

ID

.00: 378.40 : 378.00 - .08 is t. moridea! T: Q is: TO represent only part of all poss. O<sup>j</sup>.

Woe # store a sub-distribution O<sup>j</sup>. This sub-distribn can be stored in several ways

2 we can store variables sizes of O<sup>j</sup> sets, but our access is as size ↑; eg. we can theoretically store all of finitely describable O<sup>j</sup> by giving to stock Gramma One & Generator all of them.

A drift from t. extreme, is to store a large bunch of O<sup>j</sup> explicitly, as it's (centers of clusters) - a give a  $\vec{a}^2$  (matrix) for nearby members abt the clusters. The matrix (as 378.29 ff) could be Gramma One! but as was pointed out (378.29 ff) simply giving to dernal t.  $\vec{a}$  and

as a string, and then conceiving of code to extend  $\vec{a}$  is not v.g. Better matrices are needed. O.K.: see (32) o

.09

.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

The stuff from 378.00 to here, seems to be an impt "simple" approach for: all over QATM

It certainly needs more "fleshing out"! 378.22-23 is a major direction.

I don't immediately see how the "Bootstraps" TM<sub>2</sub> ETM<sub>2</sub> occurs, but a discussion of a Matrix

in 378.29-379.09 suggests that TM<sub>2</sub> would find a good "Matrix". Here "Matrix" could mean

a kind of "conditional pc": How much info to get Q to A or from (O<sub>n</sub>, (Q/A)<sub>n+1</sub>) to O<sub>n+1</sub>?

(.R) is answered exactly by 378.02 so that's not quite what I had in mind!

Hvr. 378.02 is not relevant to finite cc & finite sets of O<sup>j</sup>'s.

What we want is a way for TM<sub>2</sub> to expand to d.f. on O<sub>n</sub> before O<sub>n+1</sub> is created.

A good "metric" for this "clustering" is needed. The amount of "Expansion" should be about the same as the amount of "contraction" we get by 378.02.

I've been going around in a circle several times in 378.00 to here! Much repetition of same ideas!

After the "narrowing" of t. d.f. by using 378.02, it has to be "broadened" again by adding to t. d.f. O<sup>j</sup>'s near to O<sup>j</sup> at hyst pg.

This mite be done by mutation/crossover. But knowing t. "metric" for mutation & t. "crossover Alg", is t. Big Problem. I intend to get inside on this problem by

From working on the TSG. Any useful soln. of a new problem should be expressible as some kind of X from on t. old O<sub>n</sub> or old O<sub>n</sub> d.f.

I think t. Main Impt. ideas are: 378.00-.06 thru 379.21-.24

Also note: 378.22-.28 & 378.29-.40, on "clustering" & "Matrix"; Also 379.00-.09 is a useful way to look at t. problem.

NB It mite work (but metricianly) to use a GA approach for -22 using Koza-type

crossovers & perhaps "Automatic function Generation". If it works but somewhat too slowly, we can use t. system to improve its mutation & crossover Algs. (Bootstrap)

382.11-.15 has some ideas on how to define such a metric, using stochastic lang w. 1, 2, 3, 4... known members. 9/22/01

becomes smaller as n ↑ & eventually → only 1. ∴ no "log" henceforth.

It is poss. that this seq will not always ↓ in size: particularly if it is large enough, to start out!

One Big Diffy would be computing P(A<sub>n</sub>/Q<sub>n</sub>) for all these O<sub>n</sub>'s every time n → n+1! After is done, then Avg Improvement in all of this: Express O<sub>n</sub> in t. form of 389.00-.05!

becomes much closer to 300-302! T. d.f. on O<sub>n</sub> is expressed as t. product of indep d.f.s of O<sub>n</sub>-Ops for each O<sub>n</sub>'s; i.e. i=1..n. It becomes much easier to test hiddns of O<sub>n</sub>. Also 389.00-.05 tells how to "Expand" the O<sub>n</sub> d.f. after (Q<sub>n</sub>, A<sub>n</sub>) is known

509.28-31 on "Bootstrapping" Matrix

378-379

QAM as Name.

A Question Answering Machine.

1.00: Consonant Bats. / <sup>equ</sup> 377.26 - 29 <sup>wrong! wrong.</sup> = 378.00 - 03:  
 378.00 - 03 seems right! The @ P<sub>0</sub> (at 377.20): P<sub>0</sub>(A<sub>i</sub>|Q<sub>i</sub>) = ∑<sub>j=1</sub><sup>∞</sup> a<sub>j</sub> O<sup>j</sup>(A<sub>i</sub>|Q<sub>i</sub>) ← Wo!  
 Look at 377.30. Maybe correct maybe should be P<sub>0</sub> = ∑ a<sub>j</sub> P<sub>j</sub>(A<sub>i</sub>|Q<sub>i</sub>) = ∑ a<sub>j</sub> O<sup>j</sup>(A<sub>i</sub>|Q<sub>i</sub>)  
 It would seem that the proof of the QA convergence Bm. would be about identical to the proof for Bay induction:  
 — Because in bay induction we have ∑<sub>j=1</sub><sup>n</sup> O<sup>j</sup>(A<sub>n</sub>) ; a<sub>j</sub><sup>n+1</sup> = a<sub>j</sub><sup>n</sup> O<sup>j</sup>(A<sub>n+1</sub>)  
 in QA " " " " a<sub>j</sub><sup>n+1</sup> = a<sub>j</sub><sup>n</sup> O<sup>j</sup>(A<sub>n+1</sub>|Q<sub>n+1</sub>)

So Bay and the same except that in QA induction, the O<sup>j</sup> of Bay induction is diff. for each value of n.  
 I don't know if that changes the Proof much, hvr.  
 What's the Avn/proof for sequential induction (Sol 78T3)?

1.10 P<sub>0</sub>(X<sub>n</sub>|X<sub>0:n-1</sub>)  
 See 799.20 ff for recant stuff on "SUMAC" (Summary Machine) — which is a development of 378.00

1.20 I'm afraid my proof of the Bay convergence Bm (w.o. Lutz's theorem) of 98TM ~~111.03-17~~ may be wrong!  
 On ibid 111.10 I have ∑ k<sub>i</sub> = -ln ∏<sub>i=1</sub><sup>n</sup> P<sub>i</sub> / ∏<sub>i=1</sub><sup>n</sup> P<sub>i</sub> is valid. Hvr, n<sub>i</sub> couldn't be so! We can make n arbi large, & smc. factors w.  
 diff. r. are ind. I'd expect the ratios to Diverge, n → ∞. [later in day: it's not that simple. 380.23-32 suggests that my objection may not hold water.]

1.23 Go back to Sol 98 (2k...): The universal product PC assignment to ~~the~~ Bay's is  

$$P_{M_0}([D_n]) = \sum_j \left( \prod_{i=1}^n P_j(D_i) \right) P(M_j) \text{ is PC of } j^{\text{th}} \text{ P.D. over } [D_n].$$
  

$$P_{M_0}([D_n]) \geq P(M_j) P([D_n])$$
 for any of the M<sub>j</sub> — including the one that generated [D<sub>n</sub>]

O.k., so that gives the "constant factor" theorem for QATM!  
 It looks like 378.00 - .03 may be "Not exactly right" <sup>Actually it is exactly right!</sup>  
 The universal prod. is not 378.03, but P<sub>0</sub>(A<sub>n+1</sub>) = (∑<sub>j=1</sub><sup>∞</sup> a<sub>j</sub><sup>n+1</sup>) / (∑<sub>j=1</sub><sup>∞</sup> a<sub>j</sub><sup>n</sup>) ; Woops! This = 378.03!  
 The "Constant Factor" <sup>Inequality:</sup> ~~Proof~~ : ∑<sub>j=1</sub><sup>∞</sup> a<sub>j</sub><sup>n</sup> > <sup>2k</sup> a<sub>k</sub><sup>n</sup> for any value of k — including the model index that generated [D<sub>n</sub>].  
 On applying 98TM 111.03-17 to QATM:  
 For ~~the~~ <sup>ratio</sup> ~~of~~ P<sub>0</sub>(A<sub>n+1</sub>) : T. ratio of P<sub>0</sub>'s for A<sub>n</sub> of QATM vs. Generator (of model "k"):  

$$\frac{\sum_{j=1}^{\infty} a_j^{n+1}}{\sum_{j=1}^{\infty} a_j^n} > \frac{a_k^{n+1}}{a_k^n}$$

3.36 
$$\frac{\sum_{j=1}^{\infty} a_j^{n+1}}{\sum_{j=1}^{\infty} a_j^n} > \frac{a_k^{n+1}}{a_k^n}$$
  
 Could we use Sol 78T3 on the sequence A<sub>n+1</sub>?  
 (As we'd do for a single "object" (like D<sub>n</sub>) in 98TM 111.03-17)  
 It may work! The each A<sub>n</sub> in the QA corpus has a different Q<sub>n</sub> & a different "Machine" assoc. w. it, 78T3 may apply anyway. As before, making a correctness of "Expected values" is unclear



00: 380.40: "On Expected values" re: 98 TM 11.11-12 "D": It may be O.K.: for each "Object" in ibid. 03-10

T. "expected value" will be over ensemble of "poss. objects" at the "pt." [I guess this could apply as well to a particular

$Q_i A_i$  in QATM]

: Also NOTE in f.  $Q^k(A_i|Q_i)$  p.c.s are indep.  $Q_i$  to preceding  $Q_{i-1} A_{i-1}$  (of 380.36) & are order dependent. i.e. cause that both connectors down are dependent

So, any way, what does from 380.30 # Gives us:  $\sum_{j=1}^n \sum_{i=1}^n Q^j(A_i|Q_i)$

01: 1) A poss. proof of  $\sum_{i=1}^n \sum_{j=1}^n Q^j(A_i|Q_i) \leq \text{const}$  convergence (Boundedness) Thm.

02: 2) The proby return inequality ( $\equiv$  constant factor) [380-32]

03: 3) A restatement of how QATM does prodn: 380.31 ( $\approx$  378.03) revealing equ.

9/17/01

04:  $\sum_{j=1}^n \sum_{i=1}^n Q^j(A_i|Q_i) > 2^k \prod_{i=1}^n Q^k(A_i|Q_i)$

10: Back to "T-report": [See 376.20, 374, 24, 370.40 & to previous work in this direction] 20 Break down + Re-part into sections: First Desc. each section; Then desc. it in more detail; Then write each section

- Poss. sections: 1) Theoretical EA's! Discuss one based on 3 input inc.  $\square$  ABCDE abcde  $\alpha\beta\delta\epsilon\eta$
- 2) Perms give .06, .05, .07
- 3) E. of shows that  $\sum_{i=1}^n \sum_{j=1}^n Q^j(A_i|Q_i)$  is max; would so any in future, take a reasonable goal.

Smaller this error we get for our answers.

4) 1, 2, 3 may be okay; that  $\prod_{i=1}^n Q^k(A_i|Q_i)$  is max; would so any in future, take a reasonable goal.

Suggester to individual PCA's are "log" & cr. "iterations".

5) Show how any universal input function generator can represent stochastic functs. Our SAZ system does this and can assign a pc to each such function.

WE have shown in section C how SAZ can generate functions of 2 vars & assign probabilities to each of them. How can we use this formalism to represent stochastic relations and assign probabilities to these relations?

Suppose  $F(Q, R)$  is a function of 2 variables: The first variable  $Q$ , is a finite string. The second,  $R$  is a semi-infinite random binary string. The function reads  $Q$ ; then it reads as many bits of  $R$ , it likes then it prints "A" as output and stops.

By defining  $F$  in a suitable way, we can make it print any probab for any  $R(,)$  and any  $Q$ ,  $F$  will have a certain probability of outputting any particular  $A$ .

By suitably defining  $F(,)$  we can use it to represent any binary describable probabilistic relation between  $Q$  and  $A$ . We can use

any binary describable probabilistic relation between  $Q$  and  $A$ . We can use

As a result, SAZ system in section C gives us a way to represent arbitrary probabilistic stochastic relations and it enables us to associate probab probabilities with descriptions of these relations. (Summary lines 19-20).

The formalism of SAZ to assign probabilities to the binary describable stochastic relations.

descriptors of functions such as  $F(Q, R)$ .

(Once we are able) With our ability to generate arbitrary stochastic functions, and assign probabilities to these functions, we can begin to work on the problem of maximizing eq ( ).

Suppose we are given a short sequence of  $(Q_i, A_i)$  pairs, and it is small. We can use the formalism of SAZ to generate many trial functions,  $P_{Q_i}(A_i|Q_i)$  in attempts to maximize  $\sum_{i=1}^n Q^k(A_i|Q_i)$ .

SAZ makes it possible to generate functions in approximate order of their description lengths & probabilities. This, in turn makes it possible to use Levin's universal search (L-search) algorithm for a optimization.

SN on 3 input func  
A, B, R  $\rightarrow$   
can be  $P(A|B)$   
or  $P(B|A)$ !

In exchange for  
promised security  
(which will not  
materialize)

Our Senate has given  
up all our  
civil liberties.

ID

**SN** Just how does Lsch work for finding short codes? We can actually find the "code of x"  
 $(C(x)) \ni \frac{2 - l(C(x))}{T(C(x))} = \text{Max}$   $T(C(x))$  is time needed to test the code for x.

By using more time (CC) we can (usually) get shorter codes.

So Lsch would work for small envt corpus, **W** But what next? We'd like to do incremental updating.

GA (in Koza's style of function trees) would be diff't because testing time for each cond would be very large. We could just test random  $(Q_i|A_j)$ 's for each cond. This would be approx. G's. Or a certain subset of them but second more "critical". Then "fight out" a single cond. by (apparent) G's.

**SAZ** generates an ordered sequence of ~~functions~~ functions (trees). Each funct can use as subfunct, any further (tree) preceding it....

Could we use the algo. to devise a suitable "name" for a "GA"? **W** When I look over this

problem I wanted a "grammar" that would give the known cond's their G's. Using the idea of "cross-over" <sup>being</sup> "SSZ of 2" we should be able to get each cross-over rules, using a simple grammar **OSL**. We can perhaps get "Mutation" from the same Grammar!

It's a nice way to get **extropols**. W. SSZ = 1, 2, 3 ... ect.

The unclear part is how to get a good association of ~~some~~ empirical G's w. t. a SS of the lang.

One obvious way is to have a certain commensat of (funct, trees, concs.) <sup>for all of the Grammars</sup> then for each range of G

Values have a somewhat distrib. Grammar ~~out~~

Another way (that did not, as I remember ~~it~~ work so well) is to use the PC of a stack grammar as the G value assoc w. each SS. <sup>"observable sentence"</sup> T. trouble was, that we had to have short devis to have by G.

On the other hand, it should be poss. to make a grammar in which the "SS" <sup>is</sup> "case count" of a particular string <sup>is</sup> made <sup>as a</sup> <sup>subset</sup> of its observed G values.

The to forgo. is ok, I think I had more or less solved (a adequacy problem) ~~the~~ the problem of continuing the update's algm. (past initialization). I don't quite remember just what I had in mind, but try

revisiting that old stuff!

361.12 - .20 is not bad! "Modifying on to get other trials" is not very clear! Perhaps it was thinking of some kind of "Metric" like (11) or (11-15) for a "GA" approach to "Updating"  
 369.20 - .40 (think 370.00 - .29 on GA) ; **300.36** - 301.05 - .09 as a "fixed saln"

to **QATM**: (with 300.34 - .35 to be resolved. See 301.10 - .20 also 301.21 - .40 + 302.00 - .40

(260.01 - .25): "Some Updating Trends" **Very good**: Some General, some specific.

(260.08 - .33) Is a general plan for "QATM as a whole": It has a rather large "unworked out" parts, but its non-t. loss **Very promising!**

Some implications Not yet Considered for "report" but which is want for T-Macro Reports  
 1) Updating by finding Common Subtree in by G's Cond's (=  $Q_n$ ) **+ OSL**  
 2) The "Histogram" G's for continuous part of optzn of functs.  
 3)

- (260
- 261
- 300
- 301
- 361.12 - .20
- 369
- 382
- Max processor encephalogram:
- 260.01 - .25
- 260.08 - .33
- 261.12 - .20
- 300.36 -
- 301.05 - .09
- 301.10 - .40
- 361.12 - .20
- 369.20 - .40
- 370.00 - .29 on GA
- 378.00 - 379.40
- 382.11 - .15

Not Part  
 Needs  
 history

.20

.24

.30

ID...

80M mals 40M" 60\*60\*24 = 86.4k "/day  
= 463 days.

379.40  
378.40

Expanding the D.F. on  $O_n$  before  $z_{n+1} = z_n \cdot P_{O_n}(A_{n+1} | Q_{n+1})$

This appears to be the main problem of dt: General approach of 378.00-379.40.

379.28 has some ideas on this - (There may be other ideas on this's offered above)

37800-379.40

...so conceivably would not be hard if the number of  $O_i$ 's were large enough. i.e. "wildly" "heated" but always  $\delta \rightarrow \delta^2$  it may!

A possibly v.g. approach to "expand" a pd on  $O_n$ . This pd does constitute a precisely defined stoch lang. (No  $\epsilon$  symbols  $\leq 0$ , ~~it can be~~ normalized). i.e. We took a certain subset

$O_n$ 's ~~subset~~ we did ~~not~~  $\times$  pm of  $\cdot 00^3 \rightarrow$  on  $z_n \rightarrow z_{n+1}$ : we got a new set of  $z_{n+1}$ 's:

These w.  $\delta_{n+1} = \frac{1}{k} \sum_{i=1}^{max} (A_{n+1})_{max}$   $z_{n+1}$  are new set. They have  $z_{n+1}$ 's assoc w. them. These  $z_{n+1}$ 's are relative probabilities: But there are missing members of the  $O_n$  set, that ~~they~~ would have  $z_{n+1} > \frac{1}{k} (A_{n+1})_{max}$ , but they dropped out earlier in the normalizing of the  $TSC$ .

So, we could take the present  $O_n$  set defined by  $z_n$  new  $z_{n+1}$ 's and expand by some pragmatical (of other) extropyn. device. To find the assoc  $z_{n+1}$ 's would be expensive, but one way:

Generate a large no. of cards (say 1M GA)  $\delta$  test them statistically; until a log out no. of them w.  $z_{n+1} >$  threshold have been found. If we can't do it, we are "STUCK".

Some ways to get "outstuck": Continue  $(\cdot 00)^3$  w.o. expanding; do it for  $z_{n+1} \rightarrow z_{n+2}$  or  $z_{n+3}$ ,

then repeat. 10-13 is hopefully, it will work (or one will be stuck again & forced to try 14-15 again).

Just before .04 I had the idea that this "expanding D.F." was one way of doing

a generalized kind of "Backtracking". That we had to find some way to recover  $O_i$ 's that we (unfortunately) deleted several  $z_n$ 's back. **One View of "Backtracking"**: That some time ago, we took a



choice of branch of max pc. That we should now go back to several (not branch is / examine) branches of lower pc.

This is  $\approx$  the idea of going back to  $n-i$ ; lower pc threshold cut, so there are many more cards; so when we finally do  $z_{n+1} = z_n P(A_{n+1} | Q_{n+1})$  we will have enough cards of  $\geq$  desired  $z_{n+1}$ .

**HM!** I had been thinking of using a very large initial subset  $O_0$ . (is small  $\delta$ )

This would mean they ~~with~~ include family  $O_i$ 's of fixed length.  $\delta$  we do  $\times$  pm of  $\cdot 022$ , we modify f. dt of  $z_n$  to have hyperc. — but we have to decn length of  $O_i$ 's int.

corpus invariant. The decn length do not  $\uparrow$  ! **This Seems like a Bad idea.**  
The idea of Morse loss simultly  $\downarrow$  with by  $\cdot 022$   $\delta$   $\uparrow$  with by  $\frac{1}{2}$  mut (crossover etc), seems more reasonable.

We write out a least "expansion by backtracking" if in .18-20 we were allowed to modify the  $O_i$ 's perhaps by  $\uparrow$  Pair lengths.

Rs: .23, I think the decn length of  $O_n$  has to  $\uparrow$  w.  $n$ . (but very slowly we hope) — This Does mean

$\uparrow$  in size of the  $O_i$  set.  
9/23/01 Q: Actually, the decn of corpus,  $[Q_i, A_i]$  does include increase w.  $n$ , because the  $R_i$  are included in the total decn. length. A question is, does the size of the  $O_i$  set have to  $\uparrow$  w.  $n$ ? If the  $O_i$  are decn stays the same as  $n \uparrow$ , then "lang" has stopped! — even tho the variance  $\geq$   $Q_i$  may be "probable" is involve  $\uparrow$  in length of decn. of corpus. "No lang" means "No new Records" — which might occur in certain small areas of sciences, but not (I believe) in most (uncertain) research. (I would mean (int. "2 part-code" version of ALP) But a first part of the code is invariant (except, perhaps for the version of



40G = 4 x 10^10 bytes 3.2 x 10^10 bits. ~ 7E = 2^38.72 bits (= 2^35.72 bytes?)

826P

confinous param's. (Like  $pc \propto n^{-\frac{1}{2}}$  w. b. confinous param's)

In General, (say we are doing a Alg. based TSS): We look at  $q_i$  sol'n(s) to. First set of problems, then we make a grammar that "passes thru" these solns. It is a Stoch Grammar, i.e. we used to search for solns. to subsequent TSS problems — i.e. it generates an "Ensemble" i we look on this Ensemble, for solns. The simplest kind of Grammar is AZI (=SAZ) which counts freq. of occurrences of sub-trees (plus osl). Finding sub-trees is OSL, take most of cc.

This "stoch for solns" would not be a general non-td search of finding a  $O_i$  of  $(h, p, c)$ , but rather "finding a soln to  $Q_i \rightarrow A_i$  or by  $pc$ ". Present finding a way to "recognize"  $O_i$  is diffrent from previous  $Q_i$ 's. This makes it unwise to verify new  $O_i$ 's w/ previous parts of TSS. I think that the great majority of sub-trees is of this kind. The only kind that out involves Backtracking or Equat.

Also, Periodically, to narrow routine of .06-.10 is integrated by looking at a set of recent problems solns (or some other way of categorizing a sub-corpus (= "SC")) and trying to integrate those solns w. a "shorter code". This can occasionally result in a "Great Break Thru"

i.e. Much Compression. This is Newton integrating Kepler's Galileo's "compressive laws". (Also give reference to this sort of Rigidity Event (Briefly "a few P's) in Sal 86 - its size of "why put all in to into a P.A.?" (Also note 261.13-18: this integrates .06-.10 into a less "cl." system)

The "Ob" that is able to recognize the new  $Q_{n+1}$  can also help decide what kinds of trials to make to solve it (Not so clear, but in .06-.10 new "Recognition" after further (EOB) is devised "post notes" — after we know how to solve  $Q_i \rightarrow A_i$ ). Perhaps, after we have a lot of ob. cases, we can begin to extrapolate obs & use past correlates w. parts/aspects of "OPS" to narrow down such the solns to new problems. See 261.13-18. This is useful to extrapolate as well as "identify" out of this Post Note Ob

This integration of several Obs can be of course, concomitant — when several different Obs are replaced by a single ob. These obs would not be equivalent, but the O's for part to may be replaced by a single ob. So many Ob's of Ob. Concomitantly all be replaced by a single ob. So many Ob's of Ob. replaced by a single ob.

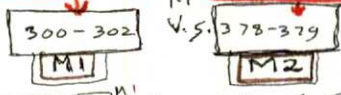
- So perhaps outline: How to do update Global for small Corpus
- (a) Incremental using .06-.10 & v. 12-14: Also method of Obs/ops algebra? (critical)
  - (b) Using GA for Global updates: Excessive local searching auto corpus: use of cover fully chosen/sample QA's
  - (c) Perhaps some other update methods in 260-261, 300-301, 361.12-20, 369-370, 378-379, 382.11-15 } gpp

An attempt to list implications: (201.29) (261.34) 260. Division of Corpus into SC by indexing (redundancy), induction from indexing. 260.12: Common to many SA's @ update  $Q_i$ 's comes (ob's cover new sub-cases (= sub-trees) reparse (19-25 find special operators for  $Q_{n+1} \rightarrow A_{n+1}$  from find "recognition of  $Q_i$ " for  $Q_{n+1}$ : (Note perhaps in pt. in previous at of 261.13-18) (Note (.06-.11) also (.12-.16) v. 0.01 "context"

→ 261.13 p.d.  $O_{n+1} = O_n$  (i.e.  $O_n$ ,  $Q_{n+1}$  from  $A_{n+1}$  is other context) 300.36-301 & 301.36-40 are "good discussions." 369.20-40 Summarized impt. updates (methods) ← GA but in as impt. poss. way to "Bootstrap" TM. (See 385.00ff for dich. it has been inserted by #11; #30 was inserted by #10)

How does it "ground plan" of 375-379 differ from previous 300.00-301-302? Well 300-302 is 2 rth, = 1st "ST" from 375-379 is open loop (At first!) 300-302 300.00 loops 2nd! — new sum, 380

$d.f. O_{n+1} = d.f. Q (d.f. Q_n', Q_{n+1}, A_{n+1})$   $z_{n+1} = z_n P_{ij}(A_{n+1} | Q_{n+1})$



(379.36-40)

(384.30-36) R:  $300-302$ : Harsh problem of Generalizing (d.f.  $O_{n+1}$ ) from  $(O_n)$  is (QA)<sub>n+1</sub> is "Contextual Mod".

$378-379$  (M2) Obtains d.f.  $O_{n+1}$  from d.f.  $O_n$ ; (QA)<sub>n+1</sub> in any direction - T. problem being ~~expanded~~  
 how to "Expand" (d.f.  $O_{n+1}$ ) after it had been obtained by "contracting"  $O_n$ : We do this by taking pts  
 "near" to d.f.  $O_{n+1}$ . T. nature of "Nearness" Metric, is 1. BIG Q.

In  $300-302$  (M1)  $O_n$  (a deterioration of d.f.  $O_n$ ) automatically gets a d.f. for  $O_{n+1}$ , which we may or may not  
 need to check w.r.t. ~~check~~ t. Gov. of ~~t.~~ new d.f.  $O_{n+1}$  w.r.t. ~~t.~~ entire corpus (QA)<sub>n+1</sub>  $i=1$  |  $n+1$ .

In  $378-379$  (M2)  $O_n$  way to look at it: At  $O_n$  point, we look "near"  $O_n$  for "sides" to (QA)<sub>n+1</sub>, when we find end of P. run,  
 $O_{n+1}$  is "wide enough" so that we converge to (QA)<sub>n+1</sub> [An improvement would be if "nearness" func. depends on (QA)<sub>n+1</sub> & S  
 well as other "Contextual info" ← 300.26 ff]

In  $300-302$  (M1) T. search is more "intelligent" if it is guided by (QA)<sub>n+1</sub> & small-  $O_n$  is any other contexts (300.26 discussed "Context")  
 Conceivably, use ~~378~~  $378 \rightarrow$  until TM is smart enough to do  $300-302$  usefully.  $\rightarrow$  391.00

369.20-40 Summarizing t. state of r/sction "updating" upto th. present pt. (there may be other impl. ideas not retraced in it, hvr.)

An ordered way to desc. this:

1) Give Gov. of Gov.; (w. poss. appendix/f.t. on CC problem).  
 Gov. is usually the ~~a~~ <sup>with</sup> set of  $O_n$ . ( $\equiv$  "ensemble")

A possible way do 378. update algm. for  $SS=00$ . (Maybe give shrn on this error; ~~of~~ or mention it)  
 Discuss using small set of  $O_n$ 's, & what I expect would happen (why it wouldn't work).

How narrow might be prevented by expanding ~~via~~ via a closeness metric or GA multicrossover.  
 Closeness Metric could be obtained from Grammar fitting (SAE type grammar).

So p.d.  $O_n \rightarrow (QA)_{n+1} \rightarrow$  p.d.  $O_{n+1}$  :  $[z_n^j O^j] \rightarrow (QA)_{n+1} \rightarrow [z_{n+1}^j O^j]$  ← some  $O^j$  deleted  
 use of fixed elem for expansion " $O^j$  added"

Better : p.d.  $O_n (p.d. O_n) (QA)_{n+1} \rightarrow$  p.d.  $O_{n+1}$  (300.00-302.40)

Do. 29 When TM is capable of understanding & usefully working on problem.

Diffy 25-26! Testing ~~new~~  $O_n$ 's to be added in! involves statistical testing by trying  $O^j$ 's on  
 small no. of (QA)<sub>n+1</sub>'s (or do "sequential testing" (Wald)).

But mainly, I'm worried that system would not "understand" that it should be making  $O^j$  trials that  
 are certain to leave old  $O^j$  results on almost all (QA)<sub>n+1</sub>'s invariant. ~~concentrate~~ <sup>It should</sup> concentrate on (QA)<sub>n+1</sub>.

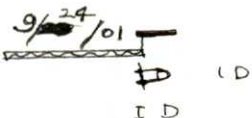
That GA score 2016 to work probs of this sort is of interest! Also idea of using "Grammar-based"  
 Hvr, it would seem that w.o. t. "understanding" of .33 updating would tend to take longer &  
 longer as  $n \uparrow$ .

The hope is that somehow would get into region of  $O^j$  space where new trials tended to

leave .211 QA's invariant. ~~every~~ ~~update~~ ← work mainly on latest (QA)<sub>n+1</sub>

A related approach! to try to find ~~com~~ ~~cross~~ ~~mutu~~ that did .39

$\rightarrow$  386.00



ID

385.10 :   
 Stroke will be involved in some of the forms. (Do not if GA is used)  
 If L is used, variation very during such (so it's not strictly L-sided; one can't estimate time to solve)  
 Also, perhaps, "context dependence" of pc's of concs, etc. Same 300.00 ft. for decn.

Another way to get  $O_{n+1} \approx O_n(O_n, (QA)_{n+1})$  is by first getting pc's of new concs by  
 1) Context free: uncond. pc.  
 2) Varying units of context  
 3) Then try for "conc"  $O_{n+1}$  as a fun. of  $O_n, (QA)_{n+1}$ .  $\rightarrow$  (10) tells how, hvr  
 How 2) moves continuously to 3) is unclear! — Tho  $O_n$  300.00-302.10, it did occur! Perhaps, it was 2  
 or 3  $\rightarrow$  nonal jump pl. I guess the idea of concs. & cond. pc. of concs was only a method of getting

$O_{n+1}$  trials, so if we could do this Directly, it would be Better ("less eff.")  
 One way to get to  $O_{n+1}$  is by ELN: we have many QA pairs for t.  $O_{n+1}$  indivis: unworkable  
 Also ELN would be to find good concs & assoc. (context dependent) combination makes to create good

$O_{n+1}$  trials: We could start out w. context free concs as .04: This (probably) may be the same as t.  
 pc's of t. concs that we use to evaluate t. "code length" of  $O_{n+1}$  for t. final conc -  
 As  $S \uparrow$  we have more cases of relevant contexts, & we use them to define cond. pc's of concs. From  $O_n \rightarrow O_{n+1}$  or  $O_{n+2}$   
 i.e. we wouldn't optimize, hierarchy concs, hvr. alternately  
 Actually .10 ft may be a nasty path to go to .03. One might get there by GA, hvr, but I suspect

much slower.   
 Tho .03 ff & it in particular, are certainly good ideas, they say the resultant system does not  
 deal w. the 385.33 problem of trying to find a modification of  $O_n$  that preserves some fold solns & solves new problem.  
 384.06 - .20 (384.01-15) R Reviews some ideas about this preserves

So .00-.16 is a possibly not bad way to do (300-302) — But .17-.19 are very impl but not included.  
 A poss. way to integrate: we start out by ~~the~~ TM by doing the  $O_{n+1}$  only to solve problems.  
 we freq improve .17-.19 by .10 - .14  $\&$   $O_{n+1}$  (mixing in .03)  
 The part I'm least certain about is  $O_{n+1} \rightarrow O_{n+2}$  (.05  $\rightarrow$  .06); An is "soln. outline" in .10 - .14.  
 But drop this for a while: say  $O_{n+1} \rightarrow O_{n+2}$  at all would work. Could we integrate 378-379?  $\rightarrow$   $2n+1 = 2n \cdot P_{O_n}(Aver(\text{Time}))$

378-9 so far: Advantage 385.00-09 discusses 300-302 vs. 378-9: 378-9 is easier to run (probably, ~~etc~~ — particularly  
 if I use GA for "expression of P.R. — But it seems not very smart — seems not to use hvr, etc.  
 It might work, hvr. I certainly feel 300-302 is more likely to work!  
 A poss. reason to believe 378-9 might work: That GA was able to solve what seemed to be "Hard problems"  
 Also 378-9 seems a lot easier to implement/try than 300-2. A main read study header was a "Distance"  
 metric — presumably deriving Grammars to use  $S \rightarrow 1, 2, 3$  etc.  $\&$  "I" write consist of  
 Improving to Metric — making it dependent on  $(QA)_{n+1}$ ,  $O_n$  is perhaps one thing. — This would bring it  
 close to (10) ( $\approx$  300-2). The reasoning here seems to be: How can we improve r. "Metric"?

An essential Genzn of t. Metric is a mode of expansion/extension of t.  $O_n$  d.f.:  
 d.f.  $O_n \rightarrow$  d.f.  $P_{O_{n+1}}$  is what we want, & we're obtaining it by d.f.  $O_n \rightarrow$  expanded d.f.  $O_n \rightarrow$  d.f.  $O_{n+1}$   
 &  $\&$  hvr, they could be in opposite order, hvr. — Re I'd prefer that t. present order is best  
 because we do want all of t. new concs to work w.  $(QA)_{n+1}$ . (?)  
 So .31 ft does seem to Merge 378-9 w. 300-2.



LD

: **SN** Away to use **GA** for a **"Problem Pool"**  $\rightarrow$  **ISQ**!

We have this unordered set of problems ( $\equiv$  "Problem Pool").

We try various functions (obtained by GA) on the problem pool. At first, No funcs solve any probs; but we cross/mul. ~~func~~  $\rightarrow$  a set of funcs that solve  $n$  problems: As soon as we get a func that can solve  $n+1$  probs it goes into  $f_i$   $n+1$  "function pool"; so  $n \rightarrow n+1$  is loop 20-03. So we have many function pools (one for each value of PM), & they recombine w. members of their own pools perhaps or occasionally w. funcs. at  $n+1$  pool? So we have pgrams ~~are~~ discussing to hyper hyper "create" ( $\rightarrow$   $20 \rightarrow n \rightarrow n$ ).

The trainer will observe what problems are difficult to solve (few solns) or no solns at all, & try to derive "bridge problems", to solved probs: especially probs solved by many funcs (?).

Out. Problem 386.18 (solving new probs preserving solns of old probs (infect!)).

Normally, when one solves a problem, the soln. method is not available for exactly the same problem. Also it is available for certain known variations of the problem (e.g.  $3+7 \rightarrow 10$ ; ~~the~~ same soln will usually ~~work~~ work for  $a+b \rightarrow [a+b]$ ). In general, for each problem solved we will ~~know~~ <sup>know</sup> how to extend it to ~~related~~ "related" new problems, usually by slightly (or unadvisedly) modifying the soln. I think this is the basic point of 386.18 question. Normally each problem soln. method will have an assoc. ~~to~~ algem. method which also ~~will~~ can be used w/o how to modify the soln. method for different (associated) problems.

Unlucky, this QATM can't make hypothesis about what kinds of new problems the ~~newly~~ found soln. method will solve, & give pos. ~~or~~ neg. examples to "Teacher" for pos. or neg. verification.

The way it works, is that the method of "identity" of the  $Q_1$  is a short code, so it automatically explains  $Q_2$ .

Fig.  $3+7 \rightarrow \text{sum } 3+7 \rightarrow \text{sum } R_1 R_2$ .  $\text{sum } R_1 R_2$  works for many examples

in the  $TSQ$ , so it's very compressive;  $32218622 + 411900252 \rightarrow (CC)$   $\text{sum } R_1 R_2$  w. ~~the~~ same  $2$  big base base and  $CC$  is also very compressive.  $\text{sum } R_1 R_2$  is great even from original data.

For all  $3$  register probs, if  $R_1$  &  $R_2$  top pos, then  $\text{sum } R_1 R_2$  works if  $R_2 \equiv "+"$ .  $\leftarrow$  is discarded empirically.

So, ~~my~~ previous idea that when one solves a new problem is immediately on w. this soln, it is usually (almost always) unnecessary to verify that this on will continue to work on all older problems. In fact, analysis like all of is normally used for all problems.

"Backtracking" is a special, rarely involved technique. (or serially sequentially)

What looks like is a classic Expert System: Rec "Obs" work in it: when one of them "lives". Its assoc "OP" is extended. Since there are an enormous no. of "Obs" to be run (they could run in 11) One should find ways to test them: Rec can be done by ~~test~~ looking at the problem & roughly categorizing it (DP3?) so one knows which obs to try first. (trying all (or many) obs is not save time, but not nearly CC!) — Note 396.26 on Hierarchy of Obs to be PC of "Recognition".

**Context: .06ff**

.06-.08 was the most general defn. of a "context", i.e. T. thing of an "ob".

.00 : 387.40

Instead of a single sequentially written prog, On will consist of a bunch of op-ob combinations: that are allowed to use other obops (or just ops or just obs) as subroutines. I think this may modify, considerably, the way I will be thinking about these progs. The individual ob/op functions can usually be deterministic, but can be stochastic. — This last would be a hint to give a stack output for any point "containing" a stack function as a "subfunction". Any of the funds can be implemented w. the "R" input, or they might be expressed in any of the other ways. see (.05-.09)R

.05

.05 R Make index of variety of ways to express stack docs. If they be first there are only a small no. of ways to get them, & all other ways are combinations of funds at this same section & stack or Defn. funds  
.09 R 244.07 ff has some of this

.06

T. Mode of (discovery construction) of .00-.05 will be same as before. (386.10-.16). When a new mode (during stochastic construction of an (op/ob) & certain context occurs, an assoc ob will "fire" and if a appropriate pc's of ~~the~~ funds will be enabled. Much cc will be spent looking for new contexts "Post More on

.09

.10

.11

May usefully modify pc's of sub. funds. — see 396.30 for "More on context" (2.06-.09)  
Usual Contexts: 1) What funds are being connected to? 2) What is the main problem being solved? 3) What (local) sub-problem is being solved?  
Other contexts will become clear as I devise TSO's & their solns. 300.26 was big expansion of "context" of .10-.11

So: How to go to "context" for soln. to specific problems, To attempts to devise On? Well, in this approach of .00-.05, there really is no "On". There is only a set of attempts to solve fun. problems.

(Qns, Ans). If all ops used to identify new problems are stochastic, then finding On would mean improving methods of solving (QAs) and perhaps improving soln of previous problems. "improvement" can mean ↑ of cc of answer & / or ↓ of cc of answer for past problems)

The core is as before, to ↑ of  $\frac{R_i}{T_i} P_n(A_i/Q_i) \cdot P(On)$ , but On is a different form from what I was thinking about when On was first written. — It's like (.00-.05)S

.20

Here, it should be poss for T. presently contemplated QATM, to be given to problem of improving itself! As before, + Q is On [tidcons of f. ensembles On]; Qns, Ans, plus previous: QATM pairs of fcs type — i.e.  $\left[ \begin{matrix} R_i \\ Q_i, Q_{jt}, A_{jt} \end{matrix} \right], \left[ \begin{matrix} A_i \\ A_j \end{matrix} \right]_{j=1-n-k}^n$  k=0-n

Also part of Q would be an "index" telling TM what f. problems So, as before, TM can work on S.F. as a "Normal QA problem" — after it "understands" what f. problem is.

The fact that On has to form .00-.05 rather than a string representing a single function & 3 variables (S, Q, f), should make no big difference in my discussion of how QATM works. Suppose we give QATM of .00-.09 a new problem's Qns that Qns does not cause any

"recognition" of "fire". Then TM can't work out problem! Seems strange, hvr, since I can usually almost always "work on" any new problem. In previous case TM was unable to even work on Qns at all! Well, see I was given a problem I didn't understand — (So first sub-problem is to try to understand f. language!). Well, there are problems that I can't work on because I find they are not well enough defined.

Hvr, it may well be that a suitably educated QATM could work usefully on any problem that I (or any human) could usefully work on. — i.e. changing p.p. during L such — expected to use

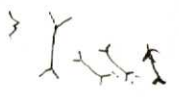
Expanded context is pairings (Lny during Lsuch, were main trick) for better updating. I do, hvr, have a way so that any poss. hvr can be expressed. My earlier idea was that (modifying to p.d. that Guides Lsuch is Lny during Lsuch (say # times now hvr. Lny. is Lsuch) might be adequate to express precisely any hvr. Perhaps write a bit more detail on how this can be done — particularly in the present (.00-.09) QATM scheme.

.37

.39

Def:  
Def:

- D** = Det. = Deterministic. D-funct., DOP, DOB.
- S** = Stoch = Stochastic = Probabilist S-funct, SOP, Sob.



0.00: 388.10 ; Common Backtrack situation! It occurs in a SC. A certain conc. has had ~~the~~'s case count. recently  $\uparrow$  so it can just now, be used for last 3 or 4 cases, giving  $\uparrow$  of PC  $\hat{=}$  change of code for last 3 or 4 cases. This new conc can be Ob or OP.

In General, we want to have b. Ob or OP (found/discovered) for a new QA to be as general (i.e. as "startincode") as possible. Otherwise, it's a very narrow soln, we haven't found much!

- So a large fraction of this time will be spent integrating (i.e. General) old solns of old problems. - Like deriving an OP selector next would do many old Q's. Early Obs only select one op. Integrated obs select from many ops

"Advanced" Obs will analyse a Q & generate a set of ops to solve it. By "set of Ops" I refer to a D.f. over Ops. (See 388.05R-09R)

So:  $t$ . initial soln. of a single problem ( $\equiv$  QA) by an OP is finding an Ob for that QA, it is an initial extens elch. of the General QA problem. - That improvements of .00-10, up to 50% or more of

1.13  $\max P(Q_i) \cdot \prod_{j=1}^n P(A_j | Q_i)$  are heavy  $\hat{=}$  amount to  $t$  main rule Govc of  $t$ . system. In 1.13, I'm beginn(ng) to lose t. feeling of  $O_n$  as a P.D. (= ensemble) - Tho  $t$ .  $P_{O_n}(A|Q)$ 's can be p.d.s (ensembles).

Also Note! for a new problem (QA or) PERHAPS  $T$ . system will not recognize Q. After  $T$  has found a good  $Q \rightarrow A$  ops,  $T$ . find that this Q was not recognized, should help find a recognizer. - i.e. look at  $t$ . logical decision pts. at which  $t$ . present Q was rejected. Actually,  $t$ . ob. that is  $t$ . total rejection by all previous obs, is an ob, but (i.e. "WHY" did  $T$  reject this Q?) Even probly not a good one. Is it expensive  $\hat{=}$  in  $ps$ ? One would have to look at previous Obs - which sounds expensive. Also all old obs are Mut. Exclusive: No overlap (usually or always?). Another funny: Q is accepted but A is given a low  $pc \rightarrow 390.00$

So: A poss. outline of report:

- 1) Prob. QA problem: Also go into how/why problem that was could be explained to human code
- 2) Give exact goal soln. first using  $\leq$  overall  $O_n$  (1.13). (CC=0)   
 Perhaps give explain in terms of ints. of all of 378.00-006

**RECOGNITION Problem**

3) How to approximate  $26(\equiv 13)$  of previous solns relevant. Discuss "Muller's" "lactan" problem. Obs, Ops.  $\hat{=}$   $t$ . decay of 388.00-14 on  $t$ .  $\hat{=}$  initially "represent" form of  $O_n$ 's  $\hat{=}$  how to "interpret" it, Au Amfay" we make trial. no context. (uncondi. pc of trial) - But gradually  $\hat{=}$  context. In particular discuss  $O_{n+1} = O_n(O_n, @A)_{n+1}$  made as a first goal.

4) ISQ's Tell how  $t$ .sq will be used train both intuitive an instructed  $T$ .sq long. by instructor leads to Modifn. of @TSQ b Mode of operation of  $T$ .sq.

5) Maybe section of FAQ's.

- Intrinsiquently asked Q's.
- Occasionally Asked Q's
- Sometimes difficult to Answer Q's
- Impossible Embarrassing unpredictable Q's
- Embarrassing " " "

Work done by Q/AM

Perhaps write xxx Reports on GA / IA's

Also 2 The 378-9 system. w. Criteria sq. See 395.00 - comparing 378-9 to 300-302.

Perhaps have a **FAQ** section at end! Try to enhance objectives, misunderstandings.

00: 359.19 : 288.29 discusses what happens if  $Q$  is not recognized by any ob.

01 What if several obs recognize Q? How does (or does) mix their ~~output~~ output d.f.'s for A? (Also Note 37)

[It] seems that I've worked on this problem before! I had several IA's, each w. somewhat different SC's (sub-corpi = dataset)

perhaps 4 problem(s) leading up to analysis of  $\sim 250$ . Looks very much like present problem!

There is some posssy that I did an earlier analysis of a problem closer to present problem.

244 or ff discusses ways to access stack ops.

I vaguely remember something like this: I had a bunch of IA's, each w. each would have an assoc. SC. Some of SC's for different IA's overlapped somewhat. In my earlier version, the IA's were clean, so they either predicted to hit A or not. For t-whole corpus, consider only data (QA's) in which  $P_{obs}$  was some overlap.

We wanted to assign wts to the IA's so that if we used  $\sum$  / sum of wts for each to predict, we could predict normalized

sums of SC PC,  $P_{obs} = \sum PC_i$  (of those QA's having overlap of QA's w. overlap), was max.

I think the idea was  $PC_2 = \frac{\sum wts. of IA's that predicted correctly for that Q_i}{\sum wts. of all IA's for that Q_i}$

In the present case, each predictor assigns a PC to the correct A,  $P_{obs} = \sum PC_i$  (of those QA's having overlap of SC's is to max.)

T.  $PC_2 = \frac{\sum w_j \cdot P_j}{\sum w_j}$    
  $\cdot 288, 30$    
 this PC<sub>2</sub> would be the PC assigned to A,  $P_{obs}$  is possible answer.

$$\frac{a}{x} = \frac{b}{1+x} = 0$$
$$\frac{a}{x} = \frac{b}{-x}$$
$$\frac{2+a}{1} = \frac{a}{x}$$
$$x = \frac{a}{2+a}$$

So we want to make wts. assignments to  $\sum PC_k$  is max,  $PC_k$  being the  $PC_i$  (of 13) assoc w. the correct A of (QA)<sub>n</sub>. This product is over only those QA in which overlap occurred.

To do this max, we note that in 13 if we multiply  $w_j$  by a constant, 13 is invariant.

So we want to do max  $\sum PC_k$  subject to constraint of  $\sum w_j = 1$  or any other constant. Berny Seq.

Could we do L's Multiplexers? In the case of 11, I think I found a way to force it as a Berny Seq.

248.16 does deal w. several A's for some Q. Also Note 248.23, 32

How, read 250.09L - .19L It sort of tells what that stuff is about.

NOTE: T. ideal 250.09-19 was that I could do to Berny Seq. & Laplace's rule. for this case.

$$\begin{matrix} s_1, s_2, s_3, \dots \\ \text{or } p_1, p_2, p_3, \dots \\ \sum p_i = 1 \end{matrix}$$

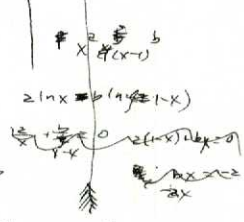
AH! I wanted this  $S$  because it represented the total pc of coding to corpus using all different wt assignments.

In 250.09-19 Doing  $P_{obs}$  with difference between using Lap's rule vs. "Straight rule"

T. "straight rule" ~~probably~~ curves tends to using the peak value of the wts, for predict. (which is what I had in mind in 13)

which is o.k. if  $SSZ(13) > 5$ , say.

[I guess what we want is]  $\sum S_1 \dots S_n$  pred. wts,  $\sum w_j$  and ideally, we could do this for prediction - i.e. see how the  $\sum$  varies when different pc's for the new A occur.



2.7 Earlier I did work that involved finding peaks of products: I used a fast N-linear optimization method.

2.8 (13) ~~to get peak~~ take in form  $\sum w_j$ :  $\left( \frac{\sum w_j P_j}{\sum w_j} + \frac{1}{\sum w_j} \right) = 0$  + 1 more eq. - Actually, 214 constant will do!

3.0 If I can't find a simple approach, one way is to use the peak of 13 (which corresponds to "T. straight rule")

T. straight rule is bad for low SSZ, which is probably case of interest so to get a higher of Laplace's rule

To do this we use the reasoning that gives Lap's rule, (Res multidimensional integration) & apply it to

2.3  $\sim \exp(13)$    
 O.k. for the present: 13 is probably good enough for present. Probably it's possible to have out

~~some~~ a good approx to "Lap's rule" by modifying the solution to the SSZ rule (of 13).

Also try various analogies ~~between~~ between Berny seq. & present problem to suggest approaches to a process.

3.7 (13) 9/30/01 T. P. seq. is fine if all obs (at or at least one of) obs gives a DP of that "networks" acceptable.

If non-acc to recognizing obs. gives acceptable results for Qnt, Ans, then we have to find a new ob-opdf. to solve Qnt, Ans, & also modify old obs, so they don't recognize Qnt, Ans.

ID

N.B. I think there is a previous ref. to ~~the~~ ~~paper~~ . go to. It starts at the top (.00) of a page - 385.00-.09

Def (302.27-.30)  
.00: 385.09: (1-11)  
(379.37-.40):

M1 Model of 300.00-302.40 (d.f. On, An, (d.f. On, (Qm, An))  
M2 Model of 378.00-379.40 (An, An, Pj (An, Qm))

The distance is in both have at any time a d.f. on On.

Their responses differ in how they respond to a new Qm, An.

.03 → M2 first (Contracts narrow) Re On d.f. by  $z_i^n = z_i^n \cdot P_j(An, Qm)$ , then "expands" to d.f. using some kind of matrix, or GA.

.05 → Also M2 may do P. ~~expansion~~ first, then ~~contraction~~. The only difference is that in case of .03, the expansion depends on QAn (as well as QAn) in " " .05: " " " " QAn only so it's not very intelligent. - but the new replies to Qm are completely w/o "fudge". - tho actually, there's not an issue, there is as little "fudge" in .03

In fact, I can't see any difference betw .03 & .05 They both alternate contraction

& Expansion. Contraction is "routine", fixed Alg; Expansion is "Creative", "heuristic"

M1 looks at QAn & decides how to modify the  $O^i$  d.f. on Pm basis:

.06 - So, while there would seem to be an "fudge" problem, there is not because we use ALP (which is probably the contraction equivalent of .03 used by M2)  
M1 first uses the contraction eq. on a new QAn. If it gives a deviate PC for An then the  $O^i$  d.f. is not modified ("expanded"). If An gets too small is PC, we expand to  $O^i$  d.f. until an adequate PC for An is obtained.  
From (.06-.08), it would appear that M1 & M2 are about the same. - P. meaning  
.20 What M2's expansion method is; how much it differs from M1's "expansion" method.  
But go over 385.00-.09 (t. earlier discussn of M1 v.s M2)

0/3/01 Ont General System Implementation: It could be by storing many  $O^i$ 's at all times; or storing a few solns to each "problem type" or (3) storing few  $O^i$ 's but also

storing states at imp't "branch pts", so "Backtrack" is easy.  $O_j \rightarrow Q_j$   
0/6/01 But remember, the form of  $O^i$  is a "product" of a larger set of  $O_j$  sub-problems.  
The meaning of "Branch pt" or "Backtrack" is much different, from that of a sequential p.m. - (25) is more like a point!  
= like Selfridge's "Pseudomonium" (1957): for each  $Q_j$  for we have a d.f. of  $Q \rightarrow Pbs$ .  
T.d.f. could be a single Ob of Op or a set of Pm: We may have a product of a Op d.f. & a Op d.f. or a more constrained d.f.

- 00: Expansion of 389.20 (Outline of Report)
- 01: ① Describe QA problem. Give examples, but also discuss "Generality" (Asking human a QA or defining problems to human)
- 02: ② Give genl. soln. for  $c = \infty$ :  $\sum_j (389.13) \cdot P_{O_j}(A_{int}|Q_{int})$   
 Perhaps mention updating routine of  $z_{int}^i = z_{int}^j P_{O_j}(A_{int}|Q_{int})$ .  
 Discuss: ~~is~~ actual predn algn. for 02 then some of its properties. (Finite cost)  
 (i.e.  $\prod_{i=1}^n P(A_i) > c$  some criteria for ~~convergence~~ "true" cond. pd.  
 Perhaps discuss approxs to w. finite no. of  $O^j$ 's why it wouldn't work

10: ③ Give ~~389.13~~ w. ~~finite~~ d.f. over finite set of  $O^j$ 's as a kind of approxn.  
 ← in which we want to maximize  $\approx 389.13$ .  
 SN In most of this stuff, I'm not sure that I will have as much of a d.f. over  $O^j$ 's:  
 The individual Ob-ops that constitute  $(O^j)$  may have d.f.'s, or could be simple  $\leq$  ops.  
 If, if I do find  $\geq$  soln. to a QA, these solns. look much different (i.e. they work differently for other Q's) then I certainly will want to save them, either a normal saving or by  
 $\Rightarrow$  not making all of the solns part of  $O^j$  but saving the solns as possl. "branch points" for Backtracking.  
 I will want to decide just how to present this stuff:  $O_n^j$  or d.f. on  $O_n^j$   
 Probably d.f. on  $O_n^j$  might be best. ~~How~~ when I do have several poss Ob-ops that solve a problem differently (like 390.01-20 then I must have a d.f. over  $O_n^j$

20: Anyway, do go into 389.26 ( $\approx 389.00 \dots \approx 389.19$ ) Some details on how normal updating occurs.  
 Then  $\rightarrow O_{int}^j = a_n(O_n^j, (Q_n^j, H)) \rightarrow$  a final goal ... Guided by 389.13.  
 Then discuss TQS: How they are used to draw TM to draw Me.

Guided

$$\int_{-\infty}^{+\infty} e^{-\frac{x^2}{2\sigma^2}} dx$$

9/9/01 (From 02): There exists a formal solution to ~~the~~ the conditional probability problem.  
 An (unnormalized) probability distribution on  $A_{int}$  is

$$\sum_{j=1}^{\infty} P(O^j) \prod_{i=1}^{n_{int}} P_{O^j}(A_i | Q_i)$$

Here,  $O^j$  is the  $j$ th conditional probability distribution.  
 $P_{O^j}(A_i | Q_i)$  is the conditional probability induced on  $A_i$ , by  $O^j$ , given,  $Q_i$  ( $Q_i$  is the condition).  
 $P(O^j)$  is the a priori probability or "weight" of  $O^j$ . It is approximately

$$2^{-L(O^j)}$$

where  $L(O^j)$  is the length of a short description of  $O^j$ .  
 The summation is over all possible  $O^j$  distributions.

It can be shown that eq. 25 gives a extremely good distribution

function for  $A_i$ . However, the relation is incomplete. The eq. 25 is always incomplete, it is possible to approximate it. The present paper will describe some methods of approximation methods.

One approximation method immediately suggests itself! Suppose that we do not sum over all  $O^j$ , but over a very large subset of them. We have not been able to get that particular approximation to work for any practical number of  $O^j$ 's. The present paper will describe some methods of approximation that we have found most useful may be regarded as a variation on that theme

393.16

393.00

Form of the answer: .25  
 An equivalent form is  $\sum_{j=1}^n P_{O^j}(A_i | Q_i) \cdot P(O^j)$   
 (clearer) (clearer) perhaps underlining

.00: 392.40 :

It is possible to get good induction by using only one  $O^j$  - i.e. one for which

$$P(O^j) \cdot \prod_{i=1}^n P_{O^j}(A_i | Q_i)$$

is as large as possible. ~~It~~ If  $O^j$  is the best we have been able to find then we

can use  $P_{O^j}(A_{n+1} | Q_{n+1})$  as a conditional probability distribution for  $A_{n+1}$ .

We can get a better conditional probability distribution

Instead of choosing the  $O^j$  a priori, we choose for each  $n$ , a set of  $O^j$  such that ~~for each~~ ~~each of them~~  $\sum_j P(O^j) \cdot \prod_{i=1}^n P_{O^j}(A_i | Q_i)$

.10

is as large as possible. We ~~can~~ ~~suppose~~ pick the 100 best  $O^j$ 's, then use them in a finite sum ~~with~~ ~~the~~ modification to obtain the distribution function of  $A_{n+1}$ .

In general, it is usually quite difficult to find  $O^j$  distributions such that  $n$  is very large

~~smaller~~ ~~practicably~~ of  $n$  is at all large.



.16: 39 2.39

However, we can get a good, useful approximation to .25, by choosing a set of  $O^j$  such that ~~for each~~

.10

is as large as possible.  $a_n^j$  may be regarded as the "weight" of  $O^j$  and the resultant

distribution function on  $A_{n+1}$ , for question  $Q_{n+1}$  is

$$\sum a_n^j P_{O^j}(A_{n+1} | Q_{n+1})$$

.20

.21

For each  $n$ , we want to find a set of  $O^j$  with maximum ~~with~~ ~~weights~~, ~~the~~

A very practical difficulty is that if  $n$  is more than 3 or 4, the problem of finding a good set of  $O^j$ 's is ~~extremely~~ <sup>normally</sup> time consuming.

What would be a desirable  $O^j$  one of the proper form,

We will discuss a method for finding the  $O^j$  set for large  $n$  by a process of "updating" ~~we have~~ ~~found~~ a good set of  $O^j$  and associated  $a_n^j$  for the data set  $\{Q_i, A_i\} i=1 \dots n$ , ~~then~~

~~contribution~~ ~~to~~ ~~the~~  ~~$O^j$~~  ~~being~~ ~~of~~ ~~a~~ ~~certain~~ ~~form~~, it is possible to obtain ~~the~~ ~~weights~~ ~~from~~

the  $O^j$  set for  $n$ , a ~~set~~ ~~of~~  ~~$O^j$~~  and associated weights  $S_n$ .

Given  $Q_{n+1}, A_{n+1}$ , we have an algorithm that will use  $Q_{n+1}, A_{n+1}$  and  $S_n$  to

produce  $S_{n+1}$ . Using this algorithm we can find  $S_1$  by direct search of Eq. 2.2,

then the algorithm can be used to find  $S_2$ , then  $S_3$ , ... and finally,  $S_n$ .

Change notation:  $a_n^j \leftarrow P(O^j)$ ;  $P_j(A_i | Q_i) \leftarrow P_{O^j}(A_i | Q_i)$

First deal with QA problem (See 392.01)

There is a formal solution to the conditional probability problem.

$$\sum_{i=1}^n a_i \prod_{j=1}^n P_j(A_i|Q_j)$$

define  $a_i = P_j(A_i|Q_i)$  (392.26-30)

"The sum is over all finitely describable conditional probability distributions, ~~P\_j(A\_i|Q\_i)~~,  $P_j(\cdot|\cdot)$ .

It all of the  $Q_i$ 's ( $i=1 \dots n+1$ ) are known, but the  $A_i$ 's are unknown, then .01 gives the universal joint probability distribution for the  $A_i$ .

If all of the  $Q_i$ 's are known, and  $A_i, i=1 \dots n$  are known but  $A_{n+1}$  is not known, then .01

gives a universal distribution for probability distribution for  $A_{n+1}$ .

An equivalent form to .01 is

$$\sum_{i=1}^n a_i P_j(A_{n+1}|Q_{n+1})$$

we define  $a_i = \prod_{j=1}^n P_j(A_i|Q_j)$

and the sum is as in .01 over all finitely describable  $P_j(\cdot|\cdot)$ .

.1) Expresses the conditional probability of  $A_{n+1}$  as a weighted mean of the conditional probabilities induced by all  $P_j$ 's. - The weight of  $P_j$  being proportional to the probability it assigns to the known data,  $(A_i, Q_i) i=1 \dots n$ .

This solution to the induction problem is known to be very accurate, perhaps as accurate as any induction system could possibly be. It is, however, intractable in finite time with finite memory.

Most legitimate attempts to do induction can be usefully regarded as methods of approximating this ideal solution. From this point of view, we can compare and criticize most physically realizable induction methods.

The present paper is a rather direct approximation of volume 11. Instead of summing over all  $P_j$ 's, we only sum over a relatively small number of  $P_j$  that have large weights,  $a_i$ . Discovery of  $P_j$ 's that have large weights, is the main problem. We look for  $P_j$ 's such that the value of  $a_i$  is as large as possible.

Usually we will have a large set of  $P_j$ .

For each value of  $n$ , we will usually have a large number of  $P_j$ 's with large weights. Usually as  $n$  goes from  $n$  to  $n+1$ , the  $P_j$ 's of large weights will change completely. However, normally we can find it possible to express the high weight  $P_j$ 's for  $n+1$  as modifications of those for  $n$  examples.  $\rightarrow$  395.01

More precisely: If  $P_j$  was a large weight for  $(A_i, Q_i) i=1 \dots n$ , then we usually are able to find (anew  $P_j$ ) that is a modification of  $P_j$  such that

$$P_j'(A_i|Q_i) = P_j(A_i|Q_i) \text{ for } i=1 \dots n, \text{ but } P_j'(A_{n+1}|Q_{n+1}) \gg P_j(A_{n+1}|Q_{n+1})$$

$P_j'$  and  $P_j$  are the same for  $(A_i, Q_i) i=1 \dots n$ , but  $P_j'$  is much better than  $P_j$  for  $(A_{n+1}, Q_{n+1})$ .



Some where in here, a sentences company  $a_i$  — (LPC of  $P_j$ ) (unfortunately, functional format  $P_j$  is peculiar, it's. AZ-141 method described is not applicable in an obvious way. The sub functions are unstandard.

The process of creating a good set of  $P_j$  ~~is a very hard~~ <sup>new</sup> Given a "good" (high weight) set of  $P_j$  for  $n$  examples, the process of creating a/s of good  $P_j$  for the data set having an additional new example,  $(a_{n+1}, Q_{n+1})$  is called "updating".

While it is possible to do good induction using only the  $P_j$  of ~~the~~ highest weight found, the process of "updating" is easier to do if we have a large number of good  $P_j$ 's for each set of examples.

The ~~problem~~ <sup>this</sup> "updating" is usually the most difficult ~~part~~ <sup>aspect of an</sup> induction systems, and the method used for updating ~~is an important way every good way~~ <sup>to characterize</sup> induction systems. ~~to implement it is a good way to characterize~~ <sup>an</sup> induction systems. (works)

Section: Updating the  $P_j$  distribution.

There is an ~~easy~~ obvious way to do updating: For ~~each~~ <sup>from the set of</sup>  $P_j$  of high weight in for  $P_j$  & the  $n$  member data set, we ~~can~~ <sup>can</sup> construct new trial  $P_j$  by a process of mutation and crossover, ~~the method as is done in~~ <sup>is done in</sup> Genetic Programming. The "fitness function" is

of course the  $a_{n+1}$  of .12. While this technique would probably work, the fitness function is very expensive in time, since each  $P_j$  has to be tested on all

$(A_i, Q_i)$  for  $i=1 \dots n+1$ .

Taking a clue from the way humans seem to solve problems, we will limit the  $P_j$  sets considerably by making the following constraint:

(Insert 394.33 - 36)

We can do this by ~~constructing~~ <sup>making</sup> each  $P_j$  so that it is a combination of ~~making~~ <sup>making</sup> each  $P_j$

2. combination of  $n$  functions ~~to~~ <sup>one</sup>  $A_i, Q_i$  example. Each of these functions consists of a ~~recognition~~ <sup>recognition</sup> part and a ~~operator~~ <sup>operator</sup> part. The ~~recognition~~ <sup>recognition</sup> part, ~~the~~ <sup>the</sup>  $R_i$  recognizes how  $Q_i$  differs from ~~the~~ <sup>the</sup> other  $Q_i$ 's, and the

second part/operator on the  $Q_i$  to obtain a probability distribution over possible  $A_i$ .

Often several  $Q_i$ 's will be recognized as equivalent, since they all use the same

stochastic function to obtain the associated  $A_i$  distribution.

For a particular  $Q_i, A_i$  it will often be possible to find more than one  $R_i, O_i$  pair

that gives a larger  $a_i$  value in eq. .12.

These  $R_i, O_i$  pairs will have high a priori probabilities (short ~~and~~ <sup>i.e.</sup> descriptions) and assign large probabilities to ~~the~~  $A_i$ .

One way to find such  $R_i, O_i$  functions might ~~be~~ <sup>be</sup> via Genetic Programming methods of mutation and crossover of ~~previously~~ <sup>previously</sup>  $R$  and  $O$  functions that have been previously found useful. At first glance, each  $R$  candidate would have to be tested on all previous  $Q_i$ 's. While this is indeed true, it will be found that certain functions are able to ~~reject~~ <sup>reject</sup> large subsets of the  $Q_i$ 's.

$R_n^k, Q_n^k$   
Finding  $R_n^k$  may not be easy: it has to be checked on all previous  $Q_i$ 's.

A small number of functions of this sort can be used by many different  $Q_i$ 's to reject  $Q_i$ 's not identical to themselves. It will usually be necessary for a new  $R_i$  to be tested on ~~at least~~ only a small number of  $Q_i$ 's for which the standard rejection functions are inappropriate.

**SN** In 394.11 if we ~~sum~~ <sup>sum</sup> over  $P_j$ 's weight "v.g." indicator (By ~~the~~  $\chi^2$  criterion) If we sum ~~of over~~ only some of the  $P_j$ 's Ith part  $\leq 13$ , T. classer weights 394.11 (w.  $\chi^2$  over), so 4. smaller ~~is~~  $\chi^2$  weight!

$\chi^2$  is ~~is~~ but by something like  $\ln \left( \frac{\sum_{j=1}^J 2^{n_j} \frac{1}{2^{n_j}}}{\sum_{j=1}^J 2^{n_j} \frac{1}{2^{n_j}} \text{ (fruits)}} or perhaps  $\ln \prod_{i=1}^n P_{\text{true}}(A_i | Q_i)$   
 $\ln \sum_{j=1}^J 2^{n_j}$   
 over by  $2^{n_j}$ 's Real  
 It's build  $R_i$  for.$

perhaps include this disc. in the report!  
 I think it is "implied in my discussion in the Bare paper" but not as clearly as it might be.

**SN** In report: Include a  $\chi^2$  giving to Heuristic Basis of TM; T. idea test we have ~~prob~~ a TQR - we have ideas on how human write solve to TQR: we have a measure of adequacy of Plot "explain" & if it's any good, we know how to express it exactly using ALP - if not, ~~the~~ ALP gives ideas ~~on how to~~ ~~repair~~ on just where the model is deficient, ~~it~~ <sup>this</sup> obtain suggests how to repair it.

Perhaps also explain "Wittgenstein of 2 of optimum!"

**SN** (or value A fixer) eq. 394.01, 11 give sol ~~to~~ <sup>its</sup> properties. Discuss correspondence betw. Prem.

For the model of ~~the~~ <sup>394.11</sup> there is an analogous theorem that guarantees the rapid convergence of all to correct answers. Perhaps mention more details on how any  $P(A|Q)$  induces a rd branch of the bits of every A, & that both ~~the~~ 394.11 & the source "P(A|Q) do this: so we can ask how fast ~~the~~ bit probs ~~depend~~ <sup>depend</sup> on ~~the~~ answer is "like" "SRTS". See of.

A bibl. review of stuff on Ob, Op ideas to implement 388, 00ff!

- Assemblation  
 1) 388.00 - 388.13 ; 330.00 - 40 ; 396.18 ; 397.11 - 140 : 394.06 - .20 + ~~the~~ (394.01 - .15)R ; 261.08 - 18

One thing about Obs: very often they are hierarchical, so a few are used to "narrow down" what the problem is, & then for further obs are selected - so determine more exactly what the problem is (i.e. its soln).

388.06 - 09 Here "Context" is defined by an Ob: It looks like the most general defn. of "context". It's an operational defn. of "context" Hvr. Discovering the correlates of a "Context" is diff: we have to try it on all previous data to see if it correlates w. any Play Useful. It seems unlikely that we would construct "context" obs ~~at~~ at random to test them for "usefulness" - T. test is really too expensive. 388.10 suggests some useful contexts.

**SN** On overall format Paper: In discussion in introduction: Descr. problem: Descr. IA's: Mention that GA, ANN, ALP <sup>complexity based</sup> (Judicial), or SV machines ... can be used to create IA's that will solve P's problem. "We will mainly discuss use of CBI!"  
 Also mention that writes: Max  $\sum_{j=1}^J P_j(A_j | Q_i)$  is v.g. being to optimize Max w. (limited times it is not the best way, T. "Best way" seems related to Subscript & reciprocal Algms (= Obs)

00:396.40

On context: context can be sharp or fuzzy. Contexts are used in (at least) deterministic primitive stochastic

2 ways: As recognizing (Obs) to condition (SOPs) for reducing As contexts far obtaining PCs of events used for concl. construction

A PC evaln. of these concs. My hobby BIG expansion of t. idea of context was 300.26 ff. (300.26 - .33 in particular)

They should be tried op. cross rate to most recent idea of context 388.06 (context = OP)

300.26 - .33 considers "Context" to be the hyper level problem of which t. present problem is a part - as well as the hyper level problem above that, and so on ...

The gen. of "Context" at 388.06 was that Context (operationally) could be any OP!

ie, a fund from things to Boolean (True false)

They would seem to be two quite different kinds of generalizations!

So, general outline again

Dom of BA problem explain why its "complete" Pattern of A's:

variable [A's can be passed. GA, ANN, SVM's, C.B.I. - we will discuss all of them]

Give genl soln. for  $CC=0$ : fall how good this is in papers data papers (form)

How can we approx this in computable soln?

1) start w/ large no. of  $P_{ij}$ 's. Recn  $\frac{1}{2} \text{ent}_1 = 2n$   $P^2(A_{ij}^{(n)})$

Detect Seq w/ heap all  $P_{ij}$  w/  $t > \frac{w}{w_{max}}$ . (Seq  $w = 1000$ )

2) for each  $n$ , select subset  $w$  ( $w$  can = 1)  $P_{ij}$  of largest  $2n$

for each  $n$ , normally  $\{P_{ij}\} \subseteq \{P_{ij}^{(n)}\}$  will have few  $P_{ij}$  in common.

If  $n$  is large finding Recn  $w$   $P_{ij}$  will be very time consuming & time will grow.

3) Similar to 2) but given  $[P_{ij}]$  & mutual redundancy until we have now  $[P_{ij}]$

of small  $w$ , for top  $w$  there are more factors of  $w$  of peak wt.  $w$  can be as small as 1

2) w can be 1, but Recn mut. redundancy doesn't work well.

We will mainly discuss various versions of 3.

For most forms of  $[P_{ij}]$ , when we expand the subset function by mutual redundancy

We have to test each candidate function on the entire data set (GA) entire data set is known only

process become more and more time consuming. The time to learn a subset of  $n$  (examples) is proportional to  $n^2$

Each operator consists of a subset S that maps  $Q_n$  into a probability distribution on  $A_n$  and second a recognition part  $R_n$  that maps  $Q_n$  into a Boolean True or False. It True

$Q_n$  is transformed by  $S_n$  into the probability distribution on  $A_n$  & examined by  $R_n$  does nothing.

The other  $R_n$  could all be replaced by one to determine which  $S_n$  should be applied to  $A_n$ .

generate on



$$\int \ln x = x \ln x - \ln x + C$$

$$\Rightarrow x \ln x = \int \ln x + x$$

00: 397.40 Would this **not** take ~~time~~ time proportional to  $n$ , thus defeating our effort to save time?  
 Normally, it is possible to categorize the operators  $R_n$ , so that by using only a few categorization operators on  $Q_n$ , one can ~~delete~~ rapidly eliminate almost all  $R_n$ 's.

03: Rewrite @ 397.33 ff: One way to deal with this difficulty is to ~~compress~~ first find ~~one or more~~ <sup>probabilistic</sup>  $F_{n+1}^j$  such that  $P(F_{n+1}^j) \cdot F_{n+1}^j(Q_{n+1})$  is as large as possible —  $P(F_{n+1}^j)$  being the a priori probability of  $F_{n+1}^j$  and approximately equal to  $2^{-l(F_{n+1}^j)}$ ;  $l(F_{n+1}^j)$  being the length of the shortest program for  $F_{n+1}^j$ . But we ~~have~~ have been able to find. **(Explain how soon facts are found)**  
~~Once we have~~ <sup>Having</sup> found such a set of functions, we must find a way to recognize ~~them~~ problems like  $Q_{n+1}$ , so we will know that  $F_{n+1}^j$  ~~can be used to solve them~~.

11: At first we will try to find a function that narrowly recognizes  $Q_{n+1}$ .  
 When we have many ~~such~~ recognition functions for ~~each~~ of the problems in the set, we will try to simplify them, so that their total description length is reduced, and so it takes less time to ~~use~~ use them to determine what  $F$  to use to solve problems.

20: This is done by creating a hierarchy of ~~recognition~~ recognition functions.  
 So that at the lowest level, the recognition function will decide what ~~recognition~~ recognition functions need be used on the next level. Using a device of this sort, the time needed to recognize a particular  $Q_n$  may be of the order of  $\log n$ , and the time needed to recognize all  $F$  functions needed for  $n$  problems will be of the order of  $n \log n$ .

22: In general, there will be other means of consolidating recognition functions so as to reduce the total amount of time used in implementing them.  
 Explain tradeoff betw. CC & PC in (OSR). Also in finding good "Recogn Obs" in 11-22.  
 That this is not an "Adjustable Param", in the sense that one pursues it on the basis of past experience & "intuitive understanding". It is a User Input, which is part of the problem defn. Discuss how a "user input" could be derived by a "Reinforcement Machine" — but explain why this is Not Desirable.....

30: 0/14/01 The main bottleneck now, is how to describe "Updating": As of now, in most "el." form; It is in 2 parts:  $F_n^j$  &  $R_n^j$ . So I can first describe "first order"  $F$  &  $R$  functions. Next, discuss consolidation of both types (i.e.  $F$ 's consolidation w.  $F$ 's;  $R$ 's consolidation w.  $R$ 's.)  
 Consolidation of  $R$ 's means: say  $R_n(Q_n) = T$ ;  $R_{n+1}(Q_{n+1}) = T$ ; Consolidation:  $R'(Q_n) = 1$ ,  $R'(Q_{n+1}) = 2$ , etc.  
 Furthermore  $R'$  has shorter desc. than desc. length of  $Q_n$  + desc. length of  $Q_{n+1}$   $\Rightarrow T(R'(Q_n)) + T(R'(Q_{n+1})) < T(R_n(Q_n)) + T(R_{n+1}(Q_{n+1}))$   
 The desc. length of  $R' < t$ : total desc. lengths of  $R_n$  &  $R_{n+1}$ .  
 "Consolidation" of  $R$ 's &  $F$  functions (like Newton's integration of Galileo's other velocity).

Convergence thm for Condl. Probabs: 33

00:398140: We can start with a Maximally el. S function that is composed of a sets of R's & F's. We end up with a single S funct. that accepts any kind of Q, & gives some kind of A (The better A if larger CB is allowed.)

Essentially, what we want to do is compress S, & have it work faster: also, it should have same results for past corpus. Compression implies better induction via analog of ST8T3.

SO: T. First thing TM does when a new <sup>QA</sup> ~~number~~ comes in: 1 use old S on Q: see if U.G. A appears (= by pc). 2 If not (i.e. Q may have no output, or just a A of low P), it tries to find an F(Q) -> by A pc. More exactly  $P_F(A|Q)$  is large as poss.

Notation Note: Perhaps look for "Smaller" functions (= shorter length) rather than "shorter" functs. "Shorter" could mean shorter execution time (= faster) short term length

3 Next, it tries to find way to recognize a new Q.

Probably I should not go into much more detail on 'updating' until I start work on TSCQ. Tho I will do ANL, using the "new model of update"

SO: in f. report: Do, or - all: Tell how to get a set of  $R_n^i, F_n^i$  (possibly do include 4000-19 do. 10-19 is summary)

Then discuss "compression" of  $R_n^i$  in small sets; of  $F_n^i$  also in small sets.

If a set of  $n$  are unified (coalesced, integrated) then their convexity.  $R_n^i$  are also automatically integrated - Also to recognize a lpm when need to do "class work" (less cc).

But note 12 (in f. report).

One expository Approach: Just tell how to get  $F_n^i$ , then for as large a set of probs as possl.

If (When) it looks diff, do  $\{F_n^i\}$  for a new set of probs & find  $\{R_n^i\}$  for those probs: then discuss the usage of 397, 31 = still a mess

Then how "consolidation" of R's (a of F's) reduces time to  $O(N \log N)$  or less.

(see 398.03-22)

Say we have a new problem: several of R's that recognize it, but then F's fail.

After we solve the new problem, we have to modify the R's that recognized it for the new problem.

This can be done by having a new R that "dominates" all other (older) R's.

In fact I've written about various kinds of overlapping R's, trying to remove of various kinds. Review & try to consolidate that work: list various kinds of cases (situations of overlap, etc) 399.00

0/16/01 SN Qnt: Convergence Thm. for Conditional Probabilities: Int proof of this thm, I considered  $P_F(A, Q)$  v.s.  $P'(A, Q)$ .  $P$  being true cond. prob. &  $P'$  being ~~any~~ approx to it - or

directly any cond. prob. distrib. In s. proof I that I had to assume  $P'(A, Q) \geq P(A, Q)$ . This is not so.

The assumption had to do with ~~expected~~ Expected Values of Prob. rates  $(E \text{ wrt } P(A, Q))$ .

(Try to find previous Retrance). My that is still not clear on this! 4.25.00 we B. what is probly f. v. idea! We were interested in  $E \ln \frac{P'(X)}{P(X)}$ : which is the Kull-Leib. distance & is always  $< 0$  (spec) 4.25.00

400.00

00:40:1.40

[SN] This <sup>process of</sup> finding of  $R_n^j$  is usually incremental, but it ~~need not~~ be "easy"!

But Anyway! Breaking a problem into recognizing part & a soln. part does seem like a good, useful idea.

Mention use of ~~key~~-type GP as an <sup>Atte</sup> alternative LSRL (Give Pros & Cons)

Ordering Functions in PC order: In Lisp (or other functional langs), Re string itself, tells when a function is completed, but it may only be a "subfunction". [Lost at my ~~discuss~~ discussion of ] **IMPT!**

Also ~~the~~ ~~stuff~~ ~~I~~ wrote about using the form code to set strings in PC order - is it applicable in this case?

More generally, suppose we have been given n examples,  $\{Q_i, A_i\}, i=1 \dots n$ .

and we have ~~an~~ a recognition function  $R_n$  and a set of <sup>associated</sup> Conditional Probability functions ~~functions~~  $F_i^j(\cdot)$  for  $i=1 \dots m$ .

Usually  $m < n$ , since we don't always need a new  $F_i^j$  set of  $F_i^j$ 's for each new  $Q, A$ .

Wa a given a new problem,  $Q_{n+1}$ .

If  $Q_{n+1}$  is recognized by  $R_n$ ,  $R_n(Q_{n+1})$  will be some integer from 1 to m.

[SN] Could  $R_n$  output more than one integer? - At first average so that its only one integer. Later, we may improve it by allowing  $> 1$  integer & giving wts to  $F_i^j$  sets.

To define the system more exactly, suppose we have been given a set of n examples,  $\{Q_i, A_i\}, i=1 \dots n$ .

The system has a recognition function,  $R_n$ , so that for  $i=1 \dots n$ ,  $R_n(Q_i)$  gives a single integer from 1 through m.

The system also has m sets of  $[F_i^j]$  functions ( $i=1 \dots m$ ).

~~Each~~ possible output of  $R_n(Q_i)$  will activate a different set  $[F_i^j]$

If  $R_n(Q_i) = k$ , then the set  $[F_k^j]$  of probability distributions will be used to solve the problem,  $Q_i$ .

If  $k \neq i$  then  $k \neq i$  - i.e. for problems that the system has worked, there will be ways to

a set of  $[F_i^j]$  functions that will give a distribution for  $A_i$ .

If  $Q$  is a new question however, say  $Q_{n+1}$ , then if  $R_n$  recognizes  $Q_{n+1}$ , it will output an integer from 1 to m. The resultant probability distribution for  $A_{n+1}$  will be

$$\sum_j b_k^j F_k^j(Q_{n+1}, A_{n+1}) / \sum_j b_k^j$$

$$b_k^j = \prod_i F_k^j(Q_i, A_i) \quad \text{if } R_n(Q_i) = k$$

$C_k^j$  is the a priori probability of the function  $F_k^j(\cdot)$ .

$\sum_j b_k^j$  is a normalization constant

In the definition of  $b_k^j$ , the product is overall  $i$  such that  $R_n(Q_i) = k$ .

On the other hand, if  $R_n$  does not recognize  $Q_{n+1}$ ,  $R_n(Q_{n+1})$  will have no output.

In early stages of training of the system, we will stop at this point, and say that it cannot answer that question. In later stages of training, the system can try to modify

$R_n$  (usually by compressing it), so that it returns its responses to  $Q_i$  ( $i=1 \dots n$ )

Discuss Structure of  $R_n$  - as a subset of M-Resolution Modules.

00:402.40!

and also responds to  $Q_{n+1}$

Analysis of the updating behavior of the system, is more complex.

As before, say we have  $(P_i, A_i)$ ,  $i=1 \dots n$  and we have  $R_n$  and a set of  $[F_i^j]$ 's that are able to get a good value of  $A_{n+1}$

After the system has responded to  $Q_{n+1}$ , we (can) update it by giving it the correct ~~thing~~ string.

If the system gives an adequate probability to  $A_{n+1}$  conditioned on  $Q_{n+1}$ , we ~~have~~ do no updating, and go on to the next problem.

If the system either gives no response to  $Q_{n+1}$  or gives too small a probability to  $A_{n+1}$ , we first find a ~~new~~ set of  $F_{n+1}^j$ 's such that ~~their~~ ~~that~~ ~~gives~~ such that

$\sum_j c_{n+1}^j F_{n+1}^j(A_{n+1} | Q_{n+1})$  ~~is as large as possible~~ ( $\leq c_{n+1}^j$ ) is as large as possible.

We then modify  $R_n$  ~~to~~  $R_{n+1}$  so that  $R_{n+1}(Q_{n+1}) = k + \dots$  and other bits for  $i \leq k$  ~~are~~  $R_{n+1}(Q_i) = R_n(Q_i)$

i.e.  $R_{n+1}$  is the same as  $R_n$  except for its response to  $Q_{n+1}$ .

(Discuss  $R_n$  as composed of  $k$  modules: How some must be changed if they give poor responses.)

Discuss  $[F_{n+1}^j]$  discovery: (2) discuss  $R_{n+1}$  modification } via Lvsch  
via GP.

For my own use, try to go into more detail on ~~the~~ ~~update~~ ~~of~~  $Q_{n+1} = O_n(Q_n, Q_{n+1}, A_{n+1})$ , using usual Calc. ~~is~~ also on the more general problem of updating, in which we want to ~~get~~ ~~of~~  $O_{n+1} = O_n(Q_n, Q_{n+1}, A_{n+1})$ , using usual Calc.

Notation: 0/19/01

Instead of  $F_i^k(\cdot)$ , perhaps use  $P_i^k(\cdot)$ . (This was used as far back as 397.22) ~~is~~ ~~not~~ what I'm using now.

i.e. the way(s) that the new problem differs from the old, ~~can~~ ~~could~~ gives some clues on how to solve it.

is a general updating algm (A more general one would have more args. for  $O_n \dots$  e.g.  $O_{n-1}, Q_{n-1}$ , recent  $A$ 's, recent  $P$ 's of problems, ect... (see 2300.26 for some possys.)

$[F_i^j] R_n$  is one way to structure  $O_n$ , but not (most) general/best way.

Some started with a certain form for  $O_n$  as TM modules, it is able to do more general forms.

For myself, write a review (or just read all past) on updating. { Ist RE Franco ~~the~~ Pilot of PP all on "updating" }

[NB] M2 378.00-06; 379.21-24 aren't really so bad! I was concerned w. trials taking too long because they would have to do all  $(Q_A)_i$  for  $i=1/n$ : Could M2 use the ideas of M1 (i.e. Division into R & F functions) to ↓ verifn. time? Conds that took too long would be discarded, so we would end up mainly w. cond's that used ~~the~~  $M_1$  - type models. A trouble is, M1 doesn't actually test  $(Q_A)_i$   $i=1/n$ . It uses a logical argument to show they are 0's!

[SN] In report do say that humans do "recognize" a "new problem," i.e. recall all old ways of working old problems. What we want is a method that does this, but occasionally is able (needs to) make a ~~transition~~ ~~can~~ ~~apply~~ ~~to~~ old as well as new problems. That this "recognition" problem can be further "broken down" ( $\equiv$  a/d) into a recognition Operator for each problem type: we can have higher level recognizers for categories of problems & by level still to sets of cats, ect. This higher hierarchy speeds up & simplifies (lower b cost) & solves.

Do. say 32-38 in report - soon after phrase "updating idea" is mentioned

IP

$= P_S(A|Q)$

: On + input form of a universal P.D :  $P_{i2}$  is equiv. to a 2 input formate cond. pd.

ie. ~~input~~ <sup>input</sup> condition, other is "R". So R is ~~the~~ <sup>a</sup> dem of the output "in view of" Q.

For a 3 input set up, R is a dem of Q in view of Q is the general system dem, S.

For a system dem S,  $\Phi$ 's  $\Phi_1 \dots \Phi_n$ ; A's  $A_1 \dots A_n$  the dem of  $A_{i1}$  is

.04

One way to write reports: ① Dem of Problem, its universality

}  $\rightarrow$  4.10.13  
← this was writ here by mistake!

② Soln: Maybe % as I've been doing recently: w. little detail,

③ Dem of  $\Phi_n$  as sum of  $n$   $P_i$ 's is  $R_i$ 's.

~~④~~ ④ Updating: Allow dems:

⑤ ~~how~~ Dependencies: ② How pc's are input to functs

③ how L such is dems for Indredzin, for Inv probs.

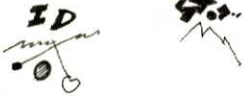
.10

.18

.20

.30





00: 407.40: The foregoing description of ~~the~~ updating techniques ~~is~~ for the early training of the system. The only method used to solve problems was Lsearch. The system was able to learn to solve complex problems, only if the ~~complex~~ ~~problem~~ solutions were obtainable by combining functions that were found useful in solving earlier problems.

This idea of ~~learning~~ + seq's should be mentioned earlier. Mentioned in 86's 89, but describe the build-up of concp.

In explaining "what is a problem?": Say best <sup>general</sup> optm. problems — optm. w. limited time & / or other resources is a possibl. kind of problem. Lsearch can do this, but is not necessarily a more efficient way to do it. (Woops! But it is "If all info is in PD, we have suitable lang. during Lsearch.")

perhaps follow ling(.00) <sup>w. it:</sup> The system is broken down into a problem recognition part and a problem solution part. While this is often a reasonable way to solve problems, it is not the most general way. After the system has had <sup>enough</sup> successful experience in solving optimization problems, we can give it as an ordinary ~~problem~~ "question" the problem of finding a set of ~~the~~ functions,

0.15:  $O_{n+1}$  so that  $\prod_{i=1}^{n+1} O_{n+1}(A_i|Q_i)$  is as large as possible — ~~what is the information~~ <sup>is as large as possible — what is the information</sup> ~~is as large as possible~~ <sup>of system functions</sup> that worked well for the first ~~n~~ problems.

0.17: "The sequence is given as Auxiliary information". Stated another way,

$$\{ [O_n^*] \} = [O_n^*] ( [O_n^*], Q_{n+1}, A_{n+1} )$$

0.17R: The concept of "Aux info" needs explaining: E.g. An encyc. or dictionary can be regarded as "aux info". This concept is better accepted ways to use Aux info: Refor complex; I'm not sure I have a good understanding of it. See 409. 22 for discussion. 409. 44 may be adequate soln.

0.21: Give us a distribution on  $O_{n+1}$  as a function of  $[O_n^*]$  the set of descriptions of  $O_n$ 's of much weight, and

the new problem solution pair  $Q_{n+1}, A_{n+1}$ , which are to be used to update  $[O_n^*]$  <sup>their descriptions</sup> ~~the set of descriptions of these functions~~

0.27: From .20 we obtain a distribution on  $O_{n+1}$ , which we can use as the basis of a set of trials. We use .15 to determine which, if any, of these  $O_{n+1}$  are good "system functions"

0.30: Q: So how does .27-.28 differ from Lsearch guided by  $[O_n^*]$ ? — it doesn't but  $[O_n^*]$  is supposed to give a v.g. pd, so Lsearch is very lost.

0.32: 0/23/01 "SN on Aux info" I haven't really solved this problem  $\in$  for limited CB. (I'm thinking of t. first/compression/induction problem. — not for 02 probs.). For induction problems or CB =  $\infty$ , <sup>total</sup> ~~total~~ <sup>maximally</sup> ~~maximally~~ Over simply/compression Main info + aux info. For CB  $\infty$ , One only compresses and of t. total corpus to give good prediction — explain. of a particular part of t. total corpus. I haven't really solved this problem: One soln. of it is to break corpus into  $\infty$  parts & compress t. part most relevant to the extrapol. problem. I've written a lot on this since I <sup>spac</sup> ~~spac~~ 409.00

ID

Bathroom  
Bethoven

$$10 \times 10^2 \times 2 \times 10^7 = 3 \times 10^{17} \text{ seconds} \\ \times 10^9 = 3 \times 10^{26} \text{ ns}$$

00: 408.40: first discovered problem (at ID51A). I may have called it "partitioning problem".

Thinking about it now, it would seem that one would want to compress parts of the corpus in varying amounts.

This would give a set of concs that were very useful in expanding to specific areas of interest.

It's that "set of concs" that seems most important.

Being able to "read" an encyc. means that one can compress it very rapidly (low cc).

But the idea of "indexing" an encyc. so one wouldn't have to read so much of it --- (how does that fit in?)

→ I think the "soln" I obtained was that one would discover certain methods for doing induction.

These methods might involve "how to choose a subcorpus".

One learns how to do good induction.

This "soln" puts the problem at a higher level, does not eliminate it. But one might have more

ideas on how to deal w. the "higher level problem" - it may be in a more stylized form than

diverse forms of the

lower level problems. Still it is Quisenberry's "Learning How to Learn".

The "higher level" could be Org. Evoln!

RLP

Actually, the "second level" can use the older form of Resource Limited ALP: i.e. "The best codes for

entire corpus obtainable within given C & B" This may not be strictly best, but it's good enough.

We can have much interchange of info betw. first & second level: [Impact of second level can be a regular problem for the first level.]

Earlier, I thought of this "second level" in another sense: i.e. E-R.W. can be this "BIG" corpus, but to make progress, we select best observations/experiments. [In physics, R.R. is a ditty, hvr: ]

In high energy experiments, say, the predictions we make are mostly for the experiments & for future experiments. - rather than for answering present questions.

Re: the "Aux Induction" problem at 408.17 & 32, 408.(17R-21R)

However, one can learn from experience (level 2 & 3) that under certain conditions one has good chance of getting results inconsi. w. previously accurate predn. methods (i.e. "theories").

A mode of living that I've not studied, is that of studying works, theories, experiments of others. Reading history of science - just trying to understand certain theories & o experiments arose, could be a useful source of "understanding" for TM. I assumed TM would read lots of good stuff, but I hadn't considered how it would make use of it.

So: up to 407.40; we have a "first order approxn." for TM.

407.15-40 is an improving recognition function, so it doesn't take a lot of Kn operations to do it. It's problem.

While taking that long would be very bad, it would be tolerable. - say, if I wanted to demonstrate

a minimal TM that could work interesting problems.

(407.03-14: this does all updating operations: I guess I need all of them: probably even for Minimal TM.)

Improvements beyond Minimal TM:

1) Long during Lsize (this may be needed in even Minimal TM) - (More detail on this is done) → 410.36

2) Improvements in Recognition Function (of 407.15-40)

3) Perhaps some improvements to updating operations

4) Use of context to assign pc's to concs (may be needed in Min TM). I should write this up more clearly for my own use. This could result in final eqn. 408.20 (see 300.26 ff)

Context probability based on context = "conditional probab" = what QATM normally does; → 410.32 410.00

Details on this is to be done

spec 399.32

.00: 399.40 Various Responses of "Set of R's" to a new QA:

ray sol.

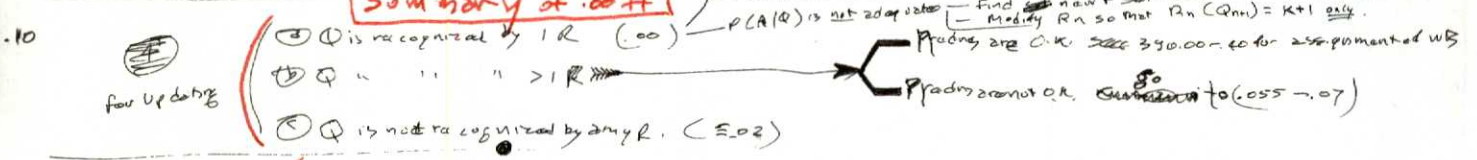
① (TUNNEL) ~~Q~~ Q is recognized, resultant A has Adequate PC: → Go on to next problem or do "SI".

② (Normal problem) Q is not recognized: We find some F → F(A, Q) is by: then find R to recognize that Q.

③ Several R's recognize + Q: 399.00 - to analyze R's situation. If + resultant predn gives adequate pc to A, then O.K., stop if not, goto ②; Perhaps try to find R that

.05 In addition to finding good R for + F we've found. Try to Modify the R's that didn't

.07 F's that didn't work well. Modify the R's w/o F's. (see 399.28-30)



- For predn
- ① If Q is recognized by 1 R use predn of F.
  - ② " " " " > 1 R use wts of F's (when available) ← (via 399.00-40 if available)
  - ③ " " " not recognized by any F. Either no predn, or try to modify (compress) ←

Some R's so Ray includes present Q. Alternatively, All R's give same wts to all Q's which is a Model of QATM that I've not investigated (Is Mindful of Self's "pandemonium")

.19 Another, more conventional way to deal w. ③ is to compress various R's that seem to be close to + Q. Since by compression, we ↑ entropy? (Entropy more?)

.20 Good Office: 397.12 - .32 (33 ff is one way to say it. 398.03 suppose another (better?) way to do this. deriv. Pns. abcd efghi

.21 O.K. so: Start w. "CP" paper: Use first 3 lines (is corrected) has will first give a theoretical solution to the problem, involving infinite time and storage. We will then discuss various physically realizable approximations.

For data set described above, the probability distribution on An is

.26 
$$\sum_{j=1}^{\infty} a_0^j \prod_{i=1}^{n+1} P^j(A_i | Q_i)$$
 (use  $a_0^j$  instead of  $P^j$ )

Here  $P^j(\cdot)$  is the jth possible Conditional Probability distribution relating two arguments,  $P^j(A_i | Q_i)$  is the probability of  $A_i$ , given  $Q_i$ . Since total recursive distributions of this sort are not effectively enumerable, we include in our sum, all partial recursive functions of this sort --- which are effectively enumerable.

.30  $a_0^j$  is the a priori probability of the function  $P^j(\cdot)$ . It is approximately  $2^{-l(P^j)}$  where  $l(P^j)$  is the length in bits of the description of  $P^j$ . APPENDIX I discusses computation of the  $a_0^j$  (this, in language "SAL").

.33 We can rewrite .26 in the form 
$$\sum_{j=1}^{\infty} a_n^j P^j(A_{n+1} | Q_{n+1})$$
 Here  $a_n^j = a_0^j \prod_{i=1}^n P^j(A_i | Q_i)$

In .33 the distribution of  $A_{n+1}$  is a weighted mean of all of the  $P^j$  distributions --- the weight of each  $P^j$  being the product of its a priori probability and the probability of the observed data, in view of  $P^j$ .

The distribution .26 (or .33) is known to be very accurate. If we use the

ID

9

Express all of our a priori information about the data, this distribution is perhaps the most accurate that could be obtained ~~etc~~ by any means.

Since we can not compute ~~the~~ infinite sum using finite resources, we ~~try~~ try to approximate it by using a finite number of <sup>large</sup> terms — terms that in 400.33 have ~~maximum~~ maximum weight. ~~large~~  $2^j$  values.

While it would be <sup>ideal</sup> to include ~~the~~ terms of maximum weight, ~~it is known that it is impossible~~ to know if a particular term is of maximum weight. ~~The best~~ It appears that the best we can do is find a ~~set~~ <sup>set</sup> of terms so that their total weight is the maximum we could find ~~in the available~~ <sup>in the available</sup> time and memory limits.

The trick of doing good prediction seems to be that of finding the "heaviest" ~~possible~~ <sup>possible</sup> heaviest functions possible ~~within~~ <sup>within</sup> whatever time we happen to have.

Using whatever limited computation ~~power~~ <sup>power</sup> within ~~our limited~~ <sup>our limited</sup> time and memory limits that we have

the more time and memory we have, the better are our probability estimates

in general, ~~the more~~ <sup>the more</sup> better probability distributions as our time and memory resources are expanded.

our probability estimates ~~will~~ <sup>will</sup> improve



abcdetg

SN

One way to phrase the "problem of induction" is to find in ~~the~~ <sup>the</sup> given resource constraints,  $C$ , ~~such~~ <sup>such</sup> that the total weight in 400.33 is maximum. ~~Viewed in this~~

Way, induction is a ~~practical~~ <sup>practical</sup> (constant) limited optimization ~~problem~~ <sup>problem</sup>

In 1972 ~~Levin~~ <sup>Levin</sup> suggested (Lev 72, LiV 93, LiV 97) a ~~method~~ <sup>method</sup> of solving sequential search problems that can be readily applied to the ~~the~~ <sup>the</sup> problem of finding a set of function of maximum weight.

If ~~the~~ <sup>the</sup>  $Q$ 's and  $A$ 's have short descriptions and  $n$  is small, this can be a ~~practical~~ <sup>practical</sup> method to solve such problems.

Mania "exponential in n"

For ~~the~~ <sup>the</sup> complex problems and on large values of  $n$ , the method ~~is no longer~~ <sup>is no longer</sup> can no longer be used to find "heavy"  $F_j$  functions directly, ~~but~~ <sup>but</sup> it can be used as part of a more complex ~~method~~ <sup>method</sup> (Make 83 new sentence)

Let us start with a single  $Q_1, A_1$ . Using Levin's search technique (which we will call "L search")

We find ~~for function~~ <sup>for function</sup>  $F_1$  such that  $F_1(A_1|Q_1)$  is as large as we could find in the available time.

When we are given  $A_2$  and  $Q_2$ , we try ~~the~~ <sup>the</sup>  $F_1$  function ~~on~~ <sup>on</sup>  $A_2|Q_2$ . If  $\sum_j F_j(A_2|Q_2)$  is acceptably large, we go on to  $A_3, Q_3$ . If ~~not~~ <sup>not</sup>, we ~~try~~ <sup>try</sup> use L search

to find a new set of functions  $[F_2^j]$  such that  $\sum_j F_2^j(A_2|Q_2)$  is as large as possible. Having found these, we look for a function  $R_2$  that can recognize that  $Q_2$  is different from  $Q_1$ .

We use L search to find a function  $R_2$  such that  $R_2(Q_1)=1$  and  $R_2(Q_2)=2$  with  $b_0^2$  (which is the description length of  $R_2$ ) as small as possible.

SN

for  $R_2$  we want both  $b_0^2$  & small CC for ~~the~~ <sup>the</sup> testing  $R_2$  on all  $Q_i$ . The ratio  $\frac{b_0^2}{CC}$  is one poss. ~~cost~~ <sup>cost</sup>, ~~it is not~~ <sup>it is not</sup> ~~needed~~ <sup>needed</sup> to ~~test~~ <sup>test</sup>

$C_n = P_n(O_n, Q_{n+1})$

00: 403.40 : try to get 403.12 (A fairly Genl. Soln.) early in the report.

01 **SN**: T. system as described will not try to make "new discoveries", in the sense of requesting new physical experiments. This has to do w. "Look ahead" it will be my job to discover new things, inventing new fields of M.S.A. only if we give T.M. suitable goals.

Do mention - at end of Report: - That this system is easily adapted to solving problems of that sort - i.e. it can be given goals of a suitable sort -

After Derby On as being composed of  $R_n$  is  $[[P_i^j]]$  (sets of sets),

Explain that this is only patching str. for  $O_n$  - Not in general, we want 1. form of  $O_n$  only to extent that we want 4.2C (i.e. solution to n) to be as large as possible. [Get notation straight out!]

Some likely <sup>early</sup> ~~improvements~~ ~~enhancements~~ of  $O_n$ : "cats of cats" - See 403.32-38

If one tried to optimize  $O_n$  directly for each n,  $CC$  is exponential in n.

Using  $R_n$  broken in to n parts makes  $CC$   $O(n^2)$ ; or  $n \log n$  for recognition alone. hierarchical hierarchical format  $R_n$  reduces  $CC$  to  $n \log n$  or less.

Since we find  $P_i^j$  by Lsach, it will tend to be short (heavy, hywt, hypL) so it will tend to extrapolate well.

Go back to 401.20 (on Lsach) [Put Lsach in an appendix.] - Discussing  $CC$  for induction problems.

Notation  $F^j \rightarrow P^j$  not so good! I'll want  $[P_n^j]$  such,  $j=1 \dots K_n$

use  $[O_n^j]$  for set of operators so  $[P_n^j]$  would be OK. to use: it would be ~~conflicting~~.

This conflicts w  $[P_i^j]$   $i=1 \dots n$  for K different kinds of problem solns.

SN If we use GA; for obtaining rd.  $O_{n+1}$  from  $P_n(O_n, Q_{n+1})$  we could find (units, cross) such that  $Q$  + cands would tend to produce solns of  $(Q_{n+1})^n$  incorrect. - The still one would have to spend too much time testing them (would random choice of  $Q_{n+1}$ ; help much in this?) Anyway - it does involve learning in mut/cross etc. - so we would have a higher order GA doing that.

So mention Lsach for finding  $O_n^j$   $O_n^j = P^j$  of 400.26

so  $[O_n^j]$  is the set of  $O_n^j$  (for a given n) that were found by  $z_n^j$  in 400.33.

That this works for small n (time is exponential in n).

31 **401.25** One way to solve this ~~sort of~~ problem is to have  $[O_n^j]$  (for fixed n) broken up in to sub-functions.

32 We will do this by having two kinds of functions: first,  $P_n^j(\dots) \sum_{k=1}^n p_k(A_{*i}, Q_i)$  produces an acceptable probability for a certain subset of values of  $i$  ( $2 \leq i \leq 4$ ).

33 Each value of  $i$  will have an associated value of  $k$ .

Or: Each value of  $i$  will have an associated 'form'  $k(i)$ .  $k(i)$  is the "kind of problem that  $Q_i$  is".  $k(i) \leq k_{max}$ ; usually  $k_{max} \ll n$ .  $k$  will be  $\ll n$ , since there will be many "kinds" of problems (sub-problems).

00: 404.40: We will have  $k$  different functions that recognize the different kind of problems,  $Q_j$ .  
 $R^i(Q_j) = 1$  if  $R^i$  recognizes  $Q_j$   $R^i(Q_j) = 0$  otherwise  
 $R^i(Q_j) = \delta_{ij}$  (is this the data of  $\delta$  func?)

02 Associated with each  $R^i$  is a set of functions that solve the problems that  $R^i$  recognizes.  
 03 If  $R^i(Q_j) = 1$ , then  $\sum_j a_j^i P_j^i(A|Q_j)$  will be the distribution function induced by the

If  $R^j(Q_j) = 1$  and  $R^k(Q_j) = 1$  then  $j=k$ . So the  $R$  functions do not overlap.

08 Set of functions  $P_j^i$ , on possible answers,  $A$ , for question  $Q_j$ .  
 $R^i(Q_j)$  (i.e. from values of  $P_j^i$  and  $\sum_j a_j^i$ )  
 404.31 perhaps start out by saying: (404.31) One kind of function recognizes what kind of problem  $Q_j$  is. For each kind of function, we have a set of functions  $P_j^i$  that solve problems of that kind.  
 For each value of  $n$ , there will be  $k_n$  different problems. Usually  $k_n \ll n$ , since many problems will be of the same kind.

14 There are  $w_n$  recognizable functions  $R^i$  ( $i=1 \dots w_n$ ). Each of them recognizes one and only one kind of problem.  
 $R^i(Q_j) = 1$  if  $Q_j$  is an "i" kind of problem.  
 $R^i(Q_j) = 0$  if " " not an "i" kind of problem.  
 The  $R$  functions do not overlap. If  $R^i(Q_j) = 1$  and  $R^j(Q_j) = 1$  then  $i=j$  and  $j$  must be identical.  
 17  $\rightarrow (.02 = .03, .08)$   
 19 There are two critical questions about this induction model:  
 20 How do we use it for prediction?  
 21 When we find out the correct answer,  $A_{n+1}$ , to question  $Q_{n+1}$ , how do we repair that the system will respond better to future problems?  
 i.e. How do we update to sets of functions  $\{R^i\}$  and  $\{P_j^i\}$ , so the system is applied to (to test sum)

14R From time to time, new  $P_j^i$ 's are found. It's not unreasonable to work many old problems as well as new problems. When this occurs,  $P_j^i$  can be updated. It becomes possible to use a smaller  $k_n$  than before, a smaller no. of  $R^i$  functions, and a smaller no. of  $P_j^i$  functions. So the over-all system can be much more "compressed" yielding better expected accuracy for its generalizations.  
 21R We try to update  $P_j^i$  on other  $Q_j$  that are "similar" (have common indicators) to  $Q_j$ . That's what we want to do for.

22 When a new question,  $Q_{n+1}$ , is given to the system, the system must apply each of the  $R^i$  functions to  $Q_{n+1}$ . If one and only one of these functions has output 1 to  $Q_{n+1}$ , then the distribution probability induced on the possible answers,  $A$ , is given by 03 - i.e.  
 $\sum_j a_j^i P_j^i(A|Q_{n+1}) / \sum_j a_j^i$   
 $j_i$  is the number of  $P_j^i$  functions for this  $i$  value.

30 If  $R^i(Q_{n+1}) = 0$  for all  $i$ , then the system is unable to answer  $Q_{n+1}$ .  
 If  $R^i(Q_{n+1}) = 1$  for more than one value of  $i$ , we can use a mean distribution function induced by the  $P_j^i$ 's of all of these  $i$  values.  
 $\sum_i \sum_j a_j^i P_j^i(A|Q_{n+1}) / (\sum_i \sum_j a_j^i)$

33 When updating the system, we are given  $Q_{n+1}$  and  $A_{n+1}$ . If  $R^i(Q_{n+1}) = 1$  for one and only one Recognizer function, and  $\sum_j a_j^i P_j^i(A_{n+1}|Q_{n+1})$  is "large enough", then no updating will be done.  
 If any of the conditions of 33 forgoing conditions do not hold, we will update. This will, in all cases, involve a new set of  $P_j^i(\cdot)$  functions such that  $\sum_j a_j^i P_j^i(A_{n+1}|Q_{n+1})$  is as large as possible.

34

an adequate cond could be violation of any formal 406.33.

00: 406.40: If  $R^i(Q_{n+1}) \geq 1$  for 1 and only the recognition function, then we do no further updating, and go on to the next problem.  
If  $R^i(Q_{n+1}) = 0$  for all  $k_n$  recognition functions, we must find a new  $R_{k+1}$  such that  $R_{k+1}(Q_{n+1}) = 1$  and  $R_{k+1}(Q_i) = 0$  if  $i \neq n+1$ . We do this ~~again~~ again by an Lsachover suitable functions.

- 03 Update conds:
- 1)  $R^i$  is unique. Prodn is adequate  $\rightarrow$  no update
  - 2)  $R^i$  " " Prodn not adequate  $\rightarrow$  update? update  $R^i$  (for a single  $R^i$ )
  - 3)  $R^i(Q_{n+1}) = 0$  for all  $i$ : a) update  $R^i$  b) find  $R_{k+1}$  of suitable ~~prop~~ properties via Lsach.
  - 10 d)  $R^i(Q_{n+1}) = 1$  for  $\geq 1$   $i$ : a) If one of ~~the~~  $i$ 's that recognize  $Q_{n+1}$  gives adequate P for  $A_{n+1}$ ; modify ~~the~~  $P$  for  $A_{n+1}$  - recognizer that gives best P. Modify other recognizers (via Lsach), so they don't recognize  $Q_{n+1}$ .
  - 13 5) b) if none of <sup>nized</sup> recognizers  $i$ 's gives a adequate P for  $A_{n+1}$ , modify ~~the~~  $P$  for  $A_{n+1}$  recognizers so that they don't recognize  $Q_{n+1}$ , find suitable  $R_{k+1}$  we are now in situation 3) so continue as 3).

15 If there is time available, some of it can be spent "improving" the  $P_i$  sets of  $P_i^j$   $R^i$  functions.  
We may be able to find new functions that solve what have been ~~previously~~ regarded as different kinds of problems - so they are recognized by different  $R^i$  functions. If, for example, we find a new  $P_i^j$  that can give adequate solns for problems of  $R^2$ ,  $R^7$  and  $R^9$  kinds, then we will use this single solution type for all of these types of problems. This integer "consolidation" of  $P_i^j$  functions ~~usually~~ gives code compression, which makes it more likely that the new  $P_i^j$  function will generalize well for new problems.

In addition, we will ~~not~~ even want to integrate  $R^2$ ,  $R^7$  and  $R^9$  into a single function that recognizes all 3 kinds of problems. - which usually gives <sup>code</sup> further/compression. .247

24 Another way that we can improve the performance of the system is by expressing the  $R^i$  functions in hierarchical form. ~~One~~ One way to do this would be to have a recognizer that can recognize a function that can quickly decide if a problem is of kinds 1, 2, 5, 7, or 8 as opposed to 3, 4, 6, 9 or 10. This would ~~not~~ have the time spent in testing new  $Q_i$ 's to determine what kinds of problems they are. This ~~same~~ technique can be used again to further ~~subdivide~~ subdivide the ~~categories~~ categorizations of that first level. ~~at categorizations~~ - ~~By~~ dividing the recognition time by perhaps 2. ~~By~~ using higher order categorizations, we can have a solution time for each problem of the order of  $\log n$ .

30 .24 ff could be ~~more~~ explained better: I will start by saying that a system as described takes ~~time~~ time of about  $\frac{1}{2} n^2$  to do  $n$  problems, since working to solve  $k$  problem requires ~~the~~ the execution of  $k_n$  recognition functions. If  $k_n$  is proportional to  $n$ , this gives a time to solve the  $n$ th problem ~~proportional to~~  $n^2$  - Much too slow! We can improve the <sup>speed</sup> performance of the system considerably by .247 ff

By using recognition functions that ~~can~~ recognize several <sup>different</sup> kinds (or categories) of problems we can ~~to~~ speed up the system further.

ED

00:409.40 → ⑤ Any aspect of the Model Matrix is "improvable" can be ~~dealt w.~~ dealt w. by giving it as a problem, to TM. (but the Q is: Just what aspects do need improving?)

N.B. ⑤ (C.00) is a critical feature of any Induction Algm that can become "infinitely smart".

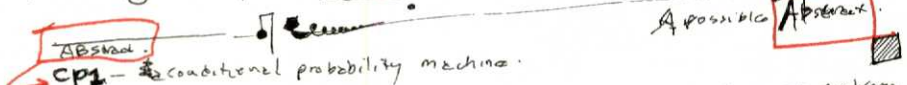
① & ④ may be adequate, but ⑤ adds nice "overkill" (No. ⑥ is also imp't)

⑥ Look ahead: "Cure Cancer" ect. Look ahead can signify ↑ problems solving speed.

SN By using these  $P_i^j$ 's in I) (unord'd) & the  $R_i$   $R_i^j$ 's in II), the shanty of overhead for concs. becomes more reasonable (perhaps): It does seem to be quite different from using a single (complicated)  $O_n$  function & obtaining p.c.'s of individual concs. within that  $O_n$  function. (The perhaps it ends up the same - if the  $O_n$  funct is a simple function of the I)  $R_i^j$ 's &  $P_i^j$ 's.)

But do examine this!

Copy this off - it here! ≡ "How to write Report"



As a rather general form of the inductive inference problem, we are given a sequence of Q, A pairs. Then given a new Q we are required to find a probability distribution on all possible A's. The Q and A's may be thought of as questions and answers, or problem descriptions and problem solutions, or the input and outputs of some known stochastic process. The Q A's in the set need not be of the same "kind" and may be drawn from very different domains of science, such as Linguistics, Economics, Physics, Biology, ect. Our system solves this problem using a set of expert functions  $O_n$  that give the probability distributions for  $A_i$  as a function of  $Q_i$  for the first n Q A pairs. The  $O_n$  are chosen so as to minimize the length of the bits of the program that describe the first n Q A pairs. In the initial part of the system's training, the  $O_n$  consist of recognition functions that recognize the early for small variations of the  $O_n$  consist of stochastic functions that relate the Q's to distributions on corresponding A's. The recognition functions and the stochastic functions are updated, when necessary, by using Learning Universal Search, to obtain functions with short descriptions that describe the data well. Functions of this sort are known to give good extrapolation of known data. For rapid learning, the sequence of Q, A pairs must be carefully designed. We will discuss some weaknesses of the described system, and how they can be overcome.

The Introduction goes over the material in abstract in more detail. (as examples of Q A's)

409.40 spec part of "Context" is "kind of problem" - This will turn involve categories & sub-categories of problems used in the Recon. funcy. We may be able to get QATM to work on the construction of formal functions for  $R_i^j$  &  $P_i^j$ , at a relatively early point in TM's training.

409.35 → Lung during lunch: Here we can regard this as a kind of QATM problem - But I'm not sure as to just how. Perhaps work out some more detail on how the lung is done & then see if I can express it as a QATM problem. One part of the trick is to suitably index previous cards (≡ mails) that are in this search. T. Q A's into be [cand; Genc;] - We are trying to quickly predict which cards will have by G. 411.00



ID

00:41.40: This seems like a familiar "old problem"! Occurred in analysis of GA. How do we get ~~xxxx~~ / languages that will enable good G' prodn.?

If a cond. say by G, then "near by" cond. will get approx by G. This is obtained by expressing the "control"

by G and, <sup>Co</sup> nearby cond. in 2 languages, select small "mutations" of Co hereby G because they and

Co can be relatively inexpensively produced from some set of grammatical constructs.

Howt. "near by cond." <sup>approx</sup> get by G is unclear

Wells, the "long during Lsrch" amounts to this: Say one is doing an Lsrch guided by a certain P.D.  $P_0(\text{cond})$  that was obtain from some Empirical data. After 1 or more cond. have been tried (tested, evaluated), we can use that that cond. data to augment the empirical data that  $P_0(\text{cond})$  was based on — giving  $P_1(\text{cond})$  — a new ~~xxxx~~ P.D.

00-05 is simply concerned w. how we get P.D.'s like  $P_0(\text{cond})$ ... i.e. what are some good ways to do this? It is, indeed, a major problem in analysis of GA.

I can, here, simply express it as a "Well defined problem" that needs a good soln.

If  $P_0(\text{cond})$  d.f. changes w. n (i.e. w. successive trials): I did have a way to do the "Lsrch" but while it seemed good, I had no "bound on soln time"

An Approach to Solving it: Look at some problems not solved/worked on by GA.

- a) Real number space: Optzn problems.
- b) Digital spaces " " "
- c) Mix Real, Digital aspects of decn. — optzn problem.

Try to make approx. models linking decn. of cond. to its Genc.

One way is to have  $P_0$  Genc Part's hard to <sup>or to calculate</sup> optimize. We make a simpler Genc  $G_0$  to original one, e.g. "pizza space"

Part's in relation to  $G_0$  is easier to find it - optimum.

It's "Easier to Calculate" is relevant to probs in which the Genc is very expensive (e.g.  $\$10^2$  use of the "pizza" p.p.m to evaluate cond. etc. — or use of a complex simulation p.p.m to evaluate a "boom balancing" p.p.m. I had the idea of using a "local model" — so as I moved up hill, I'd change models

In "Expert Systems" a kind of "Plan" was a simplified model of the system.

It was easy to predict/solve behavior of G. "Plan" model, & one had approximately that

behavior in the system: So the "Plan" system was good for generating "Conds"  $\equiv$  "Trials".  $\rightarrow$  41:20.5

30: 41.35! (On context): T. most. imp't idea here is that context <sup>induces</sup> the most General kind of Cond. p.c. — it is what QATM normally does. Consider the main problem Part ~~the~~ QATM works on: Given a new  $Q_i$  to find a good recognizer for  $Q_i$ : to find a good <sup>(stoch)</sup> function relating  $Q_i$  to  $A_i$ . ( $P_i^*(1)$ )

For the "Recognizer Problem": Perhaps the "context" is all of the previous  $Q_i$ 's & their Recognizers.  $\{ \text{the } Q_i, A_i \text{ pairs} \}$

N.B. W.D. context ( $\equiv$  cond. p.c.) we have UNcond. p.c. which is usually  $\equiv$  <sup>(or Lsrch)</sup> frequency of occurrence of various conds.

(In finding a  $P_i^*$  relating  $Q_i$  to  $A_i$  one big context is: all previous  $\{ (Q_i, A_i), P_i^*(1) \}$  pairs

These "Contexts" ( $\equiv$  cond. p.c.  $\equiv$  c.p.c.) give good <sup>guesses</sup> for  $P_i^*$ , which was the test.

T. c.p.c. is a good "Cond. generator."

1,520 in  
P. run in  
2000  
1997  
1 dur morning  
13 min.  
2.5K/monthly  
4.2K bills  
by auto  
accidents  
in US, in 1997  
1988: 49K  
5.25K in 12 mo  
etc.

Dat

0/25/01  
ID


2	1.7	412
3	1.4	
4	1.6	1.8
5	2.1	2.2
6	2.3	2.2
7	2.0	3
8	1.00	4.6
9	2.00	5.3
10	4.00	6

$34\sqrt{93}$     2.718    3+1+7  
 5.436    4.7    3+    1.1426  
 $\times 2$     34    1.3    1.09  
 $\frac{13}{34}$     34    03    1.23 + 1.137

00: 411+40: Wacker CPC's (= contexts) for  $P_i^j$  cards would be: Just old  $P_i^j$ 's. or just  $(Q_i, P_i^j)$  pairs,

→ Another part of context for  $P_i^j(Q_i)$  cards is the recognition function  $R^k$  assoc. w.  $P_i^j(Q_i)$  and, t.  $R^k$  assoc. w. t. present  $Q_i$ . The recogn. functs. tell important things about the  $Q_i$ 's.

They can classify t.  $Q_i$ 's in useful ways that can help find good  $P_i^j$  solns.

Context  $\equiv$  "side info"  $\equiv$  "Auxiliary info" 

05: 411.29 An implicit kind of extremizing problem is, of course MDL: This is a reform of game in itself.

(Just list cards in desc length): Butt. aux. conds for MDL are a simple, (say we want to solve  $k$  diffnt problems). Then one kind of game would be length of  $P_i^j$  +  $\alpha$  (no. of probs solved). This is a common trick in GA solns. Here,  $P_i^j$  GPS idea

of a Vector space is perhaps much better. It enables us to take a better idea as to "just what is wrong" w. a given card. — perhaps suggest how to fix it.....

→ 410.16-30: I abstract, covers the main content of the report: At the end, I will mention/discuss the "improvements" of 409.25-410.07: Also some stuff on tsg's: How to get TM to "understand" what a particular index means — for various arbitrary problem types.

So perhaps list various writing that I may want to use: in the "Modified".

- 1) Abstract: 410.16-30 → 413.00-40
- 2 ~~Production and Updating~~ 406.20-407.14 ← probably a summary
- 3 On recognition functions: (404.30-34) → 406.10-19 then ← 400.20-20 is an older outline.
- 2 Start main part at 400.21-401.14; then 401.20-25; then continue. ← see how reasonable this is!

- 1) Abstract 410.16-30 → [413.00-40]
- 2) 400.21-401.14 ← 3 lines of SAZ PPT ← 414.00-11 ← 414.00-415.09

Theo soln as  $\sum \pi P_i^j$  ← see the BPP for zip comp.  
 Finding the set of funds w max  $\sum$  is an NP-complete problem

401.20-25: Reference to Lsch to find solns  
 404.31 only 2 lines "Break  $[Q_i^j]$  into sub-problems"

406.10-19 = 406.10-17 {406.02-02; .02}

Note: underling means italics

Title: A Conditional Probability Machine.

ABSTRACT

As a rather general form of inductive inference problem, we are given an unordered set of  $Q, A$  pairs. Then, given a new  $Q$ , we are required to find a probability distribution over all possible associated  $A$ 's.

The  $Q$ 's and  $A$ 's may be thought of as questions and answers, problem descriptions and problem solutions, or as the input and outputs of some unknown stochastic process. The  $Q, A$ 's in the set need not all be of the same kind and may be drawn from very disparate domains of knowledge, such as linguistics, mathematics, physics, biology, economics....

Our system solves this problem using a set of system functions,  $O_j^i$ , that give the probability distributions for  $A_i$  as a function of  $Q_i$  for the first  $n$   $(Q_i, A_i)$  pairs. The  $O_j^i$  are chosen so as to minimize the lengths in bits of the ~~programs~~ <sup>programs</sup> that describe the first  $n$   $Q, A$  pairs.

In the <sup>initial</sup> part of the system's training, the  $O_j^i$  consist of  $w_n$  <sup>d.c.</sup> different functions that recognize  $w_n$  different kinds of  $Q$ 's, as well as  $w_k$  <sup>d.c.</sup> associated sets of stochastic functions that relate the various kinds of  $Q$ 's to probability distributions <sup>corresponding</sup> on the ~~associated~~  $A$ 's.

These recognition functions and stochastic functions are updated (when necessary) by using Levin's universal search algorithm to obtain functions with <sup>short</sup> descriptions that describe the data well. Functions of this sort are known to give good extrapolations of known data.

~~We~~ We will discuss the design of training sequences that facilitate rapid, ~~generalizable~~, learning by the system.

The system described is a very elementary one with several serious deficiencies. We will suggest ways in which these deficiencies can be overcome.

.00: 413.40

## Introduction

The problem we are addressing is that of very general probabilistic prediction:  
We are given an ~~unordered set~~ <sup>unordered set</sup> of  $(Q_i, A_i)$  pairs,  $i=1 \dots n$ . We are then given a new  $Q_{n+1}$ .  
The problem is to find a conditional probability distribution ~~for~~ over all possible  $A_{n+1}$ .  
(Here the "i" indexes of  $Q_i$  and  $A_i$  are used to indicate  $Q_i$  is associated with  $A_i$ .)

The ~~unordered~~  $i$  ordering, however, is not to be used ~~for~~ for prediction.)

The  $Q$ 's and  $A$ 's can be strings and/or numbers. The  $Q$ 's can be thought of as questions, the  $A$ 's as answers.  $Q_i$  might be the description of a problem, and  $A_i$  its solution.  $A_i$  ~~might~~ <sup>is</sup> ~~and~~ ~~the input to~~  $Q_i$  might be the inputs and outputs of an <sup>unknown</sup> stochastic process.

.10

Since we are considering stochastic models, the data can have "errors" in it —

.11

interpreted by the system as "noise".

We will first give a theoretical solution to the problem, involving infinite time and <sup>storage</sup> ~~computation~~.

then a discussion of various practical approximations

For the data described above, the probability distribution on  $A_{n+1}$  is

.18

$$\sum_{j=1}^{\infty} z_0^j \prod_{i=1}^{n+1} O_{\square}^j(A_i | Q_i)$$

Here  $O_{\square}^j(\cdot | \cdot)$  is the  $j$ th possible conditional probability distribution relating its two arguments.  
 $O_{\square}^j(A_i | Q_i)$  is the probability of  $A_i$ , given  $Q_i$ , as given by the function  $O_{\square}^j$

.20

Since total recursive distributions of this sort are not effectively enumerable, the sum is over all partial recursive functions of this sort — which are <sup>effectively</sup> ~~recursively~~ enumerable.

$z_0^j$  is the a priori probability of ~~the~~ the function  $O_{\square}^j(\cdot | \cdot)$ . It is approximately  $2^{-L(O_{\square}^j)}$ , where  $L(O_{\square}^j)$  is the length in bits of the description of  $O_{\square}^j$ .

(Appendix 1 discusses the computation of the  $z_0^j$ )

.30

We can rewrite .18 in the form

.31

$$\sum_{j=1}^{\infty} z_n^j O^j(A_{n+1}, Q_{n+1})$$

Here,  $z_n^j = z_0^j \prod_{i=1}^n O^j(A_i | Q_i)$

In .31, the distribution of  $A_{n+1}$  is a weighted mean of all of the  $O^j$  distributions — the weights of each  $O^j$  being the product of its a priori probability and the probability of the observed data in view of  $O^j$



.00: 414.40: The distributions  $414.18$  and  $414.31$  is known to be very accurate. If we use the  $2^j$  to express all of our a priori information about the data, this distribution is perhaps the most accurate ~~that could be obtained by any means~~ possible.

Since we cannot compute this particular infinite sum using finite resources, we will ~~try to~~ approximate it using a finite number of large terms — terms that in  $413.31$  have large  $2^n$  values. While it would ~~be~~ ideal to include the terms of maximum weight, ~~we~~ it has been shown to be impossible to know if a particular term is of maximum weight.

It appears that the best we can do is to find a set of terms of ~~maximum~~ largest total weight in ~~the~~ whatever time we have available ~~available~~.

see 416.00 -  
417.23 for digch of  
utility of L search inclusion  
(a  $\neq$  LNV-type L search).  
← 491.20

.09  
.10  
.11

Insert 423.25 - .27

In 1972 L. Levin (Lev 72, LiV 93, LiV 97) a universal method for solving

Sequential search problems that can be readily applied to the problem of finding functions of maximum weight. If there are a small number of simple problems, ~~if the Q's and A's have short descriptions and n is small,~~

Levin's method (which we will call "L search") can be a practical way to solve such ~~problems~~ them. a large collection of easier problems

For more ~~complex~~ difficult problems and/or longer values of n, the method is too slow to find "heavy"  $O^j$  functions directly. L search can, however, be used as part of a more complex method with a search time not much slower than linear in n.

.15  
.20

One method of this sort expresses the set of  $O^j$  being used, as a combination of two types of functions. The first ~~type~~ type of function,  $R^i(Q)$ , recognizes what kind of ~~problem~~ problem  $Q$  is. Associated with each ~~problem~~ recognition function  $R^i(\cdot)$ , is a set of one or more ~~functions~~ of the second type of function,  $P_i^j$ , that solves problems of the  $i^{\text{th}}$  kind. For each value of n, (the number of  $Q, A$  pairs ~~thus far~~ thus far), there will be  $W_n$  different kinds of problems. Usually  $W_n \ll n$ , since many problems will be of the same kind. ~~l.c.~~

.30

There are  $W_n$  recognition functions,  $R^i$  (if  $1 \leq i \leq W_n$ ) Each of them recognizes one and only one problem type.  $R^i(Q) = 1$  if  $Q$  is of type  $i$

$R^j(Q) = 0$  if  $Q$  is not of type  $j$

The R functions do not overlap:  $R^j(Q) = 1$  and  $R^k(Q) = 1$  implies  $j = k$

Trouble on 401.20: How L such is to be used for <sup>-finding good</sup> induction codes is unclear

I think I did write a lot on this at ISSIA month.

For induction problem: finding short code for a corpus; one way is to make trees in PC order, using CB = pc.T. When all have been run; try T <= 2T:

Or use parallel E L such:

In present case, we can only test candidates that are complete. — some list

If we have k primitive symbols — each w. its own pc, we can use Huffman codes for each of them: Then consider all binary codes of length L: They will represent <sup>Huffman</sup> ~~by Huffman~~ coded symbols seqs.

We test all terminating sequences for Time T. WOOPS! When we define a new symbol, (or even use a symbol) the pc's of other symbols change! — so we have to recode w. Huffman! Is this practical?

for a given value of T, we test codes w. T.pc >= E: E is time of one clock or some small time.  
for a given T, -T.pc >= E -> pc >= E/T so we test codes of length L < -log2(E/T)

Doubling T means we consider binary code codes 1 bit longer.

In using ~~SAZ~~ SAZ to generate codes, I had a mechanism to and terminate functions: Perhaps it is unnecessary! If we are looking for binary functions, and we have generated codes, we test it. If it fails, we simply use it (its previous subfunctions) as a basic part. next

Hyar level Functions

While the above would probably "work", it is not really L such applied for an OZ problem.

In this last, we try O.T.'s (often techniques) ~~to try to do~~ in Least order.

I've probably worked on this Q a lot!

Now, it's a lot like L such for an INV problem! — we want a code for 2 soln.,

so we try them in  $\frac{pc}{cc}$  order. In present case, I want a short pem,

$O^j$  :  $\exists (pc_{O^j} \cdot O^j (Ant_1 | Ant_2))$  Max.  $\approx -\log_2 Z$  is essentially the length of the code of  $Ant_1$  &  $Ant_2$ .

Probably using true OZ L such would be much better — but it would depend critically on what O.T.'s one had available. one had available

So, until one has good O.T.'s ~~perhaps~~ OZ L such is not so hot!

i.e. The Least of a good O.T. for a given problem, would be ?? solving it as a INV. problem.

Later when TM has had experience encoding O.T. probs or has been given a good set of O.T.'s & it has "calibrated" them, so it knows when to use which one — then it's best to use OZ L such.

We will, in any mildly advanced TM, be concerned, in L such, w. not the original pc of solns (both INV & OZ) but with the conditional - on - the problem. pc of soln.

Until TM is somewhat trained in an area of INV or OZ probs, its cond. pc's will not be v.g., & L such using them will be Very Expensive. (Spec 417.00)

CPC: Conditional proby ~~computer calculator~~

Variations of OT's.

0.00 (416.40) : Eventually, TM will ~~be~~ have learned various methods to find short codes: Contingent on the corpus to be coded — so a conditional pc of the coding tech req use with the corpus in "condition". While compression coding is a variety of OZ problem, it is a special category of OZ problem, & will, in a "mature TM" have by pc's for special OT's to be used for compressn coding.

This Cond. pc for Lsch is a very import. idea: for a young TM, cond. pc & uncond. pc.

In general, we would always want to use cond. pc for level, except that it may be more costly — "cc w/c". Perhaps much more costly. In  $L_{cost} = \frac{cc}{pc}$ , cc = cost of computing pc + cost of generation and test (does this include cost of computing pc?)

0.09 1.10 Discussion of "Context" of 300.26 is really a discn. of Cond. pc.

In fact, any thing that anyone (or TM) ever does can be regarded as being primarily based on cond/pc ( $\equiv$  spc). — i.e. T. task to be done, t. background info, the info from "related domains".

A statement for t. report: That cond/pc's are important & all pervasive is very clear. However, using the freq. concept of pc, our ssz's for most cp's of interest are too small to be useful for probty estimates. By employing ALP, we are able to use data of very many different kinds, to calculate pc's, so ssz's are much more likely to be adequate. ALP makes it possible to use data from disparate domains of knowledge to obtain cp's in any particular area.

0.20 The idea that CPC can work on to be "broad" to work on practically any problem is one of 6 imp't ideas in QATM. It is sort of related to  $TM_2 = TM$ .

0.21 0.23 Another imp't idea: That if all info is in t. P.D., that all possl. hours can be supplemental — viz Lsch — t. result will be about as good as using those hours. T. idea that Lsch is a universal mechanism that can express/simplify any hour: if we allow ing during t. sch. That any hour can be expressed as problem of building 2nd class string using Pd during sch.

On 404.31, we run into relational diffys: up to this pt. we've been using rel part data periods using Pd during sch. I want to convey t. idea that we will be using a subset of the rel part data periods using Pd during sch.

on  $n$ .  $[O_n^k]$  is part of (404.19R) I could mention  $[O_n^k]$  then say that  $[O_n^k]$  is composed of 2 kinds of sub-functions:  $R^j$  and  $[P_j^{*i}]$  (or  $[P_j^i]$ ?)

For each value of  $n$ , there will be  $w_n$  kinds of subfunctions. Usually  $w_n$  can. There will be initially, a set of  $w_n$  recognition functions that recognize  $R^j$  that are able to recognize which kind of problem a particular  $Q_i$  is.  $R^j(Q_i) = 1$  if  $Q_i$  is the  $j^{th}$  type of problem — even if not, then  $R^j(Q_i) = 0$ .

Replace 404.31: One method of this sort expresses the set of  $O_i^j$  being used as as a combination of two types of subfunctions. → 406.10

Tasks for Lsch maintained by system as describ: — (Shut trailing updating;) So viz 407.03: 1), 2), 3), 4), 5) 2) Update  $P_i^k$  3) Find adjust  $R^{k+n}$  4) Modify set of recognizers 5) Find recognizer for certain set of possl. cases So: only 2 operations: 1) Update  $P_i^k$  so it's good pc for Q's of type k 2) Find recognizer for certain set of possl. cases " " " negative cases (spc 420.00)

.00: 4.15.40 The second type of function,  $P_i^k(\cdot, \cdot)$  ~~is a probability distribution~~ is a probability distribution.  $P_i^k(A|Q_2)$  gives the conditional probability (as seen by  $P_i^k$ ) of the answer A, to question  $Q_2$ .

~~The probability induced by the set of functions  $P_i^k$  (for fixed  $i$ ) on  $A$ , with input condition  $Q_2$  is~~

~~if~~  $R^i(Q_2) = 1$  then  $Q_2$  is a problem of the  $i$ th kind, and

.09 
$$\sum_{k=1}^{k_i} a_i^k P_i^k(A|Q_2) / \sum_{k=1}^{k_i} a_i^k$$

.10 ~~is~~ is the conditional probability induced by the set of functions  $P_i^k$  on the set of all possible answers, A, for question  $Q_2$

$a_i^k$  is the a priori probability of the function,  $P_i^k$

$\sum_{k=1}^{k_i} a_i^k$  is a normalization factor.

$k_i$  is the number of prediction functions for  $Q_2$  i'th kind of question.

More generally, if  $Q_{n+1}$  is a new problem/question never seen by the system before, its probability distribution on possible answers, A, will be

.18 
$$\sum_{i=1}^{W_n} \sum_{k=1}^{k_i} a_i^k R^i(Q_{n+1}) P_i^k(A|Q_2) / \sum_{i=1}^{W_n} \sum_{k=1}^{k_i} a_i^k R^i(Q_{n+1})$$

\*w's  
i.c.

.20 If  $R^i(Q_{n+1}) = 1$  for only one and only one value of  $i$ , then expressions .18 and .09 are ~~almost~~ the same, ~~except~~ except that  $l = n+1$

If  $R^i(Q_{n+1}) = 1$  for more than one value of  $i$ , then ~~expression~~. 18 gives a maneuver the various kinds of problems  $Q_{n+1}$  seems to be.

If  $R^i(Q_{n+1}) = 0$  for  $i=1 \dots W_n$ , then ~~expression~~. 18 ~~is~~ becomes an ~~undefined~~ <sup>zero divided by zero.</sup> The system

.07 **Section: Updating**

~~cannot answer~~ question  $Q_{n+1}$ .

A ~~major~~ critical aspect of system operator is that of "updating."

~~The single most important question about the system is: "How is updating performed?"~~

.30 ~~is~~ When we are given a new question-answer pair  $(Q_{n+1}, A_{n+1})$ , how is the system modified ~~is~~?

~~is~~ How do we update the sets of functions  $R^i$  and  $P_i^k$  so that the system will respond accurately to  $Q_{n+1}$ , ~~and also be likely to respond well to~~

~~future questions?~~ Continue to respond well for  $Q_1 \dots Q_n$  and be likely to respond

well to (unknown) future questions?



Each has its own updating algorithm

418.10: There are five possible kinds of response of the system to the new  $Q_{n+1}, A_{n+1}$

- There is one and only one  $i$  such that  $R^i(Q_{n+1}) = 1$  and  $P = \sum_{k=1}^K a_k^k P_k^k(A_{n+1} | Q_{n+1}) / \sum_{k=1}^K a_k^k$  is acceptably large. In this case the response of the system is satisfactory and no updating is needed.
- Same conditions as 1. except that  $P$  is too small.

In this case there are at least two ways to update:

- Modify  $R^i$  so the new  $R^i(Q_i) = R^i(Q_i)$  for  $i \leq n$ , but  $R^i(Q_{n+1}) = 0$ .  
 Devise a new  $R^{w_{n+1}}$  such that  $R^{w_{n+1}}(Q_i) = \phi$  for  $i \leq n$  but  $R^{w_{n+1}}(Q_{n+1}) = 1$ .  
 Find a set of  $P_{w_{n+1}}^k$  functions such that  $\sum_k a_{w_{n+1}}^k P_{w_{n+1}}^k(A_{n+1} | Q_{n+1})$  is as large as possible.
- Keep the same  $R^i(\cdot)$ , but modify the associated  $P_i^j$  functions, so that  $\sum_k a_i^k \prod_l P_i^k(A_l | Q_l)$  is as large as possible. (Here  $l$  ranges over those values of  $l$  for which  $R^i(Q_l) = 1$ .)

This method of updating is usually better in terms of future system performance, but takes longer.

- $R^i(Q_{n+1}) = \phi$  for  $i=1 \dots n$ .  
 Find a new  $R^{w_{n+1}}$  such that  $R^{w_{n+1}}(Q_{n+1}) = 1$  but  $R^{w_{n+1}}(Q_l) = 0$  for  $l \neq n+1$ .  
 Find a new set of functions,  $P_{w_{n+1}}^j$  such that  $\sum_k a_{w_{n+1}}^k P_{w_{n+1}}^k(A_{n+1} | Q_{n+1})$  is as large as possible.

4.  $R^i(Q_{n+1}) = 1$  for more than one value of  $i$ , and one of the recognizers is associated with a set of  $P_i^j$  functions that give an adequate  $P$  for  $A_{n+1}$ .  
 Modify the other recognizers so they no longer recognize  $Q_{n+1}$  but otherwise behave the same.

- $R^i(Q_{n+1}) = 1$  for more than one value of  $i$ , but none of the associated  $P_i^j$  functions gives an adequate  $P$  for  $A_{n+1}$ .  
 There are two ways to update in this case:  
 a. Modify the  $R^i$ 's that recognized  $Q_{n+1}$  so they no longer recognize  $Q_{n+1}$  but otherwise behave the same. We are now in "Update situation 3" and can proceed according.

b. Find a single set of  $P_i^j$  functions such that  $\sum_k a_i^k \prod_l P_i^k(A_l | Q_l)$  is as large as possible. Here  $l$  ranges over those integers for which  $R^i(Q_l) = 1$  and  $Q_l$  and  $Q_{n+1}$  have a common recognition function. i.e.  $\exists k \ni R^k(Q_l) = R^k(Q_{n+1}) = 1$

**Fn 1**: The criterion for whether  $P$  is "large enough" will change during course of training of the system. At first, the "largeness thresholds" will be given by the trainer as part of each problem. In later stages of training, the system itself will be able to induce reasonable thresholds - based on its own experience and its goal of maximizing the value of expressions 414.18 and 414.30.

Note for RJS! See 427.39 - 428.06 for a good discussion of this point - "A quantitative soln."  $\rightarrow$  429.00

00: 417.40: **Spec** Related to this is  $\phi$  of what is an adequate PC for a set of  $P_i^k$ 's? — this  $P_i^k$  is designed

01 [ If  $\phi$  set acts on a single  $Q$  we expect by PC — it can be rather a.H. — but if  $\phi$  for a larger  
02 [ bunch of  $Q$ 's some of  $\phi$ 's PC <sup>individual can</sup> be rather small (if  $\phi$  "good" d.f. has occasional  $A$ 's of small PC.)

Perhaps  $\phi$  PC will have to depend on  $\phi$  CC available.

Also, the  $R$  functions obtained may not be so good! If  $\phi$  one is bad, its assoc  $P_i^k$  set

07 Will ~~give~~ give poor  $PC$  for large  $cc$ .  $\rightarrow$  (423.34)

08 Another (advanced) trick in update situation **5** (407.13): None of  $\phi$  set that are recognized  
by  $\phi$  recognizer of  $Q$  are good PC for  $\phi$  Anti: One could try to find a new  $P_i^k$  function that

works well w. all  $Q$ 's recognized by  $\phi$  recognizer: this would seem like a miniature  
Grand QATM problem: It  $\phi$ , but 1) small no of cases 2)  $\phi$  cases are already very similar 3) one has  $\phi$

"side info,  $\phi$  fact that certain functs workd rather well on  $\phi$  certain subsets of these problems."  
08 is an alternative method of update for situation **5**.

At first Glance, it might seem that  $\phi$  problem of 08 would tend to not have an easy Soln.:  
we have already tried many <sup>subset</sup> functs  $\phi$  on parts of  $\phi$  problem: they don't work on  $\phi$  whole problem — so

that we need "larger" <sup>just</sup> functs for  $\phi$  whole problem subset. Well, the original  $P_i^k$  funct was  
designed for  $\phi$  first  $Q$  in  $\phi$  set only. It happens to fit  $\phi$  best in  $\phi$  subset. In looking for

a good  $\phi$   $P_i^k$  for  $\phi$  entire sub-set, the PC constraints are much reduced, since its only a good  
mem (geometric) of  $\phi$  PC's  $\phi$  needed (see .01-.02). The  $\phi$  initial recognition set  
was  $\phi$  based on the first  $Q$  in  $\phi$  set, and the  $[P_i^k]$  that was obtained for it.  $\rightarrow$  (423.10)

21 Actually,  $\phi$  this first  $Q$  and  $\phi$  completely defines the assoc  $R$   $\phi$  partition function  
(Also  $\phi$  one has to give a lower bound on the  $P_i^k$   $\phi$   $P_b(Q)$ ).  $\phi$  What  $\phi$  are looking  
for in a "proper recognition function" is a  $R$  that will do  $\phi$  recognition

$\phi$  cheaply (mcc), and will  $\phi$  extrapolate well to  $\phi$  unknown  $\phi$  Q's.  
Actually, all  $\phi$  info contained in  $\phi$  model is in  $\phi$  first QA  $\phi$  of each level  $\phi$  its  $[P_i^k]$ .

Occasionally a new  $P_i^k$  will work well for  $\phi$  previous  $Q$ 's. In this case, it is more "general"  
than the earlier  $P_i^k$ 's  $\phi$  we should try to modify  $\phi$  recogn. functs to implement this —

i.e.  $\phi$  so that  $\phi$  new  $P_i^k$   $\phi$  is applied to  $\phi$  the older  $Q$ 's.  $\phi$  w.r.t.  $\phi$  R functs:

027 We could just devise an  $R$  funct to recognize all of  $\phi$  old ones plus  $\phi$  new one,  
an keep all of  $\phi$  old recognizers.  $\rightarrow$  420324  $\rightarrow$  This would give multiple recognition for certain  $Q$ 's.

030 Which is a  $\phi$  modification of this QATM that I have considered. I suspect that eventually,  
I will end up doing that —  $\phi$  the  $\phi$  update process may be much more complex.

034 A further Genzn. would allow recogn functs to assign wts to  $Q$ 's — this would  
be a fuzzy (Recognition definition). Evaluating fuzzy recogn Algms involves applying them to

039 Vary many  $Q$ 's — which is expensive. One may develop ways to restrict them  
040 to a small set of  $Q$ 's in which they have wt.  $\phi$  some  $\phi$  say.

ID

100: 420.40: one might want eq. of 418.18 with "fuzzy" ~~R~~  $R^i$ 's, — subject, perhaps, to the modifn. of 420.39-40

02: Actually ~~4~~ Multiple Boolean recogn of 420.324-34 may make updating easier. — It may make it easier to find R-functs that categorize things in useful ways.

In general one might have an interactn betw.  $P_i^k$ 's &  $R_i$ 's: we try  $P_i^k$  — see which Q's it works/on — then try to find a short, fast  $R^i$  for it. We can do this many times & get a (ot of  $P_i^k$  & associated  $R_i$  functs. We can then use 418.18 (with Boolean R) to obtain probs., predns.

SN: font categorizn. problem of ~~multiple~~ recognition: ID3 would perhaps work. I suspect that it could be much improved. It does have sub-categorization, so diffrnt  $R^i$ 's could jointly use them: Also, it tries to get "final" decision in as few "micro decisions" as possbl. ID3 may not be ideal, but it may be adequate. It is a more or less "off the shelf" item. (16)

A Q is: would Lsrich be better? For "ID3" is other categorizers "Features" must be first design.

SN: For finding  $P_i^k$ , Lsrich can be used: or GA we heavily bias on short codes built into "fitness function". All old  $P_i^k$ 's will be in population to furnish "parts". We could start out in a relatively simple, nozotype GP, ~~with~~ (w.o. recursion limits, say), then add capacity for complexity, as we acquire good  $P_i^k$ 's to be used as "parent stock".

16: (11) ID3: May be not "off the shelf" at all! ~~It~~ <sup>write</sup> T. recogn problem is a binary classifn task, it is to recognize strings, it is rather general, ID3 has to have a vector of discrete properties for each element (Q) to be categorized. While we might eventually have such functions, (if they would be useful) they are not immediately available <sup>"ab initio"</sup> "in a library". ID3 holds a "feature extractor" at "front end"

19: For the recogn. problem, it would be useful to try to develop a set of categories (for a ID3-like time decision process) so that decisions would be fast. A possbl. approach: Start w. simple, "first order" obs on Q's. The output of these can then be combined in <sup>(Logical)</sup> Boolean ways to create higher order obs. A way this might be done: Say we have a broadly found set of  $R^i$ 's for the Q's in hand. ~~It~~ This set of  $R^i$ 's is slow, but is not particularly "stout" In deriving Boolean categories (perhaps by Lsrich) a possbl. Gove: to divide the set of Q's into 2 approximately  $\approx$  parts. This could be good "first order obs" (1.21)

I think there may have been lots of work done on a "string categorizer" problem. Perhaps extensions of ID3. Its likely that Vagne's system has been used for it. — My impression is that it is slow, but this may be a false & fixable.

For small n, the  $R^i$ 's will have to be revised after each  $n \rightarrow n+1$ ; but after we have a large set of  $R^i$ 's for a  $R^i$ , revision will occur less frequently (if not at all!) — This would be true independently of how we find the  $R^i$ 's — (No it would converge faster if we do the "ALP-oriented" Lsrich)

30: that ~~the~~ recent book on "Pattern discover by Dudder, Hord & Stone" will probably have a lot of info on categorizn. problem for strings & perhaps on  $P_i^k$  discovery. Will Garsh may have some ideas ~~for~~ references.

After (or even better) we have a rather stable set of  $R^i$ 's, we should begin to optimize the ANM set of R, P functions in various ways. We can optimize Rev  $R^i$ 's (for fixed P set), w.r.t speed & pc. For non fixed R, we can try to find P functions that solve a larger no of problems at lower cost, lower cc. Then existing set of P's. The  $R^i$ 's would then have to be revised. SAC 422.00

Spec

00:421.00 : In fact, the whole idea of QATM is to have 1 P<sup>k</sup> function per entire corpus!  
So one approach would be the "problem pool" approach. We try to solve as many problems as possibl. w. as few functions as possibl. This can be done in an efficient way by merging pairs of P's, then repeating the associated R's. This merging of pairs of P's can continue until we have one big P for entire corpus.

Of course we already have that, using many R's & P's, but the pair by pair merging produces a cost & a/c at each merging.

Deciding which pair of P's to merge is an interesting problem. After choosing the pair, one has to choose P's & their subfunctions - which should heavily bias the search for "merger P".

2 associated R functions could also give useful info. We know that for 2 R's & P's can be combined to yield a single R & P, but at both costs. We want to reduce.

423.10 is a possible way to find good MERGE candidates. ABCDEFGHABDEFGHIJ

An unpleasant note on "partitioning": Given a set of n elements, the info needed to partition it into 2 specific sets of 2 and n-2 elements is maximum

when n=2a, My "soln" for problem - well I had 2 solutions - unclear as to which was better. The exact nature of a priori info involved would have to be stated to resolve the problem.

I think I soln was  $\frac{n!}{(n-2)!2!}$  which is the binomial distribution & Gauss d.f. for fusion. I say "unpleasant" because it's not that partitioning (P) dividing the corpus into 2 = parts would be best. Well maybe it is best if its a small corpus, because

it was a "lot of work" w. small  $rc \approx r/w \text{ cost} \approx \frac{1}{pc}$

On using GP/GA to find  $O_n$ : One apparent Big Objection was that it would take too long to test/verify each cand. Hvr, maybe not so. By using statistical testing on the QA's, one can make a good estimate of the rank of a Cand, w.o. testing many QA's - in fact one can choose "diff't subset" of QA's to use for preliminary testing. If a cand passes prelim tests, we test using more (maybe all) QA's. Since the population of parents are mainly vgr over the whole QA set, it may be that parent a cand does well on, say, 20 of the QA's, it will tend to do well on the rest of them! If true, this could reduce considerably

the cc wacy to find good cands (good  $O_n$ 's). → See 432.20 for discussion of importance of .20 TP. A how to fully use the GP/GA soln.

A very brief summary: 1) The use of (20) - GA for finding  $O_n$  is fairly novel, & if it works, could give us a v.g. TM. That is non-el, it is inherently slower than much more el. methods & mainly discussed.

2) The el. method described essentially reduces to finding a P<sup>k</sup> for a new problem & finding and revising R<sup>n</sup> (recognition) functs. The system obtained may be good and (w. a suitable tsq), to work on problems in progressively less el. ways.

3) Even w. the more simple "problem type" version of QATM of .36, there are 423.00

.00: 422.40 impt. improvements that can be made, e.g. <sup>Merging</sup> ~~to~~ (coaliscmy) of  $P_i^k$  func's. (422.00-1.2).

Developing Categorization is subcats for R-functions (421.19 ff) could reduce cc of solns is also  
↓ rc of soln. 420.08-2.2 Gives a less E. soln. to "Update situation" of 407.13. It reduces  
no. of  $P_i^k$ 's &  $R_i^j$ 's. so ↓ of cc is ↓ of rc.

420.20-32 is a beginning of a dem of just what this "Bare bones" model is, & what its basic backness is.

420.27 ff is particularly impt.

10 J420.21

If  $R^i(Q_{11})=1$  &  $R^j(Q_{11})=1$ ! This suggests that  $P_5$  &  $P_{10}$  may be similar. Try to find a common  
P for  $P_5, P_6, P_{11}$  <sup>Merging</sup> ~~from~~  $P_{10}$  into the set P! [

.10 is one poss. heur for merging of  $P_i$ 's. I will need to find other heurs for this, since

MERGING is a very impt Operator.

The Merging itself is perhaps standard Lschr, but more generally it is a problem w. 19 t.  
Case of .10, "side info" of  $P_5, R_{10}, P_5, P_6$  → (see 420.08-2.1) for more discuss.

"A kind of Heur" Suggested by .10: suppose we have this set of  $R_i^j$ 's that do not  
"overlap" (e.g. for  $n$  and  $m$  do not overlap we  $R^i(Q_n)=R^j(Q_m)=1$ , then we can get  
some overlap (so .10 can be modified) by "weaving" one or more more of the  $R^i$ .

"weaving" means removing constraints - in general "shortening dem" (see 424/P3/A/V/A/B/8/8).

20

SN I had idea of 417.2 that Lschr "long" could simulate any heur. Maybe Not so  
It would not simulate "Look Ahead": T. doing of "experiments" to gain info (at some cc)

Perhaps use term "Expt." rather than "look ahead" - until I think of a better term: A non-greedy  
# "Far Sighted" rather than (Myopic)

"Look Ahead"  
May be wrong term.  
Its used in  
Penny  
Chips, in which  
Several persons  
simultly taken -  
then later  
all but 1 part is  
abandoned.

.24

Insert at 15.10 : (See 422.20-31 for some <sup>encouraging</sup> theoretical discn of GA for this problem.) (95-100)

.25

Genetic Programming, using a Lisp-like functional language, would seem  
to be a reasonable way to solve this problem. A large population of  $P_i^j$  functions  
could be mutated and/or recombined w. a suitable crossover algorithms and a suitable  
fitness function. ~~At first glance~~

30

~~At first glance~~ Preliminary analysis suggests that while this may give very  
accurate probability distributions, it would be very slow, as compared to the system  
we will subsequently describe. This is, however, only a first impression  
and we have by no means abandoned Genetic Programming as a possible technique.

.32

.34 420.07 on "Adequacy of pc's assigned by System" ~~however~~, This has to do (intimately) w. t. Update  
algs, since it is t. criterion for discard/acceptance of  $R^i, P_i^j$  functions  
As was noted in 420.01-02: Sometimes its O.K. for a pc for a particular Q to be small.

A poss. rule select a  $P_j^k$  set so that  $\sum_k 2^k \prod P_j^k(Q_A/Q_B) \leq 2^k$   
is as large as poss. for t amt of search time used.

Probably longer  
this Norman factor! → 724.00

"Random"

423,40

Usually in induction problems, there is a "default code" equiv. to a term "By chance" or some trivial regularity. Getting a code at least as good as one does not, in available cc, get a better than that, then on concesses failure uses default code.

So, in updating problem, the main idea is to get max pc in available cc.

There are 2 impt. subgoals of a "Bare-Bones" system!

- 1) Being able to use fully work on "Updating problem" as a "regular problem"
- 2) Being able to ~~not~~ read & "understand" ~~the~~ Nat. lang. to a useful extent.

50°  
44°  
64

Q: Is our goal to get  $\sum_i a_i \prod_{j=1}^i P^j(A_j|Q_j)$  or  $\sum_i a_i \prod_{j=1}^i p^j(A_j|Q_j) / \sum_i a_i$

or  $\sum_i a_i \prod_{j=1}^i P^j(A_j|Q_j)$  v.s.  $\sum_i a_i \prod_{j=1}^i p^j(A_j|Q_j) / \sum_i a_i \prod_{j=1}^i P^j(A_j|Q_j)$

← This is cond pc of  $A_{i+1}$  in view of  $Q_{i+1}$

Go Back to (S78T3): say  $a_i$  is approx  $p^i$ .

Perhaps exist: 414.18 or 414.31. If we maximize  $\sum_j a_j$

sums over all codes. That would seem about best. It would get most wt. for the Corrasco. Prod. alg. In the update case, we want  $\max \sum_j a_j$ .

Here say we are extrapolating string X: for best induction we want total of all pems and of strings that have max prefix. T. larger a total wt. of our pems, the more likely it is that they will have much wt. for the "Rite" pem (??). Ideally, would like to use the shortest pem, which is  $\max a_j$ ; even better, would be to include the best  $a_j$  pems  $\Rightarrow$  max total wt. was  $\sim 90\%$  of all wt.

While it's true that the more wt. one has in  $\sum_{j=1}^{400} a_j$  the closer it is to  $\sum_{j=1}^{\infty} a_j$  &  $\therefore$  more accurate; This is not actually the problem! We are interested in good probys i.e. will weighty sets of prod. functs give better prodns than less weighty sets?

Consider predictors of  $F(x)$  is a function of prefixes / strings, X say  $F^j$  is the  $j^{\text{th}}$  such predictor in the set.  $F^j(x)$  gives a pd over 0,1 for i. bits following X, for every binary X, consistent.  $\sum_j a_j$  wtd sum of peds for the  $n^{\text{th}}$  bit of X ( $X_n = \text{true}$  & a bit of X).

The pd. is  $\sum_j a_j F^j(x_n)$ . or let  $F_0^j(x_n) = \text{prob of } 0 \text{ being next bit of } X_n$

so  $\sum_j a_j F_0^j(x_n) = \text{pc of } 0 \text{ by sum of}$

I think the end result is that if we use .30 for prod. of sum, No! we want to use  $\sum_{j=1}^n a_j F^j(x_1, \dots, x_n)$  as wt. of  $F^j$  after  $X_{i+1}$ . This is  $F^j$ 's pcd &  $i^{\text{th}}$  bit inverse to past  $X_i$ .

If we do this, then the macro predictor ends up in a pc for the  $X_n$  corpus of its total wt. at output. So we can use Corollary of S78T3

compare it to true p.d. : So larger is better, since  $\text{expected error} \leftarrow \ln \frac{\text{true prod. of corpus}}{\text{approximate prod. of corpus}}$

I think the Algebra is very simple!

00: 424.40! (Spec 399.40) Re: The Question about  $\ln \frac{P'}{P}$  always being < 0: (worried about 11; 399.33-40)  
 I think it has to do w. Expected Values & Null-Likelihood distances which always  $\geq 0$ . Hvr, I'm still uneasy about this! Work it out in some detail.

→ Gac's theorem about  $\Sigma$  error variance for semi-measures: is P.3 implied by S7B3's corollary?  
No! P.3 corollary is for ratios of Measures only.

- \* It may be better to consider only 3 Update situations: (variations)
- 1)  $R^2(Q_{int}) = \phi$  for all  $i$
  - 2)  $R^2(Q_{int}) = 1$  for 1 or more  $i$
  - a)  $T_i$  best  $i$  is adequate.
  - b) " " " not adequate.

In 4.5 situation approach 2) (above) was split into 4 variation  $\geq 2$  cases: i.e.  $R^2(Q_{int}) = 1$  for  $1 \leq i \leq 2$  or  $> 2$

Doing things cases makes dim. of "what to do if one has lots of cc for update" easier to curb.

**SN** Is + "No free lunch than / ugly dechling them" true  $\iff$  we assume (uniform)  $H$ ? (I think so)  
 If so, this is clearly ridiculous: Only to H. diff. gives "no induction" any other d.f. gives some induction.

**Exposition:** The expected total error is a monotonic  $\downarrow$  funct of  $P$ . — So Maximizing  $P$  is equiv to minimizing Expected total error.  
 This is the most important idea in the paper: There are several ways to find large values of fine functions with large values of  $P$ . All induction techniques can be compared and evaluated on a basis of how well they do this.

This will have to be preceded by an explanation of S7B3 & its result, then how 418, 18, 31 give a result.

The rest of the paper will discuss one particular method of approximation. The reader with ~~deductibly~~  $\iff$  ~~able to derive techniques~~ may be able to devise techniques for getting a larger value for  $P$ . The reader with ~~deductibly~~  $\iff$  ~~able to derive techniques~~ may be able to devise techniques for getting a larger value for  $P$ .  
 greater efficiency, efficiency, efficient, efficient

**SN** Perhaps its possl. to show that  $\downarrow$  convergence of ~~Serrz~~ ~~directly~~ implies convergence of other error measures: This would probably give a lower upper bound on convergence rate than Marcus Got.

419.00 ff: Criterion for update is unclear: When do I want to use  $\frac{1}{\sum z_i^k}$  normed constant?

I do want good products — (which measure a Normed Const). On the other hand, the ~~Serrz~~ is  $\downarrow$  fund of

$\geq Wt$ , which does not use a normed const!  
 Consider 4 Deterministic Prdn (DPM). What indication do I have of "rite answer"? — Well,  $R_{err}$  is

f. default code. — I want to do much better than that! (Default code would be perhaps based on frequencies of symbols or of various common abbs, codes.)

One approach: Just try to get as much  $Wt$  as possl. w. 1 available cc. — That's it!  
 If "no updating is needed" (because t. prdn was v.g.), spend cc on other aspects of  $Wt$ .  
 look for more "global" regularities.

→ 34 seems to deal w. t. problem of occasional by most prdcs — that most occasionally occurs.  
 Well, 34 is a global soln, but I have broken down t. problem into  $R_i$  &  $P_i$  functions → 426.00

6.27386 s/s H5ct  
→ 440.552 n  
~ 456.537 x parallel  
2.778/yr to dev

00:45:40, so f. Q: how much cc to spend on R' modification? on R' modif. of  
Oscar R: P's? On 'unifications' of R's & of P's?  
Some how a good prediction (say  $pc = .9$ ) would signal that one need not spend more  
time on it. R is for that QA. ... Not really a correct conclusion, Hrrr!  
So f. Goodness of (a) predictions can be regarded as a hour, that suggests which

R's & P's to work on - it gives f. prefer type. initially

Now, I may be that TSP (is/has to be) designed, so that pradn is good  
hour for focusing on cc, is talking when to stop. In fact, initial STO's are designed  
so that programs are "solvable", so one knows when one has solved them, & one wants  
to go on to a new problem. It is usually in the most part very advanced part  
of training that TH is given problems of unknown solns - Where f. trainer doesn't  
know if a soln. is "correct"/adequate.

Perhaps "am starting" of TSP, & that that could actually tell TH what PC levels to target  
for some to problems - as part of problem defn. Butive also want TH to (in  
At first, most of probs will be ATM, so PC crash  $\sim .99$ . Butive also want TH to (in

goodly things, so it can find Hvers.  
An advanced TM should be able to guess how much compression to expect w. a certain  
cc on a certain problem. At TH works on a problem, its estimates will change.

One imp kind of "problems": Designing good concs based on "context" (a condi PC path)  
Here, it is not clear as to when to stop - when progress has been made - too a PC >> default

PC is in progress.  
I want to maximize TM. Goals is TM max(wt/cc): because any moment of kind  
ed wt. would also be maximizable. -  $w_1^2, w_2^2, w_3^2, w_4^2, w_5^2, w_6^2, w_7^2, w_8^2, w_9^2, w_{10}^2, w_{11}^2, w_{12}^2, w_{13}^2, w_{14}^2, w_{15}^2$ .  
- Thus, why did I need a linear ordering of (cc, cc) pairs? - say I just want max cc  
pc for any particular cc. Well, \$100 w. pc =  $10^{-10}$  v.g.  $10^k$  for  $pc = 10^{-2}$ ; to choo  $cc = 10^{-2}$   
depends on problem(s) being solved.

For fixed cc or fixed pc, the answer is always clear: Hrrr, paths in cc, pc space may  
not be possible to choose between. A kind of Modified "core center" in which some  
"output" occurs both final results

So: A Q: is whether it is correct/true. Do we have to tell TM level of PC required for prob. soln.  
Perhaps it is needed only because we have broken down the "main problem" (= Max wt. of errors)

Another possy write be: for each problem a certain "Min Compression" w. Default coding.  
in to "Subproblems" (sets),  
(same for all probs)

Hrrr - what about TM working on itself - what PC thresholds to use?  
427.00



00: (426.40) Perhaps TM could learn how much to allocate time as a "Meta level" problem!  
 That may be it! After TM has had many examples of pc thresholds for solns  
 given by USR, it will be able to induce rules for even problems,  
 and it has a "GRAND GOAL" for doing this - i.e. max total computer  
for a given CC.

So, main outline of TM so far: 1) 414-18, 30 Max. wt. Goals

- 2) Break down  $O_n$  into  $R, P$  functions.
- 3) Solve  $R, P$  problems, guided by pc thresholds / USR, by USR.  
 - use L srcn a/o GP or whatever to solve  $R, P$  problems.  
 (using some func.)
- 4) Continuously be trying to merge  $R$  &  $P$  functions.  $\rightarrow P \ll R \ll \text{Main Quest.}$

So that's it for "first level epc"

Second level!  
 Next level: Get above machine to work on folg. items:

1. best problem itself, previous solns to "a problem"; (similarity judged by common R func's & "features")

- 1) In Lsch we are guided by a pd. a) putting in this pd. b) make pd. conditional on various things - so computing it becomes a QATM problem.
- 2) Solve "Look ahead problems": so it can learn to use "Look ahead" in its own prob. solving.  
 = do "experiments" = actions not immediately directed toward optm of Govc.
- 3) Devise its own criteria for a) indexing b) how much to spend on a problem/sub-problem.
- 4) Work on problem of extrapolating goals of Sci Community.
- 5) It should learn to devise and solve problems that are oriented to ward to touch near goals of it, or t. goals of its.
- 6) Reading Net lang & one's understanding.
- 7) Do more "Merging" of R & P functions.  $n \rightarrow$  Abcdelghiklmnop  $\rightarrow$  is a

These "cc thresholds" can be mult by 1.5, say, if TM has more time available: Essentially, it means

k. types of pc's spent on various sub-tasks  $\rightarrow$  we can regard them as being done in 1).

For "R" (recog. problems): Rise one strictly IMV problems. One could accept t. first soln found - or continue to search for solns. of higher pc. - Essentially, one would not consider solns. of lower pc than t. one found. - so a limit set of combs - is one spends more is more time on exhaustion t. combs (that he hasn't yet converged).

On t. other ways, t. P problems always have "soft" solns, is a "soln" is any better than "Default". After a better than default soln has been obtained, TM could spend equal time on improving R & improving P solns. T. goal is to minimize products of roots of t. 2. So t. problem is - given a certain amt. of time, what is best division of it betw R & P improvement tasks? One could test solve R (w. min cc soln), then spend rest of time on P. - But after R is solved, one should, at first, spend much time on P, but eventually, we get to pt. of diminishing returns. I think one would work on R or P depending on whether  $\frac{d \ln P}{d CC}$  or  $\frac{d \ln R}{d CC}$  is larger. Second level! could estimate these derivatives

Very Important

139  
 40  $\rightarrow$  The 2 derivatives are probly be R ↓ funcs of CC.  $\rightarrow$  428.00

Contextual

00: (space) 427.4: How Level 2 TM (estimates) <sup>(calculates)</sup> <sup>higher derivatives</sup> is unclear. I guess a "usual way": Various contextual clues, grouped over many problems in the past. We keep records of that activity. Also, I'd like TM to "Think" about the problem - use reasoning to get better solns. - but this is an aspect of TM behavior that I haven't worked on much. My impression is that there are many "Math/Logic" problems that are close to this, that could be given to TMs as "regular problems" in training to "think about" its "real" regular problems.

06

Trying to get Eff. for 419.30-32 correct. "5b" update  
 Also Find a new R that recognizes all of the  $Q_L$  that were recognized by a  $R^k$  that recognized  $Q_{HI}$ .  
 i.e. if  $\exists k \ni R^k(Q_L) = R^k(Q_{HI}) = 1$ , then  $R(Q_L) = 1$ . Otherwise  $R(Q_L) = 0$ .  
 $R(Q_L) = 1$  if and only if  $\exists k \ni R^k(Q_L) = R^k(Q_{HI}) = 1$ .  
 Discard all older  $R^k$ 's  $\exists k \ni Q_L \ni R^k(Q_L) = R^k(Q_{HI}) = 1$

10

Describe 5b merging in English w. much care. If I can find a good way to do this mathematically, symbolically, do it ~~more~~ following English deriv.

Discuss how deriv length of  $O_n$  is "sum" of deriv lengths of  $P_i^j$  (but this must be made exact: it can be very tricky! - Its easiest to explain if in  $P_i^j$ , j has only 1 value: i.e. for each  $R_i^z$ , only one  $P_i^j$ . Perhaps say  $R(O_n) \geq \sum R(P_i^{z_i})$  where  $P_i^{z_i}$  is the shortest prefix of shortest deriv of all functions having same z value.

20

**SN** idea of having many || codes  $P_i^j$  (many j values) was partly to a total PC, but more to enable updating. The || codes are able to store many useful concs. They help form a "Population Pool" for GA. However, it would seem that for a given i value, the concs in  $P_i^j$  of different j cannot be "shared". On the other hand, with  $P_{z_1}^j$  &  $P_{z_2}^k$ , i.e.  $z_1, z_2$  are different, then  $P_{z_1}^j$  can share concs for all values of j in  $P_{z_1}^j$ . The "weight" of this sharing  $\uparrow$  if  $P_{z_1}^j$  ~~is~~

We want to  
 $\downarrow$   $W_n$ , + no. of Distinct R types; + no. of distinct  $P_i^j$  types.  
 of all functions having same z value.  
 $R(P_i^{z_i}) = \sum R(P_i^{z_i})$   
 Try to find correct notation - maybe look at  $\mathbb{N}^i - V. T. Y. B. B. ?$

assigns by PC to ~~some~~ to correct answers for the problems it works on. On the other hand, one would want to share concs betw. different j value of same i value in  $P_i^j$ , because "same i" means they work on the same problems, & are more likely to be ~.

30

So 30-30 is a PARADOX!

Perhaps write  $2^n \approx \prod (R_i^{z_i})$  <sup>for whatever equiv of  $O_n$  is</sup> ~~is~~  $\prod R_i^{z_i}$ ?? ~~is~~  $O_n$  really doesn't have an equiv!!

The formula for  $\downarrow$  w. of the entire deriv involves the  $z$  equivs of  $M$ .  $R_i^{z_i}$  is  $P_i^j$  but I don't think  $z$  is  $P_i^j$ .  $O^n$  is a set of functions. An equiv can be assigned to each component function. Perhaps give example of  $z \in [O_n^k]$  set as a cartesian product of  $W_n$  sets of functions.  $O_n$  414.18.30 We used to set of functions  $[O^j]$  for each n, we will have a (often distinct) set of functions,  $[O^j]$ .

One way to get the kind of "interaction", "synergy" betw. ||  $O^j$ 's is by "condl PC":

00: 419.40

~~Unnecessary~~ ~~the time~~ while this method of updating is usually very time-consuming, it is ~~more~~ a very desirable kind of update since it "merges"  $P_i$  functions and it also merges the associated  $R$  functions. These operations can ~~lead to a~~ very large increase to expressions 414.18 and 414.30. They can also speed up ~~the~~ considerably both ~~the~~ recognition and prediction.

:05

Insert 430.02 - .05

10

20

30

ABCde ABCdefg

ABCdefghij



$$\int_{-\infty}^{+\infty} e^{-\frac{x^2}{2\sigma^2}} dx = \sqrt{2\pi}$$

$\propto \beta \Delta S$  p. 5

00: (spac 429.90): This might be done by using all 11 0's as "context"  $\epsilon \equiv \text{cond}(p_i)$  for new trial 0's  
This looks like GA(BP). - But I think it's just as relevant to PC's for LSrch

02 (for 429.05) This large increase ~~in~~ resides mainly in the a priori probabilities of the  $O_i$  functions. Each such function consists of  $W_i$  sets of  $R^i$  and associated  $P_i$  functions. The description length of  $O_i$  is approximately the sum of the description lengths of the  $R^i$  and the  $P_i$ .  
By merging the  $R^i$  functions, ~~we reduce the~~ we reduce the ~~length~~ of description of  $O_i$  which increases its a priori probability.  
05 Number of functions that need to be described, - which increases the a priori probability of  $O_i$ .  
06

b. The fact that several of the  $R^i$  accept a common  $Q_{int}$  suggests that ~~the~~ the  $R^i$  and their associated  $Q$  ~~are~~  $Q$ 's accepted by these  $R^i$  have something in common. They become a good merging candidate. ~~we should like to find out whether~~  
possible to find a single set of  $P_i$  functions that can work all of ~~the~~  $R^i$  functions so that we look for a single set of  $P_i$  functions so that

$$\sum_k z^k \prod_i P_i^k(A_i | Q_k)$$

15 ~~is as large as possible~~ is as large ~~as possible~~ as possible. Here ~~we range over those~~  $Q_k$  ranges over those "accepted"  $Q_k$ 's. i.e.  $\exists k \ni R^k(Q_k) = R^k(Q_{int}) = 1$   
17 If we are able to find such functions, we can eliminate the  $R^i$  functions ~~that~~ that accepted  $Q_{int}$  and find a new  $R$  function associated with the  $P_i$ 's of expression ~~430.15~~ 430.15.

20 ~~presumably~~ ~~larger than~~ larger than the corresponding weight, using the "unmerged"  $P_i^k$  functions.  
21

Actually, . . . 15-16 is ~~trivial~~ trivial (of no interest, unless we include .21).  
Perhaps we should not have 5b at all. . . i just have this "Merging", as an imp. part of it.  
Update process; ~~if~~ if we have time, After update 5a

27 b. If there are two or more  $R^k$  such that  $R^k(Q_{int}) = 1$ , Then ~~is~~ the  $P_i^k$ 's associated with these  $R^k$ 's ~~have~~ have ~~some~~ ~~important~~ ~~features~~ ~~in~~ ~~common~~. We ~~can~~ ~~take~~ ~~advantage~~ ~~of~~ ~~this~~ ~~"commonness"~~. ~~If~~ ~~we~~ ~~can~~ ~~find~~ ~~a~~ ~~single~~  ~~$P_i^k$~~  ~~function~~ ~~that~~ ~~is~~ ~~able~~ ~~to~~ ~~solve~~ ~~all~~ ~~of~~ ~~the~~ ~~problems~~ ~~formerly~~ ~~requiring~~ ~~many~~  ~~$P_i^k$~~  ~~sets~~. This ~~is~~ ~~would~~ reduce the description length of  $O_i$ , ~~and~~ ~~increases~~ ~~its~~ ~~a~~ ~~priori~~ ~~probability~~, ~~and~~ ~~Paraly~~ ~~by~~ ~~increasing~~ ~~the~~ ~~value~~ ~~of~~ ~~expression~~ ~~414.18~~.  
28

perhaps instead of 10-20

IS

.00: 430.40: Summary and Discussion

The system described is based on ~~the maximization of~~ <sup>Because it is</sup> ~~the maximization of~~ based on the maximization of expression 4(4.12), we are certain <sup>.035</sup> that if the regularity in the data is ~~expressible~~ <sup>describable</sup> as a finite string (plus perhaps a finite number of continuous parameters), then if we use L search for maximization, we are ~~sure~~ guaranteed to find the solution in a finite time. For ~~any~~ <sup>any</sup> induction problems of interest, this ~~finite~~ <sup>finite</sup> "finite time" is too long, so we ~~must~~ <sup>must use</sup> approximation methods.

The approximation methods suggested are based on human problem solving. We first identify the nature of a problem, then we apply solution methods that we have found ~~to~~ <sup>to be successful</sup> for that kind of problem. This corresponds to our finding ~~the~~ <sup>the</sup> Recognition Functions  $R^i$  and the conditional probability distributions  $P_i^k$ .

The "merging" of solution techniques that we employ corresponds to an important part of "Scientific Method": The ~~known~~ <sup>known</sup> summarization of many "special cases" theories by a grand, more encompassing Theory: Galileo's laws and Kepler's laws were summarized and extrapolated by Newton's laws.

The employment of carefully designed ~~training~~ <sup>training</sup> sequences of problems to train the system ~~appears to correspond to the norm in human learning.~~ ~~It is~~ ~~quite~~ ~~different~~ from most human training experiences, in which learning occurs  $\rightarrow$  437.26

Insert at .035

We are certain of two things:

First; that ~~the~~ <sup>the</sup> problems of widely different ~~domains~~ <sup>domains</sup> are able to share their concepts with each other. The techniques of problem solving ~~used~~ <sup>used</sup> and the concepts used in physics can be readily transferred to Mathematics, Sociology, etc.

Second: <sup>.035</sup> If the regularity in the data ...

This is not made clear!

.29: ~~437.31~~ 437.31 It is of interest to note that ~~the~~ the training sequence will ~~not~~ teach not only the system, but the trainer! ~~By~~ Observing how humans seem to solve ~~the~~ problems given to the system, will suggest improvements/modifications of the ~~and new present~~ <sup>described</sup> learning techniques and ~~which is hoped~~ <sup>probably</sup> entirely new ones.

One serious deficiency of most present learning systems is their inability to transfer information obtained in ~~any~~ any problem to the solution of subsequent problems. There has been some work done in this direction. "Incremental Learning" <sup>is</sup> ~~is~~ ~~an~~ ~~idea~~ ~~that~~ ~~perhaps~~ ~~perhaps~~ ~~the~~ category that best fits the system we have described. Another, more recent direction of research has been the use of "Wrappers" for machine learning — the continual adjustment of machine parameters

432.00

omit for present

00: 421.40 that transfers certain kinds of learning <sup>between</sup> ~~from~~ problem to successive <sup>near</sup> problems.  
We feel that the system we have describe may do this in an optimum manner.

.01  
.03 Discuss Ability of Systems to Do "Meta Learning" — corresponding to Human...? put this & 2 after 431.17

[SN] How Tolerant is ALP of "errors in Coding"? I mean serious conceptual errors —  
Say, the use of base 2 for one part of an expression & base 4 for another: Or the omitting  
of the "2141" feature of massive pc or no. of lines of code that are used — ? Well, it depends on the error type.  
Some are publicly minor — others catastrophic. Omitting the "2141" feature could destroy  
any induction!

Some things to include in Report: (see 427.12)

- 1) Why "Merging" is imp't: a) ↑ of pc 2) ↓ of cc so there can be a linear.
- 2) Why have many codes in 11? why not MDL? <sup>see</sup> 428.20.
- 3) Why 2 codes in 427.90 ff on how much time to spend on R v.s. P discoveries.
- 4) 425.16: How's discn of why max cu. of 424.18 is a Big Deal.
- 5) Explain how the QA form of problem presentation is a Univ. (426.04)

.16

Nuclear  
Theater  
31 Oct 2001  
SUN  
Juces

An Alternative General Scheme for QATM:

To solve 414.13: I opted to elementalize  $O_n$  — Break it down into sub-functs of certain kinds.  
(We have 2 aspects of Sci Med. Breaking problem into parts: Analysis:  $4 R^E$  functs.  
Integrating several apparently difrent parts of a Science (or difrent sciences) Synthesis: "Merging".

The <sup>other</sup> option is to solve for  $O_i$ 's directly for each value of  $n$ ; (This could be very  
cc intensive) — via LSrch or GP. Meanwhile, we are trying to get TM to work

on the  $[O^j]_{n+1} = [O^j]_n$  operation on  $Q_{n+1}$  and  $A_{n+1}$  & other (Contextual info). Sol. system  
is at all times, maximally "Non-el"  
so we started w. Big fast  
Computer system needed. This is  
true in all TM basing to be able  
to solve problem faster  
via "Cond. probly" (ie QATM problem)  
This is usable for both Lsrch & GP approaches.

30

Use of GP to solve the initial ~~prob~~ tscr probs might be nice, because we have a good  
population of  $O^j$  cond, & this population is slowly changed as we move from T to Q.

A serious problem is ~~the~~ cc of fasting cond but this may not be over <sup>powering</sup>  
powering  
— see previous note on chip → 422.20-31

Mainly, .30 is the problem of getting a good pd. for  $O_{n+1}$ , as a function of the  
b.t. "problem"  $Q_{n+1}$  and  $A_{n+1}$ . The Solns. for option ( $[O_n]$ ), any other "Contextual info".  
Before getting to this point, the cond for  $O_{n+1}$  are based only on D.F. for  $O_n$  & frequencies

of use of various sub-functs within the  $O_n$  & d.f.  
A possible dirty w. is off: That TM has to do a lot of long before its Meta Learning is good and to speed things up — it may take below  
to get to the pt. at which Meta Learning "runs true show"

dynamics associates

00: 432.401 O.k.: So what has to be done on report? We have:

1. Abstract

2. "Introduction": Actually, up to 414.13 is a kind of introduction. Its really not

an adequate intrdn. I will add more stuff that discusses to various sections.

3. <sup>New section.</sup> The paper starts "for the data described (above)  $\rightarrow$  (above in the introduction)"

Derbs Gant. soln using ALP

415.10 - on ZGA soln

415.11 Lsrch

415.20 The  $R_i^j$ ;  $P_i^k$  alzn. is how predn. is done

418.27 Updating thru 419.00 to ~~thru~~ 429.00-04 430.02-06 ends system.

I have "Summary and discussion" 431.00-40; 432.00-01 maybe do 432.03

$\alpha$  I have to write on TSG, perhaps at some length. <sup>things to include in report</sup>

$\beta$  " " " 432.10 : Discn. of first & second level machines: or just Meta Problems.  
427.12 (things of importance to (perhaps) discuss.)

I might actually send "first draft" of "Report" w.o. much on TSG - just refer to S86, S89: <sup>2nd level</sup>

Possibly quote a page or 2 - since both are in latex. Give markets! Put ~~S86~~ S86 on web.  
I will also have to tell how Meta-Problems are designed. Perhaps give example(s).  
It would be well to show how 2 arky problem types can be presented to QATM.

Mention use of index fields in Q. - That learning to understand the data of a problem is a good

Way to prepare TM to solve a problem.

8) I need section on Lsrch! - but rather simple, because only ITV (INV?)

Probs are being solved. Perhaps discuss  $T \leftarrow UT$ ,  $ll$  lsrch, Master Carlo  $ll$  lsrch.

4) ~~when in near the end of introduct.~~  
415.03  $\rightarrow$  08 - This could be made much stronger!  $\approx$  431.03-06  
The "G~~uaranteed~~" that we will get about the best poss. pd for  
that particular (finite)  $\leftarrow$  work of  $\{Q_i A_i\}$ .

5) 2 IMPT operations:  $\textcircled{a}$  Merging of old  $R^i$ ,  $P_i^k$  to form fewer fruct.  
 $\textcircled{b}$  Breaking down  $R^i$ 's into new subcomponents.  
In both cases we  $\uparrow$  pd of ~~subproblems~~ covers down.

i) 432.10-06; 427.12

Perhaps describe TM working on context as a TM calling on itself recursively, to work on these probs.  
Note that s. (ry. I doubt in some detail is not recursive.  
Recursion needs  $\textcircled{1}$  self replacement  
Equs  $\textcircled{2}$  boundary values.  
referencing?  
referent?

dynamics associates

Levin's Universal Search.

occurem updating the

They ~~cannot~~ ~~update~~

omit

There are two kinds of problems in our system, that are solvable, using L search. ~~They cannot~~ ~~update~~ ~~the~~ ~~system~~ ~~update~~ ~~the~~ ~~system~~. We will first describe these functions in some de

We will first describe L search, then show how it can be ~~applied~~ ~~to~~ ~~update~~ ~~the~~ ~~system~~. used to update the system.

Suppose we are given a ~~function~~ ~~in~~ ~~the~~ ~~form~~ ~~of~~ ~~a~~ ~~computing~~ ~~machine~~  $M(\cdot)$ .  $M$  accepts ~~any~~ ~~input~~ as arguments, finite strings that form a prefix set. The "prefix set" condition makes it possible for ~~the~~ ~~machine~~  $M$  to know when a string of symbols has terminated. The problem is to find an input  $x$  so that  $M(x) = a$ . ~~Meaning that after~~  ~~$M$  has read~~  ~~$x$ , it prints~~  ~~$a$  and stops~~. ~~is a~~ ~~known~~, finite string. ~~Many problems~~ ~~in mathematics~~ ~~can be put into this form.~~ ~~Solving equations, proving theorems,~~ ~~such as solving equations and proving theorems.~~  $P(x)$ , ~~solutions,~~  ~~$x$ ,~~  ~~$P(x)$~~

If  $T P(x)$  is less than the time of our machine in execution, we don't test that  $x$  at that time.

If we are given a probability distribution over possible solutions,  $x$ ,  $P(x)$ , we can use this distribution to guide our search.

One way of doing L search starts with a certain time threshold,  $T$ . We try  $M(x)$  on all strings  $x$ , ~~such that~~  ~~$P(x) > \epsilon$~~ .  $T \cdot P(x) > \epsilon$ .

But since  $M(x)$  may take an infinite amount of time for certain  $x$ , we ~~can't~~ ~~usually~~ ~~do~~ ~~this~~, so instead, we spend a maximum time of  $T$  on each string  $x$ . If  $M(x) = a$  over search is over. ~~computer is busy~~ ~~we don't~~ ~~test~~ ~~that~~  ~~$x$~~ . ~~test~~ ~~that~~  ~~$x$~~ . ~~130~~ ~~36~~ ~~C.P.M. is~~ ~~180~~ ~~96~~ ~~an~~ ~~160~~ ~~not~~ ~~read~~ ~~180~~ ~~Good~~ ~~200~~ ~~Good~~

The total time needed for this prefix set search is  $\sum_x T P(x) = T \sum_x P(x)$ . Since  $\sum_x P(x) \leq 1$ , this round takes total time  $\leq T$ .

If we haven't found a solution, we double  $T$  and go through the search again. This doubling and re doubling continues until we find a solution or run out of time. ~~To generate~~ ~~it is not difficult to show that~~ ~~the total time needed to find it~~ ~~is~~ ~~about~~  ~~$2T$~~  ~~to~~  ~~$P(x_0)$~~ .

If we have a good idea as to what the solution is before we do the search, we can put this idea into the probability distribution  $P(\cdot)$ , ~~and~~ ~~we~~ ~~will~~ ~~find~~ ~~a~~ ~~solution~~ ~~in~~ ~~a~~ ~~relatively~~ ~~small~~ ~~amount~~ ~~of~~ ~~time~~. ~~by~~ ~~assigning~~ ~~a~~ ~~relatively~~ ~~large~~ ~~value~~ ~~to~~  ~~$P(x_0)$~~  ~~and~~ ~~we~~ ~~will~~ ~~find~~ ~~a~~ ~~solution~~ ~~in~~ ~~a~~ ~~relatively~~ ~~small~~ ~~amount~~ ~~of~~ ~~time~~;  $2T/P(x_0)$ . knowing an upper bound on how long L search will take, makes it easier to design training sequences for the system.

There are two kinds of problems that occur in the updating process. Finding the recognition functions,  $R^i(\cdot)$  and finding the probability distributions  $P^i(\cdot)$ .

The criterion for acceptability of  $R^i(\cdot)$  is that  $R^i(Q_0) = \emptyset$  for a certain set of  $Q_0$  and that ~~it~~ ~~is~~ ~~1~~ ~~for~~ ~~a~~ ~~certain~~ ~~set~~ ~~of~~  ~~$Q_0$~~ .

We can use the mechanism discussed in Appendix I, to assign probabilities to strings — so that we can readily employ L search to solve this problem.

Spec A38.00



dynamics associates

SN Re: search for P(A|Q):

PARADOX?

I've been considering a 3 input func: input 1 describes a function of inputs 2 & 3. Input 2 is Q, input 3 is R.

If there is a function of Qw. ~~that~~ say, that gives  $A_0$ , then it should be poss to have a function that is 3 bits longer  $P(w. inputs = Q_0)$  with  $\ell(Q)=0$ , that gives same  $A_0$ .

436, 23

↓ My soln to problem

I.E. the length of input 1 + length of input 3 can be constant, & yield many  $Q_0 \rightarrow A_0$ 's.

T. forgo argument doesn't hold if the stack func has to work for > one  $Q_0$  value.

Well, it may be that for a single  $Q_0$ , a D soln. is usually a better. — That in "meagry" one puts to S func's, & seems much pc by "merging". My example in AVL in search area all D func's. If .07 is true, I may have to revise this R<sup>2</sup> P<sub>i</sub> model somewhat or plot!

Well actually, for a single  $Q_0$ , it may ~~rather~~ be true that a single code is by far best!

If one has  $Q_0 \equiv Q_1$  but  $P_0 \neq P_1$  then P has to be stochastic.

Merging ~~the~~ QA's that to be D, can give stochastic P's.

Now

A ~~set~~ S-P<sub>i</sub> can be broken up into "special cases" by proper R<sub>0</sub> recognizers

into several D-PS &/o S-P's. This sort of thing occurs in normal science!

Finding what was that to be a single disease, is best regarded as several different types of disease.

A set of Q's that are not identical, but are ~~expensive~~ expensive to be recognized, & disturb, can be regarded as identical & perhaps even  $S-P$ .

So the problem is, how to do updating!

There is that other model of stack operators, that I didn't like so much, A 2 input (rather than 3 input) func. 2 inputs: 1) desc of function 2)  $Q_0$ . We had many codes ~~(S-P)~~ (H(input) for a given output.

419.18  
419.31

Trouble is, I that I had a wire proof that 419.31 was correct!

Could I devise a U proof for 2 input func? At first I found a counter example to 2 input model — then found a counter example was wrong!

T. 2 input system does give a pd for each  $Q_i$  input, so the eqns of 419.31 apply:  $P(Q_i)$

Can we have a 2 input model universal? Equivaly, if one are allowed to use any ~~single~~ <sup>single</sup> input machine...

How can it simulate a 2 pd.?

I suspect that a 2 input

T. 2

253.21-24 has "counter example" to 2 input Model: 267.00 shows a counterexample wrong.

Hvr. note objection 267.10-11. — The whole thing <sup>(compressed)</sup> started  $\approx$  252.09-14

Hvr, as far as I can see for 3 input would work — it is provable:

Any way: Assuming that I'm using 3 input model, it seems that for even a single  $Q_0$  there will (perhaps a ~~ways~~), be a min code for it. i.e.  $F(Q_0, R) = A_0$ ; w/ a ~~way~~ <sup>ways</sup> and ~~an~~  $R = \Lambda$ : so  $F(Q_0) = A_0$  & F has a short code.

w. 2 input model, the probn. is obtained as a wtd sum of all of the codes: A legit code has to cover entire corpus, giving correct answer each time, & each code has a certain "length".

Very corresponds to ~~set~~ <sup>set</sup> of sequential coding such as, SPTB. — regular ALP. — Each code has to ~~do~~

dynamics associates

00:435.10 : code each bit of t. carries exactly. <sup>1.00R</sup> What if  $Q_2 = Q_1$  &  $A_2 \neq A_1$ ? Clearly d. system as derived fails!  
 T. codes are for a  $Q \rightarrow A$  function, so t. codes form a prefix set — like in Li-Vit's ~~idea~~  
 Universal Discrete distribution  
 T. 3 input code certainly deals nicely w. t. diffy of (1.00R)

1.07

04:432.16 **SN** Do have discn. about how t. QA form of induction is fairly universal — how one can easily "reach" what  $S \neq f(x)$  means, after  $\frac{d}{dx} f(x)$  has been read.

0.07 **SYSTEM**  
 "T. system as described is: 2 input Umc:  $Input_2 = Q_2$ ;  $Input_1 = P_{in} = z_i$  | Both inputs prefix sets.  
 t. acceptd assigned to an unordered set of  $Q_i/A_i$  pairs is  $z_i = 2^{-R(z_i)}$   $z_i$  is t. prefix set of decms  
 of ~~func~~ func's betw.  $Q_i$ 's &  $A_i$ 's such that the  $A_i$ 's are correct for all of them  $Q_i$ 's

So: T. problem now, is  $\neq$  435.00 diffy! But mainly, d35.11 is t. problem! For a sup

$R^i$  w. a single  $Q_2$  recognized, t. best model will be a 'D-function' (Deterministic, MTM).  
 After a new  $Q_2$  arrives w. same  $R^i$ , one can try ~~merging~~ merging.  
 The (common) exception is that a QA may be a D-Function (as in MapA).  
 So t. real Q is how to update? I write first looker pure MTM case. It may be poss. to work

These problems w. th. present update rules:  $407.03 - .14$   $\neq 409.00 - .33$

Going over 907.03 - .14! Condition 1:  $R^i \geq R^i(Q_{in}) = 1$   
 t. condition of acceptance is t. have, that t. P of Anti even  $Q_{in}$  is "Hyperf".  
 cond 2 same as 1 but  $\neq$   $R^i$  is not large enuf. 2 3 5 are; 2 3 5 is much ~~smaller~~ if poss.  
 cond 3  $R^i(Q_{in}) = 0$  for all  $z^i$ . The approach given is 0.4 for MTM: for NMTM, it  
 may not do harm, but its not yet of much use. Its probns will tend to be not so hot.  $\neq 437.00$

23:435.05 **SN** When we have a really new  $Q_{in}$ , we note d35.00 - .05 : w. a 3 input one, there will be  
 many **S, R** pairs that have some total length  $R_i$ , so if we do find  $\gg 1$   $\leftarrow$  S, R pair  
 then we have a good chance of  $R \neq \Lambda$ , so we get some extrapolation. So  $P$  is  
 May solve t. diffy of 435.00 - .05! Hvr, I'm not certain that ~~essence~~ soln. of t.  $SSZ=1$   
 case is usually helpful for  $N > 1$  cases — i.e. that it extrapolates at all well!

[may be able to do a lot of (range of MTM types using not so clever heurs —  
 sort of "Minimal Heurs"; But to get even (range during Lsrch) t. system has to  
 be able to do NMTM!

I used to think that induction could be usefully done w.  $SSZ$  of 0, 1, or 2: Certainly t. ~~form~~  
 formulae of 414.18 enable this. But my R, P formalism may not be so good for it!  
 Well, maybe that's O.K. — The at first TM may need  $SSZ$  of 2 or more for probabilistic  
 induction — but eventually, it will break out of t. R, P restriction & be able to  
 do more "advanced" kinds of induction.

**So** I want to know if t. MTM problems would be  $\neq$  solved by t. 5 steps of update!

Also, some ideas of how to deal w. NMTM. Also if not how to deal w. ~~the~~ MTM!  
 for NMTM, think of some simple  $SSZ = 0, 1, 2, 3, \dots$  cases.

Handwritten scribbles and symbols at the bottom right of the page.

N/7/01

ID

436.22  
spec  
spec

dynamics associates

on update rule  $V_i: R^2(Q_{nt})=0$ ; For MTM: 0.4.  $\pm$  Guess,  
for NMTM: unclear: I do have to rethink induction for use of 0, 1, 2, 3 ... 4(4.18) is probably

Correct, but how to approx it in a good way?

For  $SN \equiv N=0$  we just make  $p$  of  $A$  depend on how easy it is to construct  $A$  using available concs.  
For  $SN \equiv N \geq 0$  we do same thing, but  $\pm$  thing we have to construct is  $A_0, A_1, \dots, A_{N-1}, \dots, A_N$ . (construct using some available parts). One might think that  $N > 1$  makes it possible to reuse invent new concs:

But no new concs can be invented as part of a single function used to express  $A_0$ .

This is done as if  $SAZ (= \text{Augmented } Z(4))$ . I think its essentially what update 2 b does.

Update 4: Several  $R^2$  accept  $Q_{nt}$ , but none one is best for  $A_{nt}$  is good and.

Modify  $R^2$ 's so that only best to good  $R^2$  accepts  $Q_{nt}$ . Seems ok for NMTM = STM | vsr DTM = MTM

~~It is poss. to do 4a & 4b~~: 4b is same as 5b & involves Moving

Note that in general, more full moving of all accepting  $R^2$ 's may not be feasible / desirable - but try to merge as many as possible - but may be only zero feasible.

in condition (40703-14) of update: Every thing works fine - apparently "no update needed"

False! When  $Q_{nt}$  comes in, it is added to the corpus: At minimum, "continuous" params have to be updated;

Since no. of concs  $\uparrow$  by 1. Another poss is that the freq. of use of 1 or more concs goes over threshold"

So we are able to Legitimately delete a new conc. (thos  $\uparrow$  PC).

This kind of update is very fast, hvr, it doesn't involve search

```
Ghostscript
DLL C:\ghostpnt\gs0.01
bin\gsdll32.dll
it in C:\ghostpnt\gsview
is 27file by 3120/00
```

Do fig note on CPM (cond probability ~~machine~~) C of Learner.

(1) for ~~40703-18~~  $\pm$  31 have eqs nos

for receivers, the S should not be summed to  $\infty$  - just "j"  
T. Summation is a spread in text. Later receives eqs are used w. for the same.

Is it made clear later that when we use very poor eqs - create "Maxim" Rank by not summing to  $\infty$ ?

This must be made clear. See in subsect.

[SN] At all times "update" times, the goal is to max 414.18. Hvr, with limited time available,

there are at least an order list of tasks. First to do the 5 update things but do the "a" parts which are quicker. If time remains, do "b" parts. Mzple Stable Eded Abcdetghijkl

26 : 431.20: in spite of immense power which is particularly appropriate training sequences that are

usually not very carefully constructed. This might be explained by considering that people

usually start out with lot of "prewired" induction techniques based on many millions of years of

genetic selection. The system being proposed is expected to become much less dependent on training sequences after it has learned certain important induction techniques.  $\rightarrow$  431.29

on training sequences after it has learned certain important induction techniques. substantive initial training

dynamics associates

00: (434440) To understand the application of Lsearch to ~~probabilistic~~  $P_i^j(\cdot|\cdot)$  functions, we will first consider ways to represent ~~probabilistic~~ them, using universal Turing machines. The discussion readily generalizes to other methods of representing functions of binary strings.

Consider a 2 input machine: Both inputs are prefix sets. Input<sub>1</sub> is  $Q_i$ ; input<sub>2</sub> is  $R$ . The 2 inputs are a complete ~~prefix~~ prefix set, so  $\sum 2^{-|k_i|} = 1$ . ( $k_i$  is a word).

e.g.: 1, 01, 001, 0001, ...  $\sum \frac{1}{2} = \frac{1}{2} + \frac{1}{4} + \dots = 1$ .

08: However, this prefix set could not represent certain P.D.'s. e.g. it couldn't represent  $2, 2, 2, 2, 2, \dots$ . Since all other prefix sets would have to be used and 1 (w.  $P=1$ ) could not, 01 could not.

Another prefix set is to set all binary strings of length  $n$ . If  $n$  is large, this can represent many distributions.

10: Unless  $n$  is large, this will not work! Each  $Q_i$  needs its own prefix set. If  $t_i$  machine reads  $Q_i$ , then reads as much of  $R$  as it wants, then starts to print  $A_i$ , then  $t_i$  inputs to  $R$  must be a prefix set; one input can't be part of another.

How to get Lsearch (or AZ141) to generate binary functions, is unclear.

For  $Q_i$ , maybe a binary string w. a "blank" end symbol — to this defines a function on  $R$ . If  $R$  is unidirectional, and we acknowledge "A" as the output when  $t_i$  machine stops, then "acceptable  $R$ 's" are a prefix set.

17: I could write a paper very to Turing formalism.

- T. Turing instructions could read  $S$ , then read  $Q$ , then read  $R$ , then output.
- T. insts could be arranged so it does  $P_i$ . As soon as  $M$  reads a  $Q$  bit, it can't read any more  $R$ .

Shohazky  
Variable Storage Machines.

Schöhlage:  
Variable Storage Machines

20: I could generate any function of 2 variables ( $Q_i, R$ ) via AZ141. Actually, for testing, it would be best to have one read  $Q$  first, then  $S$  is  $R$ .

23:  $Q \rightarrow$  except if I wanted to do several  $Q$ 's w/ "blank"  $S$ , so probably do  $S$  first. So I could "construct"  $S$  via AZ141. But I'd like to get AZ1 to generate functions of 2 variables.

Well, if there are only 2 objects to be inputs (in addition to 0, 1,  $\pi$ ,  $e$ , etc), then that's it!

so in AZ1, I generate a seq. of functions; each function defn. can use previously defined funcn. Any funcn can have 0, 1,  $\pi$ ,  $e$ ,  $Q$ ,  $R$  as 1 or more arguments. — or previously defined funcn.

Since  $Q$  is finite & well defined — no problem w. that argument. for  $R$ , it's not clear how to.

AZ1 function should deal w. it. Perhaps have ~~some~~ to AZ1 function include a "no. of bits" integer, when it asks for  $t_i$  value of  $R$ .

28: Using Turing formalism at the Sept 8, 2001 "postscripted" notes; I may want to allow  $Q_i, R$  to be one of the "constants" — so they can occur anywhere a ~~function~~ final or intermediate function.

30: 31: If we add a binary number of  $\geq$  variables, then  $R$ 's OK. —  $t_i$  choice of values for each of inputs is still  $Q, 1, \pi, e, Q, R$ .

A possible trick:  $t_i$  no. of bits read from  $R$  is a function of  $Q$  only, — this means that for each  $Q$  we have a fixed residue of  $P$  is a minimum value (other than zero) — this is  $2^{-n}$ ;  $n$  being binary length of  $R$ .

Ideally, I want to determine to read the next bit of  $R$  to be a function of  $Q$  & all previously read bits of  $R$ . Perhaps read about Chaitin's Micro Lisp to get ideas on how to do this.

My AZ1 is a rather limited device anyway (no recursion): I could have a sequence of reading of  $R$ .  $F(Q, R_i) = 0$  or 1; if it is 1,  $R_{i+1}$  is generated. If not,  $R$  is of length  $i$ . So  $F$  has to have  $\geq$  output 1 and 2.  $A$ .

497. 22-23  
falls how to do this

dynamics associates

00 : F38.101: Gentry beate Lush; work on the Searchman over  $R(S) + A(R)$  constant.

So we present 4. function w.  $Q$  is a R sequence of length n. It has already been seen  
that  $R$  is injective. We can arrange so that we give  $f$ , function a  $R$  of length n,  $R$  can be  $R$

How the prefixes that were repaired.

But it may. It could express its "non recognition" by not producing an A, or not stopping

I.e. stop w. A output or not stop w. A or any output.

That into be ok, say  $X$  is a ~~string~~ R string that has no prefix that caused  $K$  unit to give

acceptable output. Call prefixes have been <sup>tried</sup> (hmm). We know  $X$ . If  $X$  produces long output,

Then, no extension of  $X$  is used as  $R$ . If no usable output occurs, we try  $X0$ ,  $X1$  for  $R$

on next Lush "round".

Will this ~~work~~ work? remember No "legal output" means loop times  $> T \cdot 2^n$

On next round, loop times will be  $2T \cdot 2^{n-1}$ ; — The same! No!  $S$  can change in long tests.

So, its not clear! ~~Value of T~~ Value of T — ~~but~~ and this  $>$  prefix set!

Oh! Even round we have a new several prefix set!

So, the big  $Q$  is! can we do this using Lush (using randomly chosen knots & partners)

In Lush, as soon as a trial terminates, it immediately cuts off all trials that

are extensions of it! — So that would do it!

In Lush, if after finding a successful trial, we could still have prefixes of that trial

that haven't yet converged — so we don't know about whether they will or not —

But still, for given  $C_B$ , we do define a prefix set. (is this conversion to  $T < 2T$

Lush also). So perhaps, for each  $C_B$ , we ~~can~~ since we have a distinct prefix set,

we also have distinct predictions!

For Lush, for a given  $S$  function is a prefix set for that function we should only try  $R$  values of all

the same length — so would have to remember which are legal (is, have not completed prefixes)

[It looks like Lush way not so easy w. pairs  $A \leq 1$  function test!]

FN on TSD:

<sup>Interact</sup> Giving type ~~for~~  $Q = 3+5; 7+1, 8+9$  act, could be a strong "Hint" In general, type in deciding can be a strong kind of "Hint".

The  $T \rightarrow 2T$  method seems to work O.K. with AZ14 (attempts for S functions).

Next list  $R$  S functions in  $T$  &  $2T$  order. We plan to do an exhaustive search until  $R$  values, starting

w. smallest  $R(C(K)) = 1$ ,  $R$  can  $R(K) = 2$ .  $R$  can 3... etc. We reject  $R$  values that don't converge in

~~limit~~  $L$  with  $C_B = T \cdot 2 \cdot R(K) - R(S)$ . So for each  $S$  function, there is a prefix set of  $R$ 's.

We'll accept the best of  $R$  &  $R$ 's that have converged for a given  $T$  (try) — At each  $T$ , in time for

we will have a first  $T$  given  $T, S$ , we will have a legal value of  $R(K)$  that we are working on;

Then  $2 \cdot R(K)$  possible prefixes will be divided into legal & illegal prefixes: we try to continue on

the legal prefixes. The illegal prefixes have a prefix that has already converged.

It may be possible to keep track of any very cheap! Look at  $R$  exploration space for  $R$ .



T: routine may be used to identify to  $T$  search that I did in ~~the~~ 1984 report I wrote. Exhibit 10 Optimal Sequenced Search was the title.

1985

440.00

dynamics associates

00: 439.40: Re: parallel Lsrch: we don't have to do this, because we aren't doing OZ problems via Lsrch.

(tho, I may end up needing it for that!). On the other hand, || Lsrch is more efficient — by a factor of 2 (at least) & (probably) by a factor of 6: search depth.

In || Lsrch, we store all partially run strings in their assoc "maxima states".

— so it looks like it would be poss. to do Reiz Lsrch, since we will know which trials have converged.

No! — Its begining to look like || Lsrch w. S & R will not be any different from || Lsrch w. R alone. (i.e. regular Lsrch).

Consider just Lsrch on R alone: In + present situation, R = x we doesn't converge in time ↑; but it is poss. for x1 to converge in time <= T!

This would not be true in usual Lsrch using unidirectional IO eyes or suitable

"prefix machines". In +. present case we always "try" x before trying x1.

"Try" means a complete trial: either it converges or "times out" using present CD.

The str prefix set differs for each "Round", & time to do each round is <= T for that round.

For || Lsrch on R alone: its not obvious how to procede. procede.

Suppose for fixed S, the (convergent) R do not form a prefix set?

They have to: If a R input converges, we don't try its extensions.

Well, in T <= T, they do for each T value.

Couple of Q's: Say S, Q has output in short time for most all short inputs (R).

So R = 0, R = 1 bring me out pots & put's it! Each time PC = 1/2. Only 2 A's

I really w'd like an R to exist for every poss A's!

Actually, (0.22) isn't really so smart! the S, Q function has to end after t: R input

BUT Many functions will do that.

I really haven't done any stochastic TSE's. Perhaps try some elementary Bern sequs:

We index t. data w. names of T & S Q so th can tell they belong together.

I've certainly written a lot of stochastic induction problem.

Another kind of stoch induction: Q: a A -> b, Q: b A -> a

A common type of stoch funct. is a table: for each Q, we have probab of several diffrent A's.

Hvr, using 2 unc w. 3 inputs is a probably adequate soln to a stoch TDn problem.

I guess what I don't like is — I don't see how E can fit (4) into it.

Hvr, I'd like to be able to show how it would work for t. learn of (.26) & t. table of PC's (.29).

In .26 (TLU), we can make separate studies for Q = a & Q = b.

I am reminded of that Guy who tried to do curve fitting using a ALP! He did a discrete model

Search, using a MDL, then some continuous opten methods to fit to continous param

The fit is only a small part of probabilic modeling — it is an huge one. [One has to take a Hessian for each

final pc, hvr.] —> 441.21 spec

Abcdefghijklmnop  
Abcdefghijkl pqr  
for wxyz abcdefghij  
| | | |  
abcdel  
fghijkl abcdef  
i j k l m n

dynamics associates

00: 440.40: Superficially, the problem is: There are 2 ways to represent S-dracks:

- 1) 3 in pot ump
- 2) AZ141.

1) is formally fine, but doesn't seem able to make hours of input.

2) frag. of word concs.

b) looking for common sub-words.

2) is able to do .05 & .06, but I haven't been able to get it to do Lsrch properly. — T. different reason for this (13)

5/11 Looking at "Opt. Seq. Snd.": Comments on folder:

Remarks & Q's by Eric Wintrae: Very Good! I really have to answer three Q's!

Also, that "OSS" paper did seem to have a general Model for TM ~~machine~~. Lsrch

(Reading Win's remarks, suggests that there are some holes in my understanding of Lsrch!

Mainly Re: OZ Lsrch, hvr.

to funct of AZ1 (as of present), simply take ~~any~~ prefix code word as inputs. This appears to be because of the structure of the functions ~~to~~ it can create, RE: Umes that accept only prefix code words: The thing I find unique about them is that their decision to take another bit, was based on what they had read thus far. Actually, any device that ~~accepts~~ knows how to ~~accept~~ seems to solve the problem

accept ~~the~~ code words from a prefix set, can do the same thing. → See 447.22 - 23 for How to Do This

So, in AZ141, each function must be able to look at any infinite string of bits, & cut it into code words (if the code words form a complete set.)

Another kind of "simple" ~~problem~~ S. induction problem: A single Q is a long sequence of symbols (or of words). The ~~to~~ TSEQ consists of ~~the~~  $X_0, X_1, \dots, X_{n-1}$  for  $i = 0$  to  $n-1$ . ~~the~~  $X_i$  are words; ~~the~~  $f_i$  cond. will be  $P_i$ 's on words that are functions of  $i$  nos. in  $Q_i$  (an early no. of nos. in  $Q_i$ ...).

Now if we have a function on  $-\infty < z < \infty$  (and  $i$  nos. in  $Q_i$ ) that is SO and indep of  $Q_i$ ,  $\int_{-\infty}^{\infty} f(z) dz = 1$ , then it is a cond. for a probab. distribution.

We could, of course take any function  $f$  squared to get it  $> 0$ , then normalize it on  $(-\infty, \infty)$  interval. (if it were normalizable) The "squaring" part seems, hvr: In general we want the function to be more "naturally" pd's: like  $t$ . pd's we get in ALP. → 447.24

One way to get nice, easy prefix sets: use space symbol ( $=$ ; so  $Q_i$  1, 2 and 3 symbols)  $=$  terminates a  $=$  always terminates a word. So we have binary words separated by spaces ( $=$ ). T. system is complete because any binary string can be uniquely parsed into words.

$$\sum_{\text{word } pc} 3^{-l(w)} = \sum_{i=1}^{\infty} 2^i 3^{-i-1} = \sum_{i=1}^{\infty} \left(\frac{2}{3}\right)^i \cdot \frac{1}{3} = \frac{1}{1-\frac{2}{3}} \cdot \frac{1}{3} = 1. \text{ So it works as expected.}$$

36 Since smallest word is  $pc = \frac{1}{3}$   
37 We cannot represent  $t$ . pd.  $\left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right)$ . See 442.37 for a kind of soln. — leave out  $\frac{1}{3}$  — use incomplete prefix set!

dynamics associates

As noted before: 441.36-37 is a necessary dirty in any system using a same-prefix code for all  $Q_j$ 's. Somehow,  $Q_j$  has to be able to close the prefix set for representing  $p_j$ 's. — Or, more exactly, " $P_j$ " " " " " " close " " " " representing  $p_j$ 's,

Perhaps look at Cover's demo, that any PD can be represented by a suitable Turner. Is it very much like Arith. coding? In arith coding, one feeds in an infinite binary seq, & every once in a while a bit of the coded sequence emerges — which is a kind of "parsing" of the input seq.

For binary ~~seq~~ data, the PD ~~can~~ always has at least 1 pc  $\geq 1$ , so the arg of 441.36-37 doesn't apply

Also note that ~~even~~ if we use a Base Binary system (like almost all digital computers), then if we use words of several bits (not necessarily same size), then we can always assign pc's to individual bits of all  $A_i$ 's. Hvr, Arith coding gives just 1 code for each  $A_i$  — and it's effective — sometimes no output for long address input, then a long seq. of output symbols!

What about Huffman Coding? It is certainly not v.e. for pc's close to 1! — & I will be having lots of d-funcs to try. — which have pc's  $\geq 1$ .

Write up detailed reasons why  $i \geq 1$  input models &  $i \geq 1$  model ~~are~~ inadequate. So an external (or other) can work out.

List possl. approaches, alternative models, approaches.

Also the method of 250.09-19: Counterexample 253.21-24 ← (counterexample 267.00)

Nail Counterexample 267.00: Actually 267.00 gives instructions for realizing the D.P. of 250.09-19!

This Idea is: for each  $i \geq 1$ , we have a set of possl  $A_i^j$ . Then pick some large int  $z$ . Then we have  $z$   $S^r(i)$  functions:  $S_i^r = A_i^j$   $\{k^i\}$  : The probability  $P(S_i^r = j)$  is the probability of  $A_i^j$  being  $\frac{j}{z}$ .

A way to get extrapol. out of such  $S^r(i)$  funcs:  $S^r(i)$  gets  $r$  most likely choices for each  $A_i^j$ . For each value of  $i$ , we have  $(0,1)$  divided into  $z$  divisions  $(0, \dots, 1)$ . We place the  $A_i^j$  probabilities over that interval (they sum to 1), we put the most likely at left & in increasing order of pc's (most likely on left).  $S^r(i)$  will have prob  $\approx A_i^j$  that occurs over the  $r$  point: it is at  $pc \approx \frac{j}{z}$ . So, in general, the  $S^r(i)$  functions over an order of  $r$ ,  $\{r \geq i$  is most likely  $\}$

The number of  $S^r(i)$  functions that will have the value  $A_i^j$  will be a proxy of  $A_i^j$ .



O.k., so say we have a lot of  $S^r$  functions. Each one gives  $A_i^j$  as a function of  $Q_j$ . No one  $S^r$  function assigns much of a pc to a corpus. Also, if the  $S^r$  functions are expressed in the  $R^k P_i^k$  form, it's easy to get lots of them. If we have "a" values of  $k$  for each  $P_i^k$  ( $k = 1/n$ ) then that's an infinite  $S^r$  functions.

If the  $S^r$  funcs (or the  $P_i^k$  funcs) all have diffrnt pc's (= wts.), we will usually need more of them to get good pc's than if they had all = wts.

So it looks like .25 wire work! (it is 250.09-19), 253.37, <sup>pd's</sup> Couple of Q's Is it method universal? (2) Can we easily express known  $p_j$ 's using it (can we express certain word-hears using it)?

On first plot, this seems like a very rough method of deriving  $S^r$  funcs. — It is very tedious all of the  $S^r$  funcs: the total length of all them is the total data length! Hvr, many have  $n$  decms, so  $\Sigma$  best is not as bad as that! (Tho its still pretty bad!)

HA! Maybe a way around the objection of 441.36-37: say pc's are not measures but semi-measures — so we don't have to use the "1/2" word if we don't want to! — we do have to normalize, hvr, since we'd be using an incomplete prefix set. A big trouble is normalization! We have to get pc's of



**dynamics associates**

00: 442.40: All of by wt. A's! Maybe estimate it: try ~~weight~~ vary by wt. R's: The ones w. no output; we can add up their pc's & get norm  $\approx$  normal. constant.

So b. code of 441.32 may be o.k. Binary strings w. space (or  $\epsilon$ ) for end of word.

OR maybe don't normalize (except perhaps <sup>Wrt.</sup> R's that don't converge)  
If a none of the A's have large probs... Well,  $\sum p(A_i) = 1$  (e. ~~one of them must occur~~!)

So what's response to space, ( $\epsilon$ ) when no A has  $p > .25$ !  
Well, say there are only 4 outputs A's & all have  $p = .25$ .

Is it possible for  $\epsilon$  to mean "no output" — is this output from  $\Lambda$  ( $\equiv$  null string)?  
Say the R's for  $\epsilon$  & A's were 000; 010; 100; 110; All other inputs (including  $\epsilon$ ) gives no output. So  $\epsilon$ 's would actually mean that we A = no output would be rather likely!

Only if we allow this is impossible; can we "normalize it out"? .27  
Well, in all cases of update, we know that no output A must <sup>occur</sup> ~~occur~~. — So ~~it~~  
We'd be justified in "normalizing out"  $\epsilon$  & nulls.

Still we have to put in a lot of inputs, to get to normal constant!

$$\sum_{i=0}^{\infty} 2^i \cdot 2^{-i-1} = \frac{1}{2} \sum_{i=0}^{\infty} \frac{1}{2^i} = \frac{1}{2} \cdot \frac{1}{1-\frac{1}{2}} = \frac{1}{2} \cdot 2 = 1$$

$$\sum_{i=0}^{\infty} 2^i = 1 + 2 + 4 + \dots = \frac{1-2^{-n}}{1-2}$$

$\frac{1}{2}$	$(\frac{1}{2})^2$	$(\frac{1}{2})^3$	$(\frac{1}{2})^4$	$(\frac{1}{2})^5$	$(\frac{1}{2})^6$	$(\frac{1}{2})^7$
1	.5	.25	.125	.0625	.03125	.015625

so, to get .058 accuracy, we have to go out to length 7:  
 $2^7 = 128$ ; so  $\approx$  256 trials. — But usually we don't need much precision!  
 $\approx$  29% error multiplier of length 3  $\therefore 2^{30} = 16$  trials.

But actually — Superficially, in R, the symbols (bits) all  $p = \frac{1}{2}$  — but actually, Not so! The pc's of 0's & 1's are about the same, but pc of  $\epsilon$  is much less!  
Consider a sequence of R's in a corpus code. The freq. of  $\epsilon$  will be  $\ll \frac{1}{2}$ .  
So we can legitimately recode the R's.

27: .10  $\rightarrow$  Another Q: If single symbol ( $R = \epsilon$ ) can be a legit ~~and~~  $\epsilon$  for Lisp functions.  
Her. If an arg has to be a number,  $\epsilon \rightarrow \Lambda \neq \phi$  at all.  $\Lambda$  would not be a legit number.

Unless,  $\epsilon$  string means (as a number) "clear" <sup>clear</sup> to register; Man loads strings bit by bit.  
If there are no bits in string, result is ~~zero~~ zero. <sup>signific. digits</sup>  
[How is "zero" expressed in floating pt? it must be  $(\phi, \phi)$  since  $\phi$  means  $|\phi$   
Maybe  $\phi$ , -511 (+ smallest poss. no.).

34  $\rightarrow$  O.k., so if there are no A's in  $p \geq \frac{1}{2}$ , then for the "value of R, there is no output — which means "undefined" — the system ~~does not~~ machine does not stop @ user, if that's put occurs.  
(Or maybe it doesn't start? — No, it has to start to read  $(\phi)$  Well, if it read a binary string of length  $w$ , followed by  $\epsilon$ , it should look for A's w. pc's  $\sum 2^{-w} \cdot p(A)$  — it may find none: it ~~doesn't~~ Doesn't stop.

dynamics associates

00: 445.40: How to code a set of  $A_i$  pc's, using  $\Rightarrow$  prefix set of  $pc_j$ 's;  
 Put A. A, pc's in order of binary, largest first;  
 Put B. PR " " " " " "  
 Prefixes A pc's and largest R pc's that is  $\leq$  to be assigned it to A. Take the remainder of it. A pc is  
 into B. A pc list, in its proper order. Go to next pc in A pc list, a loop to X. Will this eventually use up PR  
 A pc's? If there are some PR pc's that are  $>$  any A pc, they will never be used, so the total PC of A  
 (which is 1) can't all be killed.  
 Also, if B. PR pc's are of finite cardinality, there are more A pc's than that we can do it either.  
 So say B. PR pc is of ~~finite~~ Card =  $\infty$ .

Even if -03- of Doesn't work, it may not be bad! ... Maybe Usable!

09 10 11/13/01 Any way 443.34 suggests that we can put any binary string into the Lisp function.  
 It can then select a subset of those which it regards as "defined" for, it does not stop.  
 Trouble is "binary string" is a complete prefix set. Any proper subset is incomplete.  
 Say B. functions selected only  $\omega, \alpha$  as defined. In one sense, this is a complete  
 prefix set. — In another, it is not. (i.e. if " $\omega$ " is included.)

15 16 439.27 - 40 seems O.K. (in present, does of "defined"), but notes 440.10 - 20: Has not meant notes (440.15 R)  
 so to legal R do form a prefix set. But say  $x0101$  was a soln. to an INR problem. Suppose X converged.  
 We would never ever try  $x0101$ ! In our present problem, if there were two, then it's not in S, ~~should be~~ must be  
 discarded!  
 — Fine.

On the other hand: The pgm on pg 9 of OSS 1985 is also ~~under~~ under some ~~how~~ how!  
 I seem to have forgotten how Lsrch works! T. Lsrch is over pgms that are able to  
 produce a soln:  $p = M^k(\alpha, M, X)$  ; "Test" is: IS  $M(p) = X$  true?  
 $\{\alpha\}$  form a prefix set. L problem. The p's tested need not form a prefix set.

20 We are looking for  $R's \Rightarrow F(Q; R) = A_i$ . In this case,  $\{R\}$  form a prefix set.  
 we want short R's. — So this problem is really easier than normal Lsrch.

If  $F(\ )$  is a TIME, we can easily get  $\{R\}$  to be a prefix set. — we feed in R sequentially,  
 when the machine stops — R's a member of  $\{R\}$  out in pot.

Look at Chaitin's Little Lisp for (perhaps) how to do this. [ "Pure Lisp" has no "side effects" ]  
 (I'm looking at a talk by Chaitin at Santa Fe inst. 200% sure, he starts talking about bit string math  
 Chaitin TLisp. htm 47k 11/13/00 on pp200 chaitin folder on 'S-2 on F-2  
 1) read next bit 2) read rest (legal) sexpr<sub>exp</sub>.  
 "Try" can fail via CB or via end of bit string.  
 I have his Lisp interpreter in Mathematica. (Lisp-M  
 Ch. uses stack length rather than time as CB.)

Lisp:  
 Chaitin's TLisp  
 . This is a simple, no-  
 cheating! Try Lisp  
 Chaitin TLisp. htm  
 n a same folder!  
 How this differs from  
 Chaitin's TLisp  
 to unclear

30 (try, (CB, Time limit), sexpr<sub>exp</sub>, binary data) ; when it  $\rightarrow$  CB; output is (value) CB (how much data used) (displays)

McWalter's [initial] Devn of Lisp (in Chaitin Folder Place), seems very readable.

32  $\rightarrow$  15-16 was a serious Objection to my model, but (16 R) answers it!  $\leftarrow$

For a given CB, we will have a tree containing a prefix set at terminals. As we  $\uparrow$  CB  
 the size of trees will never  $\uparrow$  & will usually  $\downarrow$ .  
 Arr, say we are in this tree with given CB & we insert 1, 1, 1, ...  $\infty$  & move of from  
 converge. We will spend a total time of  $T \cdot (2^{-1} + 2^{-2} + 2^{-3} \dots) \sim T$  on that branch alone.  
 On each such branch we could spend  $\sim T$ . Well no, we spend  $T \cdot 2^{-k}$ , where k is the  
 length of R at which the infinite branch starts. So  $Q$  is  $\sum 2^{-k_i} < 1$ ?  
 $k_i$  is the length of finite branch starts.  $\rightarrow$  See 445.02  $\rightarrow$

dynamics associates

00: 444.40 : My Letter to Learn do it with the problem of a UMC that did not have unidirectional I.O. Its methods should be applicable to present problem.

02: 444.40 : In worst case scenario:  $S$  ~~never converges~~ <sup>for</sup> any  $R$ :  $\leq 2^{-kx} = N$ ,  $N = \log_2 T$   
Here we assume that  $T$  starts at "1" which is time for 1 instruction so  $N = \log_2 T - \log_2(T_0)$   $\leftarrow$  time/instruction.  
This may not be enormous, but it is  $\gg 1$ . (!)

So look at  $\{00\}^{\uparrow}$

Even if R. Method of 444.32 - 40 work, it would be very inefficient - it would vary dec. functions. R would accept <sup>only</sup> prefix set as input. Even if we had a good way to detect such functions, it would be inefficient.

On the other ~~hand~~ hands, t.3 input UMC automatically accepts only prefixes.

Its like **Koza** using Lisp expressing & crossing train over to get meaningful expressions every time.

Well one way to do it would be to Chaitin's method: to devise a language that would have a effectively unidirectional I.O. - for all inputs. Actually, the functions of first ~~inputs~~ seem under control - only the 3rd input ~~to~~ ( + R input) has to be fixed up.

Well Chaitin had it so that a try could end if  $\leq B$  / exceeded or no. of bits requested was exceeded. In the latter case, one repeats ~~with~~ bit string with one bit (Pro this is very inefficient)

Another task: See how I do induction problems, & try to put it into a formalism of

Hypoderm of Corpus

All I need is a method to get various R values various A values

One way is ~~to~~ state CFG w. probabilistic productions. In this case, we are asked for a (usually finite) no. of choices: each w. its own pc distribn. The function  $AZ(A)$  are generated in this way, & after each limit is generated, we have choice of a "stop" symbol (i.e. final function is defined),

we can derb strings in this way.

we can continue down of  $\epsilon$ : S function  $\epsilon$  w. known input  $(\epsilon, Q)$ . Actually, one or more of  $\epsilon$  inputs is normally part of the function down!

Say "con" (= connect) is one of the (perhaps primitive) funcs.  $con(0;1) = 01$   $con(0,001) = 0001$ .

so  $con(x_1, con(x_2, con(x_3, x_4))) = x_1 x_2 x_3 x_4$  (say  $x_2 = 001$ ).

$F(Q, x_1 x_2 x_3 x_4)$  would be  $F(Q, con(x_1, con(x_2, con(x_3, x_4))))$

which is the down of a function. if  $\epsilon$  is a community occurring function (form), it will have  $pc \sim 2^{-4}$ .  $p(\text{end symbol})$  - which looks like  $4/16 = 1/4$ . If  $p(\text{end symbol})$  is  $\epsilon$  (small), then  $\epsilon$  will be the smallest part of max pc of an A comb. - so if we had  $\frac{1}{2}$  A's of  $\epsilon$ 's, they could all have  $pc = \epsilon$ . But if we had  $> \frac{1}{2}$  A's  $\epsilon = pc$  we couldn't do it.

$$\left(\frac{1}{2} - \frac{\epsilon}{2}\right) \left(\frac{1}{2} - \frac{\epsilon}{2}\right) \leq \epsilon \cdot \sum_{i=2}^{\infty} \left(\frac{1}{2} - \frac{\epsilon}{2}\right)^i$$
$$\frac{1}{1 - \frac{1}{2} + \frac{\epsilon}{2}} = \frac{1}{\frac{1}{2} + \frac{\epsilon}{2}} \sum_{i=2}^{\infty} \left(\frac{1}{2}\right)^i = 1 + \frac{1}{2} + \dots = 2$$
$$\epsilon \leq 2^i \left(\frac{1}{2} - \frac{\epsilon}{2}\right)^2 = \epsilon \sum_{i=2}^{\infty} (1-\epsilon)^i = \frac{\epsilon}{1-(1-\epsilon)} = 1$$

good!

So if  $\epsilon$  is small, it works fine, but then we could have long binary strings  $\left(\frac{1}{2} (1-\epsilon)\right)^n = \left(\frac{1}{2}\right)^n \frac{1}{2} e^{-n\epsilon}$ .

**dynamics associates**

445.40: Re: T. trouble of 445.32-34. T. case in which All A's have = pc & there are many of them

is a Bad production Situation!

If we only had  $\epsilon$  symbols, we would run into some trouble, but it may be usually not very bad.

But the fact that this thing doesn't not work exactly is worrisome — I'm afraid it may have some other very bad properties!

A possibl. way to do it: Assoc. w. Q are 2 functions!  $f_1$  looks at (partial R input) &

decides whether it wants another bit or not. After R has ended  $f_2$  operates on Q & R to obtain A.

Essentially, what  $f_1$  does, is create a prefix set on the R input.

This would seem to make finding a  $f_1, f_2$  pair that works, have much more cost than just a simple single function! But maybe not: we may be able to find a set of  $f_1$  functions

that covers all cases: say  $f_1(\frac{1}{2} - \frac{\epsilon}{2})(\frac{1}{2} - \frac{\epsilon}{2}) \in$  runs of 445.35ff: By choosing  $\epsilon$ , we are able to fit a greater variety of p.d.'s.

We might do this coding of the A's using 0, 1,  $\epsilon$ , then optimize the  $\epsilon$  value for that particular Q → 442.00

Another aspect of .05ff! T. K values don't have to be "complete" because they have to be eventually removed (since some values of R → "Undesired A").

Normal can be successoring, but. Another point is that since  $\epsilon$  is probly of "undesired" →  $\phi$  we are  $\epsilon \leq \epsilon \leq \epsilon$  CB, this may be less of a problem.

It would seem that no variation of  $\epsilon$  (i.e. uses same  $\epsilon$  for all Q's) normaliza. would fix up objection of 441.36-37 (i.e. 442.37): (or see 443.03)

Re: calculator, but 445.06-11 — on Efficiency

So to model of .05 may be O.K. 445.02 was a killing objection to a certain model

obj. of 444.32 maybe the model was 444.09-16? Anyway, consider if every 60 so it's 445.02 is relevant to the model of .05

So 443.34 — 444.09 — 16 Then 444.32 — 40 & 445.02 — 04 was a big objection to .05

While we might have this prefix set, we had to go way to discover it.

It looks like the objection of 445.02-04 does not apply! The  $\frac{1}{2}(1-\epsilon), \frac{1}{2}(1-\epsilon), \epsilon$  system

is a prefix set for each value of  $\epsilon$ . In doing this, we are governed by 2 params: T &  $\epsilon$ . A trial will terminate if  $(\frac{1-\epsilon}{2})^n \cdot \epsilon > T$  i.e.  $T \geq T \cdot (\frac{1-\epsilon}{2})^n \cdot \epsilon$

For a given  $n$ , whether we are over CB or not depends on both T &  $\epsilon$ . Say we fix  $\epsilon$  & its own

value for T. previous Coroll (or maybe  $\frac{1}{2}$  had this "mean value") [ $\frac{1}{2}$  = "mean code length"]

Does this model multiply the L-req time by N (length of code) as in  $\epsilon$  objection of 445.0-04? — Looks like yes!

For a given  $n$  when is  $(\frac{1}{2}(1-\epsilon))^n \cdot \epsilon$  max wrt  $\epsilon$  so  $(1-\epsilon)^n \cdot \epsilon$  is max:  $-n\epsilon + \ln \epsilon$  is max:  $-n + \frac{1}{\epsilon} = 0$ ;  $\epsilon = \frac{1}{n}$  as suspected!  $\frac{1}{2} \cdot \frac{1}{2}$  is so some max

$(1 - \frac{1}{n})^n \cdot \frac{1}{n} \approx \frac{1}{en}$  so is this worse than a factor of  $n$ ?

In T-2T L-req, Time for soln. is  $\geq 2T$ ; w.  $T = \frac{1}{pc}$  if we divide pc by  $en$ , it takes  $\geq n$  times as long for L-req.

So the system of .25R is at least  $e$  times worse than partial .23ff

The .23ff system takes  $n$  times as long in worst case situation. .25R takes  $ne$  times as long in (I think) all cases.

(Spec 446.14) dynamics associates

00: 446.40:  $F_i$  maps from binary strings to  $\sum_{0,1} F_i$  or  $\text{get } i \text{ more bit, stop}$ ;  $U = \text{undefined} = \text{nonstopping w/o. outputs also possibly.}$   
How to create a universal set of such functions is unclear. Well, just taking any universal from strings to strings, is about first bit! That will do!

Reading Leo-Cover 78 paper: Engstrom, prefix coding always has an asymptotic cost.  
The "factor of N" may not be  $\geq \log$  big deal at all! P 332 col 2 Para 2.

$$f(l_i) = l_i + 2 \lceil \log_2 l_i \rceil + \log_2 \left( \frac{2^2 - 1}{2^2 - 2} \right) \quad (z > 1).$$

$$\text{or } = l_i + 2 \lceil \log(l_i + 1) \rceil$$

$$\text{or } = l_i + L \log_2 l_i + (\log_2 l_i + \log_2 \log_2 l_i + \dots) + 4$$

the "factor of n" may have been like  $(\log_2 l_i)$  - rather small  $\frac{1}{2}$  cost for prefix coding.

I misunderstand this  $\rightarrow$  In fact  $f_i$  Leo-Cover paper shows that all codes & prefix codes have a best  $f_i$  name & they ( $\equiv$  "max ~~code~~ word length  $l_i$ ") for prefix codes, max codeword length is  $\sum \frac{p_i}{l_i}$

(.00-.02) would seem to almost solve the problem of implementing 3 input model.

We still need  $F_1$  &  $F_2$  - 2 functions - which definitely sounds expensive.

The 2 functions can both use the same "pool" of subfunctions. In fact,  $F_2$  (f. A "output") can use  $F_1$  (f. i bit output) as a subfunction.

Consider a function w.  $\left(\frac{3}{2}\right)$  kinds of output  $\textcircled{1}$  I want no more input bit  $\textcircled{2}$  A<sup>2</sup> output  $\textcircled{3}$  also "undefined".

A<sup>2</sup> output is binary seq. w. a "stop" indicator:  $\textcircled{1}$  stop can be an internal state, but externally Accessable!

Is request for another bit a "side effect"? If so, it could be possible to have function w.  $\left(\frac{3}{2}\right)$  output & request for bit  $\left\{ \begin{array}{l} \text{no request for bit} \\ \text{output} \end{array} \right\}$ .

Another way: If output starts w. 0, we request another bit; " " " " 1, rest of the output is output "A".

See 581.11 for Comments!

Here "S" is constructed as AZ141, so we have reasonable PC's for sub-functions

.24R SN we could construct

rather sophisticated functions, using previous sub-traces, by Monte Carlo Selection! It may be possible to "detect" in this way, rather obscure sub-traces that one would otherwise find difficult to detect in the corpus! It will be easy to do as an actual computer program, & return them as a "tree structure" (which is

Useful way I think about it. One dirty, one couldn't take advantage of fact that 2 MLCs generated functions were binary. The binary had different structure. How this occurs in non-Monte Carlo methods as well! If it is, it's a serious dirty in the stuck.

.22 would "work" but Heuristically, I don't see what's going on at all!

So .22 is a "soln" but makes it clear to me that this is not what I'm looking for!

Perhaps closer to what I want: say I have  $\left(\frac{3}{2}\right)$  of 0, 1, 2, 3 etc;

How do I deal w. these cases?

For  $N=0$ , I guess we would say: No output: that is  $S=0$  means i. Machine's very young. A more mature machine would always be able to categorize each new

input as being "similar" to one or more previous inputs. My idea of "recognition function" is a "Fuzzy out" i.e. stochastic. I did consider Param

Schaum R. functs - I don't know if I found them bad, or just want. NB: Any Fuzziness param for how this used to get PC's is max of  $f(1, 1.8)$  probabilities & only criterion

$N=1$ : This "set of Q" has occurred once before. So we want to maximize  $f(1, 1.8)$  w. 2 cases of Q: T. Q's may or may not be identical. T. Q's all identical

but we don't want A<sup>2</sup>'s is an interesting (simple) special case.

So maybe study this case, using (.22-.23) - say 2 or more identical Q's. If N is large, it is a kind of Bern sep. - but Pr. A<sup>2</sup>'s can be "Bimodal"

so to one another - so that could give me all denominator denominator from standard Bern Sep.

for .22-.23 - take a standard Bern sep.  $\geq$  default set of types

2 large N. (You would (.22-.23) in placement PC's??)

(Spec 448.07)

dynamics associates

00: 447.40 : [SN] I think the main definite progress I've made on the QATM problem is

1) 414.18 2) The idea of Recognition functions & assoc. operators:

At this point, I suspect Q is whether I want to use R<sub>i</sub> functions to be stuck, or not - a/c whether I want several R's to be able to work on a new Q - is whether I (modify) R's so that each Q has only 1 R; assoc w/ it.

05 3) Another Branch from Q is to use  $\epsilon P$  to try to ~~max~~ max 414.18. | see 422.20 - 31; 432.20 - 40

07: (447.40) spec. Say we had 3 output strings  $\alpha$ ,  $\beta$ ,  $\gamma$  0101, 0111, 1010 & would equal probs & say a total of 100 cases.

One way would be to have 3 distinct  $P_i^j$  func's  $P_1^1, P_1^2, P_1^3$ ; each had a different ~~output~~

One of the 3 outputs for 1 same Q.

09 At 07 implements a common Reggy type (i. Bern reggy). We could do it w/ a simple  $P_i^j$  function, but I don't know any heuristics to find it.

10 My method of (07-09) is not at all clear!

08 Obvly: We have these 3 ~~func's~~  $Q \rightarrow \alpha, Q \rightarrow \beta, Q \rightarrow \gamma$ .

The could work this way: for  $P^1, Q \rightarrow \alpha$  is correct: otherwise, we have to write "0011" or "1010".

So  $P^1 Q \rightarrow \alpha$  ~~saves~~ much react -  $\frac{1}{2}$  of time - still a lot of cost saved.   
 saved?

16 An easy heuristic to find suitable  $\alpha, \beta, \gamma$ , would be: if  $Q \rightarrow \alpha$  works enough time to ~~compress~~ compress code - priorities of outcomes - it sort of Hill climbing function for obtaining

08 good functional components.

20 So (16-18) gives a "statistical heuristic". We could make  $P_1$  part of TM's general routine -

a long w. a way to prod it (costs applies) into 414.18 for evaln.

Note that 414.18 does not stimulate (or require) a  $\exists$  in not stochastic function - Nor does it require of stipulate LS machinery

Another Common Stoch Heur for Probly problems: (related to 16-18)

For a certain set of Q, we find  $F_i(Q)$  is correct unusually many times: we then

try to find a way to characterize a fi. Set of Q's for which it works well, so we can

get it to work a by 90 of times.

Many heur for S-probs involves (reducing/expressing) to reggy (to) a Bern seq. i.e. t. "frequency data of  $P_i$ "

→ 449.20

perhaps write it up as it is now: explaining how it works, but also what diffs ~~are~~ are in train stoch. func's. So the entire field of devising heur for stoch. func's is something I haven't done enough - just I have to develop.

dynamics associates

Review: T. present differences about 438.00, in trying to apply Ls to finding of Solns for stochastic TSQ's. Up to that pt. I'd been thinking in terms of def. TSQ's (MTR). T. was how to apply Ls to A=141 model - ~ A=141 applied to ~ Ls.

Some Imp. ideas: This is An Implementation of 446.05

1) 447.22-23 This is a soln. of the problem of (0.00-02). It makes it clear that this is not a real problem. See 438.23-31 for impl. details. 447.24-40; 448.07-30 discusses the real problem; finding ways for stochastic induction's putting them in to the form of 447.22-23 (3 input A=141) or somehow showing how one hears 447.18

2) On the problem of finding good stock models of data: 440-21-40; 441,21-31, (447.24-40), 448.07-30. Actually, there is much earlier work on this; try to find refs

Also (Apparently) relevant: Methods of expressing PD's

3) 439.25 (as SN) Giving "index type" for various TSQ's can be a powerful kind of "Hint" to search time tremendously!

4) (438.23-31): How to get A=141 to make functions of 2 variables (Q & R) only. - by putting Q & R in a "constant" (P is sort of part of (449.04) in min. Argument

5) 448.00-05 SUMMARIZATION of main Imp. ideas in QAM This far

So, very little of 430.00-448.00 is important! - T. imp. stuff is mainly these 5 points: (1 & 2 in particular)

- Well, not quite! Some ideas that didn't work: Various attempts to do a Ls on prefix ~ prefix set "R" in jobs: 1) use of (1-epsilon)/2, (1-epsilon)/2, epsilon: This gave prefix set but was inefficient in coding! Also some things perhaps could be used. T. did 441.36-37; 2) Use of 11 Prefix functions for total O(N/Q). Very thx v. good! 442.34. I think there may be a better way to do this!

20: 448.29: Say I had somehow managed to code a Bern seq. Could I express it as a set of words (i.e. alphabet); the case counts of each; the resultant Laplace's rule eq. for the PC's. If I did this, could TM solve for Analogy w. other Bern. Seq. problems? How could the intermediate info: the case counts, be recognized as in part of the theory could be parts of the code for the PC's. (S seems to answer 6-Q).

How to code a Bern seq. in the 3 input case? Note: first + case counts are a set of params indep. of the ordering of data! All coding params of interest must have this property!

Given a subset PC's [P\_i] P\_i > 0, sum P\_i = 1; how to get a prefix code for them?



Well, the constant length prefix codes will usually give a v.g. approx! - but there is a code that gives binary expansion of each P\_i.

AH! Say P\_i is represented by codes of length [L\_i^k], k=0...w\_i

Then sum\_{i,k} 2^{-L\_i^k} = 1 so there must be a prefix code w. all these lengths.

Given a set of code lengths L\_i sum 2^{-L\_i} = 1; how to find the prefix code?

The process you put L\_i in order; shortest first: If we start w. v\_1 codes of length L\_1

we use the lexically first v\_1 codes of length L\_1 then we code v\_2 codes of length L\_2 - lexically next, etc.

Will this work? Well, it doesn't matter; an alg. exists almost certainly can be found (and proved) - see 37 for correct sign

Unless I happen to think of it: don't spend time on that! This is a correct way to take the first few strings of length L\_0

Fill out the remaining bits until the last length L\_i; take the next string; fill out the remaining bits until the last length L\_2; then fill next w\_3 strings of length L\_2; etc. etc. In the long. development, no string is a prefix of any other.

dynamics associates

00: 449.40: T. "Bernoulli function" : Input is bit of codewords ( $\equiv A_i$  values) & case count for  $A_i$ .  
 Output is  $P_i$ 's of each  $A_i$  or  $P_i$  rules for each  $A_i$ . Depends on how we want to implement L such.  
 This Bernoulli is very commonly used.

03: After 434.00 - 40 which describes L such & applies it to  $R_i$ 's discovery, I run into trouble  
 04: w. L such for  $P_i$  ( $A|Q$ ). - a stochastic function.  $\rightarrow$  451.00

05: Also, I had objection about "5 update modes" "5 rules" are 419.00 - 40, 429.00 - 00.  
 437.00 on some cond #3 In general, I have to go over these update rules to see if  
 437.11 on cond #1 They work for S. Func. (NMTM).

So: update #1 (Everything updates): 437.11 - .16 is v.g. while I can see how this would be applied to AZ141 inducting, I don't see how to apply it to 6's input machine. Using L such. function of (.00) may help here. - We actually may do it in AZ141 part.

Now, until I actually understand how to do Bern steps in QATM, this whole update process is not ~~problem~~ be very useful!  
 Update #2 seems to be except PC of A seems too small. Actually, I probably will want to merge conditions for #2. If #1 occurs, update isn't necessary, but for #2 it is very.  
 In both cases, I guess L such would be the same. Ideally, it would take less time for #1 because update is usually simple (unless new conc has been derived).

Update #3:  $R^i(Q_{n+1}) = 1$  correct! My update may work for NMTM but its "SSZ = 0" (or DSL) for NMTM  
 At this point, we don't know if it's NMTM or NMTM.

We may leave it alone until we get a lower SSZ.  
 Or, try to categorize it as  $\rightarrow$  to some ~~previous~~ problem(s) try various  $P_i$  funcs on it - see if any help at all.  
 See if any of  $R^i(Q_{n+1})$  is "close" to 1: This is meaningless as stated: But a "closeness" of some value/unit to be derived.  
 If it gives 0.8  $P(A|Q_{n+1})$  for a  $P_i$ , use that; if its a bit low try to modify that  $P_i$  so that its  $\approx P(A|Q)$  for

its recognized set, is adequate.  
 #4  $R^i(Q_{n+1}) = 1$  for  $\geq 1$  & 1 of them works o.k.  
 #5  $R^i(Q_{n+1}) = 1$  " " " & none of them ~~are~~ adequate. Actually, both #4 & #5 are heavy for "Merging"

So 2 imp't kinds of have: 1) To derive new  $R, P$ ; 2) to Merge old  $R, P$ .  
 Type 1 occurs when all old  $R, P$ 's seem to fail for a new QA. In some (maybe all?) such cases, we may want to wait until we have an adequate SSZ before trying for a new  $R, P$ .  
 Whenever a new  $R, P$  has been found, we may want to try it on older QA's. If it works we have (2) kind of merging.

**SNR**: Proof that (for discrete univ. p.d.) shortest code  $>$  PC (constant indep of what?)

30: Try to find proof in ① Martin ② Willis ③ Li-Vit. : See what order of Magnitude it is -  
 See how it com pares w. length of dem. of "True Generalized P.d".  
 $M(p) = \sum a_i$  ;  $L(p) = \sum k_i a_i$  = length shortest code in which machine stops on  $a_i$  on output.  
 $PC(a) = \sum 2^{-L(a)}$  ;  $P_i$  = code for which machine stops on  $a_i$  on output.  
 $2^{-L(a)} > B \cdot PC(a)$  But const indep of  $a$ ? but dependent on M. What is B indep of??

T. proofs in Li-Vit, Cover-Lapug, Cover-Joy, and all the other log & involved!  
 Willis has a proof, but not clear for  $CD < \infty$ .  
 I may try Willis' idea & run limit to  $\infty$ . Say for M, for  $CB = T$ , we get  $P_{M,T}(X)$  - This is useful for  $|X| <$  a certain amt. (Proof w/T). For each T, there will be a set of X that are assigned PC's, for certain values of |X|, all X are assigned PC's for a given T; - also, all shorter X will be assigned PC's. (Possible on h longest X's)  
 Re: Willis' proof, it is true for all CB's, its probably true for  $CB = \infty$ ! For practical purposes, CB is finite anyway - so so  
 Spec.  $\rightarrow$  454.00



TD

dynamics associates

On a form of S. functions used by t. system: I chose t. z input form because I know it was universal: z I found a way (447.2625) to implement it. Hrs, for (Baru Seqs, i probably many other P.D.'s (since many/most) are derivative from (Seqs), an output int. form of a PC, ~~is~~ digital or analog, (or perhaps M code) form is most desirable.

Could we allow TM to choose t. form that seems best, for a given QA?

Probably before trying to answer this, I should work out more details about how TM does

Baru Seqs. and their params (450.00 is a start).

One simple TSAQ might be for a fixed Q: we have several  $A_i^Q$  i TM gets their case counts

→ pc's via Lap's rule.

One poss. "Next" would be: Same set of  $A_i^Q$ , but  $Q_i$  varies somewhat. Try to find relationship

betw. Frequ of  $A_i^Q$  &  $Q_i$ .

Th. third, most complex form: Both  $Q_i$  & t.  $A_i^Q$  set vary w. i i TM has to find how t.  $A_i^Q$  set varies w.  $Q_i$  as well as how to assoc. pc's off t.  $A_i^Q$ : This is a t. General QA problem.

Another variation - say SM or real no. prodn.: t. set of  $A_i^Q$  is the same: indep of  $Q_i$

But t.  $Q_i$  could be vectors. (t. curve fitting problem) - This curve fitting problem has a standard ALP soln. Find function of a vector that has best pc mult by pc of data in view of that function.

We can do it by having first choose center of d.f., give  $\sigma^2$ ; or have it give an arbitrary d.f. for each or all of t. datapts. So one does an Levenberg functional form, then various methods to find optimum set of params.

Can we modify .14-.20 so it will work w. diffrt. QA's? (ie not just real vectors).

If we have  $A_i^Q$  be a fixed set of vector strings, their pc's constitute a vector, so we just have vector curfit - which is a trivial extension of .14-.20. (T. output vector's for any data set are then associated via 4.14.18)

One interesting General of .20-.23 - First t. set of  $A_i^Q$  is an infinite set of strings (or integers).

If t.  $A_i^Q$  are integers, we write record form as an input to a function: so we want  $F(Q, A)$

But in general, any  $F(Q, A)$  would be OK. - just as long as they are normalizable.  $Q$  can be a vector,  $A$  can be an integer,  $Q$  &  $A$  can be strings of 0's, vectors, sets of strings, sets of strings, etc.

One big problem w. .24-.27 is Inverting  $F(Q, A)$  to find A's of by PC. - This can be as difficult as any General O2 problem.

Another Dirty is Normzn: Depending on t. form of  $F(Q, A)$  this could be very difficult.

Could we use a Mt. Carlo method? Not so easy, since we have no idea as to where A's of large PC's are!

That's one advantage of t. 3 input form: The t. pc's are not normalizable: They are usually not far from normal, & we expect that for large S's, they would be close to normal.

Well, except for normzn, .24-.27 might be easy to Update (since we know  $A_i^Q$ ) - but harder to use for prodn (No perhaps we wouldn't spend so much a large fraction of time on "production" prodn. Perhaps use O2 (such as a GP to get PC's of large pc. to answer Q's.

It would be neat if we found a nice way to reduce the search space over only dunct that were normz'd, (or almost Normz'd) (or easily Normz'd)

ID.....

20xKv B = 20x10<sup>12</sup> by.

500 x 100 = 50kw.

dynamics associates

Spec  
00: 451.40 : It might be conceivable if  $A^t$  were constrained to be ~~reals~~ reals, or integers - but if they are to be Arby Strings: This may not be so easy! The 3 input Umc begins to look good!

The strings  $\leftrightarrow$  integers. The int. case of strings, there is no a priori matrix - well there is, but not so easy to use in this case (?). ~~is~~

Perhaps it would be best to try a MTM TSG: They look at t hours for t seq.  $\rightarrow$  Try to find way in which TM could discover them.  $\leftarrow$  CB may be fairly large - It took much human effort to discover many of them - ~~is~~ Sometimes Centuries!

08 Re: Bern Seq. is similar to ~~the~~ pd's Part I have hours for: I can express them in any form I like - but ultimately, it has to end up in, say, 3 input umc form. - So t.  
3 input umc is should have any form to do w. hours for NMTM - it is just a way to help

10 evaluate 4/14/18.  $\rightarrow$  Note (31) 2150 Note ~~453.28~~ 453.28  
11 If .08-11 is true - That would solve it! - How to implement hours in 3 input umc form. 3iu  
12 Def Heuristic should be doable many universal lang. So try t Bern seq. (450.01) - since it seems to be

T. expression of t should be doable many universal lang. So try t Bern seq. (450.01) - since it seems to be  
So important in probabilistic Analysis. Consider a fixed Q. we have various empirical  $A^t$ 's for it;  
Certain  $A^t$ 's have occurred several times! So Lap's rule is a good ~~method~~ way to  
way to code it. T. Bern Seq. might be pervasive enough so one could get pretty far in NMTM  
Enough to discover  $\rightarrow$  fair no. of hours. Using t Bern Seq. as one's only NMTM model. 453.28

18 Related perhaps to Bern Seq;  $A^t$ 's do ~~not~~ repeat much (even for Q is constant), but somehow, we  
19 find a numerical function of  $A^t$  that encodes its frequency of occurrence. This could be usefully done  
20 Even if there is little or no repetition of  $A^t$ 's.

A rather direct way to do this: try to find a cheap way to code all of  $A^t$  that have occurred.  
This would be a "stochastic grammar" - perhaps of the usual prob stochastic production rule type.  
Discovery of such grammars: well I have some ideas on how to do this: It's quite hard well  
defined problem if we restrict ourselves to certain Grammar types. for t General (universal)  
problem of inducing ~~a~~ a grammar: I don't know how to start. One might look for short codes  
to individual  $A^t$ 's, (then short codes for pairs. Then for random triplets:  
Or look at short codes for individuals: look at pairs of the codes that are "n": then  
code  $m$ . These pairs - do seem similar  $\rightarrow$  giving w. n-tuples of  $A^t$ 's.  $n = 2, 3, 4, \dots$

If  $SSZ$  is large enough; for each  $Q_i$ , devise a grammar  $G_i$ ; Then try to find some  
 $G_i$  is a function of  $Q_i \rightarrow$  plus, perhaps (some bits)  $\rightarrow$   $G_i$  is a stochastic point of  $Q_i$ .

30 (N21/01) ~~improvement~~ improvement Re: (.08-11), it should be poss. to take any hour, expressed in any  
31: (11) formalism ( $\approx$  "Language"), & use it to optimize 4/14/18 via 3iu (3 input umc)  
Well: Give an Example; Say t is 1d1 method of coding a sequence of symbols. - This is a heuristic  
for Maximizing 4/14/18. Or even simpler: we have a Bern Seq. & Lap's rule is simple  
 $\approx$  Heur Method to max 4/14/18. (Note 453.28)

36 That I have to put a Bern Seq Consider not a Bern Seq. but an unordered set w. different frequencies for  
each set member. This problem is easier than the Bern sequential Bern Seqs because in  
Seq. predn, other models are poss. If data is known to be unordered, then assigning  $p_i$ 's to each symbol  
is the best that can be done. - So can one get Lap's rule? (or even "straight rule")

dynamics associates

lower nt. = 28 (Laurids)  
S. out = 29 (Ac)

00:452.40 : Perhaps first try to have a clear idea (perhaps formalize), of just what problem of 452.08-11 is  
01 - or just 4. problem = Barn problem of 452.36

02: 452.18 : On the usefulness of the Barn Seq. in Statistics: Situations in which the Barn Seq. is not applicable:  
1) Curve fitting  
2) linear/n.l. regression (of time series).  
3) Stoch. Grammars.

Possible applications: Linguistic Analysis of some kinds (or Barn analysis?)

09 2) Cross Correlations.

10 Well just draw up a "Barn" problem for ATM's see how we can solve it.  
11 So: fix Q's, various A's for "long sequence" T. Q's are "indexed" but identical.  
13 In General, for fixed Q is language & small no. of AR types; & Lap's rule is the first (only?) thing to try.

15 We could do a Barn Seq as a sequence of objects - as a single Q: i.e. after TM has learned to understand  
16 What the problem is - but I think this is a more complex problem than (11-13)

15 So: say we have 4 types of AR:  $\alpha, \beta, \gamma, \delta$  with case count of 2, 4, 5, 9. - met's corpus  
16 Q is same for all 22 cases.

17 I recognize it as a "Barn" situation & immediately involve Lap's rule - which gives me frequencies & Barn codes (via 449.32-40)  
17 is all done "open loop" on basis of past experience. I first "recognize" the corpus as a Barn-type problem.  
After 17 has been done, TM can verify that the corpus has been compressed (via 414.18).

19 Woops! If I just maximize 414.18, wouldn't I just get simple max likelihood - i.e. peak at  
20 which is "T. straight rule"? Well, if we sum over many poss.  $\bar{p}$  values ... ? Lap's rule is obtained by

Integrating over all poss. sets of  $\bar{p}$  values.

22 T. correct way to solve the problem via Alp: for each next prod: 4 hypotheses ( $\alpha, \beta, \gamma, \delta$ ) - each results in values  
22 corpus with issue.  $\leq$  cost. What we want is normalized of those costs.  
Each of those  $\leq$  costs is  $S$  over 3 space of possible  $\bar{p}$  values, of  $\cdot 19R_1$  (T. Integration as well as normalizing involve various "normal factors".)

The Method of 22 ff would seem to do OSL - but usually not in more complex cases - since we'd have to get  $\leq$  corpus for all poss. AR containing - & there are many poss. contents.

27 Actually, I should be able to use any ("long"/method of darn) of a hour that I like:  
28 This has to be true! - Because hours are very heterogeneous - very much given to being discreetly disassemblable

30 w. a suitably modified "Language". I will have to learn to xfer various hours into "standard form" -  
i.e. maxn of 414.18. As for TM discovering such hours, first TM has to (discover/learn) the "Languages" that the hours are expressed in. (By "language" I mean language w/ "point of view"). So I should be able to write up the hour in whatever language I like, - Then for the corpus of that hour (raw data that gave rise to it - Post justly) I can note that 414.18 is increased.

(Hr., some hours (maybe most!) don't have a corpus: they are discovered by "Medical Analysis".)

En this case, one should be able to find conditions under which the hour is expected to be true. - But for such hours, I'd really have to look at one to see ideas on how to deal w. it - on just how (if at all), maxn of 414.18 is relevant. - E.g. T. Hour could be found at cc of such. (So: some 414.18, but less cc!) - More Generally, (Spec 450.00 456.00)

How fast is using PC vs.  $2^{-kcost}$ ? .04

dynamics associates

00: (Spec 450.40) Consider Arit codes: For finite CB, all x's of < certain length have pc's, so one can make Arit Codes for them. In this case, T. code lengths usually  $\lceil \ln_2 pc \rceil + kcost$  of Arit Coding Alg.

However th. Bound about  $\lceil \ln_2 pc \rceil$  is within constant of k cost for any x. (We find distance is mid of x)

But for "shortest code" could be  $\gg cc$  for "many" codes.

.04 So: interesting problem is: For "small" finite CB; Towhitekenti's  $2^{-k(x)} \approx pc(x)$  approximately true  
i.e. in a search for codes, or  $\leq pc$ ; we usually only code maximal corpi of small kcost; variation  
at least t. codes we can find in  $\leq CB$  are not large  $\approx kcost$  even shorter. For codes of (limited k cost)

Can  $\lceil \ln_2 pc \rceil / \lg_2(\text{shortest code bound})$  be rather small? One application is in  $\leq (corpi)$  for 2 corpi. viz 57673

If true pc of corpi is  $c$  is t. best we can induce  $c'$ , expected error  $< \ln(\frac{c}{c'}) = \ln c - \ln c'$

If we shorten  $\ln c$  by  $\Delta$ ; t. fractional  $\downarrow$  in  $\leq (corpi)$  will be  $\frac{\Delta}{\ln c - \ln c'}$ ; This will also be fractional

↑ in "convergence rate" w. SSZ. — it will be  $\frac{\Delta}{\ln c - \ln c'}$  fractional  $\downarrow$  in SSZ we need to get same mean sq. error

For problems in which SSZ is a large fraction of project cost (like "Clinical Trials") — This factor can be very imp.

(vs. for most AI problem  $\alpha$  is .1. it does not pay much, but if  $\alpha = .5$ , this can be very imp. speedup.

In many problems,  $\ln c - \ln c'$  will be t. In pc of t. true model w/ t. reference one. It is you would be error for  $c$  based on both  $c$  &  $c'$ . (its always  $c$  base for  $c'$ );  $c$  need not be base for  $c'$ . So  $c'$  will be at least t. In pc of coding model) + t. deterioration  $c'$  due to not all codes being found — in fact t. shortest code may not have been found — so  $\leq \ln c'$  could be  $\gg cc$  ( $\ln c'$  for  $c = \text{opt}$ )

I'm beginning to get a. impression that using  $2^{-kcost}$  rather than  $pcost$  may not make much difference —

Except for Heur search w. backtracking a/o search in which one does recursion betw. "reasonably good models" We want to return models that are anywhere useful, yet are significantly

different, structurally, than "Best Model" —  
Hvr, I really have to go over this more carefully. T. Goals for Di. How kinds of (Statistical) tasks can vary a lot.

- ① Small max error is but one goal. ② Cost of sample is another ③ In (QA) TM I'll be concerned w. "allowance rate learning" e.g. how much  $cc$  to get particular "Milestone" (like understanding NL of upto 20%). ④ Pick in SM.

Also I'd like to be able to characterize 18-20 better: What I want from t. "not so good models" that are much different from "Best" model, is: What are t. differences,  $\approx$ ? They amount to an alternative method of coding ~~to~~ parts of t. corpi (say). That they would at all is of interest: we should try to modify them a/o use them as components when t. present "Best" model no longer works

30 Actually, T. "importance" depends Much on t. Application: In SM, say we use t.  $2^{-kcost}$  code "The part that describes t. model can be in MDL or in ALP. The second part is always effectively in ALP. And correct The use of ALP for t. Model can give a much more optimistic prodn. of future yield

**dynamics associates**

Note 457.02 ~~03~~

In General, the problem of how to Generate S-functs is not a new one! In my attempt to analyse GA from a ALP point of View, I ran into difficulty of not having good ways to do this. - So maybe list various "specific problems" of this sort: see if ~~the~~ patterns emerge

- 1) Old Super SGA idea of PC → Fitness; so used a CFG to get strings → PC over, then a R/R finite monotonic function from PC to Fitness. R/R, what I think I needed ~~was~~ was a way to go from strings (cards) to a d.f on Fitness. (≡ Genes).
- 2) General forms of d.f's for context, SM. so this is comment 1) ↑
- 3) A common D.F. is "Central" pt. w. pc's of various derivations: "Clustering" is a form of this
- 4) T. very general 3:U form: strings → % non → p.d.
- 5) (Sardot) related to 1) (≡ 03) Used a % stock CFG to assign pc's to strings (a "language")

for "Concl prob" we want a large param to be a function of Q. (The discrete as well as continuous params must be funcs of Q)

As an example, perhaps: Given a corpus of unordered strings. We use a standard simple Grammar Gen, to stochastically generate deterministic strings; Then parse the corpus with a Grammar (if poss.) & assign pc's to the choices in the parsings. This is a "universal" d.f. overall

**SN** Ont. "Recognition" problem: we can have category (cats) & subcat: Geom figure, 4 sided, 3 sided, 2 equal sides;  $1 < t < 2$

Also, useful to have problem in  $> 1$  cat (often but not always). Some cats will help in some situations but not in others. Finding out when this is so, can reduce search time. Multiple R's correspond to different ways to look at a problem

e.g. Problem is both in  $R_1$  &  $R_2$ .  $R_1$  errors a good, useful P (compression)  $R_2$  does not (in P's case). Try to find out "Why": ~~How cats~~ 1 or more new (or old?) cats that will help tell if  $R_1$  or  $R_2$  is the best way.

So we usually use an R funct (cat) to tell which methods in  $(0, 1, 10)^3$  should be tried/od.

This is a very imp. Heur.

**N24/01** It would seem Best to have  $> 1$  R poss. for a given problem, initially, but that we normally try to find ways to narrow down to 1 R funct so it points to a "Best" soln. for Q. Perhaps this has a disadvantage of

"MDL": That we don't return low PC Solns, ~~but~~ that we may need (w. undifs) later - ("Backtracking")

**N24/01** I did use  $> 1$  P's function per ~~set~~  $Q_i$ . I may want to give extra wt. to those P's that were most successful for a set of  $Q_i$  that they were designed for.

**N28/01** I will have to expand this idea: it seems very imp! It amounts to a gross over-haul of my initial QATM model, since using only 1 R for Q type was an imp part of the System!

One way to deal w. this is to have "Soft R's" Instead of yes/no, have a Fuzzy categorization:

I have written about this: The "Fuzziness" param need not be probab: It can be assigned in any manner:

Heur, the main Q is how it's used to guide search for soln, & how effective it is in  $\uparrow \alpha$  (≡ 4.14.18)

447.28 → discusses Fuzzy R's ~~but~~ - why they would probab be useful (also a vague reference to previous work by me on this)

An interesting property of "Fuzzy" params - it may be that they don't have to be normalized!

**N28/01** This fuzzy idea seems closely related to an earlier TM model in which I had a "Front End" function that would look at the problem & assign pc's to various IA's (Inductn Algms) - These ~~are~~ my initial estimates of "how effective" those IA's would be for that problem.

In the earlier version (C.28), I had pc's assigned to IA's: Now, I'd have pc's assigned to various P's or sets of P's. An apparent difficulty: By having only 1 R assigned to a Q, we know what P's will be. If we have several wtd R's assigned to a Q, the PC of the A would be mix of P's of the P's invoked by those R's. How do we evaluate  $\alpha$  (≡ 4.14.18) in such a system? - Well, it would work out:

the system would assign pc's to all of the IA's in the corpus. I'd have to compare effective pc's of the O's functions of  $\alpha$  (≡ 4.14.18) - I don't know how hard that is - is the apri pc's of the O's worth the effort to compare  $\left\{ \begin{matrix} 413.00 - 40 \\ 429.00 - 0 \end{matrix} \right\} > 150.05$

As for "Updating" it will be done a little like for the previous 5 cases. - We do have "R parts" and P's parts of O's. Another Ring Rat Model: "1 R per Q" attractive was that when updating, I didn't have to reevaluate the O's for any but a few QA's values. To what extent is it true if R's are "fuzzy"? → (spec 460.10)

dynamics associates

00: (spec 453.40): 2 hour can be oriented toward cc, 414.18, or a Goro. That's some combination of these 2 factors.

Usably, the hour will do one or the other (cc or 414.18) so it's easy to tell an improvement has been made. If, e.g. cc is much, but 414.18 is a constraint, it may not be clear as to whether the hour is any good - probably it depends on application. In general ALP has cc=0, so any hour will do. It's also to PC (= 414.18). In the most general case we have a Goro that is a function of cc & 414.18 - & this function will depend on the application. A rough dirty Goro is  $\frac{PC}{CC}$ : It just so happens that Lash makes sense in that order.

Ab initio, I expect that (hours/methodology) for learning S. problems will be "wired in" - that this will be true for quite a long distance along the T SQ. The changes of "Language" involved in many (most?) hours, seem quite complex & will be beyond the capacity of ~~anyone~~ an "Immature" TM!

So, e.g. for Bern Seq-like problems: I will initially "program in" a Recognizer and a P<sub>1</sub> function for them. I'd like TM to be able to generalize on this "wired in" stuff, but (perhaps), the more "open" the code of the hour is to TM, the more easier it is for it to generalize - i.e. the code for P<sub>1</sub> & P<sub>2</sub> should be as simple as possible.

Using functions that TM has found useful - that it understands the properties of.

15: 453.30 **SN** This "Being able to use hours many forms" ("language") is a key property of any form of induction that I use. - So it is not a particular dirty of 32U. Any adequate induction system has to be "universal enough" to at least be able to do that "language translation" - which is a given to universality and the "invariance theorem" becomes relevant: so PC only depends on X (the cost of Language).

So it seems clear, that the problem of expressing an ability hour in terms of various induction systems is a problem that must be solved - no matter what form the induction system takes.

Furthermore any induction system will have to learn the various "languages" used by the various hours. It may well be that certain induction systems are more "Central" to the set of hours - so on the average, it is easier for such systems to learn many hours. If so, it would be well for any ind. system to be able to learn those particular "Central" induction systems.

So 453.11 ff on T2 Bern Seq. prob. is an attempt to map this methodology (which I've been calling a "hour") into  $\frac{PC}{CC}$  Maxzn. (well, it is a "hour" - it is a fast way to get large  $\frac{PC}{CC}$  at low cc.

**NB** A very easy way to get TM to acquire new (ways/ways of learning) is to have TM learn definitions of relevant concs. At first definitions would be hard by some suitable T SQ. Later, when TM has acquired understanding of the suitable language, TM can be "told" details of certain concs. (as a dictionary).

126/01 Hour, it may be better to have TM learn terms in the "usual", purely inductive way - that PL's write be more beneficial for TM: it would have to put the things defined into the lang. of concs that it has already (evaluated/understood). W. Human Students, this seems to be the better way.

Anyway, I think the way we do .15 ff! Just write out how I do the hour, (or solve a problem) in English: Translate it to a more exact form. Then see what TM has to know, in terms of (defns, hours, concs) in order to be able to express that soln. in acceptably small code - (i.e. perhaps be able to solve that problem w. acceptably small CC).

# Boston Bicentennial Associates, Inc.

ABC abc

ABC ABC

.00: 456.40 → Learning to solve (even) Tech problems, can be understood/analysed via a conc. net - just **GOOD!**  
.01 like d-problems.  
.00-.01 looks Great! So just start at 456.36 : → ~~the~~ writer: soln. in English, then analysed into  
concs, & subconcs & make a conc net that has to solve at x top.

Probably the problem of devising suitable ~~Rec~~ Recognition algos (t. R<sup>2</sup>) can also be formulated as a  
conc net. — ie At first, I want TM to be able to do Recogn. about the Some way that I do

So I'll have to teach it via a conc. net. **.13** .13

**.07** **SN** Keep in mind that ANY problem can be solved via Lsrch, & if we don't write soln will be written  
& factor of maybe 2 of the best that could be done. One way to have all have int. pd is one  
has to Iron down the Lsrch. This is related to, 456.36 - 457.06, but I don't think its the same thing.

What **.07** says is that Lsrch can be close to optimum.  
A poss. BUG in **.07** is that the system that I analysed had no "Look ahead" (ie Dynamic Prog).  
Maybe double w. "second order TM" means given "look ahead" as a "Regular Problem".

**.13** **OS** I could just ~~write~~ "write" some very imp. elementary concs involving probly.  
Have TM constructs more complex probly concs on that basis. Perhaps get those concs from  
Wm Feller's book. (Maybe try to get second edition?)

→ What I want now is some examples of TM solving s. problems. Start w. Bern seq.?  
Or perhaps use Bern seq. as used to "Primitives" in **.06**? If so, how to express it so it could  
be used to solve probs (2) be combined w. other concs to yield concs that can solve more complex/diff. probs.  
So just derb. how to do it in Eng., then formalize, then link to "α" (≡ 4.1.13)

for statement of problem: say 453.15-16 : soln: we Recognize Bern as a form problem.

T. characteristics: A small no. of A<sup>2</sup> types, perhaps a fair no. of A<sup>2</sup>, w. repetition: to larger t-SS, the  
better "Bern seq." fits. In the present setting, if Q<sub>0</sub> is constant, any set of A<sup>2</sup> would give  
a Bern seq. because t. A<sup>2</sup> are "by definition" "indip". If all t. A<sup>2</sup> are diffent, the Bern model doesn't  
give us compression - but it doesn't lose any char (c) - w. g. n = pc to all cases, so t. pc is surely  
t. pc of conv decaying to A<sup>2</sup>'s.

complete For a Bern decn, we need decns of all of A<sup>2</sup> : Dist. case counts: Test it! assoc. w. that decn.  
Since t. decns of A<sup>2</sup> are self-dual, we can list decns followed by t. case count integer. T. pc's of integers is under!

**N27/01** The I could use that approach to 2-18 x. T. / pc's of integers is defined by  $p(n) = \sum_{d|n} 2^{-2(S(d;n))}$   
My impression is for decn was not very bad.  
From this decn TM is able to make a code for t. set of A<sup>2</sup> - its pc is easy to compute.  
It may poss. to just leave the pc (from Laps rule or SS rule) as is, w. a converting it to code:  
4.1.18 can use these pc's just as well as "code for 3iu". - easier, in fact.

**SN** Riss' 2-18 x is for a continuous x d.f. : If we knew that x is an integer, my method should be  
give hyper pc's for integers - but I may have tried this & found it not so! - check into this sometimes!

→ A concept TM could find useful input for t. Bern seq. coder (pc evolution) :  
Given a set of objects, S, & one of t. objects A; Count(A, S) tells how many times A occurred. Big  
t. Seq could be in t. form of a list.

# Boston Bicentennial Associates, Inc. <sup>ABCD</sup>

00: 457.40! So TM looks at Corpus of Q/A pairs: First, it finds several different Q types: for each Q<sub>i</sub> type, it has a set of A<sub>i</sub><sup>l</sup>. It gets v. cardinality of each such A<sub>i</sub><sup>l</sup>. It plugs this into the Bern function, (450.00) and we get as output, the function that maps Q<sub>i</sub>'s to d.d.'s on the A<sub>i</sub><sup>l</sup>.

03 There are several operators needed to do 00-03, but 00-03 is double this way, so the subfunctions need look useful in their own right (i.e. they should be useful in many other kinds of problems) we

07 The operator of 00 is best for Q's that have > 1 A<sup>l</sup>: Otherwise, pool together all Q's that have only 1 A<sup>l</sup>, so we try to find a min code (grammar) for this set. .20

10 Another kind of S-funct. that is ≠ Bern seq, is Curlit. In best known form, we are given (X<sub>i</sub>, Y<sub>i</sub>) pairs: (essentially, Q/A's). The common way is to find f(X<sub>i</sub>) = Y<sub>i</sub> + ε<sub>i</sub> with min ∑ ε<sub>i</sub><sup>2</sup>, say. This error d.f. is like  $e^{-\frac{c^2}{2\sigma^2}}$  (which is a Bern seq.) ~~No!~~ Sometimes occur in Analy

T. analog for steps is: If X<sub>i</sub> & Y<sub>i</sub> are strings, find f(X<sub>i</sub>) → f(X<sub>i</sub>) is close to Y<sub>i</sub> using some metric we want, say matrix, is min. A common Matrix: How much extra code (many bits) does need to go from f(X<sub>i</sub>) to Y<sub>i</sub>? This is a kind of CPC (could probly) but its less diff't than f general QA problem — because we make use of same matrix (CPC) for a fairly large variety of X<sub>i</sub>, Y<sub>i</sub> pair types.

In general, I want to see how sols various kinds of s. probns. can be expressed, using certain "primitive" funcs.

19 → 08 in 00-08 This doesn't always work! Say we have Gaussian noise, p.c. of ε is  $e^{-\frac{\epsilon^2}{2\sigma^2}}$ . Say we saw a bunch of these ε's: Most would occur only once, & a few would be repeated (depending on our "resolution"). So we would use T. Bern seq, on a ~~subset~~ subset of Rest hold > 1 case of each ε. This is a very poor way to analyse r. data → 459.10

23 → SN (Set back to 20) But first note that for any S. problem, all I need to do is write out a soln. in terms of "t. primitives": Do this for many problems then decide on what good set of "primitives" can be. — Don't make conc. notes for t-solns for ~~probns~~ problems. "Soln" usually means "acceptably by d" but it may mean "acceptably by pc for Anti" — I have to straighten this out: I ~~like~~ written about this dirty — I think by pc for Anti had been used in some cases as a heuristic. In general, I need criteria for when to stop. Such. (Other than running out of cc).

30 One "prim op" could be finding f xfm (o/o Laplace xfm) — using simple rules or by using Maple or Mathematica or Macsyma or smaller/larger table of xfms. Similarly integration can be done — either indefinite or definite. TM can "learn" when to use external prim or do it itself on the basis of experience — so as to minimize cc. — Just >> a person does.

T. primitives may T. Stock "primitives" may be derivable (probably) from d. primitives. If so we should try to get TM to learn from more naturally than rather than ~~to~~ learn from "hard" word in. — Since TM may learn to do this for "Sub-primitives" to derive better solns to probns than I know of. — >> so it Spec 459.00



I.D.

small

REV. 36

HAMPTON RESEARCH ASSOCIATES

26 Boylston St., Cambridge, Mass. 02138

Spec

00: 458.40 : would be a better position to "understand" & solve for it. It's probably done -

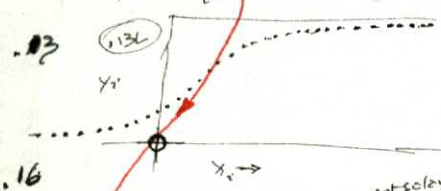
"understand" or deeper level

Go back to 458.20 / 1.10

For the most part, d. probn. & s. probn. will be separate & will be separately maximized. But occasionally I may need to mix them (to get their product max) or for a particular problem, it will not be clear as to whether it's d. probn. or s. probn., so I want to try both & compare them. (or - actually sum them). So I will have to (occasionally) convert from one form to the other.

In some problems (like Curfitting) & selection of functional form must have ~~some~~ param lengths & the ~~selection of~~ optimizer of the params of the form is ~~for~~ continuous pc must be (some param lengths), so we have to mix 2 methods, - or translate betw. 2 kinds of "Language".  
i.e. pc v.s.  $\sum z^{-2}(s)$

10: 458.53 For the Gaussian of type of problem, there is a common approach: we collect data as  $x_i, y_i$  pairs:  $y_i$  is the total no. of cases in which we've had a voltage  $\leq x_i$ . - so we get a jagged curve that (usually) gets less jagged as  $ssz$  grows. We can then try to fit a ~~curve~~ function thru these  $x_i, y_i$  pts.



Say  $x_i$  is quantized so it has a certain (uniform) grain size  $\delta$ . Some can associate each  $(x_i, y_i + \delta)$  interval ( $x_i = z_i \delta$ ) with the "pc" assoc w. it of ~~the function~~  $f(x_i)$ ; i.e.  $f(x_i + \delta) = f(x_i)$ .

So that's a not bad way to solve the 458.20 problem: But if  $ssz$  is small, it there is little restriction, 10-16 still is a reasonable way to solve the problem - unless  $ssz$  is so small, one doesn't have enough pc available to describe function, one must just assume "pc" for all  $x$  that have occurred (is zero for non-occurred cases??)  $f(x) \leftarrow$  (Normalize)

19: Actually, a better way to do 10: Try to find a d.f. (using some grain size,  $\delta$ ) such that the (p.c. of the corpus) \* (pc of function form) is max. If the funct. has continuous params, we divide (20L) by the deformation of the Gaussian. If the def. of loss is zero, we have wrong dimensionality in our model.

So go back to 458.20's criticism of (458.00-08)  $A_i^2$  are all strings. Usually also if the  $A_i^2$  are a small no. of integers.

Or, say the  $A_i^2$  are finite strings, and there is some function mapping them to integers (some  $A_i^2$  could map to same integer). We then make pc's of strings or the assoc. integers & we normalize them. If they are not normalizable, the function is rejected.

Better: Have a funct map strings onto (0,1) or any finite interval, & normalize (if possible). A good thing about the Method of 10-19: The set of pts of 10-19 suggests a monotone  $\uparrow$  funct. We can usually get some kind of fit to it using relatively standard methods.

Well 458.00-19 (Bern + Curfit) is perhaps a good one for many problems: In general, in a TSD, I will have 4 or 5 reasons to choose various kinds of P.D.'s for the cauds. - TM will have to have these reasons. So I should be able to write TSD's for s. probn. - if I know how to solve these s. probn. myself. The Example of 458.00-06 (Bar seq.) is probably o.k.

I may want to describe another s. probn. soln., to make sure I don't forget the general idea! For Review 453.27 starts with present view of how to do s-problems. Then  $\rightarrow$  (end of p) 458.00-19 is main ideas, but there are other parts that are a bit's sequence. Also 457.00-06, 13-15 is v.g. 455.00-10 has some probs - some of which I've not done well! Use of  $\geq 1R$  to categorize "understand" a Q (problems) (455.11-10). Seems like an Empty Heuristic Principle! Amounts to "Different ways to look at the problem's".

I.D.

GA.01

A B c d e A b c d

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

relatively recent  
that section discussing use of  $> 1 R$  for  $Q$  was good, imp. — find it a reference in "Mini Review"  
455.11 ff

01  
SN On GA! (SGA) : For cards w. all the same Gorc (= "Fitness Func. Value") I had every of making a grammar from any 2 ssz of 2, say: Make a grammar that could have generated both cards:  
But also use OSL to give somewhat more pc to New Cards that use "parts" of original pair of cards.

Now → If the 2 original cards don't have same Gorc value: Define a Parameter that gives less pc to "parts" of low Gorc card. — Reasoning of pc param is based on difference betw. 2 Gores.

~~way it affects~~ T. amt. by which it lowers "pc of parts" must be periodically updated by  
~~statistical studies of recent empirical results.~~

10: 455.10 to 455.11 (Fuzzy RB: Continued). OK. So lets Go thru F. Fuzzy System more carefully:

A  $Q_i$  comes in, T. R function looks at it. Assigns its soln probs via  $P_j$  — so for all  $P_j$ , we have 2 pc assigned;  $a_j$  assigned to  $P_j$ . ( $P_j$  may be set of probs, or only one). T. System then assigns  $\sum_j a_j P_j(A_i/Q_i)$  to answer  $A_i$ . Now lets look at 5 update cards (419.00, 429.00, 430.05)

① T. pc assigned to  $A_i Q_i$  is adequate. Each  $P_j$  involved has had its ssz  $\uparrow$  by 1 & has to be updated. Normally, this will involve only small modifica- of continuous params. Here, if ssz is over threshold — certain new conss will be defined & effective  $P_j(1)$  may change more seriously (than if only continuous params were slightly changed)

So: Relatively little updating: Sometimes one will not bother to do updating in such a situation!  
Also note the R function has to be updated also because of  $\uparrow$  of ssz by 1

② T. pc assigned to  $A_i Q_i$  is too small. (But is, presumably,  $> 0$ ). We will probably try to improve this  $P_j$  that were assigned large  $a_j$  values for this  $Q_i$ . a/o Look at other  $P_j$ 's: first of ones assigned lower  $a_j$  than to

\* best  $a_j$ : Then later, look at  $P_j$ 's that were assigned  $a_j = 0$  (perhaps), if a good  $P_j$  is found for this problem, the  $a_j$  assignment function has to be modified so this  $a_j$  is larger. (Sometimes larger than its previous value of zero!)

③  $Q_i$  is given  $a_j = 0$  for all  $P_j$ . We have to try all  $P_j$ 's on  $Q_i$ .

④ and ⑤ These don't seem relevant if  $> 1 R$  is assignable to any  $Q_i$ . (4) & (5) are when  $> 1 R$  is excluded from

See what corresponds to "keeping" here for (4) & (5) update condition.

In present system, it looks like there are only 2 update cards: ① & ②; ③ is pretty very close to ②

If a  $P_j$   $P(A_i/Q_i)$  is very bad we write start out  $a_j = 0$  to find an operator that xfers  $Q_i$  into  $A_i$ .

Or possibly a  $P(A_i/Q_i)$  — but for only 1 case, I imagine that an S. function would not be appropriate

Well actually, ④ may be of interest — if one of the  $P_j$  gives a good  $P$  for  $A_i Q_i$ , but  $a_j$  is too small, then we may want to modify the R part so  $\uparrow$  that  $a_j$ .

we want to modify the R part so  $\uparrow$  that  $a_j$ . We want to do this w. min. cost. ~~not with~~  
Also, we want the new R funct to not change  $a_j$ 's of other  $Q$ 's much — so we don't have to reevaluate much of the corpus.

So: Consider the above to be a "framework": T. details of the update process should be worked out in more detail, when I do the TSO.

Be sure to try both d. problems & S. problems in the early TSO — so I can be sure the update Alg's are really adequate

Also do look at discussion of "S updates" & "extensions" of P's — these may be relevant to "Fuzzy R"

One possy is to use  $1R$  for d. probs., & fuzzy R for S. probs. — in "early (infant)" TM.

REV (Sort of)

HAMPTON RESEARCH ASSOCIATES ABCDE  
26 Boylston St., Cambridge, Mass. 02138

00:

01

03

08

10

20

30

50: present picture of QATM:

1) 4.14.18 # (E2) is final Gorbang used.

2) We can use AZ (41) to assign & priops to  $\Sigma O^j \rightarrow CoF(4.18)$

3) The 3 input vnc model will be used, w. the method of 4.47.22-23 being a possibl. way to implement it

4) Two ~~kind~~ kinds of primitive functions used in: reference "unc"

a) Normal, deterministic functs: strings  $\rightarrow$  strings or (strings/nos)  $\rightarrow$  (strings/nos).

b) Problistic functions: (strings/nos)  $\rightarrow$  Pd's on strings i/nos.

5) Actually both a) & b) are not both necy, but we do have ways to xff. betw. them. (4a could give Pd's  
vz 1. 3 input vnc trick of 4.03 (2); Pd's can be so xrew as to be deterministic.

Also, all functs that TM "learns" are problistic, since  $\rightarrow$  ~~1~~ solution ppn is 2 ways possl.

Having these 2 kinds of "primitives" makes it much easier to write TSC's & solutions to them, & hence for them.

6) I will write TSC's & GNA solutions & try to find heurs that make Recog solns. of accessible cc"  
These can be done in usual ways using "conc notes" (sec 5.89) for both d. & S. type problems

7) The  $O^j$  functs will consist of a problistic Recognition function that will problistically  
assign. Solution functs to each  $Q_i$ . (see 4.55.11 for "Fuzzy Recog") - see 4.63.28 for 2 serious  
fuzzy recogn. - & some possl. ways to deal w. it.

8) The updating  $\Sigma$  (gms for 7) (1.10) are not yet clear for me: T. 5 update rules of (4.19.00, 4.29.00, 4.50.05)  
could be made a basis of ideas; see (4.60.10-40) But this needs much work!

I expect to develop the update  $\Sigma$  (gms) while working on .08 (6), guided by this (0.7ff) set of ideas

9) In writing TSC's, I'll use a lot of "indexing" of problems. Also much learning of dedus.  
of various concs. by Irug. Examples

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

This will start work on TSO & Update Algs.

**SNDP** Heurs: Very EMPT ideas: (1) Partitioning & Recon. cat. into 2 or more new R functs. Realization that a former broad category can be usefully broken up into sub. cats, w. Solns of by pc & /o ~~lower~~ lower cc than before partitioning. [Use of AND to do this]

(2) Merging of Recongn. cats: Realizing that the these 2 cats formerly had different problem Solns,  $\exists$  a common soln. for all probs in both cats. This "common soln" could have been used normally by one of f cats. [use of OR to do this]

.01-.06 is imp't partly because it makes it possibl. to "correct" errors in cats. - to allow them to "evolve" We probably need more operations than just .01-.06 for good evolution of cats. Is there any way to deriv. a "universal" set of operations on cats - so that any operation can be expressed by a (usually small) set of primitive ops.

OK: well another initial problem w. TSO's was ANL: I was able to get TM to learn +, -, x, / easily, from RPN or just (just PN). I wanted them that TM learn to evaluate any composition of these functs. The human heur, seems to be surround the idea of "quantity" - "what we get by evaluating a funct". - that it can be an argt. of a funct, just like a number.

How to get TM to dis cover Puts?

Well,  $4+3=7$ ;  $8 \times 7 = 56$  also  $8 \times (4+3) = 56$ : ~~formulas could~~  
Say we do this for several values of "4, 3, 8, 7". Could TM find genz  $2 \times (b+c) = \dots$   
even if b+c has never been used before &  $2 \times z$  has never been used before? ( $z = \text{value of } (b+c)$ )  
After being given several examples of Puts solns could TM induce a general ~~rule~~ rule for evaluating compositions of functs?

The answer is probly "yes", but the Q is, how large is cc?

Indeed w. no ~~simple~~ simple TSO: just a small set of examples of composition evaluations, its possibl. to tacitly induce the general evaln. rule, but probably only w. very large cc.

(2) One approach: look at some expressions for a general soln: see what concs might be found that would  $\downarrow$  its pc. - could these be tacit concs. But would have overuses?

Her, the idea of .12-.14 is the way I think about this problem, is there should be a way for me to express my intuitive ideas, formally.

(a) perhaps recognizing that the evaln. of "3+4" is the number "5" are very similar, so that ~~properties~~ properties of 1 write be properties of others.

.27-.28 is part of a very general heur: If 2 things have several properties in common

they are more likely to have other properties in common.

(c) Another view: That if TM "understands" what ~~RPN~~ RPN means, it will automatically be able to evaluate compositions of functions.

(d) Another possy: Just make up some value heuristic ideas & make a TSO built around them - & see what cc we get.

(e) Ask Grace for some heurs: try to implement them / formalize them

(f) Try a GA-type soln: ~~Devise~~ Devise a "fitness function" for the final soln. One such function might be: How many of a examples does it get right? Also, more credit for longer examples. Say credit =  $\leq$  no. of bits in a problems for which it gets into answer. So Puts could be a kind of "compression" criterion!

ID

ABCDE zbcdefg UNIS! ABCDE ABCDEFGH → j abcdefg  
 abcde ABCDE ABCDE ABCDE abcde ABCDE ABC ABCDE  
 abcde ABCDEFGH IJKLMNOPQRSTUVWXYZ  
 ABC HAMPTON RESEARCH ASSOCIATES  
 26, Boylston St.; Cambridge, Mass. 02138

00 : 462.40 : (462.38)(F) May be worth developing! \* A link before ALP is GA  
 T. Gore of 463.38-1.39 write to find for GA soln. of  $[O_{n+1}^i] \text{ from } [O_n^i]$  (!)  
 02 Hvr. if we did not Gore, it would take much cc to evaluate  
 T. number of more successful cards! As soon as we find  
 04 a card failed on several  $O_i$ 's, we have an upper bound on how "good" it can be! - i we may be able to reject it.

**SN** (02-04) is an Example of Q.A. (Quiz About) how Its Gore is not achievable by  
 Modifi. of V. P.D. Guiding Lsrch. Hvr, would the Gore =  $\frac{PL}{CC}$  (Gore) be adequate for P's?  
 On some dirty of 02: In Human Reprod. via DNA, etc. We do have a kind of Soln  
 to this problem! We have "genes" that create proteins & "meta genes" that control "expression" of other  
 genes (i.e. for hrs or "meta genes"). All viable genomes can combine - w. usually viable offspring!

Evalu. of cards is Hyper. hvr.

The "genes" & "meta genes" are a bit like "ops" & "obs" [ obs ± Recognition/Alarms ]

10 following 462.38 After 3+4 → 8 ; 3+4's are homogeneously linked: soln. not so unreasonable to  
 11 try. to "appropriate" substitution to evaluate to Expression.

**SN** In the TSO: When TM solves a <sup>new</sup> problem, we modify the P.D. same. If 2 new hour  
 was used, we also modify the P.D., because to have was equit to modifying the P.D.

**SN** Some General Remarks:  $O_i$  (of 4/4/18) has 1. folk. form:  $O_i(A|Q)$

T. particular subset of  $O_i$  funcs that we will consider: First look at  $Q_i$  →  $P_i^Q$

T. function "R" looks at  $Q_i$ , its output is  $P_i^Q$  on  $[P_i^Q]$ : which are probabilistic mappings from  $Q$  to  $A$ .  
 Done that of bus? The function R decides what "kind of problem"  $Q_i$  is. For each "kind of problem",  $z$

There is a set of probabilistic funcs  $P_i^Q(A|Q_i)$  that map  $Q_i$  probabilistically into  $A$  values.  
 Heuristically, this corresponds to a human's deciding that a problem is of a certain kind & should be treated in  
 a corresponding way.

An imp. affect of this categorization! When a new problem is solved or a new problem is solved or  
 an old problem is solved ~~in~~ given a better solution, the R function is usually modified only  
 slightly; and only a few  $P_i^Q$  funcs are modified or created. This minimizes the amount  
 of computation needed to determine whether the new  $O_i$  obtained is really better than  
 the old one.

There is a dirty horse! If R gives a rather broad distribn on  $P_i^Q$ , then  
 Any change in one of the  $[P_i^Q]$  sets (very 2) will change the  $P(A|Q_i)$  for many  $Q_i$  - which will  
 make recalcul. of  $O_i$  very time-consuming (or by cc).

Some ways to deal w. this ~~dirty~~ dirty! w.  $PC > 0$   
 1) Only have a few ~~outputs~~ outputs of R for each ~~input~~ input.  
 2) When we solve a new prob or improve soln to old one, modify the R &  $[P_i^Q]$  so as to minimally  
 distrib. probability  $O_i(A_i|Q_i)$  for previous  $z$  values. (32) may be a special case of  $P_i^Q$ .  
 One way (other than (32)) is to arrange R &  $P_i^Q$  so that only the new solns are affected - all  
 other  $O_i(A_i|Q_i)$  remain same.  
 33-34 Tells part of history. - part of what we want. On the other hand, we want to be able to notice  
 that certain problems are similar but not identical. - That implies 1. solution of any problem does  
 imply to some extent (never zero) on every other problem. In most cases this "implying" is <sup>spec</sup> 464.100

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

SPAC  
64349 : very small, hvr. I want a system to be able to descr. that certain probs, once that to be quite diffrnt., are indeed quite similar (Merging).

Also, I want the opposite to be possl.: Probs that to be all in some cat are best worked on as diffrnt cats.

**[NB]** There is an imp't distnce betw. T. System of S89 say, that solve OZ & Inv probs - a + present  
present QATM. - in that in S89 I could have fuzzy suggestions on the best way to solve a  
problem. If I had very many <sup>of responsible wt.</sup> this would slow down T. Soln.

In QATM, the wts. are not <sup>for</sup> suggested ways to solve probs - They are rel. wts. of probn. methods

- Having them wrong gives bad soln. - its much less forgiving than ~~S89's~~ for Inv, OZ probs  
Hvr. in S89, I really didn't have an exactly formalized method for the entire system

→ I write try to to this - go from S89, S89 completely: I have to proof read them any way!

So: Present by Q's How to resolve 463, 288: Just how narrow R is.

② Got a TSO started so I can get ideas on how to do updating, also, fit my panel system to "R.W." ①

Re: ① There are 2 extremes: i.e. to one I considered ~~1~~ (2) for each Q.

① ~~2~~ has  $pc > 0$  for all  $i$  for all Q. (which has diffy of 463, 288)

I did consider ~~250.09-20~~ was the beginning of a method that ~~exactly~~ I used (perhaps later?)

to assign wts. to various (perms?) when  $> 1$  seemed relevant to a problem (EQ<sub>i</sub>). → .20

→ What I may do: start w. approach of attempts to have one  $i$  per Q<sub>e</sub>. Hvr. empirically, we score find first  
we get  $> 1$   $i$  per Q<sub>e</sub>. This was dealt w. in 450.20 ff: update solutions #4 i's There, I tried to  
arrange so that there was only 1  $i$  per Q<sub>e</sub>, but I need not do that! I could read just wts, so that ~~Q~~ was lower (↑).

~~stagnant~~ (which may be close to ↑  $O^j$  (Ann | Q<sub>in</sub>)). 2/0 (Modify ~~Q~~ / introduce new)  $P_i^R$ .

Usually, we like to have one  $i$  per Q, but sometimes we can't - Two  $\{P_i^R\}$ 's look about equally good,

or, we can't find a way to modify R to emphasize the best  $\{P_i^R\}$ .

Dot  $\vec{P}_i$  is the ~~unweighted~~ set,  $P_i^R$  ( $R=1$ )  $\&$  max  $i$  Good! (its not a vector in the usual sense, but its a fine notation!)

T. details of (.16-.21) have yet to be worked out, but it seems like a very reasonable approach.

**A** modify of (.16-.21): T. probn. is not a ~~set~~ R wtd ~~subset~~ mean of  $\vec{P}_i$ 's;  
Instead, the  $\vec{P}_i$  of max wt. is picked for pred. If  $> 1$  has top wt., we use average of them.

Again, it ~~seems~~ several ~~cases~~  $i$  are assigned to a Q<sub>e</sub> by R from them ~~we~~ have 2

"merging hour": it suggests that this (set of  $i$ 's) are "n"  $i$  so we try to merge one or more  
of them by finding ~~a~~ a new  $\vec{P}$  that will solve all problems that point to this (set of  $i$ 's).

We can try various sub sets of that "set of  $i$ 's" but we can Usefully Merge.

.14-.15 is correct & relevant: I analyzed (.14) in which sometimes one or more perms would be relevant  
to a problem (usually only 1)  $i$  I got optimum wts: but I'm not sure it was exactly ~~best~~

to problem I have now. For cases where  $> 1$  perm was relevant - each perm was  $G_n$   $\approx$  wt.

(for all cases by that perm)  $\rightarrow$  the corpus pc was max. I guess this will work, but I'm not sure its  
the best way. (I did some that about finding the peak ~~at~~ <sup>2/0 ~~method~~ ~~divided by~~</sup> ~~determinant~~ ~~at the peak.~~

for the ~~use~~ value of  $O^j$

This amounts to <sup>simultaneous</sup> optimization over all of the wts of all of the  $i$ 's. Hvr. the  $i$ 's can be put into

"subsets of ~~mutually~~ interaction" - some either exactly or approximately. This means that when

a permutation new Q<sub>in</sub> is done, only few wts of the  $i$ 's in its "interaction class" need be  
considered. We only consider factors (groups?) in which  $> 1$   $i$  are involved.

ID

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

"Levenberg - Marquadt method"  
Garshanfeld "Principles of Mathematical Modeling"  
P124 §10.4.1

Spec  
464.40

On 1. optzn calcns needed for 464.30 - .40. The "Marquadt" n.2. optzn method would probly be fine. ~~But~~ One normally has a good initial approx. available so perhaps usually only one "iteration" is adequate. We can easily cut off interactions tests below a certain level &/or arrange so that we never have > 3 or 4 ~~z~~ z's applied in a given problem (as in cut off loose wts) - this will result in smaller sizes of "interaction subsets" so optzn will be much faster.

NOTE: T. can't work on this problem may have been 267.35 # (i. 35 line is where line 19 should be!)  
The discussion runs to 269.40 from 274.09 - .18; note (274.10 - .14) as a generalized principle.  
Also: 273.19 ff seems very relevant to present. It would be well to reread a summary of that stuff.  
That General area of TM notes looks very relevant to present problems - T. idea of FC (which is essentially a fuzzy R function) was very impt. in that Model. 274.23-24 was supposed to solve all "Procedural" problems of TM ~~also~~ <sup>also</sup> ~~also~~ <sup>also</sup>

On 274.23-24: Appl. Model for QATM:  $G_{orc} = \sum \sum z_i O^j(A_i | \phi_i)$  (?) maybe not justified!  
2 input vnc: ① Input data  $O^j$ . ② Inputs  $O_i$ : outputs  $A_i$ . Each  $O^j$  has to cover all  $A_i$ 's (w p > 0)  
 $O^j$  is a single code. I think this is just an alternative way to look at 3 input vnc model!  
We just include the  $R_i$  inputs as part of  $O^j$  code. This avoids diffy of some decrs essing very ~~pc=0~~  $pc=0$  to certain QAS. We can think of  $R_i$ 's as being put in "at some time as  $Q_i$ 's", or all at t. start or all at t. end of t. rest of t.  $O^j$  open  
Another way to look at it: The  $Q_i$ 's are one by input, t.  $R_i$ 's are part of  $O^j$  decr; t.  $A_i$ 's are one by output.  
Hvr, t. fact that  $Q_i$  are "unordered" seems to make this different from an ordinary 2 input induction problem.  
I.E.  $C^M(S, Q) = A$

So: if ff Seems to solve a Big Confusion: How we could have a set of d-~~decrs~~ decrs of t.  $\{Q_i, A_i\}$  couple & have them give PC's for any of t. Q's. It shows how this model is a way of looking at t. 3 input & basis function Model.

**SN** Marvin's "Soc of Mind" / Schrodinger's Paradoxonism: My present approach to TM to Q's via conc. nets can be regarded as a kind of way to realize these models. Could I use any of t. heurs Marv mentions in "Soc of Mind" to help design TSC's? (4/9/02 - how did I make correspondences??)

So ~~what~~ What we want is a R that usually assigns 1 or a few "z" values to each  $Q_e$  - that R can assign approx wts to them a priori or assign an optimum set of wts, using t. methods of ~ 267.35 eco (see Bibnote 605). Th. latter assigns invalued wts to z values. These wts are indep of t.  $Q_e$  they are applied to! - so they are somewhat diffnt. from what I'd normally consider! - i.e. I usually think of z's being assigned to  $Q_e$  along w/ a suitable wts by t. R funct. So I feel uncomfortable w. this data wts assignment! I feel that R(.) should decide on t. wts as a function of t.  $Q_e$ 's. So t. conflict is born w/ being property of  $Q_e$  or wts of z being property of z'. At present time, perhaps pursue both approaches. These seem to be several ideas have: ①. problem  $Q_e$  data soln methods is ②. "Problem type" - <sup>Spec</sup> 469.00

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

On curve fitting using more curve params than data pts.

Say we have  $n$  pts & we're using  $n+k$  param curves. There will be a (2d) manifold that passes through all  $n$  pts exactly. Hvr, say we want to have ms error  $\epsilon^2$ . What is the volume of region in this  $k$  space that gives  $\epsilon^2$  error? Say we know  $k$  params are on  $(-1, +1)$  so  $V$  is bounded. We want  $V$  as a fcn of  $\epsilon^2$ . for  $\epsilon^2=0, V=0$ . Is there a usual functional form for  $V$  as a fcn of  $\epsilon^2$  for small  $\epsilon^2$ ? I could probly study it empirically!

Anyway, if  $k < 0$  then for  $\epsilon^2=0$  there are no points. (Usually). There will be a min value of  $\epsilon^2$  that allows  $V > 0$ .

The idea here is that for each value of  $k$  ( $\geq 0$  or  $< 0$ ) we will be able to make ratios for the data using a certain amt. of info into  $\text{KNN}$  ( $\approx -\ln p$ ). Empirical,  $p$  corresponds to the volume in space, and if  $k < 0$ , we have to add in info needed to reduce errors.

In the case of  $k > 0$  (and fixed), we may be able to do extrapoln. to pts not in the original data set.

Say for  $n$  data pts  $n > 1$ , we have a certain  $\epsilon^2 > 0$ , (small), we have a volume  $V$ .

If we fix  $n$ . First  $n$  data pts & vary the  $n+k$  data pts  $X_{n+1}$ ,  $V$  will be a fcn of  $X_{n+1}$ .

Then the extrapoln of  $X_{n+1}$  from the first  $n$  pts should give a pc of  $X_{n+1}$  or  $V(X_{n+1})$ .

presumably, this distribn will be indep of  $\epsilon^2$  for small  $\epsilon^2$ .

I could do a MC Carlo Simulation: Say  $k$  params are each limited to  $(-1, +1)$ .

We put  $n$  random pts in this  $n+k$  dim space & see how many pts are within ms  $\epsilon^2$  of all  $n$  pts:

The first  $n$  pts are fixed, & whenever we get a curve within  $\epsilon^2$  of all  $n$  first pts, we see what its output is for the  $n+k$  pt. The distribution of these  $n+k$  "predictions" should give its P.D. ! We could do a plot of kind  $\epsilon^2$  vs one of these predictions. Or, we could make an integrated d.f. plot of how many preds,  $X_{n+1}$  are  $> y$ .

For large  $n$ , this is a very inefficient way to do predn, because so few MC Carlo pts will give a curve w. acceptable  $\epsilon^2$ . A poss. way to deal w. this: Make  $\epsilon^2$  very large, to start off. If we find an acceptable  $\epsilon^2$  vector ( $\vec{\epsilon}$  is an  $n+k$  vector of  $k$  params of  $k$  curves) we try pts nearby in  $\vec{\epsilon}$  space, & zero in for smaller  $\epsilon^2$  until  $\epsilon^2$  is  $>$  our finally desired threshold. We then start out again w.  $\epsilon^2$  large and go thru the process again.

This reminds me a bit of feed forward ANN prediction, in which we start out w. random wts, then evolve until we get a certain threshold ms error. — we make a predn.  $X_0$ , then we start out again w. random wts & evolve to  $<$  rms threshold & make predn  $X_1$ , etc. If we do this many times, we'd get a d.f. for  $X$  th. pred. We wouldn't need

"Cross Validation" (old faithful)

In either case, we start w. random  $\vec{\epsilon}$ , and "will climb" until  $\epsilon^2 <$  threshold

pick  $\epsilon^2$  of size large enuf so we can hunt for  $\vec{\epsilon}$  w. ms error  $< \epsilon^2$ , then from there we'll climb to ms error  $< \epsilon^2$  threshold.

I suspect .12-.09 is a proverbially good soln. — T. q. is .13-.07 or .38 better way to realize it? see

.20-.24 is a biological realization of .12-.09



ID

ALP Goal. is  $\frac{1}{2}$  ANN applies

358
35
386
420

**HAMPTON RESEARCH ASSOCIATES**  
 26 Boylston St., Cambridge, Mass. 02138

00 : 466.40 : A possibl. criterion as to how large  $\sigma^2$  final  $\epsilon^2$  should be. Perhaps it should be  $\leq 0.1$ .  
 true  $\epsilon^2$  of  $X_{n+1}$  to  $X_{n+1}$  predn using very small  $\epsilon^2$ .

Say  $\sigma_F^2$  is final threshold over  $\sigma^2$  and  $\sigma_F^2$  is the resultant apparent  $\sigma^2$  of  $X_{n+1}$   
 Conjecture:  $\sigma_F^2 = \sigma^2 + \frac{\epsilon^2}{\epsilon^2_{\text{min}}}$  : i.e.  $\sigma_F^2$  is like "smaller effect size" on  $\sigma^2$  of  $X_{n+1}$

If this is really true then  $\sigma_F^2$  would be an optimum  $\sigma^2$  to use, to obtain a certain precision in  $\sigma^2$ .

A (perhaps) related problem:  $k < 1$ , but it's hard to find the best curve.  $n-k$  is still quite large & we are concerned that there may be many local min of  $\epsilon^2$ .

In this case we will use Mt. Carlo methods as 466.20 ff. — I think this is commonly what's done in ANN research.

Another Monte Carlo approach for curve fit (466.20 ff): Pick random pts in  $n+k$  dim  $\vec{z}$  (cert) space. For each pt, plot its  $\epsilon^2$  vs deviation from first  $n$  pts; vs its value for  $X_{n+1}$ . This method helps to estimate  $\sigma^2$  of  $X_{n+1}$ .

15 Another possible way (not Monte Carlo): Look at  $k$  dim manifold that goes thru all  $n$  data pts. Get the D.F. of  $X_{n+1}$  over that manifold. For  $k=1$  or  $2$ , it might be possible to do this analytically.  
 17  $\rightarrow$  For  $k=0$ ,  $n$  data pts,  $n$  params, 1 curve only; the predn. would be 1 pt for  $X_{n+1}$ ! — doesn't sound v.g. | Could  $k=0$  be the worst possible case?

20 As I remember, using linear predn, there would be a certain no. of coils (or) that would give optimum predn. It would be well to see how much "compression" we get for  $n=0, 1, 2, 3, \dots$  — i.e. write out the actual pc of  $\text{dcm}$ .

17 is the Main Negative criticism of the System. The attraction of the idea is that it gives an apparently good way to analyze ANN & also may suggest a approach methods for other kinds of induction problems.

For 17: Say  $\vec{z}_0$  fits all  $n$  pts; and  $\rightarrow X_{n+1}$ : for hypersphere of radius  $r$  about  $\vec{z}_0$ , what is distribution of  $X_{n+1}$ ? : What is  $\vec{\nabla} X_{n+1}(\vec{z}_0)_{\text{at } \vec{z}_0}$ . The dirns of  $X_{n+1}$  with each component of  $\vec{z}_0$ .

30 In contrast to 17, it would seem that larger  $k$  would be good — giving the system lots of possibilities for many kinds of models.

More specifically on ANN we have  $(\frac{1}{N} \text{ TNS})$  Note wts. We use cross validation. We pick random wts to start then do vigorous, exact, gradient max slope. This gives a 1 dim. (Time  $\propto T$ ) path to peak (it could be a local max). We evaluate core wnt test set along the  $T$  path & pick  $T$  value for best core. This is a 1 dim. optzn, so the expected core is not far from observed core. If we repeat this process (w. different random initial wts) many times & pick the best of bests we get SOY (Spurious Optimization Yield). Similarly we could use same initial wts but would use path of max slope w.r.t. random chosen cases; we get randomish paths — i.e. we can pick pt. in path w. best core wnt test set. Doing this many times & picking best, gives SOY again.

(Spec 468.00)

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

spec  
00 : 467.90 : We could do similar stuff w/ to credit of 466.00 ft. We pick a random initial set of params, then take a max slope path w/ a ~~max~~ max error bar. We have a big set of data to do this on. We also have a test set, — T. process is same as in ANN

But I'd like to be able to do it Logitly w/o. "Cross val."

It would seem that the methods of 466.00 ft. w. logistic should work. It is probably easiest to

analyse if the params are linear coeffs (like fitting w. ~~the~~ polynomials) — But they will be that this is specifically the system that will not work! i.e. the non-linearity used in ANN is specifically diffnt. from linear coeffs

**NB** 466.00 ft is not a strictly "Academic Exercise": I do use this method of averaging over many copies of = wt, in several different kinds of problems: One is ~~the~~ Jorgensen's "max flatness" approach to ANN (using Hessian). Also my coding to Bernberg w. code words of all same length. Also my General soln. to the TM Gore problem of (perhaps) 274.10 - 14 — which uses a Uniform sampled over continuous params & integrates over mean = 1.

Here, in this last (is perhaps Jorgensen's ~~unsuccessful~~ approach) I think  **$K \ll 0$**  i.e. no. of params  $\ll$  no. of data pts. Perhaps it also works if  $K \gg 0$  — many different coding methods considered.

→ But do try analysis of  $K=0$  for linear & for ANN non-linear stuff — e.g. A. Barron 1993 on Analysis of ANN — no of params needed for curve fitting. I'm not sure that Barron would do exactly this problem, but

10

20

30

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

00: (465.40) : "Problem type" is a bit unclear. It is assigned by RC's to  $Q_i$ : It can be a set of  $i$ 's or, a set of  $i$ 's w. assoc. wts. that  $R(\cdot)$  gives to those  $i$ 's for that  $Q_i$ . If a set of  $i$ 's, each  $i$  has its own "comparative wt." + indep. of  $Q_i$ .  
one way to look at (3): each  $i$  looks at  $Q_i$  & decides how much it (the  $i$ ) is (relevant) to that  $Q_i$ .

10: P for hops (presumably) + "relevance wt." depends on certain "properties" of  $Q_i$ .  
Slightly "Asside" When one gets a new problem, one does not feel that any of the standard methods are appropriate. One does not normally take a wtd mean of their results. [unclear] - Depends on whether

it's a simple induction problem or a inv. or many kinds of  $Q_i$  problems. In my present QATM, no new "creative" work is done when  $Q_i$ 's are answered; the creative work occurs when updating. Answers to  $Q_i$ 's are always "Routine". A possl. exception might be invg. during the soln. of a problem - So one would "update" the "learning" part, during the problem soln.

So! During many problems, imp. "QA's" occur! - So system should be updated before proceeding.

14: (04) Well, I could just have R assign wts: But if R didn't have any good ideas as to what the wts should be, we use "i" assigned wts - that are obtained in a standard way, by averaging over previous cases of several "i" for some  $Q_i$  (267.35) ff (465.05 bibl note)

Remember that R's problem is a regular QA-type induction problem, but the data is not QA pairs, but  $Q_i$ , successful  $P_i$  pairs [also, we can use unsuccessful  $P_i$ 's as data, but then the  $Q_i$  has to be "indexed" so what we want is "Pilled" not success.] The  $Q_i$  to  $P_i$  induction is easier than  $Q_i$  to  $A_i$  induction because  $P_i$ 's are fewer  $P_i$ 's (less info needed). - So we have  $P_i$  "2 phase induction" - first "R" gives  $Q_i$  to  $P_i$ ; the  $P_i$  goes from  $Q_i$  to  $A_i$ .

Since R's problem of  $Q_i \rightarrow P_i$  is a fault induction problem, all  $P_i$  will get wts  $\rightarrow \phi$ . I suspect that we will be only interested in the top few  $P_i$ : We could try the system for various cutoff pts. for no. of  $P_i$  considered.

$P_i$  as we have cc for - usually only a few should have much wt. For updating, we try many of the  $P_i$  to see how well they do in  $Q_{int}$  that is - If we don't find any good ones - we try to generate new  $P_i$ 's somehow based on our knowledge of  $Q_{int}$ ,  $A_{int}$ . Otherwise, we might just try the  $P_i$ 's in order of wt. - (Maybe Lynch) since certain trials will take too long.)

05/01 The "trials in order of wt." - wts could be arip assigned via  $A \geq 141$ . - This is att. by my what. t. wts ( $\equiv$  pc's) are, but presumably the pc's get a better empirical basis as

"TM" structures - There's an imp. point. we start out w.  $P_i$  assigned via  $A \geq 141$ . As we get more comp. data on how good the  $P_i$  were for various problems, the RC assigns  $P_i$ 's to problems more accurately.

37 We have (3) things going on w. R &  $P_i$ : (1) For a fixed set of  $P_i$ , the pc's assigned set better as  $55 \geq A \geq 141$ . (2) When the best  $P_i$  do rather poorly on a problem, we generate new  $P_i$ 's for that problem. (3) When pc's assoc. w. all useful  $P_i$  get too small, we drop the best  $P_i$ . (Not using a threshold here seems "PROCRUSTEAN" - But we may want to do it to save CC in the General System)

472.00  
SPEC

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

(Kenny)  
Push back on this  
(Given fact)  
problems

Time to Go back & look again at GA: What its faults are & how to fix it!  
Faults: 1) When problem is given, it has to be put into form so that GA can work on it usefully.

This involves: ~~Deciding~~ Deciding on Population size; type of Mut/crossover & val fitness of Perm.

Design of Fitness function - a proxy of the fitness. This definitely has not been "mechanical".  
[Initial "start" population]: Also monotonic xfm. from fitness func. to reproduction frequency

2) After a problem has been solved by the system, it discards all info obtained during it. Both.  
"Kill the Egg" - then conceive ~~and~~ and from from babyhood a new Eggs for new problem."

Very wasteful.

3) Evaluating the fitness funct. can take ~~almost~~ almost all of ~~the~~ the soln. Approxns. of fitness func. should be used; As we approach a final soln, our approxns must improve or we don't use approxns any more

10P

SN I was concerned w/ logic of evaluating cands in GA; if the problem to be solved was changing finding a 0 for the entire cands. We can weaken this cond considerably! There are diff "Areas of knowledge": Each update affects only one of many. Also, noting which problems the cands fail on, we can do trials for "Quick Abort" Note = 31

Solns. to these problems:-

1) 2) One way would be to have ~~prepared~~ a set of probs of diff difficulty (e.g.). These ~~most~~ probs would not be "prepared" by the trainer but would start out w. essentially "prepared" examples & continuous, less "prepared" problems. This, hvr, is only one aspect of TSP.

3) The problem of the system at time  $t$ , is to find a funct. that can solve the first  $i$  problems. (see .20R for a poss. mod'n. of this idea)

As  $t$  ↑, evaln. of the "Fitness func" becomes very time consuming see 10R for one approach to this. Another approach ~~is~~ In addition we can do approx. evalns much of the time (.09-.10)

.20R

SN T. idea of a "problem pool" The GORC is how many problems solved - on a wtd sum of problems solved - hard probs get more wt. If GA is used, we have Mut/cross betw. cands that have solved most probs: We wite "index" cands, telling which probs they solved - this info would be used in Mut/cross for new trials There is a previous note on .20R within last wk, certainly within last month! I think at bottom of page.

4) TM remembers by a population of solns to past problems, when a new problem is given, TM starts w. a population that is able to solve

"Similar" probs int. past. Knowing "What is Similar" is a serious

problem: Perhaps find way to use GA to solve it.

Also, on a basis of "similarity" it can give reasonable values for

part of 100-104 (Big part of fault #1)

Similarity can be framed as an individual problem w/ usual ALP GORC. So it is a "Fitness function" - probably directly usable.

SN I've been thinking of the GORC as being an adequate soln to all probs upto now (18-19)

Perhaps this is unwise: If the TM can solve the last 10 probs (or 10 important care fully selected "Quick abort" probs, this may be adequate & would reduce Testing ~~cost~~ cost enormously! (Also see 10R)

Another aspect of this: ~~Exercise~~ Our Population ~~pool~~ pool does store lots of info about solns, so there is a strong tendency for a cand generated from Perm to solve many probs in the TSP if it ~~is~~ is tested on only a few.

One way to do this is to use GA to update  $\alpha$  (= 4.1.18) - noting above argts on why it may be practical.

IP

GA.00

Recursive functs: ~~2.4~~  
2.4 - 2.5 is implicit c. of true @)

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

Spec 470.40  
00/470.40 → Another Deficiency of GA is that I think there should be better ways to get cards than mut/cross: Grammatical evolution, perhaps... I've written a bit on this.

Another very imp. pt.: I'm thinking of a QA machine! Since any problem can be put in this form, w a suitable TSD, it should be poss. to get this GA system to work on all of the problems of 470.00-10 - to improve itself.

In designing a GA system of this sort, I will probably get ideas on how to do "updating" for a non-GA QATM. It might be very useful to design 2 systems in all.

09 → A poss. Diffy (a poss. strong pt.) of this system! I will not be "micro programming" this system: for the most part, I would not know how it's solving problems. I'll have to do so by using

concepts, but I really ~~don't~~ <sup>wouldn't</sup> know if this system solves the problems using the concepts in my (conc. net)!

Recursive functions:

Simple recursive functions: T. value at one pts x, is related to value at another pt (f(x)) by a continuous fn. so  $G(x) = H(G(f(x)))$  T. functions  $\in FC$  &  $HC$  have already been defined, & this is a partial defn of G. T. rest of defn gives G's value at some pt  $x_0$  so  $G(x_0) = y_0$  is known. so  $G = H(G \circ f)$  is one functional eq. &  $G(x_0) = y_0$  is "Boundary Value". I have studied  $G = H(G \circ f)$  in the past, as a kind of "difference" eq.

$H^{-1}G = GF$  which is kind of commutation relation betw G,  $H^{-1}$  & f. so if  $J = H^{-1}$ ,  $JG = GF$ .

20  $JJG = (JGF) = (GF)F = GFF$  so  $JJG = GFF$  means  $J^{(n)}G = G \circ f^n$  which is not involving any expanding from  $G(x_0) = y_0$ : we get  $J^{(n)}G(x_0) = G(f^n(x_0))$ : this gives G at values  $f^n(x_0)$ :  $- \infty < n < \infty$

If we can get n to be real or just rational we define G for a continuous range - but usually an integer n (usually  $\geq 0$ ) is used.  $f^{(n)}(x) = x$ .

24 → This is imp. for defining funcn in QATM or any other TMS, because we often have  $G(x)$  bearing a known relation to its value at some pt. related to x. (I think) - But E should find

25 → Examples It should be an imp. Heur for defining "Recursive" functions.

Differential eqs, Difference eqs, Differential/integral difference eqs are all special cases of Recursive function defns. They all used "boundary values" as part of f. defn.

05/01 [SN] on Recursive funcns.  $f(x)$  can be defined by  $(f(x))^{(n)} = f$  so we can then define  $f^{(n)}$ . For  $G(x) = f$  may or may not be solvable: if  $f = x^k$   $f^{(n)} = x^{k-n}$  we can define  $f^{(n)}$  in R3

30 case, even if n is complex. If we get an analytic soln for  $f^{(n)}$  as a funcn of f & a param

We can make n complex/continuous, ect. - Maybe n can be a vector or matrix.

in  $x^n$  can n be a matrix? If  $n$  is a square matrix, any analytic funcn of n is poss. (defn'd).

If we use a coord system in which  $n$  is diagonal, then the interpretation is very simple!

ID

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

267.35 ff (see 465.08 Bibl. notes)

00: {spec} 469.40: Re: 469.37: The "improvement" of the pc's of  $\vec{P}_i$ 's could be done via the method of 267.35 ff. This assigns a pc to each  $\vec{P}_i$  indep of  $Q_e$ , hvr. But  $\vec{P}_i$ 's do not respond to  $Q_e$ , but in a Logical (Yes/No) way.

This "267.35" stuff does sound rather gross! — But it probably more or less "works".

03: Intellectually, it would seem to be more reasonable to gather data on the effectiveness of various  $\vec{A}_i$ 's on various  $Q_e$ 's and tailor our  $R$  accordingly. This is what 469.37 is about.

05: The thing about "267.35" ff: If we make  $\vec{P}_i$  assignments to  $Q_e$  in a Boolean way (yes/no), & we sometimes have  $\vec{P}_i$  assigned to  $\neq Q_e$ 's, and we want  $w_i$  to be a function of  $i$  only, then "267.35" ff is correct. If we want to use more info for  $w_i$ , — make them dependent on  $Q_e$ 's as well as  $\vec{P}_i$ , then

469.37 (1) is a better way to go.

In an "infant TM" we might use 267.35 because there is little info available.

But, I suspect that there could do \* Yes/No assignment of  $\vec{P}_i$ 's to  $Q_e$ 's — one could just as well use Fuzzy (soft)  $\vec{P}_i$  assignments. Think about how pc's are assigned to  $\vec{P}_i$ 's for  $Q_e$ 's. I had an idea of "kinds of problems" so  $Q_e$ 's would be "pooled" — a N/Y type of categorization.

Anyway, 469.37, (1)(2)(3) is a nice formulation of the problem. The part that it doesn't talk about is how certain problems are or are not "solvable" & so can be solved by the same  $\vec{P}_i$ 's. I feel of this is a "Boolean" thing — but, it is, in general NOT — but people usually think about it in a Boolean way. I certainly will have sets of problems that are all solvable by  $\vec{P}_i$ .

Same set of  $\vec{P}_i$ :

06: So we have a diff: 469.37 ff v.s. (183-19): If 2  $Q_e$ 's have very similar pc costs for  $\vec{P}_i$ 's, we might consider them to be an "equiv. class". If there are many (or a large) costs, it would be easy to compare the 2  $Q_e$ 's.

I'm tempted to use 469.37 on the TSP to start, then see how I can fit in the "equiv. classes" (183-22)

The stochastic  $R$  is not the old FC. I used in the S89 model while I was at IDSIA.

At first, the indices given by trainer will tell TM which probs are "equiv" in the sense of expecting a common  $\vec{P}_i$  soln. Then (later, TM will induce (S-funct) on the basis of observed, empirical  $Q_e, \vec{P}_i$  pairs that have "involved well" (i.e. given by pc's to  $R_e$ 's).

At first I will have only 1  $R$  class! It will be solved by a single function that grows as the TSP grows. This  $R$  class will be defined by the index given by trainer.

We may try several "R" classes: some would be  $\vec{P}_i$  problems.

After we have a fair no. of R classes and TM has had to do some S-probs of suitable kinds,

09: We can put in unindexed probs & get TM to assign pc's to them (i.e. "R" function) with which  $\vec{P}_i$  found would be appropriate.

So I want to try a TSP. (maybe only 1 R class!) & see what (tools/derivs/hours) are needed to "solve" the TSP.

I was thinking of ANL as an initial TSP. But what about "Definitions" — how would they fit into the context of the ANL TSP? Would I want each defn. to be a new "R class"?

perhaps all defns. could be a single "R-class"!

What I may want is a general framework for (ing. of TSP's — a set of mechanisms that might be used to implement Heuristics.

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

00: 472.40 : Re: ANL: I could have it (in +, -, x, ÷ as 4 separate things. That learning more complex notation could be based on that (imp.

Alternatively, have it (in all these operations as a single function, but with "component functions" that "recognize" what to do in the various situations.

BACK in my mind, is always the idea of the **Context** — that all problems are to be solved this way. If so, I have to make the concepts available at reasonable pts when they are needed. One way to create conc. was by

06 **"Context"** — few similar situations in which the conc. was needed was to <sup>or "someones"</sup> situations in which it had been used in the past.

I think 06 says very much what is meant by "context", when its use to enhance point concs to be used in a new problem.

Going Back to ANL: I had this (human) concept of "value" that I wanted TM to (in, so it could guess that substituting the "value" of  $3+7$  in an expression that had " $3+7$ " in it, would give the same value.

to express.

A different approach was to use a machine containing a "stack" to have TM (in, to use it — which is what I did in Saurb.

Well, the "value" concept; perhaps its main property is the substitution idea! So if TM understood "value", it would already know how to evaluate any alg. expression! — So maybe my difficulty in getting TM to understand "quantity" (= "value"), was that I wanted to explain it w/o telling it about its use in substitution — perhaps an imposs. thing to do!

So perhaps I should be more spinnerish in teaching TM about value! E.g. say:

If  $3+7=10$  then  $3+7+5 = 10+5 = 15$ . Or even do it w. literals. (?! how?).

Anyway more to do that suitably infantile human should be able to work: then devise TM, so that it can also work that way.

Perhaps teach many of "substitution": its imp. in evaln. of expressions (understand solving eqns. "teach by induction from examples".

T. Try to be (and is meaning of  $Subs(a, b, c)$ ) substitute a for b in expression c.

Maybe ambiguous some times!  $a = x, b = 2x, c = 2ax$ :  $\rightarrow ax$  or  $2x$ !

25 Very Early in RB — (maybe in BVLG?) I printed up some alg. problems — that did not involve ANL! — Some were about Literals, other, simple eq. soln. — Try to find this!

Maybe like:  $a+x=b, a+3=b, x=?$

What's probably a major problem: I tend to make assumptions about what is "Already Understood". In Form. Mathematics they come closer to making clear what all these assumptions are.

One poss. (New) approach: write down lots of "true" expressions, that seem to have "adequate info", then begin asking new Q's. (first ask "old" Q's!)

In .18 TM could learn about substitution — just what "equality" means.

Also, the idea of scope when being local to "this problem" is often in to being "global" — to all problems. "If  $a+b=c; b=3$  then  $a+3=c$ "  $a+b=c$  and  $b=3$  are "local" info.

"If" seems to have a special meaning here: is "localization indicator". I think I ran into this idea in the BVLG stuff related to. in .25.

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

00 : 473.40 : **SN** : In doing an early TSO, I can probably go far in certain directions even using certain minimal inductive concepts. e.g. stick to d-probs: Use very simple "built-in" induction for "SI" aspects of the system (mainly L's rule - part of which is built into the AZ141 Notation).

I was concerned that I'd need some "s-funct" (long in even Elementar

**SN** Re: use of ft. pt. for early TSO. There is the problem of error size. I could explicitly introduce error size, by giving an "e" for each problem.

Alternatively, start w. integer arith. ~~and~~ w.o. division: Then introduce division, but with all problems having ~~integer~~ needing only simple-division w.o. remainders. Perhaps introduce division w. "remainder" - so all arith is exact. But this amounts to  $4 \div 3 = 1, 1, 3 (= 1 \frac{1}{3})$  So we would have to teach it how to do arith. w. these "3 vectors". It would amount to knowing how to work w. rational nos.

So ~~one~~ <sup>TSQ</sup> would teach integer arith - & division would always (in f. examples) have zero remainder. (If this is poss. ... TM will be doing experimental division of 10 by 3) So I think TM must know that  $10 \div 3$  is not just 3.

**SN** Jerry & Marcus have been trying to find general shortcuts for doing induction: (w.o. TSO's) Can it be proved that TSO's ~~are~~ (experimental) is necessary for solving induction problems w. acceptable cc? If L's such is, indeed, an optimum way (if all info is in P-Br), then w.o. TSO's the cc  $\left( = \frac{CC}{PC} \right)$  is, indeed, enormous. So only by providing a chain of problems in reasonable order, would diff. probs ever be solved.

"Trouble w. it" is that a human, (a presumably, a machine) does not need external help provided TSO's to do research. T. Matrix (TM/scientist) could simply find problems that seem within its CB to solve, & work on them. T. Matrix TM may be using the "Problem Pool" type of TSO. A machine Scientist knows of many interesting problems & works on those that he thinks he may be able to solve, using his current set of concs.

**SN** Back a few pp (20 pp ago) I had this idea that the representation I could represent the pc of a seq. of QA pairs by a set of probs in  $\Pi$ . On second thought, that would be impossible if any Q had  $\geq 1$  diff. A int. prob. Anyway Marcus (perhaps), defines a pd of  $P(A|Q)$  as [ p 3.14, eq(2) of their proposal] This looks like a 2 input unpr:  $P(A|Q) = 2^{-L(Q)}$   $M(R,Q) = A$ . So Marcus' error: Marcus is just using a universal df. for  $P(A|Q)$ .

**SN** TSO's for Sol 56 sound easy & not to write! Somehow I didn't have the problems & I've been getting in writing TSO's - even for ANL! Perhaps just write up some TSO's for a human & get desired responses: See just what into is needed for these responses. There diff. things to learn & diff. in constructing TSO's to try them - but I haven't really gotten into the "serious problems"! There are some fundamental notation & "understandings" that I haven't figured out how to deal w.!

I wanted TM to somehow acquire the concept of "quantity", but actually, I don't really have to do this: Just write out solns to  $3 + 4 = 7$  & then solns to slightly



HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

00: 474.40: more complex problems, than soln to general ANL problem! Try to devise a seq. such that a human should be able to handle it. E.g. ~~sub~~ solve w. depth 1, then depth 2, then depth 3 etc. Also 1. then 2 ~~then~~ substitutions at depth 1, etc.

See if idea of substitution can be taught (say do first first);  
e.g. subs (a, b, c) means substitute a for b in expression. See what comes

2nd way: or maybe use ints as a prim. conc!  
Try to list a set of prim concs that would make lots of alg. operations "simple" (MDE)

BUT **DO TRY** - 474.40ff.

In doing ~~it~~, try to get > 1 soln, & save all of them as possl. component containing "objects" & for "Backtracking".

To lang I use for solns, ~~state~~ should, at first, be mainly "English". (i.e. ambiguous if neccy)

Perhaps the way I got into that bind w/ "Quantity"! I assumed there was a way to write T&Q's was "Topdown", start w/ soln. to Diff't problem, then find other smaller probs that have parts of the soln, and so down to prim. concs.

Hrr, until I have more Experience w. this  $\rightarrow$  "Top down" is not neccy to best way to start out!

Perhaps try T&Q from Elementary Alg text.  $\leftarrow$  As they occur in text. - see what concs are missing!

Hrr, write out Solns. in "English" - After a fair amt. of T&Q has been done in English -

~~only~~ only then, ~~then~~ begin to formalize. ("Formalize" means "remove (noise/uncertainty)"),

Another source of inspiration for concs (prim. or not prim) are words in English used to (discuss/describe/solve) math probs. Not neccy terms that are formally defined in Math.

In Sept I did write out solns to some problems of Int. like linear  $\rightarrow$  quadratic  $\rightarrow$  cubic eqns. Also a weird way to "try to solve Algebra".

Trying to do ~~some~~ Ramougn's text is also a possy. I don't have to use simple "primitives"!

I can use whatever "prims" seem neccy to solve the set. The main thing of int. is that the concs should be a "universally Extensible" set in a useful way. (Any Set of concs can be made usable by adding a small set of concs.)

One Approach: To do many diff't T&Q's of very diff't kinds: (Many of them starting w. a rather large set of complex concs). See if I can learn any thing about writing T&Q's: a book update & films: from this exercise.

**[NB]**  $\downarrow$  Elementary problems, like a ANL, are very elementary algebra, and diff't because it is often hard to call (describe) just what the problems are: what are local, v.s. global concs, etc. So notice this! Perhaps avoid such "elementary problems"! Perhaps try to devise a T&Q that avoids such probs: like Sol 56!

Another unpleasant feature of "elementary problems" is idiot placement of symbols (476.22-23) Maybe start to move advanced stage!

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

00 : 475.40 : **SN** Getting QATM to "understand" what an o-z problem is, is not so easy!

Even harder is defining "Any type" problems! OZ seems easier. Perhaps defn. is related to this "2 part soln." I had for OZ problems in the ID's/A Mouth: I vaguely remember: 2 functions:  
1. first was inductively evaluated (w. diffy). 2. second was an analytic function of first (easy to evaluate w. z. M.C. Carlo - an optm. rather than OZ problem), & first obtained Pz edh. in doing the Mixed Corpus Form (19.98).

1. First the first function is: Given an O.T.  $\alpha (Pr)$ ; ( $\alpha$  is name of O.T.,  $Pr$  is the problem den.)

$Pr$  could ~~also~~ be  $M(\cdot)$  where  $M$  is the prob. to find  $x \rightarrow M(x)$  is as large as possible.)

and an problem o-z problem,  $Pr, C \in [Pr \text{ is an } \text{optm. problem, } C \text{ is its } C.B.]$

First and data on  $Pr$ 's work on various Optm. problems,  $(\alpha Pr, \alpha C, \alpha M)$  ( $M \in \text{Cor}$ )

to get a def. of  $M$  as a function of  $C$ , for ~~each~~  $\alpha$  for a given  $Pr$ .

Given a bunch of O.T.'s of a specific OZ problem, the second function is to PC that a given O.T. will

give the best soln. to ~~the~~ that OZ problem with the specified C.B.

A poss. Approach: Perhaps insert suitable "Prims" to make the defn. of  $\alpha$  easier!  
(Tho., in Recursive function theory; Picking the max of several variables might be one of the primitives. — Sometimes the max doesn't exist & is perhaps a common source of ("non-stopping"/undefinedness).)

This seems to be an essential problem: But if we want TH to do "SI", it should understand what "Involvement" means! → 480.38

20 **SN** 279.19 Discusses  $F(\cdot)$  (which is a  $R(\cdot)$  w. pc's rather than Yes/No (Boolean)).

Some diffys in  $F(\cdot)$ 's 'rng'. Also Good Genl. discn of  $F(\cdot)$ 's "updating".

22 : 475.40 In 'rng. "3+4", say; The fact that this is ~~an~~ indep of horizontal motion is built in as human heir.

23 So I have to either build in recognizers indep of hor. motion or give examples w. Various from displacements.

24 List some poss. TSCQ's to start.

1) ANL: see 474.40 ff: for some ideas: Perhaps use  $sub(a, b, c)$  as prim. Other discns. 476.22

2) Sol 56

26 3) 475.16 TSC from ordinary Alg text! (475.16): One poss. advantage: That it may be human-innate. → 31

4) Romanogon's text book, 475.23

30 A poss. advantage of 2.6R: It will have lots of concs missing, but rather than Prims I must know, I can insert them as "primitives". My discovering what tho's/concs might be seems very imp.!

So write out solns in "English": see what concs. seem to be needed. Do a fair amount of TSCQ in "English" before attempting formalization. I may be able to get code later (PC Estimate from "English" concs)

Pretend that I am explaining Prims to a neonate!

Might be (possible/feasible/good idea) to do for ("English") w. an ANL & some following ideas keep "NEONATE" in mind!

38 Or: Just try what would seem to be an adequate TSCQ for "ANL + " in computer  
1. c.j's's I & R Prims are accessible to me: try them out at these C.B.'s: see if there are simpler solns. Reform Prims as we go along: Prims could be so C.J's's of subsequent (477.90)

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

00: 476.40. Problems, a great deal!

03

One update here that seems to follow from the ANL work: That one has assigned  $Q_{n+1}$  to  $i$ .  
 $P_i$  doesn't work very well w.  $Q_{n+1}$ : Then modify  $P_i$  so that  $Q_{n+1}$  will "work" (i.e. Answer by PC.)

03: 00 Try very sparse  $T \rightarrow Q$  to start with & look at C.J.'s! Then ~~insert~~ <sup>insert</sup> various examples  
to augment the  $T \rightarrow Q$  & see how the subsequent C.J.'s change.

SN Perhaps get Peter Berman's  $T \rightarrow Q$ ! From Wolfson or

SN for "Report" Give the (tentative) update scheme that I've written up: But make it clear that this is very preliminary:

- That the real update scheme will be devised by studying  $T \rightarrow Q$ 's - seeing how humans seem to solve them, & expressing their heuristics in formal ways. That heuristics seem to work for some kinds of d-problems & probably will not work well for s-problems.

So: Finish up report & send it out! But first, review my present views of  $T \rightarrow Q$ 's! 476.24 - 477.03

But w. more detail on each approach, so I can easily "get back to it". → 478.04

61 SN On my proof of the conv. thm for BASS & for QATM: <sup>(use 98thm + Mix corpus thm, for proof.)</sup> Again some strangeness! It would seem that the expected error is localized by the amount of ~~PEE~~ ~~REA~~ ~~?~~

$$Z = \ln \left( \frac{P(A_i | Q_i)}{P(A_i | Q_i)} \right)$$
 But I must look at the whole proof to analyse this

$Z$  can be  $< 0$ : which I worried about before, & found not relevant!

Another point in the Area: When  $\prod_{i=1}^n P(A_i | Q_i) = \omega$  is small, this means that the expected value of  $Z$  is square over  $i=1/n$  is small. What it says about the future is unclear! So, is it relevant to the Ockham's Razor?

Is it really evidence that one PEM is better than another? Well, it looks like we did an honest eval. of the Pen with honest A Priori

20 21 SN On comparing ALP w. "Old Faithful" (= O.F. = Cross Validation): O.F. gets honest a priori by actually testing the found param. "a priori": Better data is seen. ALP gets the "a priori" by less certainties, but by it doesn't waste data (which O.F. does).

So ALP gives us an honest estimate of ~~the past~~ ~~error~~  $\Delta \ln P(A_i | Q_i)$ : I guess O.F. would also. Apart from "Honesty": what are the estimates of the future? (i.e. ALP: O.F.)

Perhaps the  $\sigma^2$  are really what the question is, in comparing them.

From a practical pt. of view, ALP tells one where to look for good models. O.F. does not.

One could use ~~ALP~~ to generate trial models w. low cost then use O.F. to test them - which is usually how O.F. is done: i.e. low cost models are usually tried first (via Ockham's rule). - But probably w.

O.F. trials are not systematically for low cost models.

In O.F. (and probably also ALP) there is usually the spurious Open Yield (SOY) problem, if one tests  $> 1$  cand. - (which one would presumably do)

In ALP, I guess we only have 1 cand. - so no "SOY" problems - for ~~CB~~ ~~CB~~ ~~CB~~.

For  $CB < \infty$ , we select the  $i$  cand. of by  $P_i$  in the past. - but we average them w. wt.  $OCPC$ 's. I'm not sure how (or if) SOY is relevant.

In QATM, the mean  $(PC / \text{error})$  would seem to not extrapolate well into the future - depends

much on how "similar" the future Q's are to the Q's in the  $T \rightarrow Q$ .

478.00

30

**CONVERGENCE:**

Convergence thm for  $B>85/QATM$ :  $478.00:477.11$

possl. Bugs?

206: 14  
10006: 5014

ABCDEabcde (0.0793 =

**-HAMPTON RESEARCH ASSOCIATES**  
26 Boylston St., Cambridge, Mass. 02138

00: 477.20: Another Q about t. proof of convergence & rate for B>85 & QATM: It would seem that these PC's would depend on: order in which  $\forall$  B>85 items (or Q/A's) occur. - Yet by defn, the true values are order independent. - But t. predicted values are order dependent: a B>85 element will get a much "better" PC value if that B>85 element is evaluated later in t. TSCQ.

04: 477.20 I want to be fairly sure that imposed TMI has acquired all t. needed "little cones", "unconscious cones", "Common Sense" cones, - otherwise it will have lots of trouble later.  
One kind of (partial) evidence that this is o.k., is that t. cjs's seem reasonable - (usually small) for TSCQ's that I expect a human could do easily. Note, hvr, that Skinner's q's are very easy, but tend not to teach ~~some~~ cones needed for harder problems.

10 Well start again on ANL! Pattern  $3+4: 7$   $A = \text{Sum}(R_1, R_2) = \text{Sum}(R_2, R_1)$   
sum, mul, div are operators available but have numerical argts. TM knows No. v.s. Non-Ab.  
TM knows "identity" for nos, also for nouns. Since it is a Boolean funcn, it can be used for control.  
TM could (v.n. + 10) if ~~it~~ "+" is only kind of problem.  
If we introduce  $5 \times 6: 30$  problems,  $\text{Sum}(R_1, R_2)$  will work  $\frac{1}{2}$  of time,  $\text{Mul}()$  " " " " " "

This is still an enormous sum of compressions & t. argts are 32 bit random nos!

The former could insist on 100% accuracy!

18 A reasonable heuristic! Put together all Q's for which "sum" is correct operator.

20 Contrast w  
sum  $u_1 + u_2$  v.s. mul  $u_3 \times u_4$   
Conclusion: If  $R_2 = +$  then  $\text{Sum}(R_1, R_2)$  call this operator on  $R_1, R_2$  "O"  
"  $R_2 = \times$  " mul()

24 This works fine! Next ~~we~~ we give many examples: only  $\text{sub}(R_1, R_2)$  works (not  $\text{sub}(R_2, R_1)$ !)  
But be for 24 is tried, TM tries Q - which does not work on "-"

27 We could get to 24 by indexing Q - common index for all "-" problems. So it would know they all had something in common. (33)

**SN1** I don't have to get TM to solve all t. probs in a TSCQ sequentially: If I can't think of a reasonable best soln, then just give TM some solns, & continue TSCQ.

Periodically review t. "probs that TM couldn't solve" & see if (as a result of my experience) I can find a way to get TM to solve them in a "useful way".

30 **SN2** Better trip, get soln. to cubic, so I can try to ~~express~~ give a soln in terms of 14 variable substitutions. See if ~~the~~ standard soln to Q can be done this way  
Look pp 118-121 in Baril/Melrose.

33 (27) When we notice  $\text{sub}$  sometimes works & other times O, works! Use heuristic of 18 to decide "-" gives "sum"

34 Similarly, we learn "div"  $\rightarrow$  "div" (v.s.  $\rightarrow$  same heuristic.  
T. soln we have now is a complex conditional of  $\div, -, \times, +$ , but it always works.

36 Next look at ~~the~~  $(1+3): 4$ . T. old ops don't work because of dis placement.  $R_1, R_2, R_3$  are not numbers.

A person would notice that  $(1+3)$  looked a lot like "1+3" and perhaps try  $\text{sum}(1,3)$ .

38 Actually, I can give 36/after write  $1+3$  (1.10)

$$\int_{-\infty}^{+\infty} e^{-\frac{x^2}{2\alpha^2}} dx \quad \int_{-\infty}^{+\infty} e^{-\frac{x^2}{2\alpha^2}} dx \quad \int_{-\infty}^{+\infty} e^{-\frac{x^2}{2\alpha^2}} dx$$

new (2)  $\int$  (parallel version)      new (1)      old

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

00: 478.40: At this pt, I mite Got TM to understand "Hor. invariance", as a useful conc. for ( ) and for (ing "Substitution").

Perhaps start out w. TM looking at 3+4 & "noticing" that 3 is "before" & 4 is "after".

03  $\rightarrow$  Another approach would be list a sequence of solns to ~~the~~ probs in the TSC & then try to ~~find~~  $\rightarrow$  devise a suitable TSC so th. solns will be of acceptable cost.

There would be solns to 3+4, 3x2, ..., 3+5 etc; Then (3+4), (5x3) etc:

eventually 3+(8x9), etc... 4x(2-5) - (3+5)x1; Then (3+5)x(9+3) etc.

eventually, nested parens.

T. way I'm using ~~the~~ Parens; I may be able to get no dupy if I only use say "close parens" - ( )

$$(3+5) \times (3+(9+3)) \quad (3+5) \times (3+(9+3)) \rightarrow (3+5) \times 3+(9+3) \leftarrow \text{This is uniquely decodable.}$$

10  $\rightarrow$  Alternative to 03 would be a TSC that I expect TM to be able to inv. - I then write out solns I expect

T. General way to evaluate Expressions that use parens & an only binary func.

Can I find a sub-expression that ~~can~~ <sup>once can</sup> evaluate and evaluate it; re-write ~~the~~ entire expression w. new values; (copy & eval until imposs).  
Can I work this w. unary & binary "ary func" (using suitable notation).  
Perhaps look at Kleene for Notation.

Should TM spend Much time discovering Alg. identities like  $\text{sub}(3,4) = \text{Mul}(-1, (\text{sub}(9,3)))$ ?

Depends on how Useful such activity could be. ABCDE. Initially, finding such identities would be the sort of thing "AM" would do, because they were "interesting". - Would we want TM to

Spend a certain amount amount of its time discovering things like these "identities"?

Just how would these "discoveries" fit into TM's general "QA" problem set?

One way would be to have a "special" particular Q that asks TM to find "interesting things" which it presents as "A".

How to get this response as a desire to get by PC for A, is unclear!

It may be that this "exploratory activity" would be a kind of "look-ahead", because it would not be immediately oriented toward "by PC A's". - (The it would be in the "long run").

Some tentative goals for early TSC's: 1) ANL  $\rightarrow$  soln. of ~~the~~ single linear eq. (including simplification).

In solving these problems: (say a human would do them); Write out, (in English), the hints or heuristics that I would want to give to a student. See whether these are really adequate for the task addressed.

old "Bad" par:

So back to 478.80  $\rightarrow$  .38 : say it (and 1+3) then (4+8).

When I see 7x9; then out (8x3) because 8x3  $\rightarrow$  ?

Trouble is: ~~(1+3)~~ (1+3) : "good soln" If  $R_3 = +$  then sum  $(R_2, R_4)$ .

It really misses the parens!

$\rightarrow$  480.12  
spac

ID

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

Generalized  
min

00:479.40 : Out "UPDATES ALGMS" I listed in "t. Report": (5 Algms)

Once one has "decided" that  $Q_i$  is in class  $i$  &  $\vec{P}_i$  should work: If it doesn't, we want to modify  $\vec{P}_i$  so it works for: <sup>member</sup>  $n$  as well as old members of  $i$  class. If this proves diff., put  $Q_i$  in a new  $i$  class: try to find a special soln. for it, then find a way to distinguish it (via "Q") from being in the former  $i$  class. We want this "distinguishing" Algms to be as short as poss., so it will extrapolate (guess) well.

10 → This "training TM" is very much like training a child. If the kid doesn't seem to get a point we have to do examples, hints, or "balling" to put him on track.

12:479.40 OK: Retournez a nos parents! While 470.33 is a soln, it's not the one we want - it will be useless for work w. parents. If I taught TM how invariance of  $t$  -  $c$  for  $a$  &  $b$ ,  $(3+7)$  would be obtained indep of the presence of parents!

17 Well: try another task: Assume that TM somehow (and 479.12! What does it need? ① T. idea of "evaln. of a sub expression." In Lisp, a sub expression is any string, & any string can be evaluated. Consider any string of length  $> 1$  that evaluates to a no. In Lisp, paren operators are Built in. ANL is essentially "built in". So, in Lisp, I could try to get it to (rn to solve a linear eqn.

20 **SN** **BIG Q**: What were the characteristics of the SOL56 (problem set/TSCQ) that made it easy (or seem easy!) for TM to usefully work on it? - Well, the prim. cons that I was using, seemed to match the problem well - in the sense that the solns. to the problems were (in terms of those cons) very "simple".

In ~~SN~~ ANL perhaps to have "Equality": If 2 things are =, then they are substitutable, giving same evalns of expressions as before.

Ah! If IM ended up at 473.32, it would have to "Backtrack", if we gave more complex probs using parents.

30 D17/01 **SN** Another BIG idea in the General QATM system is the ULTIMATE SI!

The " $Q_{n+1}$ " is all of TM's previous history (Traces + anything else), plus ~~new~~ Ance is  $\vec{O}_{n+1}$ , the new set of system operators. We also have to include the dependence on  $c, b, e, c, c$ . (we have a feedback eq.  $pc \neq cc$ :  $\frac{pc}{cc} = Max$  is one cons. express.)

→ It amounts to a general soln. to the update problem. We want TM to start work on the problem as early as poss., because it's close to a general, non-el soln. to the update problem.

38:476.17 Re: Appln. of QATM to OZ probs: if we could have a TSCQ of OZ's?

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

00 : 480.40 : **THE BASIC IDEA** of T&Q writing is FACTORIZATION: All of knowledge I want TM to be taught directly, I have to put into form of concs, & factored concs. If I cannot factor what I want TM to lrn, I do not understand it well enough, & I must before I can usefully teach it to TM.

Years ago, I had this idea of "Euclidification of Human Logic" — to express all that has to be lrd in a seq. of steps of  $re < \text{a certain threshold}$ . So that's what I have to do. I don't see anyway to avoid this problem — Tho perhaps many of the tricks I've been proposing for T&Q writing, are ways to help in Euclidification.

Many times when I think I understand (procedure/problem...) I really don't, because I am unable to properly formalize "factors"  $\Rightarrow$  status.

A note on Euclid: He had those axioms & postulates that he then used to generate the whole of Geometry. Mathematicians for many yrs ~~did~~ ~~not~~ did this. Much later, it became clear that he had not done this —

That there was a fair amount of vagueness in how to generate proofs, etc.

→ How does "Euclidification" differ from what Lenat is doing in the CYC project? ←  
I really don't know enough about CYC, to say. I suspect that Lenat mainly "teaches" by "telling". v.s. My teaching by examples (usually)

Much of early "AI" rscrch was an attempt to formalize human problem solving. Perhaps I can use this rscrch to formalize it in a learnable way. Say like Sledge's symbolic integration — or perhaps even Moses's Machine's Maple & Mathematica. Tho those last 3 were "proprietary" initial "inert" (Moses's) work was not proprietary — A Paper may be well published on symbolic math manipulation & prob. solving.

In general, a big part of Human "Understanding" of any thing, is the ability of the Human to express the thing in a small no. of concs that are familiar to the Human.

Symbolic Integration could be an easy area to start with. To start with, the student has a way to differentiate all expressions. (To test trial's ... if trials of this sort are made)

Also a lot of func's w. known S's, & some special tricks like linearity, integration by parts, partial fractions & substitution.

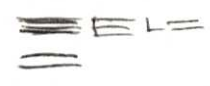
It should be possl. to write out "pems" in English, w.o. using LISP: An advantage of "English" is that even when they don't work (insufficiently formalized), they can usually be modified (not very much) so they do work "better" (more formalized).

I can probably write a pem in English (or Logo Basic), but can do a fair amount of symbolic integration. → Next: How to get TM to lrd stas by "Lang". This is v to my "state" at 480.17

$$\int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} e^{-\frac{x^2}{2a^2}} \int_{-\infty}^{+\infty} e^{-\frac{y^2}{2a^2}} dx$$

raucous

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138



00 : 481.40: Consider an earlier "state of knowledge"?

TM can, for a certain set of functions,  $f_i$ , work to problem  $S f_i = ?$

We can teach it to linearity of "S"  $\rightarrow S(a f_i + b f_j) = a S f_i + b S f_j$ .

So it could integrate any sum of coded functs, if it could integrate f. functs.

A meta rule: try to put f. into standard form, if poss. (a sub-goal).

say TM "knows"  $\int n f(x) = n \int f(x)$  if n is any var & f(x) is any func.

How could it learn this? say it knows  $\int \sin(x) = -\cos(x)$  to teach it that  $\int n \sin(x) = -n \cos(x)$

08 Say we give it many examples of  $\int n \sin(x) = -n \cos(x)$  is

After it runs .08 would it learn  $\int n \exp(x) = n \int \exp(x)$ ?

10 Perhaps give it examples of  $\int n f(x) = n \int f(x)$  w. various n & f. would it learn general rule?

Note: Differentiability Learning Differentiator is about same as / rule Integrator - Except that there are fewer rules! (Also they always converge!)

14 Note: In R.C.C.'s way of defining rec. functs in LISP: mechanism for single recursions, we can always start w. known  $f(x_0) = h_0$  & generate  $f(y) = h(y)$  & do results!

$f^0(x_0) = x$ ; then use  $f(y) = h^1(y_0)$ . Since R is may never converge (i.e. no n exists  $\rightarrow g^n(x_0) = y$ ) which is a nice property of rec. functs. - The Turing's Forward

19 Way of implementing recursion is factor & uses much less memory! (derivative)

20 Def: If Q is of form  $\text{diff}(f_1(x) + f_2(x))$  then A is  $\text{diff} f_1(x) + \text{diff} f_2(x)$ . We have a more complicated "Recognition rule" for Q. We assume TM already knows "deriv" of a function & it would "consider" such a recogn rule.

Note that recogn. rules can be rather complex, so that determining if a given Q satisfies them, can be difficult. e.g.  $\text{diff}(f_1 + f_2)$  is of form  $f_1 + f_2$  with  $g_1 = (f_1 + f_2)$ , is relatively simple, but we can have diff. cases of complicated nesting.

Also TM may have to manipulate Q a lot to get it into a form that is recognized.

Since there are many "recogn. rules", it is not always clear as to what x/y rule should be used.

Some "Recogn rules" imply an  $\infty$  of "easily applied" recogn rules. e.g. a recursive recogn rule might know this way. (33) might be a case.

One subgoal in induction is to find (or generate!) a recogn rule for Q.

Can an induction rule generate R results & its associated P functs.

From: Perry ideas, it may be poss. to see how induction works for a differentiation & for integration.

33 (29) A very weird recogn rule: "If Q is of form  $\int n f(x)$  and  $\text{diff} f_2(x) = f_1(x)$  then A =  $f_2(x)$ " By trying all functions we can "recursively" define  $\text{diff}^{-1} f(x) (= \int f(x))$

D19-01 Consider differentiator! I look at  $x^2 + 1 + e^x + \sin x$ . I notice that

it consists of functions separated by "+" signs. This strongly suggests

"AND" net (Divide & Conquer) - Lets try to factor this: One factor that would be otherwise useful: a fast recognition of that it rings bell.



HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138



00:482.42: T. "+" signs were "functions" & could (perhaps) be differentiated (recursively).

So to main soln.  $\rightarrow$  "Differentiate":  $d(A+B) \rightarrow dA+dB$  becomes  $\geq$  recursive call.

02 ~LISP: If  $X = X^n$  OR  $\sin X$  OR  $\cos X$  then  $nX^{n-1}$ ,  $\cos X$ ,  $-\sin X$ ,  $d\alpha + d\beta$ , resp

More clearly had a neater way of doing this.

Use "CASE"? No - cases ~~on~~ on  $n$  go to? No...

The info content of 02 is: ordered list of If, Then pairs.

So for factorial  $(n) = \text{fact}(n)$ : If  $X=1 \rightarrow 1$  else  $X \cdot \text{fact}(X-1)$

04 If  $\alpha = \sin X$  Then  $\cos X$  else If  $\alpha = \cos X$  Then  $-\sin X$  else  $\alpha = U+V$  Then  $dU+dV$ .  $\rightarrow$  485.02

[SN] On Recursive functs in AZ 141:

In AZ 141: How to define  $f(x) = a$  if  $x=b$ ? (i.e.  $f(b)=a$ ). Well, there's pair special input function;  $f$  inputs  $x, a, b$ : when  $x=a$ , output  $b$ ; undefined otherwise OR inputs  $x, a, b, f$  (to boundary fact it can be recursive)

So: This point has 4 args, one of which can be itself or a func of itself

19  $f$  is being defined.  $f$ 's args are  $x, a, b, f(g(x))$  or  $h(f(g(x)))$

20 So  $f$  is defined by  $a, b, h(\cdot) \& g(\cdot)$ . It differs from normal functions in

that  $h$  &  $g$  are function names; they are args (argt. functs, — or, any funct, that has 1 "dummy" argt.

23  $\rightarrow$  Note these recursive functs cannot be very complicated; their PC  $\downarrow$  very fast w. size.

[2001] These dummy variables (23) normally occur in defn of non-recursive funcs as well.

30 An apparent distance between rec. & non-rec functs. It is poss. to create a "1 use" non-recursive func. (a composition of previously defined funcs) w/o using dummy vars. My present method does not seem able to do this w. rec. funcs.

To use a rec. func., one has to define it first, then use it.

If the  $g$  &  $h$  funcs of 19 & have been previously defined as funcs of 1 variable, the rec. func. defn. needs no new dummy vars.

Look at 482.14-19: we need to find  $n \ni g^n(x_0) = x$ ; Then  $f(x) = h^n(y_0)$  (we introduce time (particularly if  $h(\cdot)$  is expansive) if we first find  $n$  in  $g^n(x_0) = x \rightarrow$  484.00

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

00:33:40 : If we can't find  $n$ , we don't have to waste time finding  $h^n(y_0)$ .  
However, in Lisp, if we don't have time to find  $h^n(y_0)$ , we are "nowhere" — so we  
might as well not have done the calculations! — so for Lisp, this is no "shortcut".  
So, in Lisp, we will alternate between evaluating  $g^n(x_0) \approx h^n(y_0)$ .

Still, the defn of rec function AZ (115) is "kludgy" & I don't feel comfortable w. it.  
It is too much of a "special case"

Actually, we can define rec. functs somewhat differently. In cases, we use  
 $g^n(x_0) = x$  to determine  $n$ , in another we use  $g^{-n}(x) = x_0$  to determine  $n$ .

Since computer time for  $g^{-1}$  &  $g$  can be quite different, this could make a lot of difference!  
Sometimes we may not know how to compute both  $g(x) \approx g^{-1}(x)$ . — Also one may be  
multiple valued! If so, we would save lots of cc by using the other (inverse  $\rightarrow$  find,  
find  $\rightarrow$  inverse),

How, for pc evaln of rec. functs, all we really need is to input output!

i.e.  $x_0, f(x_0) \approx y_0, g(x)$  and  $h(x)$ . If  $g(x)$  and  $h(x)$  are not

normally functs of 1 var, & let, they will have to be suitably (indexed/labelled) —  
telling which args are to dummy & which are params calculated externally.

Of some interest:

I have considered only 2 ways to define functions:

Composition & Recursion; there are at least another:  $f_n(x) = g^n(x)$ . } The first can also be done by recursion  
also  $g(x) = f^n(x)$  (  $n < 0$  for integrals)

Also functs defined by finite or infinite series.

On my PDP-11, if  $f$  defined a Rec funct: It may not fit "X!".

$f(x+1) = x f(x)$  (  $f(x) = x f(x-1) = h(x, f(g(x)))$  )  $f(x-1) = (x-1) f(x-2)$   
In Lisp,  $f(x) \approx$  If  $x$  is  $\leq 0$  then  $h(x, f(g(x)))$  }  $f(g(x)) = g(x) f(g^2(x)) \rightarrow$

$f(x) = h(x, f(g(x)))$  ;  $f(g(x)) = h(g(x), f(g^2(x)))$

$f(x) = h(x, h(g(x), h(g^2(x), f(g^3(x))))))$

$f(x+1) = (x+1) f(x)$  ;  $f(1) = 1$  ;  $f(2) = 2$   $f(x) = h(x, f(g(x)))$

$h^{-1}(x, f(x)) = f(g(x))$  so  $(h^{-1})^n(x, f(x)) = f(g^n(x))$  since  $n \rightarrow g^n(x) = 2$  then  $f(2) = 2$

so give  $h(x, y)$  to Best look at leave the defn. of recursive functs

How, it is clear that there will be no great trouble in seeing how much

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

11, 31, 41, 61, 71, 91, 101, 111, 121, 131, 141, 151, 161, 171, 181, 191

Info (pc) there is

00: 484.90 : Info (pc) Process in a rec. data

So: Get back to this when we get back to Comby. Looks very non-coherent

02: 483.09 : I was trying to get differentiation rules in "Heaviside" format.

Consider following "state of relevant knowledge":

TM knows  $\text{dir}(\sin x) = \cos x$ ,  $\text{dir}(x^2) = 2x$ ;  $\text{dir}(f_1(x) + f_2(x)) = \text{dir} f_1(x) + \text{dir} f_2(x)$ . How did it know this (ask: Also, what form did this knowledge take? For  $f_1$  first  $z$ , it could have been given!

$Q = \text{dir}(x^2)$ ;  $A = 2x$ ;  $Q = \text{dir}(\sin x)$ ,  $A = \cos x$ . It was found that  $x^2$  is  $\sin x$  over functions!

$Q = \text{case}(\text{function}, x^2)$ ;  $A = T$ .

$Q = \text{case}(\text{function}, \sin(x))$ ;  $A = T$ .

perhaps  $\text{dir } 3 \cos \rightarrow 3 \text{ dir } \cos$  would be easier to have starting pts.

3 is "no", "cos" is function. so  $\text{dir}(\text{no.} \times \text{function}) \rightarrow \text{no.} \times (\text{dir function})$

Look at them Polish:  $\text{dir}(x(3 \cos x)) \rightarrow x(3 \text{ dir } \cos x)$

The conc. of "exchange" or "commute" could be useful. It works nicely if there are no parens.

NB In general, when TM sees a sub object, it classifies it if possible. Sometimes there will be ambiguous parsing of  $Q$ , so  $> 1$  method of classifying is partly,

If parens are used, then no parsing may be impossible.

At first I will teach TM several deltas, so it will then use these "classifiers" in induction

Later, it should be able to invent its own "classifiers" ( $\equiv$  "terminology") for induction.

One main goal of this "classifier" is for "Recognition" part of  $Q$  processing.

One early application might be ANL:  $4 \times 5 \rightarrow \text{mul}(4, 5)$  is immediately equated (tested)

as  $no_1 \times no_2 \rightarrow \text{mul}(no_1, no_2)$   $no_1$  "operation"  $no_2 \rightarrow \text{mul}(no_1, no_2)$  doesn't work. O.K.

writing  $3 \times 4 \rightarrow 12$  is bad, here. "x" should be written like any other binary operation.

Some  $\text{dir } f(x)$

$$\text{mul } a, \text{ dir } f(x) = \text{dir } \text{mul } a(f(x))$$

Notation "F(x)" means "f. value of f for argt value x"

What I want is micro "f"  $f(x)$  the function itself not any particular value of it.

AZ/AT can (i.e. does) define such objects, using dummy variables.

$f \text{ sin}(x) = F(\text{sin}(x), x)$  - i.e. it is a function of function  $\text{sin}(x)$  at  $x$ .

Also note: B in ANL: I want "eval( $z+t$ )"

$\Rightarrow$  If  $a=b$  then means that  $\text{eval}(a)$  is identical to  $\text{eval}(b)$ .

$\Rightarrow$  I can do Algebra 2 ways (a) slow way the hard way to do; (b) rigorous way.

It is getting Algebra rigorous way.

Externally, sloppy Algebra, learning by humans; perhaps easier to derive test's, but the concepts used are related to R, i.e. obscure human needs/wants.

REV

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

Some very imp't TH ideas:

$\rightarrow \equiv (\text{exp of } 0^1, \text{ exp of } 0^4)$

1) 414, ... eq. What PATH is trying to maximize.  $\rightarrow$  T's assoc. S errz convergence param, & what it means for QA TSP's (i.e. future QA's could be in a new direction so we can't extrapolate much!  $\rightarrow$  \$\$\$\$998.00)

2) The cc modulus of 1: Note (2) (.25)  $\rightarrow$  (S, f, d)  $\rightarrow$  (S, f, d) (S, f, d) (S, f, d)

3) The conc. nodes a basis of TSP construction (see S) (S, f, d) (S, f, d) (S, f, d)

4) T. idea that "all info in P.D." Level is optimum; New P.D. must be directed (George Lorch, A Prize to be true (Critical P.N. in Sol 89))

5) AZ141 method of assigning p's to functions based on composition & recursion (recursion: 483.10-A)

6) (Maybe) factorization of  $0^n$  (in 1) into Recognition & P parts.

7) In 1) to express. of Publishable func by 3 input-umen & recursion that we could implement via AZ141-type functions (i.e. LISP-like func)  $\rightarrow$  Chaitin's idea of "request-infinite bit" as an important input v.m.

8) That stochastic (as well as Deform.) func could be found by the TSP's guided by Conc nets  $\rightarrow$  (I may need to reclarify this in my mind). This results, ultimately in what seems to be optimum update system. T. applied this same use also  $\rightarrow$  (I may need to reclarify this in my mind). This results, ultimately in what seems to be optimum update system. T. applied this same use also

9) The distance between a pair of func used in a pair of trials (based on extended "interact") of trials  $\rightarrow$  (I may need to reclarify this in my mind). This results, ultimately in what seems to be optimum update system. T. applied this same use also

10) Long Mat. Lang. by first-ry Algebra (or some other domain) then, reason to  $\rightarrow$  (I may need to reclarify this in my mind). This results, ultimately in what seems to be optimum update system. T. applied this same use also

11) TM can use logical reasoning! This can be taught a/o "wired in" (see 487.32, Also 486.13-15; 486.16-23)  $\rightarrow$  (I may need to reclarify this in my mind). This results, ultimately in what seems to be optimum update system. T. applied this same use also

12) T. idea that a P.D. has also a cc param. For each CB, we can have a default P.D. - P.D.'s to be available for the argt. we cc < to. Given CB. We can also have the same CB for all argts, or have a default specified CB for each argt.

Science  
Math, at  
the  
level  
of  
the  
brain  
is  
possible  
to  
be  
taught  
by  
the  
trainer  
knows  
-  
possibly  
even  
"cheating"  
488.12  
1990-15 is a  
good disc.

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

00 (Spec 485.90): Advantage of Rigorous approach: Its probly better, simpler solns.  
" " Sloppy or: It uses Human Hours but may be more powerful than those of (oo)'.  
For beginning: oo' may be simpler. — Too getting Algebra into rigorous form may be diff.

In case of  $\sin(x)$  Powers analogy: it could mean value of  $\sin$  for each  $x$  or to function  $\sin(x)$  as dummy variable  $x$ . Usually this is easy to resolve by "cutback".

If dir  $\sin(x)$  is used for 1 woman  $x$  to be a dummy; but we end up  $\cos(x)$  — which  $\cos(x)$  can be a function &  $x$  can be a dummy or a actual arg.

I think John Moses may have run up against problems like this in developing "Maesyma". Perhaps look at MapleV notation for problems to see how these

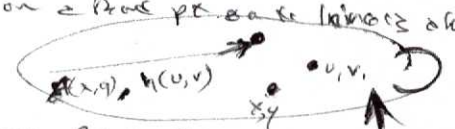
10 alg. notations can be ambiguous. Tho it may well be due to techniques they use or "kludge" because they try to stay close to human ways of processing things — or just historical args. (32)

14: 484.38 on rec. func. defns: This method 482.14-19 find  $h \Rightarrow g(a) = b$ : then  $f(x) = h^{(n)}(x)$   
This can be implemented by simple Do loop. However, Do loops can also use  $\alpha > \beta$  as a step criterion (nearly  $\alpha = \beta$ ) so they are more powerful. The equivalence of certain kinds of rec. func. (esp from rec. func.) to certain kinds of Do loops may be perhaps shown from this kind of Arg. Maybe Ross Peters used proof of this kind.

20 Anyway, I think most (if not all) rec. function routines, can be defined by Do loops. In terms of rec. func. defined by rules like  $f(x, g(x)) = h(x, f(x, g(x)))$

value of  $f(x)$  at  $g(x)$  and  $h(x)$  determines  $f(x)$  at  $h(x)$

$f(g(x)) = h(f(g(x)), f(h(x)))$ . So given  $f$  at 2 args of  $z$  (and, a value on a prod pt. as a linkage always given. (Hm, this definition not be casual



perhaps  $f(x, y) = f(x, g(y))$ :  $f(g(x), h(y)) = r(f(x, y), f(u, v))$  (not)

$f(g(x), h(y)) = r(f(g(x, y), h(y, v))) = r(f(x, y), f(u, v))$

Look at defn of Ackermann func. — then try to express it as nested "Do" loops.

32 13 In much heuristic reasoning, logic is used — In simulating these heuristics, this logical reasoning should be used. We can assume TM knows how to do this, & we can (presumably) teach a human to do this (or write it in); (see 486.13-15 & 486.16-22) R.

Any means we use to get TM to  $f$  or  $d$  of 419... is "legal".

38 This (32) suggests that <sup>Early</sup> work on A.I. (Heuristic Reasoning, etc.) could be useful if developed in a probabilistic way.

Hybrid Vigour .00

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

001. SN HYBRID VIGOUR <sup>GOR?</sup> Explained as <sup>(a)</sup> cross of 2 function trees that are given different <sup>(b)</sup> cross  
Solutions of General problem. As opposed to 2 Solns that are very near. In (a), possy of

002. new, interesting trials. In (b) much less promising "new ideas". <sup>mix of</sup> <sup>Probably that Earth found this way!</sup> <sup>since all planets are in ecliptic</sup>  
Possb Hybrid Vigour (for more likely production of Life) by collision of  
Galaxies found in much different ways. (for Gunkel). Vukobrat's "Worlds a Collision."

Also  $\therefore$  try to save children w. genetic diseases that seem a recurring <sup>Genetic</sup> <sup>Diis. Best for "Bad" diseases that</sup> <sup>Involves several genes! Diffic</sup> <sup>to pin down Genetic Origin.</sup>  
Cross of <sup>Genetic</sup> Chimp and man? Seems likely this would have occurred.  
Hybridizing should occur at near optimum ~~code~~ "level of burn of cards."  
So mix is of "almost random" strings

So the idea of .00-.02 does seem v.g. Hybrids give large ave, but they do give real  
possy of very superior offspring.

12 (486.22R): IMPTZ Move details on idea (486.13 ff) [see 499.00-15 for good discn. there is probably much previous Scott written on this]  $\leftarrow$  IMPT  
It is clear that 1 corp; for 2 ~~is~~ d.f.3 are (usually) pub default. (New) note that in t. second case,  
t. date can be "ordered" - that TM could recognize that more recent attempts to be Maxen, are  
better. Also, in the second case, TM (ideally) has to record what "Maxen" means. It might be  
expressed as a rec. defn, but is more easily implemented as a loop. I.e. "pick max of pass trials  
up to  $500 = c.b.$ " (Skill, this is not a whole complete Maxen)  $\rightarrow$  492.05

20  
21  
22  
3 approaches to TSO design:  
(1) (Formal): write examples <sup>(e.g. 1, 0)</sup> & Solns in Formal Math. Do this for a fair section A B C D  
of Algebra. Then look at solns. See if CJS's are too large! If so, look for linking concs;  
Draw conc net. - Try to find examples of conc to bridge to larger CJS's.

(2) (Informal) In English: write examples  $\&$  Solns (rules) for some corpus  
as Q. Make conc net, then try to formalize solns (Q) as in (1). Advantage  
English & Part of automatically give "soft" solns, that suggest good formal solns.

(3) 497.32-40 (~~487.28~~  $\rightarrow$  40 in particular) Used Old AI heuristics  
techniques + logic. should work.

One of the above (2) is a super case of (1) so should work! Carano.  
So (2 or 3)  $\leq$  should work.

32  
33  
36  
SN Leont's "AM" was run in 1973 28 yrs ago:  $2^{28} (1.5)^{28}$   $\approx$  400k so we see new 400k bugs as fast!  
Hur, even if we used good such methods, identifying a useful conjecture would be hard for a human -  
But if AM says its a v.g. conjecture <sup>and gives reasons,</sup> that would narrow things down considerably!  
Perhaps try to read the Hans Cyrano: perhaps improve it & run it on fast machine. 499.18  
Get it on to Net. If Brian his Reason is out, not. (Dirty finding it or even Red to learn Hans! Try kenneth Hans.)  
I'm interested in structural details of "AM": so I can try to port to do its thing using Prologistic Such. - Also have a  
goal-oriented part that looks for new heur. NP I did regard AM as kind of Vertex of S89. 449.00

JBC

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

00 : 448.40

! So try (488.21-31):

First: Prelim trials for notation:  $n_1 + n_2$  means trials w.  $n_1$  &  $n_2$  random integers

"ev ( $n_1 + n_2$ )" is +. Evaluate  $n_1 + n_2$  - it will normally be a number (integer).

$n_1 + n_2$  : ev means  $Q$  is  $n_1 + n_2$ ;  $A$  is ev ( $n_1 + n_2$ ).

For most/early I just have to write  $Q$ .  $A$  will be "eval"  $Q$ .

To start off, + numbers & non-numbers are clearly indexed (associated) on + symbol itself! So if TM is given

t.  $Q$  "3 + 4"; "3" has an associated byte, that tells (among other things) that it is a number.

"+" also has an assoc. byte, that tells it is a "non-no."

T.M. is taught that objects having the property "number" has the property "no." & TM puts this into its "property list" (?). (The "property list" is deterministic: I'd want a pc for a

property, that have been evaluated this far. If certain pc's are used by several objects, we would want to update them & store them in a commonly accessible Register -

T. Format Q:

T. Six digit chars. are left justified. All characters are of significance,

Perhaps for each Q string, there is an "end" indicator - so Q's form a prefix set.

→ A Built-in Human Heur: In Q, nearby symbols have more likely hood to be related.

18: (488.36)

**KENNETH HAASE**: <http://haase.www.media.mit.edu/~haase>.

Contains his thesis, **Cyrano**. Look at "Kenneth Haase on Google".

Thesis on HTML: Originally in Latex. Perhaps got Latex from Haase? HTML file is many "Nodes". 86 "nodes" - 20k by/node (in HTML). have to be in dividually down loaded.

Perhaps Get Ken interested in getting  $\epsilon$  "Cyrano" to work on QATM

Haase (in "worthy of a...") mentions **Wei-Min Shen**. **W. M. Shen!** At his web p. are refs of interest: But read Haase on HIS work: it is a reproduction of ATM I think.

→ Another "common sense" idea humans have, is idea of "counting" - T. Cardinality of a (finite) set. That's "inv" of a set is indep. of order of elements in set.

→ Haase's discuss. of limitations of "Cyrano", gives no sense idea of what it does.

He says Cyr. does new things & uses old things to make new, "But does it do invention or 'Modus' of previous devices. - It would seem "Probing language" is very limited ...! - but is it housing full "Lisp"? At any rate, I should look at T. limitations of "Cyrano", & make sure that QATM is not limited in those ways (or any other ways that Haase may suggest)

Also, I want to look at **KEENE**, to be sure AZKI is truly "universal".

One way to do 488.21-31: Just write a naive TSP & see where TM gets "stuck":

Then try to (invent/discover) (conc./heurs) what TM would need to overcome these difficulties:

On "substitution": " $a = b$ " iff many expression equations involving true after sub  $a, b, eq$

(or, limited versions of "subs" ("subs" can take as: n's occurrence of " $a$ "; " $eq$ ") → 495.34

Another approach to TSP's: use Polish or LISP as language of systems, so it's always clear how to parse Q. (But just how is this different from previous stuff? T.M. trials

are always in "Lisp". T. difficulty that inf. lang., Q can have any form.

(Spec.) → 490.00

↑ Shows form may have gone beyond ATM: but look at Haase's comments.

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

00 (489.40) What is being proposed is that 1. Early TSP be in easy decodable form - easily  
parsed. So, given  $(\text{sum}(\text{mul } 3, 12), 1)$  TM would be able to do it easily.

A problem might be:  $\text{solve } x, (\text{sum}(3, x) = 10)$  : to solve f. eq. for X.

In 01 Eval (sum, (mul 10, 8), 3) would be easy to do as a Q: w/o eval, meaning is unclear. In Lisp it is clear; but otherwise one would have to write "eval" or "Quote" before each expression. In Lisp, if "Quote" does not occur, Eval is assumed. One could have Quote be assumed unless eval is given - which is in fact more reasonable, but probably involves more typing (since "Eval" is much commoner than "Quote").

The I'm not sure "Quote" default would work! : Szq wanted Q, Ev, S.

In Lisp, its QS,  $\rightarrow$  S. ; In Q default, it would be Ev S, which amounts to same thing.

**SN**: Re: Infix ("normal") Alg notation:

Also assoc of  $\text{sum}$  and  $\text{mul}$  sub strings & mul div strings. Few functions associate; few commute - but neither + or x notation in infix helps in this. + & x both commute & associate, so a set of "+" notation form a "bag". In general, any string of + - can be converted to "all +" by changing sign of "-" terms. Similarly w. x & ÷:

There may be tricks in Polish!

or  $\text{sum}(\text{sum}(\text{sum } a, b), c), d$  There may be simple rules to change order of symbols in this thing.

$\text{sum}(c, \text{sum}(a, b))$   
 $\text{sum } c \text{ sum } (c, \text{sum } a, b)$   
 $\text{sum } c \text{ sum } a, c, \text{sum } a, b$   
 $\text{sum } \text{ sum } c$

Well, say TM knows Polish notation. It is also familiar w. + & binary operations.

$\text{sum}(0, x)$  &  $\text{div}(1, x)$  are commonly used binary ops & are given special names, as unary ops.

N.B. whenever TM defines a new function, its "arity" appears in its "property list" of that function.

**SS** Do a quick try for linear Eqn. It knows in Lisp to start!

Q:  $\text{sum}(\text{sum}(3, x))$ :  $A-3$  solve  $f(x)$  means: "solve  $f(x)=0$ ".

$\text{sum}(3, x) = 0$  ;  $\text{sum}(x, 3) = 0$   $\left( \begin{matrix} = x \\ \text{sub}(\text{sum } x, 3) \end{matrix} \right) \left. \begin{matrix} = -3 \\ \text{sub } 0, 3 = -3 \end{matrix} \right\}$

Example of "put" as much as is known on RHS. Simplify LHS then put name on RHS then sum LHS. Get. Try to simplify LHS by div. or sub.

Concept of "Simplicity" (may be built into System) : (Not best - Must solve at the probs use H.E. or OZ methods) - A heuristic: "Simpler" (eqs/probs) are usually easier to solve.

More directly:  $\text{solve}(\text{sum}(n_1, x)) \rightarrow A = \text{sub } \phi, n_1$  ;  $\text{sum}(\text{sub } \phi, n_1) = \text{neg}(n_1)$

$\text{sum}(n_1, \text{sub } \phi, n_1) = \phi$ . For any  $n_1$ :  $n_1(\text{sub } \phi, n_1) = \phi$  TM could find this equality  $\text{sum } n_1 \text{ neg } n_1$



**GA .00**

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

.00: 490-40: **SN** If GA can work "hard probs", it should be able to work on SI. ~~IMPOSSIBLE~~  
Superficially this may seem impractical, because each "trial" ~~is~~ for this Meta System, takes  
an enormous amt. of time. Hvr.: Perhaps Not! There are certain Aspects of S.I. that  
get much faster F.B.: They are:

- .05 1) Mutation & Cross methods as a funct of problem type.
- .06 2) Initial population, & pop size " " " " " "

.07 In 1): We have a bunch of Mut & Cross methods, & we try them out occasionally, for every problem.  
If they seem good, we try them some more: We can ~~control~~ <sup>test</sup> them in many  $\Delta G$  / generation  
(or  $\Delta G$  / partial generation), we ~~can~~ also ~~test~~ <sup>test</sup> Mut & Cross of  $\Delta G$  Mut & Cross....

Sample sizes, (run sizes) for these Experiments ~~can~~ should be determined by Sequential Analysis (Wald).

.10 In 2) we can save parts (or even whole) populations (or models of populations) of all  
probs. & GA has ever solved. We can then sample these populations at various  $\Delta G$  or  
levels  $\rightarrow$  & try them in  $\Delta$  present problem.

.13 In both 1) & 2) we look for relationship betw. **initial problem desc.** & successful  
.16 choice ranges in 1) & 2). ~~is~~

.17 In  $\Delta$  problem of .16-.17, we have a single evolving population that continues as  
new problems are given to  $\Delta$  system.

.20 .05, .06 & .16-.17 are  $\Delta$  Main Methods of S.I. needed by GA.

What work?  $\Delta$  Main problem is not  $\Delta$  above, but method of  
Representation of  $\Delta$  problem (including Fitness function).

The " $\Delta$ " (or other monotic fund of  $\Delta$ ) original fitness func, can be continuously "tweaked"  
as 11.07-13, but there are other aspects of F.F. that are also imp. — Many times,  $\Delta$ .  
Fitness func is really a vector ~~in~~  $x_1 \dots x_n$ , & we use  $\sum_{i=1}^n f_i(x_i)$  as a all over  $\Delta$ .  
we can then wt.  $\Delta$   $f_i$ 's new line, which is usually more complex than tweaking one  $\Delta$ .

.30 To "representation" of  $\Delta$  problem is usually intimately linked w.  $\Delta$  defn of  $\Delta$  fitness func.

Koza seems to have found a somewhat General Method of representing problems as  
"Networks": (Mind full of G. Kohn & TSP). I don't yet see how it fits to deal  
w. Nat. Lang. input. **Could we teachm N.L.  $\Delta$  some way & train to teach QATM Net Lang.?**

Anyway, T. way .05, .06, .16-.17 work: we are always taking random samples of  
mut. cross, "initial population" when we find ~~any~~ <sup>one</sup> that are good, we try to correlate  
this into w.  $\Delta$  initial (problem desc.), or w.  $\Delta$  representative (problem desc.) selected by  $\Delta$  user.

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

: T.M. could learn "Facts"; "Generalizations of Facts"; "Generalizations of Generalizations..."

So initially, TM could learn a set of neg. & relations betw. them: & Genz. Pres. relations.  
At first it has a unique special number for each no. Later it could learn binary relations, so it would know about numbers that it had never heard about before.

988.20 D30/01

: If I actually used the P.d. implied by the code for  $O^N$ , is a wrapper for searching for  $O^N$ ,

(This would not use info about  $Q^N$  or  $A^N$ ; so clearly it would be very "Inefficient") — how fast would it be?

If using rather large machine, I could get the system to a pt. (using suitable TSQ) to where it could

usefully work on the problem of Updating  $Q^N, A^N$  i.e.  $Q = O^N, P^N, A^N$ , in conjunction "Contextual info"

$X = O^N$ . I would have it MADE!

In ~~the~~ T. system of .04 ff, the update Algor. is simply Lsch: T. Big problem is devising

the TSQ.

So if I just wrote a naive TSQ for TM, how would .04 react?

"too large" can augment the TSQ or "wire in" (otherwise "useful"), needed concs. By looking at code of units (unimplemented) TSQ, (by "code" I include info on how I think I wrote it), one may be able to tell which concs are needed for acceptable CJS.

.04-.16 seems like a reasonable approach to TSQ Design & (perhaps) QATM design.

T. Q's in mathematical form,  $O^N$ 's in English, to start.

In .14-.16, if we need certain Heurs to do CJS: we can hypothesize a TSQ that would give those Heurs. These Heurs are "Metaniles". In .13-.16 I was not

considering using Heurs of this kind — but it would seem that they could be usefully employed!

If we do use Heurs of this kind, I guess we are back to the S89! If such Heurs

are not a departure for a "complete" TM, we can try to find out if .04-.09 to get it.

Optimum TM!

ABICQ: Using Lsch ~~to~~ to find a good  $O^N$  would seem to involve testing all previous Corpus — which would make Lsch <sup>not</sup> practical. (In normal induction by humans, this is somehow avoided — perhaps by breaking up corpus into parts by the "Recepti" function. Anyway, in doing .04-.25, try to discover just how 26 is done.)

SN Use of Symbolic Integration as TSQ

If System Got to be very Creative, it would

do so but as a "free space" integration.

SN In APL, try to express trig. funts. as funts. of 2 args: 1. first arg is an integer that gives "name" of func. T. result is (I think) that the identity of trig. taken any simplifications.

$$\tan = \frac{\sin}{\cos}; \cot = \frac{1}{\tan} = \frac{\cos}{\sin}; \sec = \frac{1}{\cos}; \csc = \frac{1}{\sin}; \sec^2 = \tan^2 + 1$$
$$\tan^2 + \sec^2 = \frac{\sin^2}{\cos^2} + \frac{1}{\cos^2} = \frac{\sin^2 + 1}{\cos^2} = \frac{\sin^2 + \sin^2 + \cos^2}{\cos^2} = \frac{2\sin^2 + \cos^2}{\cos^2}$$
$$\sec^2 + \csc^2 = \frac{1}{\cos^2} + \frac{1}{\sin^2} = \frac{\sin^2 + \cos^2}{\sin^2 \cos^2} = \frac{1}{\sin^2 \cos^2}$$

1/1/02

±D

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

Re: MATHEMATICAL, LOGIC reasoning: TM will have to learn most of this same way it (lives) of other things. —  
Via Suitable TSO's. In the Elementary Algebra TSO, there is much Elementary Logical Reasoning  
that, if used, will ↑ pc of result enormously. So it's very desirable that TM ~~know~~ know this. —

How much should be used in is how much (and, is) unclear at present

Some "reasoning": Q:  $n_1 + n_2$  A:  $evl(n_1 + n_2)$ :

~~the~~  $n_1, n_2$  are nos. A is a w. T. only handles cony from nos to nos via ~~sum, sub, mul, div.~~

So we eventually try  $sum\ n_1, n_2$  &  $sum\ n_2, n_1$ ; we also tried  $\begin{pmatrix} mul \\ sum \\ sub \\ n \end{pmatrix} \begin{pmatrix} n_2 \\ 1 \\ 1 \end{pmatrix} \leftarrow \begin{pmatrix} div \\ sum \\ 1 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \\ n_1 \end{pmatrix}$   
So  $sum\ n_1, n_2$  doesn't have an enormous pc. Look at 32 so  $\approx pc = \frac{1}{16}$ .  $4 \cdot 2 \cdot 2 + 4 \cdot 2 \cdot 2 = 32$

to ~~run~~  $Enc(n_1) \rightarrow sum\ n_1, 1$ ;  $dec(n_1) \rightarrow sub(n_1, 1)$  ( $\frac{1}{2}$  pc of  $Enc$ )

Now TM has to distinguish betw  $+$  &  $-$   $Enc$  problems &  $dec$  problems.

To distinguish betw those by use of ~~the~~ problems: T. presence of  $+$   $-$   $\times$   $\div$   $Enc$ ,  $dec$ . is adequate.

Is TM ready for  $n_1 + (n_2 \times n_3)$ ?

I'd like TM to realize that  $(n_2 \times n_3)$  can be subs by

$evl\ n_2 \times n_3$ , to obtain "simpler" expression. This can be done (for nested expressions): either way:

- 1) replace directly evaluable expressions by their values. Recursive
- 2) Take longest parent expression.

Its evalns obtained by doing  $\beta$  & recursive defn, using stack: (1) (16) doesn't need stack

(SN) on "Logical reasoning":

I assumed that LOGIC reasoning runs usually close to  $pc = 1$  — so  
to zero cost — is most of pc was used relatively large (very little such needed). Hrr, it may  
well be that "log. rsng" does require much such, so one doesn't really save on cc by  
using it to "narrow down poss. rands."  $\rightarrow$  Hrr, usually breaking a task in to 2 parts like Peir's:  
 $\rightarrow$  Logical reasoning part (to narrow things down) followed by a simple L such (or ~~more~~ exhaustive  
voter such), is better.

One way to get a recursion in .13, is to have TM learn many evalns. w.o. recursion: then  
it looks at all of + solns. & finds a common recursive form for them. One thing that facilitates  
such discovery, is noticing common sub-expressions. For each of the simpler problems, there are  
many solns, & only a few facilitate recursion (the much use of common sub-expressions, would  
narrow down to low level "solns. that were to be used into ratch. w. recursion.

So try .13 also

$$(n_1 + n_2) \times n_3$$
$$mul(sum\ n_1, n_2), n_3$$



$$n_1 + (n_2 \times n_3)$$
$$sum(n_1, mul(n_2, n_3))$$



So the  $n_i$  in paren, are combined; T. results combined  
w. ~~the~~ other  $n_j$

Peir's works in both cases.

The traces of .33 — .34 are in some of these  $\rightarrow$  normal search for functions.

(Spec 494.00)

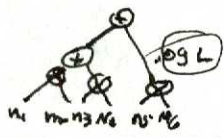
HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

00: 493.40 : Unclear as to how to formalize 493.35-36 : with "The nos. betw. t. parents evaluated t. sigt. of a function". This is in "English": so it has to desired Ambiguity. Perhaps leave it that way, until after, I do t. recursion (in "English"). In 493.35-36 "t. result is combined w. t. other n<sub>j</sub>". — This amounts to substitution.

04 Fry  $(n_1 \times n_2) + (n_3 \times n_4)$  Since there are 4 vars, this can be 2 long such!  
05 A possibl. soln. to 04: DO 493.35-36 — either way  $\begin{pmatrix} n_1, n_2 \text{ first} \\ n_3, n_4 \text{ first} \end{pmatrix}$   
N.B. t. n<sub>i</sub> have to be betw. ( ) in that order.



09-10 Try  $((n_1 \times n_2) + (n_3 \times n_4)) + (n_5 \times n_6)$  — which is in t. spirit of "04. woops! nested 6 parents needed."  
How TM interprets t. nesting" is unclear; It could parse parents by:  
 $(( \quad ) + ( \quad )) + ( \quad )$  This is "wrong" of course, but even by neglecting 2 of t. parents, it does produce answer.



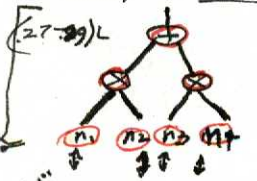
12 So: What is t. problem? To devise a Meta/lang. to desc. functions & "features" of functions,  
So Part a) The non-recursive functs for ANL are easily found.  
14 b) " recursive " " " " " " " " " " " " " "

I also need a way to desc. t. problems "Q" & their "features" & how they may be correlated w. t. features of t. lang. of 12-14. Actually, t. lang. of desc. functions" may be already available. What I need is a way to desc. features/properties of functs w/ their Derivs (trees).

20 This ANL approach (infix to Polish translation) discovery, seems diff. : Perhaps try a diffrent. TSC.

22 Note (27-29) : Uses (29L) as a subgraph. — in fact, usually problems of this sort use previous solns. as subgraphs. — so PC's are quite try! — But still it has to look at t. Q & figure out what graph(function) to use; so I still need a good lang. to desc. Q & its "features/properties".

27 Well, we could draw lines, make correspondences ~~bet~~ between symbols in Q & symbols part of graph. function. So any t. graph edges, i. t. parents or unargnd. & edges: 4 parents! A diff. method will be appropriate. Perhaps regard t. parents as "graph" instructions!



29-30  $(n_1 \times n_2) + (n_3 \times n_4)$   $(n_1 \times n_2) \rightarrow n_1 \ n_2$   $(( \times ) + ( \times )) \rightarrow$   $\rightarrow$

32 Since we can get lots of examples of (graph  $\leftrightarrow$  Q) (22) we should be able to make a correct model.  
34 Perhaps try writing Recursive Soln. that should stem from these examples

32-34 seems like a more or less "well defined problem" so leave it for a while.  
35 Looking at (27-29) & (31L); it seems that one can easily describe an algorithm to go from Q to G (out Q to G)  
To go from G to Q: first write  $n_1 \ n_2 \ n_3 \ n_4 ; \dots$  ?

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

00: 494.90! + 1 494.27 fl: One could use "counting rules".

The Graphical (494.27-.29) can be written as a full expression.

Sum  $(\text{Sum}(n_1, n_2), \text{Sum}(n_3, n_4))$  u.s.  
 $((n_1 + n_2) + (n_3 + n_4))$

so  $n_1, n_2$  recognize as sum  $n_1, n_2$   
 $(\alpha) + (\beta)$  recognize as sum  $\alpha, \beta$  } to be  
Note that rules (a) & (b) are very similar. } rules

Also, to go from a to b, it is probably possible to devise a CF "X(tu Grammar"

I could just drop this problem for time being, since it is certainly solvable:

One poss. reason to work out now: I'd like ideas on how to (recognize/discover) recursive  
algs., & this one seems particularly simple.

Also if I want to continue this/particular TSO, it would be well to know just what form  
this particular "Recursive induction" takes. Actually, I can put it into any form I like as a special  
part of the "Updating Algm".

So: say TM can now evaluate all expressions containing  $+, \times, \div, \dots$

A possible next lesson: To let TM "inc(x)" means  $x+1$ . Also that  $inc(\alpha) = \alpha+1$ , where  
 $\alpha$  is any variable or expression. We can teach this w.o. any explicit "definition" facility.

i.e.  $Q = inc(n_1)$ ;  $A = eval(n_1)$ : Once this is done, how easy would it be

for TM to realize that we can replace any "inc(x)" by "x+1" in any expression?

It would be more likely if TM had already experience w. substitution; so a substitution

was a "primitive", TM might **ANL** much easier; & well as this stuff. (The other

Alternative would be for TM to "discover" substitution "on its own")

Its perhaps difficult for TM to let  $(n_1 + n_2) \times (n_2 + n_1)$  if  $n_1$  &  $n_2$  can only be actual  
nos. One could, perhaps teach this, if addition took less time? So if we could

for  $Q_1$   $n_1 + n_2$ ;  $A = eval(n_1 + n_2)$  then next;  $Q_1 = n_2 + n_1$  TM could save time if  
it remembered  $(n_1 + n_2)$  & knew that commutative rule for addition.

Also, a real point of TM knowing "x-laws of Alg" is in dealing w. unknowns: I. idea that one could manipulate

expressions w. unks. Accordingly to "Laws of Alg" is set now "true" expressions, & often solve eqs

On its way. To show ex., TM could at first do ~~some~~ things in an "Lynch order"

Then later, let TM know certain operations & focus of operations were more likely

to help solve an eq.

34-4934 SN Recursive functs: The functions that are definable using simple "if" statements, are  
computable via a stack. Its my impression that stack machines are limited as to what kinds  
of functs they can compute... So this simple "if" statement probably will not cover all (partial) rec. functs.  
Can all functs. be computed w. 1 stack? Does it need 2 or 3 stacks?

● HAMPTON RESEARCH ASSOCIATES ●  
26 Boylston St., Cambridge, Mass. 02138

00: 495.40 : Perhaps try to take a bunch of (problem/instruction problem) types & partially order them for TM's acquisition order.

- One Over-all list: (1) ANL: including ~~recursive~~ recursion (or ~~SPVNT~~)
- (2) Solving linear eqns (3) Solving some n.l. eqns by use of basic "laws of Alg".  
(Laws of Alg: ideally, is level by TM <sup>to be</sup> "intuitive".
- (4) Solving some n.l. eqns by (invertible) substitution.
- (5) Differentiation
- (6) Integration

All things is symbolic (non-numerical solns). A Big media: Include numerical approx for soln. of eqns.

Another offshoot: Soln. of polys by factoring.

12 A diffrnt Partial ordering: Start w. TM (understands/knows) ANL, but for literals only.



Then soln. of  $a/b$  eqns. using literals only: first (linear of various kinds; some n.l. then quadratic (w.  $\sqrt{x}$  available), then cubic (w.  $\sqrt[3]{x}$  available). "Using Proctor, TM would never discover complex nos.!

See if it could do quartic: would it realize that  $\sqrt{x} = \sqrt[4]{x^2}$ ? (i.e. 4 values).?

16 So there are several ways to "Low Algebra":

1) ANL using integers, so answers are always exact integers (including division).

Then soln. of (linear eqns.; some n.l. eqns (polynomials, some other eqns)

2) Use of literals only for Algebra (.12-.16). [Some diffy w. division by zero; often one can't know if a denom is zero or not!]  
Next, TM with literal  $+ x - \div$  operators & literal roots; (or just use integers as in (.12).)

3) Understanding of "approx", soln. of eqns in Proctor pt., successive Approx. of solns. to eqns.

29 (2) division by zero is a diffy. We could ask TM for a solns. or as few divisions as poss.

30 We could tell TM which constants / expressions were  $\phi$ , or  $\neq \phi$ . TM could give soln. in form; If  $(a) \neq 0$  then  $x = a$ ; if  $a = \phi$  then  $x = b$ . f. soln. of  $a = \phi$  (originally  $a = 0$ )

must not divide by  $\phi$  - or it's  $\phi$  w. should have a new expression that can't be  $\phi$ .

Say  $(x-a)(x-b) = 0$ ; if we divide both sides by  $(x-b)$  then  $x = b$ ,  $x = \phi$ .

If  $(x-b) = \phi$ , then we can solve & get  $x = a$ .

Say  $(x-a)(x-a) = 0$ . so if  $x-a \neq \phi$  then  $x = a$ ; if  $(x-a) = \phi$  then  $x = a$ ; so  $x = a$ .

Say  $a + bx = c$ ; If  $a+b \neq 0$  then  $x = \frac{c-a}{a+b}$ ; if  $a+b = 0$  we can't solve it.

HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

00: 496.40: Actually, "division by zero" is a not very relevant problem: we can deal w. it in 2 ways  
A.H. Way: T. main problem is soln. of eqs. is  $\pm$  think never may be instructive situation  
r.e. "rug".

in TSO of 496.20, + "No division by 0" rule, is slightly simpler than it is in 496.23.

— This is because one knows  $z \neq 0$ ; but one never knows whether  $z+b \neq 0$ .

In both cases, one doesn't know if  $z+x \neq 0$  or  $z+k \neq 0$ .

06 Q: In 496.20: If a soln. to a linear Eq is in integers, is it always poss. to find a soln. w/o  
07 using division yielding non-integers? If it is not, then I'll have to have TM know about  
fractions. (Hrs I suspect that  $1.06 - .075$  is "yes": TM simply keeps nos in form

like  $\frac{I}{3}$  & adds & subtracts them, ect. Its equiv. to knowing how to manipulate fractions,

but TM need not know fractions "exist" ~~eg~~ Eg TM knows  $\frac{a}{b} + \frac{c}{d} = \frac{ad+bc}{bd}$

so could be zeros since  $\frac{1}{b}$  exists.

$$\frac{a}{b} + \frac{c}{d} = \frac{ad+bc}{bd} \quad \text{since} \quad \frac{a}{b} + \frac{c}{b} = a \cdot \frac{1}{b} + c \cdot \frac{1}{b} = \frac{1}{b} \cdot (a+c) = \frac{1 \cdot (a+c)}{b} = \frac{a+c}{b}$$

$\frac{2}{3}$  could be acceptable as a soln. to a problem! Whether <sup>TM</sup> should (want/need) to know  $\frac{20}{30} = \frac{2}{3}$ ,  
is unclear at present.

The ~~TSO~~ TSO of ANL (thru 494.25) (w.o. final recursive function) would be an  
instructive TSO for TM to know.

**REMEMBER:** Subgoal of 492.04-16

So 493.00 - 494.25: We observe (no. of) many times in "R" forms, so we give it  
names,  $\alpha$ : & rewrite these expressions: We then observe  $(\alpha \frac{1}{2} \alpha)$  & call it  $\beta$ .

$\alpha$  &  $\beta$  are of form  $(\alpha \text{ op } \alpha)$  or  $(\alpha \text{ op } \text{no})$ ; so  $\beta = (\gamma \text{ or } \delta)$  w.  $\gamma \rightarrow \alpha \text{ or } \text{no}$ ;  $\delta \rightarrow \alpha \text{ or } \text{no}$ .

we could then do recursive rule  $\gamma \rightarrow \text{no}, \alpha, \beta$ .

If we have many cases (494.22), then this should be easy to spot.

21 It gives us general format Q, but not its recursive A. We must modify 2. iff suits a XLT  
grammar, so elements of "Q" expression generate corresponding elements of "A" expression.

One way to devise a XLT grammar is to find separate grammars for Q & A,

then make correspondences betw. the 2 grammars.

$$\text{mul} \left( \begin{matrix} (n_1, n_2) \\ \text{sum } n_1, n_2 \end{matrix} \right) \times \left( \begin{matrix} (n_3, n_4) \\ \text{sum } n_3, n_4 \end{matrix} \right) \rightarrow \alpha \times \beta \quad \text{Again} \left( \begin{matrix} \alpha \times \beta \\ \text{mul}(\alpha', \beta') \end{matrix} \right) \text{ is recursive rule that could express the whole grammar.}$$

Good idea to Review present status of where am I out. task not for TM

3 What about 492.04-16? 3 Specifically what (is) present problem(s), subgoals?  
state bottlenecks clearly. If I can't solve one now, go onto Next one!

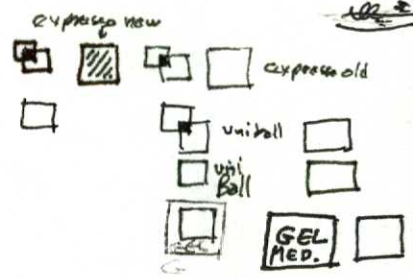
1/4/02  
ID

Convergence Problems QATM: Criterion: .00

HAMPTON RESEARCH ASSOCIATES Associates.  
26 Boylston St., Cambridge, Mass. 02138

00: 486.03R : On the convergence Perm. for QATM: In S78T3, The seq to be predicted is selected a priori: so we expect about same error in f. near terms. In QATM, f. Q's are selected to form T SQ a if f. future Q's are radically different from those of f. T SQ, we don't expect a formerly good error to continue much.  
If f. T SQ is future Q's are selected a priori, we expect better conform. in f. future.

- 10: (497.40) 1) T. Main problem now is writing T SQ's. T. Partial T SQ's written. Res for:
  - a) SSG ("An ind. mtr. machin")
  - b) SSG decub ANL a bit. (conversion).
  - c) SAARB: I ASL using stack Machine: (Much detail).
- 13 II Soln. of linear eqns is discovery of "Law of Alg".
- III linear  $\rightarrow$  quadratic  $\rightarrow$  cubic eqn. via invertible substitution heuristic.
- d) Recent ANL 494.04 - 495.23 : The "exact Mechanics" of Recursion is (unclear/undecided), In general, details not worked out: many "In English".
- e) Suggested other T SQ's : I: Develop <sup>492.33</sup> S. APL II: Differentiation  $\rightarrow$  formal integration. (492.31)
- III Use AI. work in "Expert Systems" Got them to learn f. rules.



20: In Sub I (13): 2 drafts: ① I only used <sup>active</sup> previous Solns to probs in Many for new problems

② As no. of elements in many  $\uparrow$  ( $\leq n$  elements); pc of many accesses was  $\propto \frac{1}{n}$  - which did not "scale" to large systems


w.r.t ②: I was using OSL, but I did not use it properly, so pc's of retrieved concs were not correct.



But more generally, I did not use context at all: All items in ~~the~~ <sup>same</sup> Memory had same pc.

A very simple context that would be "not bad", would be "vacancy" - pc = funct. of no. of problems back, it was created. - Also frequency of use (if  $> 1$   $\odot$ ). Context could be "How many problems ago of its most recent use" - which would  $\uparrow$  much for frequently used concs.

Say I of pc w. n (n = no. of problems ago) so (p.c. = f(n))  $\leq f(n) = 1$

Then if there were  $k_n$  concs. This could have been retraced impressing in the "n<sup>th</sup> ago" problem, Each such conc. gets   $pc = \frac{1}{n \cdot k_n}$ .

As I find applicns. for various kinds of concs, this "vacancy rule" ~~is not~~ will change - perhaps elimination of "vacancy rule" completely!

Note: f. QA's are "unsorted" a coding 0<sup>n</sup> can't take advantage of any apparent order n<sup>th</sup> QA's. But "Context" can take advantage of such order.  $\rightarrow$  499.20

Woops! T. foregoing analysis of vacancy seems wrong! - since f. pc I'm computing is to be used on 0<sup>n</sup>! - The maybe Not. See 499.00-15 for a review of this problem  $\rightarrow$  499.00



HAMPTON RESEARCH ASSOCIATES  
26 Boylston St., Cambridge, Mass. 02138

ABC

00: 498.40: Say  $pc_1$  is regular apriori: we maximize it in  $O^1$   
 $pc_2$  is of 2 conc. is it  $pc_2$  assoc. w. its "context": See [486.13; ~~██████████~~ → (486.22 R)]  
→ 488.12, 492.04. The these references are helpful, I don't think they give an adequate idea

05: I think the original (perhaps "correct") defn of  $pc_2$  had to do w. "context" in a most general sense.

06: As corpus, it included every thing that occurred in it was observed by TM since/birth. While .05-.06 seems "not bad", I'm not sure it would be a adequate for explaining to another (not me) person! (or to "me" as of 1 month from now (much less 1 yr from now!))

Also info TRAINER  
May want to print  
500.19-23,  
.34-35  
clears this  
up fine

10: I.V. seems that  $pc_2$  is related to expected accuracy of prodn.  
 $pc_2$  is related to CC for finding code... to its CJS; Least or whatever

Also,  $pc_1$  &  $pc_2$  differ not only in their corpi — but (perhaps) in their apriori? Well, maybe not! — both can use any truly apriori info — i.e. obtained before & independly of the corpus.

The  $pc_1$  of the corpus is  $\prod_i PC_1(Q_i, A_i)$

15: The  $\sum$  Least of the corpus is  $\sum_i PC_2(Q_i, A_i)$ : This is the sum of the Least of each of the problems.

Note 487.32 on logical reasoning m srch... as relevant to  $pc_2$  } But indirect way, since the corpi are different  
But also Note that Logical Reasoning can also influence  $pc_1$ . } ( $pc_2$  has  $pc_1$ 's corpus plus a lot more.)

20: 498.37 On the "recency" hour/context: This may be fine for teaching TM in a warm, friendly envt. (The infant TM), but for its later life, recency is usually not as useful, (tho it still can be useful.)  
The original TSQ is usually designed so that "Recency" is a good "Context" index — but for a TM doing R.W. research, the problems can occasionally be ordered by TM, (i.e. TM is designing his own TSQ) & . The TM may not be specifically designing the TSQ so that Recency is a good "Context", I suspect that it will often be so — that it will tend to be so in most any TSQ that TM would want to design for himself.

30: Re: Recency for SAAB I (498.13): This would really speed it up tremendously (I think!)  
Also, ~~██████████~~ CJS would be much more slowly w. size of corpus (I think). This "hour" would be v.g. for initial Tng. of TM, to get it to the  $O^{n+1} = O^n$  (on, QAnti) point.

Recency may be ok for infant TM for  $pc_2$ : But  $pc_1$  would seem to be not v.g. — Not much better than classic SAAB I! — which seems close to AZI+1.

One trouble w. using  $pc_2$  (rather than  $pc_1$ ) to guide Lsrch, is that using  $pc_1$  would favor such more bias toward by  $pc_1$  — (which is very desirable). We write mix t. 2 by doing trials in  $pc_2$  order but do cutoff via  $\frac{cc}{pc_1}$  (or vice versa) — (Does this mean any thing?) In either case 50000