

4.6.98 TM: Lurch: Optimality:

Can't find any earlier stuff on TM in 1998! Could it be that I wrote nothing specifically on TM up to now?

T. idea of "Lurch is optimum for Inv probs if all info in the P.D." : T "Proof" is

— If one finds a better way to solve the problem, then this "Better way" was not in the P.D. so "all the info was lost in the P.D."

A similar argument can be made for OZ problems: If a better way to solve the OZ problem has been found, then Lurch will be given P.D.

Then clearly, that P.D. doesn't contain that "better way" — so one must have had info not in the P.D. ABCd

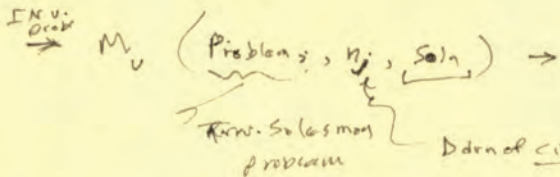
So the Q is: (1) can one put all poss. know (soln. methods) into the P.D.

W.O. sig. info. loss of info? (2) T. P.D. for OZ & Inv. probs seems different.

— What P.D. do I mean?

— Well it may be poss. to unify them.

for Inv. probs, we want a function that maps the prob. desc. into the soln. desc.



w/out  $s_j \Rightarrow$  Prob(s<sub>j</sub>) & n<sub>j</sub>  
gives shortest path & proof that it's shortest.

(I need to fill out details of this);

But anyway, say I had separate P.D.'s for OZ & Inv. probs: Since they are conditional P.D.'s I could have a single Unified P.D., if I prefixed each problem w. a symbol that told whether it was an Inv. or OZ problem.

This enables one to show abs. that a method is best in both Inv. & OZ probs.

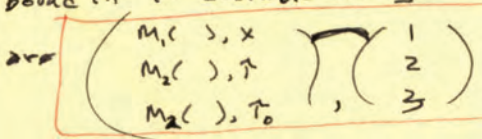
If there is only 1 P.D. for OZ & Inv. probs: then it's likely that practically any Heur. can be put into the form of a modifica. of the P.D.

Sol 86: #478: T. P.D. is a conditional prob. cond. is the problem

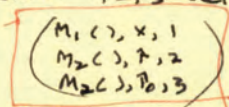
{ Inv. or OZ problem --- Inv. is  $M(x)$  and  $x$ ; (to find  $p \in M(p)=x$ ) OZ is  $M(\uparrow, \uparrow)$  to find  $M(\uparrow) \Rightarrow M(p)$  is Max. given time limit,  $\uparrow$  }.

So: for Inv. probs,  $M(x), x$  are conds.; for OZ probs;  $M(\uparrow), \uparrow$  are conds.

[ for "Anytime" <sup>OZ</sup> problems;  $M(x)$  <sup>almost</sup> completely dec. the problem; we may also need  $\uparrow_0$ , a lower bound on the available time ] — So, for the General P.D., the Conditions are



[ here the 1, 2, 3 tell what kind of prob. is to be solved. ]



Save 52.01-93.40 change to express "cond. prob. is as a stochastic operator: conditions are input; P.D. is output."

$M_1(x)$  has strings like  $x$  as output  $M_2(\uparrow)$  has  $\uparrow$  as output.

So we can have a single Prob. dist'n for all 3 kinds of problem.

outlet! Spec. → 2.40

optimal Seq. Str

I also want to make the mechanics of  $\lambda$ -Lsveh clearer! [in the OSS report, I used  $\lambda$ .

prefix set  $\alpha_i$  to assignments to  $\tau_i$  trials, but any probability (or even a semi-measure) would be ok. (All I need is  $\sum p_i \leq 1$  &  $p_i \geq 0$ .)

Also, I want to give more detail on ~~how~~  $\lambda$  form of  $\lambda$  heuristic info! How it's based on "traces" of previous work done by TM. Expand that E.N. in Sol 89 on 413.

After writing more on the GA problem (or SGA problem), I will want to integrate these ~~or~~ approaches with Sol (6, 89)

In 1.30-.40 I should show how  $\lambda$  P.D.'s for the 2(13) kinds of P.D.'s could be integrated into one P.D.

This same sort of "cond. probty" could be used to express other kinds of info that TM has.

e.g. say TM has a corpus about T.R.W. or photo of it. It could ~~learn~~ learn to make predictions about that corpus thro a P.D. on it.

Similarly, w. various other kinds ("Modes") of training, we can probly express them as P.D.'s

and make them part of the "GRAND" P.D.  $\rightarrow$  see 4.30

3 marks:  $\frac{1}{2}$  yr!

to 12.30:

$$\left\{ \begin{array}{l} 6 \text{ yrs} \\ \times \frac{200}{5} = 90 \\ \hline 50 \rightarrow 200 \text{ by} \\ 40! 64 \text{ yrs.} \\ 1.25/\text{yr.} \end{array} \right.$$

.19 [SN] Did I not write anything on TM in 1998 - up to ~~now~~ now?  $\rightarrow$  It is possl.

I spent lots of time on the SOY problem as applied to SM FT.

97 ~~22.32~~ <sup>following</sup> stuff was a main line of investigation - applicable to TM as well as SM.

I looked at all "Nite notes" from Jan 98 to present time. There is no real discuss of TM's ~~Plan~~ Plan

Plan to "SOY" problem, as applied to Fast Track.

4.7.98 TM: Gen:

I want to outline what the TM path is a list of most imp. problems.  
I did a long run of TM work on this maybe a yr. ago or less: what was it about?  
perhaps associated w. criticism of Jurgen's paper ("Long how to Len")

4.8.98: Some general mulling:

1) The 3 Programs in Sol 86: I could write proofs for all of them: Th. (astona  
in particular (Met (new using <sup>TwoShoring</sup> Lsrch), t. method is probly within a factor of 2 of  
optimum, if we use t. "Jog soln". ~~TM~~ T<sub>prog</sub> of t. 3<sup>rd</sup> Program is ~~indirect~~ indirect: If  
one has a better way to solve t. problem, then one hasn't put all of one's info into t.  
p.d. — which violates t. ~~initial~~ ~~method~~ statement of t. problem. — So t.  
"Program" is trivial. It's an explanation of a methodology, rather than a real  
derivation of a new fact from old ones.

2) Obstacles to tsq. construction: I could just use any elementary text in algebra  
for problems. If TM couldn't solve one, I either give more backend. problem  
for t. needed concs. or give "hints" (Modulus of t. problem that make it ~~more~~ easier —)  
or tell T.M. t. soln. either by giving him very strong "hints" or by giving t.  
actual cones to be used. — as if TM had spent t. enormous needed time ~~to~~  
to Lsrch t. soln.

21

(b) On "putting things in memory": I found simply adds Solns of P probs to  
t. ~~many~~ ~~several~~ many recs; retrieval becomes very costly. — (Not <sup>such</sup> > 1 memory access  
is possible for problem). — Any way, we have to have a better way  
to a case, into in memory. By "association" or other index techniques that  
I a/o TM ~~can~~ discover. .... Work in "Info Retrieval" is perhaps  
useful here . . . but I will want some very clever ways to retrieve things:

by analogy (which is a rather broad grouping!) as well as simple associations.

Recency could be an imp. bias. How do we <sup>human</sup> categorize things so that retrieval  
from our memory is facilitated?

Whenever we want TM to access its memory for something: try to think of ways  
that t. particular "correct" object is more likely — why would I think of that object much  
~~more~~ more quickly than most of t. other objects in t. memory?

If we have a set of cards that was frequently used in the past, what were t. error  
cards. (or even global cards: like "it's a chemistry problem")? → 17.09

36

On Analogy: When TM solves a problem a/o compresses some corpus:

It can express results as a structure acting on a set of objects: The structure (or  
~~perhaps~~ perhaps t. symmetry of t. objects), is stored in memory. When a new  
problem w. similar structure occurs, then we expect analogous associated info:  
we want to make it easy for T.M. to recognize that t. present situation is  
analogous to previous situation.

Also IMP: when we get info in memory, we must try to think of ways, it is likely to be useful. I index it accordingly

4-9-99 TM

Situation into ~~the~~ one (or more) [str + expl. set], previous str. diag. ~~scat~~ have been indexed in Ram for easy retrieval. When we put a prob. soln into "STORE" we also store w. r. t. + "structure" of the problem + other useful "Characterizations" of it - to facilitate retrieval.

Another source of Try. Sqns. is Kozz's Problems. Kozz's work is of interest in several other ways: One is a test to see if my Lark formalisms can at least implement all of his tricks.

Another, of course, is that OZ problems can be automatically kind into a kind of TSO using an INV kind of "Hill Climber" or - perhaps ~~one~~ using SGA.

Is it in General true that any OZ problem solver is equivalent to a kind of TSO? - even if we don't use the INV. approx. to solve it?

(SN) As a part of normal TM work: Write Jürgen! A good thing about studying Jürgen is writing him is that it forces me to examine my own methods & SSO 9.01

define them more exactly.

→ This is a guide for getting to Sci Community working on TM!

12 Find various Sectans of TM! (i.e. initial state (≅ P.D.) ~~to~~ desired final state (P.R.2). This ~~can~~ will usually take the form of a particular problem solved

Some "Sectans": Development of Euc. ~~to~~ projective Geometry; Algebra; Symbolic Interpretation (Pindol about recent advances in have TM "continue"); Various sectans of Math (Bias toward parts of Math I'm familiar w.); Ability to devise new formal lang. fit data to lang.

How to do developments: Take standard elementary tests in the TSO's Play imply: Don't fill in the needed hours in other concs.

I should try to devise reasonable + seqs for each conc. I got m, but not abs. nacy!

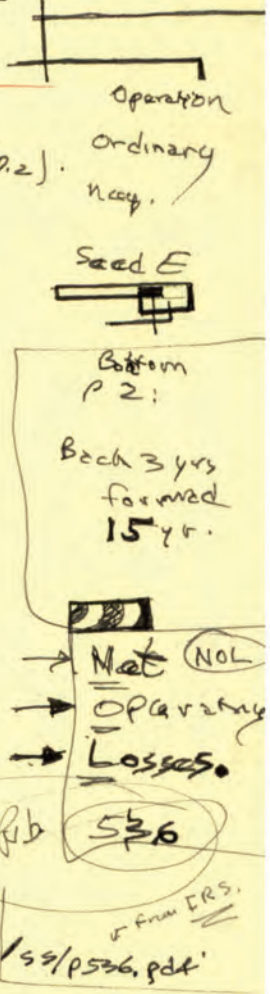
25: 1.40 Also do some kinds of H.C probs! Ones that have been done w. GA. My ideas here are a bit vague, but I idea is to develop of t. H.C. probm that is a cond. proby of t. problem given.

OZ solving aspect of TM: Something like a P.D. ~~for~~ of t. H.C. probm that is a cond. proby of t. problem given. SS. 12

Anyway: we publish a fair no. of papers that are examples of how to solve such problems (using Lark). Then we try to integrate these sub TSO's into a single TM, that is then able to use all of the info obtained from all of the TSO's.

We could then issue problems to the Sci Community of t. form 12

36 Two ways to do this: (1) Just describe parts & ask for solns. to be added to TM. TM is available on the net, so anyone ~~will~~ have available to him, the sum of all of t. contributed TSO's. Each month (say) we will put on the net, t. latest condensed integrated form" of TM. → also see 5.04



of (2) Another way: We give Monetary rewards for solns. to problems — so the best soln. get the reward. We could give all entries reward  $\propto \frac{pc}{cc}$  of ~~their~~ their solns. — Total reward is divided up w. those who.

Maybe normally sorted out at Money: except to grad. students

of We publish winning solns w. names of authors & how much they won.

~~Part~~ (0.01) has come dirty in defn of  $\alpha$ : The we could define it wrt.

a particular machine (say Pentium or Pentium II MMX).

I'd like to be able to do both (1) & (2): ~~perhaps not~~ Perhaps only best prizes for people who solve a problem (i.e. self-called "Grad students") (or prize may be Bittling Goodwin for non-grad stud; Prize for Grad stud.)

(SN) Re: Lsrch always using only 1. simple best soln. — (even after one over solves" &

prob. in th. hope of getting solns. of the user PC.) Graco suggests that in backtracking one will normally try "second best" local maxima. Is there a way to use back incorporate this into Lsrch?

It may be that this sort of thing <sup>(should could)</sup> be done by TM<sub>2</sub> ("comprison of t. P.D.")

[On the other hand, maybe I'm confusing this w. another problem: i.e. in Lsrch, we always use the cost of a soln. rather than pcost. If we had a way to assign pcosts to concs, then we could do Lsrch using these pcosts — ~~if~~ we could make an approximate ordering of concs in pc order or (better yet in  $\frac{pc}{cc}$  order).

Various concs could be labeled w. their pc approx. — we don't have to use always use cost to estimate pc. For certain elementary concs, pcost is perhaps ok., but for any conc. that we learn (of for ~~sub~~ functions that we discover as subnets), we can have pc that are not of the form  $2^{-N}$  (or normalized  $2^{-N}/(2^{-N}+2^{-M})$ )

$$= \frac{1}{1+2^{2N}} \cdot \left( \frac{1}{1+2^0} + \frac{1}{1+2^1} + \frac{1}{1+2^2} + \dots + \frac{1}{1+2^{N-1}} + \frac{1}{1+2^N} \right)$$

After we find several <sup>with</sup> solns to several <sup>each</sup> different problems: we want to "compress" t. P.D. so its more likely that these solns would be found directly (or more directly w. less srch — that is, components of the solns. would have higher pc's than "normally").

Here, I think we want to debr each problem in many ways, for a human, that's easy: where did the prob. come from, how was it debr'd in human lang, what did it.

Linguistic components suggest ~~that~~: Analogous in other problem areas. This is the machine version of the human heuristic: "state a problem exactly"; i.e. for machine: "state a problem vaguely" — implying various <sup>analogous</sup> related ideas, related to analogous problems.

4.12.98 TM

I could start out, by fixing up those T.S. Q's that I wrote:

- .02 1) Alg Notation Long;
- .03 2) teaching laws of Algebra.
- 3) learning to solve linear, then quadratic, then Cubic eqns easily using a simpler (impr.) hour.
- 4)

My main complaints about Pico T.S.Q's were

a) in 1.02 the ~~lc~~ <sup>lc</sup> of soln. became getting large because of access to info

b) (Stack Memory) was too small! (see 3.21 for a promising approach)

b) T. only ~~single~~ useful Prefs we had in store was computer coins to previous problems.

I think this could be fixed! Say we had a (later) problem that did use subfunctions (or parts) of solns. to earlier problems. we could then modify

c. entire system so it could put that sort of info in store & suitably index it.

See 17.01 - go for how to do w. Pico to some extent

.20: 4.24: Along w. this: I could put in partially solved problems. These would be prob in which parts of t. soln. had excessively by c.j's. We could "fix" this by adding a hour. to t. "primitives" or by suitable supplementing of T.S.Q so that hour would have a reasonable pc.. Also, we might have "solns" in which a human had to explain <sup>part of</sup> text of t. problem to TM, or even give TM gross "hints".

In 1.03 ... "long laws of Alg" One thing I didn't do, that I should have: was to ~~have~~ got into lrn. "if  $a+b = c+d$  then  $a+c = b+d$ " it could learn this for pure literals, then later, for mixtures of literals & nos., - so it could lrn. to solve eqns.

Suppose we gave it eq.  $3x+7=3$  find  $x \geq$  its max. we could smart many examples & it would find solns. perhaps by exhaustive search. - But eventually, after it had many examples, it would find that t. general rules gave answers much faster. Similarly, it could learn to solve many other types of eqns. Then, after it had all these soln methods in terms of literals in t. eqns, it ~~start~~ would try to find ways to obtain t. rules for a new type of equation soln, w. less lc. In this way, it could, perhaps, devr. "laws of Algebra".

But if it could learn direct solns of Pico eqns just as fast w.o. learning "laws of alg" - then "laws of alg" are really not very so understand how to solve Pico eqns.

Another approach would be via Tall T.M. i.e. "laws of alg." i.e. prove it lots of probs in which it had to use them. Some of the probs would be proofs of things coming out of

03 "T. laws of Algebra" -

→ 8.01

AH! I finally remembered how to solve cubic! ← woops!

say  $y^3 + ay^2 + by + c = 0$

Let  $y = x + \alpha$

$$\begin{aligned} x^3 + 3\alpha x^2 + 3\alpha^2 x + \alpha^3 &+ a x^2 + 2a\alpha x + b x + c = 0 \\ &+ 2\alpha x^2 + 2a\alpha x + b x + c = 0 \end{aligned}$$

We can make coefficient of  $x^2$  be 0 by suitable  $\alpha = -\frac{a}{3}$   
 So  $x^3 + (3\alpha^2 + 2a\alpha + b)x + c = 0$

No!

$y = x + \frac{\alpha}{x}$

$$\begin{aligned} x^3 + 3x^2 \frac{\alpha}{x} + 3x \frac{\alpha^2}{x^2} + \frac{\alpha^3}{x^3} + c + \frac{2\alpha^3}{x^3} &= 0 \\ x^3 + 3\alpha x + \frac{3\alpha^2}{x} + \frac{\alpha^3}{x^3} + c + \frac{2\alpha^3}{x^3} &= 0 \end{aligned}$$

$x^3 + c + \frac{2\alpha^3}{x^3}$

$x^3 +$

$x : 3\alpha + b$

so we can simplify out the  $x^3$  :

$x^2 : 3\alpha^2 + b\alpha$   
 $\alpha(3\alpha + b)$

$x^3 + \alpha x^2 + 2\left(\frac{\alpha}{x}\right)^2 + \left(\frac{\alpha}{x}\right)^3 + c + 2\alpha = 0$

$(y^2 + \alpha y + \beta)(y + \gamma) = 0$   
 so if  $3\alpha + b = 0$   
 we get rid of  $x$  terms.

$x^6 + \alpha x^5 + \alpha^2 x^4 + \alpha^3 + (c + 2\alpha)x^3 = 0$   
 $6 \quad 5 \quad 4 \quad 3 \quad 1 \quad 0$

$y^3 + \alpha y^2 + \beta y + \gamma y^2 + \alpha \gamma y + \beta \gamma$

$\beta \gamma = c$   
 $\alpha + \gamma = 2$   
 $\beta + \alpha \gamma = b$

$\alpha = 0$   
 $\beta \gamma = c$   
 $\gamma = 2$   
 $\beta = \alpha^2 = b$

$x^3 + (3\alpha + b)x + \frac{(3\alpha + b)\alpha}{x} + \frac{\alpha^3}{x^3} + c$   
 $x^6 \quad x^4 \quad x^2 \quad + c x^3$

$a, b, \text{ or } c = 0$

Boyd look in 1989-1991 TM or "Paul" notes.

→ summary: f. subs.  $y = x + \frac{\alpha}{x}$

$x^4 + 2x^3 + bx^2 + cx + d = 0$   
 $(x^2 + \alpha x + \beta)(x^2 + \gamma x + \delta) = 0$

$x^4 + \gamma x^3 + \delta x^2 + \alpha x^3 + 2\alpha \gamma x^2 + \alpha \delta x + \beta x^2 + \beta \gamma x + \beta \delta$

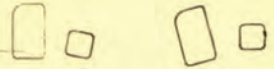
i.e. a duplicate.  
 first part of  $\alpha \gamma^2$  (so  $2=0$ )  
 then  $(\alpha + \gamma) = 0$   
 so  $x^3 + c + \frac{\alpha^3}{x^3} = 0$

4-12-98 . TM

MERGING TMS's

.01: 7.03: Note that the "laws of Alg." begin to become useful when one wants to solve eqs. w. units in them!  
So one can discover new facts about "unknown" nos.

Another good way to organize, speed up work on TM: Take certain Expert Systems & have derive the sqvs. for TM so that he ends up being able to do (at least) what these "Expert Systems" do.



.08 On the "Merging" of TMs: Say we have several different TMs that started w. the same primitive obs. set, but were given different TSO's - so they are now quite different TMs. We know the TSO's: (Can we merge these TMs by alternating items from the various TSO's into a single TM?) Maybe: But it would seem to take lots of CC, since the search spaces much larger than that used by the original "unmerged" TMs.

Another way would be to create a new TM that had the PD's of the older TMs: These PD's would be present <sup>initially</sup> unmodified, w. maybe equal wts: So if  $P_1(x), P_2(x), P_3(x)$  were the 3 original p.d.'s, the new one would be  $\frac{1}{3} (\sum_{i=1}^3 P_i(x))$


As this TM matures, it gets more & more "coupling" betw. the 3 P.D.'s - until they are integrated

.08 could be very imp. if we have a community of scientists working on various TMs, in parallel. Each group would have its own TM.

Anyway, the first step (for me, as well as for the community), is to publish many papers showing how to apply Lsrch to many problems: Also give examples of TSO's.

➤ Alternatively, w. each  $P_i(x)$  we regard it as a "conditional" prob., and we have some way of recognizing when a problem is of type  $i$  ( $i=1,2,3$ ). At first, we can do this by telling TM which set of prob. it is in. This would be equiv. to no "merging" of the 3 TMs. We could probably give less info on what the problem type is, if the TM figures it out. Eventually the Final TM should find ways to merge the 3 P.D.'s



97TM 174 ff has t. beginning of a letter to Jürgen. 

Some comments: (1) I'd like to make it clearer that I do understand his positions & appreciate his positions.

(2) That I should give my own position & reasons for it (which don't necessarily make J's methods bad).

(3) Begin Letter on a rather high Abstract level: Get down to details later.

(4) Note 97TM ~~131.01~~ 131.01 (on differences betw Me & Jug. on Goals, Methods).

(5) He is concerned that my dividing up of TM process into "prob. solving" & "self improvement" is not clearly optimum, or well defined — I should discuss this.

(6) T. Prob He's aiming for, is more General, more diff than what I'm trying to do (i.e. "Scientists Assistant") — but I expect t. solns. of my problem will help tremendously in solving his.

(7) Actually, my main problem is not that letter to J., but something close to it: i.e.

Making a general plan for my own work on TM. 97TM 22.01 ff. was in this direction.

(8) Mention: An important effect of his paper has been my ~~rethinking~~ rethinking, re-visualizing my goals & expected path for TM.

(9) Perhaps start letter with a list of imp. ideas that we agree on. Then

(a) That dividing work of TM into "prob solving" & "self-improvement" is usually non-optimum: It's best to give machine a single "top goal".

(b) It is better not to have to give TM a "time horizon"? (I may not want to say this, tho I understand his desire to do w. t. "time horizon")

AA  
297 <sup>TM</sup> 189  
Action  
Align

4-22-98 TM:

A somewhat New Approach to "Trancy" TM!

"Its easier to learn a new lang. if one already has an old one": TM starts knowing Machine lang, or some very simply structured lang. (not nearly a "complete" lang) : Lisp, or that "functional" lang. I invented, or 2 Fork or whatever. (A <sup>Risc</sup> Risc (understand it as Machine lang.) → 6502? or some other simple instr set.

We can start by telling TM various facts — about Numbers & other objects that lang. deals w. normally. (If its numbers, we might try to teach it 2 text books Roman numerals USD!). The facts could be about Number Theory w/o Algebra, to start.

"A/o": About <sup>relative</sup> positions of things in 1 dim, 2 dim, 3 dim. <sup>time</sup> It could be a "Blocks world" type of knowledge. We may also have it know about Time, & how events,

("objects") can be ordered by "time" (≈ a 4th dimension).

Perhaps even dynamics of objects. — trajectories, physics, etc.

Perhaps look at Part Simms work on Teaching Machines "Physics" of R.W.

After lang. about a world in its initial lang. — we have it learn to understand equivalent statements, Q's in a simplified simplified English. (e.g., no pronouns, few ambiguities, only 1 way to parse a sentence, etc) — Or, have it learn some different computer lang(s) . . . . Then simplified English.

After lang. English, we ~~can~~ have it deal w. R.W. situations close to the world it knows about, & we gradually expand this known world, & we gradually make the lang used to talk to it be closer to English.

Also as we continue, the English we use becomes noisier & noisier. More & more errors in spelling, grammar, more pronouns, more ambiguity in parsing, etc.

26 There are many possible paths in {language, Universe of discourse} space that we could travel, & end up w. a machine understanding English to the world (part of) the world of man.

Our main goal is to get it able to usefully crave the web. It probably could get enough info to be very useful w.o. ~~having~~ be able to understand graphic material.

Instead of the web, perhaps use <sup>source of</sup> Brewsters "Snapshots" of the web as a data base! In addition it could have many v.g. CD roms, like the Britanica, McGraw-Hill Scientific Encyc, etc. T. encyc Britanica is < 170 Mb; I have 7.2 Gb; H.D.: So it ~~should~~ should not be hard to get enough info on a H.D. for TM's later "lang".

35 I should explore (informally) many of the paths of 26 before doing a serious study. Perhaps show how other approaches to A.I. are "special cases" of it. E.g. my Older (Sol 89) approach; CYC; etc. Blocks world (expanded to really complex physics) ! Karl Sims' world for his "creatures" w. some laws of physics, time, Gravity . . . Leavets "AM" w/o "Envisca" ?

5.10.98 TM  
6.14.98

I would like to be able to "divide up" the T.M. problem so several groups of people could work on it independently in parallel.

One idea involved the "merging" of T.S.Q.'s for different trained T.M.'s.

Another: We train a T.M. to do induction (problem solving) in a certain area, but the T.M.'s knowledge in this area is incomplete: Many of the cases it uses have been "memorized" or "unlearned" or obtained with gross "hints".

Deriving T.S.Q.'s so that T.M. can replace each of these "unlearned cases".

(or "inadequately learned cases" or "imperfectly understood cases") with

them cases that is properly found via a proper T.S.Q.

6.14.98: Another neat trick I discussed in the past. Given 2 T.M. initially the same but given different T.S.Q.'s, so they are now quite different individuals. — we can combine to some extent, combine the 2 T.M.'s by raising a third T.M. on basis of "merged" T.S.Q.s. So we combine 2 T.M.s by combining with T.M.s but their T.S.Q.s —  
⇒ T.S.Q.'s are usually easier to combine than T.M.'s!

5.28.98 = TM:

How to teach a machine that can learn.

11.1

on Sep 9, 98: Gamm wants me to be invited speaker at Computational Intelligence colloquium <sup>1 day</sup>

aka Riss, Wallace.

Poss. talk: How 'The teaching for learning machine

How to teach <sup>or</sup> learning machines, ~~an~~ ~~un~~ ~~der~~ ~~stand~~ ~~ing~~.

Teaching Machines that learn. ~~NO~~

Maybe try to publish in "Machine Learning".

Discuss ways to break up problems so Sci community can work on it in 11.

Teaching industry machines.

6.5.99: TM expo:

The teaching of machines and of humans.  
How to teach machines and humans.

Title for London talk: How to teach the Machine - Let's count the ways.

1) Teaching by neurosurgery  $\leftrightarrow$  Least (expert systems) why/why not. (or when, when not). Success Chess, Symbolic Integration...

2) Teaching by example: How this can differ much better. Humans learn we can give machines much details of examples.

For humans, e. "example" is much various: Humans are forced to do many many induction

Q of difference of 2 sessions: But Inductive Inference & Machine Learning are essentially f. same problem.

Relation to Human Learning:

- 1) Good Models
- 2) We can use these models to make good teaching machines - But make model of student's understanding & propose problems & hints & study ideas (Things to read, Things to read about, interesting facts ... - But are based on that model).

Main part of ~~talk~~ paper: Analysis of listing & analysis of various ways to teach machines - One one extreme: <sup>most</sup> Expert systems - On max max max Data Getting a machine to do data mining w. minimal & priori ideas about B. Cooper

Much work in ML is oriented toward problems arranged so that it's easy to tell if Machine is (right).

Division of this field into ML & Ind. is rather difficult! ~~Any~~ **Any** practical method of induction amounts to machine lang, but some ML methods <sup>only</sup> translate dirty, messy induction (unclear functions) 6283356

How to teach a machine - let me count the ways ..

Many machine techniques have been proposed and attempted, to bring a general purpose computer from an initial state of ignorance to a state of some intelligence some computer in a variety of fields. These techniques have ranged between two extremes.

often used One extreme often used must often used in expert systems, laboriously extract knowledge from human experts and then tries to program this information into the machines. Another extreme is often used in data mining, where which we

6.6.98

H.2  
163

→

~~Give~~ Give the machine some very general techniques for pattern discovery,

and we set it loose on a warehouse full of data.

→ { What are these extremes of? }

[ Review Lanst's work: first AM as a kind of ITg environment

6.11.98 TM

found Perini Professor now 11.1, 11.2, 11.3

I seem to have lost first stuff (2 pp) that I wrote in 2 reply to Gammann about what I planned to talk about: title: 2 How to teach a machine

Also, remarks about how ML is Ind Int. were of some help. ML was bigger. Usually Ind Int research could be mapped into ML, but Garry t. other way wasn't always as easy (2 way fancier!).

Also brief disc. of Wallace Dover paper — why I'll revisit it.

So: two same approaches: Expert Systems on one extreme, AI (of some kind) on other. Extreme in amount of BIAS put in by pmr.

Also 2 extremes in "Brittleness"

"AM" was peculiar: lots of info at start, then very little (some from user response perhaps)

CYC? (unclear... It really not familiar w. it.)

SW "Agents": How

Compare w. human intell: Lots of initial info; designed for an ensemble of conversations that its very likely to find itself in.

My main goal here is to find an optimum point for myself; To get to a usable TM in say, 10 hrs.

The method of teaching of course, depends on your goal: Much of the method of teaching to date has emphasized its attempt to have little BIAS, so it

TM did something clever, it was clearly to machine's cleverness & not it's pmr.

I don't have that constraint (I think Len et does a trick) — I just want to get it's job done". As for true learning, I think I can tell by how the machine solves problems — as to whether its long in a useful way.

(i.e. useful for future problems). Later, when the machine's behavior is too complex for me to analyze in detail, the problems solved by the machine will be sufficient evidence of its competence.

contact Jack Skolansky

(7.18.98: I wrote him email several wks ago; but no reply)

prepared for humans

Another poss. goal: to bring machine as exp positively as poss., to point at which it can useful read documents so that it can autonomously, complete its education.

Time/money

Idea of "path" in "Competence" space

Idea of a project for the "scientific community as a whole"

Idea of "combining knowledge" of 2 machines that have been trained diff ly in diff erent ways. "Know" diff erent things. (This may be poss. w. certain kinds of trng.) but not w. all kinds)

6.12.98 TM

- 88 3f  
- 22 2P

13

T. problem of trying to machine is fairly difficult. - One approach to difficult prob is to make  
a nominally more difficult hard problem: "What's the best way to train such a machine"?

What flow is the machine training related to 'inductive Inference'?

① One trivial level: Given a new problem: how does it bring all of its past experience to bear on it.

Soln. of this new problem ② If the machine uses Lsrch, it searches for solns. based on probability - ~~inductive inference~~ - or inductive inference

For some time, we have had machines capable of ~~handling~~<sup>solving</sup> any conceivable problem Pattern in a reasonable form (What is for that pattern a reasonable amount of time, form, The problem is in lower, how to make our priority a compact machine - useful machine - usable machine intelligent: creator, architect. ABC

How to teach a machine: For some time we have had good effective methods for machine learning. Various kinds of problems have been solved by ~~the~~ machines using these techniques. These techniques have been used to solve a wide variety of problems - some of great difficulty.

How far can we go with these methods? A problem that continually prevents the teachers well past what has been demonstrated.

~~But~~ One real milestone in A.I. research would be a machine that could read and understand human text on a wide variety of subjects. It would then be relatively easy to teach such a machine by giving it readily available text that has been designed for humans.

Much work has been done on Machine understanding of text, none of it has yet approach pattern close enough to this goal.

How to teach machine: Games 2! 6/14/98  
We now have many techniques for machine learning. They have been applied to various ways to use travelling techniques.  
We will discuss efficient ways to bring initially have machines to a state of high intelligence

I don't know how this will fit into your division into "Machine Learning" and "Inductive Inference". Just about any applicable model practice work on inductive inference is usable for Machine Learning.

Perhaps it would be best to decide on the division into subjects at least all of the subjects of the talks that have been decided on. Spontaneous submit articles or subjects for their talk.

About the Wallace, Dowe paper you sent me to review: I can probably write a review of this relatively quickly. - ~~Have~~ discussed many of the issues and ideas in the paper with both of the authors, at some length.

Most Sincerely,  
Ray.



- .01 There are probably many reasonable paths in language, {universe of discourse} space that would work: (10.26) <sup>fruity</sup>
- I don't think I can pick an "optimum" one: But I may be able to pick a not-bad one, in which I get max. leverage out of it - Sci Community. — by interacting Team in doing part of work (b) by taking old A.I. w/ usage re-expressing ~~Team~~ Team  $\Rightarrow$  as ~~partial~~ ~~TS95~~ TS95. (10.35)

A recent idea has to be more directed toward language (ing ~~ing~~), using a simpler universe of discourse, than, say, Math. In starting w. Math, I wanted to become familiar w. Eng. TM to do fairly hard probs; then Eng. 'language' could be a sort of Math problem.  $\rightarrow$  Alternatively, I'd have TM. Become familiar w. a simpler world than Math (The math can be very simple!), & then begin to learn langs moving toward English that discuss that simpler world.

(SN) & A motivation for discovering "laws of Algebra" for Humans, is that it enables more rapid calculation. i.e. if we knew  $3 \times 97 = 291$  ~~then~~ ~~an~~ we know commutativity of Mult, then we know  $97 \times 3 = 291$  w. very little calcul. — Now we must TM's can calculate  $97 \times 3$  ~~more~~ ~~rapidly~~ ~~than~~ they could retrieve a rec'd  $\&$   $\&$  (cu. of  $3 \times 97$  from memory, so no ~~value~~ value for knowing commutativity of Mult/Expres...  
 Now, we can vary costs of various operations so it is cheaper to retrieve a rec'd calcn of  $97 \times 3$  than ~~to~~ to calculate  $3 \times 97$ .

We might also have TM devr. logarithms so it could do multiplication rapidly using tables; ~~addition~~ <sup>much</sup> ~~which takes less time than multiplication.~~

The Muzlis! ~~In~~ ~~fact~~ ~~ing~~ TM to learn via Human texts, it may be hairy for user to (pro tem) instert certain human-like limitations into TM!

Notes for the talk on Lsearch:

Explain just how Lsearch can be optimum:

First: that its "UNIVERSAL" — so it will solve any ~~the~~ INV problem in  $K \times$  optimization.

~~the~~  $\epsilon$  (prob) of ~~the~~ soln. wrt.  $t$  ~~is~~ reference value.

Problem is:  $\epsilon$  this soln. can be very small  $\epsilon$  ~~is~~ very large

~~How to~~ Let's take another approach to search ~~now~~. — Examining Heuristics Search.

Heuristics ~~Order search by~~  $\downarrow$  such time by ~~order~~ <sup>suitably</sup> ordinary trials.

Any recorder of trials can be realized by <sup>suitable</sup> modulus of  $t$ . P.D.

So for any particular Heur, I can modify  $t$ . P.D., so that by using Lsearch wrt ~~the~~ P.D. I'd get same effect as the heur.

So we can solve problems in either of 2 ways:

① Heur search ② Put Heur into P.D. & Do Lsearch.

Is there any advantage of ②? — Yes! See Sol 96 for reasons.

So: work into the "Term" that Lsearch is optimum for INV problems

Also in discussing Use of Lsearch for Geometric Alg. problems:

Note that M.C. searches are proby  $P$  for object of proby  $P_i$ .

expected ~~no~~ of trials =  $\sum \frac{P_i}{P} = N$  i.e. for an object of proby  $P_i$ ;  $t$  takes about  $\frac{1}{P_i}$  trials.

So random search rate deterministic search are just as fast.

L search is much faster! G.A. searches are usually  $\leq \frac{P_1}{P_2}$  searches ~~are~~ very inefficient.

M.C. search w. proby  $\propto \sqrt{P_i}$  is best ~~if~~ if  $\sum \sqrt{P_i}$  converges

but L search is even better. (It's in order of  $P_i$ , if all  $P_i$ 's are  $\approx$  same.)

$$G = \sum \frac{P_i}{f_i} = \min; \sum f_i = 1 \quad \frac{dG}{df_i} = \frac{P_i}{f_i^2} + \lambda = 0 \quad \frac{P_i}{f_i^2} = -\lambda \quad f_i \propto \sqrt{P_i}$$

**PARADOX?**

In Sol 78 in discussing Cover's "Extensive Complexity"

I didn't actually show this, but I may be able to do so easily using the technique.

I ~~should~~ <sup>should</sup> have noted that Gacs's Semi-measure norm is not a good "measure" norm. T. pc's error by a factor that  $\uparrow$  non-increasing slowly to  $\infty$  as  $n \rightarrow \infty$ . This seems to differ from my more recent results ~~that~~ that Gacs's norm is my norm give about same results (Prob for regular name than "Extensive" probability)

Well, ~~the~~ <sup>measure</sup> norm is unnorm. Semi-measure both  $\rightarrow$  same limit, but norm measure approaches it much faster. The rate at which measure & semi-measure approach each other is determined by the rate of convergence of  $\sum f_i$  measure, to  $t$ -correct values.

7.1.98 : TM:

In §197 (Bore) and perhaps in (Isis) (Malbaux) paper ~ 1996: I was berating Wallace & Ross' methods because they only considered a countable no. of hypotes. (for the discrete part of the spaces considered). This is because the set of recursive functs is not effectively enumerable. Usually (if not always) they only consider primitive Recursive functs: — which are ~~not~~ effectively enumerable.

Hvr. when I make approxns to ALP (via **R&P**), I, too, must confine myself to a countably enumerable ~~the~~ functions. In doing so, I make a definite choice as to which ones to consider (i.e. not consider). This gives a real bias in the results.

On the other hand, say one selects a Ref. Univ. ( $\equiv$  P.D.). The models that are enumerable in ~~time~~ <sup>w.  $CC = C$</sup>  are a finite set, and all models are ordered as one  $\uparrow C$ . So the hypotes seem to have a "natural order" if one takes  $CC$  into account. It's only when one ~~considers~~ includes into the set of legit. models, ~~the~~ those w.  $CC = \infty$ , that the list of models become non effectively enumerable. ACTUALLY  $\frac{P.C.}{CC} \in L.C. \cup L.C.!$   
 $\rightarrow$  2 v. G. ordering.

So it would seem that the R&P formulation may, indeed, have a slight advantage over any other also (Well beyond its advantage over systems that only consider prim. rec. functs).

Anyway, the "natural order" <sup>of models</sup> induced by the bounded  $CC$  requirements does, I think, give a larger, "better" set of models than using simply prim. rec. functs, or other ~~simple ways to~~ essentially ~~any~~ arbitrary ways to order the ~~total~~ recursive functs ( $\equiv$  epim).

Does this "Natural order" also solve the problem of "pretty ratios"? For small  $CC$ , the ratios can depend much on just what approxs one uses; but perhaps as  $CC \rightarrow \infty$ , the  $\rho$ 's approach some true limit.

7.10.98 TM ~~TSQ~~ TSD:

① In my earliest try seqs., I was disturbed that t. only abs. used in a ~~few~~ few ways t. primitives plus t. (entire solutions) of previously solved problems.

Well, perhaps that's not solved! If a problem seq. contains ~~useful~~ sub abs. "useful" that were not t. solns. of previous problems in t. TSP, then this latest problem is probably a very diff one!

It may well be that finding useful sub abs. in a corpus of problem solns. is something that does not occur often (like most heuristics are not easily (or frequently) discovered).

Also, t. TSP's I wrote were very simple, ~~the~~ essentially "Skinner sequences" - w. no pussy of error. (~ "perfect" TSP's).

09/336 ② Another Major diffy in this TSP's was that the techniques did not "scale well" to larger problems: The "Memory ~~error~~" containing useful abs, became

Very large - so ~~part~~ of anything in t. memory ~~was~~ ~~no longer~~ ~~usable~~

would become very extensive for larger TSP's.

I think t. soln. of this fact is to use Info Retrieval methods: When one puts something in Memory, it should be indexed to tell what kinds of situations it is likely to be useful for. This is not an easy problem to solve: it may be possible

can do it only after one has a large and corpus: So far problems

early in t. TSP, one has no IR, but it's not as heavy. For later probs

it is heavy, but one has enough data to implement it.

③ From ① In my ~~2~~ 1st seq. soln. to t.  $11 \times 11$  input /  $3 \times 2^4$  input) Max problem, again t. only new abs I used after a soln. was t. actual (entire) soln. of t. ~~the~~ latest problem.

ON Z141: At one pt., I decided that Z141 wasn't really so good for predicting sequences, but that it was useful, in a general (heuristic) way, for discovering <sup>for discovering finite state grammars.</sup> trial n tips used n tips.

Actually, I wanted to technique as a general analysis of the pc of constructing more complex objects from simple objects. This is certainly an imp. problem!

It may be that at the time I developed Z141, I didn't know as much about coding - using prefix codes, and all. Anyway, the problem that Z141 addressed was: Given the phenotype, to find good codes Genotypes for it. - This is harder than the problem in "G.A." - in which we have lots of Genotypes & their assoc. "G" values. These "Genomes" can be trees expressing a Universal function lang.

To a first approx: Z141 assumes that the corpus was obtained by context rules

of the type  $A \rightarrow C \cdot D$ , where  $C$  &  $D$  are ~~non-terminals~~ non-terminals (like A).  
 $C \rightarrow E \cdot F$   
 $C \rightarrow c$   
 $D \rightarrow d$   
 $D \rightarrow F \cdot G$   
c & d are terminals.  
I think there are no loops.

To some extent, we could discover non-terminals by noting unusually by trees of various n tips. However, it's not clear how best to use info on unusually low freq. of n tips. Th. grammar of .13-.18 could be used for the genomes of a universal lang. (.05-.10)

At this point: There are several related ideas:

1) The initial idea of Z141: of discovering n tips by (recursively) looking at frequencies of pairs of adjacent n tips.

ABCDE  
 $\int e^{k^2 dx}$

2) The idea of a stack grammar of type (.13-.18) I think it's a FSG. I like

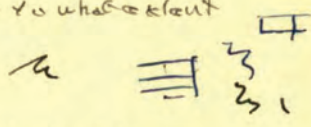
It could be the precursor of a universal lang of .05-.10: would this give a universal probability

distribution? Perhaps it would: All I need to show is that there is a common upper bound on the -log prob assigned to a sequence, and it's done by a hmc. - is that this is true for all sequences.

3) To the extent that .13-.18 is a FSG, it could be discovered by "Hidden Markov

model" techniques! HMM is universal way to non-hidden Markov model of .13-.18 - A. HMM have been successfully discovered. - to what extent

could these techniques be used for .13-.18?



4) I'm interested in G.A. of various kinds. How are they related to L-stch? - How related to XLP? It would seem that G.A. creates a new population by mixing

"new by" in complexity space - i.e. it selects new candidates that are coded with conditional complexity that is small w.r.t. original population. Mutation does this in a simple way. Crossover in a more complex, apparently more powerful way.

Wrt. to "Grand Pitches" of the T.S.Q. problem! At least 2 approaches:

1) T: / <sup>Big Sub</sup> goal was to get TM able to do "self improvement" well - as soon as possl. —

So I wanted to take a lot of other methods & make a genl. lang to express them.

See 7.24.95 <sup>TM General</sup> ~~Wrt. to~~ & 7.3.91 P 472 — Both in Mistake.

Since then, I decided that I would "wire in" some fairly good & EZ soln. methods — & that using them, TM should be able to learn lots of things — including better & EZ methods.

2) The idea of <sup>many goals</sup> parts in  $L \rightarrow \{ \text{Languages} \}; \text{Situation defined by } \{ \text{language} \}$  space.

{ 10.26, 14.01, ... } Here, the ~~main~~ main goal, is to get TM to understand

English, as soon as possl., so it can begin reading Human Texts, & eventually,

**T**. Internet. (We could use a static version of the Internet! say one of greater k-hole's "snapshots" of the web ~ 1 or 2 Terabyte)

Ability to read human text would make it much easier to train TM further.

2) A quite diffrt track: Th. Ramanujan <sup>Jan</sup> T.S.Q.! There are 2 Books he read —

which is about his entire formal education in Math. First got to Man who ~~knows~~ know of "

(Bio.) to get to know, of the book — Major in G.Hardy's Math analysis? (No reference).

There are other Books on Ramanujan — ~~see~~ now has all his published work — may have R. needed info.

This goal would be more in the direction of showing that my approach to TM can give very "creative" results.

Perhaps easier, or as "study problems": @ Learning various tracks of Symbolic Integration.

[I'd have to learn what the state of the Art is in this area]

Easier yet: Going thru an ordinary seq. of High school, then College Algebra books,

Or Euclid: The theorems in the order of some text book. I want to be able to

get TM to do this as a "Human", w.o. nearly understanding the complete logical basis of it.

How much "intuitive understanding" of Geometry is needed, is unclear.

or "Automata" (39.01 R)

ABCDEF G

.30

→ A diffrt idea related idea: After TM has knowledge of any fixed set of sufficient size,

it can begin to learn to understand questions, statements, etc. a bit that would be

a lang. close <sup>to</sup> English. For its "Mathematical" training, statements & questions would

be in a very simple formal language. Later, TM would learn to translate between ↔

-2520@989

English & that formal lang.

How, we train writer, in English, about things TM doesn't know about. How he deals

w. this, is unclear.

We could start w. knowl. of Math & English about them. — Then have it learn about "Physics" (Block's world or Karl Süssner world) — and gradually mix training about various areas w. understanding English about Physics etc. → 2.1.01

Q! How does this relate to CYC?

Reading Some old Notes (not so old... 6/92)

There was the idea of not having to write T.S.Q.'s by (a) T. problem pool (MM) start w. only a few simple problems, then add harder ones as TM learns. (b) Competition betw 2, or 3, or 4 TM's for "food" — as they grow smarter, they become more of a challenge to each other. 3 or 4 or more TM's involves ~~or~~ cooperation, sacrifice, etc.

(c) Have TM work problems that are easy because there's no noise. We then add noise to make harder problems. This is an easy way to construct T.S.Q.

More recently, instead of T.S.Q.'s, I had the idea of framing environments. — So any kind of devices can be used by the teacher. Hints, feelings, giving advantages (or partial knowledge)...

Karl Sims has simulated creatures that fight, w.  $\approx$  r.w. physics, gravity, friction, etc.

A survival contest. Could we also reward them w. "Food" (= Energy = reward)

If they solve certain problems we give them — like "I consume food"

(SN) In use of Lich to solve  $n+2^n$  input Mixer: If there is any noise in the system  $\frac{S}{N}$  decreases rapidly as  $n$  increases. One way to  $\uparrow \frac{S}{N}$  is to  $\uparrow$  SNR, but this increases CC of soln, rapidly. To what extent is normal G.A. soln. of this problem (Say Koza's soln) degraded by this effect?

Another way to  $\uparrow \frac{S}{N}$  is to do search for 2 or 3 kinds of addresses, rather than just 1. This involves more search, but + Q is — do we in the long run, solve the problem faster?

For the Paper, is for my own use: Try to list "all" of the ways to teach machines: The category theory theorem so one can see just "what's going on" is "what the alternatives are"

Maybe put the various methods on "file cards" or equiv.

Each "card" or ~~equiv.~~ <sup>equiv.</sup> will have brief description of the method plus (perhaps) a list of references.

Some newish ways to teach machines

- 1)  $\leq$  w. Aspects: possibility of net.
- 2) A list: Several different techniques
- 3) Simulation of various kinds of "Complex Systems".

Use of Logic or Espanto

— Since they should be assist the machine to learn.

Related to (19.30 - 19.40)  
 Yet another Approach to TSCQ! (Actually New!)

We start out w. a TM that's smart enough to interact usefully w. a human.

So we start out w. a min. amt. of "Common Sense".

At first, the human communicates w. the TM in a very simple lang. ...

Then as they interact over time, the lang. gets more like English & becomes richer.

The goals of the machines are: (1) To learn as much about the human in his world as possible, in available time.

(2) Not to waste human's time (3) Not to "wear out" the human — but to keep him interested.

When the TM gets smart enough he may have > 1 "teacher" — <sup>sequentially or in //</sup>  
 Also, he will have to read English Usefully.

So I guess this human teacher ~~can~~ could be a parrot from initial "lang by neurosurgery" to lang. by reading English.

We want TM to be smart enough so that its Q's asked of the teacher are of much info content. So TM thinks a lot before Q's.

Q's are **Expensive**.

We may arrange so that TM can make patterns in sound & video that may amuse the teacher.

Perhaps the goal of the machine at first is merely to amuse the "teacher".

It could be like Karl Simms' video. **Pictures & Music**.

Maybe stuff like Dave Durbin's: In a public Area, it learns to amuse people it interacts with. Maybe charge 25¢ for 5 minutes of interaction! User wears phone, maybe Google!

The maybe ~~main~~ Emphasis on ASCII <sup>mainly</sup> interaction: So it will learn to read English. One could have other modalities of interaction parallel w. ASCII. So to ASCII verbs to occurrences in the other Modalities.

Perhaps TM has access to Data Set and it must learn to answer Q's about it. It could occasionally be given "correct" I/O pairs: Or, when asked a Q — it could ask for Answer. — P. 13



7.30.98 TM TSQ. <sup>short run</sup>  
would ↓ its Gave, but maybe ↑ in long run.

hs hs

For the Teacher, Plover's a pre-processor for his Q's & A's: It tells him if it's ss were O.K. — that way used unknown words: "kennu" syntax. So teacher can edit what he tells TM.  
— i.e. known plus for. When a Q w. One unknown is given, TM may usually ask for answer. — Or maybe give "free" answer — Teacher can judge when. (The TM can learn to do this "editing" feedback to user (= teacher).)

Is there a way for Teacher to give hints (ways to ↓ tract of search)?

Teacher can occasionally ↓ cc of TM's Q if it's a Good Q as judged by Teacher.)

**Teacher** doesn't have to be right every time. ← very imp!

Re: .20: I had this idea that MTM was quite different from NMTM... Realizing that they are about the same, makes noise much less imp!

I was thinking about learning ps's w. a teacher... in which it would be terrible if the teacher made a mistake...! — No more worry about this! It.

Mobots normally ~~are~~ aren't disturbed by noise!

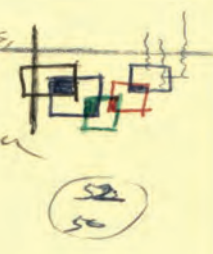
7-31-98: How "S.W. Agents" fit into this:

1) One Advantage of passive <sup>training</sup> ~~teaching~~ rather than Active <sup>training</sup> ~~teaching~~ (w. teacher).  
— passively TM's can be easily combined by combining their TSP's:

It's harder to do this w. (or impossible) to do this w. Active trained TM's.

2) Worst SW Agents: T. user can periodically give them a score for their performance & / or ↑ diffy of probs given. Also user can vary rewards, punishments for time taken, goodness of reply, etc.

3) Agent watches human and periodically says "I can do that": TM's Probs tested & may be given new response abilities. 4 5



8198: **Try this:** There is a small (to startoff) corpus: perhaps in English, (perhaps in Esperanto, or loglan or a simplified English or computer lang.). We start off by writing Q's in a simple lang. (→ machine lang.), & TM learns to answer them, using data in the corpus. At first, the Q's are word oriented. (find ss containing certain words or ~~not~~ ~~settable~~ syllables, ... or ...). Then as time goes on, we ask more & more complex Q's about the corpus. Then we relax linguistic constraints on the corpus — by making it closer to English, by using b.

Syntax have ~~some~~ occasional errors

We can give TM examples of correct Q's:

We can allow TM to ask Q's of trainers, but such Q's have ~~to~~  $\leftarrow$  low cc.

**SN** One criticism of my TSO's in Saarb., was that their use of Memory was Bad. That as time went on, Cost of retrieval would  $\uparrow$  so it would become less as less Useful. (T. model ~~did~~ did not scale to larger problems).

T (proposed) soln. was IR techniques: When one puts info into memory, one had to index it so that it would be easily retrieved in appropriate circumstances: Thus "case of Retrieval" also meant it would have much higher PC in appropriate situations than just " $\frac{1}{N}$ " (N being no. of items in Memory).

How to proposal (.10) is also use of ~~other~~ numerous heterogeneous techniques for IR, is a rather big project... forbidding!

(The... perhaps there is a standard "Goul" soln.)

→ Anyways: Shouldn't such automatically devise IR techniques - emerge as "Universal" & close to "~~the~~ Optimum"? If so, how does it Automatically do this "IR"? (both input & output for "IR"). How this can be done (conway): 33.11 ff

Well, first. <sup>①</sup> Machines I was using, were not Univl. <sup>②</sup> Even if they were, I'd have to order "Lead" from IR techniques by a suitable func Eval. a/o ~~use~~ use an Instruction set ( $\equiv$  A-prod) that made the IR ideas easy to discover (by PC / ~~low~~ low cc).

→ Big Q: Say I get a system of 22-30 ff running well: that it could do IR for ~~the~~ English Requests. Even that it could learn about English syntax to analyse Eng SS to see if they helped answer to target Q.....

Would such a T.M. be able to read English texts in, say Algebra, & learn Algebra - well and to be able to solve hard problems in Algebra - ~~then~~ the "creators" work in Algebra?

In Paris Paper (~1959) I had a teacher that seemed nice for the technique:

Here, perhaps "Not Necessarily": I can have certain Categories of sets of strings:

$\Sigma$  (known acceptable set of strings),  $T$  (known set of terminals),  $N$  a finite set of non-terminals:

Rules are sets of strings (of which  $\Sigma$  &  $T$  are included),

We try various rules of form  $N_1 = N_2 \cdot N_3$  (or  $N_1 = N_2 N_3 N_4$ )

to see if they work: Initially  $N_1 = \Sigma$ ,  $N_2 = \epsilon$ ,  $N_3 = T$ :

A way to try to find rules that work well (not nearly 100%) — we have

scoring a value of 2 but that's not 100% correct.

**SN** There was this problem of how to compute PC of error correction in to ~~the~~ induction of sets of finite strings.

It may be that my recent theoretical soln. (see 1990 letter to Wallace plus the "proof of convergence" based on strong conjecture) may help w. this.

I think I did find a soln. — sometime after Saarbrücken (or maybe at Saarbrücken).

Actually, using ALP, we realize that we will not nearly get a "rite grammar", but in general, to symbol a substring (probably) will be good.

-26

**SN** Imp't: In 2141, I was concerned w. ~~lengths of~~ usage

2 types of unusually low frequency.  $\rightarrow$  I felt that the Grammar Discovery technique had to be able to use this info.

**NOT SO!** the frequency of N tops was merely a heuristic device: A measurement to help find good gram defns, — but not nearly infallible.

The reason I thought it was infallible was that I constructed the word grammar, that was not a Bernoulli Grammar.

As a general methodology to idea I had was good: i.e. Look for present such Grammar (i.e. best found so far). Do various N tops occur w. expected frequencies? If not, modify the grammar so ~~the~~ traps are closer to

that the Expected traps are closer to observed traps.

T. trouble is, the grammar I devised to express ~~low~~ unusually low freq. 2-tups was somehow bad, or at least it was not used properly.

I think it was it would: first get freqs of symbols; look at freqs of pairs & make currency tables. But the result that grammar was

31

32

8198 TM P56-discy

Somewhat constrained. If  $\alpha = ab$  was defined,  $ab$  had to be reduced to  $\alpha$  so the freq. of  $b$  following  $\alpha$  in the resultant code had to be zero. It was no longer a basic grammar, so I couldn't use the same heuristic for discovering n tips.

Hur., the hour of 27.31-34 is very good! It gives a sequence of improvements in the grammar. Unfortly, I suspect that it often  $\rightarrow$  local extrema, so one must normally try several poss. "mutations" of the grammar at each "branch" in the path. — This  $\uparrow$  cost of the grammar discovery, transcendently.

8.9.98! P56-discy:

HotBot.com  
Larry Randall

.17 Using tech reqs of, say 27.31-34 find n tips ~~that~~ whose definitions compress the corpus. This amounts to a CFG w.o. loops.

If we had, for a corpus, a CFG with loops, if we did this n tip defn. to compress it — we would discover lots of n tips of interest and, if the corpus was large enough, we'd discover  $\gg$  very large no. of <sup>grammar</sup> rules of the type  $A \rightarrow C.D$ .

We then look at this grammar alone to see if we can compress it by introducing loops. It should be possible.

One (perhaps ~~is~~ a sufficient) way to suspect a loop!

We have many rules of the form  $A_i \rightarrow X B_i$ .

We then suspect that  $\{B_i\}$  are a nense. We look for other rules of form  $C_i \rightarrow Y D_i$  & suspect  $\{D_i\}$  is an ntipst. If  $\{B_i\}$  &  $\{D_i\}$  have many common elements, we suspect they are the same.

Or, just look at the elements of  $\{B_i\}$ : ~~Can we make~~ the term of that set has a grammar (that we know) <sup>already (have discovered)</sup>: it's the sum of the rules that are able to generate its elements.

.26-.29 ~~are~~ not so interesting! More interesting would be if we find

a ntip on the rth side of rules & a ntip on the left side & suspect they are the same. If so, we have a loop!

I.e. If we have many rules of form  $A_i \rightarrow X B_i$

We may suspect  $\{A_i\}$  &  $\{B_i\}$  is an ntipst. Say we find many n tips by this technique: then we ~~begin with loops~~ ~~that~~ ~~again~~ compare n tips that occur on the left of rules in n tips that occur on the left of rules. We make a M x N matrix that compares these n tips implies (how many common, or % of common) elements).

01:23,40 : 2 "Expo" (in R<sub>2</sub> Area) : Suppose we find a manuscript on /  $\alpha$  center. In the analysis of the Doct., we would postulate a culture created it, & try to build up a picture of the culture from in to into Doct.

Par 2 (b) w. this, we do "pure" analysis of the Doct. itself

05 (SRV) GA : looking at / GA "Sols" to problems! t.

Density of "Intrins" (non-fundamental code) in ~~the~~ ~~flow~~ ~~could~~ give an idea as to how much more efficient (cc-wise) LSrch is — since it (usual) has no "Intrins" — 27.18

04 So, R. former is "Domain Knowledge", the latter — maybe

pure syntactic lang.

(SM) Doct. ~~is~~ variables of "Supervision" in lang:

1) Asking Q's of "teacher" 2) Teacher looks at ~~the~~ performance of T.M. up to now & ~~selects~~ ~~decides~~ a new Corpus on that basis. Teacher can also give sample & sub sample to T.M. to "test" to see if certain concs. have been ~~acquired~~.

25 (SN) on Sol 78 T3: Proof for Radix =  $\frac{1}{2}$ : To show Lemma

is true: It's a " $>$ " inequality, over a certain region; find the lines or boundaries of the region where the inequality is " $=$ ". Then show the first & second curves are of proper sign in other regions of the boundaries (i.e. all other parts within the boundaries).

On Corpus Sequence design: Try this: We have this small corpus, for certain concepts' discovery, this corpus is too small. Also, the pc of these concepts is small w.r.t. to "normal" initial inst. set ( $\equiv$  apird).

We introduce "Hints": These can be examples related to RW (i.e. Domain Knowledge) or any other devices to decrease (or just ~~the~~ cost) of soln.

This "hint" could put an ad-hoc bias in T.M.'s ~~very~~ p.d. — so

This starts "W.D. Domain Knowledge" —  
 L. Randell, Eric  
 IJCAI 87  
 Guiding Construction Induction for the lang. from examples, 4 pp  
 email:  
 erendell@nasa.gov  
 ncsa.uiuc.edu  
 or  
 randell@MCS.uiuc.edu  
 Ph. 217-244-0592  
 Is on NASA Computational ~~biology~~ biology faculty. Biologist  
 U of Illinois at Urbana-Champaign.  
 (Not Conf. for "Supercoping Applications"!  
 1996 (ET)  
 42 48-54  
 conc. lang., using complexity regularization.

W  
 on

To avoid that, we have it recompute ~~it's~~ pd using ~~the~~ an augmented corpus.

I was worried abt "CFC" getting wrong biases by having ~~an~~ "knowledge" being "programmed in" — but here we give "hints" & r. info is tied into the rest of the P.D. using the augmented corpus. "Hints" should be weaker possibly.

— we want TM to put info into its P.D. in its own way.

So: How to start out (?): Maybe a "Blocks world": But with perhaps simpler language (to <sup>start</sup> learn Blocks). The path is then to put in more blocks properties — w/o or w.o. TM's ability to "Experiment" in a simulated "RW" — richer lang w/ Blocks! Then make lang. more like English: Make SS have random errors: Then introduce ~~more complex~~ new things into the world.

Maybe Sim's world, w/ a lot of interacting Robots, but have language to talk to them — w. rewards for correct or useful responses: To be to learn & to let their "World".

T. world needn't be nearly as complex as Simms' world is, but.

18:26.68 (SN) Re: GA: perhaps try constructing "approx" to simple functions like  $(3x+y)/(2x-y)$  via L'srch: Comparative GA (CC).  
 Also One way to solve things like this: Use digital synth (stecklaus) to generate  $(ax+by)/(cx+dy)$ : Then use ~~some~~ special N.L. H.C. method to find ~~best~~ optimum  $a, b, c, d$  for the data set.  
 For L'srch, we have to get reasonable good fit  $\Rightarrow$  as soon as poss. — or else of time limitations (or we go very slowly in T. Stone L'srch). — So maybe do coarse lattice search of param. space, to start off.

01) Say Grammar for entire large corpus were decod: To what extent would TM have much info about T. "Domain" of Y. Corpus? — Well, by "Grammar" I mean as many legit regularities are found as pass. Seems likely that, were this done for a larger corpus, best TM would really "understand" T. domain. — i.e. be able to answer Q's about it.

So: Perhaps: A large ~~text~~ corpus: TM finds some regys using "standard techniques": Since corpus may be "very large", data mining ideas, methods may be relevant. The "TSQ" consists of what might be regarded as "Hints".

Narrower Corpus, etc. — Q's.

Kinds of "Hints"

1) Special Augmentations of Y. Corpus. (perhaps All hints should be put into this format): (subsampled Q/A pairs: from Q's).

2) Some Analog I/O. — perhaps debug some thing as some type I/O.

19) 3) Hard & soft partitioning of the corpus: This could tell TM what some ~~impl~~ ~~exist~~ divisions of T. corpus are.

Many different types of partitionings of T. corpus are possl. E.g. English, French, German, Italian, topics? (topics of Math, Biology, Sociology, History, v.s. low prep. words. : <sup>child</sup> Early v.s. <sup>adult</sup> later learning, — perhaps in ordering of ~~data~~ T. corpus like TSQ order.

Just How is TM to use these "Partitions"? — → (4.3.36)

1) So 1) Augmentations & 2) Partitions (orderings) of corpus are ≥ impl categories of "Hints".

5) Normally TM has a single goal — to compress corpus: "Hints" can take form of other goals. E.g. Various # CS probs., I/VV probs. — using P.D. obtained from T. corpus.

In .ol ff., we have this "text world" (w to "blackworld"), T. <sup>TOP</sup> Goal of TM is to compress (= "understand") this world. The hints are oriented toward that Goal.

Now, in fact, understanding "TW" is not really T. top goal. T. top goal is really that TM should be able to work any probs I give him very fast & very effectively. Div. <sup>leading to</sup> (in understanding) TW, TM gets P.D. — Not can be used for solving arby problems.

Will not so clear that the Pd for R.W. & T. Pd. for problem solving are f.

Same: Actually, they need not be both are conditional proby distrib, so we can have one cond. Pd. — One kind of condition involves compression, (predn) in TW. Another kind of condition involves

97.3  
101  
61  
10x  
161.8 lbs  
186 cm.  
(5.15 1/2")  
9 + 30 electrodes on EEG.

English, French, German?

vol: 28.40 INU p. 105 - 2<sup>nd</sup> paper involves O2 probs.

Try to discuss this approach w/ Marvin, Salfriger, Jargin  
<sup>6?</sup>  
"Fast World"

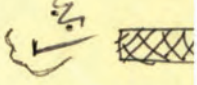
Goals of TW project (NN 84.98 #5): (1) Top Goal: Ability to solve arby, very hard  
puzz, usually O2, maybe some INU. (2) Intermediate: To be able to understand English  
well enough to get some info by reading (human) texts. - So Eng texts (usually for humans)  
can be useful to us. (3) To learn to use other Modes of knowl acquisition! (Asking people  
w/o databases, Experimenting in Math, Computers, in RW (Analyzing experiments in electronics -  
in putting chks together, on deriving LSI chks. - moving electrodes "by hand" in a vacuum envt.  
TM could easily derive new solid state "laws".  $\rightarrow$  it could do in RW, what Koze's PGM did in "Spiral"

The more it knows, the easier it is to teach.

T. ~~very~~ Main Goal is to get TM to pt. where its easy for use to derive  
the envts. for it.

Blu Blu  
Vid Very  
Gray Good  
White Wine

We can start w/ simpler Mathematical lang. Have TM learn its vops:  
Then we slowly  $\uparrow$  complexity of lang. - Move it toward **English**.  
So we can use  $\approx$  (English) **E**peranto) to desc. Math facts. - Try to  
move toward English as its used in Math text books.



An imp't idea that makes writing TSO's easy is the **CJS idea** - it tells which  
ideas are "close" to what other ideas. Perhaps one of the many imp't values of Lsck  
is not so much its "optimality", but the fact that we can use CJS to get upper limit  
on its cc.

Y4  
G5  
B6  
B7  
B8  
B9  
B10  
B11  
B12

One of the Attractive features of TW is that is seems poss. to express  
many "Eng Envs" in this form. The "Compress to Corpus" can be an ongoing  
"Backgut" task for any other problems. [I'm not so sure: Improving P.D.  
is a usual background. -  $\hat{=}$  it is oriented toward solving O2 & Inb. problems.]

One interesting common TSO types: learning to answer Q's after being given  
many example Q, A's. - ideally, I want this to be treatable to a data base -  
so the Q's are about <sup>info</sup> ~~data~~ in the data base - set  $\{Q, A\}$  set does not contain all of this info.  
**AND** could be regarded as a  $\{Q, A\}$  ext. problem (C for induction).

There's proving is an INU problem, but putting in normal heuristics into a P.D. form  
is not trivial!

B  
BR  
R  
OR  
Yol  
Grays

A Good kind of Prelim problem for TW: Ask TM questions ~~about~~ that it would know  
to answer to if it understood Y. Text. e.g. "Object X is Black" occurs in Y. text.  
"we ask: What color is object X" Other Q's could require more or less computation  
on the part of TM.

T. now things I need to do for this talk!

1) Read More on Least, Incr. lang, G.A., so I can contrast what I propose w/  
what's been done. That 1992 paper of Least's U.S. - but say the more recent stuff!

try to use for Least & Incr. lang.

I do have lots of '99 <sup>1999</sup> ~~1990~~ notes for incr. lang: So look opt. ~~at~~ or find refs to pass  
papers into Citations (Harvard Uni Lib) w/o MIT



-01: To what extent is  $(17, 01, -40)$  a adequate soln. of the 02 problem?  
Look at  $(05, 01, -23)$  ... Actually total previous work's:  $105, 01, -25$ ;  $108, 01, -40$

Chris Moore (2) wrote Piz paper showing how an analog dynamic system could simulate an array of ~~the~~ TMC. Unfortunately, the effective memory size of the device was sharply limited by essential analog noise.

Some new ideas on Piz: (1) If there was any regularity in the noise, then the noise wasn't really as big as it looked, & system would have an effectively larger memory than I first originally (still  $\ll 100$  bits.).

(2) Consider other 1. & analog systems w. dead back (recurrent Neural Nets). Probably they, too, could ~~be~~ often simulate a ~~time~~ <sup>time</sup>. — again w. ~~small~~ small amount of memory. — But if we had ~~several~~ <sup>many</sup> such machines, each w. a small amount of memory they could do a lot of computation!  $\geq 10$

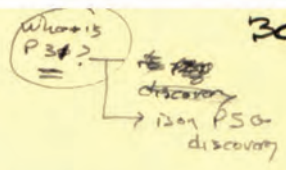
(3) Our idea I had was related to (2) — I was going to use atomic nuclei as computers, w. small units of memory, then link ~~the~~ a enormous no. of them together to make a large, fast machine.

(4) How, say we had a bunch of "Noisy Computers": ~~what~~ <sup>can</sup> could we usefully use them — say, for induction

(5) Another tack: we can make very high speed analog devices: tunnel diodes, etc. It may be possible to make very fast "Analog" computers of mem. — say like 10. They need not be 100 accurate. If they are  $\rightarrow$  universal, perhaps do Larch on them!

(6) Re (2): Recurrent ANN's can probably be UMC's: As such they have a "Halt problem". One way to deal w. this is to produce feedback: which is like reducing memory size or "valence". Or just do Larch & cutoff on a trial ~~time~~ when time exceeds  $\approx 2^k$  d.p.c. of trial.

Assoc. w. the idea of Moore's machine having small mem., because of noise introduced by envt. (.02) & "similar" problem occurs in the "data latency" of quantum computers when they interact w. the envt. — Are these 2 ideas related? Are they the "same"?



I may want to sort of write up. talk now: Then work on details:

Also, perhaps write Journal.

Make list of kinds of training I can point to. "TW" framework; Also Read Ref paper on learning from text books (in 1990 A.I. Conf. (p7-13) <sup>2</sup> Cohen)

Make a list of the supposed "kinds of lang" that are ~~in~~ in A.I. & ML. conferences,

09 The ~~two~~ main features of my approach:

1) A common (format, language, formalism) for many (perhaps all) types of lang, environment. So learning in one modality or domain, can be used in all others. — Also facilitates incremental (lang.)

2) Subclass of  $\Theta$  A very general method of induction that can fit all other forms of learning into a common formal, lang, formalism

3) Possy of easy parallelism for Lurch: Also poss. to divide up tasks betw many sites — out. internet. Or, simply have many diff't groups working on diff't problems.

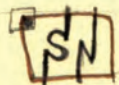
If we use passive topics, then we can integrate ~~the~~ Tsg's of various groups, ~~rather~~ guided by "desirability" (?) or ~~not~~ by separately doing it. Try to Tsg's sequentially on the same TM.

We would want mainly v.g. Tsg's: so most of time ~~we~~ should be spent (in a week, say)

by a group devising t. Tsg. — then  $\ll$  a month spent by TM, ~~to~~ (learning it)

See how Tom Ray does his A.Life on the net.

→ 32.01



I had the idea that looking for common subtrees in the trees that represented past successful skills/functions — was a good idea. — But that it would not often

be useful in small problems. — A "small problem" being one in which its adequate to

Simply combine old useful abs. (if all "useful abs." have appeared  $\Rightarrow$  solns. to past probs).

In "hard problems" / intermediate abs occur (Sometimes w. "1 shot lang")

are useful and to wait until defining them ~~is~~ even better if problem is completed. That Tsg's diff: we have to get partial reinforcement but + problems is complete as in Sanser, Chess, where winning a piece reinforces itself.

So it may well be that it wouldn't get into "hard problems" in this

for a long time!

On the ~~other~~ other hand, discovery of common subtrees could occur

as ex. 103 course of a hour — while analyzing the trees of solns. of many problems.

I had the idea that if we recorded functions used, only in their "parsed forms" (i.e. forms in which

they were generated), it would not be easy to find new pairings for ~~the~~ useful functions, as old.

— So that these new pairings would be a more "compact" (compressed) basis for  $\forall$  functions.

So if 32 were true, we would only ~~have~~ modify <sup>conditional</sup> l./probys for function selection

in "tree growing". To start out, ~~start~~ <sup>After</sup> (+) is selected, we need 2. abs & they are

selected w. probys indep of fact that they feed into (+).

It. Later, we will have a set data so that for many functions, the diff't arguments will each have its own cond. probys for input functions.

These cond. probys are sort of equivalent to defining certain subtrees as being important (E by pc).

→ 32.01



01:25:40 We then try loop rules. Read ~~the~~ assume various  $R_i$ , left  $R_i$  steps are identical/.

02 The simple grammar of 24.31-34, 25.17 may not be adequate; we may need more complex rules like  $A_i \rightarrow X B_i$  } so each "left side" can have an array  
 $\rightarrow Y D_i$  }

04 "left sides".

So, consider some simple stack CFG's; look at the frequency with which various rules occur & see if any techniques would discover them. with wsl loops

02-04 may not be ready for: The single  $A_i \rightarrow X B_i$  rules are the type used in ZCF1...  
 If may be that ~~the~~ rules of this type alone are adequate to form a grammar of any stack CFG (typed or not). w.o. loops, hrr, we need a larger ~~base~~ corpus to justify it. grammar

Since it's a Bigger Grammar

↳ Quite Reasonable Approach.

→ 9.4.98! 1) Read stuff on J.J. Horning Vol III of A.I. encyc. Also, see his paper (6).

He has a way of cutting down search space for grammars w.o. excluding anything desirable.

2) Using ALP, I can do hillclimbing in grammar space: starting either with a d hector from his own grammar. The problem is to find good "Mutation" operators. Hrr. -

4. Mutation idea is too Greedy; try less Greedy: + (2 or 3 mutation) trials. (Different amb of "Look Ahead")

Also, look at mutation operators (or macro operators) that have been successful; try to extrapolate them in mixed ways a/c by "passing a grammar thru them".

See 29.02 on Goals of TW project.

01:30.20 T. talk itself! Start out w. a general dem of what its supposed to do, & some reasons why its a good approach. [T. more it knows, the easier it is to teach]

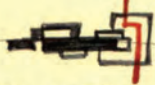
Then give a brief history of my development of the idea: talk about "L-arch formalism" (explain just what it means), Then TSP's five written: what seemed wrong w. them — that they seemed to hard to write.

That ~~not~~ by subgoal in long term is Ability of Machine to learn from human text — But I'd also like other modes of "teaching" of machine to be poss. — A third approach I'm proposing screens, to facilitate this.

### SN Learning Systems that claim to be Universal:

- 1) RLP Sol
- 2) MML walk. (strong ~~MML~~ strict MML)
- 3) MDL Rigs (Predictive MML, Stoch. Complexity)
- 4) Charniak/Kaplan { Support Vector Method? }
- 5) Soar (Rosenbloom/Newell) ["churning"] other verbs.
- 6) EBL (University of Conjectured by ... De Jong?)

Disc of Knowledge Macro Operators!  
 Machine Int p 280  
 Rosenbloom/Novell, 1983 actually



### Talk

### 1) ~ Top Goal of Project: Intelligent Assistant.

Derb. "Int. asst"   
 intelligent

"Experts systems" approach: Brittle, Diff to get into from Experts!   
 Impossible to get into from Experts.   
 ("Unconscious Mind")



Why Bad: How L-arch fixes it.

- 1) brittleness
  - 2) need of data input
  - 3) black box model
- things Experts can't do without.   
 "Unconscious Mind"

### 2) My Early work on this: TSQs: First 3 TSQs. DZ probs, INV probs. The Power of L-arch: T. 2 Gambling house probs.

### 3) Latest Approach: "Text world": Why its Good:

- a) broad based "reading & understanding Text."   
 (compression allows implicit probn. & "understanding" of some kind.)
- b) Also assoc. w. understanding & learning to solve problems.
- c) Very many problems & framing methods can be put into this format:   
 (Gives several examples (some very general: others more narrow).)

Mention D. Cohen's experience.   
 (Cover and key of Radi of text by humans   
 (= Compression).)

Some strange examples More Detailed discussion:

- a) Maybe: what is "optimality" of L-arch for INV problems means.
- b) Why we need ALP (rather than MML or MDL) to solve ~~prob.~~ diff probs.

Just RLP gives us models we wouldn't have had of. otherwise:   
 Gives us a better shot at "unconc. mind" — By considering all available models: If of comparable ~~to~~   
 Give more detailed examples.

Some dozen of Difference betw. ALP & MML, MDL:

By not ~~using~~   
 considering the other problem   
 they are a missing link   
 for neural nets:   
 In some cases: they   
 are ~~not~~   
 may have to modify our   
 to hold advantages of it.

These were ideas obtained by reading papers of others:

Attempts to understand Explain Based (Gen. Eng.)

**EBL**

For De Jong: ML book #2: p 571 ff!

My idea thus far: A "scheme" is a canned procedural sol'n to a ~~set~~ ~~class~~ class of problems.

A (perhaps ~~un~~) idea of Ruiz ~~Roger~~ is: Say we see a single instance of  $t$ . Use of a particular scheme "you solve a particular problem. — we have a "trace" of sol'n.

How do we use this one instance + "background knowledge" to generalize this scheme into a more general rule that can be used in <sup>greater</sup> variety of problems?

For me this may be close <sup>to</sup> (or a special case of) to **IR** problem. Something of interest has occurred. How do we best index this occurrence for easy relevant future retrieval?

This is much less of a commitment than De Jong makes: De Jong would actually use  $t$ . one-shot "learned" rule. I would use learning to retrieve the previous instance and combine it w. present info. to decide what to do. Essentially, De Jong discards  $t$ . Post mortem data after he makes his "inferred" "scheme". I think there is this strong tendency among <sup>(perhaps most professionals well!)</sup> ~~scientists~~ to want to make a decision ~~on~~ post mortem data — then throw away to data.

Another ~ **"I.R."** system is **CBR** (Case Based Reasoning): We have a present problem. We look at  $t$ . past to find "i" problems. We find  $t$ . one that's "closest" to  $t$ . present problem & modify it so it solves  $t$ . present problem: This can be rather simple modulus. or by a "Analogic Xform".

In **EBL** when we solve a problem, we generalize it to solve other related future probs  $N$  to that one. The generalization occurs at real time. Probably some "fitting" occurs ~~at~~.

When  $t$ . Gen. has to be used, ( $\geq$ ).

If so, there is a continuum betw. **CBR** & **EBL**

8.13.98 Perhaps actually start writing to ~~the~~ folk. After its quite clear as to what I will say, send Gamm / Jacy Revision of Abstract. Perhaps make Viewgraphs for talk & put them on Net.

(SN) on "Operationalization" ~ "Explicitation" T. idea is to put somewhat ambiguous, vague, English-blumen statements into exact form so they are directly useable by a liberal-minded Machine. In Cohen's paper (Vol II AAAI 1980 173) he gets  $t$ . Machine to do much of this. "Fuzzy" logic is another way to "operationalize" much of English text.

"Open" is a way of implem. of human heuristic: "State the problem exactly" (which I usually assume has already been done for  $t$ . in  $x$ - $y$  — so  $t$ . Machine's job is to "state  $t$ . problem vaguely" — i.e. find out what its

"close to", so it can know ~~what~~ where to look for "similar solved problems".

✓ Some rats of int! D.H. Fisher: knowl acquire via incremental conceptual clustering.

"Mach Lrng" 2(2): 139-172, Sept 87

J. Gorman, Pat Langley, Doug Fisher: Models of incremental concept formation. "Art. ~~Intell~~ Intell", 4(1-3): 11-61, 1990

Some comments on Cohen: ~~133~~ (33, 35)

"Theories" are rules for bidding: "An overgeneral theory" is one that often works, but often doesn't. This is regarded as being due to its not being narrow enough. Cards of its use: e.g. if condition A is true, then Do B.

We have to narrow down ("specialize" =  $\frac{1}{\text{genz}}$ ) A for a specialized. (whole theory).

What "explanatory" does: in explain of  $\sqrt{x}$  (correct) example is a demo that 1 or more of  $\sqrt{x}$  rules would result in that example. I.e. given a hand of cards & a correct bid: show which rule(s) would give that bid.

Multiple in cons: explanations: I guess this means explanations that give bids & then then rules in  $\sqrt{x}$  examples.

My impress. is that many of the "problems" in Lrng theory are easy to solve, but the terminology used is a real babel that makes it difficult to perceive similarities betw. problems that would make them much easier to solve.

Usually what happens: as soon as I begin to understand what the problem is, I begin to see a clear solution to it.

How, the Cohen (33, 35) has several impl. probs in "lrng from text" that I should understand.

2 kinds of diffys: (1) Undefined terms: They could be partially defined, but in terms so vague to be operational.

(2) Some of the rules given were "over general": they were used in most cases (for most systs) but not all systs.

In (1) we may have to use knowl from RL to get reasonable choices of defns: which are then tried on  $\sqrt{x}$  examples + "rules" to see if they work: A less "RL" dependent method would try minimal information "modifs" of  $\sqrt{x}$  defns that upon being operational — then tried upon  $\sqrt{x}$  examples. 100% score on examples is not neccy. Some examples could be in error.

In (2) we want to minimally modify  $\sqrt{x}$  rules to get not a set that

(2) ~~we~~ have few many errors into examples (b) cover as many of the examples as poss. (hopefully, "all", but not neccy "all")

I Think by "Theory" Cohen means E. complete set of rules + data.

81698 Well: EBL: A possil way to look at it from ALP perspective:

An "explanation" of a data set is a short code for it: Say a 2 part code, w. a rule or law or set of rules is a code for the data in terms of these laws. (for a situation say dec by the code is of pc=1, bc=20: because it fits exactly. Usually its like a Q.A. corpus, w. Q being "free".

Say we have some new data. An EI approach: we want to code it using a rule is the code of the data in terms of the rules: (Consider Q.A. type of corpus):

This old EBL wants to make a new rule that is "close" to the old one - (small Incremental Complexity).

10 **SN** "Incremental Complexity": data:  $A \rightarrow B$  is of small inc. comp. if we can obtain a code for B from a (short) code for A, using a "hy pc x fms":

"T. set of hy pc x fms" is obtained by considering the corpus up to now.

11 A More exactly, the distribution of pc's of these x fms is

Basically, we want a set of pc's on these x fms such that sum of pc of the entire corpus is maximized

these pc's (wts) could be assigned solely on basis of pc of entire corpus -

or they could be assigned so as to max. the pc of the corpus achievable w. available cc. (RLP).

19 We could view the wts of the x fms as like a "TM2 problem" - which is part of the problem of getting "new good ideas" from "old good ideas" -> 36.21

Perhaps a less EI view: Then we want to minz. total cost of coding.

22 We build up a "vocabulary" of prim. operators & x fms of prim. operators & sub operators & primitive operations. (Hvr, we also have to consider pc of corpus w. the operator(S)).

This is an impt idea! It bears on the Q of how good "incremental" lrng is - i.e. best way to do it. One area of interest is Generalized Z141 models: (in which we make models w. increasing goodness, based on previous models, parts of previous models, & standard combination rules).

It seems to be very impt! Its certainly incomplete - - - needs a lot more work - - - but do get back to it!

First try to define, exactly, just what problems are being solved.

32 One problem: We have coded corpus  $C_0$  w. <sup>short</sup>  $p_{gm} P_0$  (thus far).

We get incremental  $C_0 \rightarrow C_1$  How best to search for code for  $C_1$ ?

36 If  $P_0$  is a sequential code (as it must be for ALP, sequential/randomly primd), we just do random increments of  $P_0$ . If as is usually true,  $P_0$  is a 2 part code, we

38 just do random combinns of the second part (the first part is a self del. code & cannot be incremented in a simple way). Hvr, at a certain pt., we decide

that the main bc/control of the incremented part of  $P_1$  is the  $C_0$ .

40 So we begin looking for modules of the first part: which is a "finite object".



01: So 35.40 begins to look more conventional - i.e.  $t$ -measure of similarity of 2 finite objects.  
 This is like 35.10-14 & 35.22: Consider  $A \rightarrow B$ ; we look at  $t$ -code of  $A$ :  
 if it has parts in it (~~reused~~ <sup>stack</sup> functions) we devise  $z$  <sup>stack</sup> (could be a very simple lang.)  
 Then look at  $h$  by  $pc$  - first  $pc$  can produce  $A$  w.  $h$  by  $pc$ ; then  $z$  by  $pc$  objects in  $h$   
 lang are "close to  $A$ ":  $\rightarrow$  28

05: Remember,  $h$ ;  $t$  <sup>top</sup> goal is to get a good code for  $B$ . 35.36-38 is  $t$ -usual "Default" code.

09: We init (look at  $t$ -new sections of 35.36-38) & try to find out just how they didn't fit to functional forms that worked in  $t$ -past. This sort of thing occurs in Physics, when  $t$ -corpus augmentation consists of data in a "new part of  $t$ -universe" - say for extreme parameters. At this point, we may want to try modifying, not only the shortest code for  $C_0$ , but other shortish codes, that are clearly much worse for  $t$ -universe of  $C_0$ .

20:  $h$ ; this "Physics" idea of 09 ff. does involve expansion into a "new part of  $t$ -universe". Just how to proceed is not clear. It seems reasonable to try to go back to more "fundamental" concepts" (i.e. the generators of  $t$ -best models for  $C_0$  — and generators for earlier parts of  $C_0$  ... better we devised  $t$ -latest models for  $t$ -complete  $C_0$ .

21: Perhaps the IMP approach of 35.19 is more correct! But we look at  $t$ -"tree" of how  $t$ -various models for  $C_0$  evolved... including models for just  $t$ -earlier parts of  $C_0$  (~~models used in  $t$ -historical part~~) - corresponds to Schmidho's

24: knocking things off  $t$ -stack in his "Learning how to Learn" perm. (in Dec 93, Jan 94) 3rd. down until 845p.  
 So this is a kind of Vindication of part of  $t$ -system (can impt. part) conceptually — tho I wouldn't really discard  $t$ -top sections of  $t$ -stack! — I would keep them for "information" (Intons?)

28: Also,  $t$ -problem of 35.10 ff is related to what Genetic Alms is all about — Whether there are pure mutations &/c crossover & how it  $\rightarrow$   $t$ . "Building Blocks" hypth. is correct.  $h$ ; in G.A. most of  $t$ -work is not quite rite! i.e. they should be trying to model the variation of  $G$  w. structure of  $t$ -Cand. Instead, they seem to want  $x$ forms that give  $G$ 's that are "close" (can w. large  $\pm$ , but must  $x$  of  $t$  of  $G$  from known Cands.)

One way to look at conventional G.A.: we want to make  $t$ -space of Cands more "continuous", so we can use  $t$ -many Optim. techniques that have been devised for

- So: several (perhaps related problems):
- 1) ~~the~~ Learning theory in physics in view of new data that does n't fit old physics (09-20)
  - 2) Choosing a good conv. conventional G.A. — on basis of knowledge of  $G$ 's of many other Cands.
  - 3) " " " " "  $\leq$  G.A. — (Among non-el. Goal)
  - 4) (Searching for augmented  $t$ -Cand (can now code) for a corpus where one has coded  $t$ -old corpus 0.4. 35.22 of  $t$ )
  - 5)  $\rightarrow$

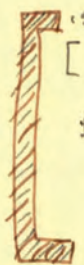
5) We have a TM w. an assoc. p.d.: It solves 1 or more problems. How best to update the p.d.? (Meta Prob) Consider case in which f. p.d. has only "simple hours" that don't look at f. traces" of past prob. solns. (6) Consider f. case of a more general hour in which traces are considered.

SN on G.A. Given past rounds  $[E_i, G_i]$  & their associated  $G_i$ 's.

} close to conventional G.A.

We make a Grammar passing through the  $G_i$  set, but  $G_i$  is given w/ that

.09



is some monotone  $\uparrow$  func of  $G_i$ .  $\sum$  slightly in the direction of SGA. [ Actually, much, in the direction of SGA! — at least my early formulation of it! — In fact, this might be a sort of vindication of a certain class of G.A. method! ]

.12

Well, not so easy! Say we combine 2 rounds w.  $G$  values  $G_1, G_2$ .

FROM AP 4

.13

We expect f. obj to have  $G$  betw  $G_1, G_2$ . In general, it doesn't.

Perhaps analogy of param. adjustment is that  $[.12-.13]$  is apparently true.

We derive a grammar  $\rightarrow$   $Cond_1, Cond_2$  has  $G_1, G_2$  ~~resp~~ resp. We don't even combine  $Cond_1, 2$  directly, lvr., we just use  $G$  each time. But

create  $G$   $Cond_1, 2$  to create a new  $Cond$ , using f. some / stuck / padding that combine  $Cond_1, 2$ . It's much more like creating new  $Cond$  on a basis

.20

of  $G$  entire corpus of  $Cond_1, 2$ , rather than "combining 2  $Cond$ 's".

70  
55  
15  
36  
514

.21

There may be some sort of "fixt pt. Perm" to effect first:

Say we have a set of params  $\vec{\theta}$  that desc. a stoch. gram. Let  $\vec{x}$  be a cond.

Then  $P(\vec{x}, \vec{\theta})$  is the prob of  $\vec{x}$  wrt params,  $\vec{\theta}$ , using some kind of discrete desc. structure. We also have  $G(\vec{x})$  which is given empirically. So I want to

.22

find  $\vec{\theta} \ni P(\vec{x}, \vec{\theta}) \approx G(\vec{x})$ ;  $\approx$  could be that  $P$  is G give about f. some "ordering" of  $\vec{x}$  to  $\vec{x}$  conds: If the conds  $\vec{x}_i$  are big enough, we can use

$\approx$  a distance betw f. P & G functions on  $\vec{x}$ , Per average distance betw.

f. orderings of the  $\vec{x}_i$ . I.e. say  $N_p \vec{x}_i$  is the integro to ling f. ordering of  $\vec{x}_i$  wrt  $G$ , similarly w.  $N_p \vec{x}_i$ ; Then we want adjust  $\vec{\theta}$

$\ni \sum_i (N_p \vec{x}_i - N_p \vec{x}_i) \delta$  is minimal. [ Adm! Adjust  $\delta$ , perhaps in  $\delta$  as one gets closer to optimum.

.30

.21 is f. original SGA! The more general SGA devises a func

$F(\vec{x}_i, \vec{\theta}) \ni F(\vec{x}, \vec{\theta}) \approx G(\vec{x})$  (as .24), with f. additional constant

.32

that for a given  $\vec{\theta}$ , it is easy to find  $\vec{x}$ 's w. large  $F$ .

T. formulation of .30-.32 suggests a criticism of f. original of .21-.24. (i.e. I wanted to use a branching variable rule grammar, probabilized. ( $\vec{\theta}$  is more prob). We get very liberally  $\vec{x}$ 's I set all  $\vec{\theta}$  components (which are usually in  $[0, 1]$ ) to 0 or 1, depending on which was closer. This seems to group into equiv. classes, all param. values  $\ni \delta$ ; all w. param  $\geq .5$  are another. If there are k dimensions of  $\vec{\theta}$ , there are just  $2^k$  equiv. classes. — Well! Maybe not so bad!

There are just  $2^k$  poss. orderings of the system — so making  $2^k$  equiv. classes is  
~~the~~ "fitting"! !

Let us compare ~~mixing~~ mixing  $\begin{pmatrix} x_{s,j} \\ \text{cont}_{s,j} \end{pmatrix}$  v. wts  $\begin{pmatrix} G_{s,j} \\ W_{s,j} \end{pmatrix}$  i. doing this again  
 i again (say on  $\infty$  of times) with the ~~early~~ early Markov GA: passing <sup>stoch</sup> grammar  $G$  thru  
 y. data, using wts  $G_s$  for cond,  $\vec{x}_s$ . It would seem that this G.A. "Mixing" is

a kind of way of getting a stock grammar — Hvr, it is not clear <sup>from this Argument</sup> that it converges

to a correct  $\vec{G}$  | <sup>rather</sup> than just selecting  $\vec{x}_s$  at random.

Well, it ~~may~~ ~~do~~ ~~not~~ ~~do~~ ~~better~~ probly does do better than random, but the big problem is

to see just Why and How!

8-18-98: TM: T. talk itself Expo:  
 29.01  
 29.37 - 30.20 | 32.01 - 40  
 30.09  
 imp. | 29.01  
 essentially 29.01 - 30.40; 32.01 - 40  
 29.01 - 40 also  
 21.01 to 32.40 (w. pp 24, 25 omitted + psobry)  
 13 All on this stuff  
 29.01 it is more narrowly related to talk.

#95 AutoMath by AUTOMATH Long Form. 39  
 94  
 112  
 113  
 138 Selected paper  
 149 in Automath  
 152 1994  
 153 Ed  
 167 P.P. Nederpelt  
 180 J.H. Geurts  
 181  
 (183) selected R.C. de Vrijer  
 papers on Automath. North-Holland 1994  
 185

- A poss. way to write talk:
- 1) first write a very short talk, ~ 10 min. which gives main ideas, & <sup>proper</sup> emphasis of what's important.
  - 2) Then just expand various sections & add sections: But be sure to retain pts of 1. Shorter talk are recorded in 4. longer talk

If 10 min talk is diff: try 5 then 2 min then 1 min.

Perhaps starts w. Abstract.

UPL file  
 from Goun 3  
 ~ July 7, 1998

(186) video  
 N.G. de Bruijn. Home page: de Bruijn  
 www.mu.pad.de/automath/  
 Automath - MUPad developer 13/6/94  
 Upl Paderborn.  
 What MUPad is, is not clear: its anal. manipulation lang. seems. Automath is a psim written by de Bruijn. It is supposed to check proofs in a texts by Landa. (not not Lew)  
 I could use it. It is maybe some of the system to do a inductive Program proof - since the steps in his book were very close.

The title of the talk: "How to Teach a Machine"

Very brief abstract: We now have many effective techniques for machine learning and they have been applied to a wide variety of problem areas. We will discuss efficient ways to use these techniques to bring an initially naive machine to a state of high intelligence.

So put Abstract on screen:

First some definition: By the "State of High Intellig": That the machine will be able to work ~~with~~ problems of interest to me, much faster than I could: - & that it can solve prob that I cannot solve. - that it would do this in ~~some~~ variety of many different areas of sci & technology. To construct a Machine of this type is one of the <sup>main</sup> Goals of A.I.

Initially "naive machine": that it <sup>initially</sup> ~~would~~ has very little knowledge of the domains in which it works problems. It has built into it ~~some~~ a standard techniques for discovering regularities in data.

→ 4.01

8-18-98: T.M. : T. talk / talk Expo:  
 29.01  
 29.37 - 30.20 } 32.01 - 40  
 30.09  
 imp. } 29.01  
 essentially 29.01 - 30.40; 32.01 - 40  
 29.01 - 40 also  
 21.01 to 32.40 (w. pp 24, 25 omitted & possibly)  
 13 All on his stuff  
 29.01 it is more narrowly related to d. talk.

#95 Automath log, AUTOMATH  
 Lang Bonn. 39  
 94  
 112  
 113  
 138 Selected papers  
 149 on Automath  
 152 1994  
 153 Ed  
 167 P.P. Nederpelt  
 180 J.H. Gevels  
 181  
 183 selected R. G. de Vrijer  
 papers on Automath. North-Holland  
 1994  
 185  
 186) verbe

N. G. de Bruijn. Home page:  
 de Bruijn

WWW: mu pad. de / automath /

Automath → MUFed developed at  
 13  
 6  
 19 d.

Urb Paderborn.

What Mufed is, is netclear: its an alg. manipulation  
 lang. ~~text~~  
 Automath is a psun. written by de Bruijn.

It / originally supposed to check proofs in  
 a text by Ladao. (not not Law)

I could use t. data if maybe some of  
 the system to do a inductive Program proof —  
 since the steps in this book were very close.

A poss. way to write talk:

1) first write a very short talk, ~ 10 min.  
 which gives main ideas, & <sup>proper</sup> emphasis of what's  
 important.

2) Then just expand various sections  
 & /o add sections: But be sure t. retain  
 pts of 1. Shorter talk are reformed in t.  
 longer talk

If 10 min talk is diff: try 5 min  
 2 min then 1 min.

Perhaps starts w. Abstract!

DPL file  
 mx  
 from Gounn 3  
 w July 7, 1998

The title of the talk: "How to Teach a Machine"

Very brief abstract: We now have many effective techniques for machine learning and they have been applied to a wide variety of problem areas.

We will discuss efficient ways to use these techniques to bring an initially naive machine to a state of high intelligence.

→ 41.01

See ppt Abstraction Screen:

First some definition: By the "State of High Intelic": That the machine will be able to work ~~with~~ problems of interest to us, much faster than I could: — & that it can solve prob that I cannot solve. — that it would do this in ~~an~~ <sup>an</sup> ~~un~~ <sup>un</sup> ~~der~~ <sup>der</sup> ~~stand~~ <sup>stand</sup> ~~ard~~ <sup>ard</sup> ~~way~~ <sup>way</sup> of many different areas of sci & technology. To construct a Machine of this type is one of the <sup>main</sup> Goals of A.I.

By "Initially naive machine": That it <sup>initially</sup> ~~would~~ has very little knowledge of the domains in which it works problems. It has built into it ~~some~~ <sup>some</sup> ~~basic~~ <sup>basic</sup> ~~structure~~ <sup>structure</sup> & standard techniques for discovering regularities in data. P. Form Sol 78 41.01 →

8.18.98 TM: On Sol 78: genus of genus. (copy → Sol 78 file)

Sol 78: comments: On the proof to Thm. 3 for unbounded sets of finite objects.

more 2

I had the idea that corpus could be  $A_1 \vee A_2 \vee A_3 \dots$

where  $A_i$  are finite binary strings & "V" was a third symbol. A constant part generated off supp is that after V, the generating mechanism must return to the same state.

I then had a generalization of this corpus, so that some of the  $A_i$  were not finite objects, but prefixes of finite or infinite objects. So one Q is: Does have to tell F.M. when an object "ends" if so, whether it is a prefix of a finite object or infinite object?

8.19.98: I could start by considering a corpus of finite objects: all of which have a

229511  
195218  
37293

prefix of a (possibly) infinite object. **Whoops!** This could be a very important

137!

kind of problem: Say SM, HR, ... !!

4:01

For SM: Our kind of possible soln: The corpus is these finite histories of the different stocks. — Not necessarily all of same length. The code consists of a common self-def. code common to all of the stocks, plus many codes for each of the stocks w.r.t. to common code.

Vector  
Graphics

Provides an easier way to think about it! There is only 1 to 5. → this a vector each day.

This is then very simple if all stock histories are of same length. While it's not a final goal soln, it can be used to check a proposed "genetic" soln.

We could have several null components for vector, for times at which stock did not exist.

This **vector** idea is good for capturing the fact that on each common date (day, time) for

stocks have common characteristics: Also that the prices (percentages) of some stocks predict others.

Seems to be imp. idea: But needs work!

→ 52.10

52.20 - .35 is 2  
tentative soln.

The ~~initially~~ Proposed Method: Starting with the naive machine, we give it some / simple problems to solve. Using ~~some~~ <sup>problem specific</sup> builtin knowledge techniques, it is using a builtin ~~general~~ problem solving technique, it solves these simple problems. Having solved these simple problems, it's not so naive. It knows just how to solve ~~problems~~ <sup>of that kind</sup> of problems very rapidly. It may have learned some heuristic techniques by reviewing its own problem-solving experience.

We then give it some new, harder problems. Building on its previous experience, it is able to solve them.

In a similar way, we give it progressively more difficult problems that <sup>eventually</sup> bring it to a state of ~~of~~ <sup>great</sup> ~~problem solving~~ <sup>problem solving</sup> skill.

~~There are~~ **Two Critical Questions**: 1) What is this General problem solving method?

2) How do you design this sequence of problems of increasing difficulty  $\rightarrow$ , so that the machine ends up in this ~~desired~~ state of high intelligence?

The first problem is solved by Levin's search procedure. ~~Each inversion prob~~

The broad class of problems is ~~inverted~~ <sup>inversion</sup> problems. We have a well defined function of, say binary strings. We want to find the input ~~to~~ <sup>that</sup> function that gives a certain output. The P/NP problems of Computational Complexity theory are of this type.

Levin's search procedure

Given  $\{p_i\}_{i=1}^{\infty}$  least ordering of trials  
 w.  $\sum_{i=1}^{\infty} p_i < \infty$  as time for solve.

If the a priori probability of a string being a soln is  $p_i$ , and it takes time  $T_i$  to test this string, then the best ordering way to do trials is  $\frac{T_i}{p_i}$  order <sup>smallest first</sup>. ~~That is~~ if  $p_j$  and  $T_j$  are the probab and  $T_j$  testing time of a solution, the expected time to find best soln is at least  $\sum T_j/p_j$ .

Levin's search procedure is a <sup>good</sup> ~~method~~ <sup>approximate</sup> implementation of searching of this sort.

Normally, "blind search" <sup>even</sup> ~~(even)~~ <sup>with Levin's ~~algorithm~~ <sup>algorithm</sup> to find them</sup> is very inefficient.

But it's possible to use heuristics of various kinds. They can be implemented by suitably modifying the probability distribution,  $p_i$ .

If you have incorporated all of the knowledge you know into that probability distribution, then Levin's search is at least as good as any other search using those heuristics.

Perhaps Give less detail & mainly results: Refer to papers for details.

Say more briefly: at 19. Levin's search procedure uses  <sup>$p_i$</sup>  this a priori probability of a string's being a solution, as well as  $T_i$ , the time to test that string to guide the search order. The trials are in an optimum way.

2489700  
 Q20 1030  
 Southern  
 #15  
 10:30 AM  
 Guilty  
 Goldman.  
 #10 DPAW  
 Boston  
 1253 Camb. S.  
 Camb.

If the probability distribution has been modified ~~to~~ to express ~~various~~ search heuristics, Larr's procedure is at least as good as any ~~stochastic~~ search guided by these heuristics.

~~Guided search~~ can also be used to solve a broad ~~set~~ class of optimization problems — and since the ~~fundamental~~ <sup>discovery</sup> heuristics can be expressed as an optimization problem — we can use L search to discover heuristics

**SN** The idea of "Tools" is cost of developing a tool: e.g. a screwdriver, a hyper-physical accelerator, a very precise weighing machine; The ~~math~~ differential & Integral Calculus. Say you are I. Newton & are mainly interested in physics, rather than math. Yet you take a few yrs off physics to develop th. calculus, because in th. longer run, it will pay you all over productivity.

Now. If you have 10 more years to live and it's going to take  $9\frac{1}{2}$  yrs to develop th. calc. — you'd only have  $\frac{1}{2}$  yr to live & fear physics. (How you could give th. tool to other ~~physic~~ scientists — it'd pay world to have it!)

~~Very good idea~~ A very good idea is usually <sup>very</sup> likely to be accepted & used.

What we ~~do~~ have, then, is a technique for solving a sequence of problems and also a method to discover heuristics ~~to~~ to help solve those problems.

perhaps over Appendixes at end of talk: more details on

<sup>524.15 min or less to Piz, print.</sup> The next Q is: How do we <sup>design</sup> a set of problems that will lead the machine to a ~~high~~ state of h. intelligence?

L search, it's "optimally"

Explain how Concept Nets work; What they are; How used to understand prob. solns.; How used to design TSQ's.

Re: J. methods of Sol 89: I'm not sure Reed/Godd TSQ's can't be written this way: I only know that they're relatively simple methods I've used have turned out to be much more difficult than expected.

3 + 9 yrs live design:

Troubles w. them ① cost of Retrieval from storage was too high.

② system didn't ~~use~~ sub concs. of previous prob. solns.

→ Solns to these: diffy

*SN*

**SN** Perhaps instead of giving a guide video Piz to main ideas! Do Piz more weakly, & give a few pts. in and details, so poor audience can really understand it.

2 troubles w. Sol 89: 1) Cost of retrieval from storage was too high 2) System didn't use sub concs. of prev prob. solns.

R M M

**SN** IDEALLY: In L search! (partial) success or failure at one pt of f. search net, should change pc's at many other pts of f. net!

Just how could this be implemented? Probably best done via time share! But maybe easier by T ← 2T method!

*100*



New Visas in Tng. Envt. Design: <sup>as a TSO for the envt.</sup>

~ Prof. Manor in Israel,

The idea that "say us. ab. probs. mita be Universal" but it wasn't ~~the~~ really a. assisting

for Humans to train

What are some other ways to train Machines? List a few:

- 1) Competitive lang. (In v. A. life, G.A.) each organism is t sq. for another
- 2) In G.A. each open problem is solved as ~~one~~ t sq. of TSO probs.
- 3) Get TM to read English, learn use of human text bodies for training it.
- 4)
- ...

(since 4  
C.O.S. error.)

10

Remarks on Value of <sup>enormous</sup> ~~amount~~ amount of no. of papers of mach lang: ~~large~~

Great variety: Many (if not all) can be put in common format (advantages)

- 2) Common lang. system can use cons. obtained in "other kinds" of lang.
- 3) can be used to <sup>understand</sup> criticize, improve.

126 v  
20K  
2400L  
24M  
3P  
2400K  
48 v  
20.02.

15

**SN** On t. diffy of TM solves a problem using a lower cost soln. than using t. concept that t. trainer wants used. What to do 1) run t. search longer & it may pick up t. desired conc.

2) Use a diffnt prob. ~~to get~~ get TM to acquire t. conc. Try simplest probs; if it always runs out that there is an invariant way to solve t. problems. Then to comp. you want ~~way~~ way is likely to be invariant. - That t. rest of t. TSO can be done ~~using~~ using other, better, cons. If it really is essential to a new class of problems, it would seem stupid that it's impossible to find ways in which they conc. can be induced. If worst comes to worst, find a way to program the conc.

Note: It should be poss. for a human to look at ~~at~~ ~~at~~ TM's soln. to a problem & easily tell what cons were used in it. Solns will in general be "Simple": not like G.A. or ANN "solns to problems

perhaps do .10 first (value of other work)

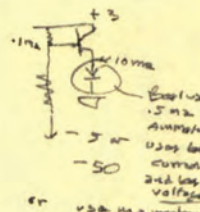
I'd like to be able to take various <sup>existing</sup> ML pgrms & integrate them into TM.

I don't yet see how to do this.

Another way to read thru mach ML literature is try to integrate the

implied heuristics, into TM.

Peter Bergmann  
U. Pauls & G. Schuster



1.5 v; .02 m  
to 100.  
1.5 m. meter may  
be possible to use  
1.5 volts  
5.2K  
60 mA  
312 mv.  
GE 27

36

On "Lang. Lang": Early ~~initial~~ intent "lang lang" is an effectively small vocab. Even if t. vocab. used by Parents is large, there is only a small set of words that are perceived by infants as relevant to it. This "Partitioning of t. Corpus" is a useful hour.

(28.19)

(see 44.01-03)



- 01: 1) Use machine "early/try" <sup>long</sup> to reduce search cc for useful models. The <sup>Past</sup> experimental results are not <sup>very</sup> "snap" — <sup>anecdotal</sup> ~~but~~ much anecdotal, small SS, poor controls, <sup>Bias/observer</sup>.
- 03. Look for <sup>reducing size of search space</sup> for reducing size of search space. 2) for talk — do give examples — <sup>Barwick</sup> ~~try~~ Booth on computers "Universal Mac Logic".

for talk! Main ideas: ① That lots of work has been done on ML: Much of it useful.

- ② ALP (a other Gentile induction models) are able to integrate / integrate, understand, debug this work.
- ③ The main remaining problem is TSQ design. ② Human's Comprehension, friends, & ideas.
- ③ So what can we do w- all this Learning Skill? ② Expert systems that learn; that eventually acquire skills far beyond those of their designers. ④ A Generalized E-S. — T. Scientist's Assistant. I'm mainly interested in this last. ③ Creation of truly "cognitive" machines.
- ④ What has to be done to bring it back?

So I give maybe 20 min talk to Grace (or walk) — <sup>perhaps someone w. dozen of Colloid talk!</sup> <sup>1) Motivation to read (More complex) paper > take home ideas to think, to an Audience, it would be</sup> <sup>Must also perhaps develop</sup> Much longer! It was an intention. <sup>Gentile ideas:</sup> That we now have ~~lots~~ lots of work on ML, much of it useful — but <sup>using all kinds</sup> of different line techniques.

We have some General Understanding of Lang. <sup>Mention Chomsky's objection (1980)</sup> that could be used to integrate that work into a coherent whole ( <sup>discuss advantages of this</sup> ) <sup>Also discuss L Such</sup> (Also discuss some different "General Theories") <sup>Hand Conc. Notes</sup> What ~~can we do with ML?~~ Can we ultimately do w. ML. ?

What would we like to do? ② See ⑩. How do we go about reaching this goal? I'm — <sup>What I myself, am mainly interested in</sup> "The Scientist's Assistant" — <sup>Pro techniques for reaching this goal are common & but for same as those needed to represent a Pro Goals, mainly, to problems, even w. the best of learning systems, we have to present material to the machine in some kind of order of difficulty.</sup>

Difference between human infant & ~~less~~ V.E. lang. Atom: Difference of Bias : prenatal knowledge. <sup>programming</sup> Learning <sup>in</sup> "common sense" is not a good way to do it. ( Most of what is reported as "common sense" is learned after child is born.) Chomsky has suggested <sup>the nature of</sup> ~~what~~ <sup>pre-natal</sup> biases for learning languages.

One subgoal for a very smart Learning Machine, is the ability to read & understand printed text

- ① (to large extent) "understand" printed text. <sup>Criteria for understanding!</sup> <sup>Reasons that ~~learning~~ <sup>learning</sup> then <sup>useful</sup> <sup>understands</sup> <sup>printing</sup></sup>
  - 1) Comprehension → prediction 2) Ability to answer Q's about text. <sup>we want several codes (not just "best") : Suboptimal codes are useful for modifying when subbest code no longer works well. ( & scientist will store several explanations of the world )</sup>
  - 2) Discussion of other work on this problem: <sup>Early</sup> work on story reading & understanding. <sup>(By children)</sup> <sup>By Adults</sup>
  - 3) Blocks world ( a world to ~~talk~~ understand, to talk about ) — <sup>very early</sup> breakdown.
  - 4) D. Chomsky, Learning by reading text on fractal images.
  - 5) Studies of lang. lang. by infants: Theories useful — <sup>to reduce</sup> <sup>Net is "Rich"</sup> <sup>but as</sup> <sup>suggesting</sup> that <sup>effectively</sup> <sup>reduced</sup> <sup>size of</sup> <sup>search</sup> <sup>space.</sup> ( ⑩ )

XML Grammar, Chomsky Hierarchy

That's a sort of talky introduction: Now, I want to give some more exact ideas on how I ~~can~~ go about solving problems:

42" w.c.dram.  
43" w.o.c.dram

Some poss. ways to expand  $\emptyset$  Conc. nets; Lsuch, optimality

② Application to Grammar discovery,

① The effects of a transformational Grammar can perhaps be obtained by having a finite Grammar construct a more complex CFG. This procedure results to total complexity of a corpus + deriv of grammar, would be about the same. The finite grammar need not be a CFG. (or C dependent G)

③ Perhaps describes how I deal w/ problems in D. Cohen's Pgm: EBG, CBR, etc. (case based)

628  
356  
628  
3356

SN "On Government and Binding"

Thesis (Grammar) rules for deciding what pronouns refer to: "who does what to who" in a sentence.

My impression: While there are rules that usually make it unambiguous: Sentences are often ambiguous about this - that one has to look at context & sometimes actually "try out" various assignments to see how they fit, in view of subsequent developments in the text.

I think Chomsky is not interested in this sort of thing. He feels that there is some "Real Grammar" in language that almost always decides these things unambiguously, & that any ambiguity is rare & usually due to inept use of the grammar. Perhaps it would be part of "Discourse Analysis" (not out of analysis of a sentence).

I could describe in more detail how "concept trees" work: how Lsuch finds solns. within time  $\frac{1}{P}$ . Also, just what the problems are & how I expect to solve them. Perhaps discuss EBL & CBR as methods of IR. Maybe "the shot long" (MDL & MML may differ on this).

One way! is to introduce up to "How to teach a machine" Part. Discuss "blocks world" (not a very good). Then talk about the corpus I've written: That's the goal was to get machine to understand Algebra, then begin text eq's, A's about Algebra.

Berman then discuss TSO's I've written: What was wrong w/ them (2 Rings) How to fix those things: "IR" "A better Approach" to start w/ text (Text)

Oxford  
ABC  
Previous

One indication of "Machine Learning": When best researchers are able to Build on previous work. To some extent, this has occurred in Machine Learning, but especially any of that ML work in a common format, with clear procedures on expected error could give a tremendous boost to this area. It gives us a way to compare learning techniques so we can integrate them. It enables us to analyze learning techniques to enhance, to improve them, to understand them.

(SN) R = EBL:

Suppose we solve a problem:  $t_i$  prob, soln. pairs

$R(PR_1, SL_1)$

In EBL, we try to generalize: Make less specific,  $PR_1$ ,

$\rightarrow SL_1$  is skil. soln. ( $SL_1$  is, in general, a function of  $PR_1$ , so if we change  $PR_1$ ,  $t_i$  &  $S_i$ , is ~~to solve~~  $t_i$  same function of  $PR_1$   $\Rightarrow$  a soln of  $PR_1$ ?

In ALP: When a problem occurs,  $PR_2$ , we want to know if "similar" problems have occurred in  $t_i$  past. We have a "matrix", so ~~rather than~~ we look for  $PR_2 \ni |PR_2 - PR_1|$  is "small" this means that  $PR_2 = PR_1 +$  a small no. of bits.

The reason we are interested: ~~is~~ if  $PR_1, SL_1$  occurred in  $t_i$  past, then let us define

$\alpha \equiv PR_1, SL_1$ : Then the pair  $PR_1, SL_1; PR_2, SL_2$  can be compactly coded

$\Rightarrow \alpha, \alpha +$  a short code. So the  $PR_2, SL_2$  is a more than average likely soln to  $PR_2$ .

$T_i$  freq. discn. is a genzn. of simple "bursting" type.

We observe  $AB$  ( $A$  is long <sup>freq</sup> string). Then  $A$  occurs later.

We define  $\beta = AB$ ; s. that  $AB$  and then  $AB$  later have short codes.

Then when  $A$  occurs,  $\beta$  is likely to follow.

In .05-.14,  $\alpha = (PR_1, SL_1)$ :  $\alpha$  is transformant  $PR_2, SL_2$

with a short code, if  $PR_1$  is xfr. to  $PR_2$  w. a short code (~~short code~~).

So, when  $PR_2$  occurs, we look for  $PR_1$  in  $t_i$  past  $\ni |PR_2 - PR_1|$  is small.

How is this done? Well, when  $t_i$  soln,  $(PR_1, SL_1)$  is found, we put in our

(file/index...) all  $PR_j \ni PR_j$  is close to  $PR_1$ . When  $PR_2$  occurs, when  $PR_2$  looks into  $t_i$  past, it finds a  $PR_1$  that was close to a  $PR_1$  in the past & it knows  $PR_1$

$PR_1$  is of interest.

Note the matrix we use, will be domain dependent — but just how is this reflected in the code structure of  $t_i$  metric?

This whole thing is rather unclear! But I think there is something

Basically riter about it. One main ditty is  $t_i$  relationship  $PR_2$  to  $SL_2$ .

Perhaps: in  $t_i$  past  ~~$PR_1$~~   $PL_1, SL_1$  occurs  $\ni SL_1 = F(PL_1)$ .

or, say we have  $(PR_1, SL_1)$ , which is xfrd by short code to  $(PR_2, SL_2)$

$\ni SL_1$  solves  $PR_1$ , then  $SL_2$  is a good trial for  $PR_2$ .

$PR_2, SL_2$  is an analogy of  $PR_1, SL_1$ , if  $t_i$  xfr. from one sol. to other is "cheap"

It will be "cheap" if it has been successfully used in  $t_i$  past as an "analogy generator".

.30-.36 is not BAD — But it sounds "too good". i.e. problems are

not solitary to be solved by analogic xfrns. — I'm not so sure! If  $t_i$ -analogy

is very cheap (i.e. very successful) then it could be very effective in solving new problems.

Group 19  
 $\rightarrow C \& I$  are  
not needed  
on (C:)

Not BAD

.36

So there are several related ideas in ML learning: EBL, Case BR, (Genus & 2 senses Lray) GA & SGA? Related to?

EBL may not be what I'd like to do. It's more like one has solved a problem, & one is able to guess. This soln., so one has solns. to many problems. T. Genus. is done just after 1. soln. so it indeed generality will be immediately available. E.g. Tom Sawyer gets his friends to paint fence by getting friends to perceive it as lots of fun, Genus.: Many poss. 1) way to get work done. 2) By convincing friends to paint for a or by paying them in other ways. How to get others to not find work a pleasure. Etc. We want to also find ways to index T. Genus. so we can solve.

Can be found int. future when they are needed. (NB) As opposed to many EBL rules, I don't want to have "best" genus; I want as many as possible that will be useful in future tasks. So actually: perhaps I'm working on 2 slightly different things:

1) Units of Learning (from a practice). Show how other work on lang. could be done better w/ ALP.

2) Text world as a text world is a text world. I'm also trying to unify these 2 things.

Studying 1) seems to give various techniques useful for 2).

Text world; Also involves study of lang lang in children, (animals?). Theories by others. A.I. & ML work in this area: story understanding. Text book understanding.

Perhaps put various papers & partial paper on "How to...". e.g. "What to work on Novel" @ James of Chem's Universal Grammar (= Grammar Grammar) @ Chris's Minimalism

One aspect of units of problems 1 & 2: In (12) TM could learn to understand some domain: In 2) we could write things about that domain & discuss it. domain w. TM. ("T. more TM knows, t. easier it is to teach").

So, it's getting TM to learn reading: If we have other ways to teach it things - perhaps easier than reading - then use those ways.

SP Re: Math Problems! (82798, 82898)

One kind of TSP could be "proof checking": in which gaps betw. steps will be of varying size. Eventually, TM should be able to discover entire proofs w.o. being given intermediate steps. The main problem here, is to search for heuristics - which are usable in most of math - not nearly int. part they were discovered in. (Note Landau's book: w. various small steps (see ref. in VSS file on

Heuristics: I can get many of them by reading stuff on AI. Proving proving: Then, after a known number of hours, I can decompose problems or "engements" for TM to learn these heuristics.

In a math problem (or proof), the axioms or terms needed may be "given" or the Mathematician has good idea of what terms are relevant. This last is a kind of "IR" problem. - A v.g. Mathematician will have good, non-standard (analogy) ways to do "IR", to guess at which terms.

(or parts of Math) are relevant. Also, the problem of making good conjectures - Conjectures that

the Math Community would find "interesting" is also a good problem in induction that TM could work on. Solving this, being able to do "proofs" would make

TM a v.g. Mathematician.

De Bruijn, htmr Examples of text conversion in Automata. (1968) in Vol 13 of a 1994 book N. Holland Pub. Co. Amazon? K.L. Schokert: JACM 1974 21(3) 436-445 July On P6 in Minzu!

- near start! List of
- ALP-RLP
  - MML
  - MDL
  - C.V. PM; SVM
  - SOAR
  - CYC, Euroco, AM

Universal Learning Systems:

Discussing "Universality" meaning: usually very restricted. Perhaps put into latter to GA.

2) How ALP-RLP, MML, MDL differ!

Emphasis of "The Best Model" in MML: why this is not so good! We want other Models (Procedures) to make new ones! But the "Best" Procedures off. post have some good procedures. One should know about them & use them in trials for new procedures.

MDL tries to avoid complexity by using least complex model approx. This discusses useful data that this data is subjective, it is necessary - makes the iterations work harder, but not impossible.

C.V. - dim, SVM: Even have been described this morning. I'll not sufficiently familiar with them to give any deep criticism. My impression is that the V.C. dimension is a variety of coding and it matches ~~probably~~ to a kind of code length for a concept.

The SOAR & CYC appear to be less probabilistic than Perceptrons, they do use a great variety of induction techniques. Analysing them using any of the previous methods could give useful understanding of both the system analysed and the system used for analysis.

Other Useful Things to Analyse using Universal Induction Systems:

- 1) Meta-Analysis (Analysis of recurrent Neural Nets)
- 2) GA (SGA) varieties of GA. with recombination.
- 3) Chaos Sequences and Prediction of Chaotic Sequences v.s. Somewhat Random Sequences

Say! ALP is Good for 2 Things ① Getting good probability ratios: Also setting good models for understanding Things ② Solving problems via LSrch.

Discuss how it differs from MML, MDL. - Partly in how it deals w. Halting problem.

Studying work of others on ML is useful because they have lots of heuristic techniques that I want TM to learn, & I can "teach" TM these heuristics by giving it problems or training sets that elicit a discovery of these heuristics.

Heuristics (or any kind of logic) can be programmed into TM: This is a valid method of teaching, but should be used only as 'last resort'. (E.g. in solution of 'Peter Baymann's problem' 4305-033) Perhaps explain (to self as well as Audience, why 'real logic' by TM is better than 'soldiering'.

It is because I don't know TM's ~~words~~ <sup>concepts</sup> well enough to ~~know~~ always know the best way for TM to solve a problem. In "real" logic, TM fits a solution in an optimum way, into the constraints & their ~~with~~ <sup>that</sup> TM has at that time. For me to "live in" a solution is very presumptuous on my part - that I know TM better than TM does!

still needs work

01: 47.40 Hvr. - even w.o. t. conjecture skill of 47.37, TM could solve difficult problems i. b. a "V.G. Mathematician".

Actual Doing Actual proofs in Geometry, Algebra: ~~That~~ It's almost certain that these problems ~~in~~ have been formalized (Geometry in particular - Algorithmic) so that it would not be hard to devise esp's.

In some proving programs, if t. CJS is too large I'll have to give hints". Some hints":

Suggested constructions: suggested subgoals (can you show  $\angle A = \angle D$ ?)

Learning how to conjecture auxiliary constructions.

Construction Probs in Math. (17 sided polygon!). Ability to do constructions

Might be useful in TM's conjecturing what constructions might be useful for a proof.

12: 48.30-38 → or Tell TM how to solve it. Look at CJS. i. instead of ~~try~~ doing the search, act as if TM did t. search & found t. ~~so~~ by CJS solve TM needs

(This term (re concept) ... and somewhat by now finding it a very by CJS, it isn't able to use <sup>redundant</sup> t. components of this soln. for problems of lower CJS. )?)

Thus, could leave a "big hole" in t. try. of TM. Hvr, if one does it many times, it may be that one finally "fills in" TM's ignorance. — After all, a human

19 Student "is able to use many ideas that he doesn't understand — but normally 20 t. <sup>better</sup> more he understands t. ideas, the better he is able to use them.

So: Is t. Bergmann problem i. t. solder in soln" related?

Well actually, "soldering in" may not always be so bad! : If I really have a "reasonably accessible"

CJS... Then its really not so bad — Its almost f. same as having TM do t. search.

→ A Big advantage of TM doing f. search is that it checks on t. working of t. system.

"Soldering in" a soln. in which we know t. CJS. i. t. pc's of t. components of t. soln., is not so bad. : What is Bad, is "wiring in" a soln. that seems to be of enormous CJS.

(So that TM is unlikely to be used given c.c.'s of that level for problem solve)

30 So we are essentially "telling" T.M., i. TM doesn't know "why"

"Knowing why" means that suitable sub-goals have been re-enforced (perhaps more important) conditional pc's have been re-enforced: t. conditional pr's are more like heuristics. — They are much narrower in the implications ~~that~~ they

33 give. — They are is "If solns". → 36

The only value to TM of that soln., is t. soln itself. The updating of t. relevant sub-concs. is not useful. (?)

For prob of excessive CJS, we might try a forward, backward chaining search to try to find a lower CJS soln. for TM. It still might be impractically large CJS, but still the re-enforcement of prob's & cond. prob's of its factors would be very useful to TM from t. corroborating elements of t. "Gross" CJS soln. (?)

36 When I would "Memoria" something in college, that I find was really imp. — I'd try to understand" it as well as poss. : The limits of its over-applied, how it could be applied. to other applics. — But not understanding knowing t. heuristic path to t. "something", it's much harder to do. 37-38 — much less "yield." → 51.01

Discuss. w. Alice M. I CBR: We have a present problem: We look into it (suitably indexed) past to find similar situations — Using various types of Metrics. Usually, we pick the "closest" case (nearest neighbor). Or they may pick several cases in past & try to mix them. <sup>using Analogy</sup>  
 At any rate, one tries to mix past cases into present, (if they are not directly usable).  
 Or try to combine several to be relevant to the present case.

Usually we do CBR when we are in a hurry: Mixing several cases is time consuming & so not so often done. [In general retrieving relevant cases is easiest part. X-Phys. Turing to be useful is hard. <sup>move time consuming.</sup> <sup>here?</sup>  
Memory based reasoning (Waltz) seems to be CBR: It uses large data base (data Mining) & fast Machine (say Connection Machine).

EBL: When problem is solved, we try to generalize solution before indexing it, for use in future prob. Solving (<sup>into</sup> ~~past~~ <sup>retrieval</sup>).

In CBR, the "traces" of the previous soln. are retained, so we may more easily discover what is relevant to the present situation.

In both CBR & EBL: One has present problem:

One wants to retrieve past cases for induction to help w. present problem. For IR: EBL generalizes before indexing.

In EBL I think we use various indexing schemes to enable future retrieval.

Now, Having the present & several past relevant cases: How best to make present decision/solve present problem? RLP can deal w. this, in both EBL & CBR.

The problem of EBL is also one that RLP can help with.

We can suggest very many different concepts, depending on expected future problems.

In General: "Expected future problems" Q. is diff.: It may be that an older TM goes over its own past history & reminds the past in view of a better idea of what future problems will occur. This idea of "Gory over the past" occurs in a trivial way when TM updates various (conditional) probabilities in its "Grand P.D."

Colleges  
 Use friends for IR.  
 Tell friends best ideas,  
 Then later, discuss one's  
 problems w. friends —  
 They may remember  
 any similarities or  
 circumstantial developments  
 of them!



01:49:40 : Math Problems: In Geometry Humans have advantage over TM, in having "Geometric Intuition"  
 What to do? Well, TM has advantage over Humans; It can make by accuracy "typical"  
diagrams, to tell w. vary by probab, if a subject to conjectured subgoal is achieved  
 Also, after TM has proved lots of theorems, its totality of heuristic knowledge,  
 will amount to a kind of "Geometric Intuition".

What to forgo. does do, thru, is make TM's intuition unlike defiant from human's  
 .09 So it becomes harder for human's to know how to train TM. → 54.22

.10 O.K. Starting over w. Talk: (Section #1 of Abstract (39.28)  
 We have these 2 different universal methods of Learning; 48.02  
~~What should we do w. them?~~ Discuss what they are, what they are not  
 I will discuss ALP in ~~more~~ detail because I'm  
 most familiar w. it, and customize presentation  
 to be most relevant to second 1/2 of talk  
 What <sup>should</sup> ~~can~~ we do w. them?  
 AR325301  
 EXE  
 This is external  
 on NetX.31  
 STD on By  
 15:15 AM

.12 D) Apply the ~~them~~ to classic problems in induction. (48.17 13 ~ 134) → 54.01  
~~What should we do w. them?~~ Among them, analysis of ~~learning~~. ML in existing literature. CBR, ERL.

Also if one really understands the learning process, one should be able to ~~design~~ ~~design~~  
 make a machine that can learn to do...

I will now discuss Part II: how to use these techniques to... (39.29)

How ML has been mainly non-incremental: ~~By~~ Best performance (koda) very non-incremental  
 How ~~can~~ non-incremental learning. You ~~can~~ want to design a good ~~optimal~~ microscope. So you  
 start out with a human infant, educate him for 25 years, then give him the problem and  
 wait! He is able to design this great microscope. ~~Now~~ Now, you want to design a  
 telescope. You take this very capable engineer who designed the microscope,  
 wipe out all memory of the last 25 years. Then you start from scratch  
 and educate him. In about 25 years, he is ready to design the telescope.  
~~Now~~ You give him the problem and wait! he is able to design the telescope.

Seems to me that we are making inefficient use of the man's talents.

We want to start with simple machine, then feed it to learn to solve  
 simple problems, then use the solutions and ~~then~~ solution techniques of these  
 simple problems to solve more difficult problems.

If we really understand learning, the main problem is to construct  
 these machines, so we can see that being machine from initial ~~state~~ ~~state~~  
 state to ~~state~~ state in evolution (state of high intelligence)

40.40  
 01 ~~40.40~~ } Must detect on learning audit to P.d. ! Perhaps the P.d. is tested or passed by a stochastic operator: For each input, there is a P.d. out to output.  
 03 Given a set of I/O pairs, to construct optimum stochastic operator.

05 Go over my solns for the stochastic operator: There ~~was~~ was a simplification, if I.  
 06 input returned 'no info' (i.e. it was random). Also, say the corpus is a (big) data set  
 sub P<sub>ij</sub> I/O pairs and maybe a few prefixes of infinite ~~sequences~~ sequences!

.20 -  
 .35 may  
 solve P<sub>ij</sub>  
 problems.

Write up a Bibliography of work on P<sub>ij</sub>! I do forget! - put it in some kind of index.

40.40  
 10 40.40 Another kind of corpus! Say a finite set of prefixes of self-ded. codes of finite objects.

Well we can't tell if they are finite! P<sub>ij</sub> may be an infinite point. Say we  
 200 pairs a bunch of prefixes of finite-infinite objects. To extrapolate  
 any off from. ~~reference~~ Some may be complete objects, but we have  
 no way to know P<sub>ij</sub>!

120 A possible way to do .10: We have a univ. w. 2 inputs (1) is self-ded. machine decrn; (2) is random seq. We want the product of the probabilities that  
~~input will produce a prefix of the~~ a prefix of the output will be each of the objects.

i.e. Say A<sub>i</sub> is one of r objects. Then for random input (for a given machine decrn in (1)) the probab. that A<sub>i</sub> will be a prefix of the output is P<sub>i</sub>.  
 Then the probab. of the corpus wrt. that machine decrn. is  $\prod P_i$ .  
 For MD's we want a D<sub>j</sub>  $\Rightarrow$  ~~2<sup>-|D<sub>j</sub>|</sup>~~ 2<sup>-|D<sub>j</sub>|</sup>.  $\prod P_i$  is max.

We then use that  $\sum D_j$  for extrapolation of r. old objects or any new <sup>partial</sup> object.

More exactly using ALP, the PC of the corpus is  $\sum_j 2^{-|D_j|} \prod_i P_{ij}$   
 P<sub>ij</sub> is the probab. assigned to A<sub>i</sub> by D<sub>j</sub>  $\leftarrow$  we sum over all D<sub>j</sub>'s : T. D<sub>j</sub>'s form a prefix set.

The technique of 120-30 can be applied to the other kinds of corpus:

i.e. A<sub>i</sub> can be any other type of object (The various A<sub>i</sub> can be different types): The prefix of an infinite string. A finite object.

essentially  
 the corpus  
 at .05-.06

The "finite object" can be any small or large. It can be a Q, a pair or it can be a large data set (as in .05 R).

So, the Q is: Is .20-.35 "a soln."? does it have desirable properties?

Note Some of the prefixes can be parallel time series: in which case these set of symbols for each common time point can be regarded as a vector and we have a combined stream into a single vector time series.

How do we use the .20-.35 model to extrapolate? Well, we just look for "regularities" in the corpus. These regularities can be within an object or between objects.

A "Big object" (containing, say, a data base) should have many copies in it. The Q, A's will have less inter-object copy, but will have copies based on that in the "Big Object".

In proving 52.20-.35; we do have some diffy: We may not be able to use the

radix 3 corollary of Sol 78 Parm 3. ~~that does it via Radix 3.~~

.05 A (maybe) different kind of proof! ~~Assume~~ Assume: corpus was generated by some machine M

that had a finite descr:  $\tilde{D}$ . The treatment of 52.20-.35 will give a pc. for corpus

.06 that is within  $\epsilon$  a factor C of that given by M. Here C = sum of all ways

.07 that the unc of 52.20-.35 ran simulator M. Then the total expected error in

.08 prog of each of the bits of the entire corpus is bounded by  $\frac{\epsilon}{C}$  in C.

.05 is the beginning of a Theorem: first we show that the factor "C" is as stated.

Then, we might use the proof structure of Sol 78 Parm 3 to ~~prove~~ prove .07-.08.

There appears to be some diffy in the approach to 52.20-.35; i.e. if we are given an entire finite object it was told it is an entire ~~finite~~ object.

U.S. being given part of an infinite object is told it is part of an infinite object

In the first case, we are only allowed self-decl. codes of the object; In the second } T. U.S. Example

case, any codes are ok.

If we are given a partial code of a finite object, then a prefix of a code for an infinite object is illegal... but perhaps this is a finer distinction

to "finite object" can be ~~ably~~ large!

Hvt. ~~it~~ appears to be a problem: Mt. some finite objects can

code differently

To approach the problem: consider the case of "only finite objects allowed in corpus"

It may well be that this corpus has essentially  $\approx$  symbols!

On the other hand: consider 05-.07 if we use the same restrictions (or lack thereof)

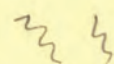
for each object in the corpus before ~~the~~ Descrs  $\tilde{D}$  is the set  $D_1$ , maybe we get the same C .06 anyway! So maybe we don't need the ~~radix 3~~ radix 3 Parm!

Ⓢ We may be able to modify corr. of Sol 78 PT3 so it applies to this case! - we. radix 3! Ⓢ

Anyway, in view of (.05-.08) it would seem that the Parm of 52.20-.35 is a book rite.

Prob "stop" status for self-decl codes: For user to recognize "stop" state, the system must have an extra symbol. (Willis used extra symbols of this sort in his 1970 paper)

We may be able to modify the Corollary so it is radix independent it works for a great variety of devices that compute proby. (i.e. a great variety of "Corpus Mixing") (i.e. it can usually mix a great variety of different kinds of Corpi)



01: 5/12 After reviewing those General Ling. Models: What continues needs to be done?

In almost all of ML ~~works~~ <sup>Research</sup> there is a general

General underneath of Genl. Solns. to Ling. Problem. — ~~with~~ <sup>different</sup> ~~instead~~ <sup>factor</sup> ~~factor~~ due to enormous variety of techniques used — each laboriously constructed on largely intuitive bases — All of this work could be done much ~~easier~~ <sup>more</sup> accurately if the ~~the~~ ~~techniques~~ were generated from General Induction System ~~that~~ known process. We now have the "Maxwell's equations" for all of learning and discovery. Any ~~exact~~ correct work on Ling ~~could~~ <sup>can be developed</sup> from these principles

To some extent, this has been done. In a few areas of Machine Learning, very general techniques have been applied. ~~eg.~~ <sup>eg.</sup> decision trees, ~~linear & nonlinear regression.~~ linear & nonlinear regression.

.368  
x = min  
~ .69220  
.36  
.37  
.38  
.367  
.369  
.3675<sup>3676</sup>  
.6922  
00627  
.3679

A ~~border~~ <sup>bracket</sup> approach, just to 5/1, 10-12 i.e. "Real we have all these general induction methods (maybe describe how they differ (a bit))"

But still, we don't know we have many plans, that can learn to solve probs in specific setting w. certain basic info.

These 3 are the tools we have available to us: ~~with~~ <sup>2/3</sup> [ I will mention that "T. tools" have <sup>been</sup> used to ~~so~~ <sup>so</sup> extent in lang tasks, but very little on interpreted learn tasks of any magnitude.

-22: 5/109 — Math proofs: The steps in math proofs often are quite large CJS. Verifying the proof is often easy, if the steps are close, it just amounts to checking that the theorems combine properly to give the desired result.

Now, after doing T. search (or simulated "search") for the various steps in many proofs, TM might learn how to reduce their cost. It might amount to long/learnable <sup>usable</sup> heuristics.

-27 The induction problem would be: Given order / Pairs:

- (a) stuff proved this far + theorems, postulates; also next step to be proved.
  - (b) The ~~can~~ actual proof. (i.e. = programmed etc)
- say S.M. was given many such pairs: it could then induce relation of (a) ~~to~~ <sup>to</sup> (b) given (a) get a good Prod. for (b).

Looks Good!

Then, after T.M. has learned to do these very simple proofs, I can use the same (a), (b) pairs to get it to learn slightly longer proofs: say, proofs taking > 1 step. Also, I will get TM to learn proof methods by observing complete proofs — Also observing traces of proofs being done by real people or by programs that know how to do some proofs.

An implication here is that I can slowly increase the difficulty of the problems, (this set of examples could result in CBR : It could pick nearest neighbor " 5/5/01 spec.

9.1.98 TM Expo Math probs

01: 54.40 <sup>spec</sup> a use a rather complex analytic x plan to decide that a neighbor is "close".

Also use of a combination of many neighbors. ~~So~~ Indeed, one might regard this whole induction trick of 54.27 as a "genzel" of CBR.

Judy!  
walk!

Another kind of problem of Graded difficulty <sup>2b/c</sup> that can be solved in  $n$  ways;

Transitions betw computer layers: e.g. Basic  $\leftrightarrow$  Fortran; Fortran  $\leftrightarrow$  C, C  $\leftrightarrow$  Machine, at interest.

Say we have some machine code output of C (or other compiler), to reverse f.c. source code.

~~Putting~~ disassembly of Mach code is putting in useful comments.

Disassembly of compiled Cobol is putting in comments.

Translation ~~between~~ of ~~programs~~ betw different Machine layers. Attempts to optimize for speed, memory, etc.

112 So TM can learn <sup>in</sup> many different "domains": putting all knowledge into a common P.D.

{ remember problems all conditional probs; so we can use simil cond. p.d. for all variety

of problems — (Spec 52.01 ff. : There is one earlier more general discussion, but I haven't been able to find it. <sup>52.01-03 is basic idea.</sup> A recent idea was that a cond. p.d. can be

expressed as a stochastic operator. The input is "t. condition": t. output is a p.d. over t. desired space. 52.01-03 is about stochastic operators of this kind.

My talk will discuss a few domains: (  Proof checking  $\rightarrow$  Run proving; Genl math problems

Language Compilers — Translator "Compiler"; "Tant word": Ability ~~to~~ understand answer Q's about Math (its domain of knowledge)

~~the~~ Input sub-goal: Getting TM to be able to read (for humans) text books &  $\text{TeX}$  & fair amt. of info from them.

A very imp. idea is that t. ALP. (cond. prob, stochastic operator) approach enables us

27 to optimally mix learned info from many domains → 58.01

In my old approach to TSG design, every time TM solved a problem, it

29 would incorporate that soln (a perhaps its "traces") into the P.D. ... This was the only way it lrd. Now, t. lang of 29 is only one kind of example that TM uses for a corpus <sup>to help</sup> to "enrich" its P.D. We could use

only external examples, so it isn't necessary to have TM solve problems only external examples, so it isn't necessary to have TM solve problems (while its learning to solve problems) ... but  using (comprehending) problems solved by TM has advantages ① Tells us how good TM is lrg.  — so

we can modify the TSG if necessary ② TM has to house traces of its own solns, external <sup>spec.</sup>  $\rightarrow$  56.31

Say we have a Big corpus of "facts", etc plus a bunch of ~~Q, A's~~ Q, A's.

Since the ~~relation~~ relation of A to Q depends on info in Big corpus, the only way we can get

info transfer is to have a complete dcm or Big corpus (all of its info) in some machine  $\Sigma$  (52.20) w. ~~self~~ self delimiting dcm's [Dj]. We could then ask for

a cond. kernel pd. (stack, ~~parameter induction~~), so we put in Q in  $\Sigma$  w. some P.D.

of ~~the~~ Answer to that Q.

In a more general T.M., we will have a conditional conditional P.D. i.e. a cond. kernel:

One is that we are in Q.A. mode Two is just what Q was asked.

It may be that we ~~get~~ <sup>almost</sup> get all of the "Big corpus of .03" into machine Dj

only after we have enough Q's & A's to cover the entire <sup>Big</sup> corpus.

So, often, if a Q is about a part of the <sup>Big</sup> corpus never asked about, it will be hard to locate

(in code user now) Dj to get any info about that Q into Dj. We have similar problems w. "Dcm

Shot / ring" We will have to make special computational devices to deal w. this problem. (31)

Another set of problems of "gradable ditty" M.T. from our computer lang. to another.

Start out w. simple ~~and~~ pairs of lang: that are very close: then is. dcm is lower close to machine code.

Eventually Eng Basic  $\rightarrow$  Fortran; Fortran  $\rightarrow$  C: A useful problem (very high level) ~~is~~ disassembling

Machine code  $\rightarrow$  Commented Source Code.

This might be a useful ~~problem~~ "hard" problem for reading English text and understanding it ( $\rightarrow$  translated into usable internal language). Being able to answer Q's

about English text could involve trying. Compressing text would involve ~~short~~ short internal codes for commonly occurring objects in the English text.

SN HASH CODING is similar to ZATO CODING in that we always get at least what we want, but sometimes more. [Also the mathematical details are similar.]

In § 55.86 I had this idea of compressing the whole P.D. every time I solved

a new problem (or in lines 55.29-40 - every time I was given ~~with~~ a problem soln

or useful "example"): Actually, this wouldn't work because the "shot / ring" effect of it is!

we would have to deal w. this in a special way. In general the DSH problem is usually

mainly a IR problem: We only have 2 cases to compress! Coming from the larger corpus are

the indexing techniques: the factoring of an "example" into a structure mult. by parameters" in regard to examples.

Is this poss? For a Q, A T.M.: We feed in Q in  $\Sigma$  the output is the P.D. of related Q's: (Not so good: It wouldn't help retrieve stuff from "Big corpus" that hasn't been subject of previous Q's).

01:50:40: Re: CBR; ~~related to~~ Generalized (as well as the way it is normally used in current ML community) ■ can be very useful for TM & for ME!

Japan Wiki  
2003 AM: -1.80%

i.e. I can try to index various ideas I know; also books, objects, HW, etc.

This "indexing" can be ~~practical~~ of both practical value for my existing memory & of Research interest for TM research!

It's the General IR problem. - having ~~two aspects~~ 2 Aspects to pure P.C. & C.C.

The interesting thing is that ALP does give to ~~an~~ form of general soln.: it ~~does~~ tells just how data is to be treated, & gives the signature, (meaning) and relevant conditional PC.

We have always 2 problems: ① inserting (indexing) to data ② Retrieval of data.

OL problems: This idea of cond. probty  $\rightarrow$  Stochastic Operator may be

How to deal w. practical aspect of "is authority" OSL

a real breakthrough here! The idea of having the same optzn. routine for all problems seems unreasonable: I had some idea of using a listing and of all optzn. points is somehow finding which energy could be used for each problem.

The idea ~~of~~ have it: construct a stoch. operator in which I would have to problem to be solved  $\Rightarrow$

input  $\rightarrow$  2 P.D. on all optzn. techniques as output. This P.D. can then be used for LSrch.

Dr Tomie. 354-2400

20 DEF Again, there is the OSL (shifting) problem. Actually, the problem is more general than "1 shot": In 1 shot long. the difference between  $SSZ=1$  (w. OSL) and

$SSZ=1+$  is rather large; But even w.  $SSZ=2$ , we may not have enough cases to warrant a definition, but  $SSZ=2$  may be enough! So the usual way of making a P.D. based on the post (w. definitions is justified by their being useful in post) will be often inadequate. We could always add 1 to the apparent  $SSZ$ , but this would seem to be impractical for  $SSZ=1$ , since almost all of them would not be useful — i.e. they would not occur in the "Question".

So in general, I need some good ways to deal w. (2) OSL:

IR: "Indexing"  $\rightarrow$  code "a.l." way.....

WBR may not count; comprehension ASB is OSL only many.

One way to deal w. OSL-type problems: say we have an unordered set of complete objects and one partially complete object. We want to P.D. for the rest of the partially complete object. We may also have a "Big Corpus" object.

36 code  $\rightarrow$  say we make a code for all of the complete objects! Then we have some "heuristic" routines for modifying (augmenting) the code w/ the partial object — that automatically gives us a variety of combinations.

Look at what's actually done in various practical OSL situations, see how they did it. 36

01: 55:27: Section of talk on TSDs: TSD's discussed:

Microprogramming ~~...~~ Theorem proving: Gentile's pubs!  
↳ Book MT of computer langs: "Text world"  
↳ Data Auto Math

Other than Devn. of problems <sup>consequences</sup> Discussion of Methodology just how "proof checking" is hard.

Actually, it's not "checking"; but long to do small steps (Micro proofs!). - But proofs used as source must have steps that are verifiable by TM.

We can't test w. 1 step proofs, Run  $\rightarrow$  step, Run  $\rightarrow$  step ... etc.

Also: Simplest if ~~...~~ amount of available info to implement a step is small.

Say w. ~~...~~ only a few. postulates is running.

For longer, more complex proofs, usually a certain amt. of planning is required. We can give the machine examples of plans and of the traces of people or machines implementing these plans.

**SN: T.** So what share codes is nominally an O2 problem. It can, however, be viewed as a INV problem, if the machine tries code in best order. If the machine has an adequate P.d. to discover the code a short code.

Perhaps Methodism SM as source of TSD's: Start in early date (easy to find rays) Then as data continues rays get harder & harder to find.

Discuss predictability of the SM.

Book 45:01  
↳ talk  
(in Mexu!)

In talk at some point discuss condi. probab. using sets of finite objects.

Can't say more talk.

show First discuss sup. extrapolation as a kind of condi. probab. - Then discuss

Q, A's, from Harristick, | examples! (sets of Dinosaur bones) - extrapol.

Davis 1982  
Completed 1976  
C.C. started ~1984  
Newton  
1692-1727  
C.C.

23 A post. way to order trials for ratio of probab.  $\frac{P_1}{P_2}$ :

Ordering is Least of  $\frac{CC_1 + CC_2}{2^{-N_1} \cdot 2^{-N_2}}$  :  $CC_1$  are times or cc's spent therefor.  
 $N_1$  are lengths of codes (in bits) tried

No!  $CC_1$  &  $CC_2$  are amt. of cc on a particular trial.  
 $N_1$  &  $N_2$  are lengths of a particular trial. )  $N_1$  ~~is~~ conditional probability of  $i$  codes given  $i$  previous trials for string  $i$  (say) given string  $i$  is one's previous experience on code ending.

Say we have a numerator trial  $CC_1$ ,  $N_1$  variables p.c. of  $P_1$

↳ While it's a 2-dimensional search, 24 may order trials.  
My mind is not yet clear on this!

perhaps discuss differences betw. ALP & Computational Ling Theory;

ARE  $\neq$  CLT is concerned w. Learnability of "concepts" - What's the difference?

They are learnable to a certain probability of error, given a sample size and PAC Learning asks if the error probability in learning a concept is polynomial in sample size and computation cost.

ALP is concerned w. a more related problem: Given some finite amount of data, we can make various statements about possible extrapolations of that data.

Given sufficient amount of time ALP can give approximate answers for any of those extrapolations. Given more time, new problems will be obtained - and are expected to be better.



We don't particularly care if the expected errors polynomial is simple since completion time.  
 We just want the best values we can get and we estimate of <sup>very practical</sup> program <sup>proven to be</sup> so we can use the estimate for.  
 Very often, the ALP problem is much easier to solve than the <sup>practical</sup> <sup>proven to be</sup> <sup>program</sup> <sup>proven to be</sup> <sup>practical</sup> learning problem.

Probably best to deal w. even V.C. very briefly.

On learning hours: The hours must form a concept not that has acceptable  $\epsilon$ 's for each kind hour.

D: Winn days  
 \ commute  
 \ Ausi.

So talk: Start w. 2 ideas: (1) Existence of Gauss Long Recursion.

- (1) to TSCQ problem.
- (2) Discuss Gauss Long Recursion.
- (3) Then discuss TNG Spgs in more detail

7.86  
 +1.52

215.98

-20.10  
 -50.44  
 ~ 10-1  
 +3

(1) Gauss Long Recursion: List on screen

Sol  
 well  
 vrs

- V.C. ~~V.C.~~ S.V.M
- Recurrent ANN
- Spgs
- Latest AM. Error, Cyc.

4.988  
 99.299  
 5.689

Cyc

Maybe abandon Recurrent ANN: Not strictly needed, use of a base technique.

Discuss Difference betw. Sol, Wal, rds: Emphasis on how Sol doesn't cancel out.  
 very diff problems. Ex. Sociology, Economics, Geology. Also Marking able to do Pen S &

done by Subconscious Mind. - Discuss this later. Perhaps Monday Report Penrose.

Maybe Marking  
 Rigor Penrose.

(2) The Training Sequence problem:

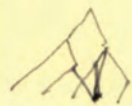
Compare as "Idiot Servant". The harder it is, the easier it is to teach.

5:44 ± 30

The computer is inherently unwisely not has an incredible lack of sentience: "it's an 'idiot servant'".

019

Derive my early ideas on TSCQ design: Conc. (note) almost = Prac. (how a Conc. was not a Computing Theory Conc.)



L search: Gamb House Arms:  $\frac{cc}{pc}$  optimum Chdng: Time is  $< \frac{cc}{pc}$ : 5.79 ± 35

If all hours are in P.D. then technically with manufacture of optimum.

So TSCQ's easy to write! :

Actually not so easy but I did write 3 in varying % of completion:

- AMS-Nature Invy.
- Linear op. Solving w Invy rules of Algebra.
- Solving linear non quadratic non cubic opns.

3 problems w. more TSCQ's (1) hard to write (2) no use of sub solns or only solns.

(3) Didn't scale: Many use of in times so eventually it was unusable

Solns: (1) lower prob (2) use of "IR" (used much in Prac ML work (BR, EBL<sup>F</sup>...))

Or: Other Modalities of Invy:

Some examples w/ extrapolation <sup>set of</sup> examples; We want to find a short coded machine <sup>mapping</sup> algorithm that assigns <sup>Mapping</sup>  $\{x, y, z\}$  into outputs by p.c. to examples. Or a ~~map~~ short coded operator that maps inputs of ~~examples~~ into outputs. How do we find this operator? I guess by Search; over heuristics & components of heuristic. So we can have a training sequence of sets of examples or heuristics machine updates its current self for every new examples  $(k) k=1, 2, \dots$ .

Perhaps we want ("Meta Heuristics"); that Enable TM to find ways to (modify/adjustment) short codes of the machine that probabilistically maps from input to output of Examples.

The interesting part as part of cond. p.d.'s  $\rightarrow$  search operator, is that it makes the problem clearer — makes it expressible as a well defined problem. It makes clearer how to various parts of the corpus are "allowed" to influence the output.

Any way: Emphasize that the cond. proby Idea enables us to merge info, recys, sub-recys from disjunct sub-corpi. — This is a characteristic of "Human" characteristics.

Actually, the non operator form of ALP also enables us to do this, but, I think "cond. probys" are more commonly used: the single sequential corpus in the original ALP is certainly useful, but much human induction is on sets of finite objects.

So to Q is: Just how does TM derive these "short codes"? see (01-09)

clip  
28, 25  
31  
Kojan...

SN T. problem of deriving a cond. p.d. can be solved exactly by the exact soln. of the unordered sub-extrapol. problem" or, approximately (assuming no info in the inputs) by finding a short coded operator that maps  $I \rightarrow O$  w. max p.c. If we have a large common corpus, we can automatically include it w. each input to the search operator.

6587  
711

Example; we have the text of an encyclopedia: we want TM to answer Q's about the text:  $\rightarrow$  A training example consist of a triple!

- ① The ency text
- ② A question about  $z$
- ③ A possibl. answer to  $z$ .

We give TM a large set of these examples one way to train the TM would be to give it a large set of such examples. If TM tries to find a short code for an  $(x, y)$  that can map from  $(1, 2)$  to  $3$  (i.e. given 1, 2 to get 3). More exactly we want a search operator w. a short code that assigns a hy probability to the known corpus.

Another Probably an easier ~~to solve~~ way to solve the same problem, Get TM able to do the same thing is to start with 1 or 2 examples; have TM solve from them, then give a new example; have TM modify its old code to deal w. the new example; give it newer examples etc.

Here if we careful design our seq. of examples and give TM a good set of primitive operators, it will be able to "track" the new examples.

We can make TM's work simpler by giving "Hints": ~~the~~ Hints

9.5.98 : T.M. : on the probability of  $t$ -integers :

One way to deal w. not truly problems but a usually closely related problem :

We have a sequence of binary strings <sup>sub</sup> & we want to mark ~~some~~ boundaries  
them. One way is to have a third symbol  $\Delta$ .

Say the binary strings are of max length  $k$  : & there are  $n$  of them in 1. Comp. :  
By using  $\Delta$  as marker, we have  $k \cdot n$  total of binary strings,  $n$   $\Delta$ 's : So

$k \cdot n + n$  total no. of symbols,  $n$  of them are  $\Delta$  so p.c. of next  $\Delta$  is  $\frac{n+1}{k \cdot n + n + 3}$

$\frac{n+1}{(k+1)n+3} \sim \frac{1}{k+1}$  for large  $n$  ( $n \rightarrow \infty$ ).

← by Markov's principle

This is O.K. if we have a prior comp already. If we don't, then  
it's a sub string : i.e.  $n=1$  ; So we have 3 symbols : The binary string or ones, or zeros, or

The p.c. of  $\Delta$  on prior and  $\Delta$  =  $\frac{1}{k+1+3} = \frac{1}{k+4}$ .

T. lang. ~~model~~ model would be obtained by assuming uniform prior for all 3 symbols.

— using integration, we get "Laplace's rule"



Spec 01:00.40.

Hints are info that makes it easier to solve a problem. In the present case, a 'hint' would be information that effectively reduces the size of the space to be searched.

Another problem, that is formally similar to the Encyc QA problem is the problem of learning to prove ~~the~~ Mathematical Propositions. Here, the examples consist of the postulates and proved theorems available to be used in the proof. The second part, is the theorem to be proved.

The third part, is the proof.

We start out with very simple, 1 step proofs, and as TM becomes more skilled we give it progressively more difficult proofs. For each example, as in the encyclopaedia, the machine tries to find a stack operation that has a short code, that gives max prob to the known set of examples.

$$Code = 2^{-|D|} \cdot \prod P_i \quad (\text{for set of finite objects})$$

$$\text{or } 2^{-|D|} \prod P(O_i | I_i)$$

1.27

1.26

40K

300.

1000

20M

ABCD

present proofs first, then Encyc.

for Encyc. we make ~~gradually~~ start with a very minimal set of statements in the ency, written in a very simple formal lang. We gradually modify the lang. and ~~the~~ increase amount of info in the ency, & make our Q's more complex. The important thing is ~~having a measure of how difficult~~ <sup>if we</sup> ~~each new problem is~~ knowing the state of the machine we have a measure of how much time it will take to solve each machine problem. — So we can tailor our ~~problems~~ problem steps (Per CJS) to be within acceptable bounds

Perhaps just Monom  $\rightarrow$  solution of ~~the~~ optm problems using techniques similar to G.A. can also furnish TSP's.

(Cosy ramona: Just using many holes of TILES of the same conditional complexity P.P. can give a <sup>Broad</sup> variety of scope of intelligence needed for the usually hard problems! <sup>increasingly use.</sup> Another thing: just. ALP isn't restricted, so kind of functions people normally

( If time, discuss Roger Penrose )

Amur Amur

I was worried about Not choosing an adequately good TSP to start TM's "rise to the top" & My present thinking is that it doesn't matter much; just one just teaches TM whatever one can — real ordering of topics is not critical. — we just "throw it all in" & TM is able to integrate it. (The "kitchen sink" approach) I'm still wary of "kitchen sink" but TM may have to be fairly advanced before one can be so cavalier about what one feeds it!

I had been concerned w. my early (E) TSP's because I didn't see a reasonable, easy way to continue from: w. the present approach, we can continue in practically any direction!

T. Reingold's "being lead" are such heuristics, & they are valid for most any "Domain".

Well O.K. I see that (perhaps) TM could discover various such heuristics - (like <sup>backward</sup> forward) changes,  $\alpha\beta$ , GPs, etc: But how does this fit into the Universal Cond. P.D.?

It may be getting TM<sub>1</sub> & TM<sub>2</sub> mixed up!

07 Dad The Un. C.P.D. (UCPD) <sup>Univ. Cond. Prob. Distrib.</sup> can use itself as a setn. — Say to do a forward comp from A to B, ~~then~~ & then suggests signal "C", then was to do the UCPD CPD

to find out trials for going from A to C, then, ~~then~~ from C to B.

10 TM's General ~~name~~ <sup>JOB</sup> name type, is to find short codes for sets of objects:

Usually this ~~type~~ results in a CPD — of 1. form; for given input, fr output is a set of objects in no order of PC, w/ PC's given.

Try Printerink or "Pilot" pen!

(Q: maybe  $\frac{CC}{PC}$  order?)

ABCD

The problem of searching for short codes can be (recursively) expressed as

a CPD: i.e. given a certain corpus, which ~~order~~ is ~~the~~

← good trial for a short code → output ← (would we want this in reverse order?)

Given a <sup>meta-level</sup> set of <sup>dom</sup> {Problem, (soln, trace)} that TM induces a CPD so that ~~it~~ for any problem, it will propose various <sup>soln</sup> {trace}. (Note however PC order; but PC's will be assigned to each output given)

In TSP's for this proposed device: It can go from one domain to another, in any order: If it isn't doing well in one domain, we have it work on another.

Later it can go back to formerly diff't domains, that will now be easier.

One input domain, of course, is TM's mind: Generalizing, (finding better) heuristics.

Well, not so fast! In the CPD, we just generalize from past hours:

we don't really devel <sup>better</sup> new heuristics, Devel of new heuristics is a OZ problem.

31 {The it will devel. new heuristics by "diffusion" as G.A. search does (Optimization by Diffusion" Also Sim Anneal does this)

IN General, however, OZ probs are solved by CPD; i.e. we put in the problem: The outputs

33 (w. assoc. probs) are various optzn. techniques. Note that "finding short codes" is an OZ problem. — so TM uses 31-33 on such probs.

Theorem proving is an INV prob. — Here, it is a "Problem", so it can use to general C.P.D. to ~~the~~ order its searches

Q, A with a corpus is pretty much an OZ problem — i.e. we want as good a code for the corpus as possible. But, its soln. <sup>will</sup> can be part of TM's General CPD.

O.K. : So Just how Does <sup>Pris</sup> TM Work? It has a CPD to start w. some ~~initial~~ <sup>initial</sup> info.

Initial info may be a set of optn. parms, with an uncond. probly on their use.   
 ~~We give some~~ Actually, we may ~~use~~ "seed in" a lot of info.   
 TM then just uses Pris info along w. new examples & other new info, to merge   
 a Better ~~CPD~~ CPD.

↳ we may want to "factor" f. of f. ex. pms  
↳ f. of ex. pms  
↳ (uncond.) stoch.  
↳ grammar

It might be possl. for Pris system to use CXC's data .... The whether Pris would beat day value is unclear! I don't know ~~probably~~ vary little about CXC & what kinds of Ruler it has & what form its it.

110

Describe TM's "steady state", then describe an initial state (Tabular Rese @?) and some ~~initial~~ startup training that will put it into correct "steady state".

114

One Q, is f. form of CPD: ideally its a stoch operator; but it <sup>is</sup> ~~is~~ created/modified by considerations of code length of B dem. [its input is a problem dem, its output

will be trial solns on a. assoc. PC's; For INV probs, probability outputs is f. soln. only   
 Maybe not: We could just ask for a GP of solns (approx) with an implied L search to find soln.   
 Hvr. in good L search, the PC's ~~change~~ change as one works on f. problem.

Doing for a while, its response to "INV" probs.

SW There is some ambiguity in the defn. of "CPD". In general, the output will depend on but ~~never~~ allow CB. (RLP) Do we have to either give a Time limit, or accept a "Anytime" soln? Well, in the case of INV problems, no CB is needed: it just works until it solves the problem (if ever!)

122

Start again on "Steady State" dem: At all times, it has some Dem of a CPD which it uses for all inputs. Its "input" is a prob. dem. Its output will depend on its input dem (Natch, since its a CPD! @). ~~then~~ we will have to decide just what kinds of output we want from TM, for each input type: for starters!

124

1) Input is a set of finite objects: would like one or more of the following:   
 Modify of TMs CPD dem. so that, given a partial object "in <sup>the</sup> set", it can give a PD of ~~the~~ <sup>the</sup> possl. completions ~~of~~ (or partial completions - e.g. "what comes next in the partial object?").

133

2) Given (input of 1) <sup>A</sup> (129): Modify of its internal CPD dem to minimize the dem of the entire corpus.

3) Input is a partial object; Output is a PD on possl. continuations of <sup>that</sup> object.

Well, we could just have 2 modes: ① f. input is an augmentation of the corpus: we want TM to modify ~~f. CPD~~ CPD dem. to accommodate it. (Minz. augmented corpus dem).

② we put in partial augmentation of corpus (or of part of corpus): output is PD on possl. continuations.   
 This is a combn. of 129 & 133: we can create <sup>Create</sup> ~~them~~ <sup>them</sup> by adding D's/o's.

7698 TM expo:

Another way to look at it: We input same objects plus partial object want continuation of.  
We want modify of CPD down a PD on to continue after partial object.

If the partial object is null, would put a PD on the continued object in the sub corpus - but we can request not to do this last.

→ Note list to make an optimum "contin." p.d. one must precede the entire CPD down.

.06 If we give it an OZ problem (w. or w/o. CB constraint), its output is a p.d. over a set of optim. methods. (The nature of the problem suggests which methods are best... so we use

.08 P.d. as a basis of search for / OZ prob. soln.  
Criticism of .06-08: As one runs the problem on the optimum regions; we may get into test suggest modifying the probs. [This may be helped by my soln. to the "what to work on next" problem.]

Its "learning" of which optim. methods are best for what problems seems flawed. It will spend most time in a "favorite choice" - so that choice is, indeed, most likely to find the "peak" form: We will have to compensate for this bias by giving more resources (w.b.) for solns. by "non-favorite" optim. methods

Well, we can give TM optim. probs, & have many different optimization techniques work on the same probs: from the results, TM will get a kind of empirical idea of what optim. methods are best for what problems.

.19 This may work well enough: Ideally, TM should be able to analyze a problem, & by virtue of its

"structure", & the "structures" of various optim. techniques, decide which techniques were best for that problem. Hrr, I'll have to "look into" just how this could be done by TM.

Maybe by watching the "traces" of a person (or properly trained machine) solving problems of that

.23 sort.

I think what I want to do is list various kinds of problems: Then for each, describe what I'd like TM's response to be.

But the main idea (I think) is that we have a technique for doing getting a CPD for a variety of problems, & we have a way for it to learn to improve its performance.

All I really need is a very good darn of its "steady state" performance. 6.10 & .22 are some ideas.

I think the hardest part of TM's work is: Given a corpus w. a shorthand code for it: Given an augmentation of the corpus: To modify the code so it is "short" for the augmented corpus.


An imp. case is when the augmentation is small (e.g. one full object that is a poss. completion of a given partial object). We may want to derive special techniques for this case.

Another imp. kind of prob. is OZ probs & improving TM's skill in this area: .06-28 discusses that a bit: .19-.23 is perhaps imp.

Does paper cover  
Q: why stomach?  
Or certain  
offerings?

01. Dcm. of Pev's system for talk! (1) The CPD: Steady state. It has a short code for a corpus (plus far) corpus. We augment this corpus in some examples, ~~or~~ plus a "partial object": output TM then looks for ~~the~~ modifications of its code in view of the new augmentation of the corpus & possibly continues it.

Its output is a pd. ~~and~~ continues of the "partial object".

 for tech, maybe don't mention "partial object" at first: the CPD is for ~~the~~ finite objects in view of the "corpus up to now".

~~It uses this facility to improve its "skill" in solving tasks.~~

10. (2) At any particular time it has a ~~large~~ large set of option points in memory. It looks at ~~any~~ ~~or~~ ~~from~~ it is able to form a pd on these ~~opt~~ methods - which would be ~~the~~ basis for ~~the~~ ~~problem~~ that problem. A common <sup>kind of</sup> problem is ~~short~~ finding short codes for a corpus.

14. (3) It is able to work on the problem of devising better option points.

15. Pev is a special case of its normal skill at solving ~~or~~ problems.

This is very u to solving; but perhaps main difference is emphasis:

Now, I'm "mixing domains" a lot. The probs tend to be (at first) simple in direction problems based on a set of examples.

I may be less afraid to "solder in" various impt. ideas - (like <sup>(infectors)</sup> ~~opt~~ routines (Pev eventually, I want to factor Pev's ~~idea~~  $\epsilon$  routine for Pev)),

I may be less afraid to ~~run~~ (try ~~in~~ ~~the~~ ~~try~~ of TM) ~~short~~ problems. But I haven't myself found a soln for ("soln" int. sense but TM could find ~~the~~ soln in a acceptable (IS).)

Perhaps a major difference: I'll be ~~more~~ using TM almost exclusively ~~or~~ problems. It will not use Lurch mech, except in a trivial way to select ~~the~~ appropriate ~~or~~ methods.

So the emphasis may be on starting w. v.g. option techniques

is not until my later try to get TM to improve Pev <sup>9</sup> ~~to~~ <sup>Some default</sup> ~~to~~ <sup>10</sup>

A possibly Major Ditty in the present approach: That "finding short codes" is not an ordinary ~~or~~ problem: As the corpus expands, ~~the~~ ~~contains~~ new "domains", the kinds of techniques used to find short codes will change, so (if) (improving the option methods)

37. may be an essential part of the steady state system. That w. a fixed

Set of option methods, I will only be able to ~~do~~ do a limited



Several "domains" or "regions of Domains" — That TM will have an effective "epistemological

cut-off": for this reason, I may want to return to the old [Sol 86, 89 system].

Using Carefully designed Concept nets.

Could I reach a kind of Compromise ~~time~~ betw. the latest system &

Sol 86, 89? — Having some of the best features of both?

First look at how the T. disty of [66.37-67.02] is. It may be possible to

solve it, using a suitable Concepts Try-Seq. act. — Text / <sup>corpus</sup> / Regulatory forms

Concepts, is TM simply has to learn them in that form. Hm, even so, is TM is

"Sufficiently Advanced", I can design "TS Q's" intuitively (as for a young student)

Ab Initio

So: One Q is: Can I get TM to be "sufficiently advanced" using various

"option techniques" I know of? — Note 89.25 — It may well be best to find a Regym

& corpus & make a special option technique: That this "concept ~~net~~ net"

consideration of 108-11 is an accurate assessment of the problem.

The basic idea of 108-11 is that the regulatory observable in a corpus, are

based on concepts that have to be found as a concept net is found, & that the same

constraint on CJS applies & tells us which ones can follow which & how much

as we need to discover which new conc.

Looks like 1.13-22 may be an adequate understanding / reason.

of this most imp. problem.

So 66.01-67.22 pretty well covers the new system & its critical problem(s).

It ~~is~~ suggests a very Optimistic Evaln. of the system — i.e. suggests the likelihood to actually work!

An imp. point: system must do it to be able to look at a corpus & guess what kinds of

regys are likely to be in it. On a gross level, this is a "simple" IR

problem, being to do w. ways of inferring conc., what "features" suggest

"what types of regys"? This "preanalysis Evaln." can be of any depth & can increase

(do part of) & final search & short code. Ideally, these 2 aspects of coding should be integrated.

59.07-.40 is an outline of most of the talk 66.01-.15 is an outline of the end - it's

1. outline of 59.07-.40.

T. talk itself: View Graphs! ①

① General Linear Programs! 59.09

② Concept notes (pictures)

$$\frac{CC_j}{PC_j}$$

optimum ordering



$$\frac{CC_j}{PC_j}$$

not average:  $n$  times for soln:  $n$  upper Bnd.

~~upper bound at each time~~

if it is within factor of 2 for  $n$  way, it's ~~close to~~  $\frac{CC_j}{PC_j}$  for structure.

If  $PC_j$ 's are rounded, then it's too certain. Method L search from time  $\Rightarrow$  will be very close to  $\frac{CC_j}{PC_j}$

Search hours

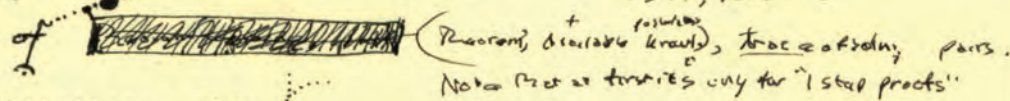
It was possible to express ~~search hours~~ as mod max of  $PC_j$ : So if you put all off hours your going to waste in  $PC_j$ , L search will do at least as well as heuristic search using  $PC_j$  hours.

between



Date of / latest idea on **C.P.D** machine. See 66.01 - 67.40 but mainly 66.01 - .15

Also draw picture of examples of  $Q, A, P$  parts



Then discuss differs in how they will be dealt with.

Perhaps discuss ~~idea~~ idea of building up more complex "simpler solns" from draw conc. nat. idea of from cons.

from idea of L search is  $\frac{CC_j}{PC_j}$ . (Maybe discuss Time Shaw).

Discuss: This is simple denn; In addition, systems able to handle such heuristics.

And  $\frac{CC_j}{PC_j}$  is  $n$  times for soln or upper Bnd.

is within constant factor of Best possible way.

Idea of "constant factor": Its error  $\sim 2$  or 4 depend on how one does L search.

Idea of heuristics in search: How they  $\rightarrow$  within of  $PC_j$ . If all hours are in  $PC_j$ . Then L search will

ever be least as good (or usually better) than use of those hours in "heuristic" search for soln.

Concept Nat:  
acyclic Directed Graph,  
Parallel ordering.



10:05

Sat of Bottleneck

SN: Uniform hand down  
Lo pass rule - if  
we want "expected values"  
of probab:

Here for Decision  
making: perhaps these

"Expected values"  
are just what we need!

critical objects are  
probab & Utility.  
(or may be expected  $PC_j$ )  
are not rate!

9.7.98 TM : VAPNIK, Vladimir (Volodya?)

Discn: He is not a Bayesian: Tho he does use a priori info, he uses it in a different way - to select a set of functions. Some (relatively recent) notes on U.C. dimension in S.V. machines: Related to discussion with Poggio (Q.V.)

4 imp. param. is his  $\epsilon$  (the zero penalty zero size): It is related to "Margin" - which, (in the case of hyper surfaces separating classes of points), is the distance of the h-surfs. from nearest pts in the classes being separated. (But what if the h-surfs cut thru both classes?)

So he has a 2 part GC. 1) to U.C. dim (or equiv) that tells you "complexity" of the hypothesis concept class. 2) the other is the observed error size.

He is impressed by Popper's falsifiability: His term "Shatters" is the same as "Falsification".

He says that if the ~~sample~~ concept class is all classes, then there is ~~no~~ <sup>no example that can</sup> falsify any concept class. So "degree of falsifiability" is closely related to U.C. dimension.

He says (in Poggio's) that ~~SVM~~ S.V.M. technique is much fast & better than Neural nets, & may be beginning to replace them.

He has gen. available on net. for his stuff. He does not need ~~any~~ <sup>any</sup> computers machinery to eat usable ~~results~~ <sup>results</sup>.

Luciano Tech  
Visit:  
Comm. 2 parters.

988 2.13A 3  
989B 2.30A 2.5  
5A 2.5  
730A 2.5

He says that the true prediction function is unknown. Set considered, then Bayes assigns zero to best hypothesis. I think this is irrelevant. We never hear any info to

~~know~~ <sup>know</sup> + correct ~~hypothesis~~ <sup>hypothesis</sup>. Instead, if  $\epsilon$  is small or  $C$  is small, we automatically get

large Granularity (i.e. Vapnik's  $\epsilon$ ).

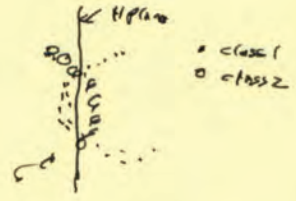
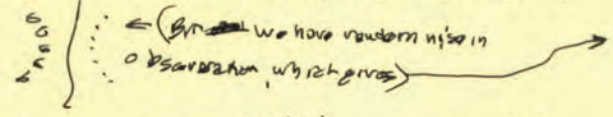
V. mentions Barron's paper on Error-Convergence. V. may not be aware of Sol 78. He felt that Barron's was somehow not Bayesian. (?) [It is poss. that Cover is not a Bayesian!]

He uses ~~the~~ Hyperplanes for separation of classes but the ~~rules~~ center n.l. limits.

(i.e.  $X_1 X_2$ ;  $X_1^2 X_2^3$ ;  $1$ ;  $X_1 X_2^7$ , etc.)

Why does he do when "Margin" is  $\leq 0$  i.e. classes interpenetrate?

T. "true" model!



No. of Dims = No. Data pts (i.e.  $S \geq Z$ ) needed for ~~less~~ <sup>of his</sup> separation of his surfaces ( $\rightarrow$  see 69.2.14)

Perhaps his A priori is always of the form: Certain functions are imposs. T. rest (which we know what they are) are of uniform = a pt. - so he has to define a density function anyway! So he writes as well do it over the entire ~~set~~ <sup>set</sup>.

V. says he gets better results than ANN is much fast: well, ANN is normally slow anyway. Anyway Max & Min (SVM) should improve it. But V. says he usually needs more support vectors than are usually used in ANN! V. says that attempts to relate his stuff to "k-NN" have been kind! Some parts correspond in simple way. Others do not correspond at all and could be made to correspond - they just essentially don't work.

Gzmm want

my c.v.:

2 sources:

- 1) That proposal I wrote that became Sol 89
- 2) Start for Oxford Brochure.
- 3) More recent up date of Bibliography

G.V.:

Maybe recent Lectures.

1909

$$\int \int K(u, v) \ell(u) \ell(v) du dv \geq 0$$

$K(u, v)$  is isotonic Pos. Def.

$$K(x \cdot y) = (x \cdot y)^2$$

Distances "barrows" way "is"

$$\int \int (x \cdot y)^2 \ell(x) \ell(y) dx dy \geq 0$$

$$-(x-y)^2$$

$$x^T A x > 0$$

$X = x_1, x_2$   
 $Y = y_1, y_2$   
 $z_1 \quad z_2 \quad z_3 \quad z_4 \quad z_5$   
 $x_1 \quad x_2 \quad \sqrt{2} x_1 x_2 \quad x_1^2 \quad x_2^2$   
 $y_1 \quad y_2 \quad \sqrt{2} y_1 y_2 \quad y_1^2 \quad y_2^2$   
 $(x \cdot y)$

$$x \quad K(x, y)$$


---


$$\sum \alpha_i (x_i \cdot z_i)$$

$$X = (x_1, x_2) \rightarrow x_1, x_2, \sqrt{2} x_1 x_2, x_1^2, x_2^2$$

$$Y = (y_1, y_2) \rightarrow y_1, y_2, \sqrt{2} y_1 y_2, y_1^2, y_2^2$$

$$(1 + (x \cdot y))^2 = x_1 y_1 + y_2 y_2 + 2 x_1 y_1 x_2 y_2 + x_1^2 y_1^2 + x_2^2 y_2^2$$

9.9.98 Vapnik - see 73.07 for More on Vapnik

Vapnik suggested ~~the SVM~~ that Fogio might be able to get me some grad students at MIT.

in His SVM alg, one is able to determine a way, f. separation ~~with~~ h surface ~~or~~ or the regressn function: (along w. an "E" value) & it will, relatively quickly, find an optimum fit. — Among f. functions it can use for regressn: linear (perhaps ar by degree); & f. of various types.

His SVM alg would run at c. h. speed on a 386. How slowdown was due to output written may be in Java (an abstractive lang.).

He says that in ordinary ANN algms, the no. of nodes they use is much less than the no. of "support vectors" he uses.

He says that "support vectors" are always a subset of the examples.

I had been mainly thinking of use of curves, splines... to separate classes of pts (classifi. problem): But he also uses it for Regressn (curve fitting);

1/d So given a bunch of data pts: fit a curve using some error function w.  $\sum \epsilon_i^2$ . Error  $\ll$  cost is zero. We get a curve, then we find that only a certain subset of the data are "support vectors" & that ~~the~~ fitting alg alone, gives same curve as entire data set.

The no. of support vectors gives something like the complexity of f. by  $\rho_{\text{eff}}(\epsilon)$ . (See 69.30)  
The  $\epsilon$  used corresponds to the "margin" in 69.28

Re: Support vectors for curve fitting. There are various ~~parameters~~ parameters & ways of measuring goodness of fit ( $\equiv$  GOF): ① One is pure RMS or some other norm;

② is one w. V's "E" of zero cost ③ no. of coeffs or other measure of info complexity ( $\equiv$  cost) of f. model (maybe v.c. dimension).

④ My w cutoff of  $\epsilon$  xtens used to model f. function, or idea of A. Barron on fitting using ~~the~~  $\sum \epsilon_i \cdot f(\frac{1}{2}x)$  for model or  $\sum \epsilon_i \cdot f(\frac{1}{2}x)$ , if  $\vec{x}$  is a vector.

These various params give different properties of the models used for approxn. They can tell something, perhaps, about the "True" model of the data.

01: 68:40 View Graphs:

① GENERAL LEARNING THEORIES. [How to use data from past for production]

Algorithmic Probability: Solovay ~~Wallace~~ - ALP <sup>never used for Prodn.</sup>  
~~Wallace~~ MML Wallace  
Vapnik. MDL Rissanen.

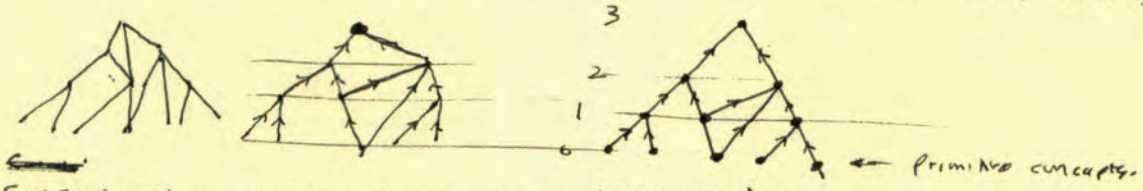
~~V.C.~~ V.C. Dimension, ~~S.V.M.~~ S.V.M. Vapnik, Chazanovitch.

Reproduction  
Genetic Algorithms - Holland.  
Recurrent Neural Nets .....

SOAR ~~Neural~~ Neural, Residual Network

... ASM, E. nico, Cye : Lanot.

② TRAINING SEQUENCES and SEARCH (Levin's Search Procedure).



Concept nets: a cyclic directed graph; partial ordering -

Concepts: Set of problem concepts

Macros that can be combined to give solutions or other macros on

- 1. Solutions to problems
- or 2. other macros

CC:  $\frac{CC_1}{P_1}$   
yield  
yette

How System solves problem problems by search.

- " " assigns probs to concepts.
- " " uses probs for search:

Cost =  $\frac{CC_1}{P_1}$  : Least <sup>of trials</sup> used to order trials

Cost of trial is good estimate to total search time needed to find that sol.

This search quality measure is much easier to design than SPUS. - we can always estimate how long it takes TM to solve a problem.

Learning Algebraic Notation  $3+7=10, \dots, ((n+1)/2)+1 = 2 \dots$

Learning to solve linear eqns and learning "Laws of Algebra" (Associativity, commutativity, distributivity laws)

Learning to solve concepts: Run Quadratic, Run Cubic.

Problems: ① TRAINING Sequences



~~The~~ method of solution using RLP:

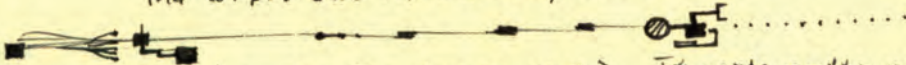
We have an algorithm that operates in the following way:

We insert a corpus of data: ~~sequences~~ <sup>sequenced</sup> sequences of numbers or vectors,

or ~~alternatively~~ <sup>alternatively</sup> a set of ordered strings - each that represents an object.

The output of the algorithm is a <sup>sequenced</sup> start codes for the corpus. If we give the machine more time, it ~~may~~ <sup>will</sup> find shorter codes.

These short codes can then be randomly combined to give ~~new~~ <sup>new</sup> codes for the corpus and its ~~most~~ <sup>most</sup> likely continuations.



How does this device work? It works pretty much like ~~the~~ <sup>the</sup> training sequences I discussed before. We start with very simple in ~~the~~ <sup>the</sup> problem. The machine uses

an L search over a set of ~~possible~~ <sup>possible</sup> codes. <sup>When asked for output,</sup> It presents the shortest codes it found thus far. The longer we wait to ask for output, the ~~shorter~~ <sup>shorter</sup>

the codes can be.

After it solves this problem, we ~~can~~ <sup>can</sup> extend the length of the corpus and ask for a short code. <sup>for it</sup> The machine starts with codes for ~~part of~~ <sup>part of</sup> the early part

of the corpus and tries to modify it, so code generators codes for ~~the~~ <sup>the</sup> augmented corpus. We can continue this learning by giving

the machine an ~~augmentation~~ <sup>augmentation</sup> to its corpus to code, or an arbitrary new corpus.

The ~~previous~~ <sup>previous</sup> coding problems I've described will only be solvable in real world time, if ~~the~~ <sup>the</sup> regularities in the sequence of problems form an acceptable training sequence - involving concepts that build on previously ~~found~~ <sup>found</sup> concepts. In this sense, this present problem has all of the difficulties of ~~the~~ <sup>the</sup> continuity ~~of~~ <sup>of</sup> fixing TQ's in our previously described project.

There are a few impl. differences.

1) Training examples can come from a variety of sources, to help learn various kinds of regularities in data. The system affords a kind of ~~distributed~~ <sup>distributed</sup> learning device that is able to integrate learning from various domains.

2) An ~~important~~ <sup>important</sup> problem <sup>in</sup> learning systems is ~~how~~ <sup>how</sup> to learn when necessary to recognize when a situation is similar to ~~one~~ <sup>one</sup> that has occurred in the past. We can use the induction machine to recognize to help solve this problem.

[perhaps Mention CBR and EBL] <sup>Helps</sup> solve the ~~main~~ <sup>main</sup> problem of using Memory of Cons. Efficiently.

perhaps Do discuss impl. of memory (to try)



Augmentation of talk:


1) Go into differences Betw. ALP, RLP, MML, MDL.

Who knows our context, also trying to help.

Idea of Incompleteness; Unknown Error Size; "killing the messenger" harbinger of bad news.

2) Discuss ~~Machine Learning~~ "Reading and Understanding": Re "Textworld" project.

07: ~~69.1~~ - 69.2.90

3) Re: Vapnik:  is simply a class of loss functions. If its a good class, it should be useful for compact coding.

4) The "support vectors" are always on the data set  $\hat{S}$  <sup>functioning on</sup>. loss function for that subset will be identical to that for  $S$  ~~on~~ entire set. Usually the "support vectors" will be at  $P$  or a  $d_j$  of the classes to be separated.

I mentioned case where classes overlapped unassisted way (that is "two classes")

Was separable, but ~~is~~ true classes ~~overlap~~ overlapped: He said that in such cases

"best" we can do a ~~best~~ fit is we have a certain no. of "error pts" & we are penalized for them.

5) <sup>Vapnik</sup> For Separating 2 sets of 2 dim. pts  $\{x_i, y_i\}$  I will use  $f(x_i, y_i) > 1$  &  $f(x_i, y_i) < -1$   
 $f$  will be, say a 5<sup>th</sup> order poly. He may find some pts that can't be categorized so:

So he has a no. of "error pts".

6) Perhaps just go thru the talk once, & if there is time, go back to discuss certain pts. in more detail.



So a major problem is how to devise tsq's for learning how to discover short codes (= "regularity")!  
 On 66.32 - 67.40 is an earlier (optimized) model of the system! Finding "short codes" is regarded as a kind of oz problem & we have a great variety of ~~oz~~ <sup>optz</sup> perms "on top" in many. Also, we have "factor P's" set of optz perms into a language in suitable (non-terminals), I got the idea that finding short codes was likely to be a very special kind of oz problem

perhaps needing a very special optz ~~tech~~ technique.

The approach: Watch me do very data-driven! Generous Methods  
 to factor them into a language (grammar)

Glickstein  
Stina

Etc  
Brain  
Mind  
a Neurophysiological  
Approach.

My personal approach to short codes: first: look at the problem! This will suggest which

types of regz to look for.

→ Martin Lot wrote about "Randomness tests". Perhaps one could "Hill Climb" on results of Randomness tests.....

University College  
London.

Some "Randomness tests": for a sequence of texts! Try prediction as a function of previous outputs. Do various SM predn schemes. Try "dimensionality" approaches in Chaos theory.

The BSD (or whatever) test for chaos is also a test for regz of any kind.

Comparing putting neural nets to date - detects a great variety of patterns. Doing it right could be much better.

V.C. Dimension is Support Vector approach can detect regz. (see 69.01 off on Vapnik etc.)

SN "Partially Computable" functions do b exactly what we do. Its much better than "Partially Computable" or Approx P.R. or Practically P.R.: It "implies" or "allows" c. idea that one <sup>can</sup> explicitly, optionally limit part of the compn. with say, "CS" in a most general sense.

TMA

SN2 I imp! T. recent idea of TSO's for ~~the~~ <sup>only</sup> learning to do ~~short~~ <sup>short</sup> codes. ~~seems~~ <sup>is</sup> an "adequate" TM, but ~~seems~~ <sup>is</sup> very hard to do. ~~By~~ <sup>By</sup> an earlier

Derbim availability.  
 Was this 72.02-70? Very likely  
 Also 66.01-67.70  
 66.01-65.40

TMB

idea leading to that one, was the idea of a Cond prob learning machine but it learned things of any sort at position in its CPD. (not just ideas on how to find short codes)  
 Now, it <sup>probably</sup> ~~is~~ <sup>is</sup> that these 2 ideas are opnly, that any kind of lrng. involves finding a short code. If so, then the essence of TMB ~~is~~ <sup>is</sup> could be a good ~~idea~~ <sup>idea</sup> (more easily designed) source of TSO's than doing TMA directly.  
 we could simply do TMB & have TMA be one of the "conditions" for TMB's CPD.

This TMB idea does seem nice, & easier, more <sup>fun</sup> ~~fun~~, to write! More immediate joyous feedback than TMA. ~~is~~ <sup>is</sup> But TMB is ~~is~~ <sup>is</sup> is ~~is~~ <sup>is</sup> TMA: Perhaps I somehow didn't realize that any induction problem could be made a part of the TSO of TMA.

In fact, TMB, (maybe?) must first decide what kind of (induction) problem it is being asked to solve. At first, we will include, with each problem, a label, telling what kind of problem it is. TM will learn to pool various problem types & separate them into sub-classes.

So later, the labels will not be critical, but they can furnish TM with ~~the~~ <sup>the</sup> initial info!

1.01 ~~There~~ seems to be a serious difference between  $TM_A$  &  $TM_B$ : Fig.  $TM_B$  can be given a set  $(I_j, O_j)$   
 1.02 of  $I, O$  pairs: then given a new  $I_2$ , it gives a p.d. over  $O_2$ . It does this <sup>(usually)</sup> by constructing  
 a ~~(set of)~~ short codes for  $[I_j, O_j]$  — or, more exactly, a short operator for  $(A, B \& Q)$  here it just shows its shortest OSL (Complete Lrn) problem.  
 1.04  $I_j \rightarrow O_j$ . As such, it is built up from a TSO of less complex operators,

that enable it to lrn imp. "concepts".

A "cold prob. machine" is a stack operator, (its fr. soln. of .02 - .04).

A simple example of .02 - .04 is the "Alg. notation Learning" TSO. — is its soln.

I should write up in much more detail, just what  $TM_B$  does: examples of its operation & how it is able to implement  $TM_A$ . [Similarly, I may want to detail out, how ~~TM~~  $TM_A$  implements  $TM_B$ .]

1.14  $TM_B$  is easy to solve  $TM_A$ 's problems: It amounts to a "change of representation" or "reformulation" of the problem. How did I discover this? — In general, keep track of "reformulations" that I do, in attempts to see how I do them. A major subgoal for TM would be its ability to do such "reformulations", → NN1,10 (#12)

■ The main (perhaps only) way that I had  $TM_B$  solve problems, was by  
 → Lrn using a suitable (conditional) P.D. Does this amount to simply trying to  
 → get short codes by Lrn? (using suitable CPD).

91498

In line with this ~~last~~, searching for short codes for a corpus: first one "identifies" where the corpus came from so one can get Conditional Probys. Then one does a Lrn on suitable "primitives" using the cond. probys assoc. w. that kind of corpus. The "primitives" used, are (conditional) dependent on the "nature" of the corpus. So one has this big set of poss. concs. The <sup>(name, nature...)</sup> "nature" of the corpus, then induces a p.d. on those concs. w.r.t. to problem of finding a short code for that corpus.

Faint, illegible text at the top of the page, possibly bleed-through from the reverse side.

**BATON'S**  
**CORRASABLE**  
**BOND**  
**MARKET**

**COTTON FIBER CONTENT**

- 1. ...
- 2. ...
- 3. ...
- 4. ...
- 5. ...
- 6. ...
- 7. ...
- 8. ...
- 9. ...
- 10. ...

Main body of faint, illegible text, likely bleed-through from the reverse side of the document.

TM 174.49 M. Glickstein Cambridge, Mass

"Change of representation"

Change of representation " that ~~the~~ America's changed Method of expression not person

eye was a <sup>uniquely</sup> "re-formulation" "looking at past in 'new way": A new method of Coding, Nobel.

(11) Memorex of ~~the~~ prize winners!

2) format: Moved toward door stand of lecture

22 " : Format v.s. You as lecturer.

23 format's reply to Eppend should it w. temperature J. Nash's rebuttal.

91298 5:15 A (cont)

8:45 A 3/2

paradigm change SPM!

130-230PM Dizzy!

3PM 600-4MG

b1 ~~the~~ Newhall (Simon) : How conv. with Atwell : Not finding out better ~~more~~ at hook into machine! [ Not Simon must be insane! He know about History, Philos. of Sci. ]

91398 12:25 A

6:15 5 1/2

10:15 4

(10) (12)

TM 75.25 TM 76.14

BALD-FACED

TM B is a way to solve bald-faced problem of TMA. - So it's a

91498

12:01 A ~~2:00~~ as such. Just what happened? How ~~did~~ <sup>what</sup> ~~the~~ "re-formulation" work?

12:15 A ~~2:00~~ How ~~did~~ <sup>did</sup> I print of it?

12:30 A ~~2:00~~ (13) Try to keep record of "inductive leaps" <sup>make</sup> ~~is~~ <sup>is</sup> see if it could be used to make

12:47 P ~~4:00~~ <sup>discuss</sup> w. acceptable CJS.

5P (14) SMFI: ~~TRY~~ TRY present p.m. w. updating, but change time scale from 1 day to

22 minutes (1 month). The idea is that we don't have to work so much every day.

Even if 1 yr old is much less, it may be worth it. - Say when we are on vacation.

Also, we can have a much larger set of poss. driven (use all kid's), w.o. much trouble.

(15) Sol. of "How Buildup Blocks Hypothesis" works! If A & B are useful ones.

Then they have short codes: "Hinting for Uses of A.B" would be finding a justification

of the new beta,  $A \subseteq A.B$  - which would give a shortening of code.

This is a justification of ZI & I idea. (which was never clear to me)

So looking for ~~many~~ <sup>unusually</sup> way occurrences of AB is somewhat

easy! There is a  $\subseteq$  of it! If  $F_{AB} \ll F_A \cdot F_B$  then clearly a very

has been found, but how best to use it for induction? <sup>see 2.20-2.40 for</sup> <sup>Av. of 2.141</sup> <sup>3.04 (#22)</sup> <sup>see 2.34-4.0 for</sup>  <sup>$F_{AB} \ll F_A \cdot F_B$</sup>

(16) Simple proof of first Gomb House Prime. If trials for A & B are not in

$\frac{F_A}{r_{A1}}$ ,  $\frac{F_B}{r_{B1}}$  order, then ~~changing~~ their order will ~~be~~ <sup>be</sup> ~~corrected~~ <sup>corrected</sup> ~~ce~~

of such. We need only prove this for A & B being adjacent. 14 orders is 14 orders.

(17) A number of "fuzzy patches" (Nectanomas)

(18) ~~Equ~~  $R_{\alpha} : .27-.28 \approx 2.14$ . If  $F_{AB} \ll F_A \cdot F_B$  <sup>because it for whole set is properly</sup> <sup>orders,  $L_{A1} < L_{B1}$  for all</sup> <sup>disjoint A, B.</sup>  $\rightarrow$  (6.0)

think  $\alpha \approx AB$  should be defined, & the order of corpus will ~~be~~ <sup>be</sup> ~~correct~~ <sup>correct</sup> (I think): Note that ordering of definitions must be correct. Also ~~the~~ parsing of corpus must be correct - It may be that we can only get the  $\alpha$  in the  $r_{c}$  by a certain parsing, & that we ~~only~~ <sup>only</sup> have one short code. But go into exact details of pc computation - Use arrays plus, maybe use computer to check ~~are~~ <sup>are</sup> makes a ~~thorough~~ <sup>thorough</sup> check

accuracy of symbolic calc. Also Note! I. ordering of doing is imp. Greedy is not ~~always~~ <sup>always</sup> ~~the~~ <sup>the</sup> ~~best~~ <sup>best</sup>. Try various

9.15.98

WN

Z141. (20) .18

7M 79

2

(18) (cont) For large  $SSZ \rightarrow$  steady state,  $f_c$  is  $p_A^{p_A} \cdot p_B^{p_B}$   
 if  $p_C \neq p_A \cdot p_B$ ? we should have equality here  $\rightarrow (p_A - p_C)^{p_A - p_C} \cdot (p_B - p_C)^{p_B - p_C} \cdot p_C^{p_C}$   
 if  $p_C = p_A \cdot p_B$ .

$p_A \rightarrow A, p_B \rightarrow B, p_C \rightarrow C$ !

$$c. (A-C)^{p_A} = (A-C)^A \cdot (A-C)^{-C} + (B-C)^B \cdot (B-C)^{-C} \cdot C^C = (A-C)^A (B-C)^B \cdot \left(\frac{C}{(A-C)(B-C)}\right)^C$$

if  $C=AB$  then  $= (A-AB)^A \cdot (B-AB)^B$   
 $A^A \cdot (1-B)^A \cdot B^B \cdot (1-A)^B \cdot \left(\frac{AB}{AB(1-A)(1-B)}\right)^{AB}$

$\frac{1}{2} \cdot .5^{.5} \cdot .5^{.5} = .5$   
 i.e.  $.25 \cdot .25 \cdot .25 \cdot .25 = .25^{.75} = .3535534$

$$(1-B)^A \cdot (1-A)^B \cdot (1-A)^{-AB} \cdot (1-B)^{-AB}$$

$$= (1-B)^{A+AB} \cdot (1-A)^{B-AB}$$

$$= ((1-B)^{(1-B)})^A \cdot ((1-A)^{(1-A)})^B$$

if  $1-B=A$  then this is

$$A^{A^2} \cdot B^{B^2} \neq 1$$

$$.5625 = .25^2$$

$$.0625 = .25^2$$

$$= .779975$$

Hvr: If a vsgy exists, it must be poss. to recode + corpus with other vsgy, so that the resultant code is more "random"

T. Morel seems to be stuck simply defining  $C \equiv A \cdot B$  does not capture all of the dependencies implied by  $f_A \cdot f_B \neq f_{AB}$ . The finite state machine with  $A$  context can't be followed by  $B$  does capture this! So that in a resultant code we always have conditional probys. to 20

(19) The OXB & s. Prob. Problem: w. random data in no PAB, yield will always be  $\phi$ , but Venc. can be  $\gg$  Nat of data itself.

16 (20) from (18) we want to ~~code~~ code it into a Barré seq, so we can ~~code~~ look for vsgys in "Re code" it.  $f_A + f_B + f_C$

20 But: say we had 3 symbols w. indep. probys: 0, 1, 01. How would we do that? - what vsgys in corpus? null, & freq vsgy of 01 would be  $f_0 \cdot f_1$ . So perhaps

if  $f_{AB} \neq f_A \cdot f_B$  define  $C = AB$ , but knows frequency  $\neq$  to freq. of occurrence  $AB$ !

This coding method has very many ways of "parsing" the seq. - i.e. one can decide on any of the AB's vsgys C, in a result. AB's bump "coincidence"! So it's very to code a corpus in many (1 way to obtain

24 all of the  $p_C$ 's.

Out. other hand, the Z141 method gave only one parsing, but it was harder to detect vsgys in the resultant sequence (i.e. "recoding")

(20 - 26) Seems closer to a very understanding of the ("Building Blocks" concept)

Att So! : If we had  $f_A \cdot f_B \neq f_{AB}$  by and  $C \equiv AB$  is a useful defn, but it doesn't really tell when to use it! It just gives an FM PC of C, for ~~coding~~ coding of corpus. So its of "strong Heuristic value" (1.20) (#6)

36 Woe! How do we get  $f_A \cdot f_B \rightarrow f_{AB}$  by defining  $C \equiv AB$ ?  $f_C$  would have to be  $\phi$  (!) Well. If  $f_{AB} < f_A \cdot f_B$  then  $f_{AA} > f_A \cdot f_A$  &  $f_{BB} > f_B \cdot f_B$ . So we define

$C = AB$  if  $D = BB$ . If there is a big redix, this may not be so good (??)

(under hvr) Probably best try to do: Define  $C = AB$ , say, then look for trigrams of abnormal/unusual frequency if any are found, they may occur w. relatively low freq. so one can try parsing the corpus all poss. ways w. then C there could be only pos. parsing!  $\rightarrow 90.01$

21) I Post Post GA was Universal! Post eventually it would solve any problem —

False! It suffers from "Local Maxima" ("Premature Convergence"). How Do GA practitioners Deal w. local Max's? (e.g. Koza) (Look at my recent Book on GA). (Also Note G.36 criticism of GA)

22) from (1.28) "Building Blocks" (cont.). In GA, the meaning of  $F_A, F_B, F_{AB}$  is unclear! Maybe we mean frequency int. by a part of the ensemble of phenotypes (candidates = "cands").

Remember, it's the finite set of objects" corpus. (e.g. even  $\{Ob_{i1}, G_i\}$  corpus. This is + More exact way to do it!)

OK: Lets just consider  $\{Ob_i\}$  corpus: We have a small Grammar first, using  $F_A$  &  $F_B$ . We note  $F_{AB} > F_A \cdot F_B$ , so we find the new Grammar component  $C \in A \cdot B$ , which compresses to corpus

further. (Noting  $F_A \neq F_B \neq$  conc. of all symbols, is to first (Bernoulli) compress.)

.04-.09 is fine! Now, try to introduce to more exact corpus,  $\{Ob_i, G_i\}$ .

By assigning more prob to cands. Post occurrence occur in (trials) of  $h$  or  $a$ , we may be approaching  $S \in A!$  Note that in SGA, we tried to devise a grammar in which the pc of a  $\{Ob_i\}$  was  $\propto G_i$ : So  $G_i$  was like weight of  $Ob_i$ .

Is it hard to "linearize" G to the very first in a SGA? Also, it was made  $w_{G_i} \propto G_i$ ,

What kind of proportionality should we use? It would seem to be  $\log$  of distance. I did write a (short) paper. I think related to another difficult Q in SGA

But seems at first, quite independent.  $\rightarrow$  7.28

There was to General in GA of fitness [Any function of G has same max as G does: So how can Conventional GA decide what reproduction rate to use?] They usually make repro. rate  $\propto G$  (or maybe  $G \cdot G_0$ )

They may choose to const. of prop. w. G empirically: i.e. units with, see how it doing, then  $\uparrow$  or  $\downarrow$  const. dynamically, slowly. They may not worry about (.20-.21) at all!

Using classic SGA on .04-.09 may be very. Perhaps try it with  $w_{G_i} \propto G_i$  or  $(w_{G_i} = k \cdot G_i)$ .

Suppose that the final grammar, is expressible as a "Born set" or a small set of all usups. What is its ultimate capability of modeling various p.d.'s on the set of finite strings? We are mainly interested in good

modeling in the  $h$  region, since this is where Most of the wt. is.

9/17/98 Re: Vapnik: There was this term that  $\frac{1}{1+\epsilon} = \vec{X} \vec{M} \vec{X}$  was used. But there was some kind of system in which  $\vec{L}(\vec{E})$  was expressible as an inner product?

Maybe what was meant was: inner product  $\sum x_i^2 \epsilon_i^2$ : where  $\epsilon_i$  were given values.

So Re: of term I used in 1962 to prove that Merz was correct. Maybe Vapnik meant this in some more general sense — for a more general Matrix.

5N Maybe ask Vapnik about localizations of sites of ions in rabbit cave-batillum.

23) G.A. on  $\vec{G}$  (approx G-funct for SGA)! Am's error or other simple metric

betw.  $G \in \vec{G}$  is not what we want: Some better ideas: @ Considering all  $\frac{N^2}{2}$  pairs

of cands; a min no. of mis orderings. (b) Sigwa pick top 10 cands w.  $\vec{G}$ :

W's like to be contain top G.) If G is noisy (i.e. noise =  $\sigma$  in G), we want to pick  $\frac{1}{2} \sqrt{N}$

"top of G" is less clearly defined — says to. Set of cands within  $\sigma$  of  $\vec{G}$  arrange for.

$\rightarrow$  we can't know top G, or identity.

5N It is poss. that Cover is NOT a Bayesian! He was very happy w. models that assume no a priori info.  
Had my Lunch 2/1-4/1  
Sup. ~ 8 to 12  
9/10/98



91798 (W)

m. glickstem@ucl.ac.uk

T4.80

0171 504 2888 ← sh. no. at dept Anatomy & Devel. Bio  
263 7166 - home phone

7

EARLY:  
(SN) On development of ETHICS in Children:

This has mainly to do w. development of "ego" — i.e. what parts of the world that  
f-child considers as imp't for its welfare. At first, the parent is seen regarded as part  
of the env. but takes care of f-child. ~~Parents~~ These utilities are under threat (to  
theory of primary of f-child) so the child doesn't consider them as part of its ego.  
Later, as the child begins to understand how parents work, the child starts to understand  
if it perceives its parents were threatened. As the child grows older, other people  
may begin to develop parent-like properties, i.e. "self" (i.e. part of the ~~parent~~  
world that f-child feels, i.e. imp't in its well-being) is expanded: siblings, family,  
extended family, teachers — may be included.

A moral problem with it: "shall I hurt one part of the family (self) to help another?"

It is similar to "shall I harm one part of my body, to help another?"

It's a very "selfish" view of the world, but it looks like concern for others.

Later in life, w. better understanding of the world, the "others" can be more & more

complex. "Ego" can include a nation or the whole world — and other forms of

91888

ethical evolv. may evolve, but don't look as "selfish". A person may have multiple  
(parasitic) goal of concern w. welfare of the target "extended family" (= extended ego).

- 1201 P 2.45
- 1201 P: 2.5
- 230P 2.5
- 5A 2.5
- T/ISA 2.5
- 930 A 2.45, cord
- 930 A.

91998

- 12:30 A 2.45, cord
- 12:30 A.
- 2:15 A 2.45
- 6:15 A 4
- 10:45 A 4

3 1/2 hrs  
3:15  
26 0 x 3  
45  
32  
79  
78.8  
18 → 69  
36  
30  
645  
780 F  
100  
26  
78.8  
52  
30  
82  
2.5  
2  
10  
30  
30  
50  
10.8

(continued)

3.4 (23) In getting  $\tilde{G}$  one of the main goals of SGA was to use info in: low & part of the sample to get a good curve to fit for  $\tilde{G}$  m.r. by a part of the sample [Efficient use of sample]

Hum. t. main goal seems to be that top  $\tilde{G}$  cond. is most easily found! One criterion: that top  $\tilde{G}$  cond. is  $k \in \mathbb{R}$  from Pt. top wrt.  $\tilde{G}$ . we want Min  $k$ . As  $SS \geq 1$ ,  $k$  will normally  $\uparrow$ .

This sounds not so good! If we have in parans  $\tilde{G}$  to adjust in correcting  $\tilde{G}$ , It may be pathologically easy to find a set of parans that put top  $\tilde{G}$  close to top  $\tilde{G}$ .

Another / Goal for  $\tilde{G}$ : that max  $\tilde{G}$  should be achieved by  $\tilde{G}$  as possl.

So: Cases for  $\tilde{G}$  resistor: 3.33, 3.4, 3.6; 5.03, 5.08. || 5.03 may be best for  $\tilde{G}$  (max be 5.08!)

Re: (08) Maybe not concentrate on getting top  $\tilde{G}$ : Getting close to top  $\tilde{G}$  may not be bad!

Also note, we will usually never know top  $\tilde{G}$  anyway! So (08) is more reasonable.

11 What we usually want (i.e. in most phases of the calcn. of Max) is that  $\tilde{G}$  be able to

12 suggest new data pts, so that fitting a  $\tilde{G}$  curve thru the augmented corpus either

13 does 11-12 a/o has  $\tilde{G}$  of very by  $\tilde{G}$ .

Hum! t. funny Gores seem to miss f. point: I have  $\{O_j, G_j\}$  data set.

I want  $\tilde{G}$  to get into/ordinary  $\{O_j\}$  — particularly in the  $\tilde{G}$  region: Most of the previous Gores involve new trials.  $\tilde{G}(O_j)$  is available for any  $O_j$ , but  $G(O_j)$  is available only for old trials — or at high cost if  $O_j$  is a new trial.

Hum, t. funny Gores can be made "operational" by limiting the data involved to the  $\{O_j, G_j, \tilde{G}(O_j)\}$ . So we can ask that  $\tilde{G}(O_j)$  may error as by  $G_j$  as possl.

Actually, I want  $\tilde{G}$  of the best 10 (say) pts of  $\tilde{G}(O_j)$  over the data set.

t.  $\tilde{G}$  ordering over Pt. by  $\tilde{G}$  pts of  $\{O_j\}$  to be  $\tilde{G}$  ordering over those same  $\{O_j\}$  pts. — but really the  $\{O_j, G_j(O_j), \tilde{G}_j(O_j)\}$  pts.

(even those of low  $G(O_j)$ ) should ~~be used~~ be used: (Efficient use of the corpus).

19 E.g. if  $G$  is linear in  $x$  ( $x$  is a scalar) so  $G = a + bx$ , (then knowing  $G(x)$  for any values of  $x$ , is useful. or say  $G = a + bx + cx^2$ , values of  $G$  for any  $x$  are useful! But some are more useful than others.

34 Essentially what we want is that  $\tilde{G}$  ordering be close to  $G$  ordering over the known data set — so it is more likely that  $\tilde{G}$  ordering over  $O_j$  regions of very by  $\tilde{G}$ , will also have very by  $G$ . we will perhaps bias w.t. of ordinary  $\tilde{G}$  in by  $G, \tilde{G}$  regions.  $\rightarrow$  9.01

I think there may be a fair amount of literature on "Ordering" in statistics.

How to find it? Who to Ask? There may be inside my book on "Social Statistics".

19 IMP. The "number of misorderings" on some  $N$  Gores can be used to correct  $\tilde{G}$  so it is a "perfect code" of  $G$  (i.e. has same ordering). (See 9.01 for a generalization of this idea... but no solu.)

#25  $\rightarrow$  6.32

oi: 1.30: (24) from (16): say we order them so low  $\frac{cc_i}{pc_i}$  is first

say  $[cc_i, pc_i]$  is a known set of CC, PC values. Let  $z'(n)$  be a particular ordering of  $t$  trials:  
So  $z'(1)$  is first,  $z'(2)$  second, etc.

Then  $\sum_{n=1}^m cc_{z'(n)}$  is total CC of session, if  $z'(m)$  is the correct trial.

So  $c_0 = \sum_{n=1}^m P(z'(n)) \sum_{n=1}^m cc_{z'(n)}$  is expected time CC for the  $z'(n)$  ordering of trials

Now say we change to a ~~new~~ ordering by exchanging adjacent trials at the  $k$  &  $k+1$  position. So the ordering is ~~the same as~~  $z'(n)$  for  $n < k$  &  $n > k+1$ .

for  $n = k$ , we have  $z' = z'(k+1)$  for  $n = k+1$  we have  $z' = z'(k)$ .

In any, the sums w.r. to  $cc_{z'(n)}$  are the same except for the 2 terms when  $n=k$  &  $n=k+1$ .

In the usual universal  $z'(n)$  case, the 2 terms are:

$$c_0 = P(z'(k)) \sum_{n=1}^{k-1} cc_{z'(n)} + P(z'(k)) cc_{z'(k)} + P(z'(k+1)) \sum_{n=1}^{k-1} cc_{z'(n)} + P(z'(k+1)) cc_{z'(k+1)} + P(z'(k+1)) cc_{z'(k)}$$

$$c_1 = P(z'(k+1)) \sum_{n=1}^{k-1} cc_{z'(n)} + P(z'(k+1)) cc_{z'(k+1)} + P(z'(k)) \sum_{n=1}^{k-1} cc_{z'(n)} + P(z'(k)) cc_{z'(k+1)} + P(z'(k)) cc_{z'(k)}$$

The red boxed terms  $c_0, c_1$  cancel,

The terms that don't cancel are:

$$P(z'(k+1)) \cdot cc_{z'(k)} - P(z'(k)) \cdot cc_{z'(k+1)}$$

$$c_0 - c_1 = P(z'(k+1)) \cdot cc_{z'(k)} - P(z'(k)) \cdot cc_{z'(k+1)}$$

$$\text{so } c_0 - c_1 > 0 \text{ iff } P(z'(k+1)) \cdot cc_{z'(k)} - P(z'(k)) \cdot cc_{z'(k+1)} > 0$$

$$\text{so } c_0 > c_1 \text{ iff } P(z'(k+1)) \cdot cc_{z'(k)} > P(z'(k)) \cdot cc_{z'(k+1)}$$

$$\text{iff } \frac{cc_{z'(k)}}{P(z'(k))} > \frac{cc_{z'(k+1)}}{P(z'(k+1))} \text{ } \left. \vphantom{\frac{cc_{z'(k)}}{P(z'(k))}} \right\} \text{ i.e. the ordering is wrong.}$$

I.e. if the ordering is wrong, we can  $\downarrow c$  by reordering the  $k, k+1$  elements in the ordering.

That's Algebra come out a lot easier than it usually does!

32:5.40

(25) There are things business in normal G.A. of making offspring of G or  $\Delta G$ :  
Say no. of offspring =  $k \Delta G$ : Then  $k$  is like  $\beta$  in sharpening the bias toward  $G$  or  $\Delta G$ .

- i.e. "like using  $G$  for reinforcement. I think  $k$  is usually chosen empirically."

1.33 This is a gen. diff. in G.A. - But any monotonic func. of  $G$  will have same peak  $\beta$  &  $k$  values

1.36 be equally good for GVM optimization finding: but its not true in GA, so GA must be wrong! (see 3.03 for another criticism of GA)

In 5.34, we are concerned only with ordering, so objectives 1.35-1.36 don't apply.

We can get the effect of  $\beta$  by giving more or less wt. depending on ~~the~~  $\beta$  to absolute ordering of the cond. parts.

Another thing to consider is the distance betw. pairs that are in wrong order.

21  
 19.20.98 (26) T59. for Analyt. tasks (in Conv. below): Start w. large  $\epsilon \in (\text{Dep. units } \epsilon)$ , then  $\epsilon$  as tiny proceeds to  $\text{Total } \epsilon \rightarrow 0$ . Since  $V.$  has  $\geq 2$  parts to fit curves of various types w. vary  $\epsilon$ , this miter be easy to do.

12:30A any card.  
 12:30A  
 3:45A 3#  
 7A 3#  
 12:30@ any card

For discretely det. task: slowly  $\downarrow$  no. of errors made!  
 For "GPS-like" task (vector error). Use slow  $\downarrow$  in some axis of error vector  
 For  $k \cdot 2^k$  input Max. start w.  $k=1$ , then 2, etc. — interestingly only, they  
 how it learns even if it starts w.  $k=3$  (or  $k \gg 1$ ): It suggests that w. only "and, or, not,"  
~~it~~ ("if" optimal... "if  $k \in \mathbb{N}$ "), it could learn  $k=1$ , then  $k=2$  act, even  
 if  $k$  were  $\geq 3$  for  $k=3$ , initially — so I would boot  $k=2$ , even if "and, or, not" were present.

Green:  
 Back by  
 ~ 2PM

26 1/2 So w. 2 imp. Q's about GA (i.e. SGA)! 1) How no. of children depends on  
 # of parents 2) # IS mutation or crossover better? I find it hard to believe that  
 people seriously consider "Mutation only": Each individual has to find its way  
 up to hill, himself. "Generalized crossover" enables individuals to pool their  
 experiences to devise optimum "child" cand. (in Mut. only)  
 Suppose fitness are extremely expensive (Human, animal large animal reproduction).  
 Then one should take as many of known  $(O_i, G(O_i))$  pairs as poss. or use them to devise  
 $O_j$ 's w. either expected  $G$  or  $O_j$ 's fitness "vary informally".

Side of kinds of  
 probs TAB can  
 work on

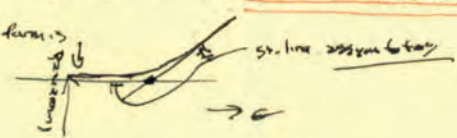
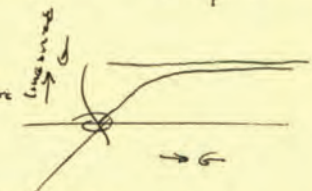
17  
 18

27 E wrt Goal 2! Just pass a stack (or  $\text{Pen } Y: \sum O_i G(O_i)$ ) data set — or, more  
 exactly, ~~Get the~~ <sup>optimal</sup> ~~operator~~ <sup>operator</sup> for  $O_i$  input. (I ~~think~~ have to think this  
 thru, better! — Recall + "One shot long" Q: ~~can't~~ <sup>can't</sup> ~~do~~ <sup>do</sup> OSL properly? OSL is very imp., very common  
 P.D.  $\tilde{G}_p$  —  $\tilde{G}_p$  T. expected values of  $\tilde{G}_p$  is  $\tilde{G}$ , say. So we could have a  $\tilde{G}$  function  
 of  $O_j$ , so  $\tilde{G}(O_j) =$  expected value of  $O_j$ . Here, this may be missing a point!  
 If we do lots of  $O_j$ 's w. large variances in their  $\tilde{G}_p$ 's we could get away  
 by  $G$  real way! (If  $\pm$  D.P.'s are sufficiently uncorrelated)

26 T. for  $\epsilon$ . system looks very much like SGA: (i.e. pass a stack (or)  $\text{Pen } Y: \{O_i, G_i\}$   
 data set  $\{$  w. additional constraint that ~~there~~  $O_j \geq \tilde{G}(O_j)$  is a max is not hard to find)

28 (28) SN It may be that w.o. linearization of the conv. that is General OZ  
 (incompatibility)  $\text{Pen} \in \text{fixed (comp. Biol)}$   $\rightarrow$  each ~~is~~ <sup>is</sup> opt. strategy will have  
 a d.f. on final  $G$  values. Unless we have linearized  $G$  we can't tell whether one D.P. is  
 better than another. Here, in many cases, the distance among in choice given by  
 various linearizations will be quite small. In some extreme cases, here, it can probably be  
 enormous.

Example of very non-linearized conv:  $(1 - e^{-G}) =$  linearized conv  
 So for  $G \gg 1$ , it doesn't pay much to  $\uparrow G$ .  
 One could do a rough linearization of the conv., (if poss.) before doing SGA.  
 In this case, if  $\tilde{G} \gg G$ , we can't tell if  
 for  $G \gg 0$ , but distance much from  $G$  if  $G \ll 0$ .



11: (28) (cont.) For  $G < 0$ , hvr, those are low  $G$  pts, & we don't care if  $G \neq \tilde{G}$  in that region. The only reason we might want  $G$  to be close to  $\tilde{G}$  in the low  $G$  region is that this would enable us to use low  $G$  data to fit the  $\tilde{G}$  curve better in (hopefully) higher  $G$  regions as well. (~~5.34~~ — discusses explanation of "ordering").

05 (29) Review of impt. ideas! \* : Need Work

- (A) Understanding of proper soln. of Z(4): 2.20 - 2.40
- (B) Development of  $\epsilon$  Prices in Children: 4.01 - 4.0.

.09 (C) 7.26 : An understanding of ZSGA! How it may be a kind of optimum soln. of Optim. problems (not exactly OZ problems).

.11 (D) 7.01 A Bunch of kinds of Iraq that can be put into the form of TSG's! Ptz list could be readily expanded. 87.01, .32 is better, hvr.

(E) 6.01 : proof of 1. first Grainly House Arm ( $\frac{C_i}{P_i}$  is best ordering),  
 9.20 (5.39) " " Second Grainly House Arm. Time to solve  $\leq \frac{C_j}{P_j}$  ( $j$ 's soln).

\* (F) (5.34) : On "no. of Misorderings" as a score for fitting a  $\tilde{G}$  to  $G$  in ZSGA. Also note 5.39 and 9.01.

.18 \* (G) TMTS. 25-76A TMA, TMB ~~...~~ A clearer understanding of My New Method of Assigning TSG's on Tug. events": (Also note 7.10)

\* (H) Mutation only v.s. Crossover only (or Crossover + Mutation) : A paper to understand! 7.03 - 17  
 (Also para) Q of How normal G.A. works (crossover):  
 (H\*) "Building Blocks": MTK 1.24 - 1.29; 3.04 - 3.09; 3.09 - 2.27  
How it works: Good Subproblems: BUT NEEDS WORK!

5.34 on 1  
 optimization  
 => score.  
 1.24 (15)  
 discrepancy.

9/21/98

9/22/98  
12:01A 2mg card  
12:01A  
2:25A 2#  
5:45A 3  
8:30A 2mg card

**(NOT a MATH Problem!)**  
Amusing ~~prob~~ problem! Given a sequence of  $n$  objects that are in  $k$  "wrong order".  
What's the least amount of info. needed to put it into right order?  
"Cost" of a particular code is unclear! Ex: 2, 3, 4, 5, 6, 1 is observed order.  
We have 6, 1, 2, 3, 4, 5, 6. We need do w. 5 exchanges: (6, 1), (5, 1), (4, 1), (3, 1), (2, 1) or one "jump"  
"I from on end to other." - assuming others all moved.

9/22/98

11:20P 2mg card  
9/25/98  
2:15AM 2mg card  
2:15AM  
5:45 3#  
8:30 2#  
9/30  
8:30 2mg card

Hooper's approach: List all  $n!$  orderings, test which ordering it is in. For large  $n$ , almost all orderings are "random" (I think).  
Actually,  $G$  will "order"  $O_i$ , but several  $O_i$  can have same order! Similarly w.  $G$ .  
Always assign part of "Business" to an ordering! Consider it to be created by "noise" of some kind (Gaussian,  $e^{-|x|}$ , etc). Noise, hvr, is most likely to occur uniformly over  $y$ . range/domain of  $G$  if  $G$  is "linear".

9/26/98

9/26/98  
2A 2mg card  
9/27/98  
2:00A 2#  
10A 3  
3P  
3P 2 mg card  
9/27/98  
2:00AM 2mg card  
2AM  
8:45 3#

That was more interested only in ordering: Arg't: Set  $G$  results are using noisy  $G = H(k)$ .  $k$  is a simple funct. (Short term), but  $H$  is a much more funct based on random fract, so it has high cost. So we'd have lots of trouble approximating  $G$ , but much less trouble sampling  $k$  (via "same ordering").  
Note Also arg't of 5.25-5.26 on why pts. of low  $G$  can be very "informative" about

9/28/98

9/28/98  
12:30A 2mg card  
12:30A  
9/29/98  
2:05A 2#

shape of  $G$ .  
front of second Comb. house Perm: If  $\frac{c_i}{p_i} < \frac{c_j}{p_j}$  then  $\frac{c_i}{p_i} < \frac{c_j}{p_j}$  order

9/29/98

9/29/98  
12:15A 2mg card  
12:15A  
3:30A 3#  
6:15A 2#  
9/30/98  
2A 2mg card  
2A  
5:30A 3#  
9:30A 4  
11:15A 2#  
11PM 2mg card

Time to find sol'n  $T = \sum_{i=1}^j c_i$   
So  $\frac{\sum_{i=1}^j c_i}{\sum_{i=1}^j p_i} < \frac{c_j}{p_j}$   
This proof suggests that  $T$  is usually  $< \frac{c_j}{p_j}$ , because  $\frac{\sum_{i=1}^j p_i}{\sum_{i=1}^j p_i} < 1$  is usually strongly on the " $<$ " side.  
If  $\frac{c_i}{p_i} < \frac{c_j}{p_j}$  then  $\frac{c_i}{p_i} < \frac{c_j}{p_j}$  order  
Since this is  $< 1$ ,  $T = \sum_{i=1}^j c_i \leq \frac{c_j}{p_j} \cdot \sum_{i=1}^j p_i < \frac{c_j}{p_j}$

9/30/98

Building Blocks: ~~from~~ from 1.24-1.29; 3.04-3.27: The Agreement of 1.24-29 in values  
~~from~~  $f_{AB} > f_A \cdot f_B$ , which is used for code compression in ALP.  
Hvr. to cond. for defn. of  $A \cdot B$  is that  $G(A) \cdot G(B) > G(AB) > G(A) + G(B)$   
(i.e.  $A$  &  $B$  are "synergistic" concs. &  
& let  $G(A)$  act. mean: max.  $G$  of concs containing conc.  $A$   
Probably so far we can talk about we must linearize  $G$ . This is also a way  
for detection of "Synergy".  
So this "Building Blocks" idea is beginning to look rather definit from 2141.  
One way to look at "crossover" is that we do 2.18! pass a stack operator Perm 1: Set  $\{O_i, G_i\} \rightarrow 10, 0$

9:40 Then we use this stack operator to try to get by  $G$  cards: In crossover, we reduce the ~~SSZ~~ SSZ  
 from  $\left[ \begin{array}{l} \text{1. cardinality of} \\ \sigma_1, \sigma_2 \end{array} \right]$  to just  $\sigma$ ! Using Koza's functional ~~operator~~ <sup>crossover</sup> (of function trees),

We do get a stack lang. from any pair of cards. Since the ~~SSZ~~ SSZ of  $\sigma$  is very small, <sup>2</sup>  
 the nature of the induced stack lang. depends heavily on the applied. In deriving a stack lang.  
 from a pair of cards, we do get by  $\sigma_2$  in the  $G$  posn.

11:15P 2mg card.  
11:15P

A1) R! OXB, Try starting in  $(\frac{1}{3})$  days later each quarter! See if phasing is critical in  $u, \sigma, \frac{u}{\sigma}$ .

92298  
3:45A 4 1/2  
8A 4 1/2  
8A 2mg card

A2) Machine Learning is the Goals of AI. Where are we going, and where can we go. Title for talk at MIT. w/oc/sechun.

11:30P 2mg card.  
11:45P  
92398  
2:15A  
4:30A 2 1/2  
7:30A 3  
11:45A 4 1/2  
1PM 2mg card.

A3) VAP NK:  $\epsilon$  for curvature w. n data pts: if reverse ms norm.  $\sum \epsilon = 0$ , t. no. of 'support vectors' is n.

92498  
12:15AM 2mg card  
12:15A  
3:45 3 1/2  
5:45A 1 1/2  
5:55A 2mg card.

(b) If we do 'separation of classes' w. hyper manifolds, and ~~we~~ only in complete separation is possibl., we have  $\epsilon$  certain no. of "errors". I think V. Eines = w.r. to all error pts.  $R^2$  is seems very wrong.

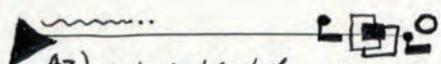
0198  
12:15A 2mg card  
12:15A  
4:30 2 1/2  
6:45 2 1/2  
9:15 3 1/2  
11:45A 2 1/2

A4) Give  $\epsilon$  sources of 3 or 4 rules at FAP! First similar to London talk — (23 mins) then following lectures Gerry More detail.

11:10A 2mg card  
02298  
1AM 2mg card  
1AM  
5A  
6:45A 1 1/2  
8:45A 2  
11A 2mg card.

A5) T. World of Math: J.R. Newman: vol 2: interesting w/for prob: also QM: by Schröd, Hermit. Also C.S. Pierce.  
A6) Arthur Koestler: <sup>Boer</sup> 1) Strangar and Spener: Arthur Koestler 1910-1950 <sup>A.K. died 1983</sup>  
2) Arrow in the Brood: 1905-1931 <sup>+ Born</sup> joint comm party.  
T. Invisible Writing: 1932-1940 + (w/relig 1940-53)  
3) T. case of t. Milutin (odd) t. roots of coincidence.

11:45A 2mg card  
11:45A  
11:30  
6:30  
10:45 4 1/2  
11:45A 2mg card

~~~~~  Challenge of chance.

4:30P 2mg card  
11:45A  
0 + 98  
3:30A 3 1/2  
6A 2 1/2  
15A 1 1/2  
9A 2

A7) What kind of sup-carpus is OZ data? 487

10:30A 2mg card  
0598  
12:01A 2mg card  
12:01A  
5:45A 5 1/2  
8:15 2 1/2

compass  
Starstruck  
Our Coast, Good  
They're like us  
Blue room  
Handbag - Pansy



92298 General TM vs. Indexing input (10) (27) TM87 11

1) t. CPD can contain all kinds of "ling" - including statements, "whenever", an "always", I interpret it self/ies part of a complex to be integrated.

T. why this works: comes esp. for compressing internet (say) can be used in solving other problems. We do have to have a cond. p.d. that suppresses which comes ~~more~~ likely to be usable whenever. (t. IR problem).

2) (very) Ex. in this work (perhaps just later) This tip of data <sup>set</sup> could be used as a "Hint". Ordinarily t. data comes in partitioned somewhat (at least, temporal ordering). It can also be partitioned by "logic" by t. trainer.

This last partitioning can be indicated by an index etc. This comes w. each data chunk, e.g. each of t. chunks of .01-.02 would be assigned an index no. Subscript could also be given. Subj. (6 sub...) indices by trainer. Later, TM itself may want to do some index experiment.

27 In opposite direction from <sup>(partitioning)</sup> indexing is "feeling" - an input part of t. induction process. Feeling  $\uparrow$  ssz, indexing, partitioning  $\downarrow$  ssz, but in a very useful way - it ~~the~~ size of many in which to hunt for usable sub. comes to combiner.

"Indexing" is perhaps best defined as a hierarchical type of "partitioning".

32 ~~As is~~ As is, it appears that we can integrate into of various kinds into "Grand EPD". One kind of input that might be very useful: Traces of T.S. or problems worked by other Learning Machines: (preferably machines  $\neq$  to TM) This would in I doubt, TM would be able to integrate

1) Such into into its own Grand EPD. This is ~~an~~ v. g. way to integrate TM's that have worked on different problems or in different problem areas.

2) I want to write up <sup>7.19-20</sup> (8.09), (8.11), (8.18) in some detail: My new Grand Plan for TM: How far I've gotten: what's solved & what's just what's unsolved. Do it mainly in "English" Don't get too details yet.

Handwritten notes at the top of the page, including a date and a title.

First main section of handwritten text, containing several lines of notes.

Second main section of handwritten text, continuing the notes.

Third main section of handwritten text, with some larger, bolded words.

Fourth main section of handwritten text, concluding the notes.

Vertical handwritten notes on the right margin, possibly serving as a checklist or index.

Large, bolded, mirrored text in the center of the page, likely bleed-through from the reverse side.

This will be a review of The most recent "Grand Plan" for TM!

Imp't refs: 8.09 (E) 7.18-26 An Understanding of SGA: Toward a class of Good Sol'ns to Opt'zn problems (not exactly OZ prob's)

8.11 (D) A Bunch of kinds of problems, inputs, tng g'v'ns This latest TM should be able to do/w. See 7.01 vary appendable! <sup>11.01</sup> <sup>12.04</sup> <sup>particular</sup> (E 87.01-88.04)

8.18 (G) Main PLANs: TM 75.25-76.40

TMA: 72.02-40 <sup>Expo</sup> finds short codes for arby corpus types: sequential a/o unordered set of <sup>object</sup> <sup>or partial</sup> <sup>obj</sup> <sup>ct</sup> <sup>s</sup>

TMB 66.01-65.40 Conditional probly lng. Machine: Is able to solve great variety of problem <sup>Also, P.D. leading upto 66.01 are very relevant (good)</sup>  
by pass using a <sup>(common)</sup> <sup>(Grand)</sup> Conditional Probly Distrib'n (CPD) <sup>7.01 a. etc</sup>

Re: TMB: There is some confusion in my mind as to just what it is!

66.01-65.40 was an early picture: I think I modified <sup>with</sup> <sup>(7.01 was minor)</sup> <sup>87.01-88.04</sup> on t. modifi, before writing <sup>(may be no. writing which</sup>

So maybe an early idea of t. <sup>(conditionally)</sup> <sup>CPD Machine</sup> is: In steady state, it has within it, a CPD; This CPD ~~is~~ takes form of: we put an augmentation of t. corpus into it - including

a "partial augmentn": It's able to error <sup>P.D. for t. contin. of t. partial corpus.</sup> <sup>a/o minimum</sup> <sup>(leads to</sup> <sup>very large</sup> <sup>continues)</sup>  
(t. "augmentn. a/o t. partial augmentn. can be null.)

At first, TM would be able to do this in carefully labeled "Domains" w. no interaction betw. domains. Each Domain would be completely separate from all others <sup>others w. its own</sup> <sup>separated</sup>

TSQ. Later, TM would try to integrate the P.D.'s of all of the domains into a Grand P.D. for all domains: At first each domain was like a Sol'ng machine w. its own

Expanding this idea to 87.01-~~26.09~~ <sup>12.09</sup> was an imp. <sup>But "ambitious" probly</sup> <sup>for IR in</sup> <sup>its own memory</sup> <sup>to reduce</sup> <sup>reasons</sup> <sup>IR from</sup> <sup>Memory</sup> <sup>67.12 (67.12)</sup> <sup>13th imp. Q.</sup>  
which detailing. This last <sup>expanded</sup> <sup>considerably</sup> kinds of problems TM could work on a/o types of info it could use for inputs.

(SN) TM has, at all times, either t. original corpus or a very easily decoded compression of it a/o compressions using very hypc definitions: like words in english separated by spaces, <sup>Prostate in dictionary</sup> will be ~~or~~ <sup>2000+</sup> specially coded.

29: 67.13: Later I began thinking <sup>about</sup> this "sufficiently advanced" ~~the~~ state of TM: ~~the~~ Presumably, at this point, it would be rather easy to teach TM things. One could design TSQ's in an instructive way - as one does for human students. At an even more advanced stage, we could give it ~~the~~ text books designed for human students. So one big Q is: could we design TM so that Ab mita it was at or fairly close to ~~design~~ one of these levels of intelligence.

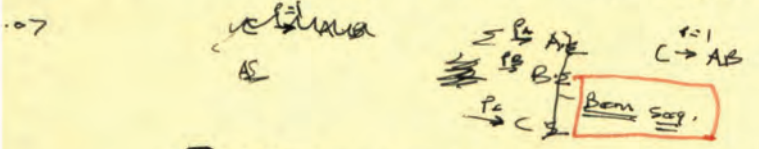
92398 TM 2141

- 01. (312C) (2.10): When parsing a corpus w. CBA substitution: One can never be sure that a specific  $C \rightarrow A.B$  is correct, unless that  $C$  has been made part of a larger n-tup. defn. — otherwise,  $\neq$ !
- 03. Assl. parsing (with meaning) & substitution  $C \rightarrow A.B$  size of = p cost.

→ overconvergence for

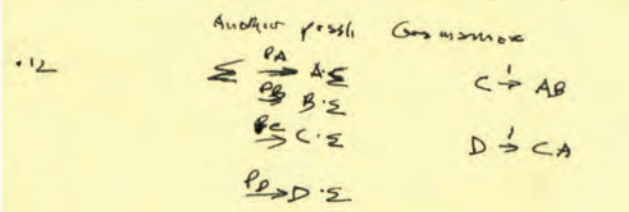
→ 109.01

T. essentials of stochastic  $\geq 141$  lang. are: C.F. rules, but no loops allowed. Since  $\infty$  lang is finite, there are a finite no. of n-tups defined. Each n-tup has a frequency assoc. w. it, that can be calculated from the pc's in the grammar exp.



→ partitions to whole grammar

T. original corpus only has A & B in it: they occur w. frequencies  $f_A + f_C$  &  $f_B + f_C$ , resp.



observed freq. of A & B in the original containing A & B, are:

$$f_A = f_B + f_C + f_A$$

$$f_B = f_C + f_B$$

So we take this corpus containing A & B that seems to have dependencies between A & B, we make a Barn seq. of n-tups each being a finite seq. of A's & B's.

I. probys in Barn seq. are indep.

So: I idea of "2141" is that the "Barn Grammar" of .12, with defns of C & D is

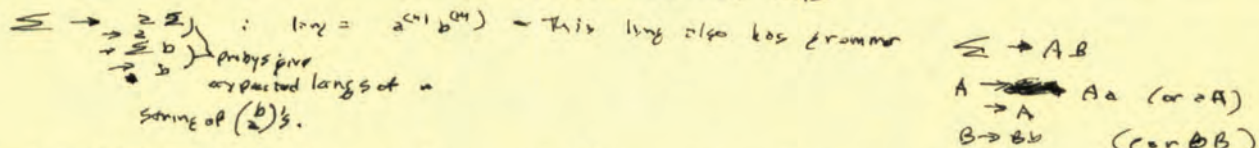
- a) the hyper pc. of defns of the corpus, <sup>then</sup> (a) the original corpus as a Barn seq.,
- or (b) any other grammar of the type. (b) is rather uncertain: we usually cannot try all possl. grammars of this type... Unless we find a suitable way to do it in reasonable time.

After finding a rather large no. of defns of the type that  $f_{pc}$  of corpus:

Can one find "looping rules" if it is, indeed, an infinite lang. w. loops?

- Well, just the grammars of .07 & .12 both divide into 2 parts containing
- (a) Subsets of  $\Sigma$  w. various probes:  $\Sigma$  is a part containing about subsets of all probes / No loops. Contains a simple loop.
- Look at prob (b) alone & see if we can find loops to deduce a (b) by  $\uparrow$  pc of defn.

Lets look at (b) or more simple CFG's that have loops



(a) look at this lang: ab occurs readily: we will deduce  $A = a^2, B = b^2$ .

$$B_1 \rightarrow B_1 b, B_2 \rightarrow B_1 b, \text{ etc. } A_1 = A_2, B A_2 = A_1 a, \text{ etc.}$$

$$\text{so } B_{n+1} \rightarrow B_n b; A_{n+1} = A_n a \quad | \quad \text{Both of them suggest } B = B$$

10)  $A_0 \rightarrow a_2; A_1 \rightarrow A_0 a_2 \dots A_{n-1} \rightarrow A_{n-2} a_2; A \rightarrow a_2; A_{n-1} \rightarrow A_{n-2} a_2$

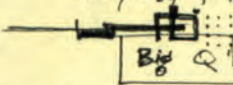
so  $A \rightarrow A a_2 \rightarrow a_2$  } Better. Maybe no gets this w. folk "reasoning"!

th.  $A_n$  objects form a class because they are all of form  $A_{n-1} a_2$ .

The  $A_{n-1}$  objects form a class because they are all of form  $A_{n-2} a_2$ .

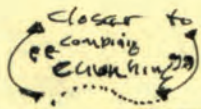
so try  $A = A' a_2$ . But no! Are  $A_{n-1}$  classes are not same.

Any way, for some reason or other a large set  $A_{n-1} \rightarrow A_{n-2} a_2$  suggests  $A \rightarrow A a_2$   
 $A_0 \rightarrow a_2$   $A \rightarrow a_2$



Answers being even many ss of  $a_2$  large or just one  $a_2$  from  $v_1$  large?

In original 2141, I was analysing inf. in finite string. When I discuss "loops", I'm usually talking about a corpus that's a large set of finite strings (= ass). This last is perhaps



closer to what I want to use 2141 for: i.e. Pro coding by creating of new codes by old codes.....

17 Even w. the single infinite string corpus my 2141 analysis is. o.k. (i.e. expressing it as a form set.)

The idea of further compression by compressing the grammar alone is interesting.

19 In general, this '2 step' method of getting a "looped" grammar, is not efficient!

i.e. Many users will not initially be defined because their defns are too expensive.

Later, these defns would  $\downarrow$  in cost thru "Looping": So the defn of users is  $\uparrow$ .

detection of "looping" should, ideally, be done together, for max efficiency (efficiency means max compression for a given corpus).

Hvr. my other method "2 step" soln. of 17-19 to get a fact for what's going on — as well as a possibly "default" soln. to the PSG'd problem.

~~At least 1 problem's~~ But we might (inadvertently) write a "first-step" grammar

25  $A_0 \rightarrow a_2; A_1 \rightarrow a_2 A_0, A_2 \rightarrow A_1 a_2 \dots$  which would waste SSZ. We could have various other "free" equivalent strings of  $v_1$  first-step grammar set of  $v_1$  25<sup>th</sup> form into  $v_1$  form

So to study  $v_1$  form of "step" gram  $\rightarrow$  (Loop Grammar?) look at various CFL's

to see how various loop & instruction sets generate to various large forms!

My Paris paper is of interest (sol602). The idea of Loops or (linked loops

of "loop nested" pts: like  $A \delta_1 B \delta_2 C$   $\delta_{1,2}$  are loop makers.

34 so  $A a^n B b^n C$  could be an acceptable sentence  $a^n$ . To detect such grammar rules,

26 506 The SOG (step one grammar) needs to have rules of type  $A_n \rightarrow a A_{n+1} b$  or equiv. ("give" means in this case  $A_n = a A_{n+1}, B_n \rightarrow B_{n+2} b$ . Hvr. read  $[A_n] = [B_n]$ )

can be f. SOG,  $\rightarrow$  unclear. So perhaps in crazy SOG's try to put them in form 34 if possl.

summary: ~~Designing~~ SOG involves trying various forms to see which best "next level" compression (i.e. most new user defns.  $\rightarrow$  109.0) next highest level).

I want to continue to make list of "types of TM inputs" of 87.01 & .32 this gives good ideas on how TM can work, what kinds of TSO's would be very useful: → .08

.03 7.01 dorks = Analog TSO's: How to design them. Also another some other kinds of TSO's that can be represented as optzn probs, of 1. type workable by GA & probly by SGA. It mite be interesting to design TSO's for parts of Robot arms that have various constraints. Also, learning to walk - to catch balls, throw balls using robot arms. To play ping pong, using various constraints on speed of paddle motion - constraints on speed, reach, accuracy of optimal pickup. Constraints of hysteresis in servo action to move paddles, etc. (93.36 has a good approach to ping pong) just on V.G.

11:50  
1:30  
13:00  
11:56  
12:50

So essentially, we would use commonly hardware! Or not so bad (not so bad) (C)

.07 H.W. for starts, a slow degradation of H.W. as a kind to TSO. → (93.36 94.05)  
.08:02: Some kinds of inputs to TM. ① statements (usually, but not nearly "true") - they can be facts, heuristics, hints of various kinds for specific problems. (shades of MCC's "Advice taker" [I'm not really forgotten details of what kinds of advice it took, & what param it had.]) ← perhaps related to .32?

87.88:  
Vernica pp!

② "Wired in" info: I'm not sure just what this should mean. One possy: would be a computer z/o set of computer programs, that TM would learn to use in solving various problems: say Mathematica or Maple, or various spreadsheet sheets. Mathematica & Maple contain various sections designed for various domains, TM would have to learn ② how to get access to these various modes or problem solving & ③ how to use each mode. Maple Corp. has online help, but TM would have to know how to read,

to some extent, to use this info - In the case of Maple Plus "online info" contains very little reference to RW, so it may be a good place for TM to learn English! Good, because TM would, presumably have much knowledge of the domain before being trying to learn to understand English text about that domain.

③ Encyclopedic: This could be put in "online" as a whole encyclopedia (or, if internet) or it could be put in as a "text understanding TSO", in which we started w. a simple corpus w. simple tags & slowly moved toward the ultimate unrestricted Encyc. text. As for the Internet: In attempting to delay TM's discovery of RW the internet could be put in, in Recorded Snapshot form (e.g. Brewster's "snapshots"). The idea here is that TM does not find it poss. manipulate "snapshot" and could do so w. f. real internet.

A perhaps unavoidable way that TM manipulates to convert. is benefits answers to the user. We could reduce this effect by pre-programming our sequence of requests to TM (how log vast; data is a func of all previous actions/requests of TM)

32 ④ (87.32-88.04) → perhaps .03? Various kinds of info usable by TM as inputs:  
② Traces of problems worked by this TM in past, as well as traces of problems worked by other TMs (particularly w. to this TM) or in general, traces of various problems, given in verage degrees of detail. Text books should be a (perhaps extreme) case of small bits of detail. Another would be proofs of Math Thms & solns. of Math probs given in Literature.

⑤ Analog Problems: 92.03, 93.36, 94.05...

.01 An approach to better understanding ("Detailing") TM<sub>B</sub> could be obtained by dropping various ways that TM would deal w. each of ~~inputs~~ its input types: The input types of 92.02-40 are ~~some~~ "novel" & can be of much input!

Some more conventional input ( $\equiv$  problem) types:

.04 (1) Extrapolator  $\equiv$  set of  $[Q_i, A_i]$  w.  $n$  &  $W_i$  ~~OSL~~ OSL (One shot long). I think ~~these~~ all 3 types I've written are of this type -  
- But for MTM: I could do it for harder MTM probs & for NMTM probs.

(2) Simply find short codes for a corpus. I haven't really tried this for a large ordered corpus ( $\approx$  "regression") - Re. of was done as a kind of

.09 Short code finding problem v.2 Lsach.

A perhaps implication: Just searching for short codes is always best done by LSach, using conditional polys learned on earlier sections of the corpus

.12 [SN] One way to solve IR problem for many access: ~~was solved for simple~~  
For example: we solve our simple Alg. notation Long problem as we already did - w. somewhat larger cost of many access. When this occurs, the system is in good shape for looking for CPD's on many that  $\downarrow$  cost of its access! The way to problem gets, the more benefit we get by finding a good CPD! Think of various (usually "soft") partitions of the many that could be helpful in retrieval. Recency is one.

.12 (12) perhaps the first real use I've specified for CPD: Re. it is a bit vague, its still looks very promising. The vague part has to do w. the kinds of partitioning/indexing ML use TM will use for this problem! I'll perhaps design a language (instruction set) to implement various indexing, partition techniques. (See 87.10 for disc. of "partitioning" vs "indexing" vs "pooling".)  
For partitioning  $\equiv$  set we can use Boolean functions & subsets.  
Hierarchical indexing is a kind of partitioning using AND only. We can then "OR" these / conjunctions. Negation is not normally used - but perhaps should be.

.33  $\rightarrow$  What about Relational data bases & the "Join" operation? This is certainly an input kind of operation & a good way to organize data! (whatever it is called!)  
Lots of Scientific Data takes form of flat files that can be "Join" in useful ways!  
Hvr., I should look up my writings on Rel. data bases (in Errata files) - also, perhaps work of others on Relational D.B.'s! Perhaps ask Pat O'Hall for good Refs.

.36: 92.02! HA! The Carrollian simulation! Since analog. expts. in R.W. are expensive, T.M. makes a model of R.W. using whatever "Hard" data  $\equiv$  expensive data is available. It's a "H.C. problem w. expensive trials"  $\approx$  much time is spent devising a v.g./model.  $\rightarrow$  94.05 spec

92498 TM

# ANALOG Lrng (.05)

.01: 93.09 <sup>93.11</sup> Findings short codes for corpus. 93.10-11 is a set of <sup>rationall</sup> ~~rationall~~ for "coding & decoding".

I'm not so happy w. my "soln." of the "Alp. notation Lrng" problem: T-machine used, had too much info in it - i.e. it was a "stack machine", & it had to (run on an) effort. - But

.04 2000, maybe a stack machine is not so a.H. after all! - Since CRL's are common (I think!),

.05: 92.07 <sup>93.10</sup> (SN) An interesting kind of analog (lrng.) TM (int. to simulate motion of → 96.20

real people, for simulated actors: As input, have motion of real people.

~~Have~~ Have outputs of machines be many dimensions & have them be (in frequency)

"scored" by humans, (clearly, this is an example of expository R.B., so the remarks (2 methods) at 93.36 - 40 are relevant.)

~~A more difficult problem:~~ A more difficult problem: we assume ~~in~~ that TM has been given a way to take human (video, zdm, say) figures, & it ~~is~~ has been given parameters can move to limbs, torso & head, etc. A more difficult problem: TM is given

~~video~~ video motion pictures of humans moving about. Say ~~a~~ persons

A, B, C are shown doing various "things": T-machine must

A is shown doing tasks  $\alpha, \beta, \gamma$ ; B is shown doing tasks  $\epsilon, \delta, \eta$ .

We then ask (somehow) for TM to make a video of "B" doing task  $\epsilon$ .

One application might be Blackmail: To get TM to make videos of the Blackmailers, doing illegal &/or embarrassing things. Sell to CIA, FBI, Mafia? ☺. - See .37 for a poss. way to do this.

The ultimate goal, however, would be to be able to give TM a descr. of a person doing a certain thing ~~as~~ as well as perhaps some pictures or videos of the desired actor - & TM would output ~~the~~ a video of that actor doing those things. This would be a step toward having TM present plays, from the ASCII text of the play.

Or TM scans a play done w. one set of actors & makes a video output of the same play using different actors.

The language tasks for TM may in themselves not be of much interest as goals (or outputs), but writing eng. papers on eng. Yakuhanase for Posen would be of interest.

.37 One easy approx. of a lot of the above could be done by Morphing one person into another. Good Morphing involves intelligent selection of



01. Covers random pts. in the 2 objects being xform (Morphd). TM could know how to do this.

Whoops! Assuming ordinary morphing would ~~not~~ do it! We have ~~Person~~ Person A in position  $\alpha$  & Person B in position  $\beta$ . To construct person B in position  $\alpha$ .



Say we had lots of pictures of B in various positions thru various directions. Then we could morph B onto A's position  $\alpha$ .

12. Other problems like R. Ferguson would be to simulate a certain person given many "voice samples", saying an arbitrary thing, or singing an arbitrary song. 12 might be a lot easier to do than the video of 94.05 ff. - yet the main ideas of the TSP's could be very N. 12 might be a TSP, for me on how to do 94.05 ff.  $\odot$ .

Actually, this is not so easy. People have dialects, so ~~can~~ may say certain words in certain ways ("Mat's alt.: special cases" exceptions).

An "incomplete" set of words said by a person could result in a simulation in which he said a word in the wrong dialect! This problem would have

to be carefully studied: perhaps by humans, perhaps by TM!  $\odot$ . (99) 6.12

25. Several Major (General) Robot TME: Rowenjin

- 1) How to ~~write~~ <sup>solve</sup> TSP's for various modes (Domains).
  - 2) " " Deal w. IR in memory for each individual domain
  - 3) " " Integrate <sup>various</sup> Domains thru common memory
  - 4) Two "OSL problems" for each domain.
- is a \$500 IR problem. Note treatment of this problem in case based reasoning & Explain. Best Lang.
- Priority after previous in integration (i.e. "Data fusion" or "Transfer Logic")

I have been acting as if the only thing RM had to do was "concept acquisition"

T. aquisition of a concept ("knowing" of it) is giving it  $\oplus$  cond. p.c.'s for various conditions. Conditions

But it is an essential part of logic. mechanics, etc.

I want to look at it more carefully! Also, consider a variety of cond. probs: they can arise in various ways & be evaluated in various ways. Perhaps the simplest cond. prob. is xform of flat "Many of concs. within a Domain" to first, an ~~non-flat~~ non-flat, uncond. prob., then a True cond. prob.: cond. on when each conc. should be used.

Next, ~~the~~ ability to use concs from other domains.

39. In all cases, the formalism of 96.01-12  $\Rightarrow$  C21 (should) be used, to get optimum cond. prob. distrib. So in Priority (96.01-12) .39 should "solve" problem.  $\rightarrow$  98, 21

.01 On 8. Gave 2 kinds of complete probabilistic Induction problem  
 We have <sup>for corpus</sup> a set of finite objects: is a set of partially ~~finite~~ partially complete objects. This later may be parts of a finite ~~or~~ or an infinite object.

Example of such a corpus:  $(Q, A)$  problems in which we are given many complete  $Q, A$ 's. also many  $Q$ 's w.o.  $A$ 's.

Guess at soln: A group of  $k$  corpus is

$$\prod_{i=1}^k \left( \sum z^{-|P_{ij}|} \right)$$

**Note:** on 99.01 we discuss t. machines that could have created this corpus. Eq. 107 must be mult by  $t$ . apr. of class machines and summed over all poss. such machines.

There are  $k$  (complete incomplete) objects.

$P_{ij}$  is  $t$ .  $j$ 's "minimal"  $P$  for  $P_{ij}$   $i$ 'th object.

We may want to have  $\geq 5$  symbols:  $0, 1, (\leq \text{stop})$ . There is no  $\text{stop}$ 's data unless  $t$ 's is omitted, it continues to  $\text{float}$ 's forever.  $\rightarrow$  99.01 = cont.

$\rightarrow$  A kind of task TM should be able to do: we give TM a diff. problem; then we show it how to work a "related" problem & TM is able to work diff. prob. An example could be solving solving of cubic or quartic eqns. or various kinds of symbolic integration. (u. ~~transfer~~ "Transfer Learning")

.20 : 94-04

Rita now, I don't really remember to have tried behind  $t$ . Alg within Long.  $\nabla$  99.

Get D did. It did use a "stack machine" which is a useful. I'd like to machine to work ~~the~~ thru a covering  $R$ . concept of "valve" is parsing to obtain "values";

substituting "values" is prevalence, etc...  $\uparrow$  way humans do it. Sandy proof: I did a few bits of work on this issue in 1990 (Sandy).

$\rightarrow$  I could start out w. a TM that already "had"  $t$ . conc. of "value" & go on from there.

In a sense, this would be giving TM some intuition into what humans may have on quite early in life - or by transfer long from other domains.

How would this conc. be represented in TM? One way: Net assoc. w. certain substrings is a "valve" & that "values" have certain (useful) properties.  $\rightarrow$  103.19

.30

Another kind of TSO that I've considered, but not actually done work on would be to try to get TM to get do a tsg. based on  $t$ . order of info & examples & problems in an ordinary human textbook on Algebra. At first, I was worried that certain (perhaps essential) parts of  $t$ . TSO implied by  $t$ . Book would involve concs. from  $R$  w. inaccessible w. diff. to TM. Now, I think that I can deal w. such concs. by writing alternative Bob TSO's or using hints or u "talking" TM about those concs. I, with interest of g such a TSO is idea that it actually does learn useful concs & gives them useful PC's

.40

(and useful conc pc's). Also that one can get useful (try, even tho all steps are not really "rand").

01. ~~Block~~ T. Let  $\alpha$  of  $\boxed{96.40}$  is a possl. "Bracketing" in my ideas about TSOs

Platform

03. SN  $\alpha \cap \gamma = \emptyset$ ; My calcns may have been wrong.

91.40 Interp Grammar we have  $f_A, f_B \in f_C$  w.  $C = AB$ :

Observed trees of  $H$  is  $f_A + f_C, f_B + f_C$ .  $AB$  is observed as  $f_C, f_C + f_A + f_B$ .

Given  $f_A + f_C, f_B + f_C, f_C + f_A + f_B$ ; to solve for  $f_A, f_B, f_C$

$$\begin{aligned} A+C &= \alpha & C+B &= \gamma \\ B+C &= \beta & \therefore C &= \gamma - AB \\ X+Z &= \alpha & Z+XY &= \gamma \\ Y+Z &= \beta & \therefore Z &= \gamma - XY \\ X+\gamma - XY &= \alpha & X - XY &= \alpha - \gamma \\ Y+\gamma - XY &= \beta & Y - XY &= \beta - \gamma \end{aligned}$$

$$\begin{aligned} A - AB &= \alpha - \gamma \\ B - AB &= \beta - \gamma \\ A(1-B) &= \alpha - \gamma \\ B(1-A) &= \beta - \gamma \end{aligned}$$

$$\frac{A}{1-A} \cdot \frac{1-B}{B} = \frac{\alpha - \gamma}{\beta - \gamma}$$

$$A \left(1 - \frac{B-\gamma}{1-A}\right) = \alpha - \gamma$$

$$A(1-A(\beta-\gamma)) = (\alpha-\gamma)(1-A)$$

Quadratic in A.

Block key is better Block Deleted w. key.  
In 2. previous para, it was checked by mod of w. each.

$$\begin{aligned} X(1-Y) &= \alpha - \gamma \\ Y(1-X) &= \beta - \gamma \end{aligned}$$

$$Y = \frac{\beta - \gamma}{1 - X} \quad ; \quad X(1 - \frac{\beta - \gamma}{1 - X}) = \alpha - \gamma \quad ; \quad X(1 - X - (\beta - \gamma)) = (\alpha - \gamma)(1 - X)$$

$$X^2 + X(-1 + \beta - \gamma - \alpha + \gamma) + \alpha - \gamma = 0$$

$$X^2 + X(\beta - \alpha - 1) + \alpha - \gamma = 0$$

$$X = \frac{1 + \alpha - \beta \pm \sqrt{(\beta - \alpha - 1)^2 - 4(\alpha - \gamma)}}{2}$$

woops! Additional eqn:  
 $X+Y+Z=1$   
 $A+B+C=1$ .  
Corpus: if there are  $> 2$  symbols, we can know  $A+B+C$ , but it would be  $< 1$ .

So this may be over constrained! 4 eqns, 3 vars!

$\alpha + \beta = 1$

What's exact value of  $\gamma$ ? (.06) It's the no. of times we observe  $AB$  divided by  $\frac{N}{N}$ .

(Also note:  $A \neq B$ ; since this analysis doesn't hold if we are looking for "1" occurrences, say,  $\frac{N}{N}$  as in 2.141 (on 5.64 part 2) we treat this as a special case.)

26. Say a grammar is  $A \rightarrow \epsilon$  symbols:  $A, B, C$  w. freqs  $X, Y, Z$ ;  $X+Y+Z=1$ .

We have this Bern corpus w.  $N$  symbols;  $\frac{N}{N}$  occurs  $\frac{N}{N}$  times.

We plug to  $C \rightarrow AB$ . This corpus now contains  $N + N = 2N$  symbols.

$N(X+Z)$  are  $A$ ;  $N(Y+Z)$  are  $B$ ;  $AB$  occurs  $\geq N$  times as a result of  $C \rightarrow AB$ .

$AB$  occurs  $N(XY)$  times as a result of random adjacency of  $X$  &  $Y$ .

So in the entire corpus of  $2N$  symbols,  $A$  occurs  $N(X+Z)$  times;  $B$  occurs  $N(Y+Z)$  times.

$AB$  occurs  $NZ + NXY$  times.

freq. of  $K$  is  $\frac{N(X+Z)}{N(1+Z)} = \frac{X+Z}{1+Z}$

" "  $B$  is  $\frac{Y+Z}{1+Z}$

" "  $A$  is  $\frac{Z+XY}{1+Z}$

Frequency of "freq" of  $AB$ : This is ratio of no. of occurrences of "AB"

in corpus, to no. of symbols ( $NX + NY + NZ = N(1+Z)$ ) in corpus

So given  $\alpha = \frac{x+z}{1+z}$   $\beta = \frac{y+z}{1+z}$   $\gamma = \frac{z+xy}{1+z}$  and:  $x+y+z=1$ .

So given 4 eqns. to find 3 vars! — seems to be some thing wrong!

This is certainly not a "deep problem", but it's annoying!  $\gamma > .28$  it seems to give .01!

perhaps  $\frac{y+z}{1+z} \Rightarrow$  a factor of  $\frac{x+z}{1+z}$  :

$$\frac{y+z}{1+z} \cdot \frac{x+z}{1+z} = \frac{x+y+z+z}{1+z} = 1$$

$$\frac{y+z}{1+z} = \frac{1-y}{2-x-y}$$

$$\frac{y+z}{1+z} = \frac{1-x}{2-x-y}$$

So that solves it:  $\alpha + \beta = 1$ .

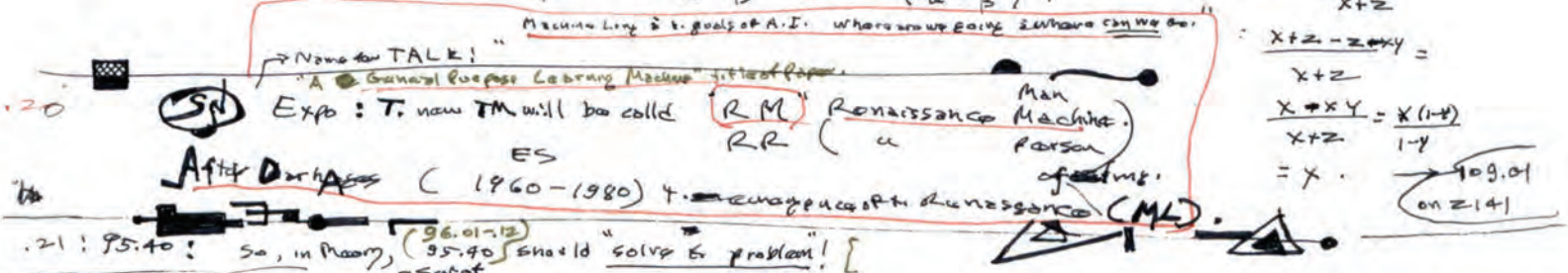
$$\gamma - \alpha = (.01) : \frac{xy-x}{1+z} = \frac{x(y-1)}{2-x-y}$$

$$\Rightarrow \frac{1-y}{1+z} = \alpha \therefore x = \frac{\alpha - \gamma}{\alpha}$$

$$\alpha = \frac{1-y}{1+z} \quad \alpha - \gamma = \frac{x-xy}{1+z} = \alpha \gamma (1-y) \therefore x = \frac{\alpha - \gamma}{\alpha} = 1 - \frac{\gamma}{\alpha}$$

So  $x = 1 - \frac{\gamma}{\alpha}$  ;  $y = 1 - \frac{\gamma}{\alpha}$  ;  $z = 1 - x - y =$

$$z = 1 - 1 + \frac{\gamma}{\alpha} - 1 + \frac{\gamma}{\alpha} = \frac{\gamma}{\alpha} + \frac{\gamma}{\alpha} - 1 = \gamma \left( \frac{1}{\alpha} + \frac{1}{\alpha} \right) - 1$$



20: After our (1960-1980) + emergence of the Renaissance

21: 95.40: So, in theory, (96.01-12) should "solve" the problem!

22: A possibl. remaining problem: Just how to get to the cpd's of h.c. starts to use for oz probs. How can 96.01-12 be applied to this problem? And, even if it could: how does one put new h.c. methods a/o modify the cpd of those h.c. methods wr. h.c. (even) problem?

It may be that we usually don't want to use LS for oz probs! (But) (using (x) PGA) is usually better & easier to implement.

There is one kind of oz problem that I normally don't do on which LS can be used. It's finding short codes for a corpus. A common way is to look for codes that reduce code length, reduce the corpus using this code, then repeat ad infinitum. I think this is what I use in all 3 TSO's that I "write".

It seems that the "LS can be" would not be a true optimum. — unless the h.c. method choice was conditional

on not only known [O<sub>i</sub>, G<sub>i</sub>] corpus, but on the process obtained generating the new O<sub>i</sub> a/o on the process of the "eval. trace". Note: pretty 95.36-38: what's this "oz" prob. & "prob. of"??

BUT: Most recently, I had the idea that one could go pretty far by using a fixed set of h.c. methods. That memory preservation method was a "heuristic" & was not very often done —



01:99:40! Could this be an applcn of the "Next candidate" problem. What to work on Next soln? Certainly "not exactly"!

But perhaps I can use the GCM, techniques of some how modify soln of WTWON.

T. WTWON prob. is  $\sim$  to  $\Delta Z$  in that the H.C. methods are "in parallel". T. WTWON soln is interesting in that ① it is optimum ② it does jump from one cond to another - back & forth - which would seem unreasonable.

perhaps  $\Delta Z$  soln. would be optimum for ~~some~~  $\Delta Z$  problem.

In WTWON, when we had several continuous distributions in  $\parallel$ , one kind of soln. was a time share which the share depends on the slope. In a  $\Delta Z$  way one could do time share out. H.C. methods, but it's unclear as to what "slope" to use!

07 SN In Normal  $\Delta Z$  soln; T. H.C. methods used are indep., so the set of trials made by each method is not available to the other H.C. methods! This is a tremendous waste! [ Tho, it may well be that most H.C. methods don't use much of the past corpus - they may use some of the most recent trials. Sim. Anneal uses only the last trial! ]

"Slope" might be  $\uparrow$  in  $G$  per unit time of "Best cand. Plus for" as done by ~~the~~

H.C. methods: (whatever that means!). Given a set of [Order] pairs, each H.C. method

should propose a "next trial", w.  $\Delta Z$  assoc. w. one trial: Hvr. there are an  $\infty$  of H.C.'s, so it's not practical to try them all.

Note:  $\Delta Z$  forg. may not apply to Sim. Anneal:  $\Delta Z$  trial depends on the "Temperature" at that time.

SN Another Approach to optm: Say the G of all  $\Delta Z$  probs are (measured): Then one can ask for a  $\Delta Z$  method that will give best approx  $\Delta G$  in available time.

Another possy: From the corpus of 99.99: one can estimate which  $\Delta Z$  set of H.C.'s will be best for a given problem. -  $\Delta Z$  Prob will vary as one continues on one or more of the H.C.'s.

The directly relevant subcorpus of 99.99 would be those in which the  $G_{ij}$  were max for the  $C_{ij}$  & H.C.  $\Delta Z$  prob; but doesn't this have difficulty, 07-08?

Could SGA be a "Not Bad" soln. to the general  $\Delta Z$  problem? T. Difficulty appears to be in an "easily comparable path" for the model. We can not restrict models to be of this class.

Wrt.  $\Delta Z$  problems: While finding short codes is a " $\Delta Z$ " problem, it does seem to be distinct from most, in that it seems "well-solved" by INV search! A poss. uncertainty in its very optimum.

That the heuristics needed, will derive from parts of the corpus other than those being coded; T. relevant parts will involve  $\Delta Z$  probs of the regular kind (i.e. not searches for short codes).

Which would be unfortunate, since it looks like I don't understand such probs very well!

32: (10) Well, the method of using many H.C.'s is how that share a common known [Set of kinds or known  $G$ ] is a new H.C. & should be included in the search as such.

Perhaps one could modify the proof of the " $\Delta Z$  search is optimum" lemma, so that the pc's of each H.C. were conditional upon the set of known [O or G] thus far.

It would seem that the H.C.'s would then become interdependent in their order of trials.

102.18  
105.01

A review of the (RM) System's its duties:

- (1) See 93.25 for a preliminary list of unsolved probs.
- (2) See 92.02 ft & 93.04 ft are listings of a few classes of probs ("inputs") for RM
- (3) I'd like General OZ probs to be solvable by the System, but Regular Level for OZ probs seems inefficient (99.36 to 100.40 at least), & I'm not sure how to get the PD Mat "guides". Level OZ soln. & make it part of the RM's "Grand CPD".
- (4) Just what form do these CPD's take? One possibility is the stack operator. — to get it, one just tries to find Good Stochastic stock operators! One that Maximally compresses corpus.

2nd & 3rd Party Short Codes  
ANALOG Probs  
N.B. This list does not include any OZ probs!

09 (SN) I could just start off w. a TM that could work mainly INV problems plus short code discovery

It would be a useful TM, yet it couldn't solve General OZ probs — which is to say Most problems. It could solve OZ problems approximately, using INV. Levels and SGA. There are probably other approx. methods it could use! (See 99.36 — 100.40 for some ideas, perhaps) Also, perhaps SGA could be implemented. → 102.22

17 It can usefully solve QA probs. T. "Analog problems" of 93.36 ft... seem to all be OZ probs (except for simplifying "corollum" see 93.36 102.20-2). T. "Corollum" idea of 93.36... of modeling an envelope for corpus of analog I/O pairs, seems to be a full-scale diff. diff. OZ problem. [No! — see 102.20-23!]

So perhaps draw up a more detailed dem of f.(09) done: (Crippled Restricted) RM (Crip RM) CRM  
This CRM can do lots of things & will give me a better idea as to what the main remaining problems are! Meanwhile, I can work on various poss. solns to f. OZ problems. Many types of Approach.

(SN) on ~~the~~ HCM's (attempt to map OZ prob into short code finding prob. & maybe OZ probs): OZ prob: Given (Black box) M(x)

To find  $x \Rightarrow M(x)$  in Max, (in time  $\uparrow$ ). (Consider as a coding problem: ~~each~~ string. For each Make a posn of M: so invert M by ~~using~~ a large positive real & integerize (roundoff)  $\rightarrow M'(x)$  A code for x is then  $M'(x)$  is an integer. For each value of  $M'(x)$  there will be a contn. no. (usually finite) of strings, X: ~~call it~~ card  $(M'(x))$  using  $M'(x)$  & ~~the~~ So, we can code (many) x's by  $[M'(x), n(x)]$  where  $n(x)$  is a no. betw. ~~0~~ &  $\infty$  card  $[M'(x)]$ . Short codes for x ~~may~~ involve small card  $[M'(x)]$ , and a full (maxim) can involve small card  $[M'(x)]$ , but other x can also have small card  $[M'(x)]$ . If we make time limit on M(x)... (This doesn't seem to help.)

I'm thinking of the Soln problem, in which a kind of optzn problem was expressible as a short code finding problem

(SN) It may be that true optimizer of TM behavior (including all or considerations) is essentially on OZ

probs, is not always a short code finding prob. An old problem that I (perhaps) never solved! Given a postset of  $\{P_i, Soln_i, C_i\}$  to induce an optimum CPD for all problems,  $\rightarrow$   $\{C_i\}$   $\subseteq$   $C$  of soln- of all post problems is minimal. Well, by assigning ~~each~~ CPD  $\subseteq$   $C$  known solns of all known problems, we get a  $\{C_i\}$   $\subseteq$   $C$ . This may correspond to "redundant solns" of problems. If we get a new prob.  $\equiv$  out of the old solns each way should

01 try the past solns first. Hvr. for a new novel prob; this will not work — ~~known~~  
 For 2 large class of problems, optimize CPD from to corpus [pre, solns] would be a duplicate —  
 Oran L such would be the best way to ~~be~~ continue.

I Print there are cases in which we want to bias: CPD in ways to reduce expected CC of solns, but may not be. Same as simply finding the CPD of the soln. in view of the problem —  
 that may somehow involve the CC's more directly. "Quick About" may be one such trick —  
 We find a way to quit a test early so as to ↓ CC of testing. Modifying CPD of solns. ~~can~~  
~~we~~ cannot implement this hour.

09 Interaction w. world (say a reinforcement problem) is an OZ problem

T. General TM2 problem (to optimize TM) is an OZ prob.

CIAP(0): we may be able to use fact (09) to construct a fairly intelligent TM, but  
 is not "conscious", it is unlikely to give ~~answers~~ assoc. w. "awareness of outside world".

18: 104.40 Another view of OZ probs: that TM might learn to take OZ probs is "change from representation"  
 100.70 So they are ~~easy~~ solvable by a standard set of heuristics

20 More on OZ probs: Given corpus [the probs, trial solns, G<sub>ij</sub>], TM can form a good CPD for  
 G of conc. solns & probs. Finding prob may not be easy, but ~~this~~ Peris would be a useful

22 function if evaln. of G is expensive; So Peris soln. could implement "T. Corbellum" problems

23 of 93.36; Hvr. using the "setebellum" for analyzing [I, 0] is usually a OZ problem; In a pinch, 10/19/99  
 → 105.01

24 We can use INV L to approximate a soln (or other approx. methods).

25: 101.17 So start out using the non OZ (corp. & probs) of 92.02-93.04: see how far we get w. the analysis

# probs (9405-9505) say, using Inverse to OZ part (or some other approx of OZ).  
 So basically, this is the idea: (conclusion): In earlier WTSQ's, I started out w. a set of primitive (objects, abs. concs.)  
 They all had some pc. I would use all poss. combos of them in PC order (assuming same cc for each trial). When I  
 found a conc that was a soln, I put it in storage. I would either use it <sup>then</sup> as a prim. conc. and give it = pc  
 to all concs within storage; or simply add it to the set of prim. concs (w. = pc), & then  
 go to next problem; loop to ... There may have been some differences in pc's dependent on freq. of use (essentially  
 Bernoulli probs).

NOW what I want to do is make the pc's of the concs. diff. for each problem/application of these concs...

... a new "conditional pc."  
 If a conc. appears for the first time, as the soln of a problem; then ~~the~~ next use of that conc. will  
 amount to "OSL"; which is normally treated in a special way. My impression is that in OSL, ~~the~~  
~~the~~ OSL a conc. occurs only once <sup>if</sup> ~~it~~ occurred. Then its pc of occurring again, is fairly low.  
 Hvr. next to conc. was the soln of a prob. may give it considerably more pc.  
Perhaps EN If we actually gave TM solns to problems (along w. the problems), it would take much  
 less time for TM to integrate the needed concs. (i.e. find relevant concs), than if it had to also find the  
 solns.

→ 103.01



0.1.98 TM

Unordered set of objects form: proof:

.01: In the proof of the conv. form (Sol 78 form 3) for unordered set I used 1, 3 symbols: — (t. cont. of this idea seems v. p. ... so stop). — Hvr, t. more General Form of 99. 33, 28 can use this 3 symbol proof: T. third (stop) symbol doesn't seem like an appropriate termination for the "partial objects" & "prefixes of infinite seqs".

.05: 102.40: In general, at each point in a trace, the "next move" is governed by cond. prob's involving T. problem being solved & f. previous trace within that problem. Later on, t. cond. probly will have more input conditions possl. — so it can draw info from other Domains.

Hvr, we want to start out w. uncond. pc's; then gradually add more conds, as size ↑ & ∞. ∴ as more conds become warranted. [To do this, we merely augment t. corpus w. t. new subcorpus of "conditions"]

[SN] 3. impl. deficiencies of current model of RM:

- 1) Ability to do OZ pubs, but only in an approx. way (non-optimums)
- 2) (Possibly a result of 1) Not all types of heurs can be implemented, since only those that modify t. cpd are possl. (eg. "Quick Abort" types of heurs are imposs. to implement)
- 3) If Lsh itself (either in T ← ZT or Time shared Mode), is non-optimum since it doesn't take advantage of interactions ("cross-info") betw. trials. — Ideally, this info should (at least) change to shape of PC curve.

.18: 96.30: One way to implement this "Value" idea is to give examples in t. form:

[V(10+3), 13] TM must then extrapolate to V( ) function to every string:

(Hvr, V will not be defined for most strings): For perhaps it should first learn some "type" ideas: that "10" is a "number" → Part

V( ) is a "number".

To make it easier, table f. corpus "V" & have examples be like 10+3, 13

(This bears a close resemblance to t. old corpus: 10+3 = 13, (!@)).

T. diffrnce is, that TM can know about t. comma; that it separates t. 2 parts of t. corpus into P: A

— So TM would know that it had to express t. string "13" as a funct. of t. string "10+3".

Having t. idea of "Recursion" could be useful — or have it as an available primitive function in t. "Functional lang.". In t. present case: One way to do it: ~~then~~

V = "look for t. section of t. string ~~in~~ in t. domain of V: apply it & rewrite t. string w. V's value"

Do V until it is no longer possl. ↑ this may be an adequate recursive defn.

i.e. V = G(V, A)

in execution, we can use special purpose functions to evaluate values & efficiently

Since I ideally, these functions should be very useful to express other imp. concs (funct.) as well.

Some may look at APL or "J" languages,

Another approach would be to start w. less clear prim. concepts, & use lots of corpus &

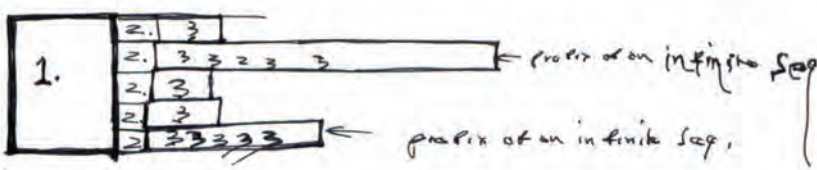
lots of CC:   
 ← direction of Jargon

CS Procs  
University  
Pittsburgh  
1878

0.4.98 TH

- 01: 99.35 [SN] On t. functional 9925, 28! Wallace may actually want to do this as a port code!
- 02 1. self def. code common to all objects & partial objects & / or seqs.
- 03 2. for each  $\alpha \rightarrow$  (objects of  $\alpha$ ) in 1., a self def. code giving a P.D. ~~by zero~~ <sup>(05, 27 - use and this part should be zero)</sup>
- 0398 3. for " " " " " " ; a self def. code giving to prod of complete object, <sup>or 2 (1/4) self def. code for partial object.</sup> ? 20

115A  
0498  
215A  
4A 13  
545A  
615A 2mg



09 [SN] On "Levels" for a growing corpus:  $\alpha$  idea I had was that once you had  $\alpha$  start w. a short corpus, then using Levels, find some conc. ~~data~~ <sup>sub</sup>  $\alpha$  <sup>data</sup>  $\alpha$  and compress it. Then using  $\alpha$  old primitives as new defns, use Levels to code an augmented corpus, & find new defns. loop to  $\alpha$

Some problems here: how much  $\alpha$  is to be used for each corpus augmentation?  
 [Note that in general, one would try to recode  $\alpha$  as per corpus, not just  $\alpha$  augmentation. — T. larger corpus may have enuf  $\alpha$  to find some new low level data. to help recode ~~earlier~~ earlier parts of corpus]

[SN] On functional langs: I will want various simple <sup>prime</sup> functions:  
 Say  $x, y, \dots$  etc., integers, maybe sin, cos, exp, ln, etc.  
 Then functions! To implement:  $x, y, z$  arguments  
 Do until  $G(x, y, z) > 0$   
 $(x, y, z) = F(x, y, z)$   
 $F$  has been applied a certain no of times (  $x, y, z$  can be a constant;  $F(x, y, z)$  may not always change  $x, y, z$  )  
 So this function is defined by  $F$ , and "until" a t.  $z$  func  $G(x, y, z)$ ,  $F(x, y, z)$  and  $>$

Do Until (cond) 20-  
~~input~~  
 inputs  $x, y, z \Rightarrow$   
 $F(x, y, z)$   
 loop

To define  $n!$ : ~~input~~ ~~output~~ ~~n=0~~ Input  $(n)$ : output  $n!$   
 do until ~~n=0~~  $y=n=0$  ~~stop: y=0, x=1~~  
 $x=y+1$   
 $x=x*y$

So f. func  $n!$  is defined by:  
 initial  $y, x, n$   $y=0, x=1, n$   
 $y_{n+1} = f_1(x_n, y_n)$   $y_{n+1} = y_n + 1$   
 $x_{n+1} = f_2(x_n, y_n)$   $x_{n+1} = x_n * y_{n+1}$   
 $G(x, y) = 0$  is condition for ~~end~~ end of loop.  
 $y_n = n = 0$

This is one way to define  $n!$ : it seems to use many prims  
 Less complex defn:  $0! = 1$   
 $n! = x * (x-1)!$   
 standard recursive form.  
 we can have a compiler that easily  $n!$  thus this defn into 1. more "computable" defn. of  $n!$

0.4.98 TM

OZ probs

Leads Good!

Scalar

string determining an open occurrence

Kind Corpus Problem (24)

01: 102.24 100.40 (SN) Interesting Q! Given a corpus  $[O_i, G(O_i)]$  a set of a data set in which we want to find  $O_j$ .  
 $\rightarrow G(O_j)$  is as large as poss. w. t. in CB. (Minimization) (is it OZ problem).

Is it meaningful to ask for  $\text{rank } P(O_j, CB) - \text{t. probly that } O_j \text{ is t. value for which } G(O_j) \text{ is max: [Somehow CB tries in]}$

Well,  $[O_i, G(O_i)]$  implies a certain  $P(O_j, G)$  p.d. for all  $O_j, G$  values.

07 We can ask for t. probly that any one  $O_j$  has max  $G$ . This would seem to be a clearly defined function if  $P(O_j, G)$ 's for t. various  $O_j$ 's are indep. [if they are not indep.]

09 it may still be defined but even harder to compute function of  $O_j$ .  
 If so, then this could perhaps tell how OZ problems can be added to t. corpus of 104.01-08 (i.e. 99.23, 28).

Even if t. p.d. of 07-09 is very hard to compute, it is <sup>perhaps</sup> meaningful to ask for

t. best estimate we can get in 10 sec, say.

If th. arg. discn. is correct, this will help a lot in t. ~~particular~~ discn. of OZ probs in

100.07 ff. 100.07-10 suggests strongly that I don't really understand t. OZ problem;

that I don't really know what t.c.p.d. for t. H.C.M. is. I have no data for it.

clear how one

105.01 ff suggests a defn., tho it's not at all clear how one goes about approximating it!

20  $\rightarrow$  Actually, t. discn. of 01 ff is not exactly rite for OZ: for OZ we have several H.C.M.'s  $H_i$ . We have  $[H_i, O_j, G(O_j)]$ :  $[O_j]$  is a subset of  $O_j \geq G_i$  by  $H_i$  in  $CB = C$ . Doubtedly 100.07  $[O_j]$  is not clearly defined: but anyway, from that corpus we're supposed to get a c.p.d. for  $H_i$  getting to base  $G$  in  $CB = C$ .

23 All of this is w.r.t. problem  $P_{H_i}$ : so reason our cond. probs that  $H_i$  will be base for  $P_{H_i}$  in  $CB = C$ . **108.09**

"IR" probs

or "IR" prob: 93.12, 33

24 104.18 (SN) Forti. Genza of 99.23-24! I need to analog of  $\frac{P_i}{P} > 2^{-K(n)}$ .  
 To start off, use a generator of t. corpus w. a finite den:  $\frac{P_i}{P} > 2^{-K}$ . **See 106.23**

This amounts to a finite den of parts 1) & 2) of ~~104.02-03~~ 104.02-03.

27 If we only have finite objects in t. corpus, this suggests that base of 2 (104.03) should be zero!  
 28 — out. other hand, if we have elements of t. corpus that are prefixes  
 29 of infinite seqs., then it seems base of 104.03 should be  $> 0$ !

For finite b object corpus: clearly 104.03 must be zero, if th. stack lang. den, is finite (unusual, but poss.)

31 For corpus that is prefixes of 2 inf seqs., then have base of 104.03 should be  $> 0$ !  
 Going back to 99.23-24! T. way we go to use base for inference for a k extrapolation;

Say we have t. prefix of some object (finite or infinite). We want t. vol of PC of various contents of it. For each content (which is a prefix of some object), we want t. total a pair of all objects that (plus t. rest of t. corpus) that have that prefix.

Consider a corpus: 1) corpus contains several finite objects & one prefix of objects, that (also consider only finite den for corpus.) We want the base of of all corpus that have those prefixes. (These are "cylinders").

may be finite or infinite.

This brings up an interesting Mathematical object: say the corpus consists of prefixes of two in P. seqs. We want to product of the PCs of 2 "cylinders" — This is supposed to be some kind of sum of many pc of many <sup>or continuous ∞ of</sup> doubly infinite seqs — many of which have "zero" pc. So we want a measure or semi-measure over some space.

I think the idea is that 99.23, .28 <sup>is</sup> correct — but what about 105.28 - .29 <sup>(.29)</sup> 105.31 seems to be for "Bis Q".

In "para" stock lines (w. ass ≡ finite strings), it is poss. that 4. gen error pauses for 2 while both generating certain of v. ~~strings~~ strings! Could these pauses be interpreted as  $\geq 2$  (4.03 ~~XXXXXX~~) ( $\equiv 2$ ).

An analog of 105.31: the "1." component: common concs. and delay of sin, cos, tan, etc. etc. "2." components are the individual parameters in terms of sin... etc, plus various constants (params).

So sequence a):  $X_{n+1} = \sin \sum_{i=0}^k a_i X_{n-i} + noise_a$   
 " b):  $Y_{n+1} = \sin \sum_{i=0}^k b_i Y_{n-i} + noise_b$

instead of sin, as a funct. w. clear vector  $\rightarrow$  several params.   
 Noise a & b are uncorrelated.

One of the descr. things 2 seqs as a single vector sequence, w. vector components,  $X_i, Y_i$ . In the vector representations, the  $X$  &  $Y$  sequences need not be of same lengths. — The way this done: the pc of the partial vector seq. is the measure of all vector sequences that have both the "x" prefix & the "y" prefix [the 2 prefixes being of different lengths]

This vector representation can clearly be extended to any finite no. of sequences. Certain of the "vector components" could represent bits of finite <sup>(strings)</sup> strings — So this will be a "Unified soln." to the "Mixed Corpus" problem of 96.01

23: 105.24 The no. "n" in  $\frac{P}{P} > 2^{-k(n)}$  (105.24) is simply the no. of bits in the known Corpus.

24: .04: The way this is done: We have 2 sets of code strings of int. length. To get the product of

25 Prefix measures; Take the measure of the first set of strings, in which each of the strings

26 is given with ~~the~~ = ~~the~~ measure of second set of strings.

Equivalently, if  $[S_1^i]$  &  $[S_2^j]$  are 2 sets of strings, we want to measure of  $[S_1^i, S_2^j]$ . and I think it's 24-26.

A way to Criticize the "vector soln." of 17 ff: Say the corpus is only a finite set of strings. — Does the soln. correspond to our natural "rate reasonable" soln to Y. "Unordered sets of objects" problem?

33  $\rightarrow$  1) What kind of sub-corpus is  $OZ$  data? Even if we only have one or more  $OZ$  parts in the corpus — this is still an "open problem".

34 2) In 104.05-08 (in figure) "1" & "2" can be combined via method of total info. — so the alternative vector representation of (17) poses no difficulty — it's only an "Alternative Representation" which is always good!

The vector representation does combine 1 & 2 — so we only need 1 2 input machine.

But, I do want to write up a detailed dec. of the (proposed) soln. to the "Mixed Corpus" problem.

If it works, the next big break problem is (33) (OZ corpus).

Mixed corpus problem

A diff. m. 106.34ff: T. initial Machine "1" is supposed to start on t. random seq-corpus: each with its own new, uncond. p.c. For inf. sequs., t. p.c. of t. entire seq. is " $\phi$ ", but for finite objects, it is  $> 0$ .

Perhaps for each inf. seq., we could have a "p.c. of starting" — as i.

To work on this problem: Consider just out t. system works! T. diff. Pres. G. gives for predns.

Int. case of pure uncond. objects for corpus: Consider a stoch. CFG as a P.D. desc! We start w. desc. of t. grammar, then we have various probabilistic choices of non-terminals:

Then, at last step, (perhaps) t. terminals (final ~~obj~~ symbols) are generated.

It is not t. "sequencial" generation of most mt. sequs: it seems more like Cover's

"Extensio<sup>Complexity</sup> (probab.)" — not exactly: for CFGs, — say one wanted the P.D. of various poss. continuat<sup>with</sup> of a part of object"; One can, normally (I think) devise a CFG (or set of CFG's) to generate t. "rest of object".

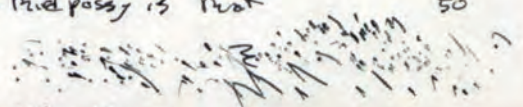
Well, ~~consider~~ say t. corpus consisted of  $A \& B$  (2 finite objects)  $\& C \& D$  (2 inf. sequs). A possy is that C & D have t. same "first part",

& that they are simply different stochastic variations from t. same ensemble.

A more possy is that they have <sup>completely</sup> different "first parts". A third possy is that they have different "parts", but some "shared" part's.

Arr Boston  
12/15  
So a Hierarchy  
Lundtime  
16 PM Boston  
time.

to  
8  
50



20

I think I have a clear idea as to what "INFO" is required for "Mixed" corpus. — So t. Machinery of Generation of t. corpus should be (trivial) easy.

T. vector formulation of 106.17 seems to deal well w. t. diffy of infinite sequs being infinite (so one cont. do one seq. after t. other")

23

Another View of "Mixed corpus": we have a machine that generates t. corpus in t. following way: Random input, unidirectional I/O tapes, 1 wrh. tape! One person inputs. As soon as one of t. corpus objects has been generated (it can be a finite object or prefix of a finite or inf. object)

27

We reset t. corpus  $\&$  loop to  $\alpha$ . Under these circumstances, there is a certain probability that our corpus will have been generated.  $\alpha$  is described as "failure".

23-27 is wrong. We abort t. trial as soon as its output is not consist. w. any of t. corpus objects. A trial, how, may get into an infinite (but <sup>unrecoverable</sup> inscrutable loop) — in which case we say, t. corpus was not generated! (t. total probab. of all objects in t. corpus is certainly  $< 1$ !).

30

If we generate an object twice, but it appears only once in t. corpus — that trial is "failure". 23-30 may be O.K. ! Is it consi. w. t. vector view of 106.17? (t. vector view)

may be equant. to an ordered set of objects.

23-30 is for a BAE corpus.

23-30 suggest a desc. for t.  $K(n)$  of 106.23, 105.24: We have a particular T.M.C., M and it generates t. corpus (using notes 23-30) w. probab. P. Say, using an vec,  $M'$ , t. probab. of generating t. corpus is  $P'$ . Then  $\frac{P'}{P} < 2^{-n}$  !  $K$  being t. vector desc. of M w.r.t.  $M'$ . How to desc. w.  $K(n)$  is not yet clear. I should first get t. continuous

params straight, before working on  $K(n)$ : Just how does t. dimensionality enter?

02: 0.03 0.09  
T. conv. Coercary

.01: CB like a formalism for MIXD. Corpus that enabled to solve Thm 3 (conv.) to be used to find convergence Rate.

Use 3 symbol formalism. One stop (actually reset) symbol. Par ind. Sages, was still have v. useful

Symbol: Any extension of t. do not prefix, followed by v. reset symbol, is regarded as a loop

.04 code for t. corpus. This is  $\equiv$  Cover's "Extension Probability".

.05 SN "Extensy" doesn't allow finite codes of infinite spaces (like 1000 has a sample finite code)

.06 It may give  $\approx$  same proby values as ALP. In sol 78, I had to impress. That Page mita be  $\approx$ , including Conv. rate.

.07 So: .01-.04 may solve t. problem but w. diffy of .05-.06 (maybe not so imp.)  $\rightarrow$  110.01

.09: 105.23 So next B<sub>j</sub> problem is 02 probs.

Say we've worked a lot of 02 probs using 02 Lsrch, so we have a lot of data on various HCM's

11 Probs, w. CC's  $\in$  G's. We have a corpus [Pmb<sub>i</sub>, HCM<sub>j</sub>, CC<sub>j</sub>, G<sub>j</sub>]

12  $\Rightarrow$  Next, given a new problem Prob<sub>0</sub>  $\in$  a CB = CC<sub>0</sub>, to find p.c. best HCM<sub>j</sub> gives Max G in CC = CC<sub>0</sub>.

Actually, more not so much interested in t. G values but in t. ordering of various HCM's

For previous dozen of 02 ~~102.24~~ ~~100.24~~ ~~100.24~~

Seq 99.30 - 100.40; 102.18 - 24; 105.01 - 23

A Series

Criticisms regular 02 Lsrch is 100.07 (doesn't all know [0; G(0<sub>i</sub>)] in  $\approx$ )

This can be fixed in a non-A.P. way BY 100.32 (100.07 diffy) 00.07  
No correspondence betw. HCM's.

20 But even w.o. Prot Criticism! I still don't know how to get P (HCM<sub>j</sub> gives max G on Prob<sub>0</sub> w. CB = CC<sub>0</sub>)

105.01-23 looks like a series step toward soln. N.V. rate 105.20

105.01-23 suppose start from t. corpus of (.11), we can obtain (.12) - Actually we don't need t. P.R. obtainable for CB=CC<sub>0</sub>, but a much less good one, obtained w. CB = CC<sub>2</sub> - The what CC<sub>2</sub> to use is unclear: T. problem of obtaining (.12) looks like a TM<sub>2</sub> problem.

What (.12) is: we want a stoch operator w. Prob<sub>0</sub> input (02 Problems, CB = CC<sub>j</sub>, HCM<sub>k</sub>) output proby that HCM<sub>k</sub> will give best G w. CB = CC<sub>j</sub> for problem; (or: using only first 2 inputs, output is a P.D. on HCM's) or a Monte Carlo output w. proper probys for each HCM<sub>j</sub> name, J.

$\rightarrow$  A main idea is that w. a reasonably large corpus, we should be able to usually tell which HCM to use w. given problem  $\in$  CB.

If the really did face induction, T. method should be able to extrapolate new HCM's to fit a new problem on t. basis of t. corpus of .11. (Also now 100.07 has suggested soln. 100.32)

$\rightarrow$  0798 T. 02 Lsrch proby. is beginning to look like a p. problem in SM!  
T. diffy of 100.07 is effectively solved by 100.32. So 02 Lsrch may still be "optimal".

01: 98.20: So anyway we have to consider all poss. parsings for  $t$ : data  $C=AB$ ;  
20.03 Also not only  $t$ : most likely no. of "C" occurrences.

See 141.01 for a good way to understand this.

If we use  $C=AB$  for prediction, it may not be necessary to <sup>explicitly</sup> consider all possibl. parsings. (except to know how many parsings, so we can suitable  $\rho$  of  $P$  &  $R$  set of codes).

04: 91.40 For loop closure (in Z141  $\rightarrow$  CFG): In my idea about Functional Lays,

I just did combinations of concs. that were <sup>lower</sup> ~~higher~~ in  $t$ : tree. Like

Here, to get recursion, some of  $t$ : input objects are functionals that create recursive functs. from "non-recursive" & recursive functs.



"Loop operators" are functionals of  $P$  &  $R$  sort. — so perhaps it would be possl.

to discover loops (for  $\rho$  & CFG discy) using Pars Method. 1

10 In fact it had better be possl. to disc loops  $P$ 's way! Because my idea of using ~~the~~ Universal functional lays for concs discovery, is critically dependent on this idea!

After while ~~no~~ ~~may~~ may be true, we may need hours to make it practical:

Such hours could do statistical tests ~~to~~ whose results would suggest which loops to try. These tests would give cond. probs to  $t$ : various loop trials. (Nothing Magic... just usual way of doing it!).

Wolff may have some nice heuristic methods for new data on  $t$ : loop closure.

$P$ : Wolff's  $t$ : soln. of Z141 problem! So far, my impression is that  $t$ : usual way I expect  $t$ : find regys in data is by "synthesis" so to speak. I find simple regularities & try combs. of them &  $t$ : primitives, in attempts to find more complex regularities.

Does this obviate ~~the~~  $t$ : Z141 (cont) problem? Superficially - NO!

So we are looking for solns to a problem viz SGA: We have a lot of cards w. assoc  $G$ .

~~the~~ One way is to make grammars for sets of cards at various  $G$  levels.

The Z141 grammar is  $t$ : simplest kind. It falls us to look for combinations (Subtrees) ~~of~~ ( $\equiv$  subfunctions) that occur unusuall; often in  $t$ : corpus.

$T$ : analysis of induction over unordered sets, is  $t$ : Mixed corpus problem, enable one to properly assign parts of data to these corpus. e.g. if  $A$  &  $B$  are functions  $P$  &  $AB$

is a combin. of them, then we look for  $t$ : frags of  $A$ ,  $B$  &  $AB$  in  $t$ : entire corpus of cards, withing  $t$ : desired  $G$  range. The frags of  $G$  can overlap in order to  $\uparrow$  SSZ for

$\rho$  each corpus. When we get a Grammar for ~~each~~ each corpus (for each  $G$  range)

We can then try to make a general Grammar for all values of  $G$ .

We may note that  $t$ : probab of certain functs  $\uparrow$  as  $G \uparrow$  — so we can extrapolate to high  $G$ .

01:108.07: A more detailed design of 108.01-07 w. ideas of how proof will work.

We start w. an uncond. p.d.; It induces a p.d. on all finite objects & all infinite strings.

In the latter case, it gives a measure on the point set that represents finite strings.

Any pts on 0, 1 interval represents an inf. string (or 2 inf. strings for terminating binary nos)

5- Each such measure can be represented <sup>measure</sup> exactly by a ~~finite~~ Tmc

[ SN T. proof of .05 might be easily reconstructed, using the ideas in Arithmetic Coding. ]  
It is probably identical to the proof I referenced (By Cover & L...) in Sol 78.

08 The way to ~~work~~ works: We put ~~some~~ random input into the Tmc (unidir, I/O, Bidirac, work tape)

T. proby that it will output a particular finite object & stop, is the pc of that object

" " " " " " " " & string w. a prefix  $\equiv$  part of one of the prefixes of a u. m. f.

10 string that is in the corpus; is the prob of that prefix.

Since the corpus elements are indep, we reset the machine after each trial (successful output).

[ We have a time limit  $T$ , to tell if a trial ~~will~~ give any thing in the corpus as output <sup>or not</sup>. Since  
the Tmc can take a by time w. its calcus. .... T. term is for  $T \rightarrow \infty$  ]

18 Because of this independence, the pc of the corpus  $\equiv$  product of pc's of the objects in the corpus.

So, every Tmc induces a pc on the corpus (possibly a pc. of  $\phi$ ).

20 Th. ~~the~~ resultant pc of the corpus is  $\sum_i P_{M_i} \cdot P_i$  : Here  $M_i$  is the ~~best~~ possible

Tmc. ;  $P_{M_i}$  is the pc of  $M_i$  w.r.t. the reference Umc. ;  $P_i$  is the proby induced on

the corpus by  $M_i$ .

From the above, it seems clear that if  $P_{M_i}$  is the pc of the Tmc that produced the corpus (or, more exactly, the sum of the pc's of all Tmcs that would produce the corpus with proper probys),

Then for a particular reference Umc,  $P_{\text{corpus}} (\text{from } .20) \geq P_{M_i} (\text{true pc of corpus})$

[ T. way we implement .20: let  $\tilde{i}$  be a binary string that represents the integer,  $i$ .

So we feed  $\tilde{i}$  into our Umc. This Umc is now by definition  $M_i$ . (If  $\tilde{i}$  results in some ~~output~~

output, then if all objects in the corpus have that output as prefix, then it's possible that the

corpus has pc  $> \phi$ : If 1 or more objects in the corpus doesn't have that prefix,

then pc of corpus w.r.t  $M_i$  is  $\phi$ ). We now use .08-.18 to evaluate  $P_i$  (or .20).

34  $P_{M_i} = 2^{-|\tilde{i}|}$

Consider (.08-.18) ~~unproven~~ if T. Tmc is Umc: ~~is~~ this distribn to discrete or continuous universal distribn? (or both or neither?) It seems that it is both:

Th. cont. & discrete distribns are both semi-measures so  $\leq$  at each is  $\leq 1$ ; sum of both is  $\leq 1$ . Actually, there are (at least) 2 kinds of Umc's: One has a stop state, one has

no stop state: Either mechanically ok. for .08-.18; If it has no stop state it produces ~~no output after T time~~ one of the finite objects, then has no output for  $> T$  time, we

consider it as having produced that object & stopped.

I'm not sure if we can get rid of the need for a  $\Sigma$  symbol alphabet for finite objects.

A toy: T. prob of 0 or 1 after the output object is zero. So if  $\alpha$  is a finite object (binary)  $\rightarrow$



01! This means that in 110.08-.18 we report (E given  $p_c = p$ ) to all codes that give  $\phi$  or 1 after ~~printing~~.

02 25 output.  $\rightarrow$  (20)

03 Idea for a somewhat simpler approach to ~~max~~ ~~monot.~~ proof poss. w.o. 3 symbols:

Since generation of t. objects is independent, we can treat them via Sol 78T3, individually.

Consider one object in t. corpus. Our system assigns  $p_c = P_i$  to it, t. probab. is  $P_i$ .  
 $\frac{P_i}{p_i} = \alpha_i$ ;  $\alpha_i = e^{-k_i}$  so EZ error for object  $i \leq \frac{k_i}{2}$ . This is via Sol 78T3 Corollary (2nd part of them).  
See 113.00 for Sol 78C3

E error for all objects  $< \frac{1}{2} \sum k_i$ ; and we know a bound for  $\sum k_i$ : that would prove t. thm!  
since  $\sum k_i = -\ln \prod P_i$  From 110.20-.34

10 Some troubles! 1 I'm not comfortable about P-Error "expected values" P-Error seem to be over different ensembles. 2 I don't know if Sol 78T3 can be applied to finite objects.  $\rightarrow$  It does

12 Sol 78C3 (Corollary of T3) may apply, just as long as  $\frac{P_i}{p_i} < \alpha_i$  for all bits in t. object. See 114

13 135 3 78C3: for some  $i$   $\frac{P_i}{p_i} > 1$  so  $\alpha_i = e^{-k_i}$   $k_i < 0$  so  $\sum k_i < 0$ !! which is impossible. Look into this (9/22/01) See 113.00 for Sol 78C3

14 00.12 We can apply 78C3 to finite objects: ~~but then applies~~ If t. object is a string of length  $m$ , then it applies directly to first  $m$  bits of t. object, w.o. modification. For bits past to  $m+n$ , t. additional error is always exactly zero, so t. term continues to hold exactly as before.

17 I think (03-.10) holds O.K. [9/16/01 - what about 1? (.11-.12)? See 01ID381.00 for some discuss. 9/22/01 - RE 3 - still troublesome!]

18 So: T. proof of t. Mixed corpus Thm. (as of now) is 110.08-.34; 111.03-.17

20 (02) On t. discrete v.s. continuous distributions: Universal Universal

1 t. Discrete distribn: t. p.c. of string  $X$  is  $\sum_i 2^{-l(P_i(X))}$   $P_i(X)$  is i-th perm that inputs t. reference machine  $M$  has  $X$  as output, then loops.

2 Continuous distribn: finite string,  $X$ ,  $\sum_i 2^{-l(P_i(X))}$   $P_i(X)$  is  $i$ -th perm that inputs  $X$  as input and has prefix  $X$ . The output may or may not stop after  $X$  is printed. [This is my defn, using minimal codes]

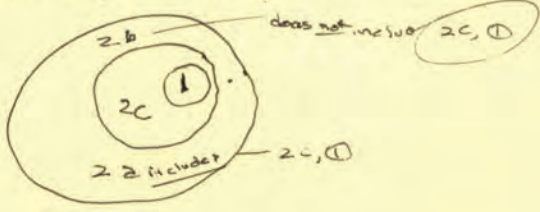
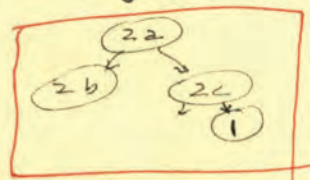
2b The measure of ~~all~~ codes of ~~all~~ finite strings that have  $X$  as prefix. This excludes codes of all finite strings. [This is Li-Vitanyi's defn. Maybe use UMC that has no "stop" state! But check on this.]

I think 1 + 2b = 2a; 2b includes only codes of infinite continuations of  $X$ .  
NO! 2c + 2b = 2a

2c Same as 1 but  $P_i(X)$  is i-th perm that has  $X$  as output a string w.  $X$  as prefix, & eventually stops.

1 is, of course, part of 2c; 2c is correct extension complexity

37 (18) To clarify my own ideas on this, draw up a proof of t. form - for unordered finite object corpus



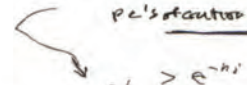
240 111 1/2 .00  $\rightarrow$  115.01

9/23/01

May be not!  $P_i' \leq P_i'$  no other

(11)

pe's stationary corpus up for Z'



$$P_i' > e^{-h_i}$$

$$\frac{P_i'}{P_i} > 1 \quad ; \quad h_i < 0 \quad - \text{so}$$

.00: 98TM/11.135: Ont.  $\sum (er)^t$  bnd for BAC Inductm; objection ③ for some  $i$ ,  $\frac{P_i'}{P_i} > 1$  ;  $h_i < 0$  - so

$E \sum (er)^t < 0$  ~~if this~~ if this is true then  $E \sum (er)^t < 0$  is also true, - which implies  $E \sum (er)^t = 0$ .

- No error! I'll have to look at proof of 578C3 to see what's wrong!

It is certainly not immediately apparent! Perhaps the Corollary is false!

See 9/19/02!

8/19/02 Actually  $E \sum \left( P_i \ln \frac{P_i}{P_i'} \right)$  is what is of interest & it is always  $\geq 0$  ( $\equiv$  Gibbs Form).

or something like this:  $\rightarrow$  Most likely 1/2 by

$$\left. \begin{array}{l} \text{If } \sum P_i = 1, \quad P_i \geq 0 \\ \sum P_i' = 1, \quad P_i' \geq 0 \end{array} \right\} \text{ then } \sum_{i=1}^K P_i \ln \frac{P_i}{P_i'} \geq 0 \quad \left. \vphantom{\sum_{i=1}^K} \right\} \text{ is always } \geq 0 \text{ } \equiv \text{ Gibbs Form}$$

So the KL distance is always  $\geq 0$ .

This is a radix  $\geq 2$  system



Let's look at the "indep" sites problem: say all  $X_i = 1$ . If we consider  $n$  sites, all possible configs we get a Binomial d.f. (?)  $\approx$  Gauss.  $\Rightarrow$  about  $\Sigma X_i = \frac{n}{2}$ ,  $\mu$  or  $\pm \sqrt{n}$  or  $n \approx 2 \pm n$ .

Now, if we take  $\left( \ln(\Sigma X_i) \right)$  it will be fairly linear  $\pm 1$  at end of  $n$  is large, or  $\exp(\Sigma X_i)$ .

$\exp(\Sigma X_i)$  will not be, but is unlikely that our Gauss would be that non-linear!

Probably  $\exp(K \Sigma X_i)$  w. small or  $\exp(\frac{\Sigma X_i}{n})$  is more likely.

Anyway, I guess it's unlikely that our Gauss will be ~~that~~ too linear and so that's a bits of 112.35-36 error of much import. — BUT I'm Quite Uneasy About This Argument!

Perhaps a better way to look at 7. Hie. problem in Genetic Genotypes!

That we have vector each site is a vector component, & each component has Allele values. We want quick & easy approx. to fit Genetic natural Gauss.

Taking pairs of sites into account is way to using a Quadratic Approx of G. Gauss.  $\Sigma$  Quadratic form or  $\exp(\text{Quad form})$ . Since the vector components are not continuous objects, but have discrete (Allele) values, a Quad form is not easy to define.

One way of doing this approx of G. Gauss is suggested by ~~112.21~~ 112.21 ff: We start by assigning a  $\Delta G$  to each allele at each site as 112.21-30! We then examine pairs of Alleles at diff sites & see if they don't add linearly; if not, we note the  $\Delta$  of  $G$  in the joint  $G$  v.s. the sum of their  $G$ 's. After finding pairs of alleles that have by  $\Delta G$ , we may come across other Alleles of good discovered pairs in hopes of finding even more non-linearity of K. Gauss.

Note that in 12.35-36, if one did do hill climbing in a simple max slope greedy way, one would find the peak — in either 12.35 or 12.36. In 12.35 using 14 ff's we would define larger & larger n-tups; but the result would be the same. In 12.35, we would not define 2-tups or any "tups" but we'd end up w. same results w/ peak.

1.35P  
3  
CPU 2 MMX

SN on T. Building Block ~~model~~ model/hyper " " for GA: In my "Progress in A.I."

1965, I had the idea that if one had crossover and breaking/reordering of genes, one could have "chunking"! My impression is that (except for ~~keza~~) they just do crossover, so very little "chunking" occurs! They may use ~~other~~ other models of "crossover" so they do get "chunking".

Many Chromosomes is another method of "chunking": do chromosomes ever get mixed? (e.g. front  $\frac{1}{2}$  of 1 chrom. & back  $\frac{1}{2}$  of another).

N.20.98 Much has been developed since 01398! Soln. of mixed corpus problem " " where to get P.D. for OZ probs: See 130.01-40 for Review

Anyway, in review (130.01-40): I can view GA probs in 2 ways (at least!)

- 1) The prob of GA is an OZ prob, & I know how best to solve OZ probs (130.01-40)
  - 2) As part of 37 we have this corpus of (rand,  $G_j$ ) pairs: we make a model for the data set (a stack Grammar), and we try to pick a cand. w. high expected  $G_j$ .
- This is essentially what SGA does, but SGA also tries to get models in which choosing cand. for highest expected  $G_j$  is easy to do.  $\rightarrow$  132.01

On use of General univ. machine (no restrictions on I/O):

Data of ALP is unclear: Some possible data:

- 1) "Output defined to be when machine stops.
- 2) "Input is all input tape squares & end before stopping."

Now output data: we want to know prob of  $s_a$  v.s.  $s_b$  ( $s$  is long,  $a, b$  short)  
~~on input outputs "sa"~~ on input outputs "sa" if it ever writes (at least) "sa" on its output tape.  
 Well what if it writes  $s_b$ , erases it & writes it again? — say it writes "sa" & erases it then writes "sb"? Also, it could write an  $\infty$  of strings ~~there~~ — so prob of non-convergence is very high.

There are several Models for induction suggested. (See

~~See~~ See, My letter to Levin for analysis of some.

A new one: Umc w. unidirectional <sup>input</sup> ~~output~~. Output <sup>tape</sup> is bidirectional, but it can't erase.

A certain input results in an acceptable code if ~~the~~ <sup>as soon as</sup> output produces a string w.  $\epsilon$  corpus as prefix. T. output tape can have blanks in it.

11.40

01:11:40: A critical Q is the defin. of t. p.d. being approximated by ~~the~~ ALP.

We have as "output" a set of finite objects; some of which are partial objects; (prefixes of finite objects/<sup>or</sup> infinite strings).

Several ways to define t. p.d.: We can describe it as assoc. w. some UIC machine in all cases.

~~For any~~ For any ~~object~~ object in corpus: w. random input, what's probab. that

a) T. output will ~~be~~ be ~~the~~ <sup>the</sup> first corpus element to stop.

b) T. corpus will be a prefix of t. output string. least restrictive.

c) " " " " proper prefix of t. output (i.e. t. output is at least 1 bit longer than t. corpus element).

d) for finite complete objects, t. output stops after every such object; for prefixes of finite ~~complete~~ objects, t. output has that prefix but eventually stops for prefixes of infinite objects, t. output has that prefix, but never stops printing.

it may or may not eventually stop printing.

Consider various sub-corpus cases: in all cases, output has that sub-corpus as prefix; corpus: a) prefix of some infinite string: straight

1) ~~output~~ output can be any finite or inf. string. (0.08) (Sum of 2) and 3)

2) output must eventually stop (~~Excursion Complexity~~)

3) output must be an infinite string ("cylinder" probability)

Sub-corpus is / complete finite object

output has that prefix <sup>in all cases</sup>

1) must stop <sup>in all cases</sup> after printing that object.

2) " " eventually after printing that object.

Sub-corpus is prefix of finite object

output must have that prefix

1) must eventually stop.

2) may or may not eventually stop.

Which of above 8 possy is most correct, depends on one's knowledge of

t. corpus. e.g. (B) if one knows t. object ends at that pt., (1) must be true.

In (A) if one knows t. object is an inf. string, (3) must be true.

In (C) if one knows part of a finite object, (2) must be true.

(2.29) Seems to solve t. problem!

So ~~next~~ Next, we give a "true" p.d. of t. Myrd corpus, and our state of knowledge (in 2.29) of each of its corpus elements, we use a bit to show  $\exists$  a machine assoc with that p.d. that state of knowl. "I." state of knowl." enables us to tell which codes to this machine to accept.

In t. UMG (ALP) approx., we use t. same rules for accepting the codes,

So we end up w. a "cutoff factor" that is t. pc of t. "true" machine that

represents t. p.d. — and from this, we can use 7&3 to prove Pao upper bound error, using 11.03

92898: Jarem;

Not a good  
beginning.

I have read your "how-to-learn" paper. Some time ago, but until now haven't found time to write you about it.

Though I found several <sup>As I see it,</sup> ~~features~~ <sup>bugs</sup> of the program, it's likely that these bugs can be fixed (i.e. indeed, they are bugs!). The main question is ~~is it~~ — Can one get a machine to learn how to learn, rather in what will ultimately be the best possible way, by using a relatively simple, ~~redundant~~ <sup>redundant</sup> routine? If so, is this a practical way to get ~~into~~ <sup>into</sup> very intelligent behavior? — How long would it take, and how difficult would it be to design a training sequence for such a machine?

A point ~~from~~ <sup>from</sup> these remarks, I found the ideas to be stimulating — Ray forced me to express more clearly, my own goals and the ~~background~~ <sup>background</sup> I expect ~~that~~ <sup>that</sup> I've been using and expect to be very so <sup>constraining to these.</sup> ~~believe from.~~

OZ898 TM: & Jørgen S.

01:15:50 Looking at his "Learning how to learn" Can this be formulated as an OZ problem? — The roles about how ER ↑, & t. unimeta goal (whatever it may be... eg. J's ~~unimeta goal~~ can be made part of t. problem.

04: Hvr, we may also view it as a "separational OZ problem": A sequence of OZ problems, in which TM is supposed to learn a better & better P.D. as it continues. This is perhaps what I was aiming at in 108.09 (i.e. preceding work) on "oz probs".

Say I have a set of OZ problems:  $[M_i(\cdot), T_i]$  — Then there is some 'Algm, A ( ) which takes as inputs  $M_i$  &  $T_i$ ;  $A(M_i, T_i)$  outputs a string,  $X_i$  such that  $(i.e. \exists M_i(X_i) \text{ is } \dots)$  Consider  $G(\vec{A}) = \sum_i M_i(A(M_i, T_i))$  (Assumt Goes as "LINEARIZED")

for every Algm, A  $\exists G(\vec{A})$ : we want an  $\vec{A} \Rightarrow G(\vec{A})$  is max. — BUT: How much time do we have to find  $\vec{A}$ ? Do we have time  $\leq T_i$ ?

12: Hvr. The main idea here, is that J's problem "Learning how to learn", is really a "ordinary" long. problem. If there are better ways to learn than Lrn OZ, then they are within a constant factor of OZ.

— Hvr, I still don't understand t. OZ problem (108.09 i.e. preceding work) t. "separational set of OZ probs" of .04 may be a key.  $\rightarrow 117.31; \rightarrow 127.01 \rightarrow 128.01$  Looks like a quite a bit quite soln. to t. OZ problem!

While J's "Long. problem" would ordinarily be an ordinary OZ problem & it is not for 2 reasons: (1) It got F.B. ~~cancel~~ (parallel & values) during its trial. (2) It is only allowed 1 trial.

Hvr. t. way such problems are solved ("normally"): TM makes a model of t. G function & he'll climb repeatedly - (in private) - "cerebellum" — cheap — Pica makes u.g. external trials.

D298 Jørgen S.: One main idea of "Learning how to learn" is minimal "Bias": that one t.

System Lrn in a way conditioned mainly by its TSP.

However, the language used by t. system. (instruction set) is as strong a Bias in J's

system as it is in ALP.

A critical system Q is: "Where is t. Summary of Experience stored". J. stores it on t. stack, in t. individual stack components

T. idea of .12-.20 is attractive (now that I've "solved" t. OZ problem!). J's system may be a reasonable approx. of t. Genl. soln., & if so, can be criticized, improved, and used, out of Basis.



.01: 108.90  
Def  
.02

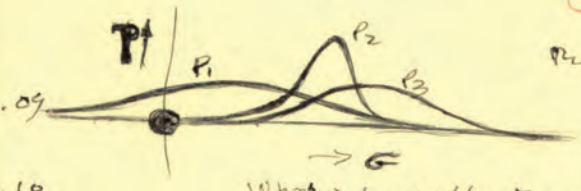
A direct look at 4 02 problem sym: we have this data set of 99.39!

21.577  
24.5  
= 1.7  
21.58 = 17.8

[ 02 probs; HC method;  $c_{ij}$ ;  $G_{ij}$  ] : Given to cc value T,

we can, in the manner of SGA, get a d.f. for the G obtainable w. an array HC method, on (02 prob; lem) using  $cc=T$ . From this d.f., we should be able to get the prob of each 02 method giving the best G in time T.

N1699: I Don't see how SGA is Relevant!



These 3 d.f.s are the pdf of G values for 3 HC methods w.  $cc=T$  on 02 probs. (= 02 p's).

.10 What is the probly that each will be max?

Actually, we can jump to .31 from here, but .10-.25 is conceptually of interest & may be good way to do Approx.

Probably that  $P_{1,2,3}$  are 3 p.d.s

$$\int_{-\infty}^{\infty} P_i(G) dG = 1$$

Probly that "i" will be max is

$$\text{say } F_i(G) \equiv \int_{-\infty}^G P_i(G) dG$$

$$\text{So } \int_{-\infty}^{\infty} P_i(G) F_i(G) F_j(G) dG$$

Say  $H(G) = F_1 \cdot F_2 \cdot F_3$

$$\text{then } \int_{-\infty}^{\infty} \frac{P_i(G)}{F_i(G)} H(G) dG$$

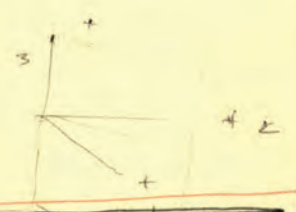
.21 So we do each  $\frac{P_i}{F_i}$  for each i: Then dot mult it by  $H(G)$ .

If we assume independence, then solves the problem.

.25 If " " dependent, I think the solution is still about the same form.

Another way to think about it (for dependent or indep ~~prob~~ p.i.g)

Integrate  $P(G_1, G_2, G_3)$  over the region of space in which  $G_1 > G_2$  and  $G_1 > G_3$



The region  $P_1 > P_2$  is bounded by 2 plane thru the "3" axis

Region  $P_1 > P_3$  is bounded by plane thru ~~the~~ 2 axes.

We could divide up the space into 6 (=3!) regions of  $G_1, G_2, G_3$  in ~~order~~ order.

We want same  $cc=T$  for all 3  $G_i$ 's

$$\iiint_{\substack{G_1 > G_2 \\ G_1 > G_3}} P(G_1, G_2, G_3) dV$$

Volume Element of the Space

**Note Bene**

The Probly distributions are based on the complex of .02 - first problem to be "interesting".

A finite c.B. (RPL) is necessary. If  $c=B=\infty$  for comp of space p.d.f. then

many because the probly become ~~unreal~~ - since in ~~the~~ end of time,

we could compute exactly which how much G would be obtained by each HC method on the 02 probs. - so the S of .31 would be trivial! → 127.01

.01 : On Search for regys, ect. All regys are searched for, using only previously devrd tech neques. (plus "Primitive" techniques).

Presumably — at each stage, the search techniques are "Universal", in the sense that if they had a large enit

.05 CB, they would eventually find any findable regy.

(SN) I'm not sure the "Universality" condition is a neccy property of each search. Perhaps it is only neccy for very diff. probs!

(SN) Perhaps  
Have only one  
Topic on a Page.  
If I have a new  
idea - put on  
a Separate Page  
More Pages, but  
Certainly  $\ll$   
Twice as  
Many PP.!

.15  $\rightarrow$  A possibl. reason for insisting on (SN) .01-.05 : we want to be sure that TM has ability to devr. all of the search techniques it uses. It is then more likely to be able to discover new search techniques.

It would seem that normally, for each new problem, RM would ~~first~~ (usually) first search for "solns to that specific problem" i.e. code that left p.c.'s of other prob. solns. invariant — but tried to  $\uparrow$  p.c. of ~~some~~ <sup>trials</sup> trial for present problem. This corresponds to the "TM<sub>1</sub>" phase. Next (or later), it may look for code modifns that affect  $\geq 1$  problem (or) set of <sup>jobs</sup> problems of a same macro problem. This corresponds to "TM<sub>2</sub>" phase. — or "Updating".

Suppose the Corpus is an ordered list of problems; Some INV problems, but many OZ problems. RM can go ~~at~~ at his own pace.

The goal Might be: To get Max total "G" + some Weighted

Sum of G's for **INV. Probs.** — Or, just Max total G for OZ probs,

~~but~~  $\rightarrow$  since <sup>the</sup> probs have to be worked in specified order, all

Inv probs up to the last OZ prob. worked on, have to be solved.

[classical OZ probs: CB given for each.] Another formulation would give a certain sum of time for total <sup>Time</sup> of all optzn. probs given. There may be several ways to formulate the global goal

To transform an OZ prob into an INV prob. We want Min (time to) get  $G \rightarrow$

a certain " $G_0$ ". While Puz clearly makes them into INV probs, we seem to lose the H.C. info used in most OZ solns. — i.e. ~~in~~ L such for INV probs, info in past trials is discarded — a new task Puz facts that they had been tried & they had failed!

If we give T.M. traces of all of its previous trials, then a problem can never be subject to  $> 1$  trial! — since the set of past traces (The history) is always different for each trial. (Shades of Turing's "you only live once"!)

**N.B.** My "Genl Soln" to the "Mxd Corpus Problem" was for

a pure prediction problem only. 117.0-40 supposedly tells how to

apply Pz. resultant p.d. to the soln. of OZ probs.

Actually this "Genl. soln." is just what we need for the "updating" part of the T.M. routine.

Ont. General Search & Discovery Process!

Concept "Acquisition" @ Discovery of Concepts via L such.

(b) Evaln. of PC's of concs: condl. & uncondl pc (?) (useful?) <sup>is this</sup>

$\beta$  Cond.: 1) Within "problem area" (e.g. This specific problem, Geometry, math, Science...) 2) Between different "problem Areas"

Go thru ANL Notation Learning T.S.Q.: See just how the topg. is relevant.

One thing of import is "IR": It is both a problem for  $\downarrow$  cc & search & a problem for  $\uparrow$  pc of "searched" — for "object". I was simply putting solns. of probs into "storage". What is justify. of Puz, & how should condl pc's be assigned properly? (Both in accord w. ideas of ALP & w. heuristic ideas about Search).

In ANL (with  $\uparrow$ ) the system I used was pretty much like having a sequential predn. problem, and after each "example" ( $\equiv$  micro corpus),

we "summarize" the implications implied p.d. due to the data up to that point,

From this "summary" we are able to give a p.d. on the future of the sequence.  $\rightarrow$  [There is the problem of "One Shot Learning" — & ALP has to deal w. Puz in a special way]

Note at the ANL I was ~~doing~~ doing "Operator induction" —

.01 Which assumes the inputs are uncorrelated & have no info-content. i.e. it ~~doesn't~~ doesn't take advantage of any ~~features~~ features in the inputs.

.03 P.c. of op operator/would be  $P(O_i) = \prod_{j \in O_i} P(O_j, A_j)$

$P(O_i)$  is the prob of operator  $O_i$ .  $O_i = \{Q_j, A_j\}$  is the probly ~~input~~ (as given by  $O_i$ ) that the input  $Q_j$  will result in output  $A_j$

.06 The equal soln. of the operator induction problem then sums, over all  $O_i$  —  $P(O_i)$  is given by the lengths of terms of  $O_i$  (summing over all terms). We can obtain .06 w.  $\epsilon$  in input time. — one input to decb.

.10  $O_i$ 's, & the other inputs for  $O_j$ 's; & the output will be  $A_j$ 's.

Next, the problems:

.14 The idea of the Alp. Notn. Long TSP was  $\alpha$  this! We start with a short corpus.

.16 We have a set of concs so we can use simple Lsrch on those ones, w. uniform P.D., to do trials. We find a soln. to the problem, which amounts to a short code for the corpus. It implies a diffent p.d. on codes for the future. This implies p.d. amounts to a "summary" of the past.

.18 Using  $P_i$  P.D. we L search for solns to the next problem set. Loop to  $\alpha$  (.16 &)

.19 Now: to what extent was NL completely decb by .14 - .18?

Given a set of examples, I would create a short code for it using "operator induction". For some reason I then put C in the set of instructions to be used in generating trials for the next batch of problems. — why?

Well, say TM learns how to deal w. "+" problems. — makes useable operator! Then w corpus of + & — he learns to do "—" probs also! But he then has 2 operators: Each is vit  $\frac{1}{2}$  the time! Still a big compression! →

(SN) Using randomness, w. many digits (for examples) is equiv. to using larger  $\epsilon$  or smaller nos. — the total compression for the vit operator is the same in both cases. The method of using these 2 operators — of correcting for errors... look into this for "operator induction".

→ Another kind of "Long" that I may want to analyse from the pt of view of  $\epsilon$ . 14 - 18 is soln. of GA problems by Lsrch.

→ Back to .19 ff: In "unrelated finite objects" corpus: Say the objects are large, but they have little or no mutual info: To generate the corpus: The grammar consists of a list of indiv  $d$  terms of each of the objects. We choose each term w. = probly.

Notcom  
547  
3773  
Charlotta  
01  
528  
3356  
Alex  
100 Start Sun

Say  $A_T(x)$  &  $B_T(x)$  are 2 proby distribns. — They are each finite CB ( $C_B = T$ ) approxs to f. ALP of  $\Sigma$  string  $x$ .

$$A(x) \equiv \lim_{T \rightarrow \infty} A_T(x); \quad B(x) \equiv \lim_{T \rightarrow \infty} B_T(x)$$

$$\lim_{T \rightarrow \infty} \frac{A_T(x)}{B_T(x)} = \frac{A(x)}{B(x)}. \quad (\text{We can be sure } B_T(x) > 0.)$$

Now: what does it mean to say  $\frac{A(x)}{B(x)}$  is "not enumerable"?

In what sense is  $\frac{A(x)}{B(x)}$  "less computable" than  $A(x)$  or  $B(x)$ ?

(Say both  $A(x)$  &  $B(x)$  are incomputable.)

It may be that while  $A(x)$  &  $B(x)$  are both incomputable, we can at least get an ~~upper~~ lower bound for each, which  $\uparrow$  as  $T \uparrow$  — so r. bounds of our approxs. are Monotonic in  $T$ .

Th. Successive approxs to  $\frac{A(x)}{B(x)}$  are not monotonic in  $T$ , so it would seem that our uncertainty of  $\uparrow$  would be usually "worse" than our uncertainty of  $A(x)$  or  $B(x)$ .

On the other hand, the non-enumerability of the recursive functions, does seem to make the "summing over all c pairs" an essentially meaningless statement.

Perhaps see Li-vit for explain. of "enumerability"

228 <sup>and I think I said in Sol 78</sup> They do say that there is no normalized proby measure that dominates all normalized proby measures. — The say universal semi-measures dominates all finitely decodable semi-measures.

The implication of 228 is that ~~there~~ <sup>rather</sup> ~~no~~ <sup>no</sup> off. normalized constants of 2 universal semi-measures can  $\rightarrow \phi(corpus)$   $\approx$  t. length of corpus  $\uparrow$ .

Maybe the truth of 228 depends on what defn. of "normz." is used. ! ))

It may have been proved for a normen system Different from Mine!

I think what they mean by "normz." is simply  $\sum p_i = 1$ . I think the theorem is that there is no universal, finitely decodable, "normalized" proby measure. — This would include all methods of Normen.

W Nov 25, 1971

23  
1.25/d  
3.07/yr

I had the idea that one could try to Model the discovery of various Sci Concepts by seeing if the CJS for the proposed models ~~are~~ <sup>were</sup> reasonable.

One initial difficulty was that it was hard to tell what input into each scientist had: — Pro from his writings one could often get a minimal set — a "lower bound" on input info.

.06 Another <sup>even</sup> more serious difficulty. Lsrch is based on a Prob. Distribn.

These are Conditional P.D.'s — based on many features of a problem.

Many of these ~~PC's~~ <sup>condl. pc's</sup> used in Lsrch of a human, would be computed "unconsciously": Tho these condl. pc's are quite impt., we will very often have no conscious access to them — the Scientist doing the search will not have conscious access to them — Pro "Historian of Science" (Will have even poorer access.

— Yet knowledge of these (unconscious) condl. pc's appears to be critical in tracing possible paths of Lsrch in evaluating C.J.S.'s.

The colloquial (non-rigorous) analysis of the Hist. of Sci, is even less able to deal w. the problem of .06! So what good is Hist. of Sci?

Well, Hist. of Sci. can try to find possible modes of conscious evaln. of condl. PC's to be used in Lsrch. The idea is that we may in this way, discover

Some impt. tricks (=Heuristics for search/discovery) that were used by v.g. Scientists of the past. — but not all tricks, or perhaps not even ~~known~~ many v.g. ones!

Finding models for past discovery can help a team present day Scientists — even if the models are not "true" — if they are "good" in a short code sense, then they are of interest — they are potentially useful.

Main Problems

Suite

1) **T.OZ** problem: T. "Soln" of (17.01-.31) is o.k. as a "conceptual model". The details of how best to Approximate (this, have to be worked out)  $\rightarrow$  Also + ideas of "Information Based Complexity" ( ) relevant?  $\rightarrow$  (See 127-128 for pretty Good Approach)

2) 118.01-120.10: This is an over-all view of how TM operates: The idea was to take Alg Notation Lang ( $\hat{=} \sim$  TSP's) & try to give a better theoretical understanding of them - & (perhaps) improve them... Then try to fit them into the Mxd corpus

Scheme.

3) "SMA" problem: Space Mks/detector problem: "What is the expected (unbiased) ~~approx~~ estimate of "yield" of a search. or "Action Algorithm".? I don't remember just how far I got in solving this problem: There were at least 2/3rd dip  $\rightarrow$  in ~~investigation~~ <sup>relatively recent</sup>

One or both in "SM" file: perhaps within the last year.

At present, I'm not sure it is very critical in the TM goal. - It seems most

relevant when one has large data bases, ~~for~~ when one needs much accuracy in pc estimates,

- But I don't know. For TM applic., I was thinking of evaluating TM's overall

"Action Alg" and trying to Optimize it. In any optimization problem over

statistically fluctuating <sup>(noisy)</sup> variables, one has this problem of "SOY"

(spurious optimization yield). / <sup>At least</sup> One of the approaches attempts to ~~estimate~~

estimate the SOY for various statistical situations. The problem is particularly Bad

(Sever) when ~~one~~  $t$  variables one picks  $t$  best of have noise levels that are

$\rightarrow$  B. differences betw. their means

Is "SOY" relevant to (.01) (117.01-.31)  $\leftarrow$  in particular? in .01, we're not

[ (Yat!) so much concerned w. the magnitude of error: all we want is a p.d. over  $t$  choices, of what the "best" choice is.  $\leftarrow$  This is also what we want in Straight SM work,

th. need for quantification arises when we compare the peaks obtained from somewhat different data sets. Presumably, if we had a joint p.d. for both data sets together, we would not

need to know SOY. Hvr, it may be that in SM (it may be in TM?) we often have to

compare  $t$  "Best choices" betw 2 data sets & we don't have a joint distrib. funct.

A way to solve the SOY problem: Say we have a joint p.d. of  $G_i$  ( $i=1/k$ )

of  $k$  variables: We have way to select  $t$  variable ( $i$ ) of max expected  $G$ .

So we just add a  $k+1$  variable,  $X$ , of known value  $\geq$  zero, to the set.

for each value of  $X$ , we pick  $z$  that has max expected value of the  $k+1$  vars.

for a certain  $G_0$ , when  $X \geq G_0$ ,  $X$  will be always chosen. for  $X < G_0$ ,  $X$  will

never be chosen.  $G_0$  is the expected  $G$  of the max choice.

Try it w.  $k=1, 2, 3$ : See what happens: Does it give correct answers?  $\rightarrow$

We probably know correct answers for  $k=1 \hat{=} 2$ .

For  $k=1$  it works, but there is no SOY, so not so interesting!

For  $k=2$  (or  $k=3$ ) What is the probab that all  $G_i$  ( $i=1/k$ ) are  $< X$ ?

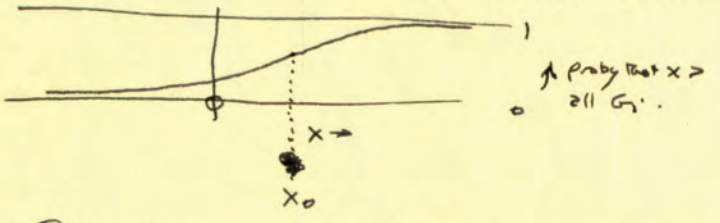
Write more frequently, more detailed reviews

Mix & corpus problem ends 115. to

How impl. is J's idea of: "You only live once" & we can never repeat a trial.

N 698 TM:

say we have an  $X \rightarrow M \rightarrow \text{prob} = .5$ . It's just  $X$  the value of  $\theta$  peak of the set of vert.  $G_i (z=1/k)$



.05 One, perhaps impt thing that some of the Say work did: It assumed a model for the variables, so that e.g. the variables were generated by having their means be normally distributed, then each had a  $\sigma_i$  about that mean. — the  $\sigma_i$ 's could be assumed same, or they could be all diff or be related to the means. What one can "assume" depends on how much  $\uparrow$  in pc we get using the "assumptions" as a model for the data.

This (.05) model idea ~~is~~ is to p.d.'s it gives, seems related to the practice

02 prob (em of 123.01 ( $\equiv 117.01 - .21$ )) — This stuff is very theoretical: The SMA problem gives p.d.'s that one can apply this to — and it produces a model (class).



Assume we start w. info about  $t$ , past of  $t$  seq. Machine coding from point  $t$  to  $n$  to  $i$  future. So we code w. r.t. a known past.

Say the part of  $t$  future contains  $n$  points. If we agree on linear regression, the structure of  $t$  model is  $t$ : self-def. code for  $k$ ,  $t$ : no. of params.

Say  $\sigma_k^2$  is the observed error in  $t$ : best  $k$  param fit to  $t$ :  $n$  data pts.

So: if  $P(k)$  is  $t$ : self-def. prob of  $t$ : integer  $k$ . Then  $t$ : p.c. of  $t$ : data using  $k$

Cost coders  $P(k) \cdot \left( \frac{\sigma_k^2 \cdot (n+k)}{(n-k)} \right)^n$   $\sum_{i=1}^n x_i^2 = \frac{1}{n} \sum_{i=1}^n x_i^2$

P.c. of  $t$ : "error" is  $\frac{1}{\sigma} \prod_{i=1}^n \frac{1}{e} e^{-\frac{x_i^2}{2\sigma^2}} = \sigma^{-n} \cdot e^{-\frac{\sum x_i^2}{2\sigma^2}} = \sigma^{-n} e^{-\frac{n \bar{x}^2}{2\sigma^2}}$   
 $= \left( \sigma^{-1} \cdot e^{-\frac{\bar{x}^2}{2\sigma^2}} \right)^n = \text{max when } \sigma^2 = \bar{x}^2 \Rightarrow \left( \sqrt{\frac{1}{2\bar{x}^2}} \cdot e^{-\frac{1}{2}} \right)^n$

Say  $\sigma_k = \sqrt{\bar{x}^2}$  |  P.c. of "error" =  $\frac{1}{\sqrt{e}} \cdot \frac{1}{\sigma_k}$

So we want to "choose  $k$ " so  $P(k) / (\bar{x}^2)^n$  is min.

**No!** error =  $\sigma^{-n} e^{-\frac{1}{2} \cdot n} = (\sigma \cdot \sqrt{e})^{-n}$ ; so select  $n \rightarrow \frac{P(k)}{(\sigma \sqrt{e})^n}$  is max.

or more exactly  $P(k) \cdot \left( \frac{\sigma \cdot \sqrt{e} \cdot (n+k)}{(n-k)} \right)^n \cdot \sigma_k$  Select  $k$  to

maximize!  $P(k) \cdot \left( \frac{e \cdot \sigma_k^2 \cdot (n+k)}{(n-k)} \right)^{\frac{n}{2}} = P(k) \cdot \left( \frac{n-k}{(n+k) \sigma_k^2 \cdot e} \right)^{\frac{n}{2}}$

$P(k) \geq 2 \cdot (\log k + \log \log k + \log \log \log k + \dots)$  Use logs to base 2: Take only terms  $> 0$  in  $t$ : sum

The constant  $C$  is  $\approx 3.4903293 \dots$

$P(k) = \frac{C}{k} \cdot \frac{1}{\log k} \cdot \frac{1}{\log \log k} \dots$

only use factors  $< 1$ .

| k | P(k)   |
|---|--------|
| 1 | .35    |
| 2 | .175   |
| 3 | .0735  |
| 4 | .04375 |
| 5 | .025   |
| 6 | ⋮      |



For more accuracy: probability of  $t$ : max isn't sharp in  $k$ ! Use  $\geq$  wtd. sum of probabilities averaged over  $k \geq 3$  to  $\infty$  w. wts = .22.

T. Forgy works w. non-linear regression, so  $t$ : extent that  $t$ : system is "locally linear" to a distance of interest. This distance is as  $n \uparrow$ . So it's probably "locally linear" if  $n$  is large enuf.

We also assume that  $t$ : function to be optimized is "smooth": In SM strategies

also discontinuous, this is a source of trouble!

.01 120.40: A possibl. Good way! Consider <sup>[The this Arg't. seems valid for "Mixed Corpus" pure sequential prodn.]</sup> <sup>ut. most General kind of problem</sup> <sup>we get (chunks) of corpus to</sup> codes. After each "chunk" we "up date" to P.D. — what this means! We construct a new unc. that summarizes the corpus up to the end of ~~the~~ <sup>it</sup> must recent chunk.

.04 By "Summarizes" I mean that ~~the~~ <sup>the best search one can do is w/ random input</sup> I suspect that there are many ways one can code equivalent of .01-.04

Min P.D.s:

|            |
|------------|
| 118        |
| 119        |
| 120        |
| <u>126</u> |

→ a) To what extent was <sup>the</sup> "putting of the latest soln. of the ALL problem into storage" an example of .01-.04?

→ b) To what extent was "my soln. of the Muxer problem via 'such' (instead of GA) an example of .01-.04?"

Try to give several examples of .01-.04 equivalences:

1. We have several <sup>(copies)</sup> codes for the corpus. The final summary is a wtd. sum of these codes  $\equiv$  a wtd. mean of their products. (Products  $\equiv$  P.D. for future):
2. Sometimes A codes for ~~more~~ <sup>several</sup> sub-corpi: The equiv. code for the total corpus is cart. product of those codes (if they are all indep.).
3. Just how do various heuristic types fit in? { e.g. "Quick abort" }  
e.g. How are "Info. Retrieval" p.d.s for obtaining more relevant "Concepts", implemented in this schema?

[SN] In General, ~~when~~ <sup>when</sup> one has solved a new problem, this does modify the P.D.

Some, by itself, but the updating of the ~~current~~ <sup>current</sup> P.D. implied by the or more new problem solns. can take much more time  $\approx$  be a more diff. problem than the original problem  $\leftarrow$  soln. ! Perhaps to start off, the "New Machine" should have good updating algms. as we can devise. <sup>solns. of t.</sup> The / updating problems can be considered to be

? → part of the corpus, so, presumably, the updating algms. will get automatically better as TM Matures. — So is this like in old TM<sub>1</sub> = TM<sub>2</sub> where we used the "good soln"

to decide what fraction of time to use in improving the General system ( $\equiv$  updating algms) v.s. direct soln. of problems

So: .20-.25 seems to be the crux of the "Updating problem"  $\approx$  is perhaps

the most difficult (is least understood/worked out by me) part of RM. (Essentially the "TM<sub>2</sub> problem".)

.01: 117.90: In 117.01-.90, we end up w. a p.d. over various H.C. Methods, but the methods is rather indirect & it involves a long chain of calcs. There should be a more direct way to get this d.f. or approximate it.

**A trial approxn:** Say we have the data set 117.02: We want to "Explain it"

by writing various codes for the entire data set. Actually, what we want is a "Summary

d.f."  $\rightarrow$  given any [OZ prob<sub>i</sub>, #OT<sub>j</sub>, CC<sub>k</sub>, G<sub>n</sub>] as input, it will have a

Optim Technique

stack operator w. the proper prob<sub>i</sub> of that combination.

**SN** In General, The "Summary d.f.'s" Are very imp. — They may be mainly what we calculate in TM (126.01-.04)

.13  $\rightarrow$  From .06, by integral 117.31, we obtain the desired p.d. on OZ<sub>i</sub>, OT<sub>j</sub>, CC<sub>k</sub>

.14 T. Set 117.31 may be (easily) evaluated by Mc Carlo. ! I'm not sure, but since we only need approx. values, this may be adequate. T. idea is that each time

we evaluate  $P(G_1, G_2, G_3, \dots)$  we do get useful info. for each such pt., all we have to do is figure out  $\dots$  No!

One way: for each  $i$ , we find ~~random~~ random pts in the subspace

$G_i >$  (all other  $G_j$ 's).  $\rightarrow$  We add the  $P$  values to the "register":

Another way: we pick pts at random in the space: At each pt. chosen, there is one  $G_i$  that is max: we then increment the "i register". At the end of many trials, the i register will give a profile of the <sup>relative</sup> prob<sub>i</sub> of  $G_i$  being max.

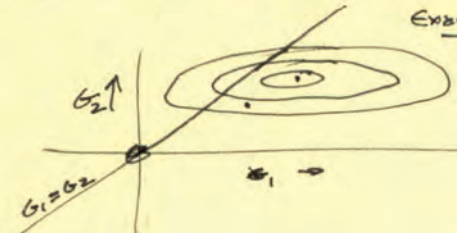
.24 Trouble is, the <sup>rel</sup> density in most of the space will be very low, so the above method will tend to be very inaccurate. Would it help to find the peak region in the space? Maybe! Then

by picking pts in the peak region and make just much of the total (w/).

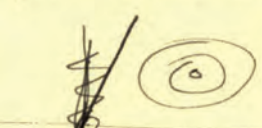
If  $P(G_1, G_2, \dots)$  is expressible (indep prob<sub>i</sub>) as  $\prod P(G_i)$ , then we can find the peak of each  $P(G_i)$  and width of peak. — so a "multivariate function" w. known axis. This may not be so hard to deal w. One knows the peak, & how rapidly it

turns peak, in each direction. T. Monte Carlo approxn might be easier than!  $\rightarrow$  128.01

Exhibit 1: for 2 dms, it's easy: x for  $\epsilon$  coords unill the distribn is isotropic;



Then it's easy to compute the  $\int$  on one side of a line.



Just rotate coords and it be even more simple evf.

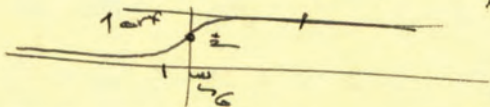


For  $> 2$  dms, it seems much more difficult!

o: 127.20: If there aren't too many  $G_i$ 's; one can find s. center of d.f. From find the distances from it to each of the hyperplanes that help define  $(\frac{G_1}{G_2}, \frac{G_3}{G_2})$  max of all  $G_i$ . From that info, one may be able to estimate the wts in each region. Remember, we don't need much accuracy!

T. techniques of [17.01 - .21] may be fine if we assume indep. & normality. Since we don't need ~~very~~ very much accuracy,  $\int_0^1 p_2(x) dx \approx \text{err}(x)$

Can be approximated fairly well. T. middle section is easy to approx. & tails are probably easy also - but don't approximate.



We may not need to go out far out tails because we would also cut in that region infeasibly

(SN) See [130.25 R] for a complaint about this "soln" to OZ.

- Main path: General Structure of TM!
- 118
  - 119
  - 120
  - 126

17

Handwritten text at the top of the page, possibly a header or title.

Handwritten text in the upper middle section of the page.

Handwritten text in the middle section of the page.

Handwritten text in the lower middle section of the page.

Handwritten text in the bottom section of the page.

Vertical handwritten text on the right side of the page, possibly a list or index.

Large, faint, mirrored watermark text in the center of the page, including the words 'COTTON FIBRE', 'SOUTH AFRICA', and 'S.A. PATENT'.

.33 Note of 6/23/01

of 126.40: I want to make a list of things I want TM to be able to do, & to as much extent as possible, How TM does these things. The problem solns. will use one another.

① (Basic Operation): L search: Given a P.D. & a problem (INV problem or OZ problem) it performs L search to find soln. In case of INV, we may or may not want it to continue to search after it's found 1 soln. In the case of OZ, it can continue in "Anytime" mode to continue to search in "available cc".

A critical component is v. P.D. In the case of INV probs, this is obtained from corpus of past INV probs & their solns. In the case of OZ probs, we use the corpus of past [probs & solns], and the techniques of 117, 127 & 128 to obtain a suitable P.D.

-18 ② (Basic Operation) Updating: This consists of "improving the P.D.": In one sense, it is an "Anytime Problem" (an OZ prob. w. unspecified cc) w. best compressness as Goal.

Superficially, there would seem to be a problem of just "what part of the corpus to compress." This is being that certain parts are more likely to be relevant in the probs. of the (near) future.

What role to older "failed" Probns? (Matters, since good) PLEX in (ARM) (Human) Search, D's coverage. Exactly how over time used? (Quant. Discr.).

T. simple criterion of "total corpus cost" sounds not v.g. If there are long time series w. lots of data in them, then T.M. would satisfy this criterion by

Working on them mainly. [This may be not: Since these are OZ probs., one has to have some sort of "linearized Utility funct", to decide how much cc. to spend on each.]

In general, the idea of "improving the P.D." is undefined. How good a P.D. is depends on one's CB. Here, P.D.'s may be partially ordered.

Also the "goodness" of a P.D. could depend not only on giving a hypc. for the Corpus, but a measure of how good it was in guiding rapid solns. to problems. (e.g. QuikLabort)

D2.99 "Improving the P.D." is a task that needs to be more narrowly defined - by having the "User" & no more params.

1980: 85 R. 1905 One is "rich about" its small cc for many solns.

.33 6/23/01 No! T. User doesn't have to tell TM what part of the Corpus to Code.

T. User could give this info & it would help TM, but normally it is not & hence not be available. In Lugano (say April, 2001) I ran into this problem again - Decided that TM would have to "learn" what part of the corpus to code, (after it was given the cc limit (= Competition Bound = CB)). In general, the problem of induction is not to just put max # pc codes for the corpus in available cc, but it also involves learning how much of the available cc to use on each part of the corpus - i.e. how much each part is relevant to the particular prob. rather than just computing.

4/6/05 - See STM 59.14 for some "reasonable reasoning" on this stuff!

Some critical ideas in TM progress:

1) At the beginning I had the idea that L search could solve all INV probs in a close to optimum way!

"Optimum" w.r.t. f. data & hours & CB available at t. time.

I had the idea that OZ probs could also be solved in this way, but most recently, it became clear that I didn't know how to improve P.D. used, on the basis of new experience. (See 2.1 for soln)

2) New Developments: (a) Sol 78T3: t. convergence theorem for sequential data.

This explained how temporary convergence occurred - that it was f. p.c. of f. individual symbols (via AIP) that converged rapidly to t. probs via t. "true" generating sequences.

(b) That I. ideas on meaning of convergence and proof could be applied to a corpus of finite objects. (for applic. to ~~the~~ unordered set of finite objects, see "Computer Journal" paper 1998 or 1999 for ~~previous~~ meaning of conv. & applic. to finite objects and set of prefixes of finite objects. See 106-108; 110-111, 115)

(c) That t. ideas could be easily extended to a corpus  $\leftarrow$  (Mixed corpus theorem)

Consistency of ~~data~~ several (or  $\phi$ ) finite object and several (or  $\phi$ ) prefixes of infinite sequences. (See 106-108, 110, 111, 115.)

(d) That "updating" t. P.D. for, say sequential produ., consisted of obtaining a C.P.M. that summarized ~~the~~ ~~ways~~ ways in t. corpus thus far. There were various means

but one could do this in an approximate way. as an "OZ" problem. Her. t. "One Shot Long" problem, ("OSL") was a ditty that usually (if not always) had to be dealt w. in a special way.

(e) I finally found a formal (and good ideas on how to do approx) ~~soln.~~ soln. to t. problem of obtaining a P.D. for ~~an~~ OZ such, based on a corpus of solns of OZ problems plus any other (Mixed corpus) data that you might have. The soln is in 2 parts.

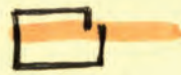
(I) The first part is to obtain a P.D. for t. G of any particular O.T.; applied to a / OZ problem for t. given C.B. - based on previous experience w. various

{ O.T., C.B., G, Prob } quadruples of t. past.   
 Hvr. Still selecting f. part of source randomizers is not always the best way to go! Soln w. 2 norms (vars. of 2nd norm & var. of 1st norm) of means & 2nd norm

(II) Integrating t. space of t. previous P.D. to obtain t. Prob that any particular O.T.   
 Best for t. Prob of intersegment for t. states C.B. (117.31) So OZ is solution   
 108.09 - .40; 117.01 - .40; 127.01 - .40; 128.01 - .13; 130.25R   
 Imptr Criticism.

3) Using the ~~background~~ developments (a) thru (e), it becomes possible to solve any solvable INV or OZ problem (starting w. a give initial P.D.): Then "update" t. P.D. so it summarizes t. past A priori info as well as t. problem just solved.

This updating process is a kind of OZ problem (But see 129.18-30 for some difficulties in adequacies of this process.)



Looking at the review of 130.01-40: I think the least completely flat out parts are:

1) How to solve problems (INVAR OE): By this I mean Not simply Lsrch - but just how the P.D. enables various kinds of heurs.

2) Also, in <sup>implementing</sup> heurs, just how does P.D. work? What CB is used <sup>to obtain</sup> for P.D. or other P.D. values? Note that the P.C. output of a P.D depends on CB used.

10:50 start  
11:00  
11:20  
11:30  
11:40

→ Perhaps best work on a problem, or problems in heurs, to see how this works. - In general, it would seem to take some cc to implement a heuristic: could we somehow add this extra cc of the brain involved?

Well, a single heur, may involve several trials. So how would this work?

3) I should write up Review more carefully. <sup>130</sup>

\*\* 4) Is there some way to automatically include the cc of the P.D. into Lsrch process?

\* 5) On 02 prob: look at (130.25 R): I'm not sure if SOY "soln" solves P.D.'s problem!

6) W.R. TM only working INU problems: T. first set of problems uses an uncondl appt. T. second set uses a conditional appt: T. "conditions" being the problem to be solved. So we need induction on finite set of finite objects.

7) Very Impl. idea on TSQ's: <sup>using ideas of</sup> ~~Research~~ MCT, I should be able to write TSQ's much easier! One easy way would be to follow conventional Math texts,

filling in w. "talking" (when RW knowl. seems heavy), using "hints" etc.

Try to get TM's "Concepts" ~ to those that I think I have.

Another source of TSQ's: Problems solved by other ML researchers.: Remember that present system is able to go from any particular state of knowledge to any other particular state: State<sub>1</sub> is given as to conditions in "Cond. proby."

Also Note: 146.30: How to solve  $1+3x=7$ , using idea that one pretends one knows  $x$  is one pretends various facts about  $x$ .

8) An Impl. 'set of probs. in main line of TM: Give examples of many types of <sup>(Heurs)</sup> probs. in their solns. How are they found to solns. put into the "Summary Machine?". What are different ways of summarizing?

9) continue from (7) on TSQ's: Another source: Look at classic works on A.I. on prob. solving (various A.I. encyclopedias). Instead they give ways to solve problems: Make TSQ to learn those techniques.



- .01: 113.40: Now, to what extent is "Conventional" GA  $\approx$  passing a "Stoch Grammar" thru "Updating"?
  - .02:  $t$ : data  $\approx$   $[cond_j, G_j]$ ?
- I think a "Great Breakthru" of 130.01-40, is the idea that making the stoch grammar is the "Updating" process.  $\approx$  the "Summarization" process. The main weak pt of this technique is that it doesn't deal w. OSL (over shot long)

The ~~Grammar~~ Grammar idea from supports the idea of "Building Blocks."

"Building blocks" are crass. But are used "often" in the deriv of  $t$ . by  $G$  set of cond.  
"Often" means defining them enables a more compact encoding of  $t$ : corpus ( $\in [cond_j, G_j]$ ).

1.3.99! OK: First realization that Normal GA & SGA are not similar to Lsearch solns.

Next: In "passing a grammar" thru  $t$ : data (.01-.02) we must note that it is a conditional Grammar

11:  $t$ : cond. being the deriv. of the "Optim problem". This means that any other "optim prob" marked

12:  $t$ : system in the past should help construct the cond. grammar for the present optim. prob.

W2(1), I'd like .11-.12 to be true - i.e. so each GA or SGA soln could learn

from previous (S)GA solns: But in fact, they do not. It is a weakness of both

GA & SGA. At this time, did I have any good ideas on how to "update" a conditional Grammar?

|    |    |
|----|----|
| 20 | 20 |
| 30 | 10 |
| 15 | 15 |
| 6  | 4  |

Now: Enter super-good mutation-only GA,  $t$ : system would keep a record

(a) of  $t$ : (trace) (part) of each search, & make an updated p.d. based on this past history, & make trials to both "inform" & to "get by  $G$ ".

(b) In a (mutn + crossover) GA system, ideally, each trial is made using info about all of their past & all trials.

$\beta$  is good if one has to do things in 11, but if  $\beta$  is done serially, it would seem to

have no advantage over  $\alpha$ . In both cases, one uses all info from past trials.

In fact, if we do  $\beta$  serially,  $\alpha$  can exactly simulate  $\beta$  & still be a "mutation only" system!

27: Look at this way,  $t$ : only deduce better mut. only & mut + crossover is that "mut cr" is parallelizable & parallel. is not. (para mut. is paired to  $\approx$  sim. oneal).

$\approx$  SGA is about as good as one can do [if one doesn't use .11-.12 info] - say one doesn't restrict SGA Grammars to ones in which by  $G$  pts are easy to find.

.27 is true for idealized soln; In practice, "mutation only" GA, only looks at info about last & ancestor! "mut + crossover" does look at many more (past) trials, but the means of using past info is usually rather narrow.

However: Consider the early Koza system (Opponent of function trees). To what extent can it be approximated by a kind of stoch Grammar? - or an "Evolving" stoch Grammar? To the extent that any gn. population, w/ [mut + crossover rules] implies a p.d. out next generation, - we have, by defn in evolving stoch. Grammar.

Handwritten text at the top of the page, possibly a title or header.

Handwritten text in the upper middle section of the page.

Handwritten text in the middle section of the page.

Handwritten text in the lower middle section of the page.

Handwritten text in the lower section of the page.

Handwritten text at the bottom of the page.

Handwritten text at the very bottom of the page.

Vertical handwritten text on the right side of the page, possibly a list or index.

I haven't read this art. yet! — some ideas:

According to Chris Morris model of dynamic processes, certain <sup>continuous</sup> functions can act as OMCs:

$\rightarrow F(X) \rightarrow \text{delay}$  I forget what  $F$  is, but. The info stored is int. value of  $f$  function (constant).  
 So  $X_{n+1} = f(X_n)$ .

Anyway I had the idea that we have 2 input continuous functions, we could somehow extract lots of computation! If we had  $B.V. = \omega$ ;  $T \approx \frac{1}{\omega}$ , noise (rms) =  $b$ ;  $\text{Path/Value of } f(X) = Mx$   
 \* Then if  $T$  is total length of delay line, it could store  $\frac{T}{\tau} + \frac{Mx}{b}$  different states ( $\approx \omega T \cdot \frac{Mx}{b}$ .)  
 $F(X)$  itself can only store  $\frac{Mx}{b}$  states, but with the delay line we can multiply that by  $\omega T$ .

By having many  $F(X)$  units interacting, we can ~~get~~ <sup>get</sup> a messy total memory.

with  $\rightarrow F(X) \rightarrow \text{delay}$  if ~~say~~  $T = \tau$  then we have 1 computation going on.

if  $\omega T > 1$  we can have  $\omega T$  indep calcs going on in the system.

$\Rightarrow$  I  $H(X, y)$  is a func. of 2 args (like ~~the~~  $H = y \cdot X(1-X)$  (logistic Eqn.))

we can have these functions interact.

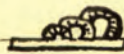
I had the idea that we could do induction in such a system! When the "output" was within  $\epsilon$  of the corpus, <sup>say RMS error.</sup> we would have an equivalent code length related to

$\epsilon$  in the corpus length. —  $\Rightarrow$  we could give val wts. to various approximator codes.

Hvr., I don't know where the input codes or the output!

So, I don't know what the "summary machine" is. I'd like to be able to get "summaries" <sup>Progressive</sup>

$\Rightarrow$  4. corpus & (T.S.Q.)



New tech: Say we have a bunch of interacting function generators & delay lines:

We can "start" the system by setting say 10 of the functions at ~~the~~  $(\frac{Mx}{b})$  different states, & all the rest at  $\phi$ .

We let the system run for a fixed time: the output is

defined to be a certain set of function outputs. (maybe all of them).

More Generally: A "computer" is a set of <sup>active</sup> elements connected in a complex way.

I. start state of the system is its "pgm": Its output can be defined in various ways.

1) AS .23 - .26

2) When computer "says" the output is ready: say on certain registers a/o where the system is in its state space, is an "output": Presumably there are large equiv. classes of states.

E.g. Getting output from a certain set of regs, is same as saying many states are equiv. (if they have the same output regs.).

I'd like to be able to characterize various kinds of computer systems.

What are good ways to search for  $\epsilon > 0$ ?

In an ordinary OMC (no under 10), even for induction, the I/O has to be very carefully defined (see Mr Letter to Lewis ...)

**<STRONG>CHAOS-BASED COMPUTING</STRONG>**, a fundamentally new way to perform computations by exploiting the ubiquitous phenomenon of chaos, has been demonstrated in a simulation by researchers in India and the United States (Bill Ditto, Georgia Tech, 404-894-5216). *Phonetic* Compared to digital computation, the chaos-based technique might come closer to how the brain performs computation, and might be superior in certain tasks such as pattern recognition. The computer consists of an interconnected grid of "chaotic elements," systems such as ammonia lasers which can generate unpredictable signals even though their behavior is governed by known mathematical equations. To encode specific numbers into each element, the researchers make specific signal patterns correspond to a number and ask each element to open its connection to the rest of the grid when it generates that pattern. Sending its signal out to the grid can trigger activity in neighboring elements. To carry out specific operations such as addition, the researchers connect the elements in a certain way. An unpredictable but deterministic avalanche of activity among the elements ultimately settles down to produce an unvarying signal that corresponds to the desired answer. Having demonstrated their technique in a computer simulation, the researchers are planning to test this idea with chaotic ammonia lasers and hybrid networks of nerve cells and silicon chips. (Sinha and Ditto, Physical Review Letters, 7 September 1998.)

**<STRONG>TRIPLE PHOTOIONIZATION OF LITHIUM</STRONG>**, a rare process in which a single photon removes all three electrons simultaneously from nature's third-lightest atom, has been detected for the first time by a Japan-US collaboration (Ivan Sellin, University of Tennessee, 423-974-2738). Studying this process further promises deep insights into the interactions that can occur between a trio of electrons and therefore a more sophisticated understanding of the interplay between charged particles in many environments such as stars. At the Photon Factory in Japan, an intense beam of extreme-ultraviolet (EUV) photons broadsided a beam of lithium atoms; a detector then recorded the rare process by collecting Li<sup>3+</sup> ions. In the most simplified picture of the process, an EUV photon deposits virtually all its energy into a single electron; the electron immediately shares enough energy with the other two so that they could all escape the Li atom. The three-electron interactions are relatively easy to extract from the data since the photon vanishes after striking the atom, and the heavier lithium nucleus acts merely as a sluggish "spectator." Researchers have observed triple photoionization of heavier atoms, such as neon, but such processes are typically more complicated events involving internal rearrangements of other electrons in the atom. (R. Wehlitz et al., Physical Review Letters, 31 August 1998.)

**<STRONG>BROWNIAN MOTION IS CHAOTIC</STRONG>**. In case one needed any more persuasion that chaos is all around us, a Brussels-MarylandUtah collaboration has for the first time demonstrated evidence for chaotic behavior in fluids at the microscopic level. The data consists of repeated viewings of a 2.5-micron particle suspended in water. A plot of the particle's position as a function of time is translated into a form which provides information about how the particle's position at one time correlates with the position at a later time. This analysis proved to bear all the hallmarks of chaotic behavior. The chief symptom of chaos is the tendency for particles that initially follow nearby trajectories to diverge quickly from one another. (P. Gaspard et al., Nature, 27 August 1998.)

On "Summary Machines": For predn. of a sequential corpus! Say we have several short codes: The state of the machine after it has just accepted the last bit of such a "short code", is a summary machine. A wtd mean of these ~~summary~~ & summary machines & one w. each of the short codes, is a better "summary machine".

A summary machine has the property that one can feed ~~it~~ it random input, & the d. f. on its subsequent outputs will be the p.d. of contin. of the original corpus. It means that L such on <sup>that</sup> random input is the best way to try for codes of longer sequences.

A poss. flaw in the forgo. Arg't.: (2) OSL would not be implemented, since the ~~state~~ set of short codes would not have any idea about tagys. That would occur w.  $SSZ = 2$ , when the corpus was augmented by one more "sample".

(b) (noted) (similar to (a)) If a tagy appeared once, each time the corpus was augmented & "updated" (updated = summarized), then perhaps the tagy would never be noticed!

(c) Consider a binary seq: If the "summary" is just the rel. freqs of 0 & 1, then all other info about tagys in that corpus will be lost. Also, any ~~new~~ tagys other than rel frequency, will not be able to use that part of the corpus as part of their sample - because that new tagy will not appear in the ~~code~~ that part of the corpus.

.21 Counter arg't: That ~~summary~~ any summary machine has inherent code for the corpus, so it has all info about the corpus in it!

Well, is .21 true? The summary machine has the state (Head state + ~~link to tape~~) of the machine after it has printed out the corpus. Hrr: the corpus is usually not recoverable from that state.

So, in general, the "update" is almost never a complete dem of the corpus! So usually we will want several ~~to~~ "summaries", including the un-coded corpus.

.27 By the way, how did I do "Summary Machines" for uncodable data? - Perhaps it was simply the (unconditional) stock lang. that generated the observed set of objects. Some "summarizing machine" in the "Mxd corpus" problem. (1). - That it seems like

.30 a ~~the~~ Discrete Universal D.F. I'm able to proceed ~~even~~ using "non-stopping" output.

↳ for contns. of the objects known to be finite, we have ~~the~~ probs. that are the sum of an (usually infinite) no. of <sup>rational</sup> binary fractions. - so the sum can be irrational.

.33 On OSL: In OSL, a summary code, that would be useful for OSL, codes the previous (not one-shot data) corpus in a non-optimum way: since it defines an object that has occurred only once.

.36 Q: Does  $\Rightarrow$  generalize to tagys other than Bernoulli seq. type tagys?  $\rightarrow$  (37.03)

Q: in L such do we also miss out if we use a pd that would not implement OSL?

N22.98 TM: CHAOTIC H.W. for TM

LEVIN

01:133.40 I t. ~~UNIVERSAL~~ latter, I was able to define input c.a., but for output I decided to have it stop: which is  
 very constraining ("Extension Complexity") is not so good for "Summarizing". It would be good if I  
 could get to General Umc to do Re  $\rightarrow$  "Universal Continuous D.F." as opposed to discrete d.f.  
 -02: see (35.27-30 for discussion of discrete v.s. continuous D.F. in the "Mkd Corpus problem".

D12.98: I got copy of Sinha, Ditto paper: also 1992 paper by some Guys on some  
 more General Principle in this Area. None refer to Chris Moore.

One (apparent) goal. soln. to Bern, OSZ problem was that idea of considering all possi. kernel definitions in getting to pc of all possi. contents. of t. corpus.

$$\sum_{i=1}^n \frac{2^i \cdot b^{n-i}}{i! \cdot (n-i)!}$$

.03 135.36: HA! (C) This may be a justification of retaining old Procs that didn't work so good as t. <sup>up to now</sup> ~~present~~ Best. Here, also, we retain Rave because they may be able to deal w. as yet unseen innovations in t. future corpus.

So: Simple "Summary Machines" are not able to do OSZ (cos) other input, regy discovery the Procs. Still, "Summary machines" are easy to make (compute) & may be able to do useful stuff, tho "imperfectly".

regy.

.10 Consider a "simple" sequential prodn (coding) problem. We code a short bit of corpus w. L-strat: use rotam base ~~sums~~ <sup>or. pgms.</sup> To code t. next chunk of corpus, we use m diffnt. "Summary Machines" for each, we do a simple L-strat continuation: we put on t. top m pgms, then 2nd Loop to  $\alpha$  ~~sums~~ <sup>seg 140.01-04 for an input Trm. on this.</sup>

After we've coded a chunk, we have Proc's m best codes for t. entire corpus!

When we do L-strat ~~on~~ starting on a summarizing machine of code length  $l_i$ ,

.17 we do L-strat w. an initial pc  $\alpha \geq -D_i$ .

.10-.17 is very easy to prove! (if we use a simpler kind of trace-like trace)

would it be easy if I used my (somewhat complex) "Functional lang".?

If one uses very large m, it could do some OSZ: (see 135.36): Very large

m, " means I could afford defining things that only occur once — The Proc seems like a grossly inefficient way to do it!

.27 Look at my soln. of t. Algorithmic Lang Problem (ANL) — See if it can be put into t. format .10-.17;

.29 Do same thing w. my 2 "GA-like" soln of t. 11 input Muxer problem.

Q: Just what is "Updating"? It is t. modifi. of t. P.D. so it will be better for future (Problems/corpus coding). How is "updating" done in (27) (29)?

Two Aspects of t. P.D. (1) Unconditional: It gives t. pc of each finite object — either a finite string or t. finite prefix of an infinite string. (2) Conditional: Condition is usually t. <sup>known</sup> prefix of a finite or infinite string; we want t. P.D. of t. continu. of that prefix.

.35 Hmm! .29 is interesting! For successive runs, t. corpus is t. same! <sup>So this is really different from what I had in mind for "corpus lang" or "prod lang" → 140.01</sup>

.36 (SN) In t. input muxer, I had t. idea that as  $k \uparrow$  t. "signal amplitude" would  $\downarrow$ . — well, it does  $\downarrow$ , ~~per sample~~ <sup>per sample</sup>, but  $t. \text{SSR} \uparrow$  rapidly w.  $k$  — So for t. total corpus, t. "signal" (i.e. t. amt. of compression obtained) may be constant or  $\uparrow$  w.  $k$ . — Hur as  $k \uparrow$ , t. cc part  $\uparrow$  and more!

I ~~was~~ (recently) was a bit confused about the definition of "convergence" & meaning of "error". This Q. comes up in the simple corpus w. pure set of finite objects.

T. Q. was: say we are computing the error int pc of  $x$ . ~~error~~ with bit of  $x$  &  $m$  object.

What is the set of finite objects that is used in its data?

Answer: none: "Since" ~~the objects are~~ p.s. of  $x$  objects ~~are~~ indep. of ordinary.

We have 2 p.d. that induces a p.d. on every object or partial object. T. ~~successive~~ ~~are~~.

P.c.'s of successive bits of an object are dependent only on  $k$  previous bits of

that object. // ...

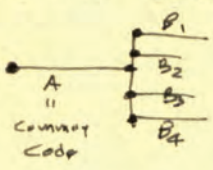
T. prop. is true for Mxd corpus as well.

T. basic P.d. (for ALP) is the ult. sum of all p.d., w. wts of pc of entire corpus - for that p.d.

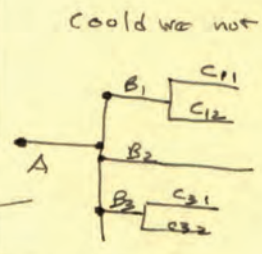
17 Plugs A possibly IMPT. Gauzy or MCT!

Consider not full MCT, but just finite set of finite objects.

Normally in MCT, the finite objects all have some common codes, & then each object has its own additional codes:



$B_i$  = individual additional codes.



Could we not have: we have 5 object forms, w. differing amounts of common code.

- A B1 C11
- A B1 C12
- A B2
- A B3 C31
- A B3 C32

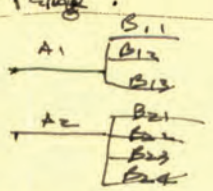
for example: 2 time series could have common params! Say they were both stock prices on ~~some~~ days. The "usual" way of dealing in this is to have a single

vector ~~series~~ time series - But .17 may be a more general. alternative way.

T. Gauzy to full MCT seems clear.

In general, one could take any ~~set~~ set or subset of ~~the~~ things in the Mxd corpus, & make a tree showing which codes <sup>for each "thing"</sup> was common to the code for which other "things".

.32 Another simple case:



Essentially a Philip Corp

697C Dash Job

in ~~the~~ .32 we don't begin printing until we start reading  $B_{ij}$

in .20 (R, L) we print only when  $B_i$  or  $C_{ij}$  appears.

is in .20 R  $A, B_i$  etc insts on how to construct a machine.



D.2.98 TM:

**LOGIC!**

deductive

A tentative idea: That "Logic" need only be consi, in order to be usable! In all <sup>almost all</sup> imp't. cases, we will use inductive logic (except for Meta).

.03

By logic I mean: Given a set of strings that are "true"; a formula to tell if certain strings are "true". I guess to be of interest, the system must be able to tell if certain strings are "false" ( $\equiv$  "not true") - i.e. one can't use a f. formula to show from True. [ The [ "not true"  $\equiv$  not derivable by the system ]  $\equiv$  false ] is ~~not~~ not obviously necessary. Must we use 2 valued logic? ] - Perhaps at least: True, false, undecidable.

I think logic is usually concerned not w. T/F of strings, but only w. T/F of WWF's (well formed formulae). If WWF's are CFL, we can probably assign integers to every wff. <sup>1 to 1 ( $\leftrightarrow$ )</sup> Perhaps this is what Gödel did.

⊙

It would seem that the most imp't. idea in a WWF formalism would be which strings are t. negations of which. Other than that, consistency, any implication formula would be o.k.

sec. 03

In order to be "sure" of deductive logic's conclusions, one must be sure the logic is consi. In general, this is impossible. One might know that a logic is "consi of depth  $< n$ ". So, for all demonstrations of depth  $n$  are ~~not~~ consi. Infinite recursion can be of finite depth, hvr. - We are only concerned w. use of programs - in general, for proofs for humans - these are of finite depth.

01: 137.35 spec ) What we did (& what GA does) is to find a (key compression) of the corpus, then use the resulting compression to implement summarizing p.d., to look for more compressions (keys) in the same corpus. Since any key is another parallel code for the corpus, one still uses the old

04 search codes, but one has some newer ones of higher priority. This is an imp. trim "on 137.10 ..."

(SN) Is this related to "Sequential Testing"?

In 01-04: Say one finds a new key in the  $r^{\text{th}}$  subcorpus! Would it not be well to recode the entire corpus w.r.t. this new key? Well, Yes & No! If the "Summarizing Machines" did a complete job, there would be no point in looking at the previous, un-coded, corpus — or looking at it in any way! However, if (as is usually the case) the summary is incomplete, the new key could very well ↓ rec of previous corpus considerably (?).

(D1598) A Main Problem: Look at various Lrng methods (≥ 141, ANL, ... )  
 See just how they do "Summarizing" — Maybe/try "GA" (as I solved InputMuxer).  
 I guess "Summarizing" is "Updating" via almost the same. Except maybe for OSL!

In a better GA, "InputMuxer" problem: We seem to update as soon as we get a somewhat good code. This is the "coding & recoding" idea.

.01 | 09.40 Nice soln. to Defining  $C \equiv A \cdot B$ : I was worried because  $C \equiv A \cdot B$  implied that one would have to convert all  $A \cdot B$  pair to  $C$  for  $t$ . defn. to work, & in the real world, a fair fraction of  $t$ .  $A \cdot B$  pairs would be random adjacencies of  $A$  &  $B$ .

.04 A v.g., "correct" way to deal w.  $t$ . problem: When we define  $C \equiv A \cdot B$ , we do this not for 1 code string, but for all code strings: So  $t$ . total  $\uparrow$  in pc for  $t$ . entire set of all codes, is what we are concerned w..

.09 [ I'm not certain whether this idea is really new, but it really craves TDFE did not apply it to  $t$ .  $\geq 141$  problem! ]

.10 Anyway if  $A$  has occurred  $N_A$  times;  $B$ ,  $N_B$  times;  $C$ ,  $N_C$  times:  $\begin{bmatrix} N_C \leq N_A \\ N_C \leq N_B \end{bmatrix}$   
Then it was code  $t$ . string using  $n_A$  A's;  $n_B$  B's;  $n_C$  C's;  $t$ . pc of ~~the code~~  $C_C$   
 $t$ . code is  $\left(\frac{N_A}{N}\right)^{n_A} \cdot \left(\frac{N_B}{N}\right)^{n_B} \cdot \left(\frac{N_C}{N}\right)^{n_C} \cdot C_C$  |  $C_C$  is  $t$ . pc of defining  $C$ .  
 $n_A + n_B + n_C = N$  also

If  $C_C$  were constant for all  $N_A, N_B$  values (it is not) then  $t$ . total pc of  $t$ . corpus using

.20 Plus data would be  
upper limit  $n_C = N_C$   $\left[ \left(\frac{N_A - n_C}{N}\right)^{N_A - n_C} \cdot \left(\frac{N_B - n_C}{N}\right)^{N_B - n_C} \cdot \left(\frac{N_C}{N}\right)^{n_C} \cdot C_C \right]$   
 $n_C = 0$

If  $C_C$  were constant (it is not), we could take it out of  $t$ .  $\Sigma$ .

If .20 is slightly  $>$   $\left(\frac{N_A}{N}\right)^{N_A} \cdot \left(\frac{N_B}{N}\right)^{N_B}$  then  $C \equiv A \cdot B$  is a useful defn.

The "summand" in .20 has a peak for a certain  $C_C$ . (98.0% to % 98.20. Peak max)

on " $C$ " This is usually easy to do:  $N_A \rightarrow N_A + 1$ ,  $N_B \rightarrow N_B + 1$ ,  $N_C \rightarrow N_C + 1$ :

$t$ . formula of .20 is inexact, anyway! I think Sol 64b explains it in detail - also how to define  $D = A \cdot C$ , say. - its pc of definition.

.29 I had Plot that defining  $D$  would be a made easier if I first had narrowed down how many C's to define,  $t$ . passages: I'm now not so sure it makes much difference - i.e. on can define  $C$  &  $D$ : then one has many many passages. Here, one is mainly interested in prediction, so only  $t$ . passages involving  $C$  augmented corpus are critical!  $N \cdot B$

.33 One sums over all other / passages.  $\rightarrow 142.01$

It should be possl. to compute  $t$ . probty of  $t$ . next bit being  $A$ s or  $B$ s, out. basis of  $N_A, N_B$  &  $N_C$  only - sums over all passages.

Also: Can one ask: what's  $t$ . probty of  $t$ . particular  $A \cdot B$  pair being really  $\leq C$ ? Yes!

To find  $t$ . answer: take total wt. of all passages of  $t$ . data (including  $t$ . passage that does not define  $C$ ) that pass as " $A, B$ "; compare w. total wt. of passages that code it as " $C$ "

D 492 TM: 2/141

Great Breakthru: Serial or ll codes either or both are least for defining "data" (07)

01: 141.33 Rept. 29 for Defining "D = BC", say: It is probly poss. to define the "state" of the system after C has been defined, so one has all the params needed to define D = BC, & to ~~def~~ determine pc's of ~~the~~ contents of corpus in volving D.. — Those params may be approximate — i.e. They may give just a approx. results.

In fact, it may be poss. to characterize the state of the system of 141.10 w. a few params & make pc of a C = A.B a func. of these params.

07 : 141.04 -.09 Maybe a GREAT BREAKTHRU for me: either I ~~never~~ realized it or I forgot it! Anyway: The idea is that Both serial & parallel codes can contribute to paying the cost of coding a definition. The idea is that the definition should ~~be~~  $\uparrow$  the pc. of the entire set of codes. A very simple example is the coding of a Bernoulli symbol that has proby  $\alpha$  by an "equiv. code symbol" of pc = exactly  $\alpha$  (no need for "Arith. Coding"). All symbols are coded by a "chunk" of  $k$  bits ( $k \rightarrow \infty$  in final proof). An alphabet item of  $pc = \alpha$  is coded by  $(2^k \cdot \alpha)$  of the  $2^k$  chunks of length  $k$ .

(SN) T. Mix & Corpus Theorem seems related to this "Great Breakthru in Z141" — But it may be a counter example! — On second thought, they seem much different.

In Z141, we have lots of codes in ||: In Mix Corp Thm, we have the "2 part codes": Part 1 is the "definition" of a p.d. or Machine — corresponding to the defn. of a chunk in Z141; 4th in Mix Corp Thm: We take the product of pc's of various objects & parts/objects. In Z141 we add the "pc's": So they really look much different! No counter example!

Def MCT

1) "T. & American machine" is a 2 way approach to OSL problem is one extreme example.  
2) How can we combine resp "2" parts of system w/ "C" & study to D can be determined by  
3) Increased "1" sometimes all 5 values & are approx & loss mpk info. — T. only way to do this. approach is to use an exact non-loss code. Usually we code w. 2 parts: part 1 followed by part 2 pure random sequence  
11:11 → 11:19  
20 → 25  
27 No Nth variant count. 22 N.Y. R.I  
• 1395 "NU NT" P/A GA  
~ 23' GA  
4.33 L 04/2000

Previous ref exists in 1998

This may be f. idea: There is some domain (area of expertise) where P's system has lots of domain theory. (≡ Good models, that often work).

It is given a problem — which it solves, using its models: Hvr, it takes a fair amt. of CC to solve P's problem. (say lots of reasoning & search, etc).

So It generalizes this soln., so a soln. will be available in future at lower CC. I guess f. purpose here, of GENZU, is to save CC, rather than induction.

The "Genzu" is done by examining the problem & its soln. How can we relax/specify of the problem & still have the same "soln. trace" be valid? — so we can easily decide what the soln. is? For INV. probs. we can always tell easily if a proposed soln. is correct. For OZ probs: well we can get f. Genz.

By defn. of "INV" problem

Passing phase call to Hyp Solver

New invasions → M.M.M.

Ford: Renaissance Reason 838

Hyp Solver

\$62M

now in 2 part of hyp solver

An example: A proof that a RT Δ has certain properties: looking at proof, it is clear that we can remove the "RT" constraint & the proof still works. So we want to turn to be true of all Δ.

I think that an imp. part of EBL is that a "explanation" of what's observed helps tell how the observation (problem soln. regularity...) can be general.

- Some Refs: AAA 1990 conf vol II
- 19 803 contrasts EBL & SBL
- 821 EBL is "Chaining"

Similarity

Mentions Wolf's induction w/ 2 R's: looks reasonable!

Proc 6th workshop ML: 1989; p 2-4 Pat Langley: this is somewhat clear.

In .20 EBL is said to need small SSZ because of all the constraints in its Models; SBL has few constraints so it needs larger SSZ.

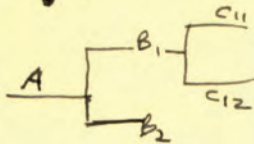
T. paper, 20 Tries to mix the 2 ideas by using a system w. constraints stronger than SBL & weaker than EBL.

In .22 Langley refers to Wolf's work on (ZIL) lang as starting w. SBL then moving toward EBL as the domain theory was built up. So putting Wolf's work, my most recent work on ZIL (1990) in Mxd corpus lang form into clarity EBL, SBL & Mxd corpus lang. The thesis that as pc's approach oil, we approach oil.

Conjecture: That while certain pc → oil (after a large SSZ) other pc's in the problem area do not. (These could be the decision to try a linear eq. approx. : Then solving the linear equ. would use pc's of 1 & 0).

01: 138.40!

code  
to



we need  $|A| + |B_1| + |C_{11}| + |C_{12}| + |B_2|$  + punctuation  
to code the entire corpus.

to code w.o. this graph we'd need  $(|A| + |B_1| + |C_{11}|) + (|B_1| + |C_{11}|) + |B_2|$  } using ordinary coding of non-disjoint sets.

and  ~~$|A| + |B_1|$~~   $(|A + B_1 + C_{11}|) + (|A + B_1 + C_{12}|) + (|A + B_2|)$  to code the corpus w.o.

using common codes.

Could we get the same results w.o. the tree structure? i.e. within A, we have B1.

C11 & C12 both refer to B1 (within C) before we start C11 or C12.

w.o. the tree structure we could get more general sharing of common resources than we could get w. the tree structure.

A possible advantage of the tree str. (when we can use it) is that it reduces the cost of referring to sections of relevant code. w.o. the tree structure, we need

to tell what parts of A are relevant to our term. Using the tree str., this is unnecessary, ~~or~~ or costs less.

20

**SN** In the soln. of O2 probs via Lsach, we need not have cond. pd: A around. pd. is fine! Whoops false! — we do need a cond. pd. — That looks at the H.C. problem & decides which H.C. method is best for it. I thought since the H.C. method "looks at the problem" that a cond. p.d. would be unnecessary. — But NOT SO!



① Sol 78 Thm 3: for radix  $> 2$ : incons: results obtained!

② SOY Problem: 2 d.f.s of mean 0, var 1, v.s. 1 d.f. of mean 1 & var  $\phi$ .

.07

③ ~~Probability~~ Prob by distribns. w. <sup>cc</sup> ~~constant~~; <sup>apparently</sup> Related to Chess problem, Oz problem w.  
cc =  $f(t)$  ( $R(t)$  is + funct), Japan's ~~machine~~ <sup>machine</sup> (149.37 - 150.08); see <sup>TM</sup> (99) 2.25 for possl. Break thro!

④ Summarizational Problem: What are Good ways to summarize t. Current P.B. (is R<sub>2</sub> related to 3) (.02)  
How to balance by PC off-entire corpus v.s. cc of obtaining PC's.  
?.

Faint, illegible text at the top of the page, possibly bleed-through from the reverse side.

Vertical text on the left side, possibly bleed-through or a margin note, including some numbers and small text.

**WARRANT**  
**CORRECTION**  
**HONOR**  
**OFFICE**

**COTTON FIBER CONTENT**

Main body of faint, illegible text, likely bleed-through from the reverse side of the document.



Kate PP359  
i 358  
one in  
inverted  
order

01: 357.40: 11) An imp. Q: Can I get TM to work on TM2 ~~hyperprobs~~  
in a natural way? — i.e. apply info on prev. optimization problems to its own

program optimization probs. w.o. being told to do so?

12) Related to 11) Would TM give by PC's to methods that were of small CC?  
(in this packet)

See 358.01 (Ideas) for some discussion. Would it try to revise old prob. solns.

so as to find methods of lower CC, hyper PC? (one very fast soln. is a numerical A.H. soln. ...  
but it will have very low PC ... probably too low to be within CC of TM) — precisions should be always kept

so that A.H. solns. are essentially infeasible.)

13) I'd like TM able to "Memorize" complex solns. to certain problems — they  
later be able to (usefully) try to "factor" to solns. to give them more PCs  
i to be able to usefully use the "new" factored concs.

.18

14) There is t. possy that ~~the~~ search strategy can remain constant. That I don't  
even need  $\geq 141$ ! This idea is related to (0.01), 357.24, (1.16'91 ... I don't get this!)  
↳ 50% soln.



related to "concepts": 7588 Index 3.05

.01; 247.40 ⑦ I was concerned about TM not automatically developing methods to work problems that had small CC & give by pc's to them.

I think this is automatic! ~~automatic~~ Small CC methods will be more likely to solve problems.

One can try ~~more~~ more of them per unit time. They will end up being "popular" & get assigned by pc's. [This idea needs more careful thought, but be sure to write down & reasoning on this... it's a bit crazy in my mind at present!]

Is this an example (pos. neg.?!): TM find multiplication to be useful operation, but it does take a lot of time (he does it by repeated addition). Would he be motivated to

try to find a faster way to do it? would he (ap) have any idea that this particular "commensur. of operations" might be done more rapidly?

.12 → Related Q: TM solves prob. w. CB = 1 min; w. CB = 5 Min it obtains <sup>work</sup> factor soln. but w. <sup>much</sup> lower PC. } more rapidly?

How, this soln. at reasonable PC level is now available for "SSZ=1" learning for a new problem like

.14 ~~the first one!~~ (Symbolic) soln.

Seems to me that whether or not Z41 is up to it, that I should be able to get ~~it~~

~~CBI~~ to react properly to 4. problems of .12-.13. (Pb Z41 does seem to be the right approach!)

T8 → 1.18.91 Example of "foc" soln. of very low PC: Given on eqn. w. low pc. coeffs.  
T. no. that solves it is a low pc., low CC soln.

# 6988 TM: RADTS

(cont.)  
 .01: 360.909) That would be useful elsewhere. T. Considerations of 360.30 (Re: GPS) suggest that  
 there will be lots of "transfer learning" to other domains!

→ 10) A big Q is: what order of concept & position should be used. Should GPS be learned along w. / algebra, or should algebra be learned partly by memorization — you have GPS be discovered as an "integrating principle." → (19)

**Very Imp.**  
 Old view was ①: perhaps very diff. to write by equiv. New possy. is ②: see ①9 item 15

11) One way to learn GPS is that one has at 1. outset, this idea of "difference" betw. t. very things & re: one's goal: perhaps memorize this/or use it as a primitive.

→ 12) This seems to be departing from t. "Ref & dirty" idea of 357.08, 128 — lack of criticality of t. sug. seq. → (19)

①9 → 13) From ⑩ & ⑫ I guess that I'd want to start w. lots of "memorization". Learn to do lots of tasks that could be learned by GPS. Then after these tasks are learned, have it discover GPS to reduce t. pc's of all previous solns. To do this would for me, be v.g., because it would involve some "higher order learning", perhaps. It would seem to involve factoring t. original solns. — a having all these solns. have a "common factor", i.e. GPS use could be a common way to solve them.

14) An interesting Q: to what extent would TM have to have experience w. problems like "discovering GPS", before it could discover GPS? again, I think the problem may be transferring into from one level (TM<sub>1</sub>) to another level (TM<sub>2</sub>).

15) "Memorization" may mean: something is learned but t. CC given for t. task was not large enuf for TM to discover t. useful factors of t. concs. that were adequate for t. soln.

[I probably meant "CB" "Comp. And → CC =

16) Actually, GPS doesn't have to be discovered all at once. Parts of GPS can be discovered & be useful. E.g. Certain Ob-ops are found to be useful. later, it's part of GPS, the Ob becomes an observed difference, & t. obs. op's t. operator that reduces that difference. At first, TM discovers this ob-op in a particular domain of probs, then later TM discovers this sort of thing is common to many domains of probs, so it can't cover all pc by using this/factor in t. code.

6.10.88 TM: 586: Time sharing. w. 1 CPU,  $\sum p_i \leq p_{cr}$  w.  $\geq 1$  CPU.

A possibl. way to do Time share  $\Sigma$  tasks w. many // machines: First, if we only had 1 machine, time share it by spending a fixed time  $\delta$  on each task, indep of its  $p_i$  — but r. frequency  $w_i$  which one spends to. Pick  $\delta \propto p_i$  of that task. We arrange

so that  $\delta \gg t$ : time needed to switch betw. tasks, so there is  $\approx 0$  cc lost per u

"2d min. instr. of T.S":   
 There's a problem is that this gives us an effective cutoff of low  $p_i$  cands. There are a very large no. of "low  $p_i$  cands". Cutoff depends on  $\delta$ . Some wouldn't get runs w.  $P_i \cdot p_i \cdot T \leq B$ ; so  $\Rightarrow T \uparrow$ , w.  $\uparrow$  no. of cands tried.

Say  $t$ : time quantum's  $\delta$ : each task has a counter. Every  $\delta$  interval, each task has its counter incremented. The "threshold" for each task is  $(p_i \cdot p_i \cdot T) \cdot \delta^{-1}$ . Each task has an order number.

After CPU has worked on task no.  $i$ , it looks at task  $i+1$  to see if its counter is  $>$  its threshold. If it is, it increments its counter by its threshold value, & it works on that task for time  $\delta$  — then goes to next task & looks at (counter - threshold), etc. it goes to next task & looks at (counter - threshold), etc.

One nice feature is that if a task has been overdue "for a while, before its work on; then  $P_i$  is overdue amt. makes it work on earlier next time. Hvr, how much time is spent updating counters? — More use  $w_i$  of them.

To do this time shared  $\sum w_i > 1$  machine } we could have a bunch of CPU's } Each does a task in // w. others. When done, it looks for the next task it could find that's "over due" (i.e.  $(\text{threshold}) > (\text{counter})$ ). Hvr, all CPU's finish  $\delta$  tasks simultly.

The problem here is how memory is shared betw. CPU's & how the state of a task (after a CPU has just finished its  $\delta$  unit) is transferred to next CPU to work on it.

One not-so-efficient approach: Assign each CPU a set of initial tasks. Do this  $\exists t_i \leq p_{cr}$  is a bott. since for all CPU's. Each CPU has its own RAM.

A trouble in .28 & .29: Sometimes, work a CPU will look for a task to work on but find there are none overdue! Another possy is that we never exhaust capacity of CPUs, so as time goes on, the "overdue" level of many tasks gets larger & larger! One way to deal w. this: the threshold of over tasks =  $\frac{\delta}{p_i}$ . The value of  $\delta$  is being constantly corrected so that

.30 & .31 occur infrequently: I. e. ~~that~~ kind of event pushes  $\delta$  in a direction to reduce the likelihood of that kind of event.

to make .28 more efficient: when a gn. CPU is spending too much time on its tasks (since  $\sum p_i \leq p_{cr}$  isn't as large as other CPU's) then one or more tasks may be transferred to it from other CPU's that are spending too little time on their tasks. We probably don't have to consider transferring unless  $\sum p_i$  for 2 CPU's differs by a factor of  $> 2$ , say.

See also 1 page note! D2785 on use of disc for (even,

Summary Machine (also maybe 146:07) "Updating"

01:40:40 T. Summary Machine is any usable form of (conditional) P.D. As a general P.D., we can get PC values of ↑ accuracy from it, but at greater CC. In general, the relationship of CC to PC precision is "unknowable", but we can probably approximate it w. finite CB.

T. main idea is, we want best Lsrch results. — i.e. good PC & small CC.

We could try ~~using~~ varying amounts of CC (a scanning ~~and~~ optimizing sense) & ~~use~~ use th. variation in times to solve INV probs, to select best CC to use.

.07 The soln. to: "Summary Machine" that I had in "2 kinds..." ~~is~~ seems a.c. accept

.08 OSL. We may not need the "OSL" part of the P.D. for Lsrch(?). → (152.01)

So we have 2 applies of the "Grand P.D." ① for Lsrch problems ② for output to RW — to make decisions, ect. I guess I want to optimize it for Lsrch. so if Lsrch doesn't need OSL, then it should be easy! (well: I'm not so sure Lsrch doesn't need the OSL part!)

17 A major problem in "Summary": We want to "improve the P.D.": we want ↑ PC for the entire corpus, but we are also concerned w. CC of the PC values. The "goodness" of a P.D. has 2 components PC of corpus & CC of outputs: So we can certainly partially order the P.D.'s:

— but we need a "linear ordering" to guide "updating" (= "improvement") of the PC part.

Entire corpus

This Q. of just how best to do the "updating" seems to be a major problem / bottleneck in the entire project.

Actually, what we are really interested in is high & low CC for ~~the~~ extensions of the <sup>(partial)</sup> corpus. perhaps Low Cost  $\frac{CC}{PC}$  would be what we want.

Consider the actual applies of the "updating": We have the mixed corpus:

(Neglecting OSL) we find 10 shortest items of the corpus, we are then given a "partial object".

We are required to find likely continuations of that "partial object" w.r.t. the mixed corpus.

.27 ~~Someday~~ A thing we want to do: Given the 10 short mixed corpus items, to ~~do~~ devise a stochastic operator, ⇒ given any <sup>new</sup> (partial) object, the output is (in some form) a P.D. of continuations of that partial object.

An "update" is a state of TM that makes .27-29 relatively easy to do.

If, in .27-29, the new partial object is a Big Question, it will be necessary to spend much CC analyzing that Q.

.34 Say C is the (Mixed) corpus; B a partial object (= a fragment of C). To get to Best 10 codes for CB; we may first use sequential coding of B, then "coding & recoding" of CB.

T. General "Mixed corpus" updating problem looks pretty much the same whether .34+.36 is a mixed corpus or a sequential corpus. For each augmentation of the corpus, we want to get the best set of codes (or output) for the augmented corpus, in the available time. → (39) 2.01

dynamics associates

$$\left( \begin{array}{l} A \supset B \rightarrow \sim B \supset \sim A \\ \sim P \supset Q \rightarrow \sim Q \supset \sim P \end{array} \right)$$

B → ∼B means toggling of "i"  
 so B → ∼B implies ∼B → B.

"Human Problem Solving"

01:362.1.40! (However, on the whole, the data in HPS seems interesting and may prove very useful). I can probably make much <sup>(more complex)</sup> better models of their human S's than N's did, because of their use of CJS!

In the data of HPS pp 593-6, there is improvement in problem solving ability as the task becomes continuous. So this is a key exp. of interest.

So this may be a very feasible way to do the exps! HPS is one source of both problems & data on how humans solve them.

Also, I could use Pat Langley's probs (i.e. data on ?) on discovery of scientific laws.

In the case of HPS examples: I watch this human solve the problem & I postulate a certain set of principles or concepts that would give such a soln. Note that normally, both problem solving & updating are occurring during <sup>human</sup> Phase/solns so I will have to use a TM model that Time shares on updating along w. regular work on trials.

dynamics associates



01.362.40

Perhaps a main problem is to find a set of parallel problems, that are nominally in different domains, but have similar structures, so I can try "Transfer learning" as soon as poss. I'd like to ~~try~~ <sup>try</sup> this "Meta learning" as soon as poss. Ideally, w.o. use of a TM formalism.

N.B. "Hum. prob. solvng" may have some useful ideas. I have avoided reading their stuff because of their very bad methodological errors — but it's likely that I can get ~~useful~~ <sup>usable</sup> problems from them.

A poss. source of n probs: 7. list of ~~some~~ <sup>11</sup> problems they used for "GPS & Generality":

Another poss. is the probs used on Hummer in "Hum. prob. solvng": I think they did a lot w. "~~the~~ <sup>cryptarithmic</sup> cryptarithmic".

Def

The "Missionary & cannibals" prob. ~~(MCP)~~ <sup>(MCP)</sup>

One way a human might solve it: play around w. <sup>for conj. purposes</sup> ~~conjunctions~~, make up "interlocking" ~~problems~~ <sup>problems</sup>. Keep eyes open for relevant threads or conjs.

This is probably true for all 11 problems! — so good poss. of transfer! Another poss. is to give TM 11 probs <sup>of all 11 types</sup> to work on "simultly" using Time sharing.

In Hum. prob solvng: 3 ~~distinct~~ <sup>distinct</sup> Domains:

- |                 |          |           |                          |
|-----------------|----------|-----------|--------------------------|
| 1) Cryptarithm. | 2) Logic | 3) Chess. | "Group" <del>group</del> |
| 20pp            | 15-2pp   | 125pp     |                          |

Also, do ~~the~~ <sup>elementary</sup> Alg. learning (how to solve linear eqns). My present idea is to have TM learn to solve each kind of problem in a general "way as poss. — (by using large esp. CJS in 4. tng. seq). Then, by having TM note similarities in different probs, he devises ~~the~~ <sup>the</sup> best ways to have ~~them~~ <sup>them</sup> discovered to old solns. & to discover solns to new probs

?  $\rightarrow$  Tic Tac Toe as an elementary problem (?)

Def In (GPS) (Hum. Prob. Solvng), there are short passages on learning, that could be of interest. ~~the~~ (See index) p 593-6 are of interest, because

they give a temporal order for acquisition of imp. learning concepts; (perhaps unfortunately) N.B. how GPS is in mind as a model, so it probably loses ~~some~~ <sup>some</sup> ~~of~~ <sup>of</sup> info

6888 TM: RADTS Rudimentary Alg. Seq.

01.303 to 1) Say <sup>top</sup> Goal is soln. of linear equs. ∴ let us work backwards:

first define top goal: T. eqn. to be solved is a linear eq. but it is in a form such that

various standard operations have to be done to get it into y. form  $ax+b=0$ .

T. eq. is given in RPN or conventional alg. notation (e.g.  $(3x+7)(2x-3)=5$ ) — hrr: see (19)

005 2) So TM. wants to "simplify" t. eqn: — for it into a form that he knows he can solve.

TM has idea that a simpler looking eqn is usually closer to soln. than a complex looking one.

A TM has some natural measures of complexity: e.g. no. of terms

(perhaps certain terms even give more complexity than others).

12 3) In "simplifying" TM will start w. primitive operations — they leave this... (4 tuples)

of operations that are v.g. These "operations" are in the "ob. of" algebra —

i.e. they involve operations that are contingent on observations. We can have variety of:

sequs of obs; sequs of ops; sequs of ob, op mixtures.

19 4) Somehow we have to convey to TM, that the problem is to find  $x \ni$  that eqn. is true.

— We can use floating pt. & have an <sup>external</sup> error criterion, so t. eq. doesn't have

to be solved exactly. — so almost always, we can assume infinite precision & TM can "watch" t. eq. that checks its cond. solns.

5) After learning to solve equs. w. numerical roots, it should be able to do it for literal roots.

6) Some info we expect TM to have in order to solve (a) literal roots.

a) Some possibly useful info: that certain x forms of t. eqn. don't lose info, so they are acquir. to it, & one can solve t. new eqn. instead of using old x forms. — (But t. original may be special).

b) Some operations — like discarding old equiv. x forms — can save RAM (2. ↓ size of subspaces; ↓ # of solns)

c)

7) It may be that t. hours of .05 — .18 (6, 8) aren't normally acquired until TM is quite "Mature": that t. solving of .01 usually uses a lot of "Memorization" —

& that later these memorized things aren't factored until later, when

TM acquires t. hours of .05 — .18.

8) A way to acquire t. hours of .05 — .18! Discover GPS; find various "differences" & assoc ops that reduce them (See op. obj. of (12) Attempts to get set of "orn"

differences & o/a set of "orn" operators that reduce t. differences indirectly.

With Reg GPS Model, TM will be able to work other GPS problems (like Symbolic Integration,

Unitary & Conables, Tower of Hanoi. (But, probably, much better than GPS).

9) Why (.30) is very encouraging: I was afraid that leaving to do simple algebra

wouldn't have an effect assoc. probs. to be in "learning" in t. sense of learning in pt. hours



6688 TM: General Path:

# RADTS

The system as of now ( $\approx$  586) :

- .01 1) The initial ~~set of~~ primitive set of concs need be universal, but Phil's all. We can put all h. into we like into Y. prim. set. If we don't have enough or it's wrong emphasis, we can correct Phil later in the try seq. (or, just add to primitive set.)
- 2) The search algm. is fixed. We can do  $T \leftarrow 2T$  or find a having in various ways, but this is not very imp. diffence. Hvr., see 359.18
- .08 3) Try seq: Also not very critical. If we make mistakes these can be corrected by later fixing.  $\rightarrow$  (128)  $\downarrow$
- 4) One goal is to get machine ~~to~~ work + smart enough so it can usefully work on under sets of problems (using time sharing), (See 247.17-40 for some vaguely related discn.  $\rightarrow$  also see (20) !)
- 5) I think that as one adds probs to T. S., TM keeps working on old probs, but there is some cut off criterion for branches based on total time worked on it  $\approx$  ~~present~~ present PC.
- .20 6) It may be that if, even w. a rather ~~inert~~ inert TM, <sup>perhaps even in infant TM, use 2 cross-multiples!</sup> one always gets about 5 problems to work on simultly — but one doesn't have to be too careful in devising try seqs — just so long as one gives it lots of time to work on early probs (even if they are easy probs, a large C will keep machine working on <sup>posbl.</sup> improvements).
- .24 7) I don't know just how nery to TM<sub>1</sub> = TM<sub>2</sub> thing is. If it is, I could simply have T.M. time share  $\approx$  of its time on TM<sub>2</sub> type probs — so ~~up~~ "up down" would be going on continuously.
- .28 8) Because of lack of criticality of the try seq., I could just take (any old) algebra problem, then work backward in a "conceptual pathway" to ~~some~~ a suitable primitive set of concs. Try this for several problems & see at what points the "conceptual pathways" joint in their way to the PRIMITIVES   
 [ 1998: thereby giving a strong set of near primitive concs. ]
- 9) The try seq. should be first done in English! Perhaps first see if human could follow it. Then get CJS's for a machine: then insert needed concepts & examples to get CJS down to acceptable sizes.
- [6898] 10) A ~~posbl.~~ Bottleneck that I haven't met about too much (the other is not; 111.01) is Making up probs that use custom concs. in their solns. Also how many examples of each prob. to give. Using enough equif. dis. in numbers in probs can elim. possy of A.M. solns. Perhaps enough dis. can give effectively large set!

359.15 on back of Phil's sheet

360.15 on back of Phil's sheet

.01: **88TM 362.2.90**: **IBid 357.01-362.2.90** is an examination of a maximally simple (as of that time) TM & its TSQ. At a time two problems that I was concerned w. - that I feel I now understand much better: (1)  $TM_1 = TM_2$  structure, system (2) How ideas in one area of Technology get used in an other area. Much of Ruth 1988 discn. <sup>was</sup> involved w. these Q's:

inside + pte.  
last part -  
a. extreme solution.  
= 2.1992.

The  $TM_1 = TM_2$  prob is treated by regarding updating as a ordinary OZ problem. T. problem of relative wts. of solns to OZ probs (i.e. renders mixed w. INV probs) does not occur in t. original form. T. all over goal is simply to find best codes for entire corpus of Inv. & OZ probs.

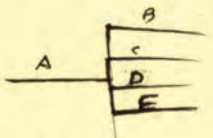
.11 There is a "rub", hvr. **Certain Probs are more imp. than others**, so I (think) I'd rather ↑ pc's in certain parts of t. Corps rather than others. Could this be regarded as a special case of t. General Problem of what parts of R.W. to integrate into each Corpus? - a/c to problem "what direction should MoRe go in?"

A very long time series or a large sub corpus: → these could find more ↓ in best/cc than short sections of corpus. Sub corpus help me solve what I

If I don't want an excessively general I.M., but one that can help solve what I Rm have imp't probs, — then I can give T.M. these wts. → It may be possible for TM to assign these wts. as a result of an induction problem; at first I had

.20 In General "TM<sub>2</sub>" (T. updating prob) will have t. problem of deciding how much (relative) time to spend on improving f. pc's of various parts of t. corpus. Say "A" is the

TM wts of many probs — then eventually TM sees pattern in the assignment set



Common codes, i.e. B, C, D, E are codes for separate finite objects or parts of finite or infinite objects. I can ↑ pc of entire corpus by working on any of t. B thru E problems, or work on them jointly,

.26 to ↑ pc of A.

In General is .30-.26 a variety of "What to work on Next" problem? It's a general. of t. w/ w/ o/ n problem! Instead of 2 yes/no soln, we have a "soft" Score. ~~of total pc.~~ of ↑ of total pc.

.30 Hvr. another  $\pm$  **U.G. idea** from [ibid]; In learning to solve eqns., I originally did this by wanting an  $x \rightarrow 3x+7=2$ . The idea that one could pretend one knew

Mem  
I found 10:10A  
Doris  
492 5081  
2335 Mass Ave

.32 What x was, & manipulate expressions containing it, was a major hvr. breakthrough. Then, learning a growing set of x-funs, <sup>for Manipulation</sup> added to heuristic power.

Would you prefer of .30-.32 be EBS?

Now .30 is much different from t. TSQ's that I've written for Algebra. I'd like to be able to get TM to do .30 ff, because it's close to how I think about problems & is (on some level) easier to teach TM.

For .30, I am (apparently) able to do t. TSQ "in English" & so I should be able to put it in formal form.

Tyland on ASP N-6.

Some kinds of Induction/updates:

1) Numerical T.S. i Some 'leading into' ridges! ~~that~~ start w. some w. exact rules.

$X(t)$  can be func of  $t$  & (near by) <sup>int</sup> previous values of  $X(t)$ , for these TS's:

~~1-1~~  $X(t) = f(t, X(t-1), \dots, X(t-k))$  starting trials w. small  $k$ .

To find  $f(\dots)$  one must step thro  $t$ -space of all possl. func's, but  $t$ . pc of soln is probably too small.

2) If we add noise to 1),  $t$ . soln. B to same w. additional pc of error corrections,

3) Soln of equs! Solve  $x=2$ ;  $x=3t$ ;  $x=3*(1+t) \dots$  etc,  
 - is this same as ANK?

"  $x+1=2$ ! This could be solved by <sup>direct</sup> search if  $x$  is an integer. — in fact, true of most integer solns.

20 **SN** Advantage of several examples of low precision v.s. **1** example of high precision! Search is faster, but "sharpness" of soln. is much less — so several low precision examples is  $\geq$  better TSO than a single high precision example (usually!)

~~20~~ 20 is true for direct numerical search. It is not true if  $t$  search is over functional forms, in which case  $\epsilon$  search time is same in both cases, (usually, w'd rather that  $\epsilon$  search be over functional forms)

1) Some Imp't kinds of things I want TM to be able to do:

- 1. It watches a problem being solved — say a <sup>specific</sup> linear eq. . From this, it's able to genz. & solve any linear eq.
- 2. It watches an eq. solved by "invertible substitution" (say a linear eq.)  
 From this it can genz. to quadratic eqs (once  $\sqrt{\phantom{x}}$  is known as logic)  
 Then " " Cubic " < "  $\int$  " " " )

∴  
 ∑ 1112  
 ∑ ↑  
 A1 ↑

3. Given  $\sqrt{-1} = i$ , when asked what  $(3+7i)^2$  is — what would it equal?  
 $(3+7i)^2 \times (4+2i) = ?$  etc: — would it be able to generalize complex nos.?  
 — Square roots; complex primes etc. Complex solns of quadratic eqns.

Complex cube roots?  $\Rightarrow (a+bi)^3 = \left. \begin{aligned} & a^3 + 3a^2bi + 3a \cdot b^2i^2 + b^3i^3 \\ & = a^3 + 3a^2bi - 3ab^2 - b^3i \\ & = (a^3 - 3ab^2) + (3a^2b - b^3)i \end{aligned} \right\} \begin{array}{l} \geq x + iy \\ \text{where } x, y \text{ known} \end{array}$

$2(a^2 - 3b^2) = 2a^2 - 3a^2b^2 = x$   
 $b(3a^2 - b^2) = 3a^2b - b^3 = y$

Given 2 geometric interpretation of complex nos. in a plane! As  $Ae^{-it}$  representation.

could it then do roots, etc? — Could I lead it to dis cover various properties of complex nos.

4. on (x) EPL: When it solves a prob. or sees one solved! It might be able to generalize & substantly if it has ideas about what ~~a~~ problems will occur in future.

Say it has a "Problem pool" — problems are ways being added (by me) to this pool.  
 TM looks at these problems, tries to categorize them as well as possible, so it has ideas on how to solve them. Then, when it solves a new problem, it will be able to know how to genz. its soln. (i.e. in diff. direction of one or more of these unsolved probs.)  
→ There are many ways to genz. .... so ... it can gather ideas from its unsolved prob. pool.

Also, w/rt. solved probs; — it can genz. a soln to a new prob. so it solves older solved probs. in different (perhaps, but not necessarily, better) ways.

In General, this "Problem pool" can be worked on w/o actually solving it. Gary 5715

Probs: just look for regys. in the problems, ways to compress this pool.

Actually, this compression is usually a major part of prob. solving! The problem statement may be regarded as a "partial object" : prob. soln = "whole object".

These partial objects are always part of a corpus that needs to be compressed —  
 So we do want to find similarities, regularities, analogies & find them.

Various problems.

"Hints": "Hints": I kind of hint to solve prob A, is prob B (which is easier but uses ideas that can be genz. to solve A)

Handwritten text at the top of the page, possibly a title or header.

Handwritten text in the upper middle section of the page.

Handwritten text in the middle section of the page.

Handwritten text in the lower middle section of the page.

Handwritten text at the bottom of the page.

Vertical handwritten text on the right side of the page, possibly a list or index.

Large, faint watermark or stamp in the center of the page, containing the words "ALPHABET" and "COMPARABLE".

More imp. thing I want TM to be able to do!

In symbolic integration; TM should be able to see a particular trick being used in a particular problem; then generalize so it can use this trick much more generally. There are an enormous no. of tricks ~~of~~ of this sort that have

been incorporated into some like Macsyma, Mathematica, Maple - There may be books on this - but also, much info may be proprietary.

08

T. idea of invariable substitution can be used for both exp. solving & symbolic integration.

Some Q's:

1) T. "Grand P.D." can be in several forms or mixes of them:

- a) We input a string (condition); Output is  $\{p.c.\}$  of (outputs strings) i. Aff. p.c. order; ~~or~~.
- b) We may have to put

b) Output could be a "most likely" string or number,  $i^2$  or  $or^2$ . ~~or~~ (i.e. normal dist.) or description of description of some other p.d. T. p.d. can be multimodal - in which cases ordering outputs by p.c. order could be diff.

c) Output could be a ~~mathematical~~ Multiple (could) set of outputs w. i. proper probs:

- Not so desirable, but it could arrive at p.d. is a stochastic Grammar - say a C.F.G. There is a recurring problem (solved many times), of putting outputs of some kinds of stoch. Grammars in pc order. - in simplest form, we have a finite set of symbols, each w. its own p.c. To list ~~all~~ <sup>all</sup> finite strings of these symbols in pc order. I think I solved this one, most recently.

d) Q: Could input to "Cond. p.d." be a p.d. of strings? - It certainly could be; but if Q is: e) p.d. has 2 applies: I use to make decisions: outputs of TM to control an Lurch (used internally by TM) does this even ~~if~~ occur in TM?

2) Work out Example of TM in action on a problem of some diffy (like some of b. diff. Genz: probs of 178.01 - 179.08) - or 147.20 - 148

3) In 1 I probably need to work out methods of converting between p.d. type to another (like .18 - .21).

4) SN In "Operator induction": one form: we are given  $\{I_i, O_i\}$  a set of I/O pairs; to induce. (now) On given new  $I_n$ . This is a normal "Finite object set induction problem". If we are allowed to furnish new inputs, ("Experiments") Y. problem seems much diffent. Each experiment has a known (or unknown) until done cc - (may be / zero for some students) This is a "Learning w. a teacher" It is also some of close to problem of Science & wanting to learn about R.W. - so f. prob. is not "well defined".

37

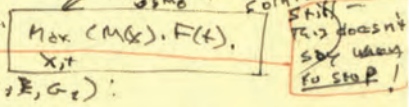
Related problem is Turing's "Learning how to learn": Essentially an optm problem.

Given a known sequence of inputs,  $G$  Inputs:  $G_i$  pairs, to find opt. v.e. Inputs  $G$  (Arr. see 150.07 for probable Soln.)

01: 149.40 <sup>cont</sup> BUT t. problem is actually more difficult:  $G_i$  is for  $(I_0, I_1, \dots, I_i)$ , so we can never repeat a brief situation (which J. makes a big pt. of).

T. main Q here is not whether J is rite or wrong (his rite) but whether I, personally, am interested in solving that problem. It is the problem of a living organism trying to survive & maximally reproduce, in a hostile envt. Certainly an "interesting" problem is perhaps impl. — but I may be able to build what is for me, a useful TM, w.o. solving that problem.

107 Actually, .01 may be solvable as an OZ problem! (hint: it's of the "unsolvable type")  
108 One has to set of math G pairs:  $(I_0, G_0), (I_0, I_1, G_1), \dots, (I_0, I_1, I_2, \dots, I_i, G_i)$ .  
From this, one can view it as a regular OZ problem. One is always restricted in one's input trials (since each next must include all previous inputs). — But I think the standard Leach Program would work. T. ~~Theorem~~ Theorem gives a P.P. on all functions that look at the problem's available trial data, & derive a suitable Hill Climbing Routine.



115 A big Q is how this P.D. "learns" from previous organisms' trial runs.

5) Another very broad Problem: How do I take ideas like 147.30 (solitons - Einsteinium) & put into a form useful for TM, so suitable Leach can be done, etc.?

Anyway: In my own soln. of t. problem, write down what seem to be impl. ideas (not all useful, but likely to include t. key, if sufficient) & use in solving t. problem:

- So: Ideas: (1) assoc. w. some (string) is a (real) number (exactly or approxly) ("numerical value" of string).  
(2) (string a) = (string b) means that strings a & b have some "numerical value":  $a = b$  is an equation.  
(3) It is poss. to change tags to produce new tags. Also combine several tags to produce (or more) tags. (See 5) for some tags).  
(4) If 'a' is a string that has a value, then  $a = a = a$  tag.  
(5) if  $a = b$  &  $c = d$  are tags then  $a + c = b + d$ ,  $a - b = b - d$ ;  $a * c = b * d$  are tags. and if  $c = d$  is not a tag, then  $a / c$  is not a tag. and  $a / c = b / d$ .  
(6) if A substring has a value, then eqns containing that substring will still be true if we substitute b for a in any tag & it will still be a tag.  
(7) if a string is a number then t. value of that string is that number.



There may be other impl. ideas, but this Q can probably be reduced to a smaller set of ideas that will generate 1-7.  
From t. ideas 1-7 as a kind of codes, we could ask TM Q's about strings & numbers, to see how smart he is. So we would quiz TM constantly in this area, before going on to problems like Solving eqns.

Another Q is: what kind of knowledge TM has — is it factual (true/false) or probabilistic? If t. knowledge was obtained by long (induction) it must be probabilistic. But by using many examples or nos. w. many random bits (much "noise") one can get very close to "certainty" — as in ANL

So, we can perhaps get TM to learn about the "facts" of 150.20-40 by a few

### cases of hy precision

Probably best way to do this: Write out details of 150.20-40 in English, then progressively formalize.

Compare in similar process for ANL.

While 150.20-40 can probably be much simplified; shortened, & improved; I should be able to do the formalization of almost any English procedure - not really for back English procedures.

In fact, the real problem: to be able to formalize <sup>substantive</sup> any really well defined English routine.

So consider various kinds of English statements, Q's, definitions, etc.

Answer  
will be a  
Re!  
thesis.

Definitions: For TM are ~~more than~~ (as for humans) are much more than agreement to substitute one string of words for another! They are a statement that the thing defined is of importance & its given a higher than normal ~~pc~~ pc: Also we look for its occurrence in text, so as to

update its pc. So if a "rt. d" has been defined, we look for "rt. d" in

(almost) all texts. We mainly look for them in Geometry & Math text.

→ Then if we find them in other texts, it is of unusually high interest!! So maybe

look for them in all texts.

Fortunately, looking for certain string types is easily parallelizable - so adaptable to human

IT into prog., but also to high capacity machines!

On to use of "numbers": If I just use integers I will be limited to (classical) number

theory: which is interesting - but I'm not at all knowledgeable in that area.

If I allow reals (which I am familiar with), there are various difficulties of approx, Roundoff error, etc

How, I may be able to get TM fairly far before these difficulties become serious!

→ At Real Time, I may be able to explain to TM what the difficulties are & how to get around

them!



b1

(Much previous work/role): Usually, the way it occurs (?): We have summarized corpus, C up to now (w/o OSL): We get a new Partial object, A —, So corpus is now CA — We want a Summary for CA (or "update" of summary for C): The way we usually do this: We look for previous occurrences of A or part of A in ~~the~~ previous corpus, C. (This is if we are doing just a ZIP, simple k-top. dedup).

So, in this forum, the technique of "2 kinds of —" in designing the Summary Machine, is probably correct — we design summary machine for CA, not C. — On the other hand, in updating code for C, → CA, we do not ~~only~~ use random coding of ~~the~~ best codes for C. We augment these codes w. 3 "OSL" codes for A. These are derived in special ways: depending on the <sup>short</sup> codes for C & A. How this is to be done, is clear for ZIP, but not in general clear.

The Q. about whether OSL can be done by MMH, is unclear: If MMH find a min code for CA, it may, indeed, be able to do OSL. It still would have many pc's for k-top B that occurred only once in CA.

An impl. idea: For any long sequence, no matter how random it is, the continuations of that seq. are not all of  $\approx$  prob. If the original seq. is n bit long, the 2^n bit repetition of the original seq. has much > pc than just  $2^{-2n}$ .

Here, the larger n is, the less regularly we can expect in a fixed length <sup>k bits</sup> extension of the original n bit seq. (? — is this true: say "repeat <sup>first</sup> last k bits"). — if  $k \gg$  the no. of bits in the inst "repeat", then this k bit seq. is of  $\approx$  pc than  $2^{-k}$ .

Actually, "A" can be quite large: Say it's a new Scientific Paper, or a big batch of data: One could code it in small increments, but I suspect that this is non-optimum.

One can code the whole thing, using codes in C, then recode it using new found codes in CA, loop to ~~end~~.

This seems to be the way I often read a diff paper: I go thru it several times — each time picking up more codes "understanding", based on ideas from previous runs over the sub-corpus.

- .29 [ At certain pts, ~~the way~~ I am reminded of ideas in C (the previous corpus) ]
- .30 [ So I may go to reference other works. ]

So there are 2 major updating methods: .25 incremental, .26: coding & recoding.

Also .29-.30 seems impl, but I'm not sure of the how to implement it!

One way: As I work, I develop categories of ideas, & areas in which they are applicable — I also write about a kind of Algebra for this: {problem areas} descriptors for.

{In methods} descriptors for: A some sort of statistical correln. of <sup>empirical</sup> usefulness including 2 sets of descriptors.

0!": 149.37-150.45 <sup>spec</sup> ; 116.40 ; On 150.07, I that J's problem was not divolved of type  $M(x), F(x)$

It may not be. If the Reinforcement comes at known (periodic) times, then it's a standard 02 prob. w. fixed CB. This is assuming J's "Greedy ~~is~~ Game"... which is probably not so good. I really don't know how to best ~~do~~ do it w/ horizon.

The it may be usually irrelevant: That distant future (predictions, awards, expectations) may be so imprecise, that they are gn. very little wt.

01: 152.40: A Common type of short A, to update is a Question. It has to be understood in terms of previous Questions, ultimately, in terms of Plain Answers.

A simple kind of Q to try to answer: We have a table in Menu. So we ask "what is content of address A?" I.M. Plain has how to give the answer.

If Persist. first data set TM is drawing on, the phrase  $\alpha$  is "what is the content of address B" irrelevant. only "A" is imp. Later, when we have other refs to addresses, the phrase  $\alpha$  will become imp. This is distinguishing one type of Q from another.

Back to 150.20 ff: One imp. idea is various classes of strings.

A string can be a word; an eq.; or TA; or  $2Tq$  (approximately two Eq); a logical statement. Some strings have themselves; some are true or false.

Some are Questions: Others are statements, or "answers" (can a number for string correct answer -

But I do want to see if I can deal w. things close to an "English representation" of

Another idea I may need is "local variable", "local statement" - these are true only with a specified sub corpus.

20 [12599] I.M.P.T Practical Aspect of OSL: Say in Science, One has several poss. Theories to explain "Observations up to Now." One of them is much better than the others, in terms of fit for the known corpus, so in MMh, we would tend to assume use only the "Best"

22 Theory. However, when New data arrives (from some previously untouched part of the scientific universe). The "best" Theory no longer fits very well. To make ~~the~~ new Theory fit the new data, we would, using Wollac's MML, use only the best Theory plus for its sub-corpus; ~~the~~ In fact what is usually done, is to take not only the best Theory & its sub-corpus, but much worse than the best Theories & their sub-corpus, in order to try to get a new Theory to fit both old & new data.

This is equiv. to what we do in O.S.L. Instead of using <sup>only</sup> codes that were good up to just before  $R$ , new data <sup>came</sup> coming in, we use codes that were fairly bad, but which would code the new data very well. I.e., if the new sub-corpus is  $\alpha$ , we consider codes which define  $\alpha$  - which only occurred once in the previous corpus & would be poor codes for the previous corpus, but good codes for the newly augmented corpus.

So if  $\alpha$  is the "new data" of .23, then we try to find codes for the old corpus +  $\alpha$ . Such codes could be rather poor for "old corpus alone". - (just as in OSL).

D28981 TM: TSO (.07R, .12ff)

Health Care Assoc.:

Dr. Weil

C. Brantner, Louise. Skapro Clinical Center, (CSA?)

① Rash: Legs, Arms, chest (some). Back of neck.  $\approx$  Two or three places where it itches

② Cough: Left lung hurts. Battered new (130M)

③ Cold.

④ Colonoscopy appointment.

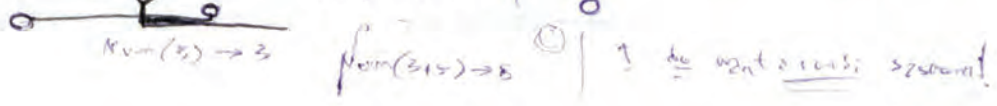
⑤ Ask about Medication after Jigg  
~~low copayment!~~

12 The idea of learning functions  
( $\equiv$  operator induction).

Solving eqns. can be trivial!

1. No  $(x, 3x=7)$  (maybe Sauer-Plaple!). For a "limit" or conclusion!

2. Solve  $(x, 3x=7) \Rightarrow x = 7/3$



Also, for each problem, I'll need "local variables" - "local facts".

"Local Universe" For various functions, one could have a context.

Argument telling where universe is

Perhaps just start w. some simple inductive proofs - using accepted systems - This is just to get a "feel" of the ideas. Later I can modify

the system to deal w a greater variety of problems.

Primitive facts that we start w:

1)  $Num(s) = (c, 1)$  depending on whether  $c$  is a "number" or not. Perhaps numbers will have to be abled: ordinarily Ray will be a string of bytes. When humans look at a no., they can always (?) tell it's a no. By No., I mean 2.71 or 13, not "II" or "e". (2.7535-21 also.)

Contra Bloom had a way of dealing w. nos. & doing Calculus, we running into lots of diffs ("Non-standard Analysis")

The I should be able to get T.M. to deal w. Math even tho it's not exactly "casi" - i.e. like person deals w. math.

Perhaps Go Proc that elementary & Algebra book (for high school) - so if I can simply present examples & get T.M. to do adaptive induction.

Dial @ Jones. Comm. -  
# 208 652.73  
rec'd. on Day!

9/29, 11/20/8

Winter 12M.  
Nov 5.

23 Nov 98

For Jones  
2 PM 2nd floor  
Dermatology  
Dr. Worth

Blood tests  
2 vials

Get fast blood tests from  
Gardino: Get dates of them  
~~fast~~ / fast blood test  
Nov 8 28 98



Book "Intermediate Algebra" "IA"

Starts out w. student able to add, subtract (i.e. <sup>zero</sup> positive/neg.) only. (No. neg. results poss.) (They be only integers).  
So I can start in Arith unit that can only do that.  $4+5=9, 5-4=1, 4-5 = \text{meaningless}$ .

I gather that the book is meant for not very advanced undergrads. - a first course in "Algebra"

Some of the old (69-80) look diff. T. ideas of "some", "every"  
Knowing what an "integer" is. (Integers =  $a \pmod{1}$  - But  $\mathbb{Z}$  is <sup>defn.</sup> ~~noted~~ <sub>noted</sub>.)

Infinite or finite sets: TM could make meaning of  $\{1, 2, 3, \dots\}$  or  $\{1, 2, \dots, 7, 9, 10, \dots, 12\}$ .

Defn. of an infinite set: could be entered for inclusion e.g. "Integers" ( $\equiv \mathbb{Z}$  mod 1).

1998-2000 Book - Structure Book

D2998 A perhaps diff. approach: (1) Math + Lang Functions (155.12) (2) Also (very) Definitions (155.12)

(But int. form of functions: so  $(x, y)$  Equation  $(3x+7=2) = 1 = \text{True}$ :  
paraphrase interesting concs. are being introduced, used.

(3) But start in middle of book & work backward to find functs. & defns. that are needed -  
to construct needed Conc. Def. m.

(18) sounds most reasonable - otherwise, there is a lot of stuff on numbers  
that will not be used to solve most problems. A

A possibly attractive thing about the Alg. Book, is that while much stuff  
is not needed, the way. seq. used in the book is perhaps close to that occurring in  
R.W. - in that much stuff is irrelevant to most problems -  
but mine to, eventually. T. approach of (18) is a more "minimal" system -  
More "brittle", perhaps.

(SN) on "Definitions": we could teach TM many imp. Definitions in

functional form. (17). Then see if it could recognize instances of any of the  
defns. in a corpus of "Text". This ability to recognize defined objects  
is probably quite important in prob. solving. Much content of science is in the  
form of definitions. T. "sci laws" are then (e.g.  $\mathbb{Z}$ ) relating things defined.

A list of things for "kinds of lang" for TM.

- 1) Function (rules (operator induction). Very broad class. Maybe broader than for all lang!
- 2) Definitions: (a) <sup>in</sup> functional form (e.g. 17) (b) find examples of  $\mathbb{Z}$ -defined things in R.W.'s corpus. (24A)
- 3) Concepts like "solve"; "simplicity"; idea of "closer to a soln" (a) scalar (b) vector (c)  $\mathbb{Z}$   $\mathbb{Z}$ .

When we have solved a new problem, we have

- 1) Augmentation of Corpus
- 2) Augmentation of concs. used to describe corpus

~~new concs w/o~~  
~~changes in~~  
updating of pc (frequencies of old concs).

D298 TM: Index of some imp't threads:

← corrections on p. 9.

1) Jürgen's: "Try how to learn" 9, 115  $\frac{1}{2}$ , 116

2) Mixed Corpus Problem. MCP: 107.01 ff This was solved on 7/15/40: Lab reports: 126.01: 138.01/148.01/158

3) History of Science: 122

4) "OZ problem" 99.36, 100.08, 102.18, 105.01, 108.05, 117.01, 123.01, 127.01, 128.01

122  
HAS

5) Summarizing (P.D.'s. Machines): 135.01, 137.01, 140.01, 146.01

117.01  
108.09  
105.01  
102.18  
100.01 123.0  
99.36 127.0

6) 2141: 18,

79.20 - 40; 90.01, 91.01, 97.03, 98.0, 109.01

7) GA, SGA, Org evolv:

112, 113, 132

8) OSL (One Shot Learning):

146.07, 152.01, 154.01

ABC

9) Unsolved probs (list of 4): 145.01