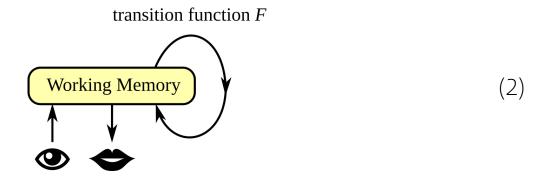
## 逻辑化 AGI 基础

这是强化学习最基本的 setup1:



状态转移函数 F 负责 **更新** 工作记忆 (WM). 如果 F 在一个 **闭环** 内训练,它似乎可以"解释"(或 预测)输入的讯息。

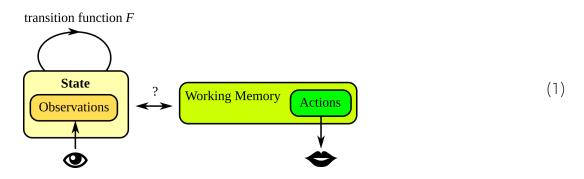
In classical logic-based AI, **inductive learning** means searching for a **theory** T (= set of logic rules) that "explains" or **implies** positive examples but not negative ones:

$$T \vdash e^+, \quad T \not\vdash e^-$$
 (3)

While logic learning is powerful, it relies on **combinatorial search** and was too inefficient, which caused "AI Winter".

我理论的重点是:透过 F 在 RL 闭环内的训练,可以做到 <mark>发现</mark> 逻辑理论 T 的效果。

<sup>&</sup>lt;sup>1</sup>The above diagram is somewhat inaccurate as there are subtle differences between the notions of "state" and Working Memory. In RL, state usually refers to the external environment through observations. WM is slightly different in the sense that an internal belief may be wrong about the environment – eg, mistaking seeing something that doesn't exist. There is on-going research as to the relation between RL and WM. I have not fully resolved this issue.



The essence of the standard model is just to identify a Working Memory as the "state" of the AGI system. One benefit of our theory is that it relates Transformers / BERT / GPT to AGI

The "standard model" is a way of thinking, that may help us better under-

逻辑化 AGI 基础

systems. These language models are phenomenally intelligent, yet many people criticize them as not "truly" intelligent. The standard model suggests that they are indeed linked to AGI.

This is the simplest form of a dynamical system: transition function F $\int$  state x = "working memory" (4)When we add a "control" or "action" variable a to it, it becomes the most

## F(x,a)

basic control system:

towards AGI systems.

from it.

1

**Memory Tape:** 

in the memory matrix M.

students of the Transformer:

stand the general theory of AGI systems.

Reinforcement learning

which is the setting for Dynamic Programming or Reinforcement Learning.

The optimal solution for such systems is governed by the Hamilton-Jacobi-Bellman equation 1:  $V_t^* = \max_a \mathbb{E}[R_{t+1} + V_{t+1}^*]$ 

(6)TO-DO: It would be worthwhile to find the brain mechanism that approximates reinforcement learning and use it to help the design of AGI. Recently, Yann LeCun's Energy-Based Models offers a way to circumvent

the problem of learning probability distributions over actions, when the ac-

tion space is hugely high-dimensional. This seems to be an important step

I call this the "standard model" because of the extreme simplicity of this

biological brain

(9)

(10)

(11)

(12)

(13)

(14)

(15)

(16)

(17)

(18)

(20)

(21)

↑ attention

input #3

value

(5)

The following diagram shows how the standard model relates to several other important areas, so we can reap profits from their interactions: "attention" mechanism

neural Turing machine

**BERT/GPT** 0 0 1 1 0 0 0 1 (7)**Turing machine** 

Neural Turing Machines and Transformers

The attention mechanism was first proposed in the "Neural Turing Ma-

Recall that a Turing machine is a Finite State Machine augmented with a

memory

reinforcement learning

chine" paper by Graves et al [2014].

Finite state machine

(S2) **Turing** (8)machine 0 0 infinite memory tape In Neural Turing Machines, Graves et al proposed the attention mechanism for an RNN "Controller" (playing the role of the Finite State Machine) to read and write from a Memory Matrix (the tape), using a content-based addressing method.

The Memory Matrix M consists of N items, each of constant size. The dis-

creteness of the address would introduce discontinuities in gradients of the

output, hence we need an Attention Vector to focus on a specific location

The Attention Vector  $\vec{a}$  is calculated via the following formula, familiar to

 $\vec{a} = \operatorname{soft} \max_i \{\mathcal{D}(K, M_i)\}$ 

where D() is a similarity measure between the key K and memory item  $M_i$ .

This then evolved into the **Self Attention** mechanism used in all Transform-

The key K is emitted by the Controller as the value that it is looking for.

## ers. Now let us refresh with this diagram illustrating Self-Attention (redrawn from a blog article on the web): output #1 output #2 output #3

The research done by Olah et al, in their 2021 paper A Mathematical Framework for Transformer Circuits, is very helpful towards understanding Transformers and Self-Attention. For example when we say "all men are mortal":

 $\forall x. \, \mathsf{Human}(x) \Rightarrow \mathsf{Mortal}(x)$ 

any object instantiated as x (eg. "Socrates") would have to be **copied** from

value

input #2

cat!

feature vectors

The result is the emergence of **disentangled features**. There is now a lot

of research papers on this topic; Personally I first learned of this from Marta

Garnelo and Murray Shanahan's paper Shanahan2019. We can think of

this as a first step of symbolization, in which objects are recognized by

In the cortex, neuronal populations are organized into "columns", with lat-

eral inhibition among themselves. When one population is activated, it

suppresses the activation of nearby populations. This is likely to be the

neuronal populations

It is remarkable that the **softmax** in the Transformer / Self-Attention seems

to be an abstract implementation of this winner-takes-all selection mech-

Moreover, the cortex is organized into layers with widespread recurrent (ie,

This bi-directional architecture may be applicable to AGI architecture

(see also §3 on abductive reasoning), possibly replacing the current uni-

directional model of feed-forward networks and the back-propagation al-

essence of deep learning? I think the answer lies in two words, "hierarchi-

cal" and "learned". As a counter example, decision trees are hierarchical

glass", this too can emerge out of disentanglement of features, because

it is a very **economical** / efficient representation of a complex scene:

lateral

inhibition

forward-

projections

mechanism that enables disentangled features to emerge:

cat!

## As is well-known, the brain does not use back-prop. The bi-directional innervation is a very significant brain architectural feature that has not yet been incorporated into current deep learning techniques. In order to find an alternative to back-prop, we need to ask: What is the

be based entirely on it. One remaining question is how to represent symbolic data in a "neural" manner. A general form of symbolic data may be as a tree. Taking inspiration from the cortex (14), we may perhaps represent the tree / symbolic data as hierarchically organized neural feature vectors:

Remember that in the Transformer, symbols are organized as sequences,

for example: "spoon · inside · glass." It may be desirable for AGI to have

multiple levels of features, such as "spoon" and "glass" on a lower level,

Juxtaposed side by side, the Transformer and the cortex seem to have many

 $\xrightarrow{\text{features}}$  high-level

Every human can recognize this as "putting a spoon into a glass", a sym-

bolic representation. Many researchers may have under-estimated how

much the brain uses symbolic reasoning, and my proposal is that AGI can

function (19)working memory Based on this understanding, we need to figure out how to design the next version of Transformer and incorporate it into our AGI architecture.... Abductive reasoning Abduction has been relatively neglected in AGI research, which had focused

on forward inference. Recently there is a call to study this important aspect.

In logic, abduction means finding the explanation for some known facts.

An explanation E is simply some propositions that imply the known fact F,

For example, why do we think a certain actress, say Marilyn Monroe, is

"sexy" <sup>3</sup> ? That's because we recognize she has some features (visual or

otherwise, no need to enumerate them explicitly) that we consider sexy.

So,  $E_1 \wedge E_2 \wedge .... \Rightarrow$  Sexy. Those conditions **imply** she is sexy, and they are

Why is abduction important? For example, when a waitress says "The Ham

recipient of the

**Award** 

**ACL Lifetime** 

Achievement

Also recall that our reinforcement learning model consists of just the state

transition

Such resonance behavior can be viewed as the system seeking to minimize an energy, ie, trying to find the "best explanation" to a set of facts.

picted in diagram (14). We can further abstract this with the following diagram, where F and G are not functions but **optimization constraints**:

When a system has both forward and backward connections, it forms a loop and its dynamics is likely to produce "resonance". This harks back to the ART (Adaptive Resonance Theory) proposed by Grossberg and Carpenter beginning in the 1980s.

 $F \left( \begin{array}{c} \\ \\ \end{array} \right) G$ 

If the input  $\vec{x}$  produces the output  $\vec{y}$  after some iterations, then it is likely

that the output  $\vec{y}$  would produce  $\vec{x}$  in the inverse direction. In other words,

we have a **neural** mechanism that implements a function f and its inverse  $f^{-1}$ . The significance of this (from the **learning** point of view) is that we only need to learn the function f and we get  $f^{-1}$  for free. In logic, if forward inference is denoted as  $\vdash_{\mathbf{m}}$ , where  $\mathbf{m}$  is a set of logic rules, then abduction is  $(\vdash_{\mathbf{k}})^{-1}$ . Abductive interpretation is basically a constraint-satisfaction process that uses inference rules in both direc-

This is also corroberated by neuroscientific evidence: areas in the cerebral cortex are replete with both forward- as well as back-projections, as de-

Dealing with assumptions

**SC**Viterbi So abductive reasoning is basically just bidirectional inference.

tions.

↑ attention attention query value key

input #1

the LHS to the RHS. Relation to the biological brain There are two distinct aspects in the brain: • Short-term or Working Memory is the electric activation of neuronal populations. • Long-term memory is stored as synaptic strengths, established by synaptic formation and strengthening. The transfer from STM to LTM is called memory consolidation.

cortical "laver"

Bi-directional connections in the cortex

forward and backward) connections <sup>2</sup>:

backprojections

Alternative to back-propagation?

symbols.

anism.

gorithm.

structures that are learned, but the learning algorithm is too slow because it uses combinatorial search (reminiscent of NP hardness). But the brain must have a roughly equally powerful learning mechanism as back-prop. A likely candidate is resonance. In figure (14) we have a hierarchically connected cortical structure. What we need is some sort of "infinitesimal" learning rule. Hierarchy of features If we consider relations between objects, for example, "spoon inside a

Softmax corresponds to lateral inhibition. The Transformer has many layers because it unfolds along the time axis the training of a recurrent network part of the reason why the Transformer is very efficient. Each hidden layer of the Transformer can be construed as a "stage" of logical inference: input  $\vdash$  stage<sub>1</sub>  $\vdash$  stage<sub>2</sub>  $\vdash$  ....  $\vdash$  output.

ie,  $E \Rightarrow F$ .

and its transition function:

similarities:

and "inside" on a higher level.

**Transformers** 

Sandwich left a big tip", Ham Sandwich here refers to the customer who ordered it (an example of metonymy). The AI knows the plain facts such as that someone ordered a ham sandwich, and then it abduces that the most likely interpretation of the phrase "Ham Sandwich" is as the person associated with it. This is the basis of Abductive Interpretation of Natural Language proposed by Jerry Hobbs:

the explanation for her sexiness.

4