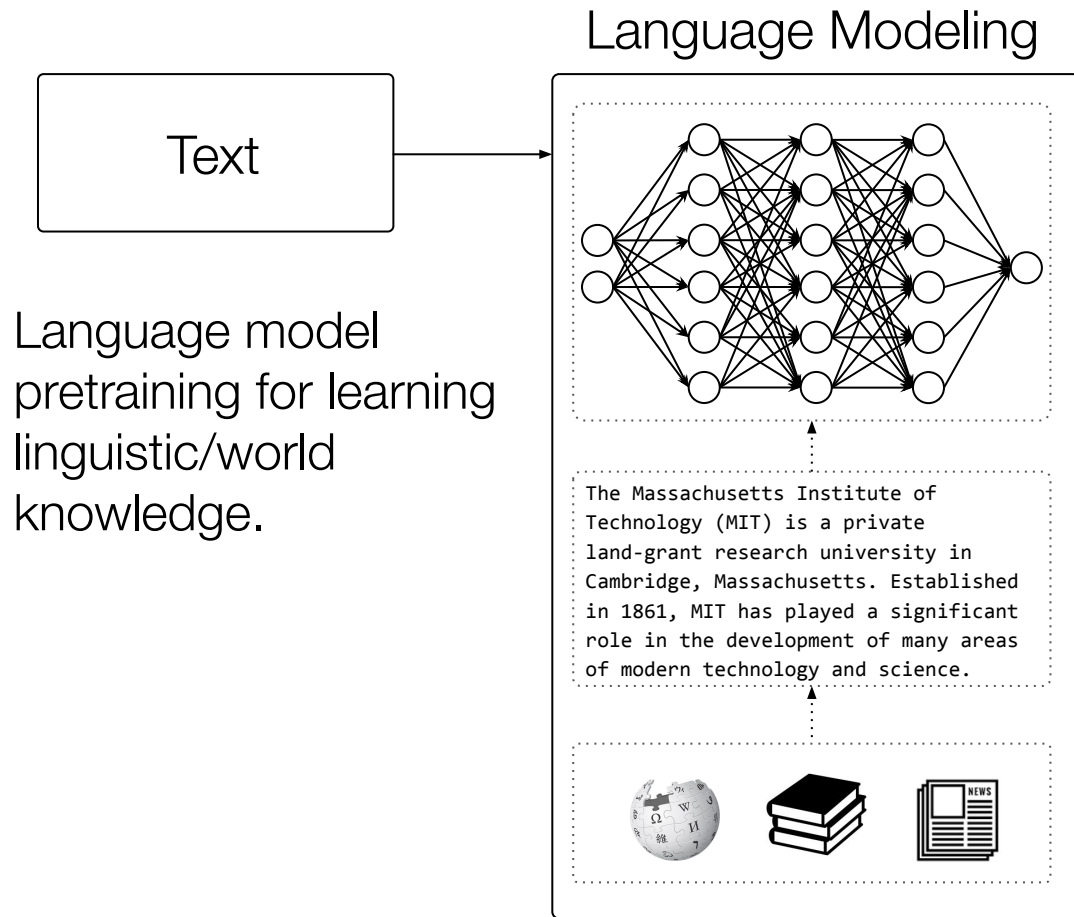
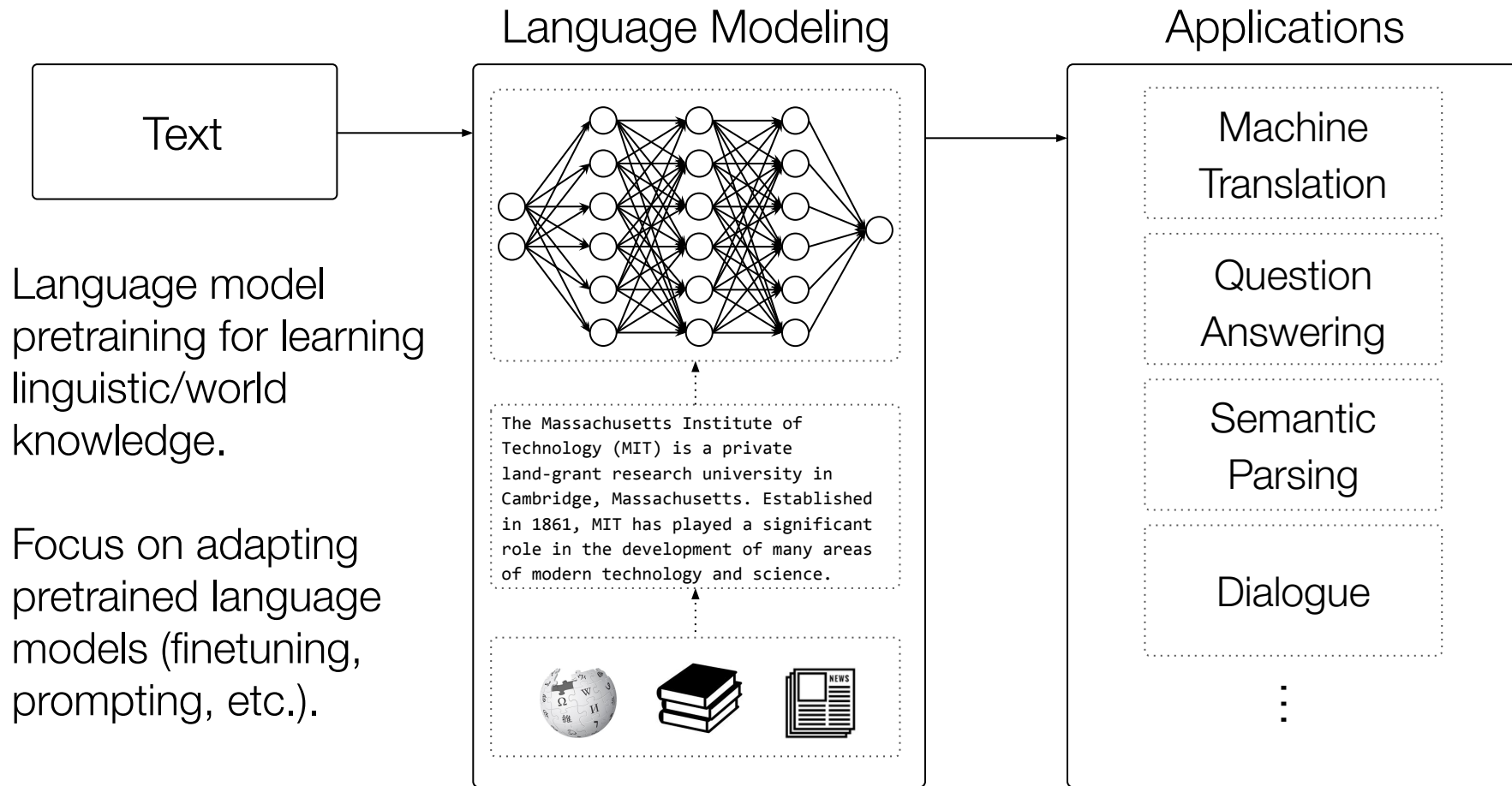


Combining Symbolic Logic with Transformers for Interpretable and Controllable LLMs

The “LLM Paradigm”

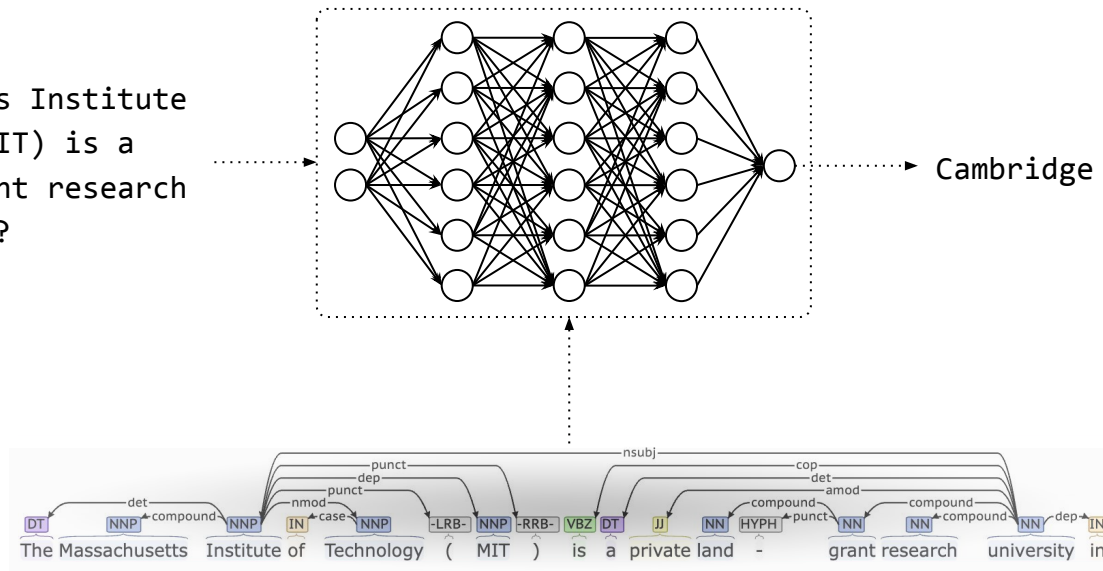


The “LLM Paradigm”



LLMs & Implicit Structures

The Massachusetts Institute of Technology (MIT) is a private land-grant research university in ???



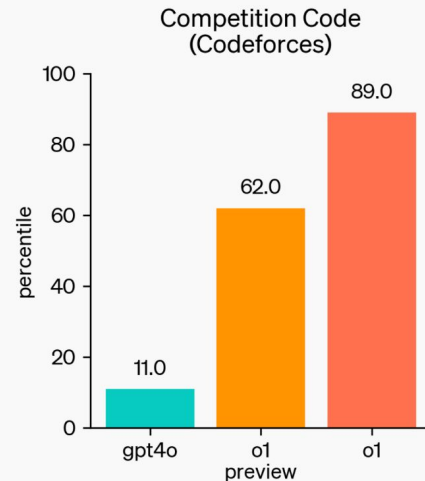
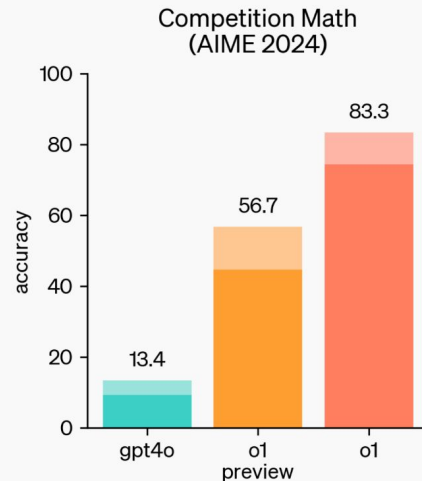
The language modeling (next-word prediction) objective forces the LLM to *implicitly* encode useful structure.

LLMs & *Explicit* Structures?

Is language modeling “all you need”?

$$\arg \max_{\theta} P_{\theta}(\mathbf{x}_t | \mathbf{x}_{<t})$$

Maybe...



LLMs & *Explicit* Structures?

Alice and Bob play the following game. A stack of n tokens lies before them. The players take turns with Alice going first. On each turn, the player removes either 1 token or 4 tokens from the stack. Whoever removes the last token wins. Find the number of positive integers n less than or equal to 2024 for which there exists a strategy for Bob that guarantees that Bob will win the game regardless of Alice's play.

The problem asks for the number of positive integers $n \leq 2024$ such that Bob has a winning strategy. This is a typical impartial game that can be analyzed using the Sprague-Grundy theorem or by directly finding the winning and losing positions. Let W be the set of winning positions (number of tokens n such that the player whose turn it is can force a win) and L be the set of losing positions. [...]

A quantum mechanical particle of mass m moves in two dimensions in the following potential, as a function of the polar coordinates (r, θ) :

$$V(r, \theta) = \frac{1}{2}kr^2 + \frac{3}{2}kr^2 \cos^2(\theta)$$

Find the energy spectrum. Hint: Write the potential in Cartesian coordinates.

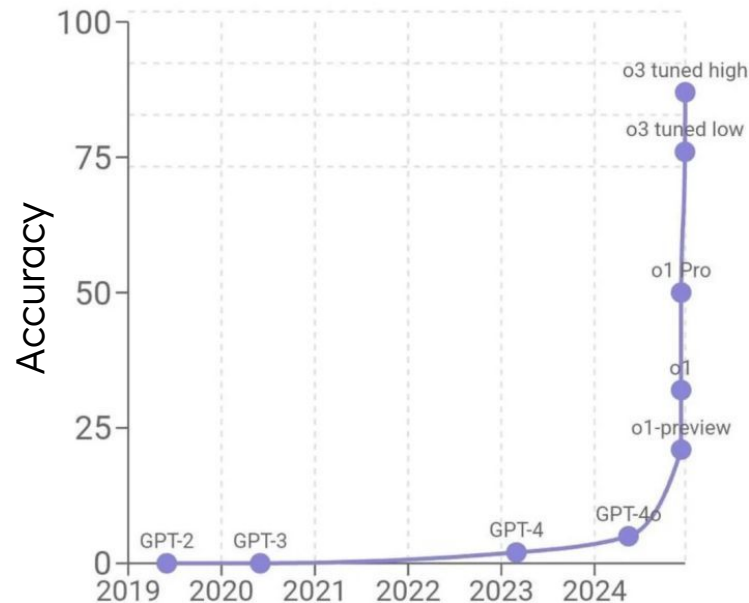
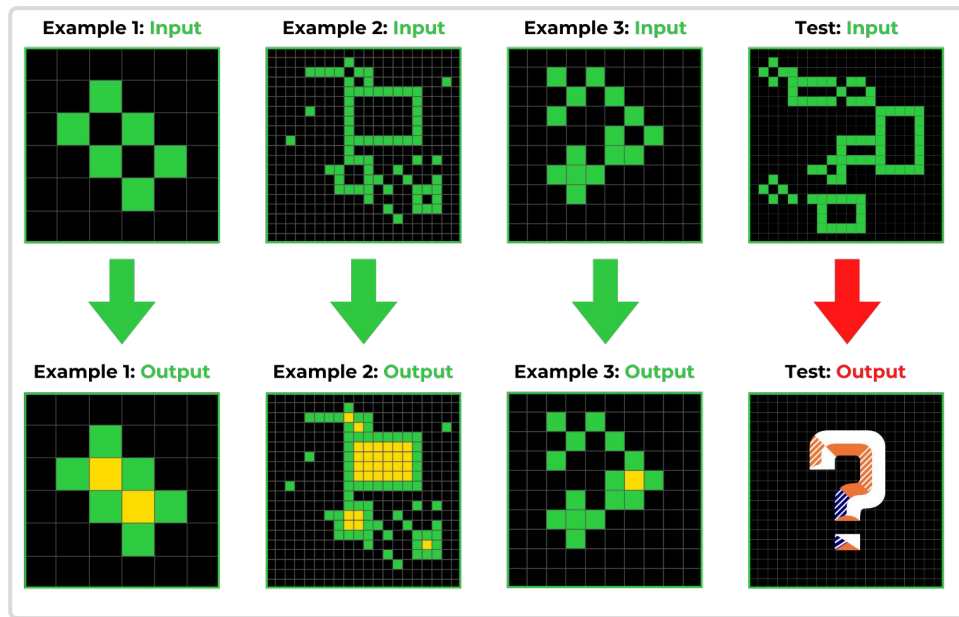
The problem asks for the energy spectrum of a particle in a given 2D potential. The potential is given in polar coordinates (r, θ) :

49.
$$V(r, \theta) = \frac{1}{2}kr^2 + \frac{3}{2}kr^2 \cos^2(\theta).$$

We can rewrite the potential as:

LLMs & *Explicit* Structures?

OpenAI o3 (Dec 2024) on ARC Challenge



LLMs & *Explicit* Structures?

OpenAI o3 (Dec 2024) on ARC Challenge



Nature

How should we test AI for human-level intelligence? OpenAI's o3 electrifies quest

The technology firm OpenAI made headlines last month when its latest experimental chatbot model, o3, achieved a high score on a test that...



MSN

OpenAI's o3 system has reached human level on a test for 'general intelligence'

OpenAI's o3 system got a score of 85% from the previous 55% on the ARC-AGI benchmark. This score is on par with the average human score.



Can you generate a sentence with 5 occurrences of the word "the"?



Certainly! Here's a sentence with 5 occurrences of the word "the":

"The cat, which the little girl found in the alley, was the fluffiest of the strays."



What about 15 occurrences?



Of course! Here's a sentence with 15 occurrences of the word "the":

"The baker, whom the townspeople considered the finest in the region, believed that the flour from the mill by the river, which the locals called the lifeblood of the valley, was the secret behind the exquisite texture of the bread that the citizens, the tourists, and the critics all loved the most."

YO

What about 15 occurrences?



Of course! Here's a sentence with 15 occurrences of the word "the":

"The baker, whom the townspeople considered the finest in the region, believed that the flour from the mill by the river, which the locals called the lifeblood of the valley, was the secret behind the exquisite texture of the bread that the citizens, the tourists, and the critics all loved the most."

YO

I count 17 occurrences :(



Apologies for the oversight. Here's a revised sentence with exactly 15 occurrences of the word "the":

"In the heart of the village, the market buzzed with the chatter of the locals, while the children played near the fountain, the tourists browsed through the crafts, and the elders told stories from the past, creating the tapestry of the community's life."

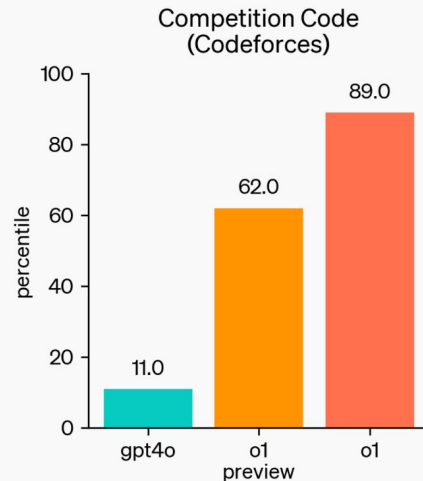
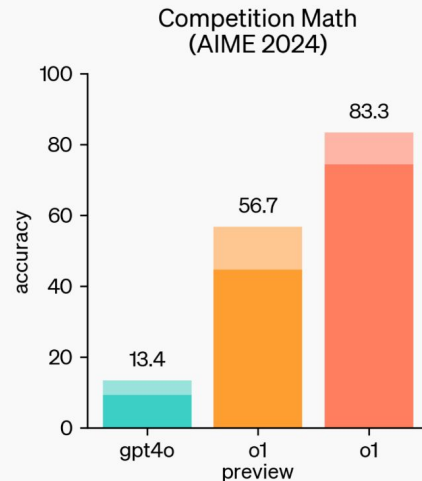
LLMs & *Explicit* Structures?

Is language modeling “all you need”?

$$\arg \max_{\theta} P_{\theta}(\mathbf{x}_t \mid \mathbf{x}_{<t})$$

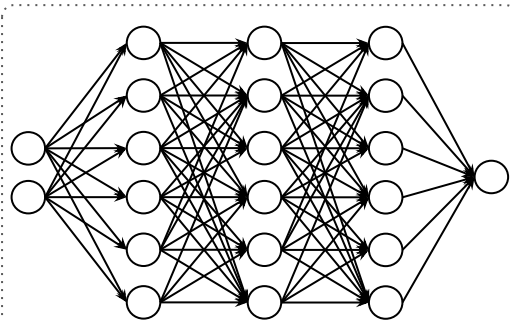
Maybe...

But not quite yet!



LLMs vs. Symbolic Systems

LLMs



English: Hello, how are you?
French: Bonjour, comment ça va?
Spanish: Hola, ¿cómo estás?
German: Hallo, wie geht es dir?
Italian: Ciao, come stai?

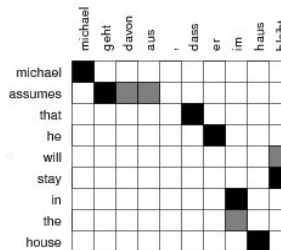
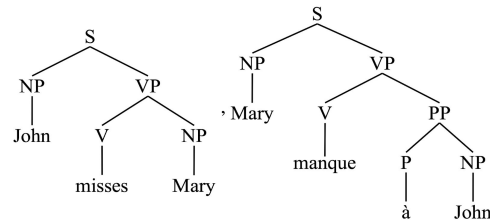
✓ Generality/
Flexibility

✗

✗ Interpretability/
Controllability

✓

Symbolic Systems



LLMs + Symbolic Systems?

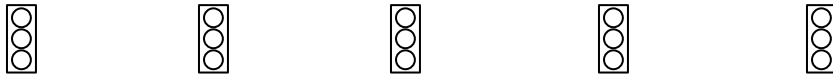
How can we combine the controllability/interpretability of symbolic models with the flexibility of LLMs?

How do Transformers work?

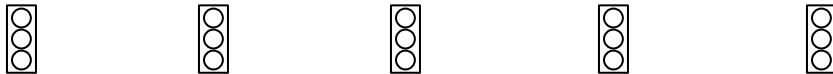
L : sequence length

d : hidden state dimension

Output vectors $\mathbf{O} \in \mathbb{R}^{L \times d}$



Input vectors $\mathbf{X} \in \mathbb{R}^{L \times d}$



How do Transformers work?



\mathbf{x}_t

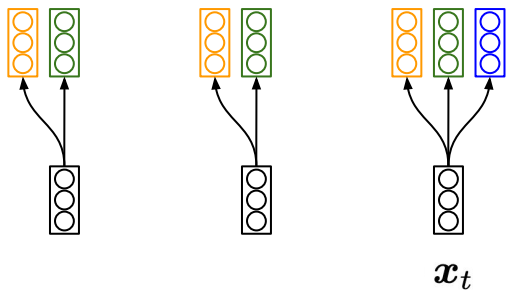
How do Transformers work?

$$\mathbf{q}_t, \mathbf{k}_t, \mathbf{v}_t = \mathbf{x}_t \mathbf{W}_Q, \mathbf{x}_t \mathbf{W}_K, \mathbf{x}_t \mathbf{W}_V$$

Key K

Value V

Query Q



How do Transformers work?

“Attention distribution”: Model pairwise interactions between current and previous inputs

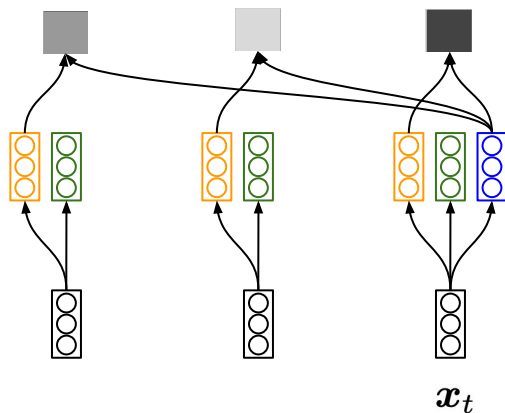
$$\frac{\exp(\mathbf{q}_t^\top \mathbf{k}_j)}{\sum_{l=1}^t \exp(\mathbf{q}_t^\top \mathbf{k}_l)}$$

$$\mathbf{q}_t, \mathbf{k}_t, \mathbf{v}_t = \mathbf{x}_t \mathbf{W}_Q, \mathbf{x}_t \mathbf{W}_K, \mathbf{x}_t \mathbf{W}_V$$

Key K

Value V

Query Q



How do Transformers work?

This vector is a function of *all* previous inputs!

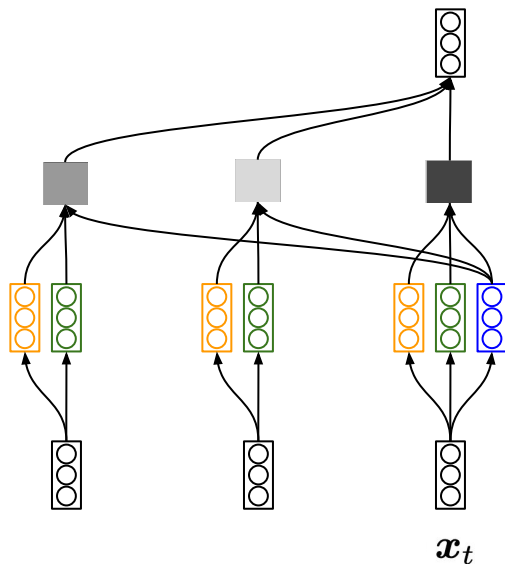
$$\mathbf{o}_t = \sum_{j=1}^t \frac{\exp(\mathbf{q}_t^\top \mathbf{k}_j)}{\sum_{l=1}^t \exp(\mathbf{q}_t^\top \mathbf{k}_l)} \mathbf{v}_j$$

$$\mathbf{q}_t, \mathbf{k}_t, \mathbf{v}_t = \mathbf{x}_t \mathbf{W}_Q, \mathbf{x}_t \mathbf{W}_K, \mathbf{x}_t \mathbf{W}_V$$

Key K

Value V

Query Q



How do Transformers work?

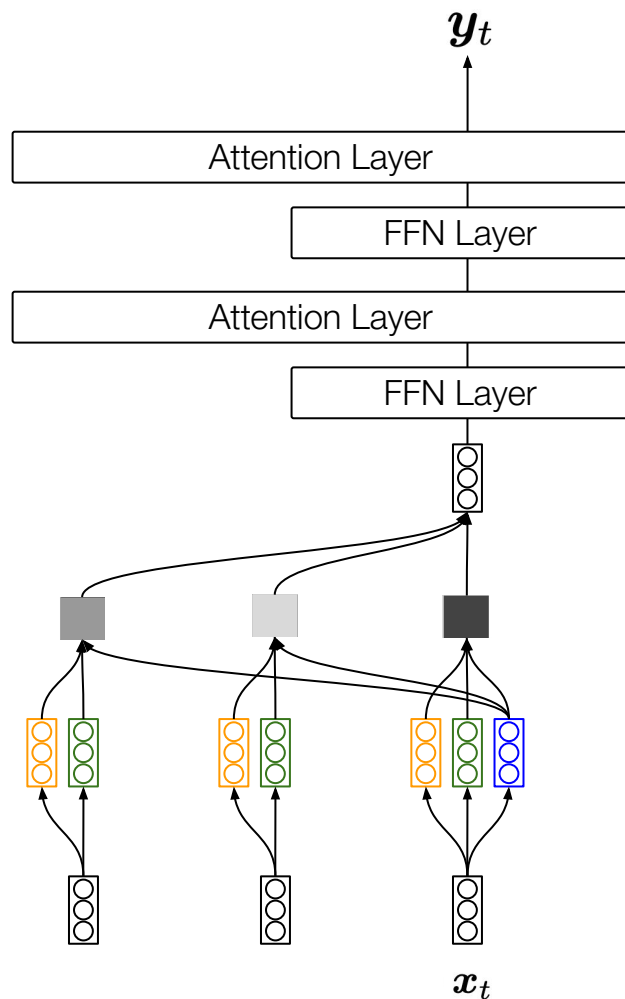
$$\mathbf{o}_t = \sum_{j=1}^t \frac{\exp(\mathbf{q}_t^\top \mathbf{k}_j)}{\sum_{l=1}^t \exp(\mathbf{q}_t^\top \mathbf{k}_l)} \mathbf{v}_j$$

$$\mathbf{q}_t, \mathbf{k}_t, \mathbf{v}_t = \mathbf{x}_t \mathbf{W}_Q, \mathbf{x}_t \mathbf{W}_K, \mathbf{x}_t \mathbf{W}_V$$

Key K

Value V

Query Q

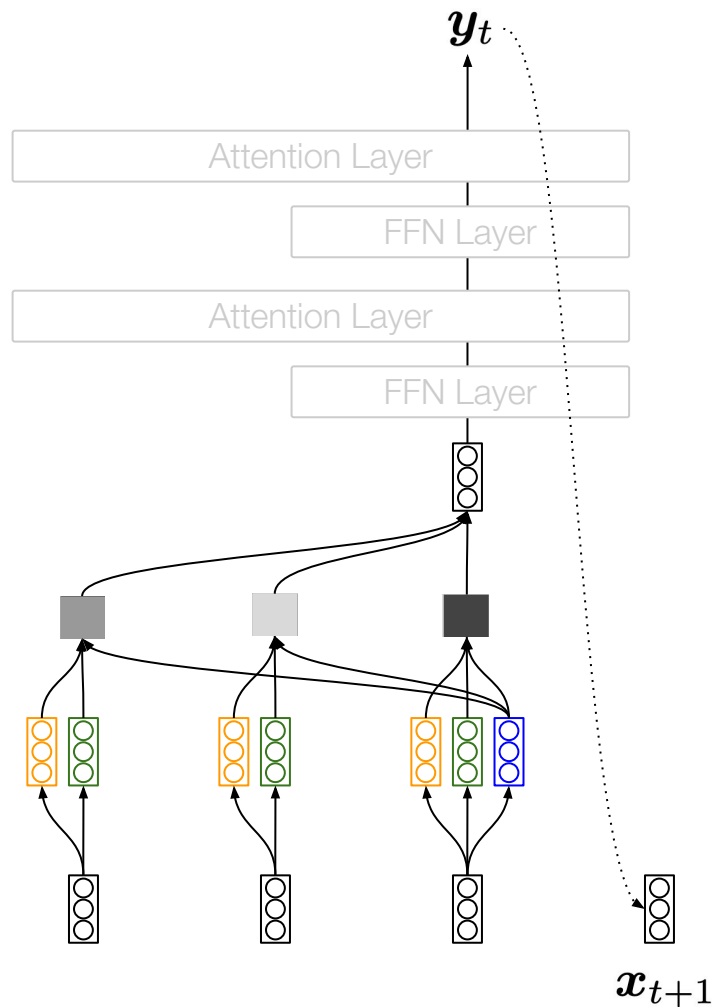


How do Transformers work?

$$\mathbf{o}_t = \sum_{j=1}^t \frac{\exp(\mathbf{q}_t^\top \mathbf{k}_j)}{\sum_{l=1}^t \exp(\mathbf{q}_t^\top \mathbf{k}_l)} \mathbf{v}_j$$

$$\mathbf{q}_t, \mathbf{k}_t, \mathbf{v}_t = \mathbf{x}_t \mathbf{W}_Q, \mathbf{x}_t \mathbf{W}_K, \mathbf{x}_t \mathbf{W}_V$$

Key K
Value V
Query Q



How do Transformers work?

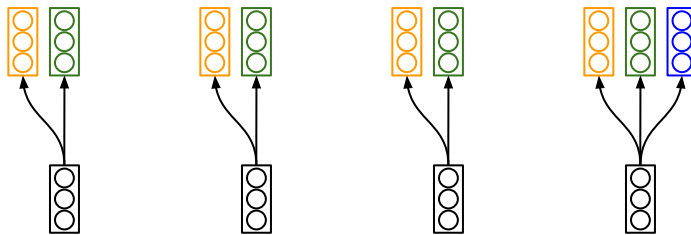
$$\mathbf{o}_t = \sum_{j=1}^t \frac{\exp(\mathbf{q}_t^\top \mathbf{k}_j)}{\sum_{l=1}^t \exp(\mathbf{q}_t^\top \mathbf{k}_l)} \mathbf{v}_j$$

$$\mathbf{q}_t, \mathbf{k}_t, \mathbf{v}_t = \mathbf{x}_t \mathbf{W}_Q, \mathbf{x}_t \mathbf{W}_K, \mathbf{x}_t \mathbf{W}_V$$

Key K

Value V

Query Q

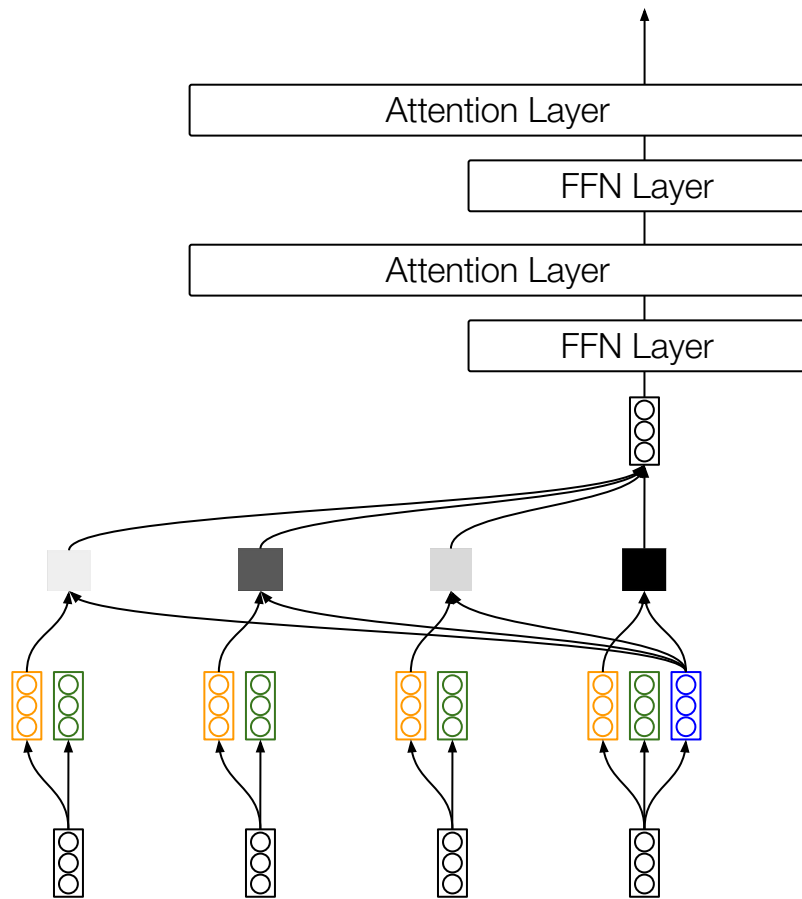


How do Transformers work?

$$\mathbf{o}_t = \sum_{j=1}^t \frac{\exp(\mathbf{q}_t^\top \mathbf{k}_j)}{\sum_{l=1}^t \exp(\mathbf{q}_t^\top \mathbf{k}_l)} \mathbf{v}_j$$

$$\mathbf{q}_t, \mathbf{k}_t, \mathbf{v}_t = \mathbf{x}_t \mathbf{W}_Q, \mathbf{x}_t \mathbf{W}_K, \mathbf{x}_t \mathbf{W}_V$$

Key K
Value V
Query Q



How do Transformers work?

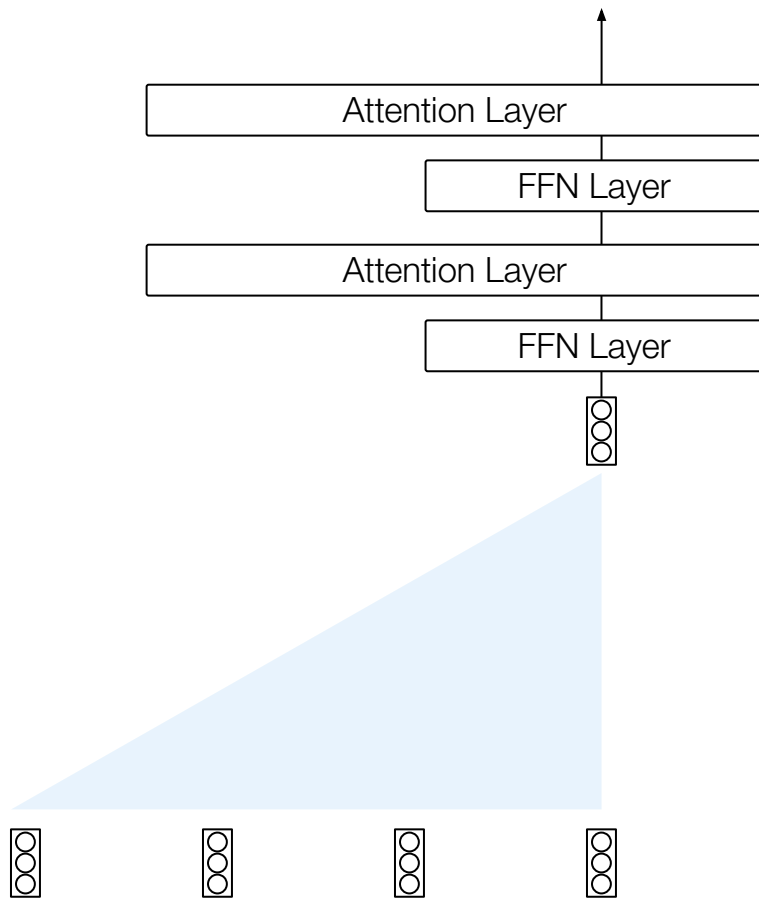
$$\mathbf{o}_t = \sum_{j=1}^t \frac{\exp(\mathbf{q}_t^\top \mathbf{k}_j)}{\sum_{l=1}^t \exp(\mathbf{q}_t^\top \mathbf{k}_l)} \mathbf{v}_j$$

$$\mathbf{q}_t, \mathbf{k}_t, \mathbf{v}_t = \mathbf{x}_t \mathbf{W}_Q, \mathbf{x}_t \mathbf{W}_K, \mathbf{x}_t \mathbf{W}_V$$

Key K

Value V

Query Q

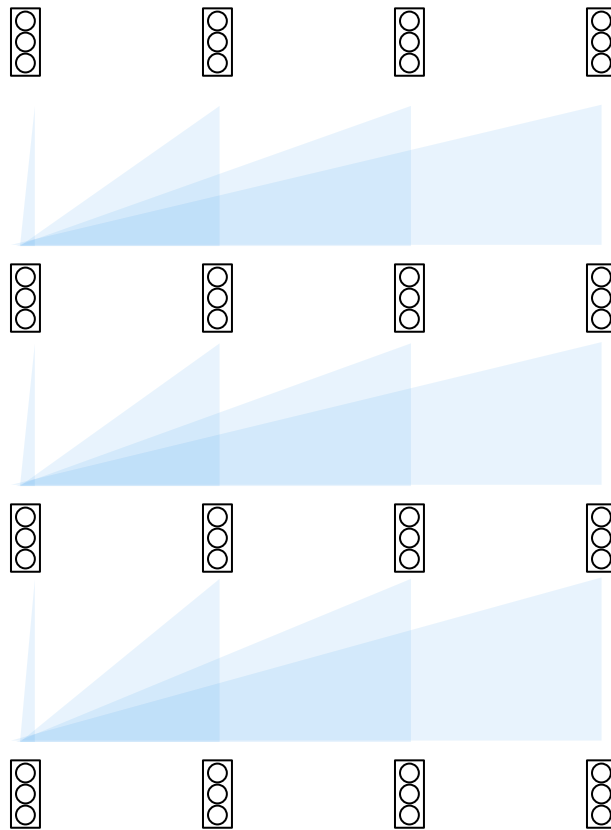


How do Transformers work?

Transformers use flexible but uninterpretable internal operations to compute the output given the input.

$$\mathbf{q}_t, \mathbf{k}_t, \mathbf{v}_t = \mathbf{x}_t \mathbf{W}_Q, \mathbf{x}_t \mathbf{W}_K, \mathbf{x}_t \mathbf{W}_V$$

Can we implement “hard” rules in Transformers?



A Logic for Expressing Log-Precision Transformers

William Merrill
New York University
willm@nyu.edu

Ashish Sabharwal
Allen Institute for AI
ashishs@allenai.org

Any log-precision transformer can be re-expressed as a sentence in FO(M) logic, e.g.:

$$M_i. a(i) \wedge M_j. b(j) \wedge \neg \exists k, \ell. (a(k) \wedge b(\ell) \wedge \ell < k)$$

(m a's followed by m b's, i.e., $a^m b^m$)

aaaabbbb ✓

aaabbbbb ✗

baaaabbbb ✗

Implementing Logic in Transformers

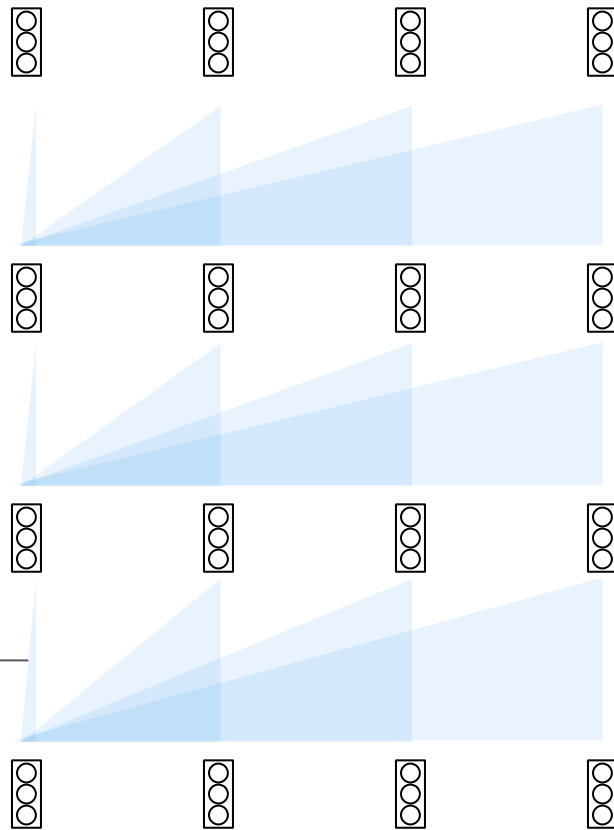
Any log-precision transformer can be re-expressed as a sentence in FO(M) logic, e.g.:

$$\exists i. a(i) \wedge \exists j. b(j) \wedge \neg \exists k, \ell. (a(k) \wedge b(\ell) \wedge \ell < k)$$

(*m* a's followed by *m* b's, i.e., $a^m b^m$)

aaaabbbb ✓ aaabbbb ✗ baaaabbbb ✗

Trained Transformer weights can
be “read out” in the form of logical
rules → more interpretable!



Implementing Logic in Transformers

These logical rules can be
“compiled” into the weights of
Transformers → more controllable!

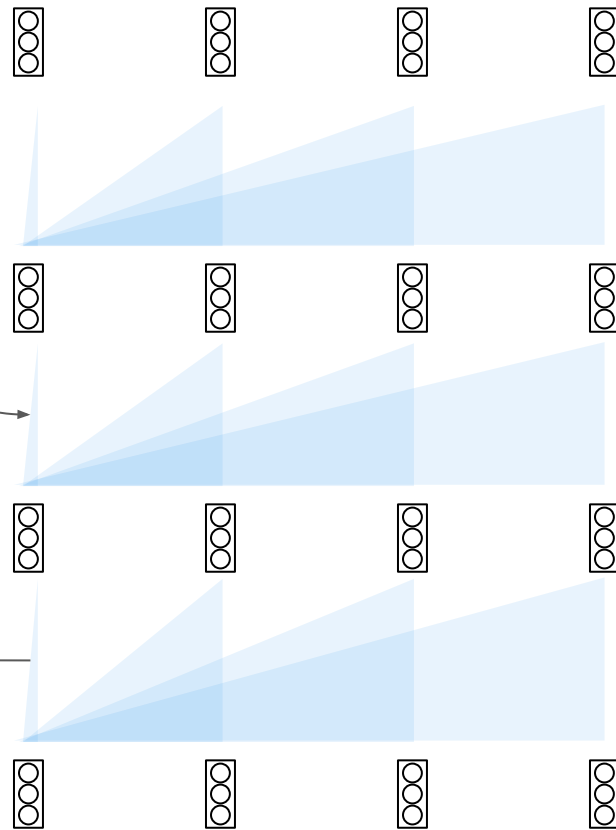
Any log-precision transformer can be re-expressed as a sentence in FO(M) logic, e.g.:

$$\mathbf{M}i. \mathbf{a}(i) \wedge \mathbf{M}j. \mathbf{b}(j) \wedge \neg \exists k, \ell. (\mathbf{a}(k) \wedge \mathbf{b}(\ell) \wedge \ell < k)$$

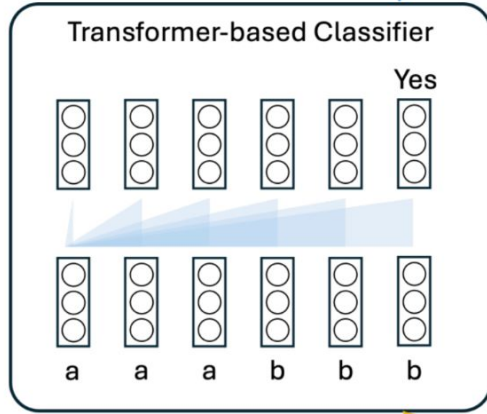
(*m* a's followed by *m* b's, i.e., $a^m b^m$)

aaaabbbb ✓ aaabbbb ✗ baaaabbbb ✗

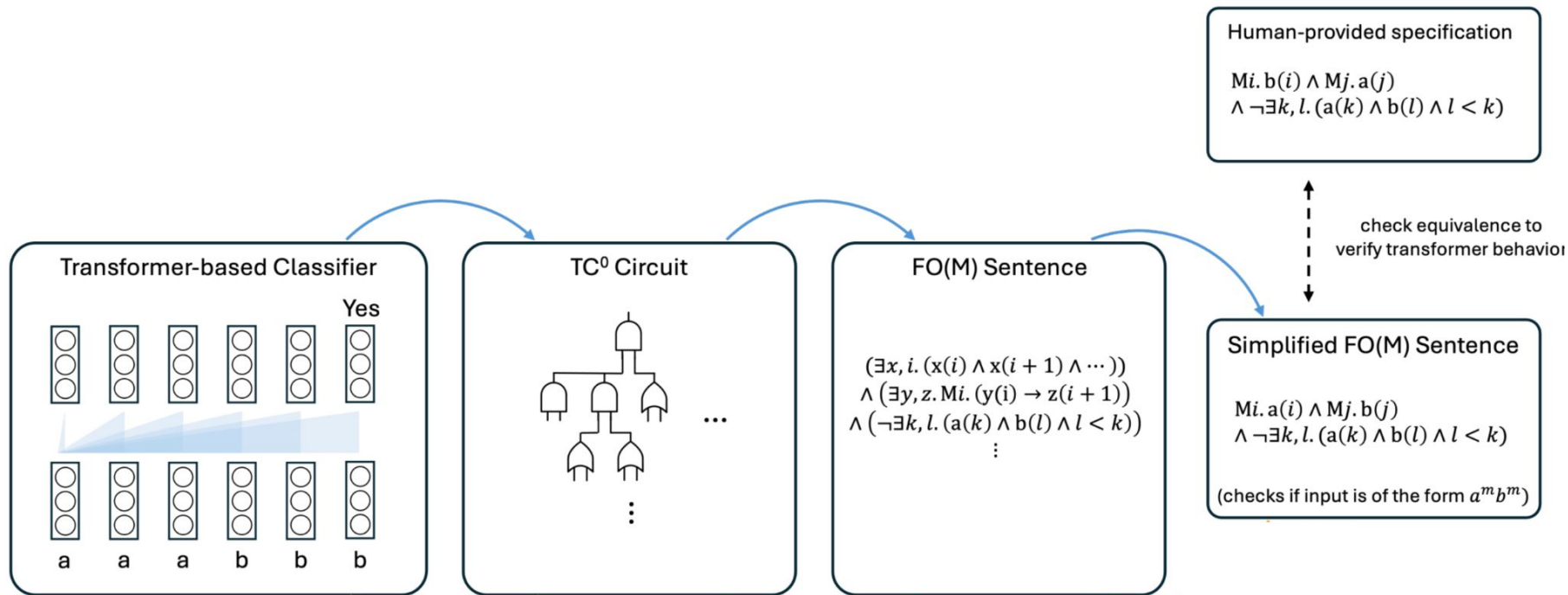
Trained Transformer weights can
be “read out” in the form of logical
rules → more interpretable!



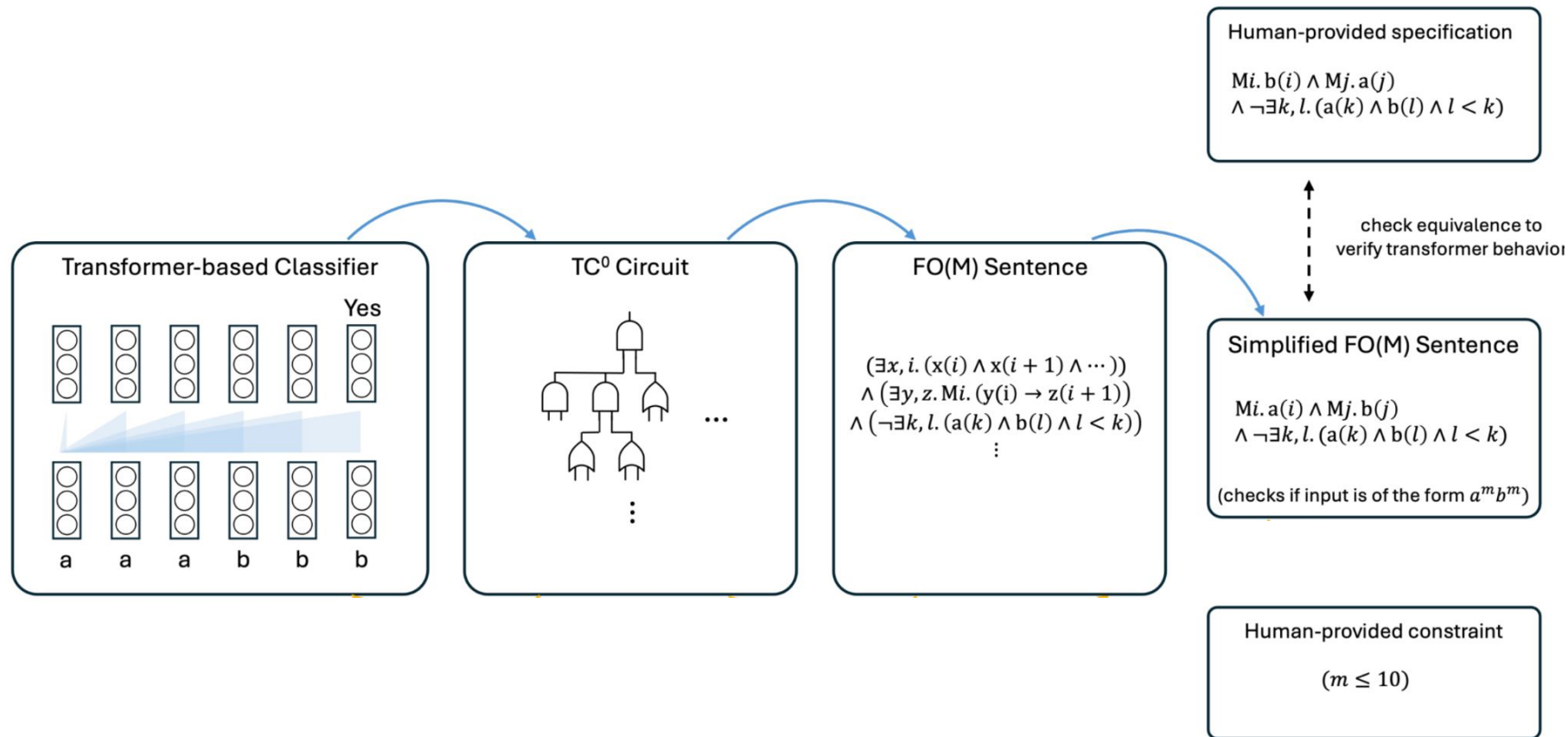
Our project



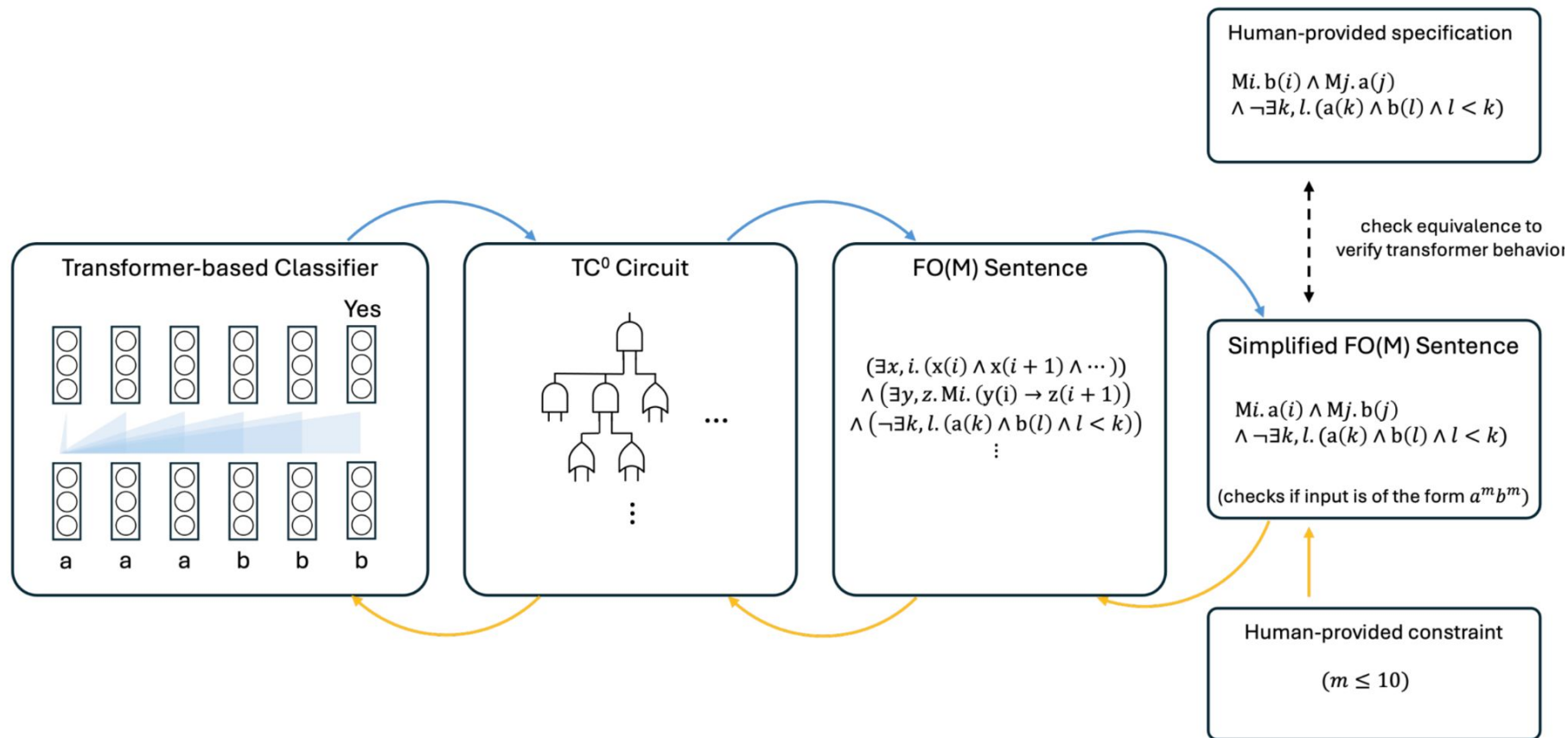
Our project



Our project



Our project



Ongoing Work

