

AGI standard model

— a proposal

YKY

March 16, 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.

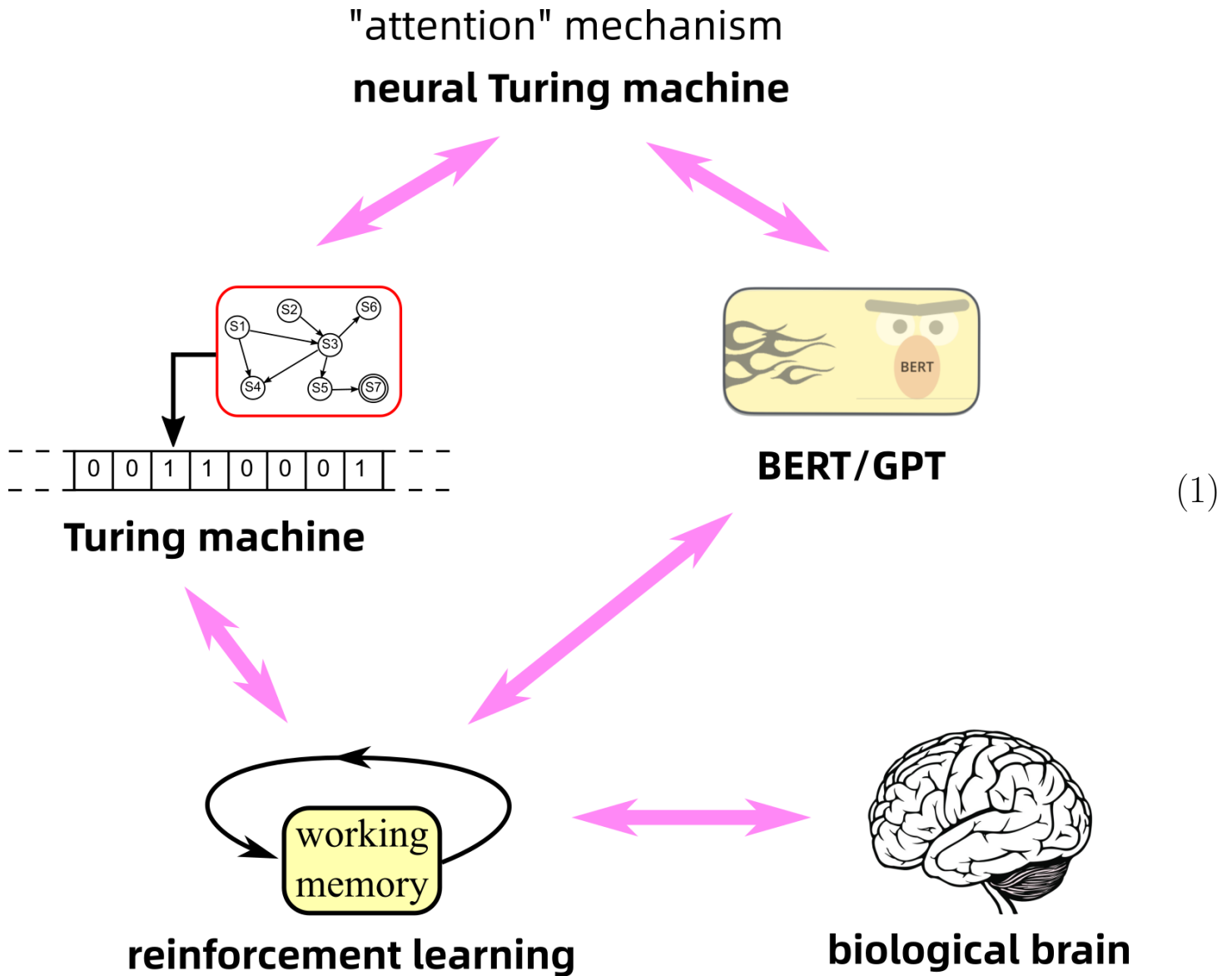
1 Reinforcement learning

The simplest form of a **dynamical system**:

When we add an “action” or “control” variable to it, it becomes the simplest **control system**:

which is the setting for Dynamic Programming or Reinforcement Learning.

2 Standard model



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

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”.

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 ψ would be the application of M to N , denoted as $@(M, N)$ or simply MN .

The assumption rule (Ax):

$$\Gamma, \phi \vdash \phi \quad (\text{Ax}) \tag{2}$$

for example can be written as:

$$x : \phi, y : \psi \vdash x : \phi \quad (\text{Ax}) \tag{3}$$

which is why we say that an assumption is a **λ -variable**.

Discharging an assumption (ie, using the \rightarrow I rule) results in a λ -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, making an assumption.