

What are Transformers?

- Introduced in "Attention is all you need" [Vaswani et al. @ NeurIPS 2017]
- Basic model used in recent Large Language Models (like ChatGPT, Gemini, ...)
- Applications in natural language processing, computer vision, language recognition, time series analysis, ...

Verifying Transformers?

- Sometimes their output is just wrong:

Could you please count the number of "R"s in the word "arbitrary"?

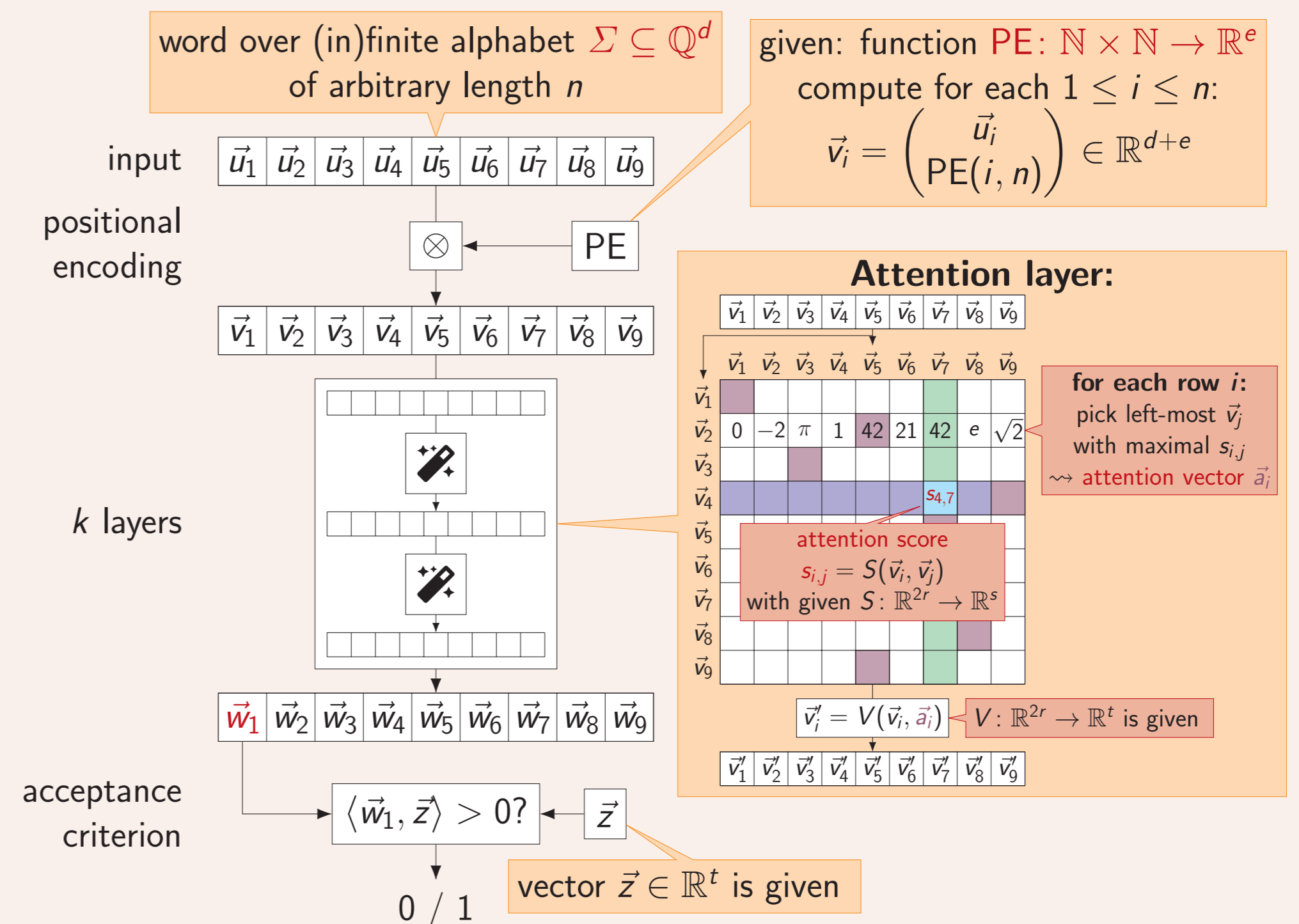
The word "arbitrary" contains 2 "R"s.

(Spoiler: "arbitrary" contains **three** "R"s)

Questions:

- What can be expressed via transformers?
 - e.g. $\text{PARITY} = \{w \in \{a, b\}^* \mid |w|_a \text{ is even}\}$ is accepted by a transformer [Chiang & Cholak @ ACL 2022]
 - What can transformers learn?
 - e.g. PARITY is not trainable via known algorithms [Bhattachamishra et al. @ EMNLP 2020]
 - Can we verify the (in)correctness of a transformer? If yes, how?
- Already a lot of work done for transformers on words.
 - Not (much) covered yet: transformers taking complex input data (like pictures, videos, voice, ...)

Unique Hard Attention Transformers (UHAT)



UHAT_{fin} = all languages over finite $\Sigma \subseteq \mathbb{Q}^d$ accepted by a UHAT
 UHAT_{inf} = all languages over infinite $\Sigma \subseteq \mathbb{Q}^d$ accepted by a UHAT

UHAT vs. Circuit Complexity

Circuit Complexity Classes AC^0 and TC^0

- AC^0 = all languages accepted by family of circuits of
- constant depth,
 - polynomial size, and
 - Boolean gates with unbounded fan-in.
- TC^0 extends AC^0 by majority gates.

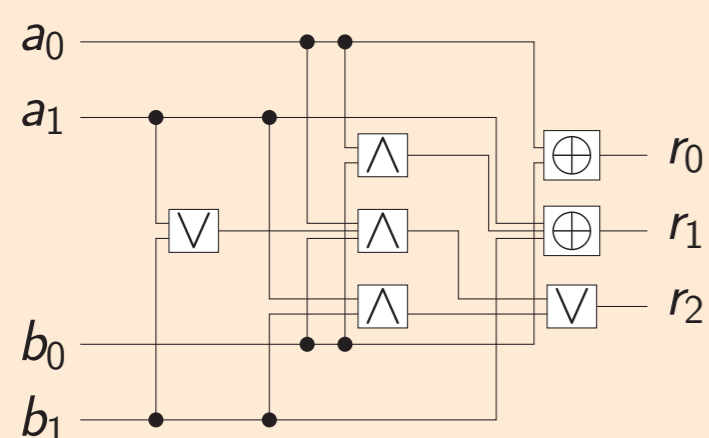


Figure: Circuit for $a_1a_0 + b_1b_0 = r_2r_1r_0$. In general, addition is in AC^0 .

Well-known: $\text{AC}^0 \subset \text{TC}^0$

Known Result

$\text{UHAT}_{\text{fin}} \subset \text{AC}^0$

[Hao et al. @ TACL 2022; Barceló et al. @ ICLR 2024]

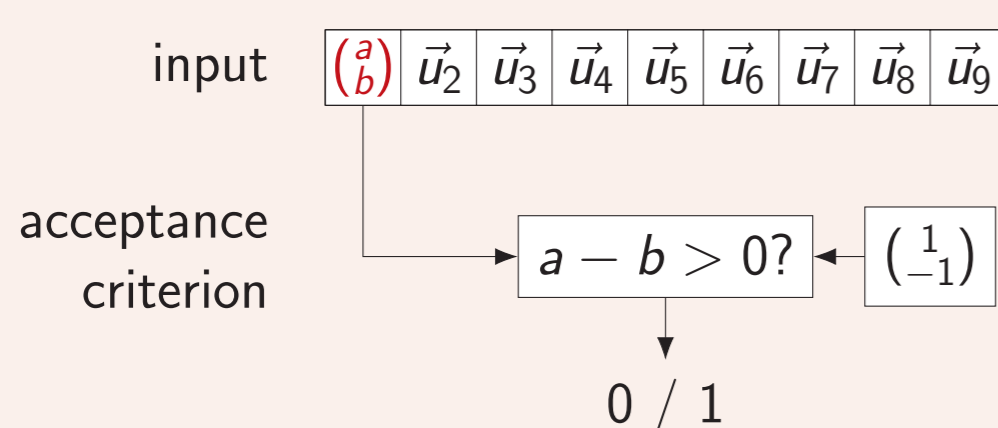
Theorem

$\text{UHAT}_{\text{inf}} \subseteq \text{TC}^0$

Theorem

$\text{COMPARE} = \left\{ \binom{a}{b} \cdot w \in (\mathbb{Q}^2)^* \mid a > b \right\}$
is in $\text{UHAT}_{\text{inf}} \setminus \text{AC}^0$.

Proof. The following UHAT (without positional encoding) and 0 layers accepts COMPARE:



UHAT vs. Regular Languages

Known Result

UHAT_{fin} (without positional encoding but with masking) = $\text{StarFree} \subset \text{Regular}$ [Yang et al. @ NeurIPS 2024]
 (StarFree = languages constructed from finite languages using Boolean operations and concatenation)

Symbolic Automata and Regular Data Languages

- Symbolic automata = NFA with arithmetic constraints as edge labels

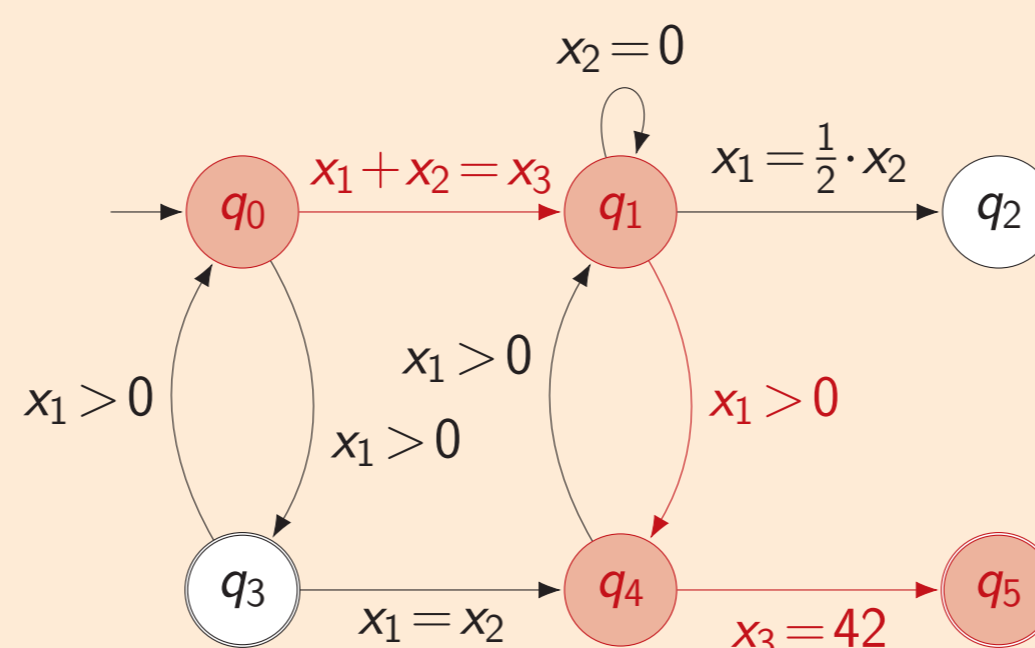


Figure: A symbolic automaton. Colored edges mark an accepting run of $\left(\binom{1}{3}, \binom{0}{0}, \binom{10}{42} \right)$

- $p \xrightarrow{\vec{x}} q$ if there is a transition (p, ϕ, q) with $\vec{x} \in \llbracket \phi \rrbracket$
- $(\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n) \in L(\mathfrak{A}) \iff I \xrightarrow{\vec{x}_1} \xrightarrow{\vec{x}_2} \dots \xrightarrow{\vec{x}_n} F$

Theorem

Define $\text{DOUBLE} \subseteq \mathbb{Q}^*$:

$$\{(x_1, x_2, \dots, x_n) \in \mathbb{Q}^* \mid \forall 1 \leq i < n: 2x_i < x_{i+1}\}.$$

$\text{DOUBLE} \notin \text{Regular}$, but $\text{DOUBLE} \in \text{UHAT}_{\text{inf}}$ (without positional encoding but with masking)

Proof. Non-regularity is shown by Pumping Lemma. DOUBLE is accepted by a UHAT with two layers:

- For each position i choose the $j > i$ maximizing $2x_i - x_j$ and let y_i be this maximal value.
- Check whether all y_i 's are non-positive.

UHAT vs. Logic

Linear Temporal Logic (LTL)

LTL describes words via following syntax and semantics:

	a	b	b	a	c	a	a	b
Atoms	a:	t	f	f	t	f	t	f
Boolean ops.	$a \vee b$:	t	t	t	t	f	t	t
NeXt op.	Xa :	f	f	t	f	t	t	f
Until op.	$a U b$:	t	t	t	f	f	t	t
Finally op.	Fc :	t	t	t	t	t	f	f
Globally op.	$G(a \vee b)$:	f	f	f	f	f	t	t

$L(\phi) = \{w \in \Sigma^* \mid \phi \text{ holds in first position of } w\}$

Well-known: $\text{LTL} = \text{StarFree}$ [Kamp 1968]

$\text{LTL}(\text{Mon})$ extends LTL by positional predicates.

Known Result

$\text{LTL}(\text{Mon}) \subseteq \text{UHAT}_{\text{fin}}$ [Barceló et al. @ ICLR 2024]

Locally Testable LTL (LTLTL)

LTLTL extends LTL(Mon) to alphabet $\Sigma = \mathbb{Q}^d$ by adding arithmetic constraints of the form

$$\langle (\vec{x}_i, \vec{x}_{i+1}, \dots, \vec{x}_{i+k}), \vec{a} \rangle > b.$$

Theorem

$\text{LTLTL} \subseteq \text{UHAT}_{\text{inf}}$

Example: 7-day Simple Moving Average

- Check for an "uptrend" in a time series
- Describes time sequences where the value at each time t is above the average of the week ending at t .
- $G(X^6 \text{ true} \rightarrow 7x_{t+6} > x_t + x_{t+1} + \dots + x_{t+6})$.

