(x_2, x_3)$  we write the definition of  $F_1$  followed by

the definition of  $F_2$  in terms of  $F_1$ :

$x, sum, mul, \phi, 1, F_1, sum, mul, x_1, x_1, mul, x_2, x_3, F_2, mul, x_1, F_1, x_2, x_3, \Delta$

N.B. F and x all have subscripts

To define more complex functions compactly, we will define a sequence of functions and constants culminating in the complex function we want to define.

The definition of each function or constant in a complex definition can only refer to symbols that have occurred earlier in the definition sequence.

At each point in a definition string, only certain symbols are legal. Using illegal symbols results in a meaningless (uninterpretable) string.

Except for special cases which we will discuss in the next paragraph, the probability of each symbol in the definition string is given by the formula, some symbol,  $\alpha$ , occurring at a particular point in the definition string is



00: 372.40

Table 1

Symbol	Legal possible symbols at this point.	Probabilities
X	X	1
sum	sum	1
Mul	Mul	1
$\phi$	$\phi$	1
1	1	1
$F_1$	$F_1$	1
sum	sum, mul	$1/(1+1)$
Mul	X, sum, <del>mul</del> , $\phi$ , 1	$1/(1+2+(1+1))$
X	X, sum, mul, $\phi$ , 1	$1/(1+2+2+(1+1))$
subscript 1	1	1
X	X, sum, mul, $\phi$ , 1	$2/(2+2+2+(1+1))$
subscript 2	1, 2	$1/2$
Mul	X, sum, mul, $\phi$ , 1	$2/(2+2+2+(1+1))$
X	X, sum, mul, $\phi$ , 1	$3/(3+2+3+(1+1))$
subscript 2	1, 2	$1/2$
X	X, sum, mul, $\phi$ , 1	$4/(4+2+3+(1+1))$
subscript 2	1, 2, 3	$1/3$
$F_2$	$F_2, \Delta$	$(2-1)/2$
Mul	<del>sum</del> , mul, $F_1$	$3/(3+3+1)$
X	X, sum, mul, $\phi$ , 1, $F_1$	$5/(5+2+4+(1+1))$
subscript 1	1	1
$F_1$	X, sum, mul, $\phi$ , 1, $F_1$	$1/(6+2+4+1+(1+1))$
X	X, sum, mul, $\phi$ , 1, $F_1$	$6/(6+2+4+(1+1)+2)$
subscript 2	1, 2	$1/2$
X	X, sum, mul, $\phi$ , 1, $F_1$	$7/(7+2+4+1+(1+2))$
subscript 3	1, 2, 3	$1/3$
$\Delta$	$F_3, \Delta$	$1/3$

.10

.20

.30

00: 3 29.40 It will be noted that we have described a technique for assigning probabilities to functions of any finite number of variables. These are deterministic functions — but on the function of interest is a probabilistic function of one variable. There are several ways to convert deterministic functions to probabilistic functions. We will describe one very general method.

To obtain a very general stochastic function of variables, Consider all deterministic functions of 2 variables — ~~the first variable is the~~ <sup>u.c. Math</sup> The first variable is the string that is a member of an infinite prefix set, ~~the second variable is a member of an infinite prefix set~~ <sup>a string and is</sup> If we have  $Q_1$  in our "QA" machine, The second variable is a member of a prefix set (not necessarily the same prefix set for both variables). The first variable is  $Q_2$ , the input question, and the second is  $R$ , a random number. If  $F(Q_1, R)$  is allowed to be any function of two such variables, we can find such an  $F$  so that a random distribution on  $R$  will induce a distribution on  $F$  that is identical to any finitely describable stochastic function of  $Q_1$ . By repeatedly giving different random  $R$  values for a given  $Q_1$ , we get Monte Carlo distribution of  $A$  for that  $Q_1$ .

We may regard the foregoing description method as equivalent to a 3 input, universal Turing machine with unidirectional output. ~~The machine has 3 inputs and 1 output~~ <sup>standardize!</sup> The first <sup>input</sup> describes the function  $F$ , as does our ~~AZ1~~ <sup>AZ1</sup> register. <sup>what parameters used for AZ1 or SAZ or whatever.</sup> The ~~other~~ <sup>two</sup> ~~inputs~~ <sup>are</sup>  $Q_1$  and  $R$ .

Perhaps dec. 17-20 first, then say that AZ1 register is a way of implementing it. 376.32

24: 370.40 Any way: Where Am I now? → 370.30: I guess I've just finished part of SAZ on ~~some~~ <sup>some</sup> ~~of~~ <sup>of</sup> ~~the~~ <sup>the</sup> ~~functions~~ <sup>functions</sup>. a prop evaln. of funcs. next

~~SN~~ My dream ~~of~~ SAZ didn't go into its dream of Recursive Functions. ~~I doubt if this will be a bottleneck of a chrl lang data or for Exposition.~~ (354.20)

Topics: How UMC's do  $Q, R \rightarrow A$ :  
 problem: we need  $P(A|Q)$ : This Monte Carlo form does not readily give it. [375.20ff has a lot on this Q]  
 HVR, in practice, it's almost always poss. to map this "input form" into any other (as it V.D.)  
 - out size issue: Huh! — changing input into non-MC Carlo forms is not so easy!

I think I did write late about this — How to get useful pds in useful forms:  
 300.35 ff ... General pp purposes to solve the update problem using Monte Carlo L's etc. — My impress is that it would be too large! — Unless only a few new QA's had to be evaluated. This may be a trick!

Certainly all the QA's pds can't be re-evaluated for each new Q on trial! (258.24 Discusses this)  
 Perhaps notes were 390.35 was all about! Just I ~~at~~ <sup>at</sup> ~~found~~ <sup>found</sup>  
 Also Note "Perfect Update" of (368.16)! 271.31 considered that various QA's would have different formats for p.d. outputs. MC Carlo, list of A's in PP order. U.S.

On ANN

00:37.40

ANN: I only have a vague idea of ~~what~~ how I expected to get IA out of ANN. I usually don't use cross validation, I think I did once a way that was able to use it in an EA - in general it ANN can either use small no. of nodes & edges & use Max Fitness, or use lots of nodes, edges & division of data into exp. & test sets.

Small updates & modulus of wts: layers, add, delete wts, Nodes, & o/a RANN concerns.

Partners add, delete memory / change memory type. Change Ranges, Squashing function.

T. goal is to get ANN to work. to output QA outputs. We want to update algo. to lock at present state, recant QAs, latest QA, & proper modulus of Next.

Testing would be a really statistical on subset of QA's.

Some poss. diff. in this Approach: in RANN we usually ~~don't~~ are not able to localize just what is being important in a particular QA. IN AZI, this is done by noting frequencies of use of sub-functions.

Hence, to an extent, to prog. discerns a bit "premature" T. main methods of "update" will be obtained by studying human ~~transit~~ search ~~usage~~ & finding ways to express/guess them using AZI, GA, ANN or any other IA's that may be suggested. Though, to implement a learn, one usually doesn't need an entire IA (?).

On  $P_n(A_i/Q_i)$  as <sup>Universal</sup> stock op. (Other than  $\epsilon$  standard 3 input unc). 244.07 had some good ideas. Also Note 244.02-20R

ABedefghij

One way is via Stock Grammars - that are derived from Dataflow Grammars by using probabilistic rewrite rules: to find ~~unambiguous~~ a pe of a  $\epsilon$ s, ~~fast~~ parse it all! Also note 376.07: on 3 input unc w. R as "clusters" to central  $A_i$ . That's a kind of "clustering" approach to ~~inductive~~ coding.

Another way 250.09 + f, 257.37 - Discussion on 269.40: on which ~~is~~ the universal: (They may be, but still not so good ~~historically~~!) 267.35 ff

Another Note! The 3 input unc is not very slow if it's thought-for  $A_i$  is not by PC! Usually this will be the case. Also, there may be several equivalent  $A_i$ 's that have Saki's factory

I think there are several notes on counting P.D.'s on strings, act. 331.26 How to count  $\epsilon$  into Map into a Map. (unclustered) 340.31 On "clustering" as a kind of stock d.f. 21505 = 376.07

Perhaps one <sup>very</sup> ~~must~~ earlier idea: That's I observe hours in human prob. solving & devise TSO's that could have given rise to them, I will discover many ~~new~~ ~~input~~ forms of stock operators. that will depend much on insider for evolving  $\epsilon$  v.s. TM.

- 9/10/01: Some diffys in in report list for: (A in my plan for COM, in general) 1) As dec'd, ~~the~~ f-SAZ (uz is not Univl. - it didn't do recursive functs 2) T. 3 input unc for  $\epsilon$  general stock operator! Seems to take too long to evaluate -> 376.10

$P_n(A_i/Q_i)$ : In general;  $\epsilon$  3 input modal may be very slow in evaln/used. main

Replies 1) This is not diff to fix, T. p. 8 I do have a way to assign a priority to decns of ~~all~~ (all) finitely decrable functions. 2) @ to may input update schemas (unlike most) Not many  $O_n(A_i/Q_i)$  have to be evaluated for new trial  $O_n$  see 20-30 for some ideas on other forms of P.D.'s (229-30) is my main idea (expect'd so (1/hope) seems very reasonable, very practical!

(Also Note 225-228)K

00 375.40 : A Better Way to organize the report:

1) Desc. problem (as in introduction) of  $\{Q_i A_i\}_n$  input  $\rightarrow Q_{n+1}$  given  $U(O_n, Q_n, A_{n+1}) \rightarrow O_{n+1}$

2) A system that we will term IA a system that is able to do this to any extent.

Also Discuss to Gorb (cap. 1)

That for a large class of IA's  $U$  is Governable & Applicable. Updates always like  $U(Q_n, Q_{n+1}) \rightarrow O_{n+1}$   
 $U$  need not be a TA. i.e. it need not be linear  
In general,  $U$  remains not really bijectively possl.

3) There are 3 aspects of a IA that characterize its behavior: 1) initial  $O_1$ , 2) TSA 3) updates Alg.

4) Explain that int. input  $U_n$ , the  $R$  input can be regarded as "Correction" to a "Control"  $A_n$  output which occurs when  $R$  represents  $\neq 0$ . (It is not  $R \equiv 0$ , i.e. it is the representation of zero in the prefix code being used. "Clustering" is one example of this kind of coding.  $R$  value represents the "Matrix" of the Cluster)

10 375.33: 3)  $O_n(Q_n, Q_{n+1}) = d.f. \text{ of } O_{n+1}$  is incorrect:

$O_n$  (d.f. of  $O_n'$ ,  $(Q_n)_{n+1}$ )  $\rightarrow$  d.f. on  $O_{n+1}$  is correct. It is the d.f.'s of  $O_n$  that are

recursively defined. eq. 11 is easy to write out exactly for  $cc = 0$ , but not so easy to implement for  $cc > 0$ .  
Hr, it would be worth looking into as an approx for  $U$  putting d.f.  $(O_n)$  is defined or developed recursively.

I may have done 2.13 when I was considering "Min Backtrack Algs". My impression at that time was that I ~~could~~ couldn't get much out of it. With present view, however, it might be possl. to Revive it. "Revisits" it.  
Also Related is the idea of a univ. set of op being Expressible as a simple P.D. on all  $\mathbb{Z}$ -ops. (rather than a P.D. on all stack ops - which is normal input  $U_n$ )  
I don't remember exact reference. I think it was at IDS or Post Ids.

20 374.40 Where I am now in report: Just finished how telling how PC of SAZ functions are Evaluated.

Brief mention of how 3 input machine works

22 Actually, I have to rewrite (374.00 - 20) on how we use SAZ to get a stack of  $U_n$ 's

New approach: Discuss 3 input ~~ops~~, ~~Algorithm~~ ~~about~~ ~~the~~ ~~program~~ about how good Program. They discuss Gorb: Global soln, Rec Incremental soln.

The discuss SAZ as a practical way to get around the effect of the universal Turing.

Note that the formalism used for <sup>assignment</sup> ~~per assignment~~ is similar to Loh's rule and ~~to~~  $R$ -symbols in the desc somewhat like ~~the~~  $R$ -symbols in a Perm seq.

30 While 22 ff would indeed be clearer & have much more useful info in it, it would be longer & take more time. what I want to do now is write finish a "first draft" as soon as possl.

32 374.00-20 Section: Probabilistic Operators:

We have shown how a method that assigning generates arbitrary functions of any finite number of ~~variables~~ arguments and a way to assign a priori probabilities to these functions.

We will now show how to obtain very general stochastic operators ~~and~~ associated a priori probabilities from these functions.

First Consider the var general





