

Combining LLM and RL, and Logic Transformer

King-Yin Yan^[0009-0007-8238-2442]

general.intelligence@gmail.com

Abstract. This paper has two main ideas: The first is to explain two types of LLM (large language model) + RL (reinforcement learning) architectures, which are probably known on the internet but not explicitly stated. The second idea is developed based on what we call “Type L” architectures. Under this perspective, formal logic is essentially the same as natural language, and so LLMs (language models) are the same as logic processors. Thus we introduce the Logic Transformer which may have advantages over traditional Transformers.

Keywords: AGI · large language models · reinforcement learning · neural-symbolic integration

Part I. LLM + RL architectures

For “string diagrams” there are usually two conventions: 1) data are nodes, functions are edges: $\textcircled{x} \xrightarrow{f} \textcircled{y}$ or alternatively 2) functions are nodes, data are edges: $\xrightarrow{x} \textcircled{f} \xrightarrow{y}$. In the following, I make explicit nodes for both functions (grey) and data (white), whereas edges merely represent linkages.

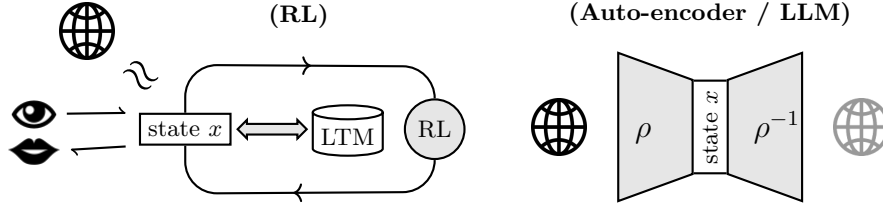


Fig. 1. The eye represents observations and the mouth (speech) actions. Because RL has to maximize rewards, its internal representation (the state) must eventually approach a good approximation of the world. LTM = long term memory, which works by associative recall and condensation, but will not concern us in this paper. The auto-encoder, of which LLMs are a special case, works by compressing world-data (via ρ) and de-compressing (via ρ^{-1}) to re-construct the data (grey world).

First let’s recall the **fundamental forms** of RL (Fig.1 left) and auto-encoders (right). The question is how to combine them, for which we can ignore the LTM

module. Both RL and LLM have the “state” x in common but there are subtle differences as to how to interpret these states (Fig.1). Initially my research favored **Type W**, where the LLM is interpreted as a **world model**. This approach suffers the problem that we don’t really understand the internal state (marked “?”) of the LLM, which consists of many layers of Transformers and their intermediate token outputs, so we don’t know how to “merge” the RL state with the LLM state.

It seems that **Type L** would be easier to work with, and it is also the “mainstream” approach, of which RLHF seems to be an instance. Here, the LLM models the “spoken thoughts” of human thinking, as found in text corpora. These thoughts are interpreted as internal states of the RL. Ever since Richard Montague’s success in converting a fragment of English into formal logic [Montague], the line between formal logic and natural language is blurred. We should not obsess with the idea that “logic” must be some kind of cryptic, undecipherable code. This leads to a novel idea: *the “representation” of human internal thought is just natural language*. Type L also has the advantage that we can directly examine the internal state of an AGI. What are currently called “prompts” is just **Working Memory**.

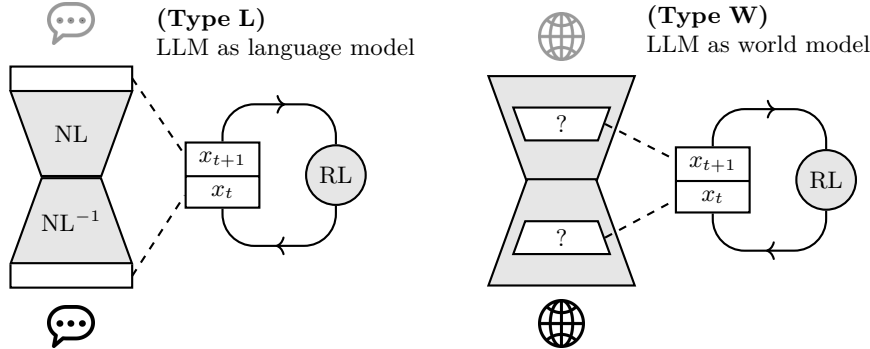


Fig. 2. Two types of RL + LLM architectures.

Part II. Logic Transformer

Having convinced ourselves that LLMs can be a model of *internal* thoughts, we wonder if the Transformer can be re-formulated as some kind of logic engine or general **rewriting system**, defined by rewriting **rules** that can be learned via gradient descent?

In order for such rules to be **differentiable**, we can’t have N rules at one moment and then suddenly $N + 1$ rules the next moment (that would not be differentiable). The only solution I can think of is to maintain a *fixed* number,

M , of rules. Every rule may potentially “morph” into any other rule in this space of rules. The set of rules would look like a *rectangular* matrix:

$$\begin{aligned}
 \text{rule 1: } & \boxed{X_{11}^1 \dots X_{1I}^1} \wedge \boxed{X_{21}^1 \dots X_{2I}^1} \wedge \dots \boxed{X_{K1}^1 \dots X_{KI}^1} \rightarrow \boxed{X_{01}^1 \dots X_{0I}^1} \\
 \text{rule 2: } & \boxed{X_{11}^2 \dots X_{1I}^2} \wedge \boxed{X_{21}^2 \dots X_{2I}^2} \wedge \dots \boxed{X_{K1}^2 \dots X_{KI}^2} \rightarrow \boxed{X_{01}^2 \dots X_{0I}^2} \\
 & \dots \quad \dots \\
 \text{rule M: } & \boxed{X_{11}^M \dots X_{1I}^M} \wedge \boxed{X_{21}^M \dots X_{2I}^M} \wedge \dots \boxed{X_{K1}^M \dots X_{KI}^M} \rightarrow \boxed{X_{01}^M \dots X_{0I}^M}
 \end{aligned} \tag{1}$$

For a rule to be applicable, its atoms must match with propositions in Working Memory. An atom may contain constants (such as “Socrates”, that must be matched correspondingly) or variables (such as x , that matches any entity). To achieve this, I introduce a trick called **cylindrification factor** (Fig.3).

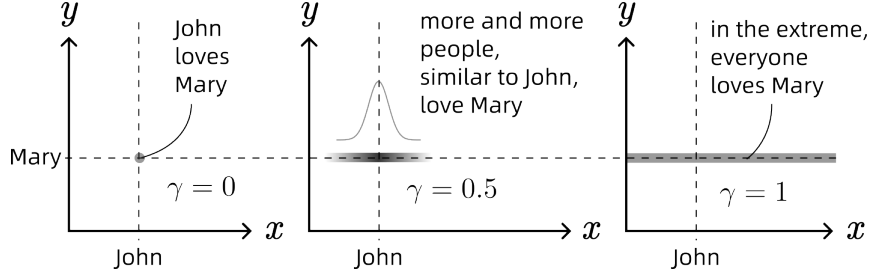


Fig. 3. Illustration of the cylindrification factor γ .

To make **variable substitution** differentiable is even trickier. Consider this logic rule:

$$\text{father}(X_1, X_2) \wedge \text{father}(X_2, X_3) \rightarrow \text{grand-father}(X_1, X_3) \tag{2}$$

We have to imagine some “slots” where we “select” variables using weights and **softmax**:

$$\text{father}(\text{---}) \wedge \text{father}(\text{---}) \rightarrow \text{grand-father}(\text{---}) \tag{3}$$

The actual implementation leads to formulas like these:

$$\begin{aligned}
 & \begin{array}{c} \boxed{\hat{x}^1} \quad \boxed{\hat{x}^2} \quad \boxed{\hat{x}^3} \\ \begin{array}{c} w_{ij}^k \\ \text{---} \end{array} \end{array} \\
 \text{father}(x_{11}, x_{12}, \bullet_{13}) \wedge \text{father}(\bullet_{21}, x_{22}, x_{23}) \rightarrow \text{grand-father}(x_1, \bullet_2, x_3)
 \end{aligned} \tag{4}$$

$$\hat{x}^k := \forall i \in \underset{\text{predicates}}{\quad} \forall j \in \underset{\text{arguments}}{\quad} \left\langle \text{soft max}_{ij} w_{ij}^k, x_{ij} \right\rangle \quad (5)$$

which is similar to the **Self-Attention** of Transformers:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{soft max} \frac{\langle \mathbf{Q}, \mathbf{K} \rangle}{\sqrt{d_k}} \mathbf{V} \quad (6)$$

except that softmax does not commute with inner products, so they are not exactly equivalent, but qualitatively similar. Let’s pause for a moment to reflect that we tried to substitute variables (some kind of syntactic manipulation), and ended up with something like Self-Attention. Also recall that Self-Attention originated as a kind of **content-addressable memory** for Neural Turing Machines, invented for its differentiability.

A final obstacle concerns the limitations of first-order logic. As Ben Goertzel pointed out during an online discussion, my earlier design lacked mechanisms for handling **Skolem functions** and thus existential quantifiers. One way to get around this problem is to define \exists via a set of **higher-order logic** rules. Without explicitly specifying how to do so (though I believe this can be done), we simply use a more powerful “higher-order” logic in the sense that rules admit substitutions in all (predicate or argument) positions, and hope that machine learning will learn the required axioms.

Armed with these ideas, it is not too difficult to work out a new kind of differentiable “Logic Transformer”, the string diagram of which is Fig.4. A detailed explanation is contained in my on-going research thesis [**<empty citation>**].

Conclusions

The significance of combining RL and LLM is such that the system can *think in loops*, eliminating thoughts that contradict accepted truths, thus achieving genuine, logically coherent understanding, instead of just parroting human speeches. This would also solve the problem of “hallucinations”. The remaining question is how to train the Type L model *efficiently*, but this is more of an engineering rather than theoretical problem.

A very interesting question concerning Part II is how to handle variable-length propositions?

