

Computability and Transformers: A Circuit Complexity Revival?

Michaël Cadilhac



Joint work with Andy Yang and David Chiang (Notre Dame)
Appearing at NeurIPS 2025

→ THE T of GPT

Computability and Transformers: A Circuit Complexity Revival?

Michaël Cadilhac



Joint work with Andy Yang and David Chiang (Notre Dame)
Appearing at NeurIPS 2025

Goals of this presentation

- ▶ Succinctly present transformers (as recognizers)

Goals of this presentation

- ▶ Succinctly present transformers (as recognizers)
- ▶ Tie them to two notions of computability:

Goals of this presentation

- ▶ Succinctly present transformers (as recognizers)
- ▶ Tie them to two notions of computability:
 - ▶ Transformers are Turing-complete *when allowed to “think”* for an unbounded number of steps

Goals of this presentation

- ▶ Succinctly present transformers (as recognizers)
- ▶ Tie them to two notions of computability:
 - ▶ Transformers are Turing-complete *when allowed to “think”* for an unbounded number of steps
 - ▶ Transformers are in small circuit-complexity classes as classifiers

Goals of this presentation

- ▶ Succinctly present transformers (as recognizers)
- ▶ Tie them to two notions of computability:
 - ▶ Transformers are Turing-complete *when allowed to “think”* for an unbounded number of steps
 - ▶ Transformers are in small circuit-complexity classes as classifiers
- ▶ In the latter case, tie transformers to temporal logics...

Goals of this presentation

- ▶ Succinctly present transformers (as recognizers)
- ▶ Tie them to two notions of computability:
 - ▶ Transformers are Turing-complete *when allowed to “think”* for an unbounded number of steps
 - ▶ Transformers are in small circuit-complexity classes as classifiers
- ▶ In the latter case, tie transformers to temporal logics...
- ▶ ...and a derive depth-hierarchy result for them: the deeper the transformer, the more it can express

Goals of this presentation

- ▶ Succinctly present transformers (as recognizers)
- ▶ Tie them to two notions of computability:
 - ▶ Transformers are Turing-complete *when allowed to “think”* for an unbounded number of steps
 - ▶ Transformers are in small circuit-complexity classes as classifiers
- ▶ In the latter case, tie transformers to temporal logics...
- ▶ ...and a derive depth-hierarchy result for them: the deeper the transformer, the more it can express



OUR CONTRIBUTIONS.

Transformers

What can transformers compute, if we give them time to think?

What can transformers compute in one evaluation?

Conclusion

Transformers

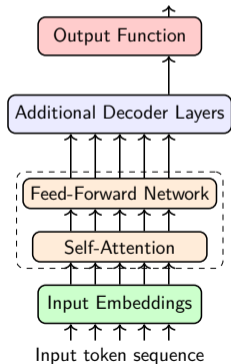
- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations

Transformers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)

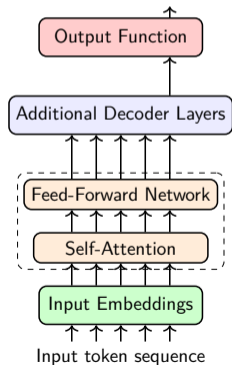
Transformers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)



Transformers

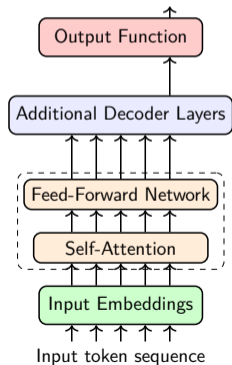
- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)



$$E: \Sigma \rightarrow \mathbb{R}^d$$

Transformers

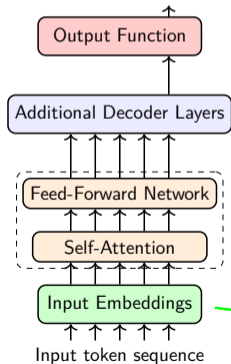
- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)



$$\text{lay}^{(0)}(i) = E(w_i) \in \mathbb{R}^d$$
$$E: \Sigma \rightarrow \mathbb{R}^d$$

Transformers

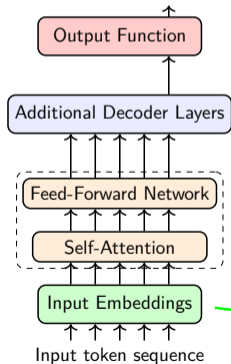
- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)



$$\text{lay}^{(0)}(i) = E(w_i) \in \mathbb{R}^d$$
$$E: \Sigma \rightarrow \mathbb{R}^d$$

Transformers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)



$$k^{(1)}(i) = K^{(1)}\text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}}$$

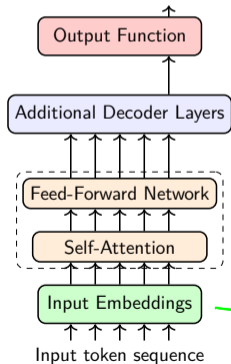
$$q^{(1)}(i) = Q^{(1)}\text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}}$$

$$\text{lay}^{(0)}(i) = E(w_i) \in \mathbb{R}^d$$

$$E: \Sigma \rightarrow \mathbb{R}^d$$

Transformers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)



$$k^{(1)}(i) = K^{(1)} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}}$$

$$q^{(1)}(i) = Q^{(1)} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}}$$

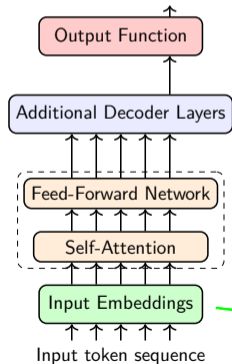
$$\text{lay}^{(0)}(i) = E(w_i) \in \mathbb{R}^d$$

$$E: \Sigma \rightarrow \mathbb{R}^d$$

$$\in \mathbb{R}^{d \times d^{\text{hid}}}$$

Transformers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)



$$\text{att}^{(1)}(i, j) = q^{(1)}(i) \cdot k^{(1)}(j) \in \mathbb{R}$$

$$k^{(1)}(i) = \underline{K^{(1)}} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}}$$

$$q^{(1)}(i) = \underline{Q^{(1)}} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}}$$

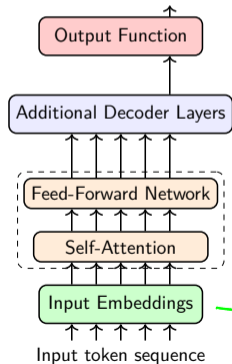
$$\text{lay}^{(0)}(i) = E(w_i) \in \mathbb{R}^d$$

$$E: \Sigma \rightarrow \mathbb{R}^d$$

$$\in \mathbb{R}^{d \times d^{\text{hid}}}$$

Transformers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)



$$v^{(1)}(i) = V^{(1)} \cdot \text{lay}^{(0)}(i) \in \mathbb{R}^d$$

$$\text{att}^{(1)}(i, j) = q^{(1)}(i) \cdot k^{(1)}(j) \in \mathbb{R}$$

$$k^{(1)}(i) = K^{(1)} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}}$$

$$q^{(1)}(i) = Q^{(1)} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}}$$

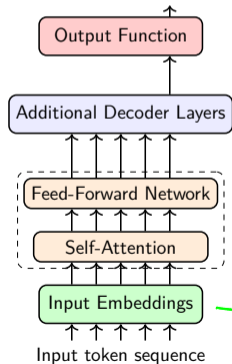
$$\text{lay}^{(0)}(i) = E(w_i) \in \mathbb{R}^d$$

$$E: \Sigma \rightarrow \mathbb{R}^d$$

$$\in \mathbb{R}^{d \times d^{\text{hid}}}$$

Transformers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)



$$v^{(1)}(i) = \underline{V^{(1)}} \cdot \text{lay}^{(0)}(i) \in \mathbb{R}^d$$

$$\text{att}^{(1)}(i, j) = q^{(1)}(i) \cdot k^{(1)}(j) \in \mathbb{R}$$

$$k^{(1)}(i) = \underline{K^{(1)}} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}}$$

$$q^{(1)}(i) = \underline{Q^{(1)}} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}}$$

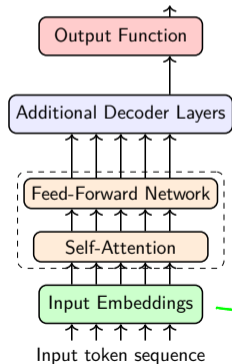
$$\text{lay}^{(0)}(i) = E(w_i) \in \mathbb{R}^d$$

$$E: \Sigma \rightarrow \mathbb{R}^d$$

$$\in \mathbb{R}^{d \times d^{\text{hid}}}$$

Transformers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)



$$\text{pre-lay}^{(1)}(i) = \frac{\sum_{j \leq i} \exp(\text{att}^{(1)}(i,j)) v^{(1)}(j)}{\sum_{j \leq i} \exp(\text{att}^{(1)}(i,j))} \in \mathbb{R}^d$$

$$v^{(1)}(i) = \underline{V}^{(1)} \cdot \text{lay}^{(0)}(i) \in \mathbb{R}^d$$

$$\text{att}^{(1)}(i,j) = q^{(1)}(i) \cdot k^{(1)}(j) \in \mathbb{R}$$

$$k^{(1)}(i) = \underline{K}^{(1)} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}}$$

$$q^{(1)}(i) = \underline{Q}^{(1)} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}}$$

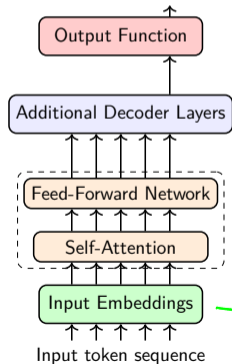
$$\text{lay}^{(0)}(i) = E(w_i) \in \mathbb{R}^d$$

$$E: \Sigma \rightarrow \mathbb{R}^d$$

$\in \mathbb{R}^{d \times d^{\text{hid}}}$

Transformers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)



$$\text{pre-lay}^{(1)}(i) = \frac{\sum_{j \leq i} \exp(\text{att}^{(1)}(i,j)) v^{(1)}(j)}{\sum_{j \leq i} \exp(\text{att}^{(1)}(i,j))} \in \mathbb{R}^d$$

$$v^{(1)}(i) = \underline{V}^{(1)} \cdot \text{lay}^{(0)}(i) \in \mathbb{R}^d$$

$$\text{att}^{(1)}(i,j) = q^{(1)}(i) \cdot k^{(1)}(j) \in \mathbb{R}$$

$$k^{(1)}(i) = \underline{K}^{(1)} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}}$$

$$q^{(1)}(i) = \underline{Q}^{(1)} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}}$$

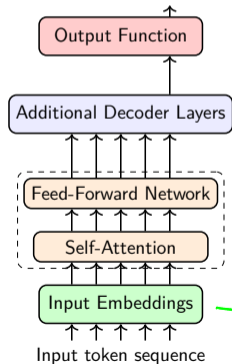
$$\text{lay}^{(0)}(i) = \underline{E}(w_i) \in \mathbb{R}^d$$

$$E: \Sigma \rightarrow \mathbb{R}^d$$

$$\in \mathbb{R}^{d \times d^{\text{hid}}}$$

Transformers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)

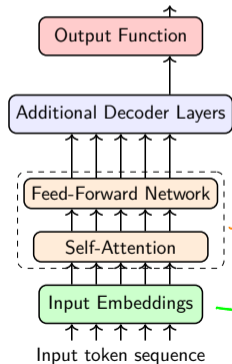


$$\begin{aligned} \text{lay}^{(1)}(i) &= \text{ffnn}^{(1)}(\text{pre-lay}^{(1)}(i)) \in \mathbb{R}^d \\ \text{pre-lay}^{(1)}(i) &= \frac{\sum_{j \leq i} \exp(\text{att}^{(1)}(i,j)) v^{(1)}(j)}{\sum_{j \leq i} \exp(\text{att}^{(1)}(i,j))} \in \mathbb{R}^d \\ v^{(1)}(i) &= \underline{V}^{(1)} \cdot \text{lay}^{(0)}(i) \in \mathbb{R}^d \\ \text{att}^{(1)}(i,j) &= q^{(1)}(i) \cdot k^{(1)}(j) \in \mathbb{R} \\ k^{(1)}(i) &= \underline{K}^{(1)} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}} \\ q^{(1)}(i) &= \underline{Q}^{(1)} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}} \\ \text{lay}^{(0)}(i) &= E(w_i) \in \mathbb{R}^d \\ E: \Sigma &\rightarrow \mathbb{R}^d \end{aligned}$$

$\in \mathbb{R}^{d \times d^{\text{hid}}}$

Transformers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)

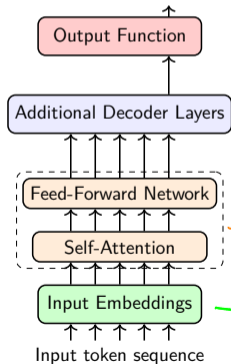


$$\begin{aligned} \text{lay}^{(1)}(i) &= \text{ffnn}^{(1)}(\text{pre-lay}^{(1)}(i)) \in \mathbb{R}^d \\ \text{pre-lay}^{(1)}(i) &= \frac{\sum_{j \leq i} \exp(\text{att}^{(1)}(i,j)) v^{(1)}(j)}{\sum_{j \leq i} \exp(\text{att}^{(1)}(i,j))} \in \mathbb{R}^d \\ v^{(1)}(i) &= \underline{V}^{(1)} \cdot \text{lay}^{(0)}(i) \in \mathbb{R}^d \\ \text{att}^{(1)}(i,j) &= q^{(1)}(i) \cdot k^{(1)}(j) \in \mathbb{R} \\ k^{(1)}(i) &= \underline{K}^{(1)} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}} \\ q^{(1)}(i) &= \underline{Q}^{(1)} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}} \\ \text{lay}^{(0)}(i) &= E(w_i) \in \mathbb{R}^d \\ E: \Sigma &\rightarrow \mathbb{R}^d \end{aligned}$$

$\in \mathbb{R}^{d \times d^{\text{hid}}}$

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)

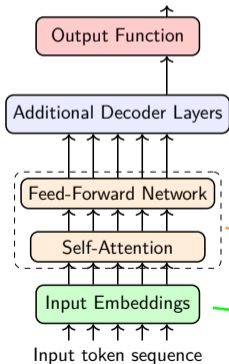


$$\begin{aligned} \text{lay}^{(1)}(i) &= \text{ffnn}^{(1)}(\text{pre-lay}^{(1)}(i)) \in \mathbb{R}^d \\ \text{pre-lay}^{(1)}(i) &= \frac{\sum_{j \leq i} \exp(\text{att}^{(1)}(i,j)) v^{(1)}(j)}{\sum_{j \leq i} \exp(\text{att}^{(1)}(i,j))} \in \mathbb{R}^d \\ v^{(1)}(i) &= \underline{V}^{(1)} \cdot \text{lay}^{(0)}(i) \in \mathbb{R}^d \\ \text{att}^{(1)}(i,j) &= q^{(1)}(i) \cdot k^{(1)}(j) \in \mathbb{R} \\ k^{(1)}(i) &= \underline{K}^{(1)} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}} \\ q^{(1)}(i) &= \underline{Q}^{(1)} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}} \\ \text{lay}^{(0)}(i) &= E(w_i) \in \mathbb{R}^d \\ E: \Sigma &\rightarrow \mathbb{R}^d \end{aligned}$$

$\in \mathbb{R}^{d \times d^{\text{hid}}}$

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)



... other layers ...

$$\text{lay}^{(1)}(i) = \text{ffnn}^{(1)}(\text{pre-lay}^{(1)}(i)) \in \mathbb{R}^d$$

$$\text{pre-lay}^{(1)}(i) = \frac{\sum_{j \leq i} \exp(\text{att}^{(1)}(i,j)) v^{(1)}(j)}{\sum_{j \leq i} \exp(\text{att}^{(1)}(i,j))} \in \mathbb{R}^d$$

$$v^{(1)}(i) = V^{(1)} \cdot \text{lay}^{(0)}(i) \in \mathbb{R}^d$$

$$\text{att}^{(1)}(i,j) = q^{(1)}(i) \cdot k^{(1)}(j) \in \mathbb{R}$$

$$k^{(1)}(i) = K^{(1)} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}}$$

$$q^{(1)}(i) = Q^{(1)} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}}$$

$$\text{lay}^{(0)}(i) = E(w_i) \in \mathbb{R}^d$$

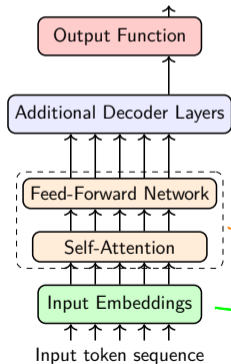
$$E: \Sigma \rightarrow \mathbb{R}^d$$

$$\in \mathbb{R}^{d \times d^{\text{hid}}}$$

Transformers

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)



... other layers ... → "DEPTH"

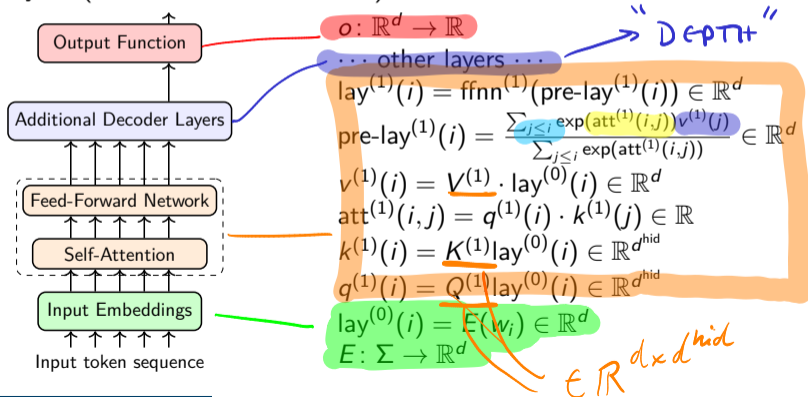
$$\text{lay}^{(1)}(i) = \text{ffnn}^{(1)}(\text{pre-lay}^{(1)}(i)) \in \mathbb{R}^d$$
$$\text{pre-lay}^{(1)}(i) = \frac{\sum_{j \leq i} \exp(\text{att}^{(1)}(i,j)) v^{(1)}(j)}{\sum_{j \leq i} \exp(\text{att}^{(1)}(i,j))} \in \mathbb{R}^d$$
$$v^{(1)}(i) = V^{(1)} \cdot \text{lay}^{(0)}(i) \in \mathbb{R}^d$$
$$\text{att}^{(1)}(i,j) = q^{(1)}(i) \cdot k^{(1)}(j) \in \mathbb{R}$$
$$k^{(1)}(i) = K^{(1)} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}}$$
$$q^{(1)}(i) = Q^{(1)} \text{lay}^{(0)}(i) \in \mathbb{R}^{d^{\text{hid}}}$$
$$\text{lay}^{(0)}(i) = E(w_i) \in \mathbb{R}^d$$
$$E: \Sigma \rightarrow \mathbb{R}^d$$

$E \in \mathbb{R}^{d \times d^{\text{hid}}}$

Transformers

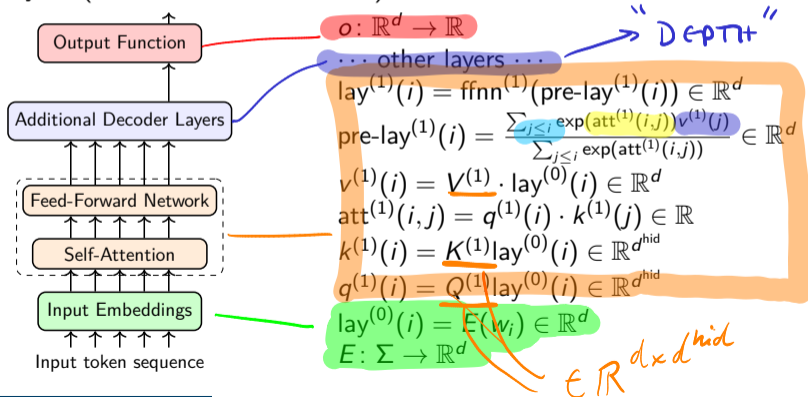
For each pos i , compute score against j , weigh by a value, then take avg (softmax)

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)



Transformers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)



Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformers as Recognizers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformers as Recognizers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)
- ▶ Language-theoretic objects:
 - ▶ *Chain of thought language*: let T “generate” tokens, until it generates “accept/reject”

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformers as Recognizers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)
- ▶ Language-theoretic objects:
 - ▶ *Chain of thought language*: let T “generate” tokens, until it generates “accept/reject”



Attention layer

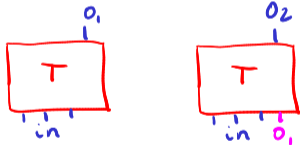
For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformers as Recognizers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)
- ▶ Language-theoretic objects:
 - ▶ *Chain of thought language*: let T “generate” tokens, until it generates “accept/reject”



Attention layer

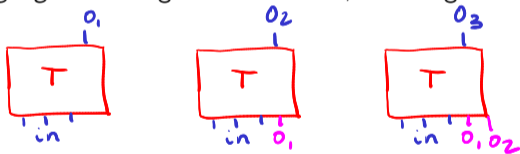
For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformers as Recognizers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)
- ▶ Language-theoretic objects:
 - ▶ *Chain of thought language*: let T “generate” tokens, until it generates “accept/reject”



Attention layer

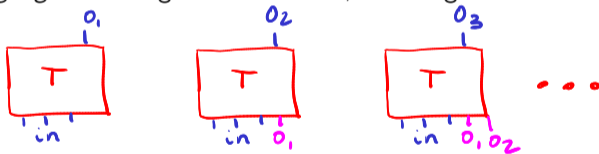
For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformers as Recognizers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)
- ▶ Language-theoretic objects:
 - ▶ *Chain of thought language*: let T “generate” tokens, until it generates “accept/reject”



Attention layer

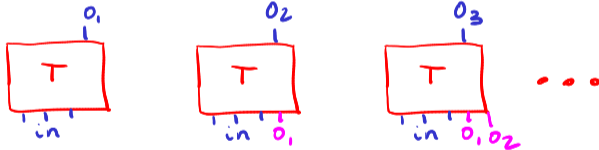
For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformers as Recognizers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)
- ▶ Language-theoretic objects:
 - ▶ *Chain of thought language*: let T “generate” tokens, until it generates “accept/reject”



- ▶ *Transformer language*: $\{w \in \Sigma^* \mid T(w) > 0\}$

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

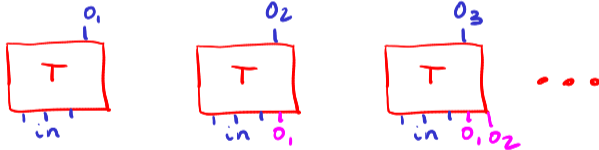
Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformers as Recognizers

- ▶ Transformers compute functions $T: \Sigma^* \rightarrow \mathbb{R}$ using a bunch of (mostly) linear transformations
- ▶ Transformers have multiple *layers*, each of them accessing the previous layers (“attention” mechanism)
- ▶ Language-theoretic objects:

- ▶ *Chain of thought language*: let T “generate” tokens, until it generates “accept/reject”



- ▶ *Transformer language*: $\{w \in \Sigma^* \mid T(w) > 0\}$
- ▶ *What is the complexity of these languages?*

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformers

What can transformers compute, if we give them time to think?

What can transformers compute in one evaluation?

Conclusion

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Chain of thought languages

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Chain of thought languages

Theorem (Pérez et al. 2019)

“Transformers are Turing-complete”

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Chain of thought languages

Theorem (Pérez et al. 2019)

“Transformers are Turing-complete”

Idea.

Let M be a TM; a transformer with chain-of-thought can simulate it:

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Chain of thought languages

Theorem (Pérez et al. 2019)

“Transformers are Turing-complete”

Idea.

Let M be a TM; a transformer with chain-of-thought can simulate it:

- ▶ Based on first character of input, output first transition of M

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Theorem (Pérez et al. 2019)

“Transformers are Turing-complete”

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Idea.

Let M be a TM; a transformer with chain-of-thought can simulate it:

- ▶ Based on first character of input, output first transition of M
- ▶ After that, at each generation step, count number of left/right moves using attention

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Theorem (Pérez et al. 2019)

“Transformers are Turing-complete”

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Idea.

Let M be a TM; a transformer with chain-of-thought can simulate it:

- ▶ Based on first character of input, output first transition of M
- ▶ After that, at each generation step, count number of left/right moves using attention
- ▶ Track what is cell content at that position

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Theorem (Pérez et al. 2019)

“Transformers are Turing-complete”

Idea.

Let M be a TM; a transformer with chain-of-thought can simulate it:

- ▶ Based on first character of input, output first transition of M
- ▶ After that, at each generation step, count number of left/right moves using attention
- ▶ Track what is cell content at that position
- ▶ Generate the next transition of the Turing machine □

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Theorem (Pérez et al. 2019)

“Transformers are Turing-complete”

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Idea.

Let M be a TM; a transformer with chain-of-thought can simulate it:

- ▶ Based on first character of input, output first transition of M
- ▶ After that, at each generation step, count number of left/right moves using attention
- ▶ Track what is cell content at that position
- ▶ Generate the next transition of the Turing machine □

- ▶ Crucially relies on chain-of-thought!

Transformers

What can transformers compute, if we give them time to think?

What can transformers compute in one evaluation?

Conclusion

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformer languages

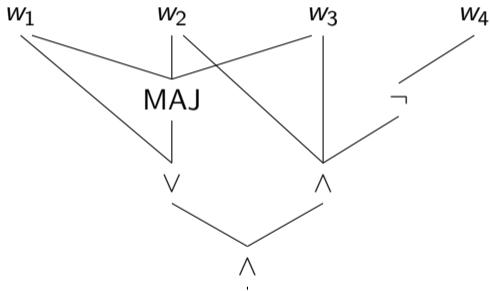
Preliminary 1: Circuit complexity

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

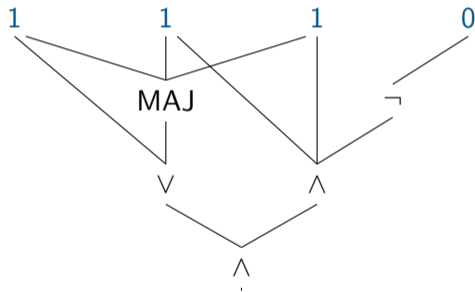
Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate



Transformer languages

Preliminary 1: Circuit complexity



Attention layer

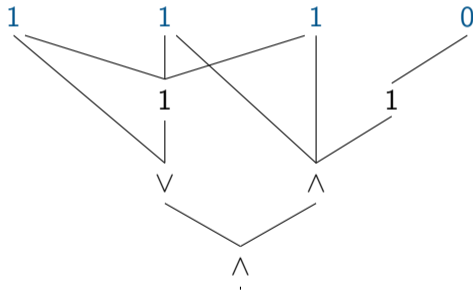
For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformer languages

Preliminary 1: Circuit complexity



Attention layer

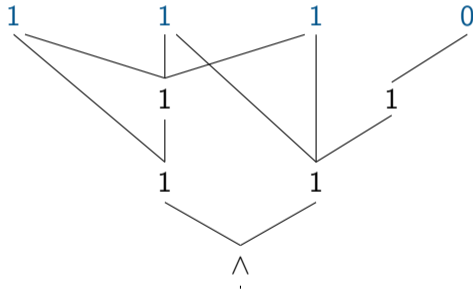
For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformer languages

Preliminary 1: Circuit complexity



Attention layer

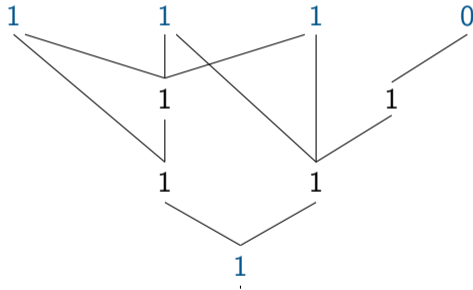
For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformer languages

Preliminary 1: Circuit complexity



Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformer languages

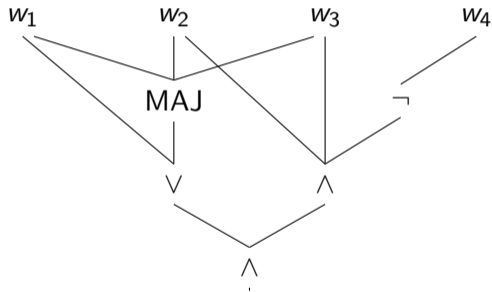
Preliminary 1: Circuit complexity

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate



$$L_4 = \{1110, \dots\}$$

Transformer languages

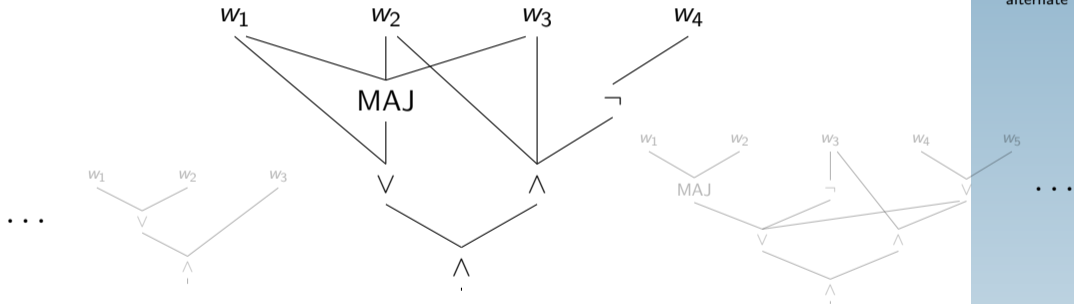
Preliminary 1: Circuit complexity

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate



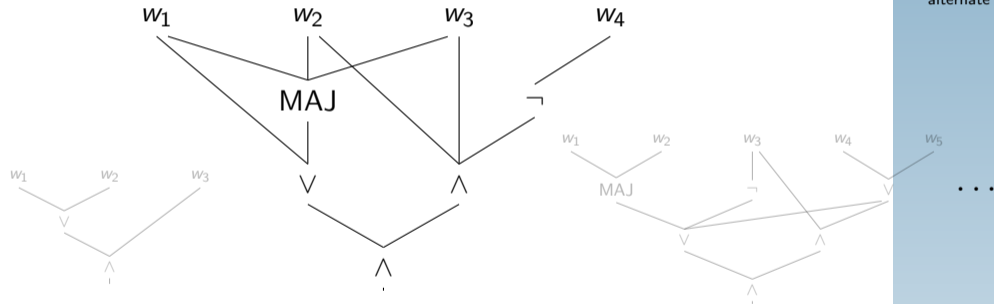
$$L_4 = \{1110, \dots\}$$

$$L = \bigcup_{i \in \mathbb{N}} L_i$$

Transformer languages

Preliminary 1: Circuit complexity

$\text{size}(n) = \text{number of gates}$



$$L_4 = \{1110, \dots\}$$

$$L = \bigcup_{i \in \mathbb{N}} L_i$$

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

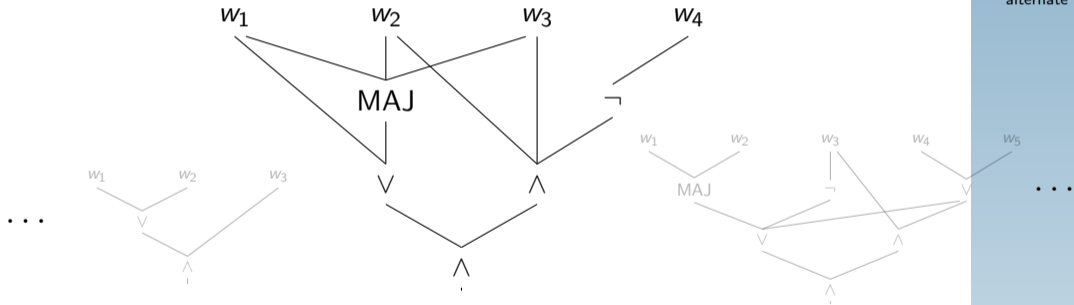
Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformer languages

Preliminary 1: Circuit complexity

$\text{size}(n)$ = number of gates
= **poly**(n)



$$L_4 = \{1110, \dots\}$$

$$L = \bigcup_{i \in \mathbb{N}} L_i$$

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

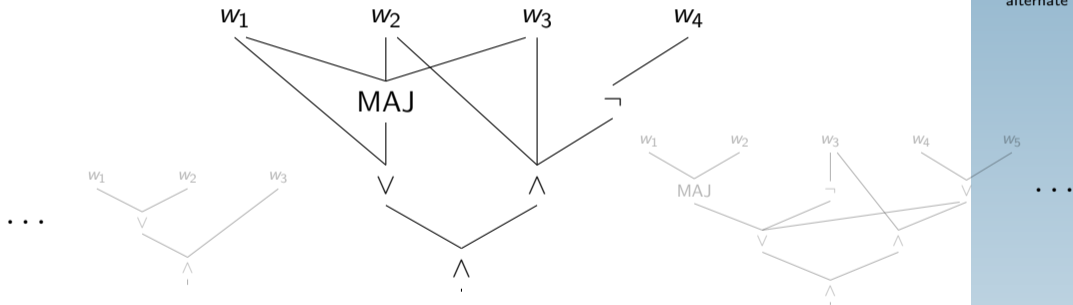
- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformer languages

Preliminary 1: Circuit complexity

$\text{size}(n)$ = number of gates
= **poly**(n)

$\text{depth}(n)$ = maximal length of a path



$$L_4 = \{1110, \dots\}$$

$$L = \bigcup_{i \in \mathbb{N}} L_i$$

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

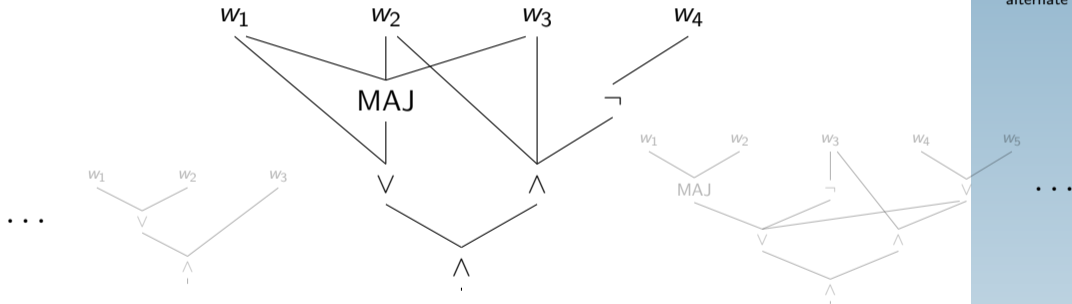
- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformer languages

Preliminary 1: Circuit complexity

$\text{size}(n)$ = number of gates
= **poly**(n)

$\text{depth}(n)$ = maximal length of a path
= **constant**



$$L_4 = \{1110, \dots\}$$

$$L = \bigcup_{i \in \mathbb{N}} L_i$$

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

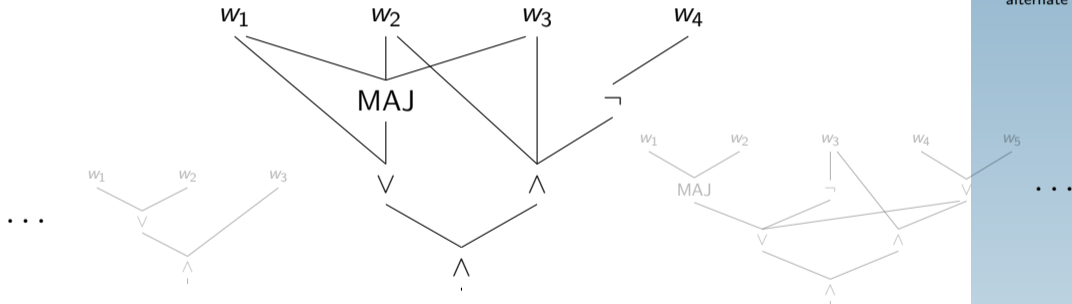
- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformer languages

Preliminary 1: Circuit complexity

$\text{size}(n)$ = number of gates
= **poly**(n)

$\text{depth}(n)$ = maximal length of a path
= **constant**



$$L_4 = \{1110, \dots\}$$

$$L = \bigcup_{i \in \mathbb{N}} L_i$$

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

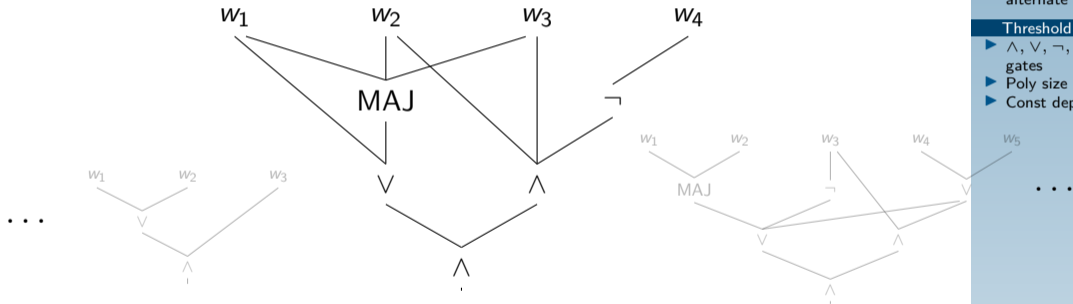
- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Transformer languages

Preliminary 1: Circuit complexity

$\text{size}(n)$ = number of gates
= **poly**(n)

$\text{depth}(n)$ = maximal length of a path
= **constant**



$$L_4 = \{1110, \dots\}$$

$$L = \bigcup_{i \in \mathbb{N}} L_i$$

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Transformer languages

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Transformer languages

Theorem (Merrill and Sabharwal, 2023)

Any transformer language is implementable as a threshold circuit family

Idea.

- ▶ This requires bounding precision (but not to a constant)

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Transformer languages

Theorem (Merrill and Sabharwal, 2023)

Any transformer language is implementable as a threshold circuit family

Idea.

- ▶ This requires bounding precision (but not to a constant)
- ▶ All the operations used in “executing” a transformer are implementable by threshold circuits:

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Transformer languages

Theorem (Merrill and Sabharwal, 2023)

Any transformer language is implementable as a threshold circuit family

Idea.

- ▶ This requires bounding precision (but not to a constant)
- ▶ All the operations used in “executing” a transformer are implementable by threshold circuits:
 - ▶ Addition, multiplication, division, max (for FFNN), exp, iterated addition

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Transformer languages

Theorem (Merrill and Sabharwal, 2023)

Any transformer language is implementable as a threshold circuit family

Idea.

- ▶ This requires bounding precision (but not to a constant)
- ▶ All the operations used in “executing” a transformer are implementable by threshold circuits:
 - ▶ Addition, multiplication, division, max (for FFNN), exp, iterated addition
 - ▶ On floating point numbers (precision $\sim \log$ of input length)

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Transformer languages

Theorem (Merrill and Sabharwal, 2023)

Any transformer language is implementable as a threshold circuit family

Idea.

- ▶ This requires bounding precision (but not to a constant)
- ▶ All the operations used in “executing” a transformer are implementable by threshold circuits:
 - ▶ Addition, multiplication, division, max (for FFNN), exp, iterated addition
 - ▶ On floating point numbers (precision $\sim \log$ of input length)
- ▶ Circuit depth correlates with transformer layer-depth □

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Theorem (Merrill and Sabharwal, 2023)

Any transformer language is implementable as a threshold circuit family

Idea.

- ▶ This requires bounding precision (but not to a constant)
- ▶ All the operations used in “executing” a transformer are implementable by threshold circuits:
 - ▶ Addition, multiplication, division, max (for FFNN), exp, iterated addition
 - ▶ On floating point numbers (precision $\sim \log$ of input length)
- ▶ Circuit depth correlates with transformer layer-depth □
- ▶ Threshold circuits are in LogSpace, which is in PolyTime

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Theorem (Merrill and Sabharwal, 2023)

Any transformer language is implementable as a threshold circuit family

Idea.

- ▶ This requires bounding precision (but not to a constant)
- ▶ All the operations used in “executing” a transformer are implementable by threshold circuits:
 - ▶ Addition, multiplication, division, max (for FFNN), exp, iterated addition
 - ▶ On floating point numbers (precision $\sim \log$ of input length)
- ▶ Circuit depth correlates with transformer layer-depth □
- ▶ Threshold circuits are in LogSpace, which is in PolyTime
- ▶ What about a converse? Or something exact?

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Transformer languages

Preliminary 2: Temporal logic with counting

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Transformer languages

Preliminary 2: Temporal logic with counting

By examples with $w \in \{a, b\}^*$:

- ▶ $w \models a$ The *last* letter of w is a

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Transformer languages

Preliminary 2: Temporal logic with counting

By examples with $w \in \{a, b\}^*$:

- ▶ $w \models a$ The *last* letter of w is a
- ▶ $w \models (\# \overleftarrow{[a]} = \# \overleftarrow{[b]})$ As many a 's as b 's

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Transformer languages

Preliminary 2: Temporal logic with counting

By examples with $w \in \{a, b\}^*$:

- ▶ $w \models a$ The *last* letter of w is a
- ▶ $w \models (\overset{\leftarrow}{\#}[a] = \overset{\leftarrow}{\#}[b])$ As many a 's as b 's
- ▶ $w \models \overset{\leftarrow}{\#}[\text{true}] = 1 \wedge a$

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Transformer languages

Preliminary 2: Temporal logic with counting

By examples with $w \in \{a, b\}^*$:

- ▶ $w \models a$ The *last* letter of w is a
- ▶ $w \models (\overleftarrow{\#}[a] = \overleftarrow{\#}[b])$ As many a 's as b 's
- ▶ $w \models \overleftarrow{\#}[\overleftarrow{\#}[\text{true}] = 1 \wedge a] = 1$

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Transformer languages

Preliminary 2: Temporal logic with counting

By examples with $w \in \{a, b\}^*$:

- ▶ $w \models a$ The *last* letter of w is a
- ▶ $w \models (\overleftarrow{\#}[a] = \overleftarrow{\#}[b])$ As many a 's as b 's
- ▶ $w \models \overleftarrow{\#}[\overleftarrow{\#}[\text{true}] = 1 \wedge a] = 1$ First letter is a

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Transformer languages

Preliminary 2: Temporal logic with counting

By examples with $w \in \{a, b\}^*$:

- ▶ $w \models a$ The *last* letter of w is a
- ▶ $w \models (\# \overleftarrow{[a]} = \# \overleftarrow{[b]})$ As many a 's as b 's
- ▶ $w \models \# \overleftarrow{[\# \overleftarrow{[true]} = 1 \wedge a]} = 1$ First letter is a
- ▶ $w \models \# \overleftarrow{[\# \overleftarrow{[a]} < \# \overleftarrow{[b]}]} = 0$ No prefix has more b 's than a 's

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Transformer languages

Preliminary 2: Temporal logic with counting

By examples with $w \in \{a, b\}^*$:

- ▶ $w \models a$ The *last* letter of w is a
- ▶ $w \models (\# \overleftarrow{[a]} = \# \overleftarrow{[b]})$ As many a 's as b 's
- ▶ $w \models \# \overleftarrow{[\# \overleftarrow{[true]} = 1 \wedge a]} = 1$ First letter is a
- ▶ $w \models \# \overleftarrow{[\# \overleftarrow{[a]} < \# \overleftarrow{[b]}]} = 0$ No prefix has more b 's than a 's
- ▶ Language of ϕ : $\{w \mid w \models \phi\}$

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Transformer languages

Preliminary 2: Temporal logic with counting

By examples with $w \in \{a, b\}^*$:

- ▶ $w \models a$ The *last* letter of w is a
- ▶ $w \models (\# \overleftarrow{[a]} = \# \overleftarrow{[b]})$ ~~As many a 's as b 's~~
- ▶ $w \models \# \overleftarrow{[\# \overleftarrow{[true]} = 1 \wedge a]} = 1$ First letter is a
- ▶ $w \models \# \overleftarrow{[\# \overleftarrow{[a]} < \# \overleftarrow{[b]}]} = 0$ No prefix has more b 's than a 's
- ▶ Language of ϕ : $\{w \mid w \models \phi\}$

\rightarrow WELL PARENTHESIZED
 $a(b)$

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Transformer languages

Preliminary 2: Temporal logic with counting

By examples with $w \in \{a, b\}^*$:

- ▶ $w \models a$ The *last* letter of w is a
 - ▶ $w \models (\overleftarrow{\#}[a] = \overleftarrow{\#}[b])$ ~~As many a 's as b 's~~
 - ▶ $w \models \overleftarrow{\#}[\overleftarrow{\#}[\text{true}] = 1 \wedge a] = 1$ First letter is a
 - ▶ $w \models \overleftarrow{\#}[\overleftarrow{\#}[a] < \overleftarrow{\#}[b]] = 0$ No prefix has more b 's than a 's
 - ▶ Language of ϕ : $\{w \mid w \models \phi\}$
 - ▶ $\text{TL}[\overleftarrow{\#}]$: temporal logic with $\overleftarrow{\#}$
- Handwritten notes:* \wedge → WELL PARENTHESIZED
 $a(b)$

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

Transformer languages

Preliminary 2: Temporal logic with counting

By examples with $w \in \{a, b\}^*$:

- ▶ $w \models a$ The *last* letter of w is a
 - ▶ $w \models (\overleftarrow{\#}[a] = \overleftarrow{\#}[b])$ ~~As many a 's as b 's~~
 - ▶ $w \models \overleftarrow{\#}[\overleftarrow{\#}[\text{true}] = 1 \wedge a] = 1$ First letter is a
 - ▶ $w \models \overleftarrow{\#}[\overleftarrow{\#}[a] < \overleftarrow{\#}[b]] = 0$ No prefix has more b 's than a 's
 - ▶ Language of ϕ : $\{w \mid w \models \phi\}$
 - ▶ $\text{TL}[\overleftarrow{\#}]$: temporal logic with $\overleftarrow{\#}$
- Handwritten notes:* \wedge → WELL PARENTHESIZED
 $a($ $b)$

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\text{TL}[\overleftarrow{\#}]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\overleftarrow{\#}$

Transformer languages

Preliminary 2: Temporal logic with counting

By examples with $w \in \{a, b\}^*$:

- ▶ $w \models a$ The *last* letter of w is a
 - ▶ $w \models (\# \overleftarrow{[a]} = \# \overleftarrow{[b]})$ ~~As many a 's as b 's~~
 - ▶ $w \models \# \overleftarrow{[\# \overleftarrow{[true]} = 1 \wedge a]} = 1$ First letter is a
 - ▶ $w \models \# \overleftarrow{[\# \overleftarrow{[a]} < \# \overleftarrow{[b]}]} = 0$ No prefix has more b 's than a 's
 - ▶ Language of ϕ : $\{w \mid w \models \phi\}$
- WELL PARENTHESIZED
a (b)
- ▶ $\text{TL}[\# \overleftarrow{\cdot}]$: temporal logic with $\# \overleftarrow{\cdot}$
 - ▶ $\text{TL}[\# \overleftarrow{\cdot}]$ cannot express $(a + b)^* aa(a + b)^*$ (Behle & Krebs 2008)

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\text{TL}[\# \overleftarrow{\cdot}]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\# \overleftarrow{\cdot}$

Transformer languages

Preliminary 2: Temporal logic with counting

By examples with $w \in \{a, b\}^*$:

- ▶ $w \models a$ The *last* letter of w is a
- ▶ $w \models (\# \overleftarrow{[a]} = \# \overleftarrow{[b]})$ ~~As many a 's as b 's~~
- ▶ $w \models \# \overleftarrow{[\# \overleftarrow{[true]} = 1 \wedge a]} = 1$ First letter is a
- ▶ $w \models \# \overleftarrow{[\# \overleftarrow{[a]} < \# \overleftarrow{[b]}]} = 0$ No prefix has more b 's than a 's
- ▶ Language of ϕ : $\{w \mid w \models \phi\}$ ^ ^ → WELL PARENTHESIZED
a (b)
- ▶ $\text{TL}[\# \overleftarrow{\quad}]$: temporal logic with $\# \overleftarrow{\quad}$
- ▶ $\text{TL}[\# \overleftarrow{\quad}]$ cannot express $(a + b)^* aa(a + b)^*$ (Behle & Krebs 2008)

Proposition (Behle & Krebs 2008)

$\text{TL}[\# \overleftarrow{\quad}]$ languages are all implementable as threshold circuits of linear size

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\text{TL}[\# \overleftarrow{\quad}]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\# \overleftarrow{\quad}$

Transformer languages

An exact characterization

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\overleftarrow{\text{TL}}[\#]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\#$

Transformer languages

An exact characterization

Theorem (Yang, Cadilhac, Chiang 2025)

Transformer languages are exactly those of $\text{TL}[\overline{\#}]$

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\text{TL}[\overline{\#}]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\overline{\#}$

Transformer languages

An exact characterization

Theorem (Yang, Cadilhac, Chiang 2025)

Transformer languages are exactly those of $\text{TL}[\overline{\#}]$

Idea.

- ▶ A more careful study of the behavior of attention

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\text{TL}[\overline{\#}]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\overline{\#}$

Transformer languages

An exact characterization

Theorem (Yang, Cadilhac, Chiang 2025)

Transformer languages are exactly those of $\text{TL}[\overleftarrow{\#}]$

Idea.

- ▶ A more careful study of the behavior of attention
- ▶ Rounding/precision is done in a slightly different way

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\text{TL}[\overleftarrow{\#}]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\overleftarrow{\#}$

Transformer languages

An exact characterization

Theorem (Yang, Cadilhac, Chiang 2025)

Transformer languages are exactly those of $\text{TL}[\overleftarrow{\#}]$

Idea.

- ▶ A more careful study of the behavior of attention
- ▶ Rounding/precision is done in a slightly different way
- ▶ The nesting of $\overleftarrow{\#}$ corresponds to the depth of the transformer



Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\text{TL}[\overleftarrow{\#}]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\overleftarrow{\#}$

Transformer languages

An exact characterization

Theorem (Yang, Cadilhac, Chiang 2025)

Transformer languages are exactly those of $\text{TL}[\overleftarrow{\#}]$

Idea.

- ▶ A more careful study of the behavior of attention
- ▶ Rounding/precision is done in a slightly different way
- ▶ The nesting of $\overleftarrow{\#}$ corresponds to the depth of the transformer



Lies, a bunch of lies! Certainly transformers can do $(a + b)^* aa(a + b)^*$!

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\text{TL}[\overleftarrow{\#}]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\overleftarrow{\#}$

Transformer languages

An exact characterization

Theorem (Yang, Cadilhac, Chiang 2025)

Transformer languages are exactly those of $\text{TL}[\overleftarrow{\#}]$

Idea.

- ▶ A more careful study of the behavior of attention
- ▶ Rounding/precision is done in a slightly different way
- ▶ The nesting of $\overleftarrow{\#}$ corresponds to the depth of the transformer

Lies, a bunch of lies! Certainly transformers can do $(a + b)^* aa(a + b)^*$!

- ▶ The *main* lie: *positional encoding*
- ▶ Some other lesser lies (fixed precision at some places)

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\text{TL}[\overleftarrow{\#}]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\overleftarrow{\#}$

Transformer languages

A depth hierarchy

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\overleftarrow{\text{TL}}[\#]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\#$

Transformer languages

A depth hierarchy

- ▶ L_k : all of words over $\{a, b\}$ that alternate k times between a 's and b 's

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\overleftarrow{\text{TL}}[\#]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\#$

Transformer languages

A depth hierarchy

- ▶ L_k : all of words over $\{a, b\}$ that alternate k times between a 's and b 's

$aaabbb\ aa \in L_3$

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\overline{TL[\#]}$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\#$

Transformer languages

A depth hierarchy

- ▶ L_k : all of words over $\{a, b\}$ that alternate k times between a 's and b 's

$aaabbb\ aa \in L_3$

Theorem (Yang, Cadilhac, Chiang 2025)

L_k is in depth- k $\overleftarrow{\text{TL}}[\#]$ but not in depth- $(k - 1)$ $\overleftarrow{\text{TL}}[\#]$

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\overleftarrow{\text{TL}}[\#]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\#$

Transformer languages

A depth hierarchy

- ▶ L_k : all of words over $\{a, b\}$ that alternate k times between a 's and b 's

$aaabbb\ aa \in L_3$

Theorem (Yang, Cadilhac, Chiang 2025)

L_k is in depth- k $\overleftarrow{\text{TL}}[\#]$ but not in depth- $(k - 1)$ $\overleftarrow{\text{TL}}[\#]$

Idea.

- ▶ Depth-1 $\overleftarrow{\text{TL}}[\#]$ languages are commutative (letters commute)

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\overleftarrow{\text{TL}}[\#]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\#$

Transformer languages

A depth hierarchy

- ▶ L_k : all of words over $\{a, b\}$ that alternate k times between a 's and b 's

$aaabbb\ aa \in L_3$

Theorem (Yang, Cadilhac, Chiang 2025)

L_k is in depth- k $\overleftarrow{\text{TL}}[\#]$ but not in depth- $(k - 1)$ $\overleftarrow{\text{TL}}[\#]$

Idea.

- ▶ Depth-1 $\overleftarrow{\text{TL}}[\#]$ languages are commutative (letters commute)
- ▶ A reduction lemma: if L_m is in depth n , then L_{m-1} is in depth $n - 1$ \square

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\overleftarrow{\text{TL}}[\#]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\#$

Transformer languages

A depth hierarchy

- ▶ L_k : all of words over $\{a, b\}$ that alternate k times between a 's and b 's

$aaabbb\ aa \in L_3$

Theorem (Yang, Cadilhac, Chiang 2025)

L_k is in depth- k $\overleftarrow{\text{TL}}[\#]$ but not in depth- $(k - 1)$ $\overleftarrow{\text{TL}}[\#]$

Idea.

- ▶ Depth-1 $\overleftarrow{\text{TL}}[\#]$ languages are commutative (letters commute)
- ▶ A reduction lemma: if L_m is in depth n , then L_{m-1} is in depth $n - 1$ \square

Corollary

Transformers of depth k can express L_k but not L_{k+1}

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\overleftarrow{\text{TL}}[\#]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\#$

Transformer languages

A depth hierarchy in Practice: Expressing vs Training

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\text{TL}[\overline{\#}]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\overline{\#}$

Transformer languages

A depth hierarchy in Practice: Expressing vs Training

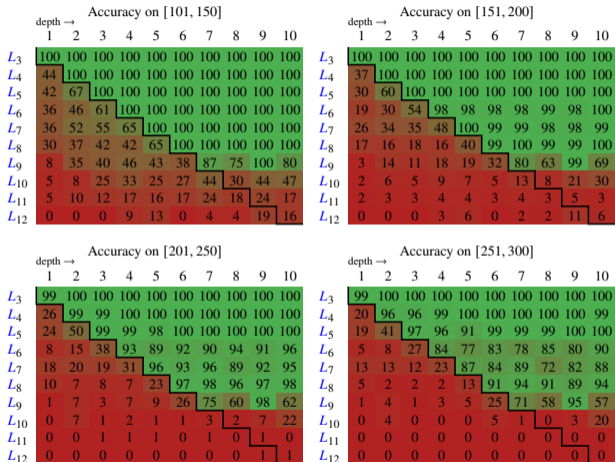


Figure 4: Experimental results. Corollary 5.2 predicts that a transformer with depth k can recognize language L_{k+2} but not L_{k+3} (demarcated by the black line). Up to L_9 , this closely predicts our experimental results (shown as numbers and colors).

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$TL[\#]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\#$

Outline

Transformers

What can transformers compute, if we give them time to think?

What can transformers compute in one evaluation?

Conclusion

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\overleftarrow{\text{TL}}[\#]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\#$

Conclusion

We saw:

- ▶ Transformers = just a bunch of (mostly) linear operations

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\overleftarrow{\text{TL}}[\#]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\#$

Conclusion

We saw:

- ▶ Transformers = just a bunch of (mostly) linear operations
- ▶ Allowed to “think:” Turing-complete

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\text{TL}[\#]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\#$

Conclusion

We saw:

- ▶ Transformers = just a bunch of (mostly) linear operations
- ▶ Allowed to “think:” Turing-complete
- ▶ As language acceptors:

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\text{TL}[\overline{\#}]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\overline{\#}$

Conclusion

We saw:

- ▶ Transformers = just a bunch of (mostly) linear operations
- ▶ Allowed to “think:” Turing-complete
- ▶ As language acceptors:
 - ▶ Expressible using threshold circuits

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\overleftarrow{\text{TL}}[\#]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\#$

Conclusion

We saw:

- ▶ Transformers = just a bunch of (mostly) linear operations
- ▶ Allowed to “think:” Turing-complete
- ▶ As language acceptors:
 - ▶ Expressible using threshold circuits
 - ▶ Exactly equivalent (more or less!) to $\text{TL}[\#]$

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\text{TL}[\#]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\#$

Conclusion

We saw:

- ▶ Transformers = just a bunch of (mostly) linear operations
- ▶ Allowed to “think:” Turing-complete
- ▶ As language acceptors:
 - ▶ Expressible using threshold circuits
 - ▶ Exactly equivalent (more or less!) to $\text{TL}[\overline{\#}]$
 - ▶ A depth hierarchy can be deduced

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\text{TL}[\overline{\#}]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\overline{\#}$

Conclusion

We saw:

- ▶ Transformers = just a bunch of (mostly) linear operations
- ▶ Allowed to “think:” Turing-complete
- ▶ As language acceptors:
 - ▶ Expressible using threshold circuits
 - ▶ Exactly equivalent (more or less!) to $\text{TL}[\overline{\#}]$
 - ▶ A depth hierarchy can be deduced

Questions:

- ▶ Main question: Learnability vs expressiveness

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\text{TL}[\overline{\#}]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\overline{\#}$

Conclusion

We saw:

- ▶ Transformers = just a bunch of (mostly) linear operations
- ▶ Allowed to “think:” Turing-complete
- ▶ As language acceptors:
 - ▶ Expressible using threshold circuits
 - ▶ Exactly equivalent (more or less!) to $\text{TL}[\overline{\#}]$
 - ▶ A depth hierarchy can be deduced

Questions:

- ▶ Main question: Learnability vs expressiveness
- ▶ Fix the lies (positional encoding!)

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\text{TL}[\overline{\#}]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\#$

Conclusion

We saw:

- ▶ Transformers = just a bunch of (mostly) linear operations
- ▶ Allowed to “think:” Turing-complete
- ▶ As language acceptors:
 - ▶ Expressible using threshold circuits
 - ▶ Exactly equivalent (more or less!) to $\text{TL}[\overline{\#}]$
 - ▶ A depth hierarchy can be deduced

Questions:

- ▶ Main question: Learnability vs expressiveness
- ▶ Fix the lies (positional encoding!)
- ▶ Can transformers contribute to circuit complexity?

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ \wedge, \vee, \neg , MAJ gates
- ▶ Poly size
- ▶ Const depth

$\text{TL}[\overline{\#}]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\#$

Conclusion

We saw:

- ▶ Transformers = just a bunch of (mostly) linear operations
- ▶ Allowed to “think:” Turing-complete
- ▶ As language acceptors:
 - ▶ Expressible using threshold circuits
 - ▶ Exactly equivalent (more or less!) to $\text{TL}[\#]$
 - ▶ A depth hierarchy can be deduced

Questions:

- ▶ Main question: Learnability vs expressiveness
- ▶ Fix the lies (positional encoding!)
- ▶ Can transformers contribute to circuit complexity?

On the circuit lower bound barrier. It was suggested in [Hah20] that unconditional lower bounds against *encoder-only* transformer would imply breakthrough circuit lower bounds against linear threshold circuits.¹ Our result avoids this barrier by exploiting the information bottleneck and autoregressive nature of *decoder-only* models. Namely, our proof crucially depends on the fact that in decoder-only models, each token cannot attend to any token after it.

Attention layer

For each pos i , compute score against j , weigh by a value, then take avg (softmax)

Transformer

- ▶ Input encoding
- ▶ Attention layers/FFNN alternate

Threshold Circuits

- ▶ $\wedge, \vee, \neg, \text{MAJ}$ gates
- ▶ Poly size
- ▶ Const depth

$\text{TL}[\#]$

- ▶ Temporal logic
- ▶ Evaluated at end
- ▶ With counting in the past $\#$