

# AGI standard model

## — a proposal

YKY

March 26, 2022

## 0 Introduction

The “standard model” is a way of thinking, that may help us better understand the general theory of AGI systems.

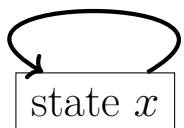
The essence of the standard model is just to identify a **Working Memory** or “state” of the AGI system.

One benefit of our theory is that it relates Transformers / BERT / GPT to AGI 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. There are other benefits.

## 1 Reinforcement learning

This is the simplest form of a **dynamical system**:

transition function  $F$



= “working memory”

(1)

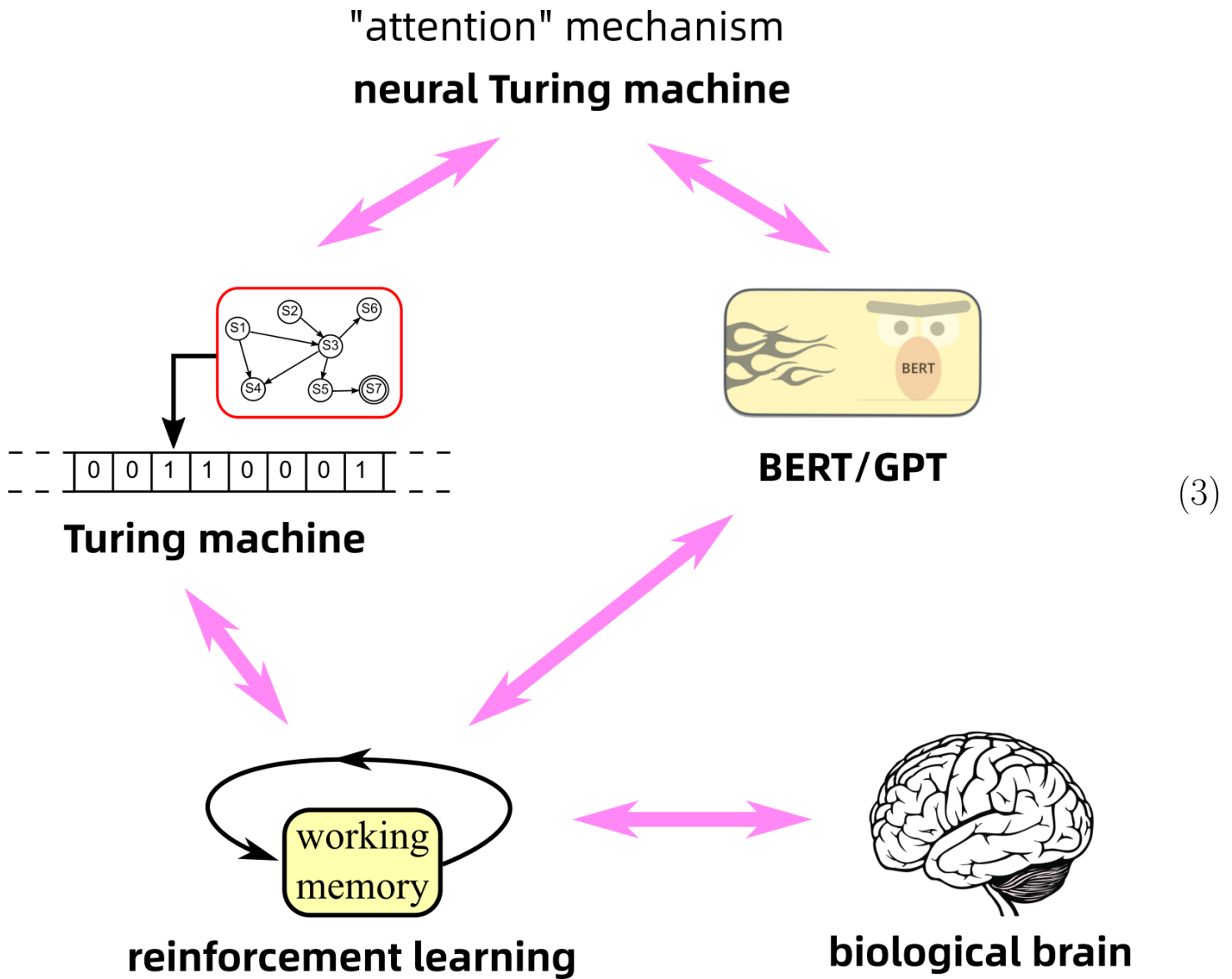
When we add an “action” or “control” variable  $u$  to it, it becomes the most basic **control system**:

$$\begin{array}{c} F(x, u) \\ \curvearrowright \\ \boxed{x} \end{array} \tag{2}$$

which is the setting for Dynamic Programming or **Reinforcement Learning**. The optimal solution for such systems is governed by the **Hamilton-Jacobi-Bellman equation**.

It is because this setup is so simple, that I am quite confident to call it the standard model.

## 2 The standard model in relation to other systems



## 3 Neural Turing Machine and BERT

The **attention mechanism** was first proposed in the “**Neural Turing Machine**” paper by Graves et al.

## 4 Relation to the biological brain

One theory has it that the prefrontal cortex maintains a number of “thoughts” with sub-populations or, perhaps, with **micro-columns**. These activated sub-

populations are in competition with each other, through **lateral inhibition**. The thought(s) that win are the thoughts we retain – they “make sense”.

There are two distinct aspects: •

**Short-term** or Working Memory is the **electric activation** of neuronal populations.

Whereas **long-term** memory is stored as **synaptic strengths**, built up through the processes of synaptic formation and strengthening.

## 5 Abductive reasoning

Abductive reasoning is basically just **bidirectional** inference.

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.

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.

## 6 Dealing with assumptions

Example of an assumption: “If I play move  $x$  now, I will checkmate in 3 moves”.

Suppose  $M, N$  are proofs of  $M : \phi \rightarrow \psi$  and  $N : \phi$ . Then the proof of  $\psi$  would be the application of  $M$  to  $N$ , denoted as  $@(M, N)$  or simply  $MN$ .

The assumption rule (Ax):

$$\Gamma, \phi \vdash \phi \quad (\text{Ax}) \tag{4}$$

for example can be written as:

$$x : \phi, y : \psi \vdash x : \phi \quad (\text{Ax}) \tag{5}$$

which is why we say that an assumption is a  **$\lambda$ -variable**.

**Discharging** an assumption (ie, using the  $\rightarrow$ I rule) results in a  $\lambda$ -term.

An AGI needs the ability to place an implication  $\phi \rightarrow \psi$  into working memory, and to prove it using the ( $\rightarrow$  I) rule, ie, by making an assumption.

## 6.1 Tic Tac Toe example

Assume the current board is  $\begin{array}{|c|c|c|} \hline \text{O} & \text{X} & \text{O} \\ \hline \text{X} & & \\ \hline & & \\ \hline \end{array}$  and it's **X**'s turn to play.

**X** can do the “double fork” by playing  $\begin{array}{|c|c|c|} \hline \text{O} & \text{X} & \text{O} \\ \hline \text{X} & & \text{X} \\ \hline & & \\ \hline \end{array}$ .

But how can an AGI know (or prove) this?

If the current board is  $\begin{array}{|c|c|c|} \hline \text{O} & \text{X} & \text{O} \\ \hline \text{X} & & \text{X} \\ \hline & & \\ \hline \end{array}$  then a double fork exists. We need a predicate to detect double forks.

We need to reason that even if **O** plays the “blocking” move  $\begin{array}{|c|c|c|} \hline \text{O} & \text{X} & \text{O} \\ \hline \text{O} & \text{X} & \\ \hline \text{X} & & \text{X} \\ \hline \end{array}$ , **X** can still win.

We can easily express the conditions for **X**-can-win, but the difficult part is to make the assumption in **red**, in other words:

$$\text{red-move} \rightarrow \text{X-can-win} \tag{6}$$

The difficulty lies in that the LHS is **not true** under the current facts. This conditional statement must be proven by, first, assuming the LHS, and then deriving the RHS. Then the assumption is **discharged** and the conditional statement is proven, via the ( $\rightarrow$  I) rule.

In the old days, in classical AI, assumptions are handled with **Truth Maintenance Systems** that keep track of inference traces symbolically. These systems can get quite complicated with the need to track multiple assumptions. For example, when we plan a bank robbery, we need to consider many possible forking scenarios.

**Algorithm:** Put the assumption  $A$  in Working Memory. Make inferences, marking all conclusions with  $A \rightarrow *$ . When enough conclusions are obtained, remove the assumption  $A$  from Working Memory.