

Chris Köcher¹ Alexander Kozachinskiy² Anthony Widjaja Lin^{1,3}
Marco Sälzer³ Georg Zetsche¹

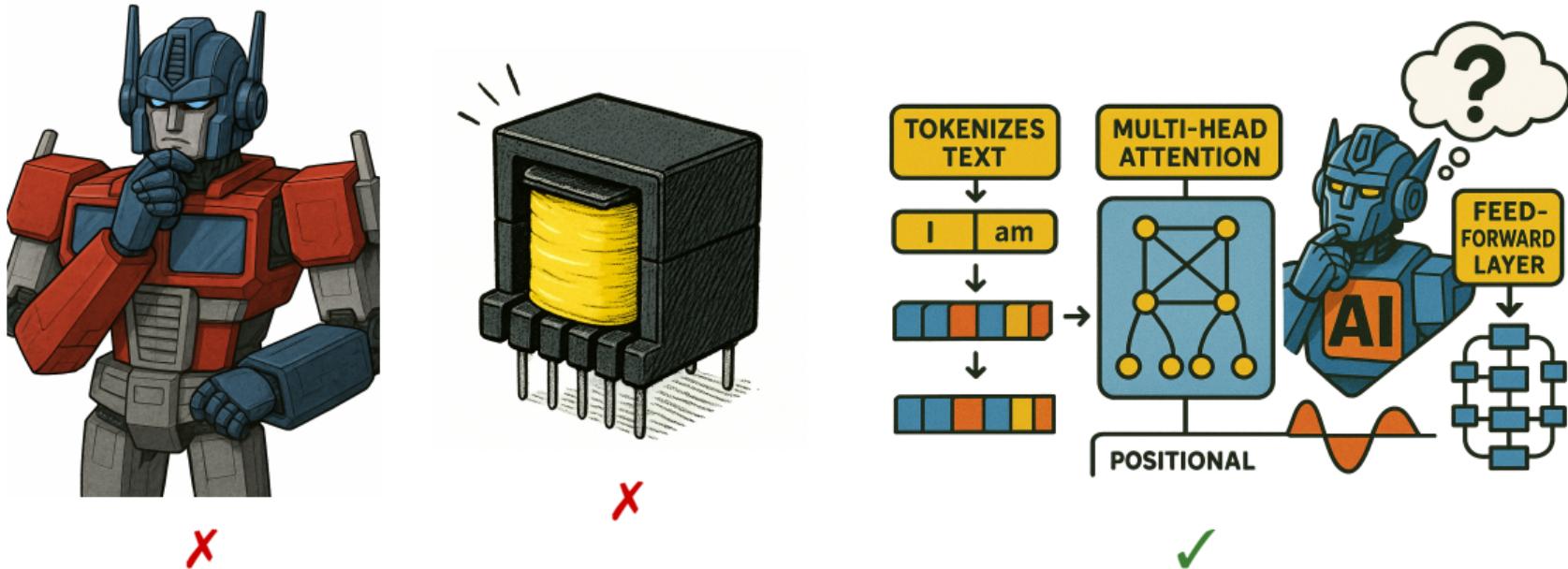
¹ Max Planck Institute for Software Systems, Kaiserslautern, Germany

² Centro Nacional de Inteligencia Artificial, Santiago, Chile

³ Rheinland-Pfälzische Technische Universität, Kaiserslautern & Landau, Germany

September 4, 2025

Transformers?



- Transformers are the basic model in machine learning used in recent LLMs



Motivation

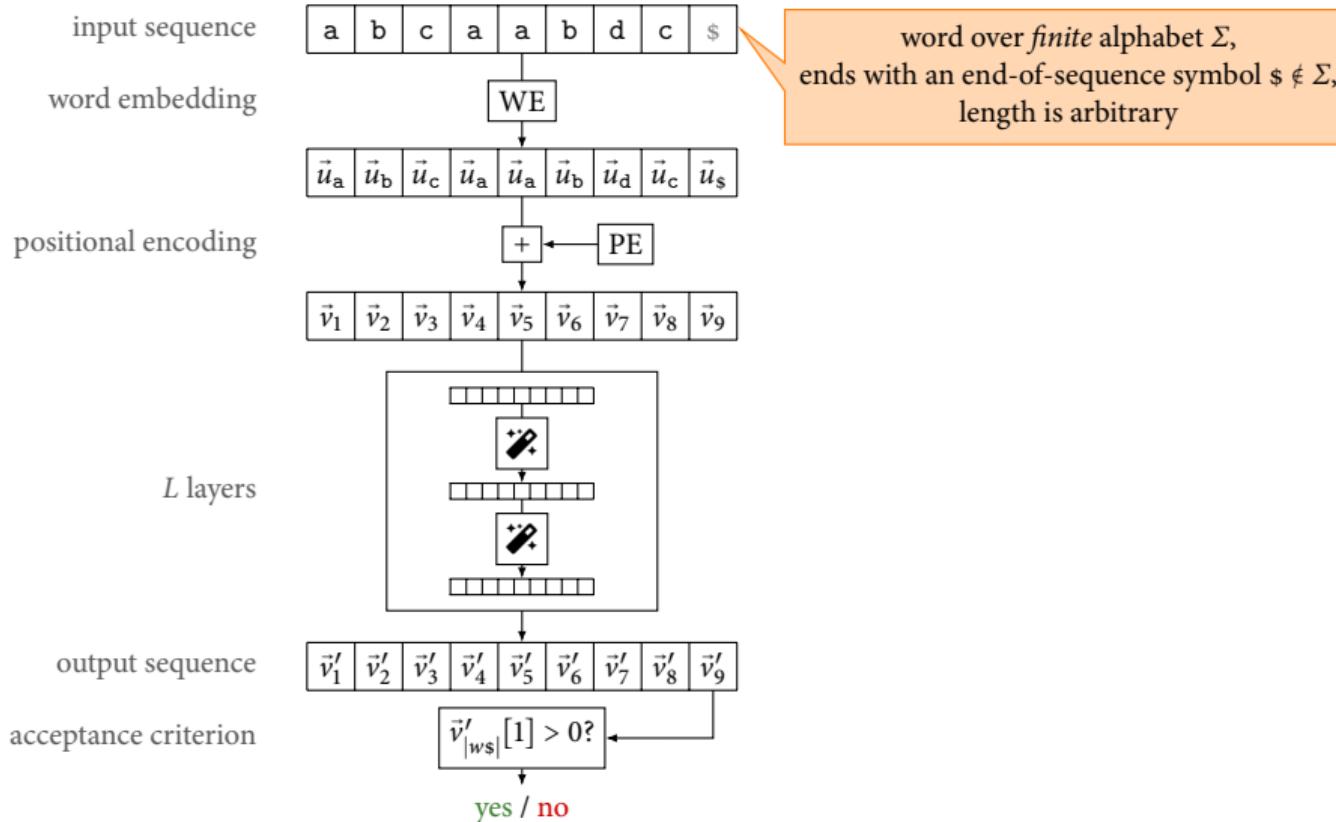
- Artificial intelligence is often not intelligent:



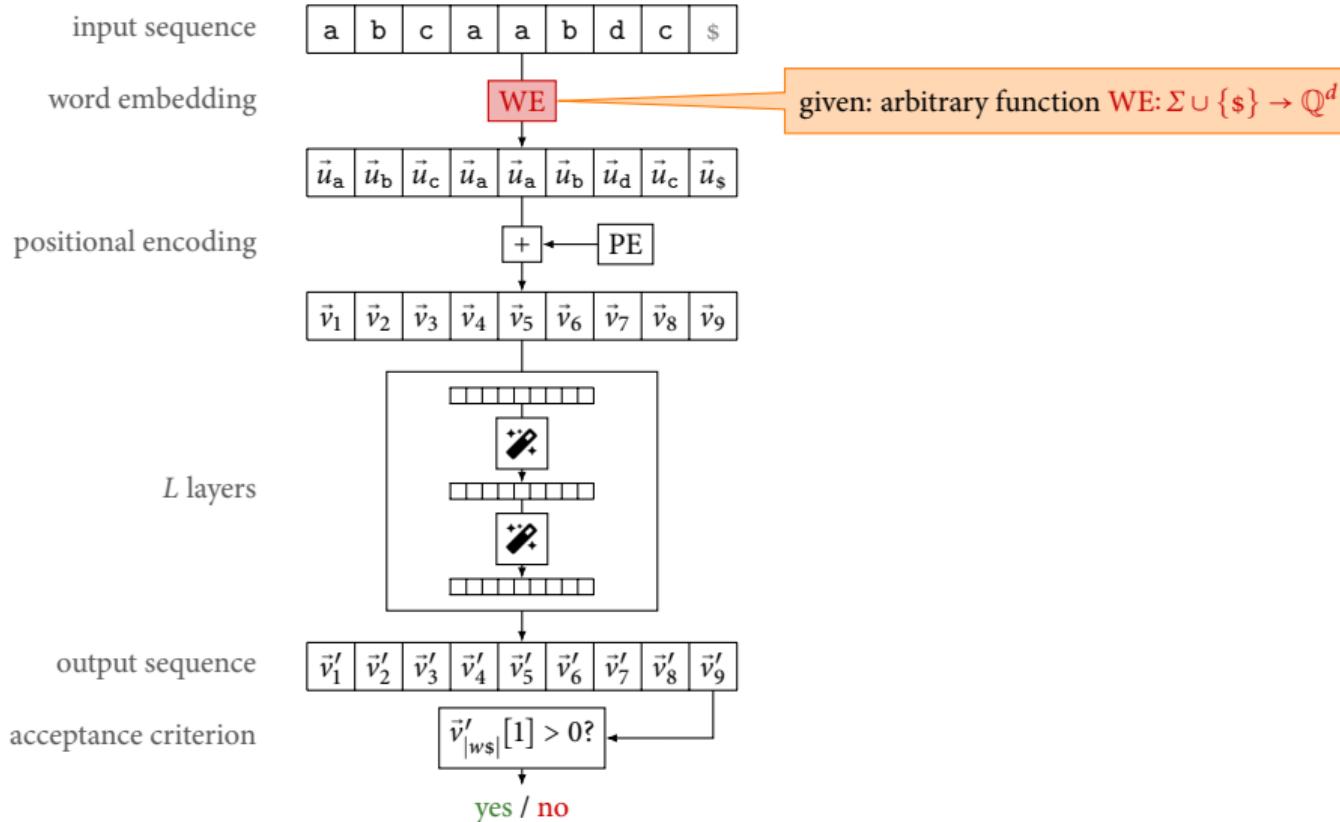
Source: Facebook, “15 most powerful passports in the world”.

- We try to understand what transformers actually can do.
 - A first step: study expressibility

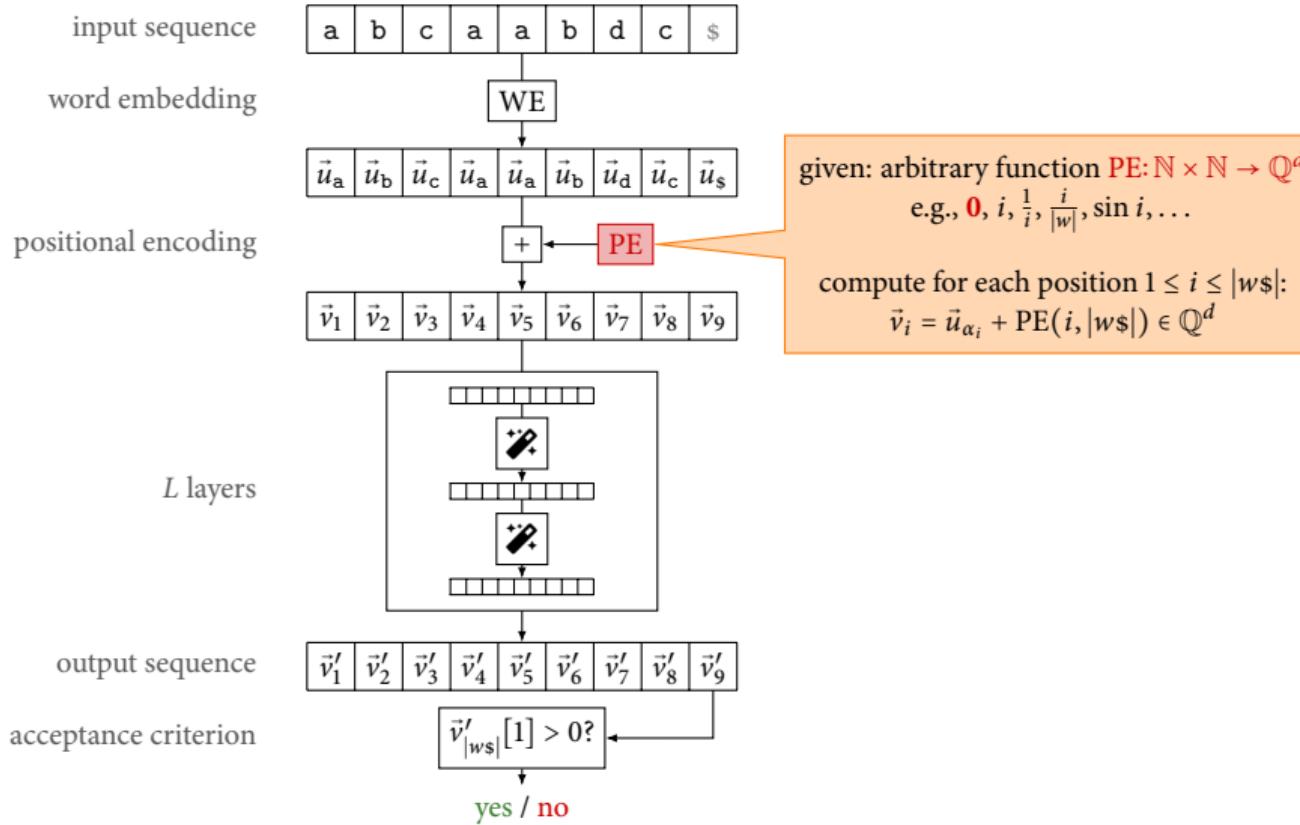
Transformers: Scheme



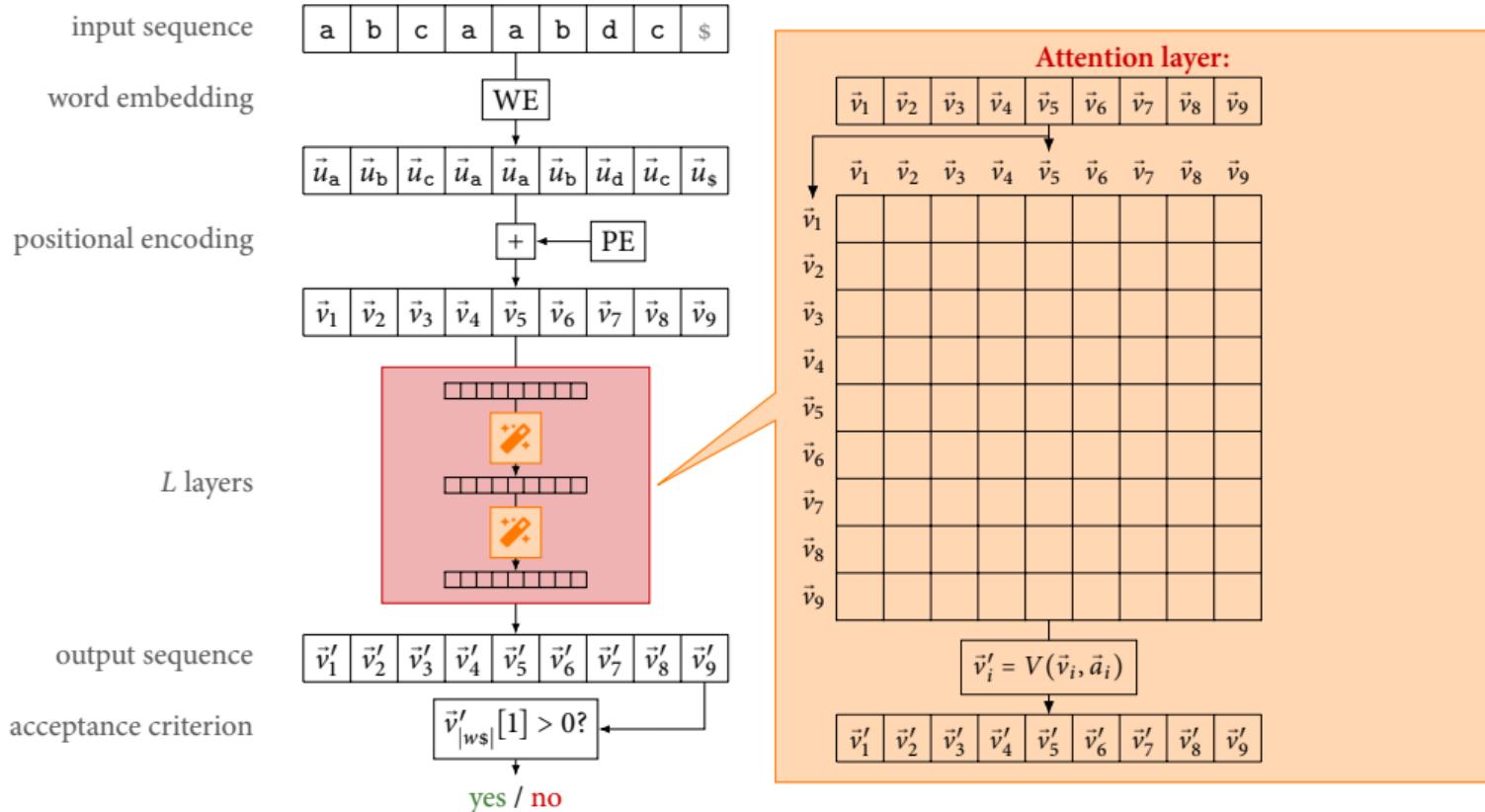
Transformers: Scheme



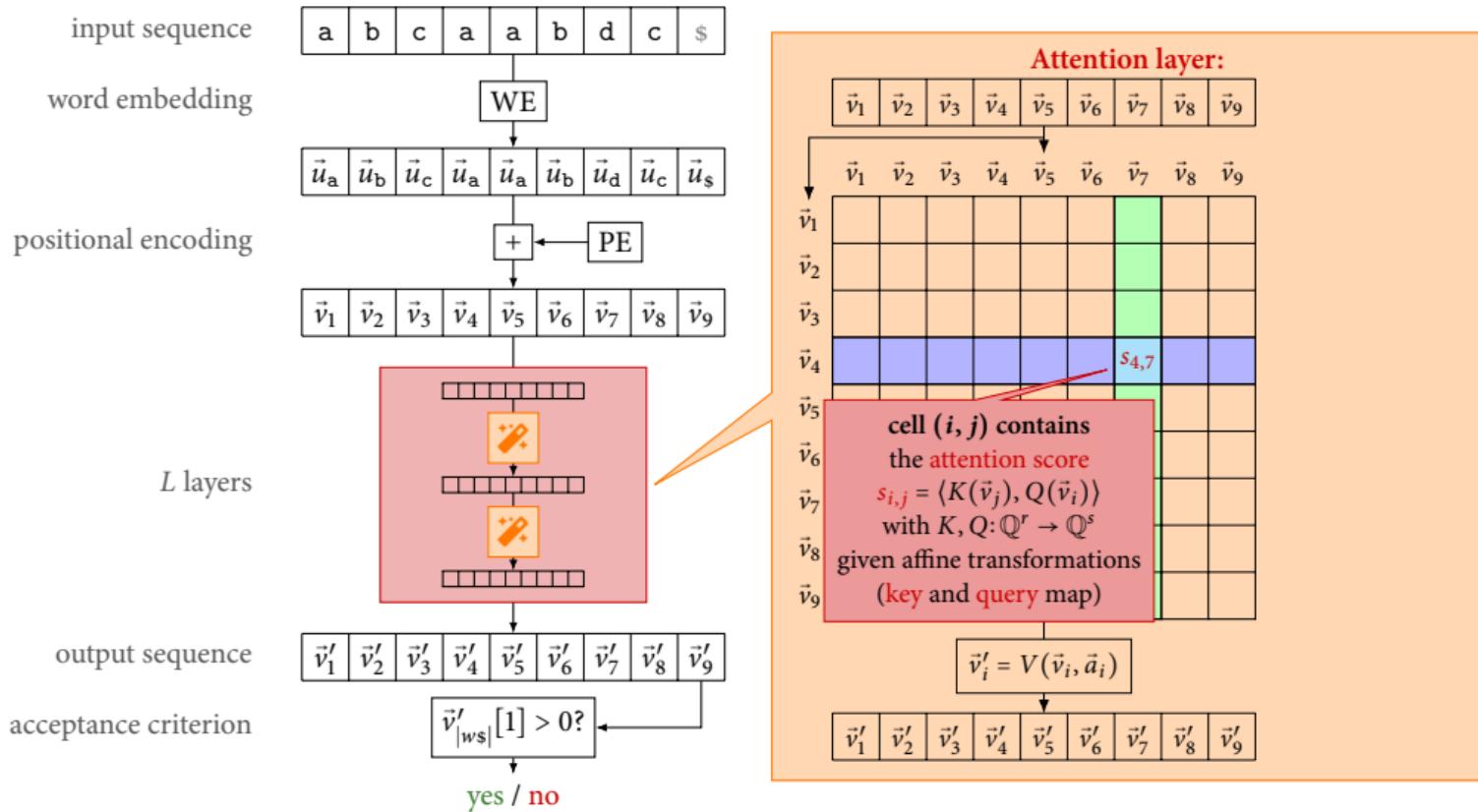
Transformers: Scheme



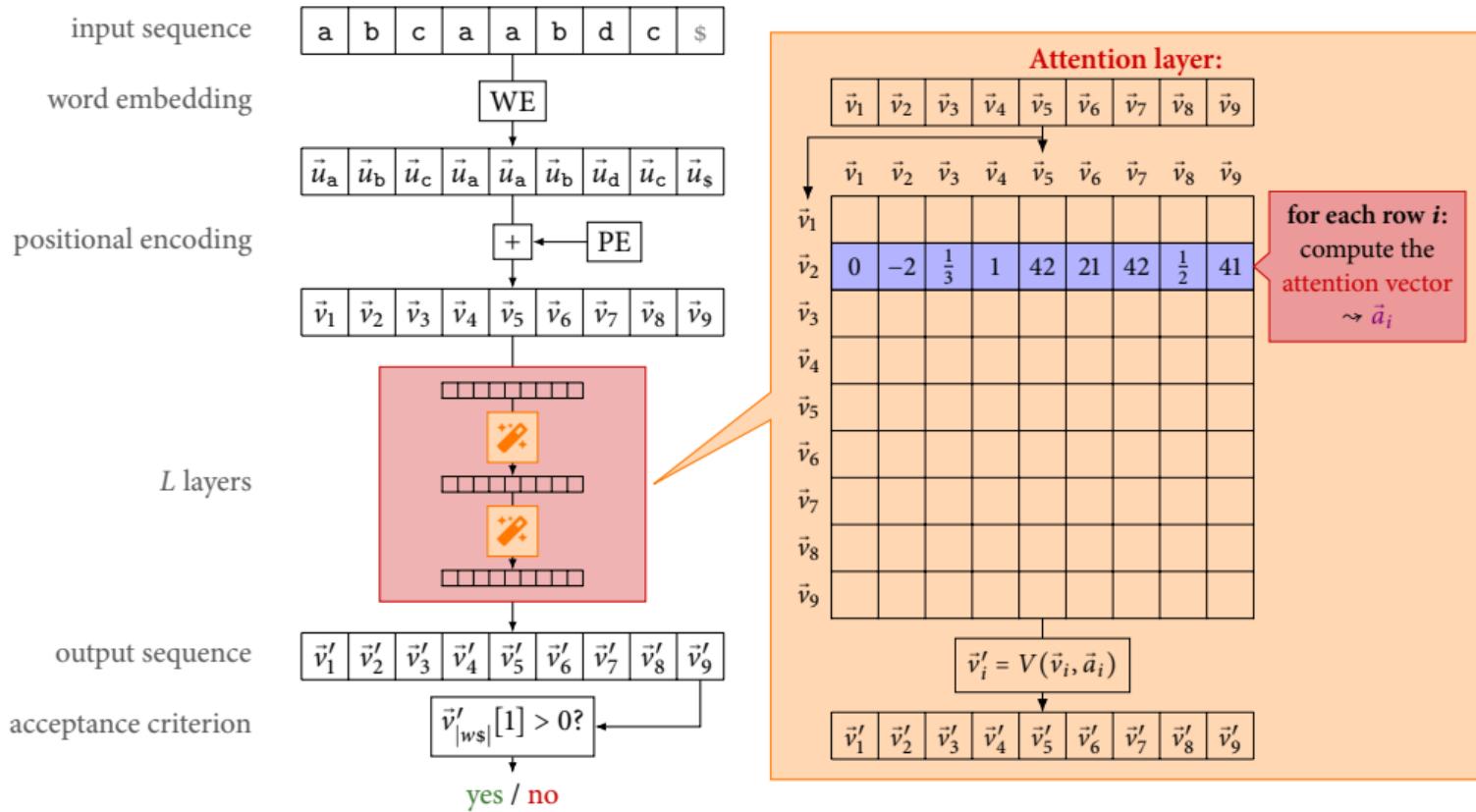
Transformers: Scheme



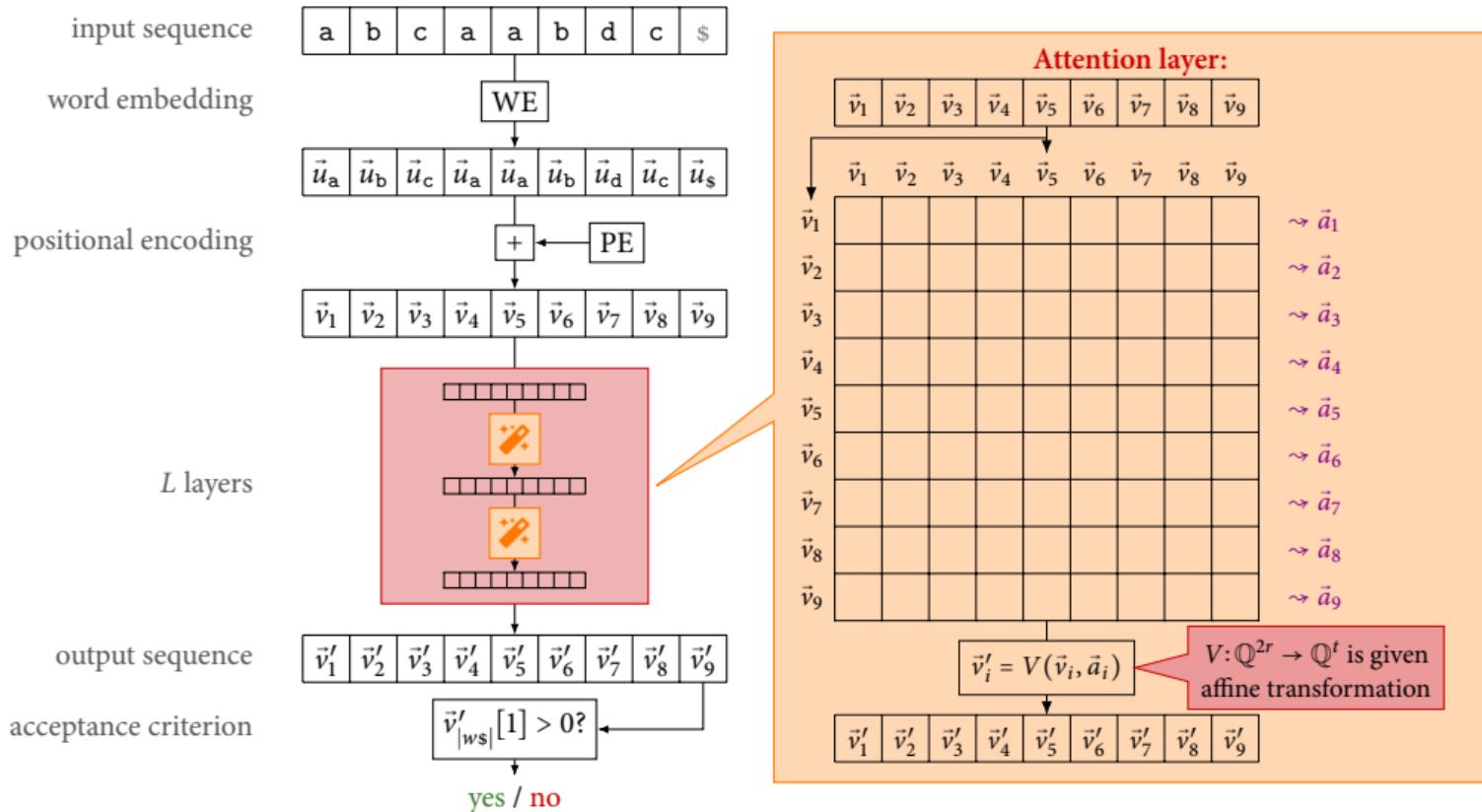
Transformers: Scheme



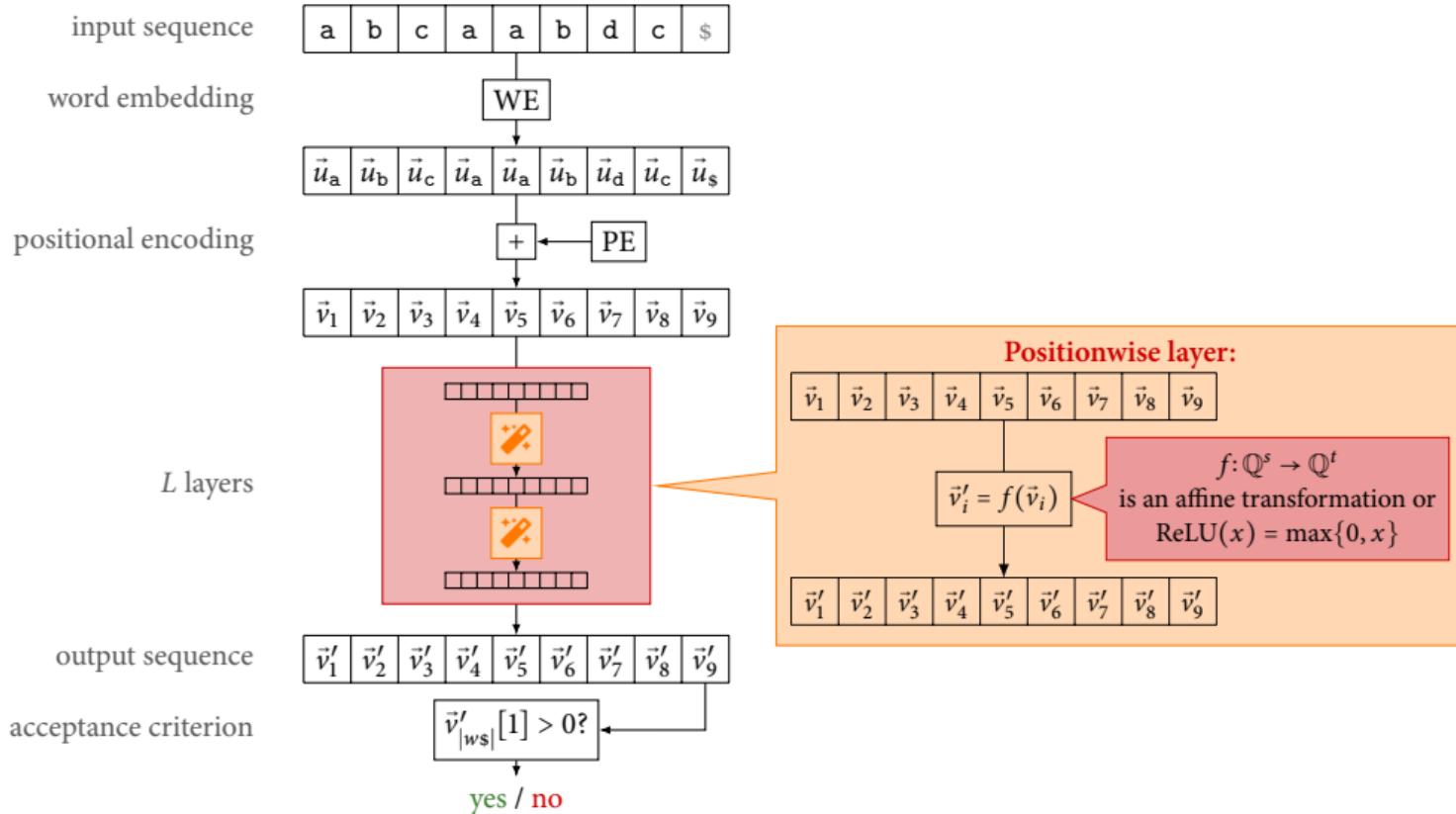
Transformers: Scheme



Transformers: Scheme



Transformers: Scheme



Transformers: Attention Mechanisms

1 Unique Hard Attention Transformers (UHAT)

	\vec{v}_1	\vec{v}_2	\vec{v}_3	\vec{v}_4	\vec{v}_5	\vec{v}_6	\vec{v}_7	\vec{v}_8	\vec{v}_9
\vec{v}_i	0	-2	$\frac{1}{3}$	1	42	21	42	$\frac{1}{2}$	41

$$\Rightarrow \vec{a}_i = \vec{v}_5$$

2 Average Hard Attention Transformers (AHAT)

	\vec{v}_1	\vec{v}_2	\vec{v}_3	\vec{v}_4	\vec{v}_5	\vec{v}_6	\vec{v}_7	\vec{v}_8	\vec{v}_9
\vec{v}_i	0	-2	$\frac{1}{3}$	1	42	21	42	$\frac{1}{2}$	41

$$\Rightarrow \vec{a}_i = \frac{1}{2}\vec{v}_5 + \frac{1}{2}\vec{v}_7$$

3 Softmax Attention Transformers with temperature scaling $\tau > 0$ (τ -SMAT)

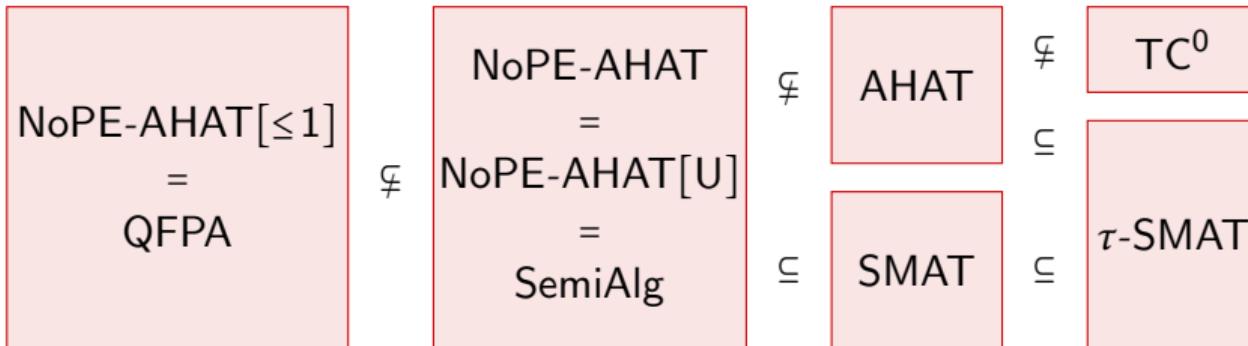
	\vec{v}_1	\vec{v}_2	\vec{v}_3	\vec{v}_4	\vec{v}_5	\vec{v}_6	\vec{v}_7	\vec{v}_8	\vec{v}_9
\vec{v}_i	0	-2	$\frac{1}{3}$	1	42	21	42	$\frac{1}{2}$	41

$$\Rightarrow \vec{a}_i = \sum_{j=1}^{|w\$|} \frac{e^{s_{i,j}/\tau}}{\sum_{k=1}^{|w\$|} e^{s_{i,k}/\tau}} \cdot \vec{v}_j$$

Transformers: Languages

- accepted language of a transformer \mathfrak{T} : $L(\mathfrak{T}) = \{w \in \Sigma^* \mid \mathfrak{T}(w\$) = \text{yes}\}$.
- AHAT / τ -SMAT / SMAT: all languages accepted by an AHAT / τ -SMAT / 1-SMAT
- NoPE-C: all languages accepted by a C-transformer without positional encoding (i.e., $\text{PE}(i, n) = \vec{0}$)
- C[U]: all languages accepted by a C-transformer such that each layer is uniform (i.e., if the key and query maps K and Q are constant).
- C[$\leq L$]: all languages accepted by a C-transformer with at most L attention layers.

Results



Definition

Let $\Sigma = \{a_1, a_2, \dots, a_d\}$ be an alphabet. The **Parikh map** is defined as

$$\Psi: \Sigma^* \rightarrow \mathbb{N}^d: w \mapsto (|w|_{a_1}, |w|_{a_2}, \dots, |w|_{a_d})$$

where $|w|_{a_i}$ is the number of occurrences of a_i in the word w .

Results

$L \subseteq \Sigma^*$ is semi-algebraic (in SemiAlg) if there are polynomials $p_1, p_2, \dots, p_k \in \mathbb{Z}[\vec{x}]$ such that $L = L_{p_1} \cup \dots \cup L_{p_k}$ where $L_{p_i} = \{w \in \Sigma^* \mid p_i(\Psi(w)) > 0\}$.

QFPA

SemiAlg

\subseteq

SMAT

\subseteq

τ -SMAT

$L \subseteq \Sigma^*$ is semilinear (in QFPA) if it is semi-algebraic such that all polynomials are linear (i.e., have degree ≤ 1).

\Updownarrow

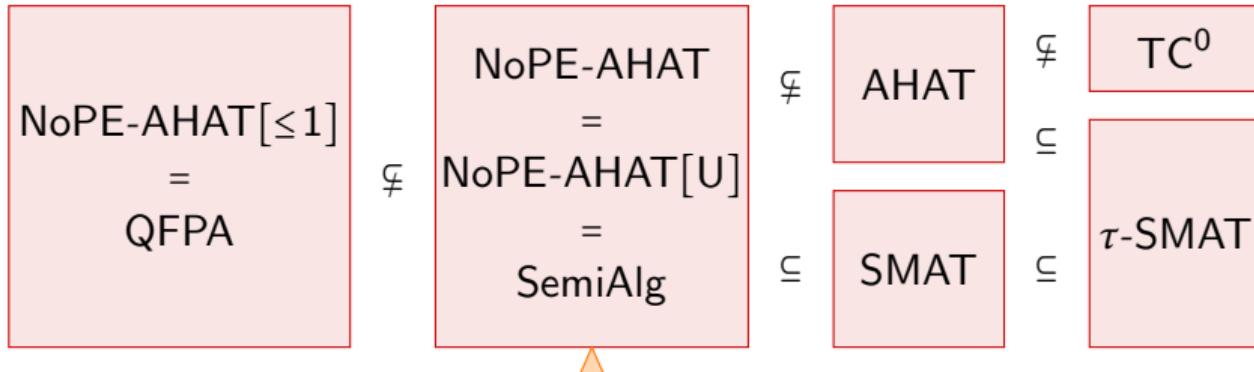
L is described by a Presburger formula.

defined as

$$\Psi: \Sigma^* \rightarrow \mathbb{N}^d: w \mapsto (|w|_{a_1}, |w|_{a_2}, \dots, |w|_{a_d})$$

where $|w|_{a_i}$ is the number of occurrences of a_i in the word w .

Results



Definition

Let $\Sigma = \{a_1, a_2, \dots\}$

Corollary

$$\text{Proj}(\text{NoPE-AHAT}) = \text{RE} \cap \text{PI}$$

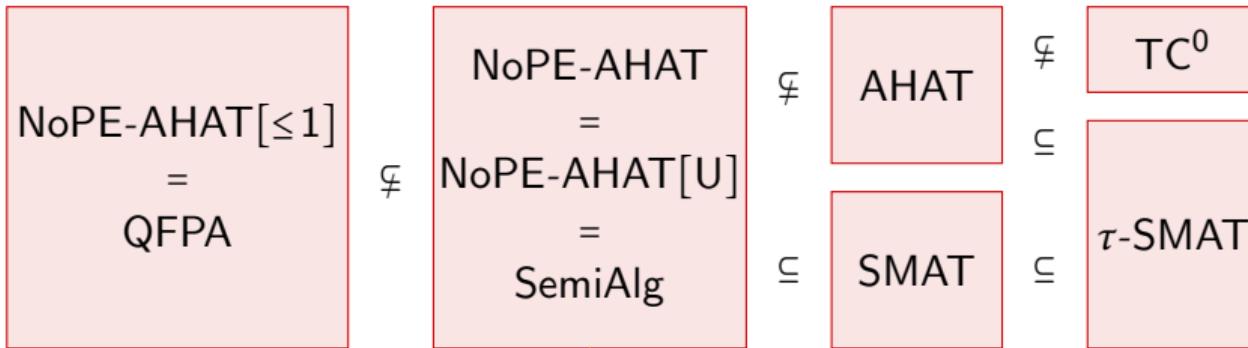
ned as

where $|w|_{a_i}$ is the n

Corollary

The emptiness problem is undecidable for (NoPE-)AHAT with two (uniform) layers.

Results



Definition

Let $\Sigma = \{a_1, a_2, \dots\}$

Corol recursively enumerable languages

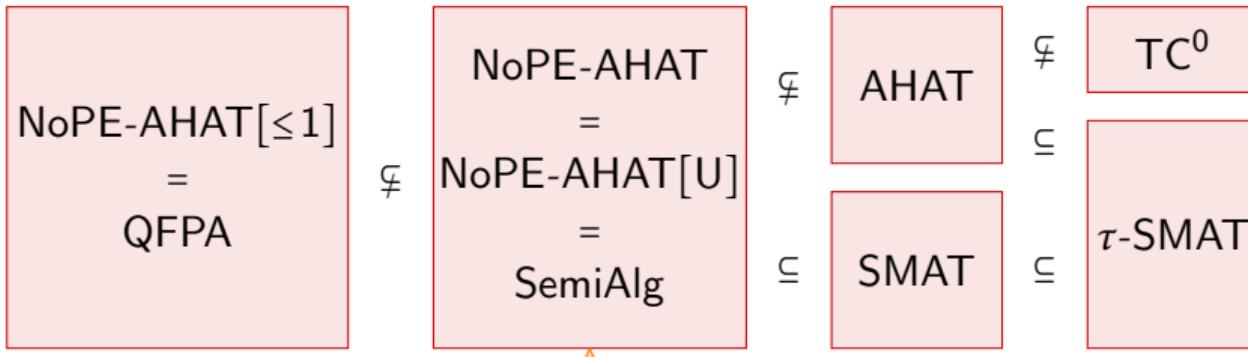
$\text{Proj}(\text{NoPE-AHAT}) = \text{RE} \cap \text{PI}$ permutation-invariant languages

languages $\pi(L)$ for $L \in \text{NoPE-AHAT}$
wh and projections $\pi: \Sigma \rightarrow \Gamma \cup \{\varepsilon\}$

$|_{a_d})$

is undecidable for
(NoPE-)AHAT with two (uniform) layers.

Results



Definition

Let $\Sigma = \{a_1, a_2, \dots\}$

Corollary

$\text{Proj}(\text{NoPE-AHAT}[U, \leq 2]) = \text{RE} \cap \text{PI}$

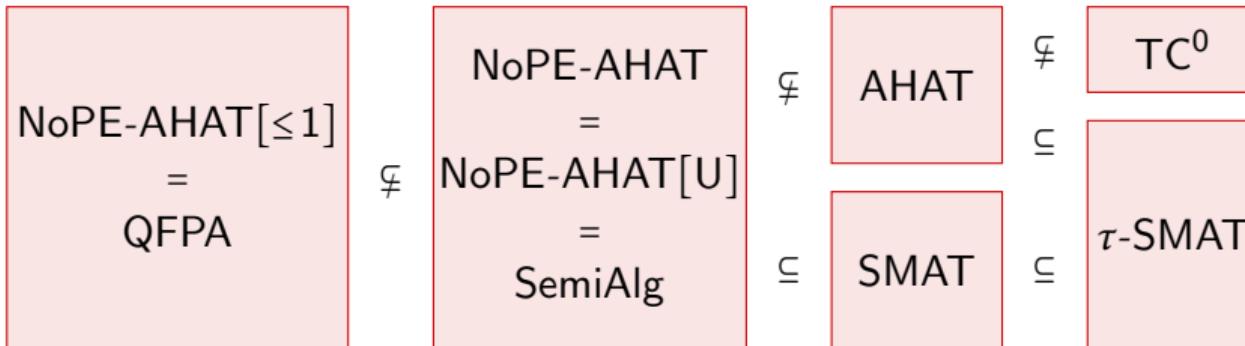
ned as

where $|w|_{a_i}$ is the r

Corollary

The emptiness problem is undecidable for (NoPE-)AHAT with two (uniform) layers.

Results



■ Open Problems:

- Are NoPE-AHAT \subseteq SMAT and AHAT, SMAT \subseteq τ -SMAT proper inclusions?
- Is NoPE-AHAT $[\leq L, U]$ equal to semi-algebraic sets with polynomials of degree $\leq L$?

Thank you!

