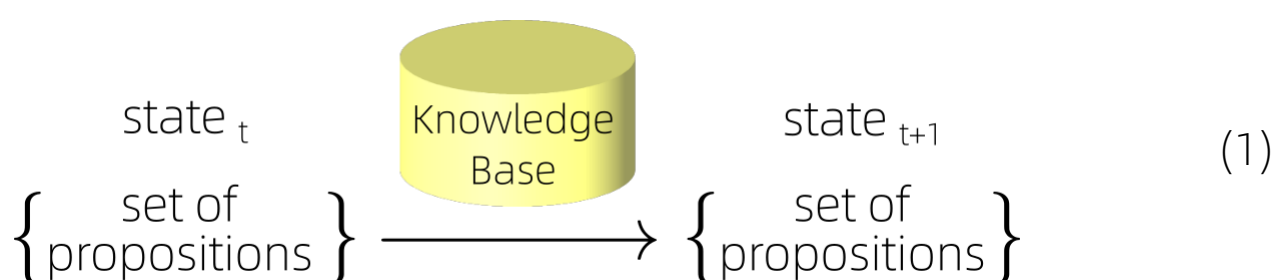


①

Comparison of Logic AI and Deep Learning

This is the most basic operation of classical logic-based AI:



It contains two algorithms:

- **matching** (unification):
Logic rules are conditional propositions involving variables,
eg: $\forall x. \text{human}(x) \Rightarrow \text{mortal}(x)$.
Unification determines whether a rule can be applied to a proposition,
eg: $\text{human}(\text{Socrates})$ can unify with the left side of the above rule.
The goal of Matching is to get an instantiated proposition (ie, specialized, does not contain variables).
- **forward- or backward-chaining** (resolution):
Deduce new conclusions from known facts, or conversely, judge whether a given conclusion can be proven.
eg: $\text{human}(\text{Socrates}) \wedge \text{human}(\text{Socrates}) \Rightarrow \text{mortal}(\text{Socrates})$
From which can be deduced: $\text{mortal}(\text{Socrates})$.

The special thing about deep learning is that it can imitate this inference process:

$$\text{state}_t \vdash \text{state}_{t+1} \quad (2)$$

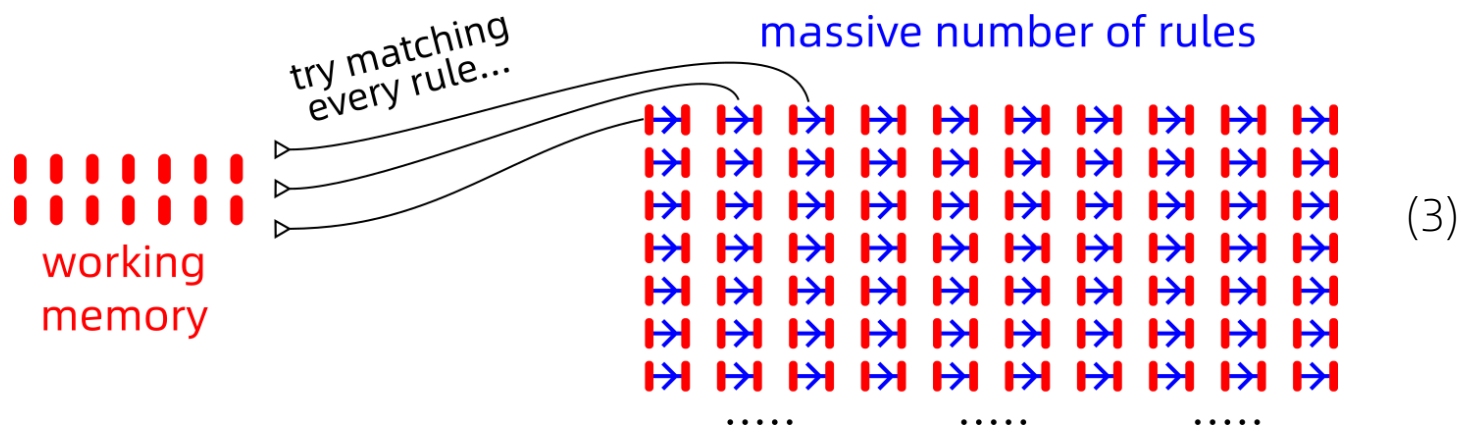
with a highly complicated non-linear function (ie, a deep neural network). By doing so, the logic rules are “mingled” together so that it’s hard to tell them apart. But it is precisely because of this “mingling” that a deep neural network compresses a huge number of combinatorial logic rules into a smaller number of parameters (network weights). It can perform both learning and inference. This simple and crude method is actually extremely efficient, and it is not easy to surpass its speed!

We know (or speculate) that an intelligent system should possess the structure of symbolic logic. Can this insight be used to constrain or accelerate deep neural networks? The answer seems to be yes. The current state-of-the-art CNN (for vision) and GPT (for language) both have specialized internal structure, instead of just being **fully-connected**, and that internal structure is suited to the structure of the data being processed. We have reason to believe that logical structure can be used to constrain deep neural networks to accelerate logical learning.

②

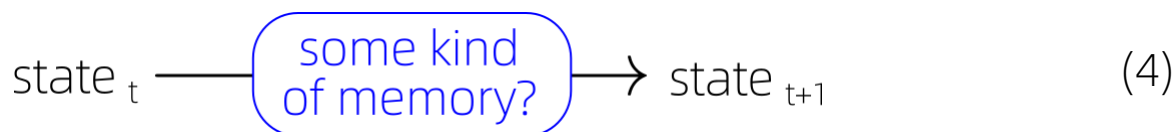
Next, let's look at the logic-based AI architecture in detail.

There are a huge number of rules in the Knowledge Base, and the system needs to match these rules one by one against propositions in the system's state (= working memory):

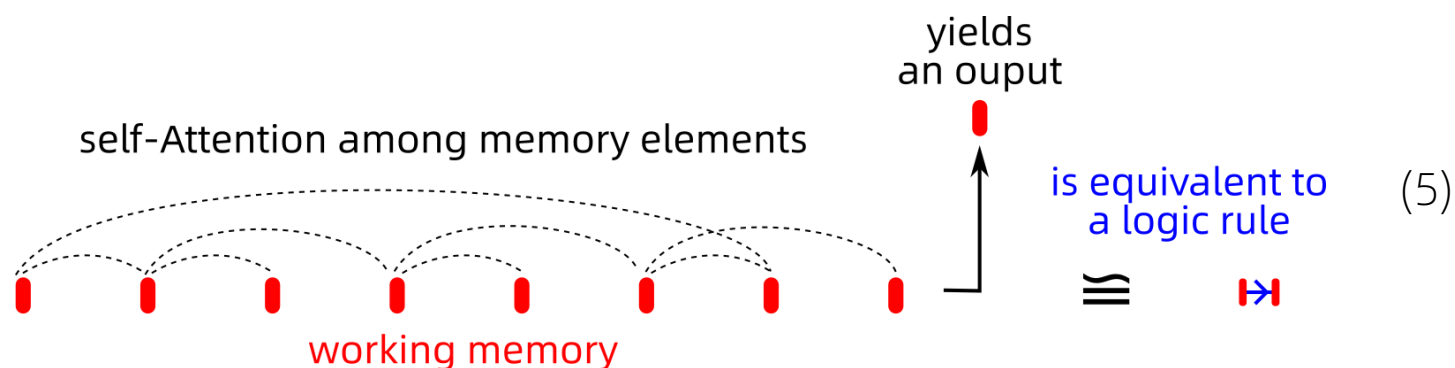


Successfully matched rules generate new conclusions that can be added back to the state / working memory.

This complicated process is entirely replaced by a neural network. Or more abstractly:



For the Transformer, this is a kind of memory stored **between** input elements (stored as the Q, K, V matrices), and it **implicitly** plays the role of logic rules:



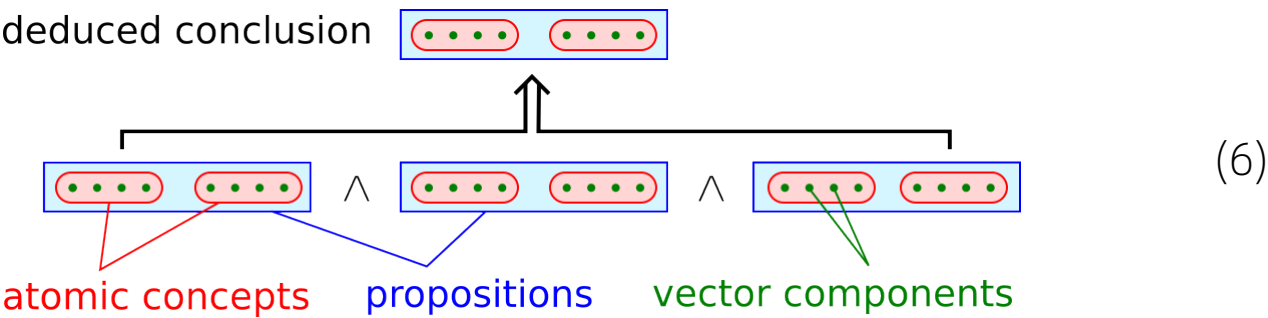
In other words, there are some sort of "distorted" logic rules inside the Transformer. Naturally, we want to find out more structures of logic / logic-based systems. That is, what kind of algebraic structure constrains (4)? To answer this, we can take insights from categorical logic and classical logic-based AI.

We wish to formulate, in algebraic terms, the constraints for (4), but for now it is easier to describe in words:

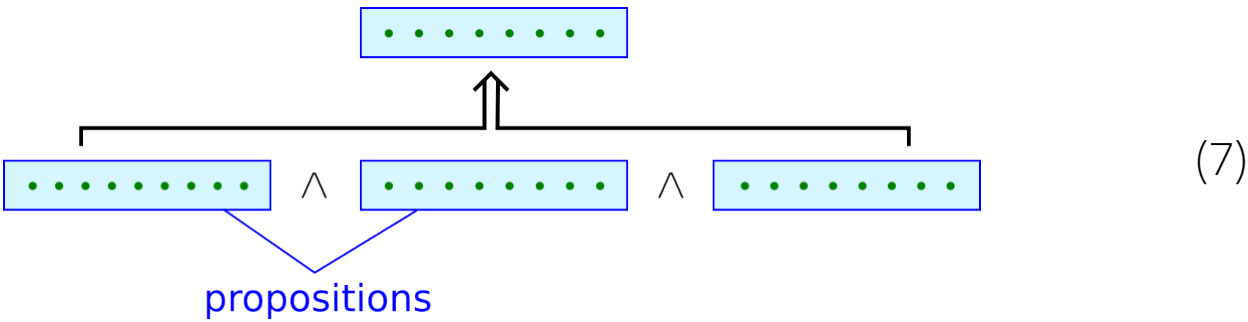
- The AI state is **granular**, as a set of elements, whose order can be permuted, corresponding to equivariance of the Transformer. (Note: Transformers are equivariant, but equivariance does not necessarily require Transformers)
- **Deep structure**: in the sense of having many layers, as composition of functions. The Transformer also has deep structure, with many layers of Self-Attention stacked up.
- Logic has granularity at the **proposition** level and at the **sub-propositional** level. The latter is the structure of **predicate** logic, eg: *loves(John, Mary)* can be represented as a **product** in an **algebra**: *John · loves · Mary*, also called a “word”. The details are unimportant. Our focus is on how to impose this 2-level granular structure onto deep neural networks.
- Each step of logic inference produces **one** new conclusion (or a probability distribution over conclusions), and this new conclusion is added to the old state as an element in a set of propositions, and the old state also needs to **forget** some old propositions, otherwise infinite memory is required. This is slightly different from the Transformer which always outputs the **same number** of tokens as its input. (We are also unsure whether Transformer tokens correspond to propositions or to predicates / atomic concepts).
- A logic rule usually depends on some premises where other premises are **irrelevant**; For example: *talk ∧ dark ∧ handsome ⇒ attractive to women*, where *rich* or *poor* are irrelevant. The Transformer’s **softmax** seems to be a mechanism to exclude irrelevant tokens.
- (There may be other structures.....)


4

In my theory, the ideal logical structure is something like this (the numbers of elements may vary):

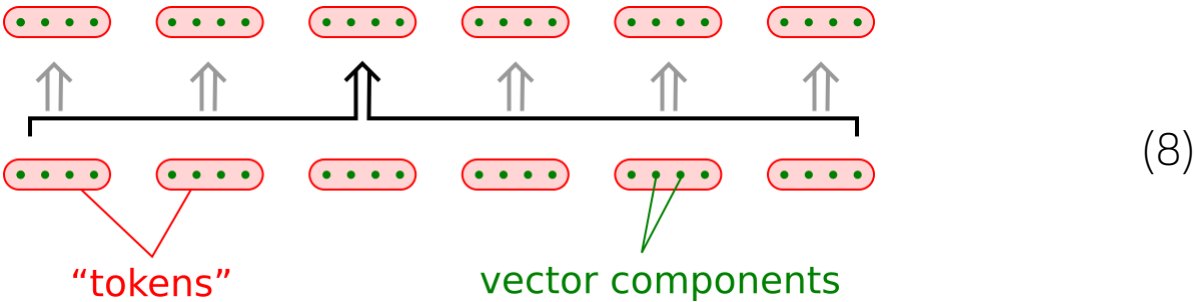


In contrast, the algebraic relation $p \wedge q = p \wedge q$ expresses only this structure:



Compared to figure (7), figure (6) is additionally constrained by the  structure. But how can we express this constraint algebraically?

And this is how the **Transformer** handles propositions and atomic concepts:



It does not represent propositions (= sentences) explicitly, but it uses a special “stop” token to signify the **end** of sentences, and there are other “tricks” such as **positional encoding**. It seems to be a rather *ad hoc* design, we should be able to improve it.

5

Now we try to answer the crucial question: how to express algebraically that “propositions are composed of conceptual atoms”?

That is to say, what is the difference between the following two structures? How to express this difference algebraically?



This is like asking the difference between $0\dots9 \times 0\dots9$ and $00\dots99$ (they are isomorphic).

Similarly,

$$\{ \text{John, Mary} \} \times \{ \text{human, god, worm} \} \quad (10)$$

and the $2 \times 3 = 6$ propositions

$$\{ \text{John is human, Mary is human,} \} \quad (11)$$

are also isomorphic. But the former is a composition of two different concepts, where components can be individually quantified by \forall or \exists ; The latter is propositional logic, where propositions are indivisible and cannot be internally quantified.

But since $\dots \in$ a non-commutative free group (ie, the group with the least structure), it does not possess a simple symmetry like $a \cdot b = b \cdot a$.

After some analysis I arrived at the following condition for “propositions are made of conceptual atoms”:

Atomic Condition (AC). Each proposition P_i is made up of K atoms:

$$P_i = a_{i1} \cdot \dots \cdot a_{iK} \quad (12)$$

where optionally some atoms can be **copied** to other locations (with a non-linear transformation τ , if they are copied to the output layer) via:

$$a_{ih} = a_{jk} \quad \text{or} \quad a_{ih} = \tau(a_{jk}) \quad (13)$$

and the transformation τ has to accord with \forall or \exists as adjunctions to a substitution functor.

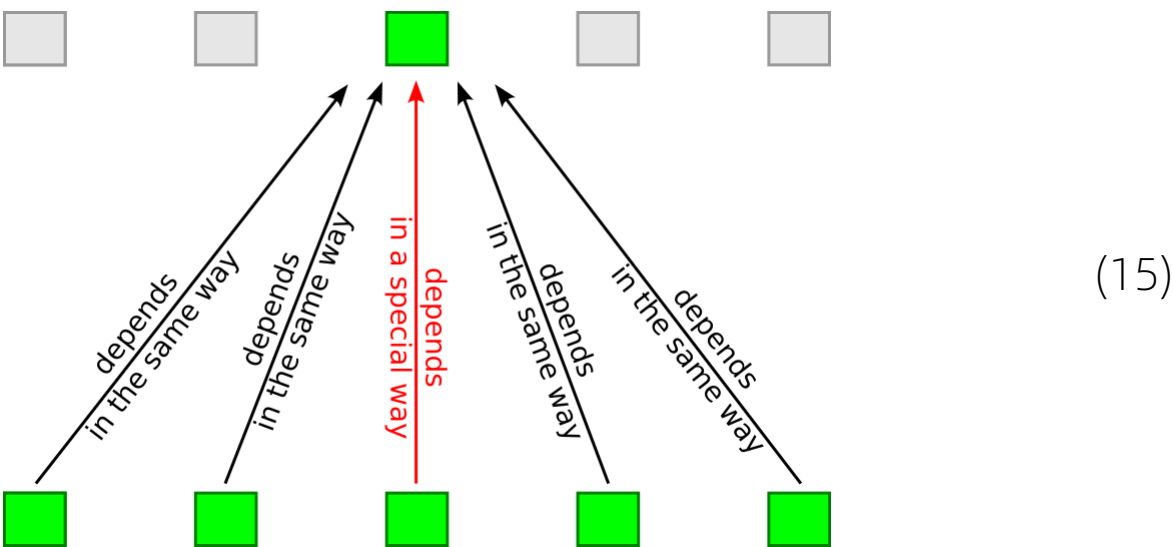
The essence of the Atomic Condition lies in the two **equations** (12) and (13), which are actually very simple. The category-theoretic description of \forall and \exists as adjoint functors is quite advanced, but not essential, and we shall explain them in appendix A. In fact, τ only needs to be a continuous function to meet the above requirement.

So where does the “=” in the equation (13) come from? In fact, it is too obvious, it is just the action of syntactically “moving” variables in a logic rule:

$$\forall X, Y, Z. \text{ grandfather}(X, Z) \leftarrow \text{father}(X, Y) \wedge \text{father}(Y, Z) \quad (14)$$

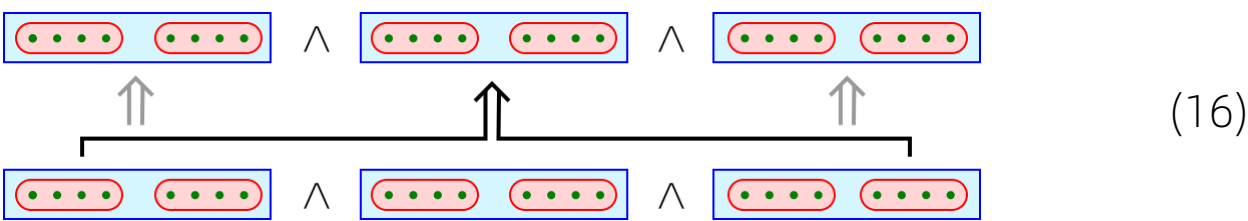
It is due to such “movements” that constitute the structure of “propositions as composed of concepts”.

The essence of **Self-Attention** can be understood as follows (abstract Attention structure):



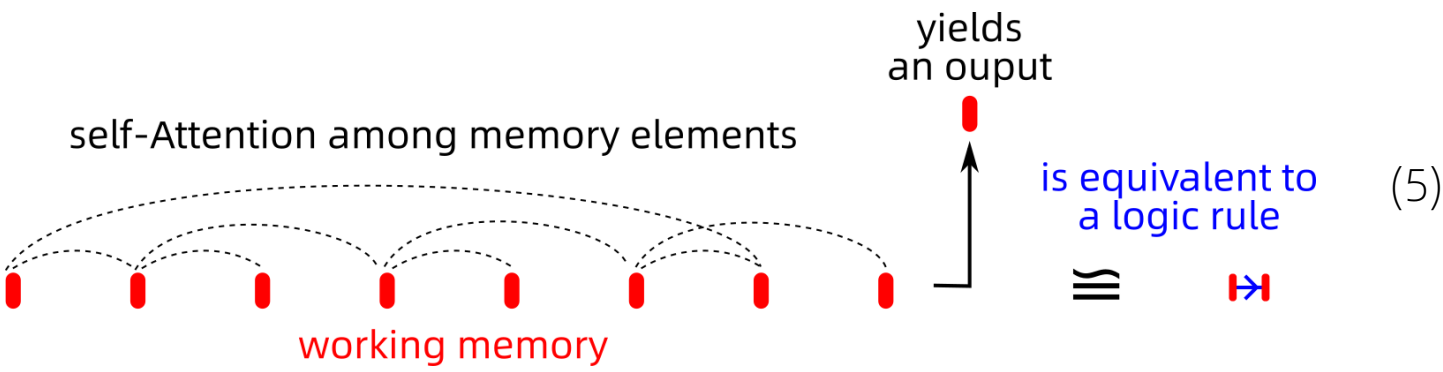
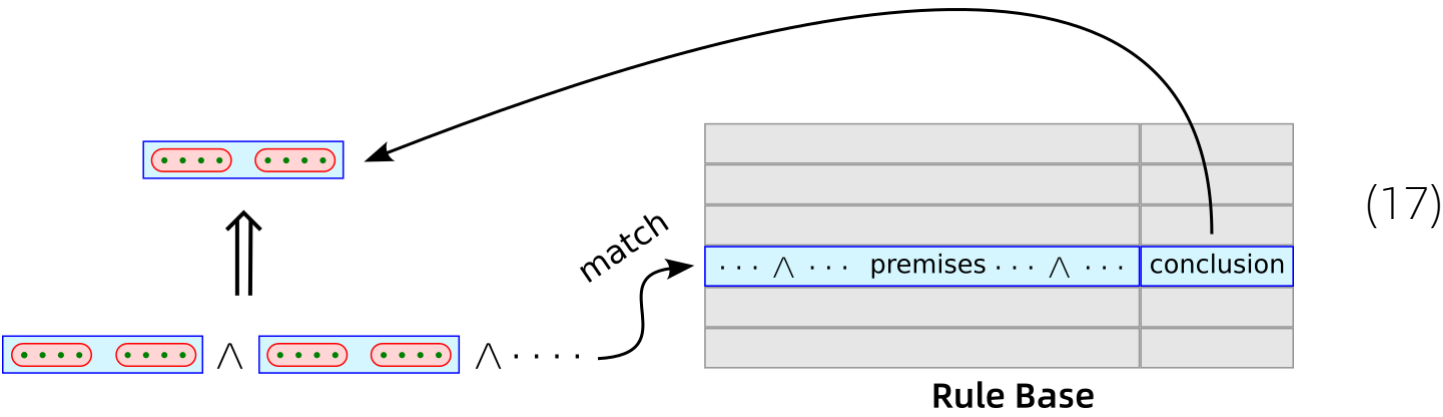
This vertical “axis” (red) is repeated per each input token. So, when the input elements are **swapped**, the output is also swapped. This is how Self-Attention achieves **equi-variance**.

We want to apply a similar method for logical structure:



Here is a very important point¹: The precursor of Self-Attention is the **content-addressable memory** from Graves *et al*’s “Neural Turing Machines”. So we have good reason to regard Self-Attention as a form of **memory**.

We want to compare two approaches. One is the naïve classical rule-base structure, and the other uses Self-Attention instead of a rule base:



You should get the sense that the Transformer is a very “twisted” way of representing rules. The correspondence is so indirect that it is difficult for us to see what the rules look like on the Transformer side. However, I think the designers of the Transformer might have had at least an inkling of its similarities to rule-based systems. In particular, look at the following logic rule:

$$\forall X, Y, Z. \text{ grandfather}(X, Z) \leftarrow \text{father}(X, Y) \wedge \text{father}(Y, Z)$$

(14)

The premise of this rule has two clauses, the variable **Y** that appears twice must be identical (red), for this matching to be considered successful. And this kind of **comparison** operation within the premises of the rule is exactly what Self-Attention can perform conveniently. But Self-Attention also ignores the symmetry of $A \wedge B$, so there may be room for improvement.

¹Thanks “Ziyu” for telling me this important information.

Some key questions to be answered now:

- According to eg. Qian Liu's paper ¹, Transformers often fail to perform certain logical and grammatical operations. Where is the problem? It doesn't seem that Transformers cannot learn the syntax at all, but that it cannot do it only with prompts. But do prompts actually have a deeper meaning, or are they just a hack? We didn't explicitly "tell" the Transformer what it should do, so is it a real shortcoming that it failed to solve those problems? I find it difficult to judge, and the research direction of prompts is shrouded in mystery.
- Now consider the naïve learning algorithm for the graph (17) that is the rule base. This algorithm is of course very slow, since two sets (working memory and rule head) similarity. Assuming that the size of the two sets is fixed at N , then $N \times N$ times of dot products are needed, and this is just a comparison of a rule. All rules need to be added with softmax. When the rule base is very large, this algorithm seems impractical.
- figure (17) may also have the problem of the old logic rule learning algorithm, that is "**plateau problem**". For example, write append in Prolog language function:

```
append(X,Y,Z) :-
```

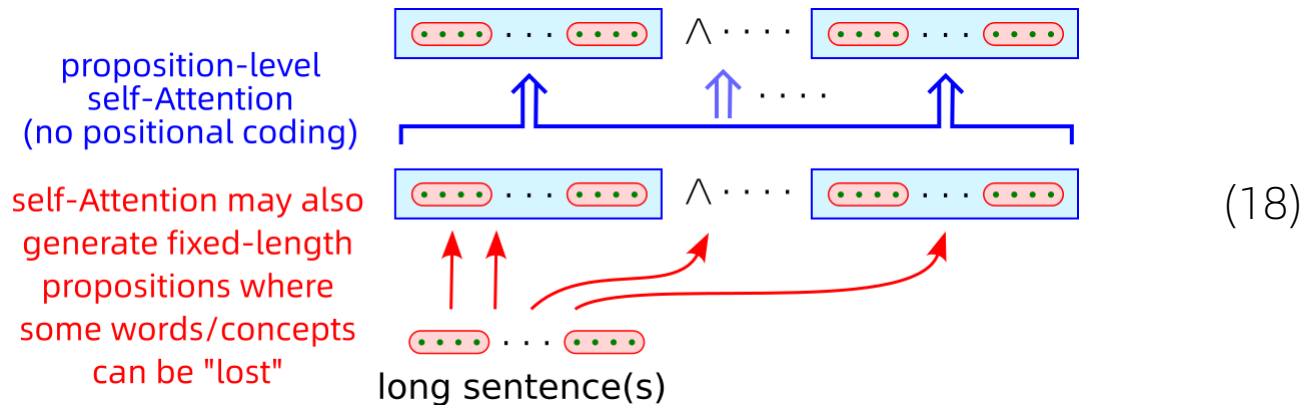
```
list(X), head(X,X1), tail(X,X2), append(X2,Y,W), cons(X1,W,Z).
```

This rule has 5 premises. When the rule is learned, the premise is added one by one, but the "score value" of the rule is always zero, until the last premise is added, the score suddenly rises to 100%. For machine learning, this situation is Terrible. And the Transformer twists the rules together. Will this approach help avoid being trapped in the local minima?

- May have an algorithm between Transformer and naïve rule base, it has a stronger logical structure than Transformer, but uses more similar self-Attention's matrix operation to speed up?

¹Qian Liu, Shengnan An, Jian-Guang Lou, Bei Chen, Zeqi Lin, Yan Gao, Bin Zhou, Nanning Zheng, and Dongmei Zhang. *Compositional Generalization by Learning Analytical Expressions*. Advances in Neural Information Processing Systems 33 (2020).

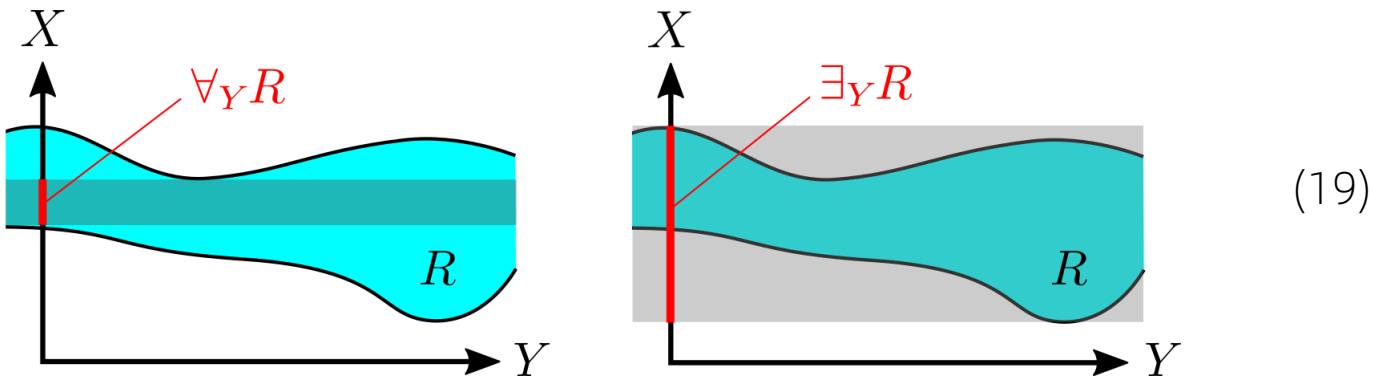
This is a more practical architecture that I can think of:



- In the **blue** structure, each proposition shall consist of fixed # of concept atoms.
- But the abstract self-Attention structure can also allow **variable length** propositions, as long as the **dot product** (similarity) between the propositions can be calculated, and the fixed-length **matrix** memory.
- **red** structure: how to generate propositions of varying number and length from sentences of indefinite length? This is the sequence-to-sequence problem discussed before. This can also be solved with "lossy" self-Attention.

In this appendix, we explain the theory of \forall and \exists in a simple way.

The following is a relation R , such as " Y loves X ", where X, Y are both sets of "persons"; two identical copies. For example, the diagonal represents loving yourself (some people don't love themselves). Note that this graph is not diagonally symmetrical, otherwise there will be no "broken love". Projecting the **projection** of the relation R onto the X axis yields the \forall_Y and \exists_Y sets. The \forall_Y set represents those X that "everyone loves", and the \exists_Y set represents those X that "someone loves him":



The category theory master Lawvere found that \forall and \exists are a so-called **weakening functor** "accompanying functor", the so-called weakening is to expand from a **discourse domain** with only X variables to There is a domain of two variables X, Y , and here Y is purely a **dummy** variable:

$$\begin{array}{ccc}
 & \xleftarrow{\forall_Y} & \\
 (X) & \xrightarrow{\text{weakening}} & (X, Y) \\
 & \xleftarrow{\exists_Y} &
 \end{array}
 \tag{20}$$

For example, $\text{Love}(Y, X)$ is a logical expression with two variables, but $\forall Y. \text{Love}(Y, X)$ actually does not have the variable Y , Because it is **\forall bound**.

The so-called **adjoint** (adjoint) means: There are two categories: left and right. You can move things from the left to the right, and do "**comparison**" in the category on the right, and this comparison can also move things to the left, and the two comparisons are **equivalent**. Here "comparison" means morphism within a category, for example in the category **Set** comparison is set inclusion.

Adjoint functors are not unique, so weakening has two adjoints \forall and \exists respectively.

Lawvere's work made the \forall and \exists quantifiers more general: the "simple" weakening functor is based on the Cartesian product $X \times Y$, but Lawvere expanded it to arbitrary **substitution** functors. (I don't know of examples of this, so I'm not sure what advantages this could bring to our application.)

In categorical logic there are **Beck-Chevalley** conditions and **Frobenius** conditions, perhaps the symmetry we need? But after a closer look, I found that the problem still cannot be solved... For completeness, I will describe it, and you can skip it if you are not interested.

Consider first the easier-to-understand **Frobenius** condition. Logically, it is equivalent to saying:

$$\exists x. [\phi \wedge \psi(x)] \equiv \phi \wedge \exists x. \psi(x). \quad (21)$$

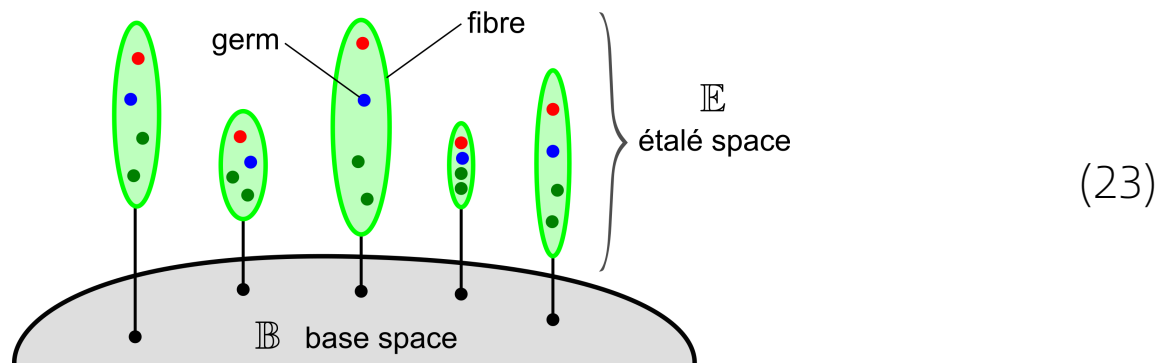
Since classical logic AI generally uses \forall and ignores \exists , I rewrite the above formula as:

$$\forall x. [\phi \vee \psi(x)] \equiv \phi \vee \forall x. \psi(x). \quad (22)$$

But the problem is, the left and right sides of the formula (22), the corresponding neural network (6) is the same (no difference). In other words, this difference may be too subtle, and it does not affect the neural network we actually implement.

As said before, predicate logic leads to **fibration** or **indexing** constructs. The Beck-Chevalley and Frobenius conditions basically say that the fiber structure is “preserved by re-indexing functors”.

Here is a diagram of the fibration structure:



This whole structure is called **bundle**, and **sheaf** is bundle plus some special topology.

fibred product of A and B over I can be defined on top of the two bundles (A, f) and (B, g) , denoted as $A \times_I B$:

$$\begin{array}{ccc}
 A \times B & \xrightarrow{q} & B \\
 \downarrow p & \searrow h & \downarrow g \\
 A & \xrightarrow{f} & I
 \end{array}$$

(24)

where $h = f \circ p = g \circ q$. This is also a **pullback**.

The **Beck-Chevalley** condition says that the following image commute:

$$\begin{array}{ccc}
 K \times J & \xrightarrow{u \times id} & I \times J \\
 \pi \downarrow & & \downarrow \pi \\
 K & \xrightarrow{u} & I
 \end{array}$$

(25)

where π is the projection representing the quantifier \forall or \exists , which are the accompanying maps of the weakening map π^* .

The Beck-Chevalley condition is not entirely hollow; it may not hold. There is a counter-example from Pitts: Consider $X \times Y$, where $X = Y = \mathbb{N} \cup \{\infty\}$ is the natural number plus ∞ as top element; but Y uses discrete order, that is, all orders are $=$ order. A is the relationship on $X \times Y$: $A = \{(x, y) \in \mathbb{N} \times \mathbb{N} \mid x \leq y\}$. Then $\exists y. (x, y) \in A$ will be the entire set of X . If we consider the DCPO category, we require the fibration of Scott-closed subsets (ordered by inclusion) over DCPO. The condition for Scott closure of $\exists y. A$ is that it is a lower set closed under directed joins; and this Scott closure condition seems to be violated, thus causing the graph (25) not to commute. (I do not understand the details of Scott closure)

First express the self-Attention structure in a functional way:

output proposition O_i is composed of atoms b_i

$$O_i = [b_1 \dots b_K] \tag{26}$$

input proposition P_i is composed of atoms a_i

$$P_i = [a_1 \dots a_K] \tag{27}$$

Self-Attention

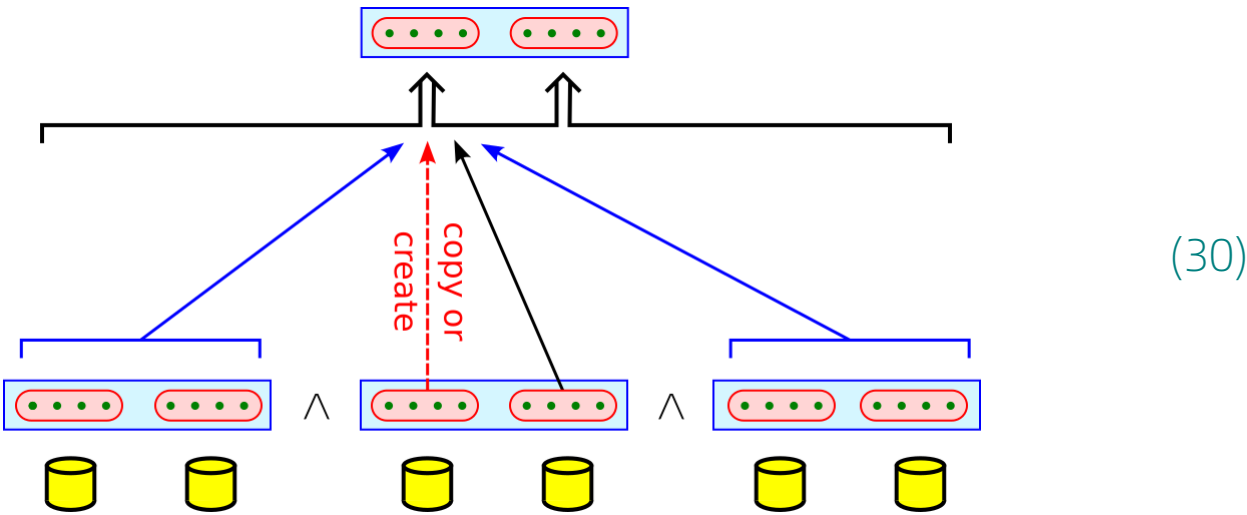
$$O_i = \alpha(P_i ; P_1 \dots \hat{P}_i \dots P_N) \tag{28}$$

$P_1 \dots \hat{P}_i \dots P_N$ means $P_1 \dots P_N$ except P_i .
 $\alpha(P_i ; \dots)$ is the function structure of the graph (15), and it can also be understood as the self Attention with P_i as the query.
How to measure similarity between (P_i, Q_i) ?

$$\arg \min_i \sum_i \min_j \langle P_i, Q_j \rangle \tag{29}$$

- First notice that the pivot structure is only useful at the propositional level, and its effect does not extend to the conceptual atomic level. However, since Transformer does not take full advantage of the exchange invariance, it becomes very efficient because it uses matrix multiplication, so from an efficiency point of view, there are also reasons to extend the axis structure to the atomic level.
- Another idea to try is: direct hard-code copying mechanism. How can this be done with Attention? It has been analyzed before, copy is not easy, because winner takes all.
- but purely using Hopfield network lacks depth. But it seems that in the RL scenario, **breadth** is also important.
In fact, as long as there is an input/output functional relationship,
- is equivalent to a KB library with logic rules. The question is what method it uses to reach its conclusions.

Self-Attention already meets our requirements, but we want to improve it. There are two main ideas: one is a more direct copy mechanism, and the other is a more detailed function dependency.



Copy:

- Copy requires no special mechanism, the output is the input. But the question is how to combine with “create” operation.
- can of course use the familiar softmax: α copy + β create, where α, β are outputs of softmax.
- Another idea is content-addressable. Use table-lookup to find executable rules. But there is a problem with variable matching.

Create:

- What kind of function is needed?