

# AGI from the perspective of categorical logic

YKY

April 2, 2024

## Basics

- The starting point of categorical logic is really the **Curry-Howard isomorphism**, without which you won't be able to understand the sequel.
- Curry-Howard correspondence is the idea of using a mathematical **function**  $f : A \rightarrow B$  to simulate or implement the process of logic deduction, specifically  $A \Rightarrow B$ .
- From this perspective, a logic **proposition**  $A$  corresponds to the **domain**  $A$  of a function. That is, a proposition is akin to something like a **space**.
- Objects in that space are so-called **proof objects**, we use  $\blacksquare$  to denote them.
- It may take a while to get used to, but in fact we see this idea in use almost every day: in **neural networks**.
- A neural network maps certain vectors to vectors. Each vector is a proof.
- The space near a positional vector (under some error margin) ought to represent the **same** concept. So we might as well think of the neighborhood space as a logical proposition.
- This way of doing things is really very obvious and natural.

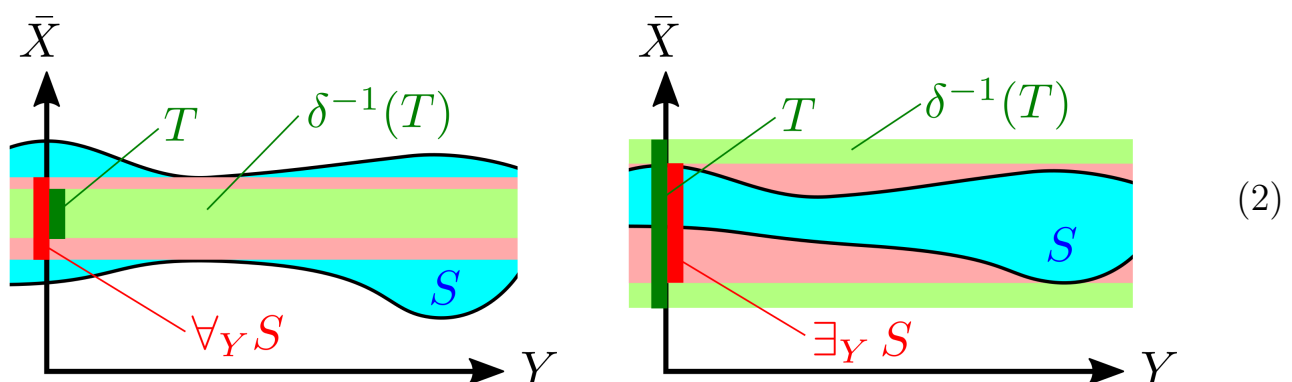
One evidence that this correspondence is on the right track is the truth table of  $A \Rightarrow B$ , which perfectly matches the cardinalities of the function spaces:

$A$	$B$	$A \Rightarrow B$	$B^A$
0	0	1	$0^0 = 1$
0	1	1	$1^0 = 1$
1	0	0	$0^1 = 0$
1	1	1	$1^1 = 1$

(1)

The goal of categorical logic is to use categorical tools as much as possible, to **describe** the structure of logic.

- “Propositions = some kind of spaces” generalizes naturally in category theory to “propositions = **objects** in a category”.
- We use **products** in category theory to express logic  $\wedge$  and  $\vee$ , and **exponentiations**  $B^A$  for  $A \Rightarrow B$ . The latter are also **morphisms**  $A \rightarrow B$  in the category.
- $\forall$  and  $\exists$  are described as **adjoints** to certain variable-substitution maps. For example in  $\forall x. \phi(x, y)$  the quantifier  $\forall x$  projects the space of  $(x, y)$  down to  $(y)$  only, so the resulting expression is no longer about  $x$ .

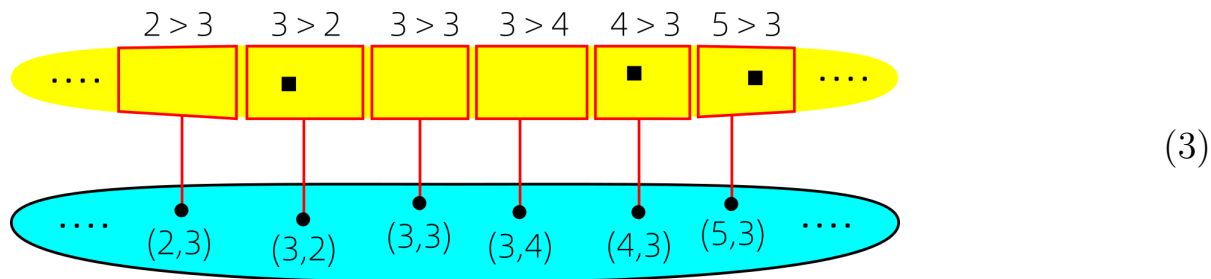


(This part is a bit complicated, but I have explained it elsewhere. The important thing is to understand the overall concept and not get lost in the details just yet)

①

# Where is GPT?

Logic **predicates** are described as **fibrations**, that is, an upper space “indexed” by the base space:



Here we are considering the binary relation “>”, so the indexing space consists of pairs of natural numbers  $\mathbb{N} \times \mathbb{N}$ . We can form relational propositions  $> (a, b)$ , and because of Curry-Howard, these propositions are “spaces”, ie, yellow squares above. Each square is a proposition, which may or may not have a proof (■).

The union of all the yellow squares above is a **sheaf**, which is the space  $\mathbb{L}$  of propositions. **GPT** is a **logic consequence operator** that maps propositions to propositions. But note: GPT is a **set-valued map**. Its domain is not  $\mathbb{L}$  but the the power set  $2^{\mathbb{L}}$  or  $\mathcal{P}(\mathbb{L})$ :

$$\begin{array}{c} \text{space of} \\ \text{propositions} \end{array} = \mathbb{L} ; \quad \begin{array}{c} \text{GPT} \\ \parallel \\ 2^{\mathbb{L}} \xrightarrow{f} 2^{\mathbb{L}} \end{array} \quad \begin{array}{c} \text{GPT} \\ \text{is here} \end{array} \quad (4)$$

Knowing where GPT fits into the scheme of things, provides some clarity. At least for me, because I’m very familiar with **logic-based AI**, I tend to understand the mathematics from this perspective.

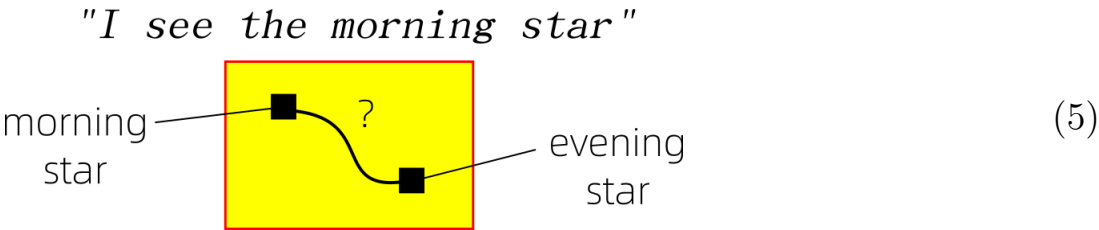
# What is HoTT?

As Curry-Howard suggest that propositions are some kind of “spaces”, can they have topological structure? This is the original idea of **HoTT** (homotopy type theory) proposed by Voevodsky.



From my limited understanding, a proposition either has a proof or not have a proof. If there is a proof, there is no difference between one proof or another. But HoTT posits that they can differ. In the internal space of a proposition, two **path-connected** points (proofs) are regarded as identical, but there can be regions that are not path-connected in this space, so their proofs are regarded as different.

An example: The ancients regarded Venus as the Morning Star and Evening Star, without knowing that they were actually the same star. This is the difference between intension and extension, which can be handled by **intensional logic**, and can be formulated using **modal logic** and Montague semantics. For details, see this article: <https://plato.stanford.edu/entries/logic-intensional/>



For another example, a group can have different **group presentations**, a situation that seems applicable by HoTT.

These spaces that are not path-connected have a **groupoid** structure, with multiple levels such as 1-groupoid, 2-groupoid, ... up to  $\infty$ -groupoid. I am not familiar with these aspects.

As the reader can see, HoTT is concerned with the interior of “truth” (ie, the yellow square), but AGI is mainly concerned with **inference**, which acts on the proposition space (ie, the entire yellow “banana” space).

This is not to say that HoTT is useless for AGI, but its impact is subtle and I cannot judge it yet.

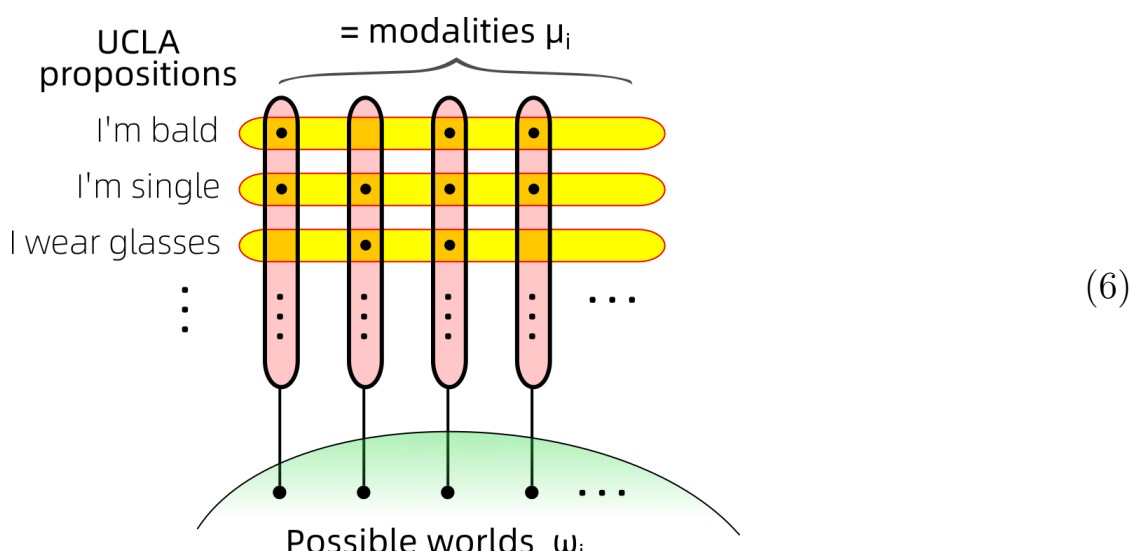
In fact, the **Transformer** has the ability to learn very complex syntactic manipulations, that is, it can implicitly learn the derivation steps of various logics such as modal logic. If so, it seems that we be “bypassed” various special logics without the need to explicitly implement them. But it is also plausible that we by imposing certain logical structural constraints that we can accelerate deep learning. These need to be verified experimentally and are promising research directions.

## Modal logic and possible worlds

Modal logic and possible world semantics are very powerful forms of logic that can handle many problems in philosophical logic. When I studied symbolic logic AI in the past, I mostly ignored it because it was troublesome to implement a modal logic engine on a computer. As AGI gets closer to reality, I feel the need to take a closer look at modal logics.

The main point of this section is that modal logic has a structure of sheaves and toposes. Topos semantics is the “latest” development in logical semantics (though it’s actually quite old, haha). It is said that topos semantics can handle almost all the logical semantics we know of (though I don’t even know what are the exceptions).

The basic picture is this:



- The green set underneath are just **indexes** to each possible world, for example  $\mathbb{N} = \{1, 2, 3, 4, \dots\}$
- Each red “**stalk**” represents a possible world. They form a **fibration** over the base.
- The yellow horizontal bars show the state of truth of each logic **proposition** for each possible world.
- **Modality** = synonym for possible world, a terminology used in John L Bell’s book<sup>1</sup>

This basic structure is laid out in the paper: *Topology and Modality: The Topological Interpretation of First-Order Modal Logic* [Awodey & Kishida 2008].

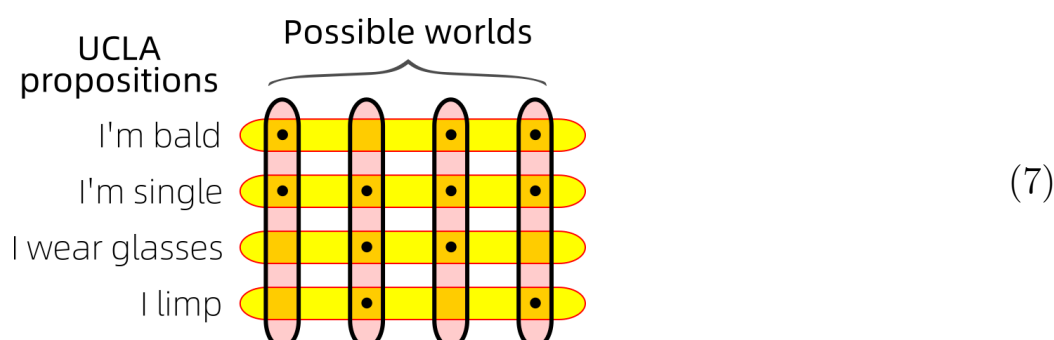
Modal propositional logic can be interpreted by a simple **topological semantics**, while first-order predicate logic can be interpreted by **denotational semantics**. The “product” of these two gives the **sheaf semantics** of first-order modal logic. (Those who have studied programming-language semantics may have heard of denotational semantics.) Let’s briefly introduce these two semantics:

## Topological semantics of modal logic

This was first proposed by Tarski-McKinsey in 1944.

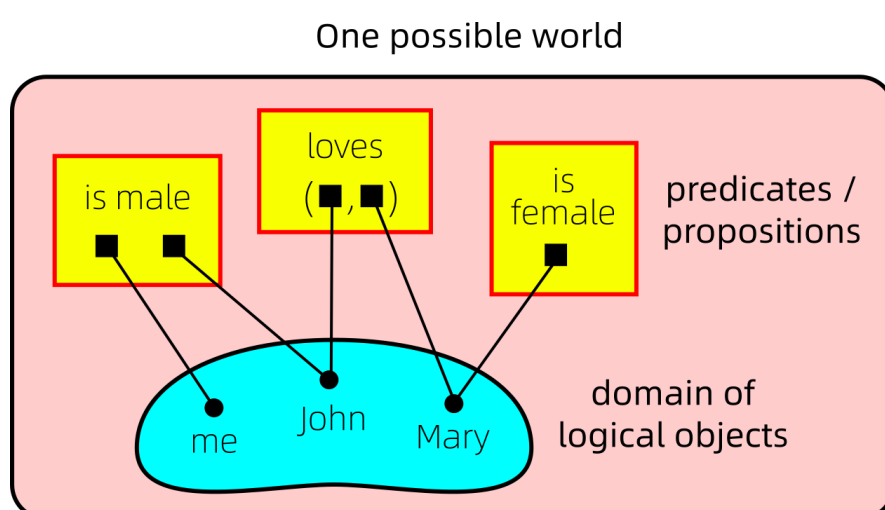
- A possible world corresponds to an **open set** in a topological space
- A proposition is equivalent to the set of possible worlds in which it is true (= set of open sets). These are known as “UCLA propositions” after researchers at the University of California
- The modal operators  $\Box$  and  $\Diamond$  correspond to the topological operations of **closure** and **interior** respectively.

For example, if we only consider the following 4 possible worlds, then “I limp” is not necessarily true, but “I am single” is. The proposition “ $\Box$  I am single” has interior = the entire domain, thus the proposition is true.



## Denotational semantics of first-order logic

This part is the more classic model theory. I don't want to spend too much time explaining it. Basically, it uses a structure to interpret logical propositions.  $D$  is the domain of a logical object, such as { apple, banana, orange, John, Mary }, etc.  $R$  is some relations or predicates, such as " $x$  is a person", " $x$  likes to eat  $y$ ", etc.  $f$  is some functions, such as " $x$ 's mother", " $x$ 's favorite fruit", etc.  $c$  are constants, such as " $c_1 = \text{John}$ " etc. A structure  $M$  contains enough information to interpret the truth or falsity of any proposition.  $M$  can also be regarded as the data of a possible world, where there is only one world.



(For simplicity, I mixed both unary and binary predicates in the same yellow space)

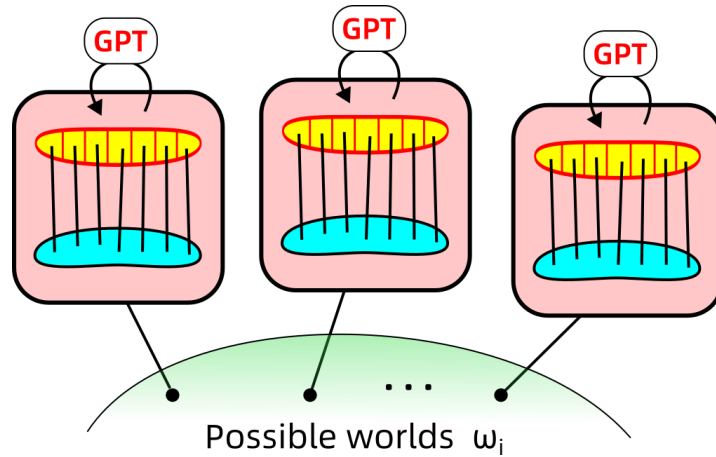
---

<sup>1</sup>Toposes and Local Set Theories [J L Bell 1996]

## Sheaf semantics

The so-called combination seems to be a product, or more simply, lining up the fibers in the direction, which is an addition:

Form a sheaf, and sheaf can also be regarded as topos. The overall image is roughly like this:



(8)

- From the perspective of the human brain, its ability to process "possible worlds" is very limited.
- The most classic example is when playing chess, each predicted move is a possible world. In general, people can usually predict only 3-5 moves.
- Each possible world is actually only one proposition different from the current world. It seems that in practice, there is no need to imagine possible worlds as "monsters".
- Possible worlds are generated dynamically when thinking. We cannot quickly train a GPT according to every possible world, so the above GPT's are copies of the same training results.
- If you want to implement modal logic reasoning, you need to expand the function of the "little GPT" above so that it can handle reasoning in multiple possible worlds. The method I think of for the time being is purely to follow the idea of classic logical AI, such as tagging each possible world with a specific proposition, and then calculating
- When the number of possible worlds is small, the concept of topological closure / interior does not seem to be very enlightening. From a computer point of view, continuous space is very "ideal" and is actually very difficult to achieve.

These are quite intuitive. I will look at them in detail when I have time, but now I suddenly feel that this direction may not be too useful...

## Some technical details on sheaves

⑤

**What's the use of all these to AGI?**