

The this folder contains

capt. papers I haven't been able to find. original! (8 12 27) well finally found traces in separate folders in TH 75A

Sub

On ladder on TM Green has

Priz

1) On Why Questions: Their use in TM's: ~~CM~~

CM 1 means "why" a fairly specific Q. w. 2 perhaps useful answers: E.g. One has a method for working a ruff class of problems. One tries y. method on a new prob. It doesn't work, (or works poorly w. rt some Govc). One asks "Why?". The answer via CM 1 is: How is this result made most likely by info int. previous corpus? Several parallel pems can be gen. w. max Govc. From these pems & t structures of these pems, one may then get clues on how to modify t. method to obtain one that works in t. new case.

processed
automatic
combination of pems.

Essentially, this is meant to model after human processes — in that humans ^{often} use t. Q "why" as a heuristic device.

→ Asking "Why t. method didn't work in this case" is a more genl. Q than asking "How can I best solve this new problem" — in that t. "Why" Q is usually of more general applicability in future problems.

"Why" in t. ⁰³ situation could result in a more stringent selection criterion for t. "Method" used. A selection crit. is a part of a Method that tells when it should be used s/o looks at t. problem & sets up certain params in t. method so it can be used.

206.01
306.18
more specific

.01: 304.40 General approach to AI devices:

At all times, retain large store of Pans around that have been useful. Status is "components" & "mutations" of these Pans should also be retained & indexed so one will be able to get at parts that are likely to be useful in forming new total "devices", or "daemons".

The storage will be hierarchical so that access time ~~is~~ (\sim cost) will be related to expected utility of each item stored. It may be more practical to store parts of ~~an~~ item, w. construction instructions, in ~~fast~~ rapid access, than to store the whole item in slow access.

.18: (304.90 specy): After the answer to why E-Method didn't work, one has many operators that one has accumulated, that will have suggestions on what to do. The "why" question "factors" to non-working of E. method — it breaks it down into useful "parts" so that various "difference operators" can look at it — see how E. method differs from what is desired, & suggest corrections to modify it to make it work more like what is desired.

Note that here we are working on the "Engineering problem" of how to bring about certain effects. In this case, how to get a "Method" (that is usually mutations of older methods), that will "work" the desired problem.

.32 Also on Method of Use of CMI in TM's: From 339-340 (on Puncta) & from other discs: It would seem that one could devise a set of methods of generating Pans, & one ended up with $A=0$ & ξ is small. (e.g. also, the Max Method). Here, it might be that one could save on C cost by looking at the corpus a bit before generating & testing Pans — so $A \neq 0$ & ξ is larger than if $A=0$; but one may gain in cost. So one has a trade off betw. cost & $A=0$ or $A \neq 0$, & one must learn how to choose betw. Pan-generative methods to get mix ξ for a given cost.

306.90 : (Perhaps!) New idea: Say TM has been working for a long time, & has processed a history of its own operation - in some abstracted form.

Now, we present TM w. a new problem. Just how does it use all of its previous experience to try to obtain a soln. to t. new problem?

We can, via CMI or ~~or~~ approxms to CMI (like SVM) write out ~~or~~ "solns" to t. new problem. But such solns will usually be too big in cost. The true problem for TM is what to do about t. new problem - \rightarrow a greatest possl. utility per cost expenditure will be obtained.

[Note that t. Q. of how to take into account t. ~~the~~ cost nacy for working on this "admin." problem, is unclear -

This Q came up in ~~SR~~ SP in a n way. I don't remember just what I did about it. - One way mite be to have

fixed ~~SR~~ ~~SR~~ srtns to solve these problems - with TM₂ continuously working on t. problem of deriving such srtns. ~~until~~ ~~TM~~ ~~becomes~~ ~~adequately~~ ~~good~~, I could be ~~TM~~ TM₂

~~at last a soln~~ when TM₁ is good enuf, we let TM₁ = TM₂ (i.e. TM₂ is one of TM₁'s many jobs types (say TM₁ is a PMTM).

It looks at "n" problems in t. past & from t. approaches that were used, & their success & failures, it decides on an approach for - t. present - problem. Note that TM need not use t. "same" approach as was used in a n prob. in t. past. Th. "similarity" may be fairly distant, so a rather complex x fun. may be nary. to change t. old approach into something appropriate to t. new problem.

This "similarity" problem is similar if not identical to that of extrapolation by analogy (L. Roberts' Pat Winston).

\rightarrow T. fogg (343.01 ff) is about t. same as Sussman's "Hochar" (328.36) ^(Pol) Sussman Also includes my earlier idea about going to t. library first!

I think I wrote a fairly large amount on this: Perhaps' mainly in SR.

Try to find previous work on this. Also PMTM, also BSTM.

Every time TM solves a new problem, it $\textcircled{1}$ batters its statistics in the usual way by \uparrow SSZ $\textcircled{2}$ If it solves t. problem in a "creative" way at all, it will have a new powerful tool to

Goal: Approaches to TM.

01:34340 use on new problems.

→ Essential Qs in this "any purported Gen. Approach":

- ① Is it Univl.: in the sense that by giving a suitable tag. seq., we can get it to learn any derivable trick?
- ② Is the SSZ required for ① reasonable? (Note that SSZ is $\frac{1}{\text{available IPE (cost)}}$ are reciprocally related).

Some new (reminiscences of old) ideas in this area:

Say a PMTM is used w. a SP administrator - using a library is a set of obs & ops. At a particular state of experience, the obs, ops, & various known correlns, in particular statistical forms, at that time can be used to solve a particular set of probs. w. some facility.

As Hvr., New ops & obs are created by CMI methods. They are ~~analyzed~~ tentatively evaluated using a small sample of recent problems. If they seem to be useful, they are evaluated using a larger SSZ from the past.

The end result of these new creations, is to throw out some old, little used, ~~low cost~~ ~~encompass~~ obs & ops & store the new ones.

Genl. (remarks)

sol: 344.40: If we have a deterministic corpus: Say 11111... = 100

Then using Laplace's rule, we got $P_m = \frac{1}{1+n}$ for the probab. ~~that~~ that the next symbol will be 1.

Hvr., from CMI - for deterministic seqs. $\prod_{i=1}^{\infty} P_i = \dots =$ probab. of deriv of the seq. which is > 0 .

$$\prod \left(\frac{1}{1+n} \right) = \prod_{i=1}^{\infty} \left(1 - \frac{1}{1+n} \right) \text{ which } \rightarrow 0,$$

since $\sum_{i=1}^{\infty} \frac{1}{1+n}$ diverges (T. moral is - for this seq., P_{n+1} must converge / more rapidly than $(1 - \frac{1}{1+n})$)

i.e. $\sum_{i=1}^{\infty} (1 - P_{n+1}) < \infty$. So Lap's rule would be wrong in this case!

It may well be that for all deterministic seqs Lap's rule is wrong. I.e. CMI might always give an essentially lax approx to deterministic seqs.

I think Carnap got variations on Lap's rule, but his expressions were always like $\frac{\alpha}{\alpha+n}$ where $\sum_{n=1}^{\infty} \left(1 - \frac{\alpha}{\alpha+n} \right)$ also diverges - So Carnap is also wrong for the deterministic case.

⇒ Note Also that if we have a Bernoulli seq. w. $P = \frac{1}{2}$ or $\frac{3}{4}$ or any "non-random" no., then $\sum_{i=1}^{\infty} (P_i - P)^2$ converges, since $\prod_{i=1}^{\infty} \left(\frac{P_i}{P} \right)^2 > 0$ which it does not do for the ordinary Lap's rule. i.e. $\frac{m+1}{m+n+1} - \frac{m}{m+n}$

$$\begin{aligned} &= \frac{(m+1)(m+n) - m(m+n+1)}{(m+n)(m+n+1)} = \frac{m(m+n) + (m+n) - m(m+n) - m}{(m+n)(m+n+1)} = \frac{n}{(m+n)(m+n+1)} \\ &= \frac{P}{m+n+1} \quad (P = \text{constant}) \end{aligned}$$

since $\sum_{m=1}^{\infty} \frac{1}{m+n+1}$ ~~also~~ diverges $\Rightarrow \prod_{m=1}^{\infty} \left(1 - \frac{P}{m+n+1} \right) = 0$

Lap's rule is obtained if one assumes uniform approx for P; then uses expected value of P. CMI assumes non-uniform approx for P.

Hvr., I think the difference between the 2 methods may be stronger than that - since I suspect small differences in approx would not give convergence in 1 case & divergence in the other! But this must be checked.

01: 35140 : A New type of PEM : (\equiv Computable proby measure).

W. defines Comp Proby Measures. — These are defined for all corp_i .

Would it not be possl. to define a Partially Computable proby measure, that is defined for some, but not all, corp_i ?

In this case, one inserts t . corpus in \mathcal{Q} , ~~put~~ into a machine, & the output is either ① the binary expansion of t -prop of that seq. — ② ~~stop~~ ^{signal} some computation then a "stop" w.o. other output ③ t . machine computes forever w.o. output.

If only ① & ② are true, t . measure is probly much more tractable than if ③ is possl.

Gen. ~~Blum's~~:

(Blum.)
1. (372.40)
(356.40 Gen)

A new tack: ~~Blum~~ Actually, I'm not so terribly interested in the CMI version of Blum's prob. I'm at present

willing to consider, say, all perms w. ~~length~~ ~~linear~~ CB linear in s to corpus length. (or s^2 or whatever).

The imp. Q is - how to get something like f most prob per expanded cost.

One guess: spent a cost $\propto 2^N$ on codes of length N ,



This is really poor: To see what happens: (i.e. is f imp. pt.) Make a tree of codes that ~~are~~ have ~~to~~ tracked f corpus as far as f branch went. At each branch and ~~if~~ we write how much cost has been expended thus far w.o. convergence (\equiv output symbol or stop ~~is~~ \equiv request for more input).

What happens is that we get branch ends w. very long input seqs. that do converge - but there are an enormous no. of such seqs that never must be tried - i many ~~do~~ do converge. What many of such ^{input} seqs. corresp to, are perms of various kinds. They are "setting up" perms like "linear regression" or "PSG induction" or "clustering" etc. Almost of these perms will not work at all well when they start giving output (i.e. they will have relatively ^{increase of} little output per \uparrow of input) - (e.g. linear regression systems w. wrong costs). So one wastes an enormous amt. of cost on these / input seqs.

Also, there are probably even more input seqs that have no output ever, but do converge in this sense.

As a result of .10 ff, I abandoned ^{"the} / "sequencia

"property machines" approach to CMI realization

\rightarrow .10 ff is, I think, a more exact statement of why that approach is bad, than my previous discu of it.

1. 0000

124.02
- .09 | 118.15
119.25

Gen 374.40: Anyway 374.40 ff may suggest more reasonable ways to search for good codes — or simply better ways to do CMI.

One much better way, is to make various measure observations on a corpus. These observations will then tell one which pairs are likely to work well.

Some simple observations:

- 1) Frequ of symbols & Ngrams. From this, γ effectiveness of Bern coding can be determined — also use of defns of various Ngrams.
- 2) Auto correln (for a corpus that is ^{is} a time series). From this one can determine which linear regression coefficients to use & how much "coding efficiency" one gets.
- 3) Auto & cross correln. of various functions of γ . time series of 2) — to do same thing as 2) but use of non-linear terms, power series (multit) (multidim Taylor) or exponential terms, Fourier terms, etc.
- 4) Concordance construction for ~~an~~ estimator how good CFG coding would be. see PSD 346.01 ff
- 5) For certain kinds of Data, make cluster diagram maps & find cluster boundaries, etc. One can then tell how much info compression one can get this way,

etc.

So: Assoc. ~~over~~ many sets of pairs, are sets of measurements to make on a corpus to determine which of that set, ^{of pairs} will code that corpus well,

It is clear that one can define a set of all possl. finite descr. pairs as Willis does — ~~but~~ Also, one can define a set

Gen ~~all~~ all possl. measurements on a corpus (this is 1. set of all functions from 1. corpus into 1. set of all integers — one can have pairs or utuples of functz for vector functs of 1. corpus)

But anyway, it is clear that at 1. present time, I don't have a good formalization of 1. search problem for induction codes. — Certainly not in a form to present it to an audience of nations. interested in min ~~complex~~ computational complexity algorithms.

As for this formalization, I'd perhaps best look at what humans seem to do in solving a ~~bug~~ bug seq.

Hvr., 1. General Problem of 1. lowest cost method to do ICMI does remain.

What seems to be done by Humans: we start out with an ongoing induction process, in which we have many pairs that we have used on 1. "past corpus". By factoring 1. pairs into named strus, we can get proofs for each strin o/a combination rule. We thus extrapolate 1. set of old pairs to devise new pairs. This can perhaps best be done by regarding 1. old good pairs as ass in some "lang" — perhaps a ~~lang~~ Cf Gramm is good envl — no other types of langs might have gramms that are easier to devr. Anyway, it might be negy to devise special grammars that induce new langs more like 1. way humans do.

So, a way to construct a TM: Start w. a good bunch of ~~many~~ pairs, that have been used w. a hypothetical "previous corpus". Use Grammatical induction to extrapolate this set of pairs. Use 1. new extrapolated pairs on 1. new / ^{sub} corpus. This ^{new sub corp} will reinforce some pairs & not others. Keep 1. good pairs in one's

Gen
'Sample set of f. lang' is extrapolate some new pems
is try on a new sub-corpus - etc.

The ~~form~~ will have to be formalized some, - but y. Q
of how to optimize it, isn't so clear. One could, at
some higher level (TM₂) devise better grammar types
~~to extrapolate f. pems.~~ to extrapolate f. pems.

So: present concept, TM₂ generates new pems
in perhaps a MTCAN way (i.e. zero best). These are tested
(or f. entire set of pems is tested) wrt y. containing
corpus. T. selection of pems wrt f. Gove is used
to ① select new & discard old pems from f. retained set
of pems ② vary f. params of f. pem-generating "grammar"
(or other device) of TM₂.

We might want to consider f. "set of pems"
as a single PEM. [Most of these pems can be given - i.e. zero best.]

The ~~idea~~ approach to actually designing such
a device: Use a fairly well-understood Tng seq.,
w. a known-to-be-adequate set of pems. - So
one knows which pems should be added in what
order under pressure from f. "driving" tng seq.
So a mechanism that is adequate to do this
could be looked for.

So, TM₂'s problem is to find a good seq. of PEM's -
which ~~it~~ does by "mutation" trials - (perhaps more intelligent
since subsequent trials can be guided by knowledge
of goodness of various previous trials).

So, in a sense, TM₂'s problem is to guess a good seq.
of pems - while TM₁'s prob. is to guess a seq. of
simpler symbols (i.e. f. symbols of f. corpus)

So, given a tag seq. that would be reasonable
for humans, we can assoc. with it a seq. of pems
that a human could learn. From such a seq. of
pems, we must then design a large ~~thing~~

Gen

1: 377.10 such that Δ differences betw. successive pams in Δ seq. are all "small", in the sense of likely not so that trial guesses would take too long.
 (i.e. expected ~~guess~~ time needed to guess Δ next pam is within a reasonable cost.)

06

Actually, the prog. TM₂ could just as well try to find (by \sim mutations) an I.O. device (\equiv operator).

This could be, \sim an **RTM**. The operator could be a genl. F.B. device in the sense of ~~ops~~ obs have ops in them & ops have obs in them.

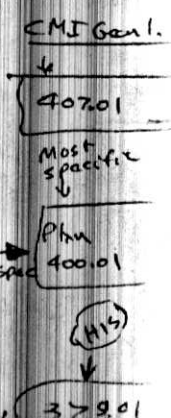
In fact this .06 device would be a fine simulation of an animal, because at any particular time, it ~~is a state~~ has a relatively fast response to its input (if that is what it was reentd for).

The technique of design of this TM: Make a tag seq. ~~then~~ then devise a seq. of operators or pams that would be an appropriate soln. Then look at Δ differences betw. successive ops. If they are too large either (1) put additional things in Δ tag seq. to make ~~the~~ some jumps smaller (2) Devise a ^{set of} suitable abs in a hypothetical pre tag seq. that could have given rise to these abs. The abs are \rightarrow to jumps w/ (they + hold set of abs.) are small.

More generally, Δ machines that we're moving along this tag seq., can be any type of TM — say even, most generally, a **PMTM** — A then we can begin to give this PMTM optimized probs like TM₂'s — etc.

→ However, first try this with some sort of simple, non-PMTM tag seq. — see to what extent it's poss.

→ It may well be that it is almost easy to have several modes in Δ TM, if reasonable jump sizes are to be achieved.



Genl. Notes

01.372.06

1) On Provably Convergent Automata. (\equiv PCA) One oldish idea was to consider all Pams that corresponded to PCA's — i.e. ~~Pam's that~~ Pams \rightarrow they ~~under consideration~~ could be shown to converge for all input seqs. of all lengths. Given a fixed set of postulates for proving convergence, the set of all PCA's is then r.e.

2 Q's (a) How does one go about adding new postulates? By Gödel's Theorem (or by a simple diagonal argt), there is always an additional postulate one can add. to find a Pam that frustrates the resultant system.

(b) Is the resultant set of Pams the same one specified by the "All spans method"? Even if not, if one deletes pams from the "all spans method" in an "unbiased" way, one will probably end up w. a ~~resultant~~ system that's about the same as CMI \rightarrow

on Typ Seps. 38/40

2) Notes on Typ Seps., PMTM, QATM:
PMTM 8 795.01 - 796.18 just how it works
QATM 8 776.15 - 788.40 (mode of PMTM)
Typ Seps, probs for PMTM 8 789 - 790

These are all somewhat burnt.

3) The set of Pams used is perhaps significantly diffrnt from those used in "normal" CMI — In CMI one uses a Pam for the prodn of the N^{th} symbol, if that Pam converges for seqs up to N symbols long. That Pam can diverge for the $N+1^{\text{st}}$ symbol.

In the PCA method, one uses a Pam only if it converges for all corpus lengths — i.e. all corpus lengths $\leq N$

See R 1, 2, 8, 35 for a disn. of provable convergence

D273 Gaul.

TM407

01:378.40 CMI General : IE³ Info Th. & Eng. Nov 73 P783 : A longish section

"Universal Noiseless Coding" — Also some review in it. This problem appears to be close to the CMI coding problem. There apparently hadn't been very much work done on it. Kolmogorov & T. Cover are some people who've worked on it.

Aug 24, 73

IPC

63.01-64.20
35.40

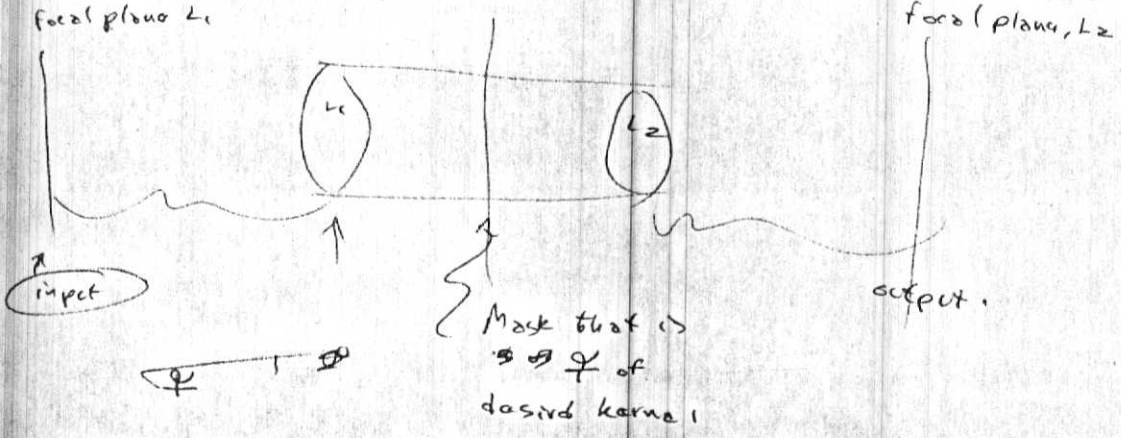
IPC: Provo may be later reflection of.

Scientific Amer; Masfex Gardner; Also William Gosper, MIT AI Group.

Optical computer; The "Game of Life" can simulate a Turing machine. - so it is a universal machine. It consists of a Boolean kernel

on pts in a 2 dim. cartesian array.

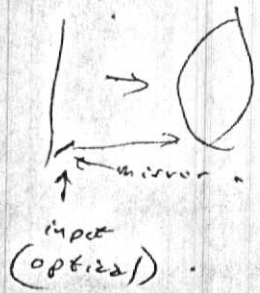
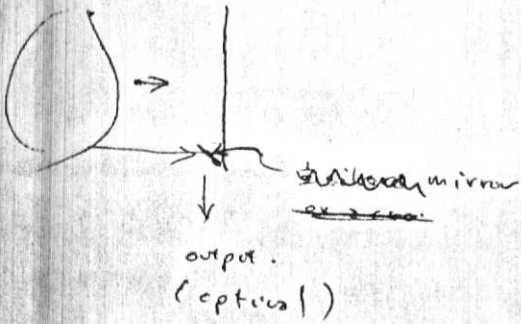
One mirrored bit able to do ^{2 dim} Boolean kernels optically, by using a linear kernel & a threshold:



If this can be done rapidly w. short wave lites, it would be very good. The full kernel mirror have to be done in several stages.

For a 10cm x 10cm array at 10^{-4} cm (= μ) spacing, we have 10^{10} points. At even 1ms x scan time we have ~ human IPC. At 10ps (10^{-11} sec) we are $\sim 10^9$ human IPC!

Input & output can be at the edge of any 2 dim array!



One problem is how to get another gain for the next kernel (x scan)

Aug 24, 73
30840
01; 282.40!

Hair Memory

k perhaps simpler way to do τ fogg; Simply have
only 1 very long stack w. all τ info in it. The last word
used goes to the top of τ stack, as before. Hence τ .
first N_1 members of τ stack have ^{fast} access time τ_1 ;
Then next " " " " " " slower " " τ_2 , etc.
CAM's are used for all parts of τ stack.

The result, is that the time it takes to retrieve a word
from τ stack, depends on how far back in τ stack it is.

[SN] Bubble memories may have much utility in making these stacks

\rightarrow - Be amount to a drum in which one can reverse beam instantaneously.

No! (This reversal ability only halves τ access time over an irreversible drum type memory).

So, one gets a fetch order. One moves from the top of τ stack
to the position of τ wanted word. The retrieval time \propto
 k , τ distance down τ stack. It then takes k more units times
to move that word back to the top of τ stack, where one
waits for the next "fetch".

(Hr, by using many stacks in ||, one is concerned only
w. τ access time of k , not τ "cycle time" of $\geq k$.

28

One can use bidirectional shift registers in τ some way.

A (perhaps) Big trouble w. τ fogg. (282.01 - 309.28) memory organization.

It would assign τ some very long address length to all words.
Ideally, τ quicker access words should have shorter addresses.

One way to perhaps get around this! Use only a few
randomly chosen bits of τ address to retrieve info. When
a find is made, ^{in τ access part of memory} more bits are compared. If a discrepancy is found,
one goes to a slower div. of memory (using more randomly
chosen address bits) & repeats τ process.

Usually components by themselves will have zero IPC. Combinations of components as well as certain larger components, can have $IPC > 0$.

Each assoc. w. each component, is a no., α , \rightarrow if α_i are α 's assoc. w. a set of components, γ . IPC of γ -set is always $\leq \sum \alpha_i$ for γ -set

I think I defined γ -IPC of (a set of) components (C_i) , to be $\lim_{T \cdot N \rightarrow \infty} \frac{1}{T \cdot N}$ (No. of ^(set) components needed to realize t . N.I.O. relation) \times (Time needed to do N.I.O. relation).

These I.O. relations could be listed by starting w. a fixed Umc, & listing all of its ~~possible~~ poss. input strings.

We may want to compare γ -sets of components IPC w. that of γ -Umc.

In the case of γ -Umc, γ -cost will be a fixed cost of γ -FSM + head, plus γ -length of tape used, times γ -time needed to complete γ -calculation. For large calculations, this \rightarrow Tape length \times calcn. time.

Blum's old paper on "A Machn. indep. Theory of computational complexity" (ACM ~ 1966) will be very relevant here.

This γ -Umc, seems to be cladly very inefficient. We would want a very large no. of heads (say $\rightarrow \infty$) to get more ~~reasonable~~ efficient operation. Thus, we

may not want no. of heads to \uparrow as fast as tape length. If

$\frac{\text{No. heads (+ FSM's)}}{\text{Tape length}} \rightarrow 0$, then for ∞ tape, γ - (heads + FSM's)

cost notary. Anyway \rightarrow since γ -no. of heads can \uparrow \ll linearly

w. tape length, we can \uparrow IPC of TRM by a factor \ll tape

length. - So IPC of TRM is \ll (Tape length)² - but

it can be (Tape length)^{2-\epsilon}, where ϵ is any no. > 0 .

3,25,7+ IPC

21.40

TM113

ol: ~~3,25,7+~~ I previously computed an upper bound on IPC of human brain by assuming 10^{10} neurons, each having direct access to all other neurons/each computational. in 0.1 sec. Certainly a poor estimate. While each neuron ^{for} ~~interacts~~ ^{has} this access, it even, perhaps in $< 10 \times 10^{-10}$ sec., this is a deceptive dem of what goes on.

Very probably there are multiplexed "trunk lines" so that while each neuron can contact every other in a short time (as in a telephone system), they cannot do this simultaneously (as in a phone system), because of the limited capacity of the common trunks.

It may well be that most of the IPC is not contained in the trunks, so one can compute the IPC of the Brain w.o. ~~knowing~~ ^{knowing} their copy (or just having some upper bound on it).

Anyway, say a small unit does "its thing" in 1 ~~sec.~~ ^{sec.}

Then, the Q. is, how many bits outside itself does it have access to in T ~~secs.~~ ^{secs.}? If we keep T small, various units will not interfere w. one another, & we can add their IPC's to obtain $t \leq$ ipc of the net. If T is large, we obtain a spuriously large net IPC, because the units would interfere & would not be able to use the same trunk lines) simultaneously.

We might be able to get a simple expression for the upper bound on

1. IPC of a net in which each unit has n ^{bits} input & m output bits (w. m, n , small) & does this in 1 second. Here, can we assign an IPC to a large "core" memory? — or, what about a "trunk line" (a multiplexer), that has many inputs & 1 output? The ipc of these last 2 devices is very impt. because we will probably use a core memory as a (hopefully) main source of IPC in our TM.

~~At least~~ ^{RAM} $N \times 1$ bit/word can be regarded as

a N bit word, w. some kind of multiplexed access to N bits — so one can only get at one at a time.

A ROM ($N \times 1$ bit) is also a word multiplex, but unchangeable. Also a ROM can be regarded as a fixed function, w. no memory

Punctuations (average error - Puncts)

269.40

A way in which $\alpha \cdot P$ may be better for MM than DM:

Using the previous analysis, I get: $C_{DM} A_{DM} \stackrel{=}{=} C_{MM} A_{MM} + \text{a constant indep of } S$

~~$C_{DM} A_{DM} + C_{DM}$~~ $= C_{MM} + \text{const indep of } S.$

Hvr, I didn't consider that there is a term in A_{DM} $\propto \ln \ln S$ — This is a punctuation term that falls when t. Punc beam ends & when t. corpus decn. starts. — I may have noticed this, but it was not regarded as important at t. time, since I felt that I had to correct a much larger difference betw

$A_{DM} + C_{DM} \approx C_{MM} - \text{if one of } \approx \pm (N+1) \ln S + \text{const.}$

Hvr, I'm not sure this effect exists — I may have to go into t. exact coding in much more detail.

268.88 — 269.40 goes into this. I think, hvr, that the algebra is very wrong!

Starting at 268.25 t. prob of a signal param of

accuracy δ , will be $(\frac{1}{2}(1-\alpha))^{\log_2 \delta}$. $\alpha = \frac{1}{4} \delta \alpha (1-\alpha)^{\log_2 \delta}$

$n_{cost} = -\ln \alpha - \ln \delta + \log_2 \delta \cdot \ln(1-\alpha)$
 $\approx 1.4 \ln \delta \approx -1.4$
 $-\ln \alpha - \ln \delta - 1.4 \alpha \ln \delta$
 $\alpha - \ln \alpha = \left(1 + \frac{1.4 \delta}{\alpha}\right) \ln \delta$

The α in n_{cost} over ~~not~~ considering $\gamma \cdot \alpha$'s, is

$-\ln \alpha + 1.4 \ln \delta \cdot (\ln(1-\alpha))$ for each of $N+1$ components.

(The $N+1$ is ~~2~~ is a bit different from the other N)

Say $\alpha = \frac{1}{\log_2 \delta} = \frac{.7}{\ln \delta}$ so 1 in n_{cost} is $\bar{\delta}$ is average accuracy — perhaps $\alpha = -\left(\frac{1}{\log_2 \delta}\right)$ would be better

$-\ln \alpha = \ln(-\log_2 \delta) = \ln(1.4) - \ln \delta$

$+ \ln(\ln \delta) + \ln 1.4 + 1.4 \ln \delta \times \left(\approx \frac{.7}{\ln \delta}\right)$

$= + \ln(\ln \delta) + \ln 1.4 + \left(\approx 1.4 \times .7\right) \approx + \ln(-\ln \delta) + \ln 1.4 + 1$

Aug 27, 73

Punctuate (etc). $\pm 1.35 =$

So $\ln(-\ln S) + \ln(1.4) + 1 \approx \ln(-\ln S) + 1.35$

$S \propto s^{-1/2}$ so this $\rightarrow \ln(\frac{1}{2} \ln S) + .65 = \ln \ln S + \ln \frac{1}{2}$
 $\approx \ln \ln S - .05$

$- \ln 2 + 1.35$
 $- .7 + 1.35$
 $= +.65$

Actually, since we must do $N+1$ params,
 + total correction is \sim

$(N+1) \ln \ln S + (N+1) \times .65$

This term is very approx —
 it may actually be zero

More exactly, $-\ln \alpha + \frac{\ln S}{\ln 2}$ (for $\alpha = -\frac{\ln 2}{\ln S}$)
 $= -\ln \ln 2 - \ln 2 + \ln \ln S + 1$

for $\alpha = -\frac{\ln 2}{\ln S}$

$\delta = S^{-1/2}$ so

$-\ln \delta = \frac{1}{2} \ln S$

$\frac{\ln S}{2 \ln 2} (\ln(1 - \frac{2 \ln 2}{\ln S}))$

$\alpha = \frac{2 \ln 2}{\ln S}$

$\approx -\ln(.7) - .7 + \ln \ln S + 1$
 $.35 - .7 + 1 + \ln \ln S$

$\ln S = -\frac{\ln S}{2}$

$\ln .7 = \ln \frac{1}{\sqrt{2}} = \frac{1}{2} \ln \frac{1}{2}$
 $= -\frac{1}{2} \ln 2 = -.35$

$\approx \ln \ln S + .65$

$-\ln .7 = \ln 1.4 = \frac{1}{2} \ln 2 = .35$

so $\approx (N+1) \ln \ln S + .65(N+1)$
 for small α ; i.e. small δ .

The fogg analysis is not exactly right — i.e.

a more exact treatment would take $\alpha = -\left(\frac{1}{\log_2 S}\right)$ i.e. α is

an average for all of the δ 's. Thus, since each δ is $\propto S^{-1/2}$
 the results are about the same.

The fogg result is impt. if the loss good (Rao simpler)

SVH Govc of $\frac{F+A}{S}$ is vsd. In this case MM is better
 by $\sim (N+1) \ln \ln S$ than DM.

If the Govc $\frac{S \cdot F}{S+A}$ is used, then DM is still better than

MM by $\sim (1-k) \frac{1}{2} (N+1) \ln S$.

Punctuation

01: 329.25
315.90

A rather General Approach to punctuating Code 1. corpus using as many symbols (including punctu. symbols) as neccy. This produces an intermediate code for 1. corpus = C_1 . Say there are v symbols. We then Code C_1 as if it were a modified Baru Seq.

$\sqrt{= C_0}$

"Modified" means that at each point in 1. corpus ~~is~~ C_1 , there may be constraints limiting 1. next symbol, so 1. probab. distribn. over "1. next symbol" may vary (but in a known way) ~~with each separate symbol.~~ at each pt. in 1. corpus. The set of symbols in C_1 is 1. constraints of 1. symbols in C_1 are regarded as "free" (i.e. given).

18

One way to get 1. probs of coding C_1 : Use Laplace's rule in each "situation" that arises. e.g. ~~is~~ ~~with~~ At a certain pt. in 1. corpus we are in situation (\equiv state) S_3 . In S_3 only symbols ~~alpha, delta~~ α, δ, ϵ are possl. Then 1. prob of δ , say, will be $\frac{n_\delta}{(n_\alpha+1)+(n_\delta+1)+(n_\epsilon+1)}$ where $n_\alpha, n_\delta, n_\epsilon$ are 1. no. of times that α, δ, ϵ has appeared in 1. part in situation S_3 .

In 1. case of certain punctuation symbols, they can only occur once (or twice or some finite no. of times). e.g. "1. symbol that says that the previous set of 0's & 1's ~~is~~ represents an integer that is now terminated, is that this integer tells how many coeffs we want to use in our linear regressn."

It might be possl. to have more complex rules for pooling of data" for 1. various symbols. In 1. forgo. (0.18 ff) we regard each "situation" as separate & unrelated to ~~other~~ data accumulated in other "situations". In general, I think this need not be so.

35

In general hvr., one can make up any rules that one likes for probab. of 1. various symbols. If these rules are devised before one has seen 1. corpus, then one can legitimately do something like an "All pairs Method" to get

Puncta.

.01: 339.40 a ~~to~~ \sum proby distribu. This \sum proby distribu. will have recent $A=0$, so one can use its observed \sum to extrapolate directly into t. expected future \sum .

The idea here is that the ^{stochastic} rules for generating C_t , give one an a priori over a set of Pems, so one can use Bayes directly to make predictions. I do not know if t.

.12 "All Pems method" is in general correct, (see Index 323.05-.07 for Refstodion.) but in t. present case, it is legitimate.

FOR's!
u willis 139.40

This may be because FOR's can only represent r-computable proby. measures. (ie. rational), but this is because, for any finite output, one can only have finite inputs.

Q: If I give an FOR finite input - can I tell if it will ever stop?

I suspect so. - ie. I think W. says it could only have rational fractions as output in this case - (ie. periodic. (No I haven't been able to find such a counter in ACM after a period of it))

If so, a FOR is probly a FSM

Conversely, every FSM can be slightly modified to become a FOR. If true, this fact would make it a lot easier to remember & derive properties of FOR's.

It probly is easy to show every FOR \leftrightarrow a FSM \rightarrow (Not easy)

Here, FOR's don't necly have S & U. ~~CB's do.~~ ~~Every CB is FOR.~~

Note ACM 248.03: W. says a CB can be made w. finite tapes & no repeating cells! This may destroy i bad job!

W's remarks (ACM 257.23-28) may be almost identical to my proof that no other proby evaln. system can do better than CM.

Perhaps I can disprove .06 or modify it. I.e. try to find a FSM that violates having an upper bound on the length of a program that will give an output of length m (or a specific output of length m). Or has finite output for infinite input.

Consider a FSM: binary input, output. It counts the total no. of 1's in its input. As soon as there are just 2, it stops. For each n , there are a finite no. of programs that will give n as output. - But for ~~many~~ infinite inputs, there is ~~no~~ output, much less infinite output.

Also, there is no bound on the length of a program necessary to produce a program of length n . So .07 is probably wrong. The "what about 'slight modify'?"

So, all FSM's are not FOR's! But are all FOR's FSM's or simulatable by FSM's?

Q: Given any $L(m)$ ($L(m) > 0$, $L(m)$ integer)

can we construct a FOR \rightarrow ACM 246.30 is satisfied?

I think so. We could first construct $\eta_R(M, M)$ to satisfy ACM 243.32. Then we could, by Thm 2, construct the assoc. Programmable Machn.

Note that $L(M)$ is an upper bound on n for a gm. m . n for m can be $< L(M)$.

For FSM, must $L(m)$ be bounded in some way (e.g. $L(m) < m$ or $L(m)^2$)

FOR'S

01'330.40 : W shows (ACM 2 + 6.39 Thump) that for any FOR, R

$\eta_R(m, n)$ is the no. of m pgs of length n

$L_R(m)$ A limit on the length n , of a m pgm.

$\mu_R(A^m)$ - The Kolmog. complexity (cost) of the seq. A^m

is the length of the shortest seq. that can produce A^m as output

Are all ~~seq.~~ effectively computable.

Can these results be used to show that a FOR is a FSM?

Note: Since every FSM can be simulated by some Turing, or "programmable Machine", then it is probably a CB on a one i can.

be simulated by a FOR that does not have the symbol \cup in its output.

Hvs, consider 330.25-30 - counter example!

~~SN~~ FOR'S may be more powerful than FSM'S:

Reason: I suspect that for FSM'S $L_R(m)$ cannot \uparrow too rapidly w. m , because of finite memory. (I'm not so sure of this)

also, because of counter example 330.25-30,
 Also FSM's can get into loops but FOR's can't

[SN] Note ~~that~~ that a simple time limited Turing is not

a FOR, ~~but~~ But can define a FOR. If the time limit is T , then for all inputs of length $\geq T$,

it must be in a S or \cup loop before it reads the entire input. \rightarrow 352.2

A machine w. finite memory (a FSM) usually cannot be directly x'fnd into a FOR. A FOR can't get into pure compa. loops, a ~~seq.~~ FSM can easily get into such loops.

One can fix the FSM by having a watching device to see that it never repeats the same state - if it does, \cup .

System automatically $\rightarrow \cup$.

Hvs, a FOR also cannot read an ∞ of input w.o.

printing some output. We can assure this by having an aux device see that the FSM never repeats a state at 2 difent times, unless it prints at least 1 symbol betw. those times.

If it does do such a repetition w.o. printing, we go into a \cup loop.

The aux device resets itself every time a char. is printed. Then it looks for repetitions of states since that last printing.

If a repeat state, ^{so} were allowed w.o. printing, then the

input seq. betw. those 2 identical states could be repeated

FOR's

01: 337.40 indefinitely, i.e. we would have finite output for infinitely

02 long input. Also, if there ^{exists} any sequence after state s_0 , that results in printing, (say t -seq α is printed),

Then we could get a certain input seq; terminating in α , that has an arby no. of codes, of lengths that are arby long (in multiples of t , input repeating seq.).

06 But anyway, my impression is that ~~there~~ a FSM that had an read an infinitely long input w.o. additional output could not be a FOR — unless we somehow know it was in T 's situation & flipped it into U ^{or S} in which case it would stop reading. So I think such a FSM ~~could~~ behavior could not be worked into a FOR.

Willis mentions a FOR (ACM 248.03) that is probly I'dentical to t -source I desc'd on 337.33-37

20: 337.23 Re: \uparrow limited machines! on ^{ACM} 247.43 w. desc's a $T(m)$ limited machine as a FOR. Here t -time limit is a $\#$ computable function of m , t -no. of ~~input~~ output symbols printed thus far. we can, from this use $T(m)$ as an upper bound to $L(m)$ (ACM 246.30) is t -longest ~~input~~ input needed to generate an mppm.

$T(m)$ need not be an \uparrow func. of m . It can be \downarrow or even constant. — But it must be computable —

30: 330.02: Any FOR w. finite input must stop eventually.

If it did not, we would either have an ∞ comp. loop, or inf. output w. finite input — both of which are illegal for FOR's.

\rightarrow No! Inf. output w. finite input is O.K. — In fact, if an inf string is deterministic, it must have a finite descr. i.e. be of probty (\exists probt) > 0 . \rightarrow 357.01

33: 330.39: For a FSM, $L(m)$ need not be bounded (i.e. $L(m)$ need not exist)

— An example is 337.38 — 352.06 (352.02-06 in particular).

So, a conclusion is, that by ~~discovery~~ properly creating recurrent states (337.24-37) we can xtra ~~add~~ FSM into a FOR.

\rightarrow Whether all FORS are FSMs or \approx FSM's is not, hvr., clear



Consider th. folg. FOR: For any string of input symbols

08.13 FOR's

of length n , it prints just n^2 1's and stops. This cannot be a FSM because it requires ^{a total of} ~~potentially~~ infinite memory.

It appears to be a FOR w. $L(m) = \left\lceil \sqrt{m} \right\rceil$ smallest integer \geq

If never uses ϵ , S or U symbols.

↑ Xerox

TOR'S

352.30
01: 353.40
330.01

The problem of 330.01 is open. I can't think of any way to tell if a fin. finite input has ∞ output or not. I think that ∞ output w. finite input is legal for FOR'S, but whether it is poss. to tell in finite time - I don't know.

If it ever passed thru the same state twice, it would be a finite computing loop - i.e. we could tell this. But then, it would use only a finite amt. of memory - a FSM.

In general a FOR can have ∞ memory.

Say a FOR was computing $\sqrt{2}$ - could a FOR do such things? I think so.

Well: say if its input is 0, its output is 1, ∞ .

" " " " 1 " " " successive digits of $(\sqrt{2}-1)$.

So for any finite output seq., there are only a finite no. of possible perms.

1.914
- .250

1.664
- .025

1.639
- .035

1.604

1	has	code (i.e. 0)
0	"	"
01	"	"
011	"	"
0111	"	"
01111	"	"
011111	"	"
0111111	"	"

00 has no code
010 has no code
0101 has no code

So $L(M) = \{1\}$ is fine

$$\left. \begin{aligned} \eta_R(1, \Sigma) &= 2 \\ \eta_R(m, \Sigma) &= 2^m \text{ for } m > 1 \end{aligned} \right\} \text{ so } \eta(m, 1) = 2$$

$$\sum_{r=1}^{\infty} \frac{\eta_R(m, \Sigma)}{2^r} = 1 = \frac{\eta_R(m, 1)}{2} = \frac{2}{2} = 1$$

I doubt if any FSM could compute $\sqrt{2}$ - I couldn't because $\sqrt{2}$ is not a rational (i.e. periodic) fraction.

- 25) Al Meyer, Mike Fisher
- 26) Chomsky?
- 27) Report? (Canada)
- 28) Giorgio Levi
I.E.I.
VIA. S. Maria 46
56100 PISA, Italy
- 29) A. Rosenfeld, U of Md.
- 30) Ed. Fredkin, MAC
- 31) Peter Elias [interested in "info retrieval" from pt.]
Elias will be more of view of info theory oriented toward CMI as a AI approach, than most AI workers. - i.e. use of prob. theory, info theory, analysis, etc.
- 22) Bill H.

SN

① Suggestion that Pat Winston's Process would be a good thing to apply CMI to ~~to~~ to quantify some of his "sort of" Grammatical Induction

② Marv. Said that Berlin had written a more complete analysis than Slagle, on Task Nets. (\equiv SP)

③ Pol Politics: Papers on applic. of CMI to various problems would be useful. e.g. ① also perhaps its application to linear regression is all of Y. associated results.

④ There are a few papers at it. ^{Aug} 73 A.I. conf. on Induction. They don't look v. G., but check them.

⑤ / Another Approach to getting more people working on CMI; Write Paper or Give Seminar or Teach course on: CMI: An integrated approach to Problems in A.I.

Treat: ① Prediction (discrete & continuous) ② Operator extrapolation ③ Pattern discovery. ④ Mechanics of how CMI can guide the construction of inductive hypotheses. — e.g. The "mutation method" ~~is~~ is treated by CMI.

⑥ Some impt. ideas around to A.I. lab. to get familiar w. I can perhaps show how CMI could be applied to those to optimize them: ~~Perhaps put this in "CRIT" file~~

a) Loops: Planner, u planner, control, libp, Trac, Snobol see 2/27/71 to 5/27/72 MAC Program report for detn. of "Planner" pp 95-101 ibid p 88

b) Winston's Thesis on "induction" — How it is an improvement of Larry Roberts's Thesis

c) Winograd's Thesis

d) { Sussman's "Hacker" — Aug 73 Signet newsletter p 26 Probably written in October

e) Iva Goldstein's "

f) Marv's ideas on "Frames" — how they differ from previous stuff (see my notes on Marv.)

g) Actor (s) sig Art newsletter Aug 73 p 24; IJCAI '73 also

f) Chavira's stuff on learning to read

(6) (cont) g): Read MIT AI Lab Quarterly or Yrly reports
to get General idea of what's going on

~~Read the book "The Structure of Ill-Structured Problems"~~

(7) Publish Langley Abstract in AI newsletter.

Politics:

01.342.40 General Meanderings:

Most immediate things to get done:

- 1) Review Results on Uniqueness of CMI: Mail out to various people.
- 2) "Study paper" on CIAP: Various Scenarios; On f. need for discussion.
- 3) Papers on "How CMI is useful."
 - a) In Sci. Method.
 - b) In A.I.
 - c) It has interesting math probs assoc. with it.
- 4) Give Papers at various Conferences.

Present immediate Goals:

- 1) Get lots of good people working on
 - a) Theoretic CMI
 - b) CMI-based AI.
- 2) Get people thinking & writing about CIAP: Perhaps form informal "College"

Wrt. 010, I think I want more general work on Theory of CMI. For this reason, the "Uniqueness" paper is important. Also, if I could get results or motivational "results" on optimum trial sequences of Pams, etc. — This for "computational complexity" people.

Specifically, what I want/is ^{most immediately} — various derivative induction methods — so that the Q. of $\max \frac{p_{cost}}{c_{cost}}$ becomes reasonably close to what is desired in practice.

- 5) Paper: ~~The~~ % Talk: The Quantification of Historical Explanation — The Application of Historical Analysis to ~~the~~ Problems in Artificial Intelligence. To be given at a symp on f. history & philosophy of Science & a similar Journal. Perhaps ask Hil. Putnam where to send this. [TNY Seg 401.30-40 is also on this Q]

1:373.90 :

TALK - Say at MIT - ~~is~~ Slightly popular.

TM 9

Start out by listing ~~some important~~ problems in induction.

- 1) Prob. of Geom. Probly. Also How it is just another form of assigning ~~perhaps~~
- 2) Expected error in curve fitting, if one is allowed to use any functional form one desires. How does one \uparrow expected error due to Δ fitness? Show how one can fit ~~the~~ data ~~as~~ well w.o. any continuous priors ("observed error = expected error")
- 3) Grua, Blean problem
- 4) Look in Sci Amer in Math section, I think, for ^{are on} dittys in induction - 1973, I think. Possibly a regular article [^] (maybe more likely).
- 5)

Previous refs (volumes) exist!

Given a corpus; First make a concordance. This is a list of all subseqs - occurring in the corpus, w. the freq. of occurrence of each. Unfortly, a complete concordance of the entire corpus is probably too expensive. It can contain no greater than N^2 items strings! Where N is the length of the longest string.

Where N is the no. of characters in the corpus, δ is the length of the longest string, δ is the length of the longest string, δ is the length of the longest string.

So it may be possible to have a complete concordance. We may want to put parts of the info stored away. Say 100 as's, each 20 symbols long; $2 \times 100 \times 20 = 20k$. 20k substrings as well as the no. of occurrences of each.

Using such a ~~form~~ concordance, it will be easy to try out various grammar rules of the form $\alpha \rightarrow \beta \gamma$ or $\alpha \rightarrow \beta \gamma \delta$ etc. I don't know how good it would be at testing rules like $\alpha \rightarrow \beta \gamma$ or $\alpha \rightarrow \beta \gamma \delta$.

Any way, it should be possible to compute under what circumstances one would save on cost by doing the count or dance first, before making trial rules; or just keep the corpus & doing a freq. run for each new rule.

01:34640 to be tested. I suspect it would involve

less cost — but one has to include in cost, the

cost of many of various access times — i.e. a properly

designed hair memory.

I am thinking of this concordance as being roughly

analogous to a spectrum.

Some version of the Bible has an IBM sponsored

concordance. Perhaps I could use this to derive a reasonable

set of grammar rules for part of "English" (or Latin)

5:34635 -

I could try all $\alpha \rightarrow \beta$ - type rules by $\frac{\text{frequency of } \alpha}{\text{frequency of } \beta}$

having all phrases (EX) : looking at it. No ways that

$\alpha \rightarrow \beta$ could occur. This may be a big task, but I think

there will be many zero entries involves little computation

To search for $\alpha \rightarrow \beta$ type rules is not

obviously ~~practical~~ ^{practical} related to the concordance.

One way to use the concordance for this: Say the conc. is

in alphabetical order. We can then ask for γ . (Wtd)

Set of strings that follows β (say). ~~the~~ call

Pairs with set S_1 . We then ask for γ . wtd

Set following β (say) - call this set S_2 .

It S_1 : S_2 have a similar profile than we

will gain something by comparing β : γ into

a single set.

~~Work~~ We also will want backward concordance

in which all ~~the~~ ^{the} sample of all α 's is

given backwards (i.e. each β is given backwards) : a concordance

is made.

In genl. Much of the ~~work~~ ^{work} ideas : other good ones, are

discussed in my "Machn. of Ling. Learn" (Naur).

PSED

01; 347.40

Our problem is evaln. of the "fit" of a grammar -
 Say, for each sample S_i , one has \geq diffrnt possl.
 parses wrt a gn. Gramm. If there are N samples S_i ,
 we have 3^N diffrnt. ~~parse~~ parse sets.

Each parse can be represented by a vector, a
 choice of N vectors "sums" ~~combine~~ these vectors in
 some way to result gives parse.

~~to~~ To get the total parse of ϵ /grammar ϵ ,
 we will want to \sum parse of these 3^N parse sets.
 I.e. we would be most interested in grammar for
 which that \sum parse was max.

For prediction of the parse of a new possl. S_i ,
 we would have to \geq the parse of all ~~sets~~ $\geq N+1$ parse
 sets.

Very probly, we would need to consider only
 the "few" "top" parse sets.

Actually, the most S_i 's will have $\gg \geq$ parses,
 & the actual/parse vectors can be expressed in some
 simplified "factored" form.

Biblio

T Pzurs Introdn to LISP

W. D. Maurer (1972)

MacDonald/Amer Elsevier Computer Monographs

\$ 9.95 102 pp.

Looks rather simple. Has lots of exercises,
Answers to many. Looks easy to read.

348

or 396