

Marco Sälzer¹ Chris Köcher² Alexander Kozachinskiy³
Georg Zetzsche² Anthony Widjaja Lin^{1,2}

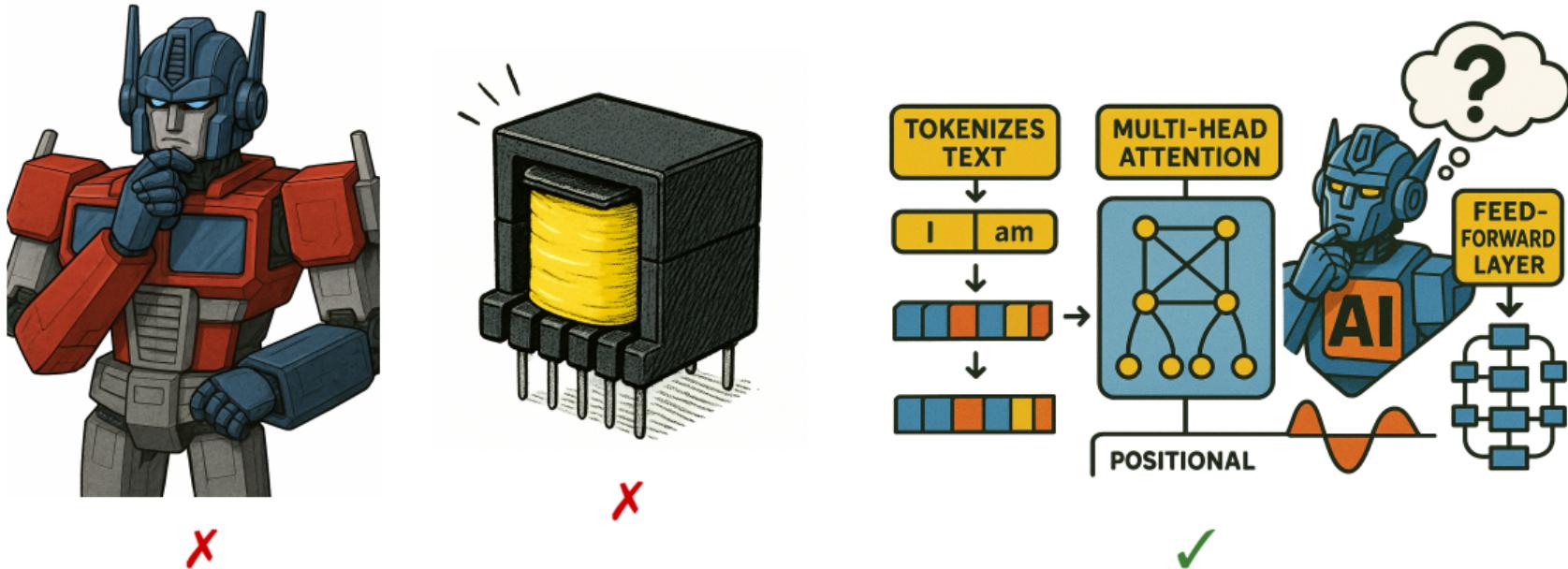
¹ Rheinland-Pfälzische Technische Universität, Kaiserslautern & Landau, Germany

² Max Planck Institute for Software Systems, Kaiserslautern, Germany

³ Centro Nacional de Inteligencia Artificial, Santiago, Chile

October 2, 2025

Transformers?



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



Motivation

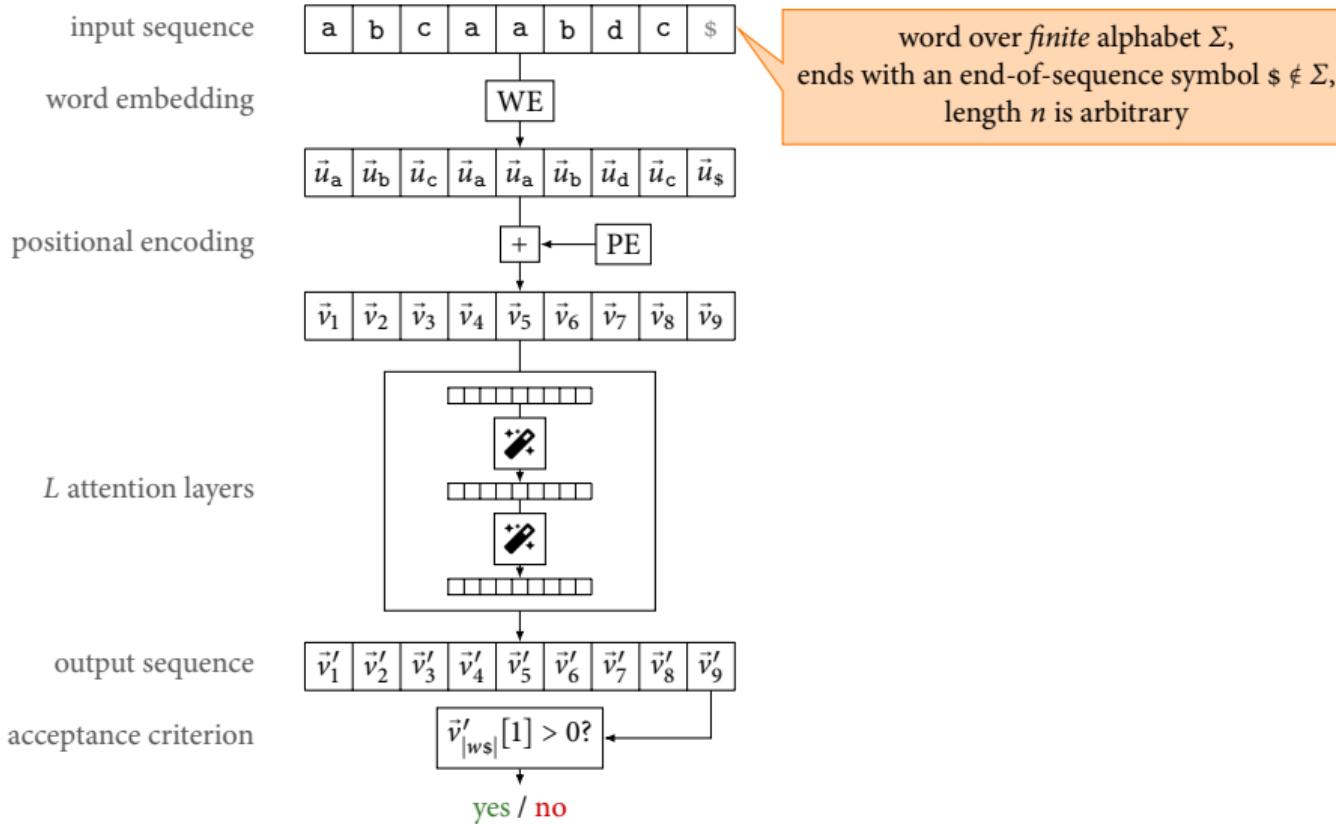
- Artificial intelligence is often not intelligent:



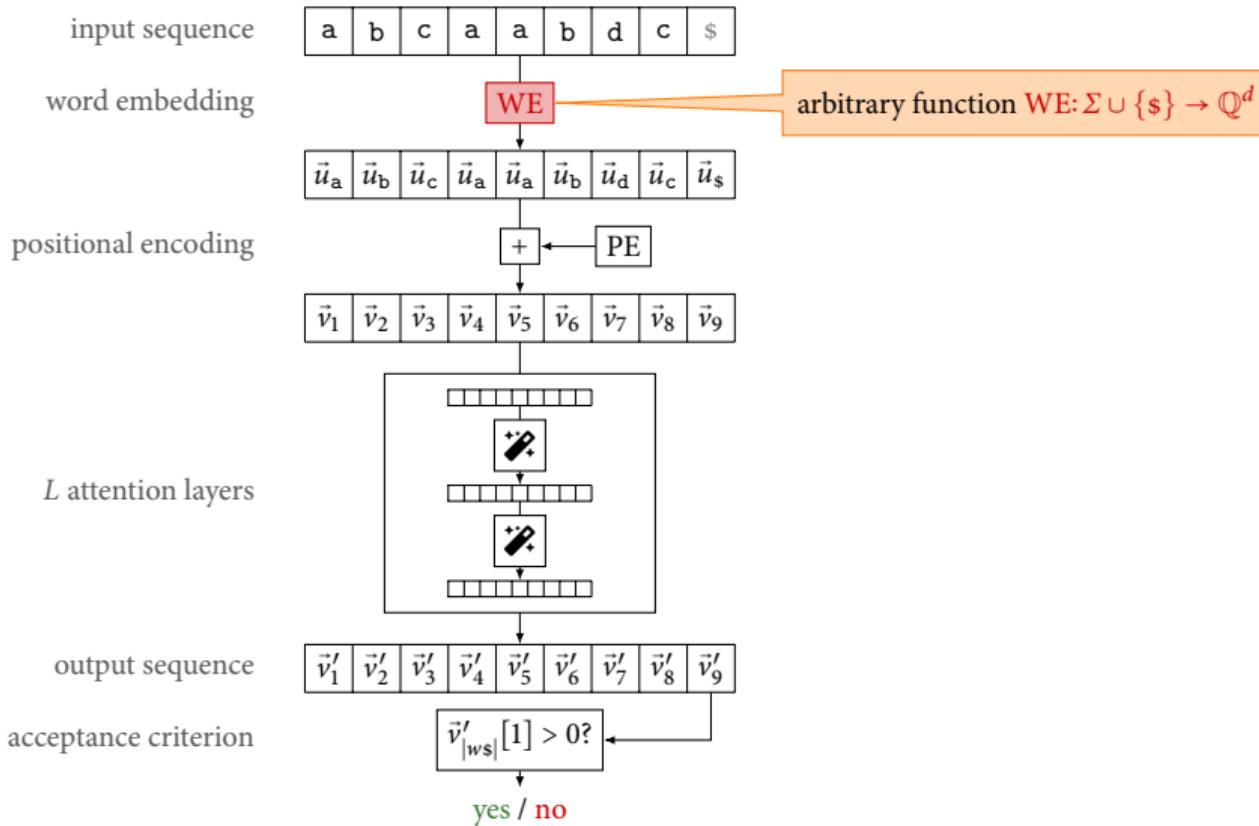
Source: Facebook, “15 most powerful passports in the world”.

- We try to verify transformers.
 - A first step: study expressibility

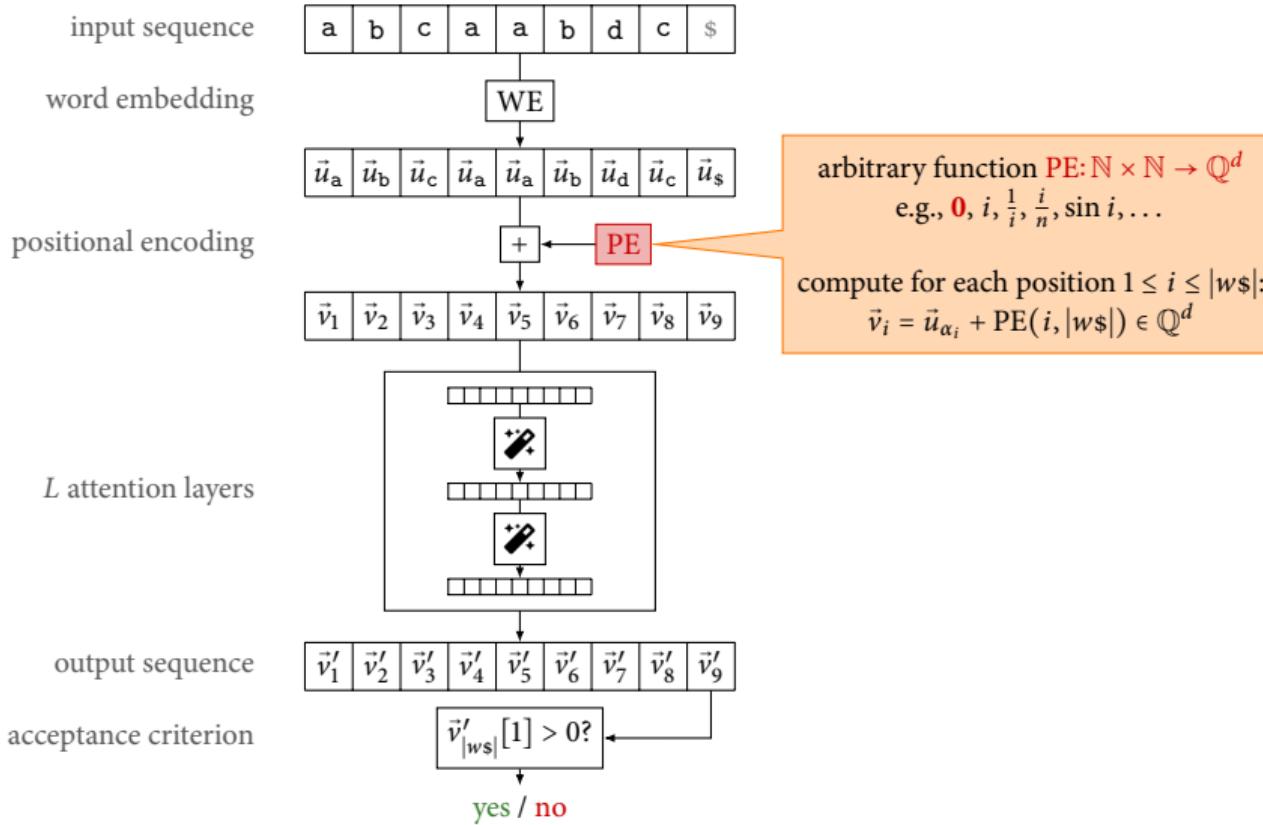
Transformers: Scheme



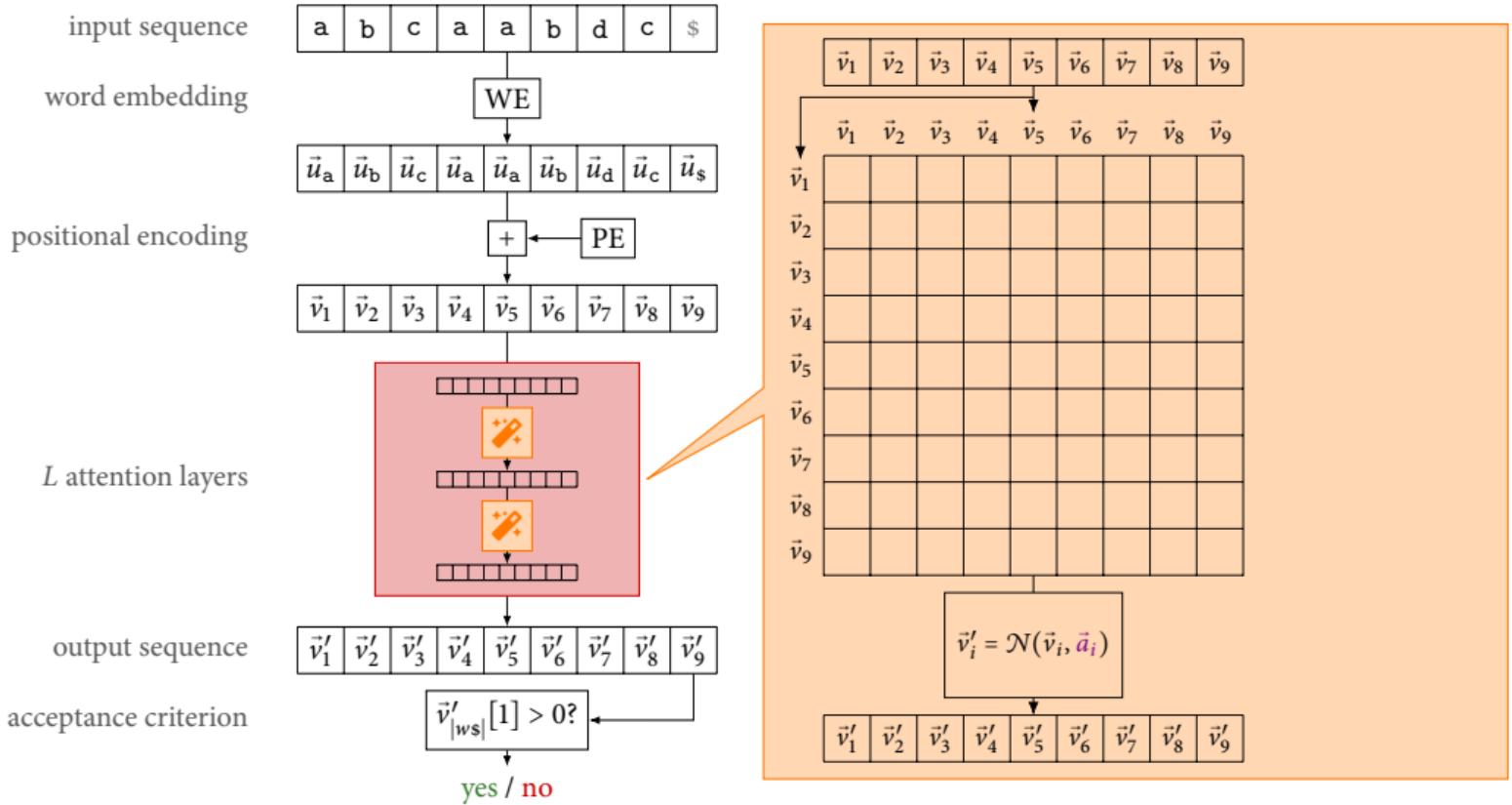
Transformers: Scheme



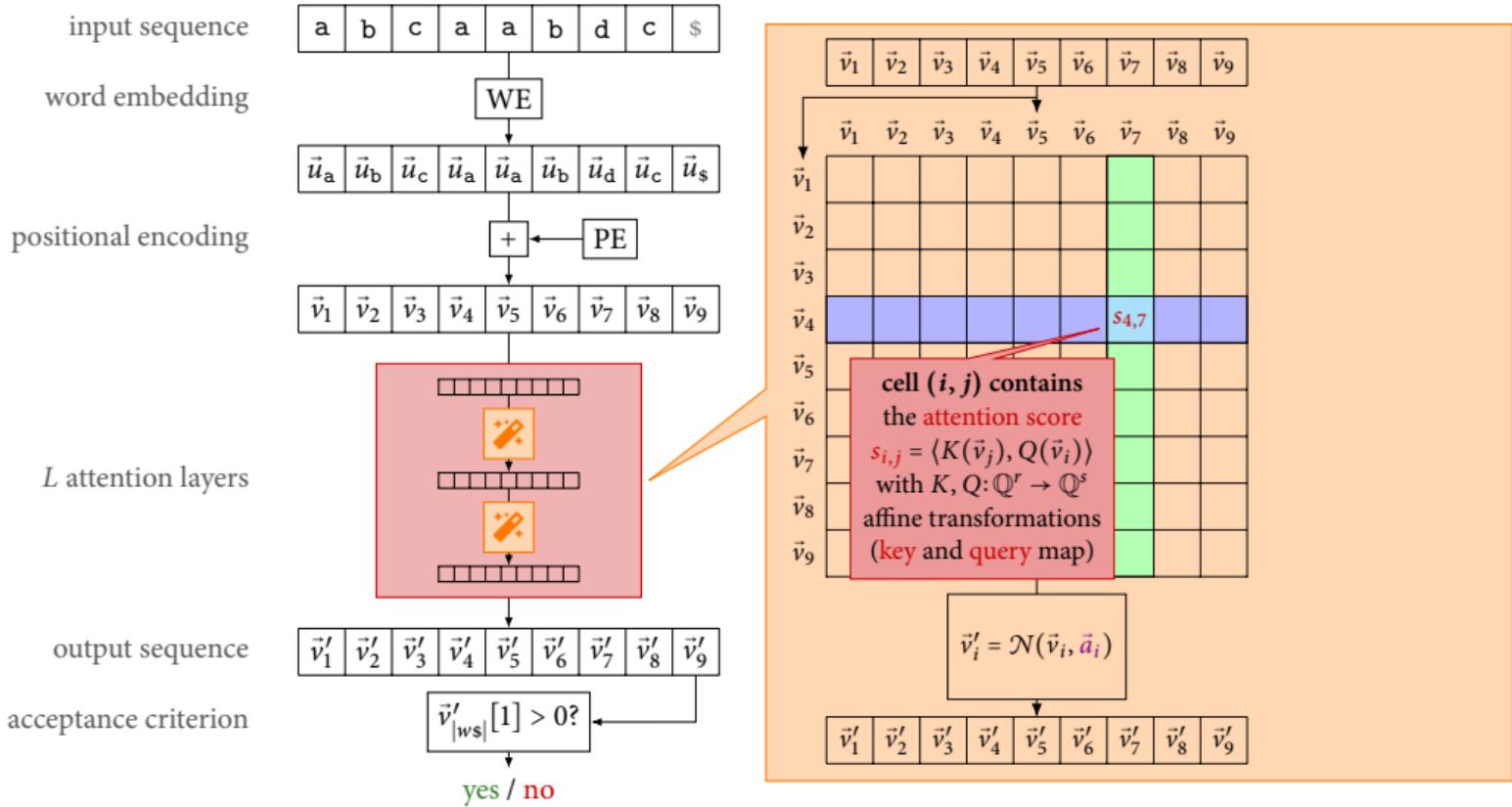
Transformers: Scheme



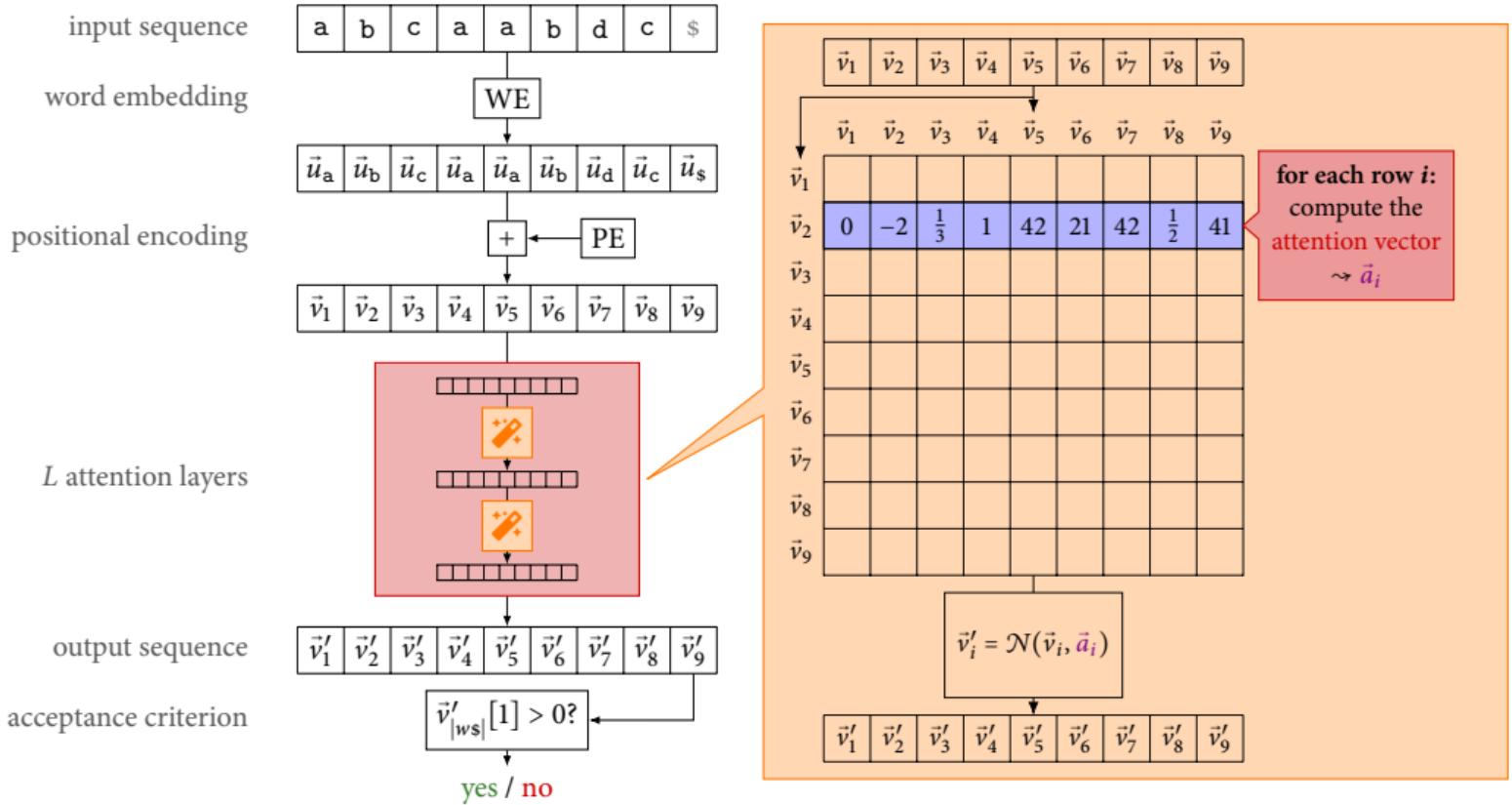
Transformers: Scheme



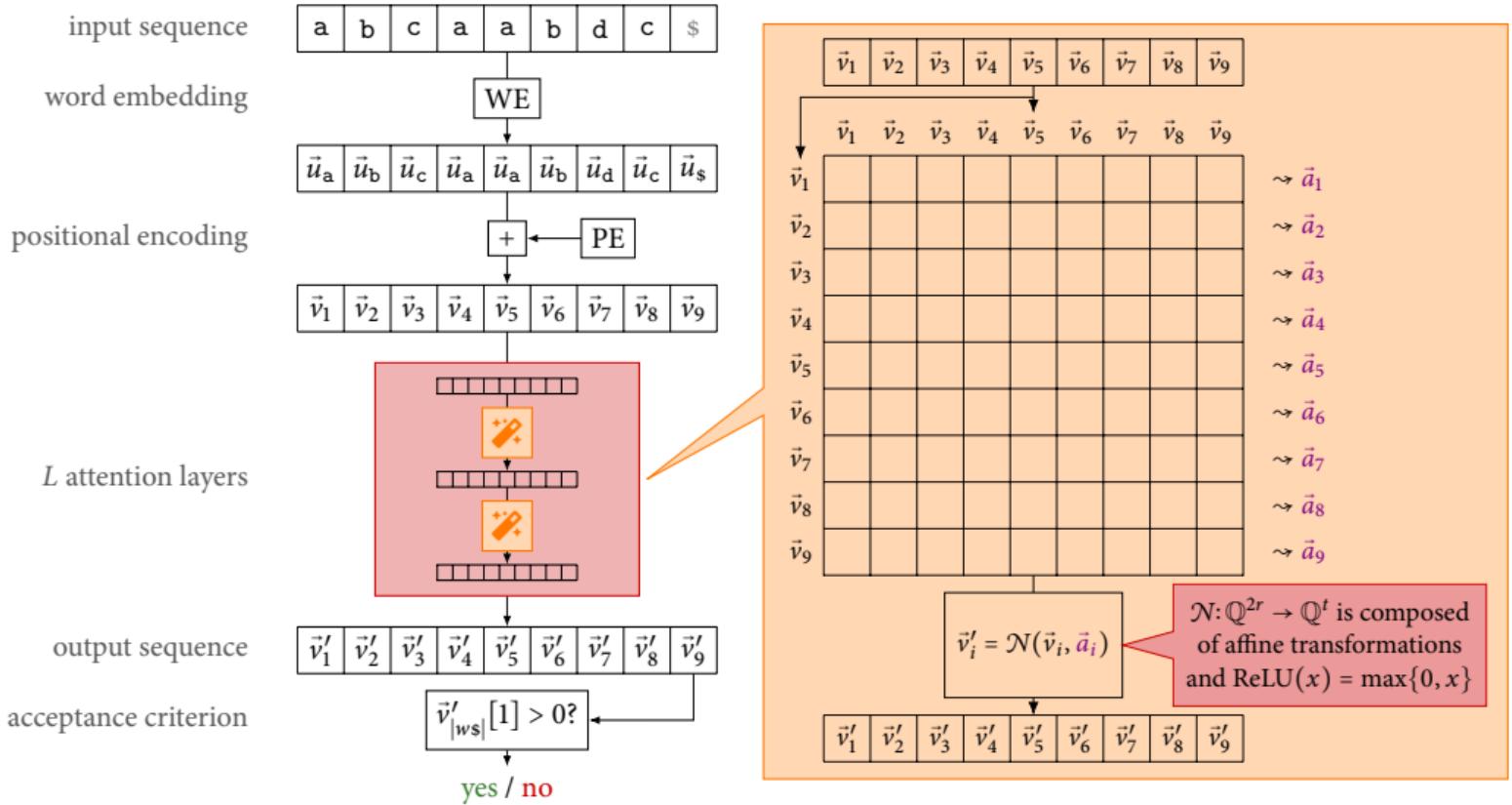
Transformers: Scheme



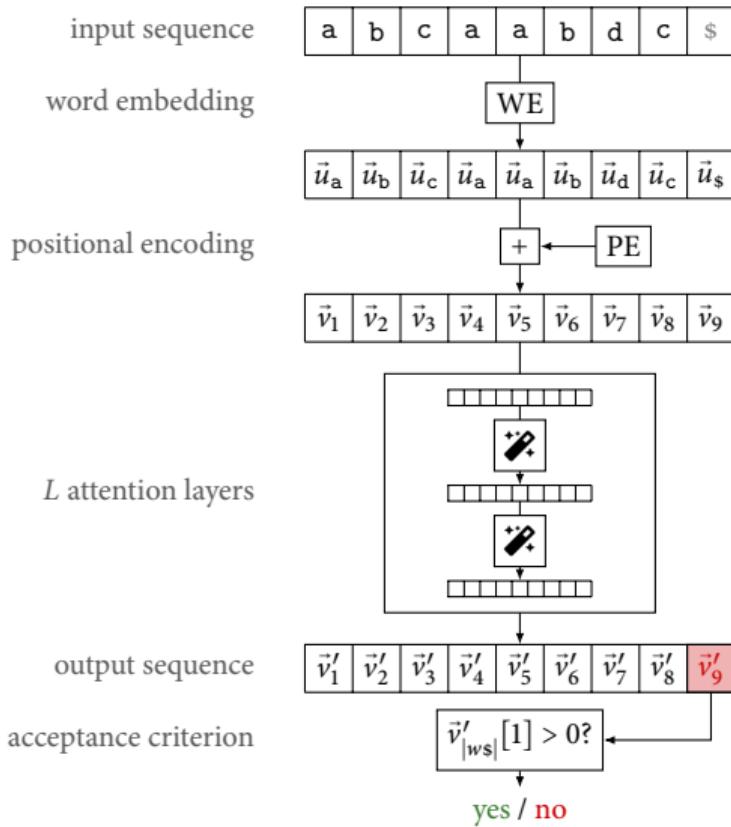
Transformers: Scheme



Transformers: Scheme



Transformers: Scheme



Transformers: Attention Mechanisms

1 Unique Hard Attention (UHA)

	\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 (AHA)

	\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 (SMA)

	\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

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

Transformers: Languages

- accepted language of a transformer \mathfrak{T} : $L(\mathfrak{T}) = \{w \in \Sigma^* \mid \mathfrak{T}(w\$) = \text{yes}\}$.
- **UHAT / AHAT / SMAT**: all languages accepted by a transformer with UHA / AHA / SMA mechanism
- **NoPE-C**: all languages accepted by a C-transformer without positional encoding (i.e., $PE(i, n) = \vec{0}$)
- **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.

Theorem (Yang, Chiang, Angluin @ NeurIPS 2024)

NoPE-UHAT[Masking] = Star-Free

Overview

$$\begin{aligned}\text{NoPE-AHAT}[\leq 1] \\ = \\ \text{QFPA}\end{aligned}$$

$$\begin{aligned}\subsetneq \\ \text{NoPE-AHAT}[\cup] \\ = \\ \text{SemiAlg}\end{aligned}$$

$$\begin{aligned}\subsetneq \\ \text{AHAT} \\ \subsetneq \\ \text{SMAT}\end{aligned}$$

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 .

Overview

NoF $L \subseteq \Sigma^*$ 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

\subsetneq

SMAT

Definition

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

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

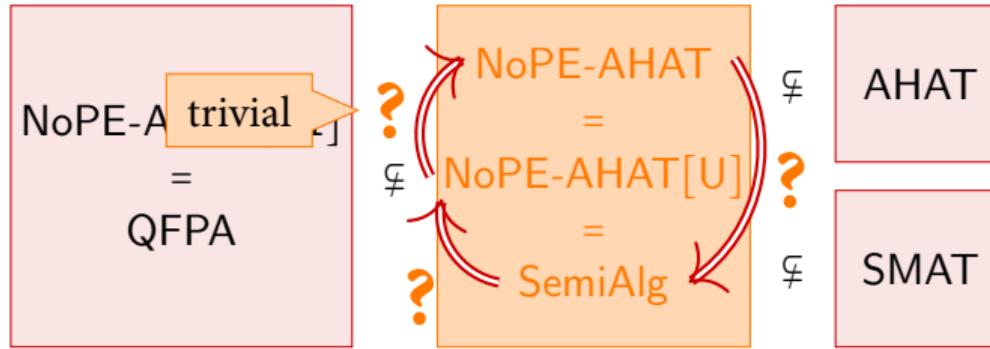
\Updownarrow

L is described by a Presburger formula.

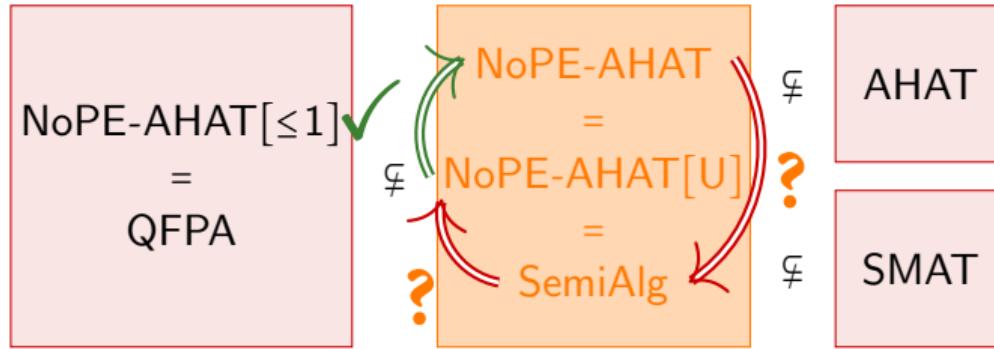
$$\Psi: \Sigma^* \rightarrow \mathbb{N}^d: w \mapsto (|w|_{a_1}, |w|_{a_2}, \dots, |w|_{a_d})$$

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

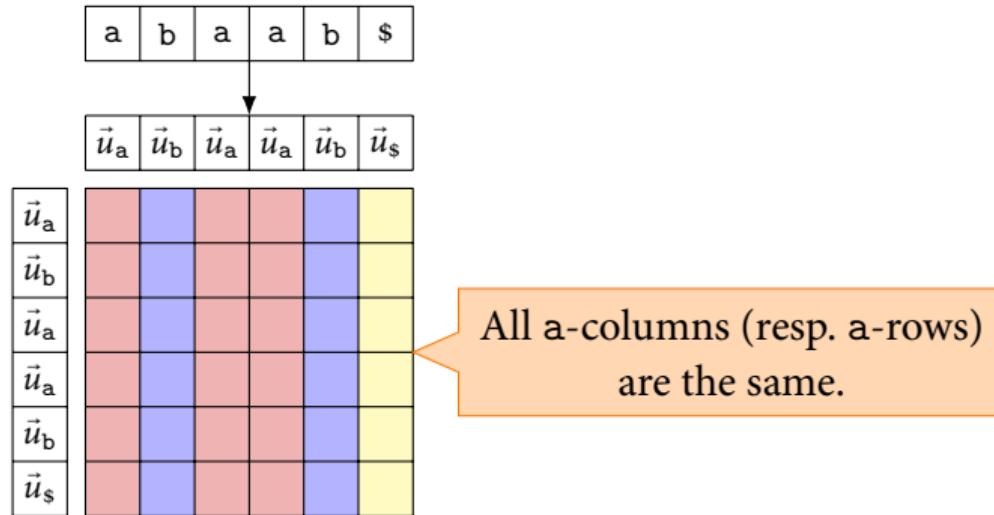
Overview



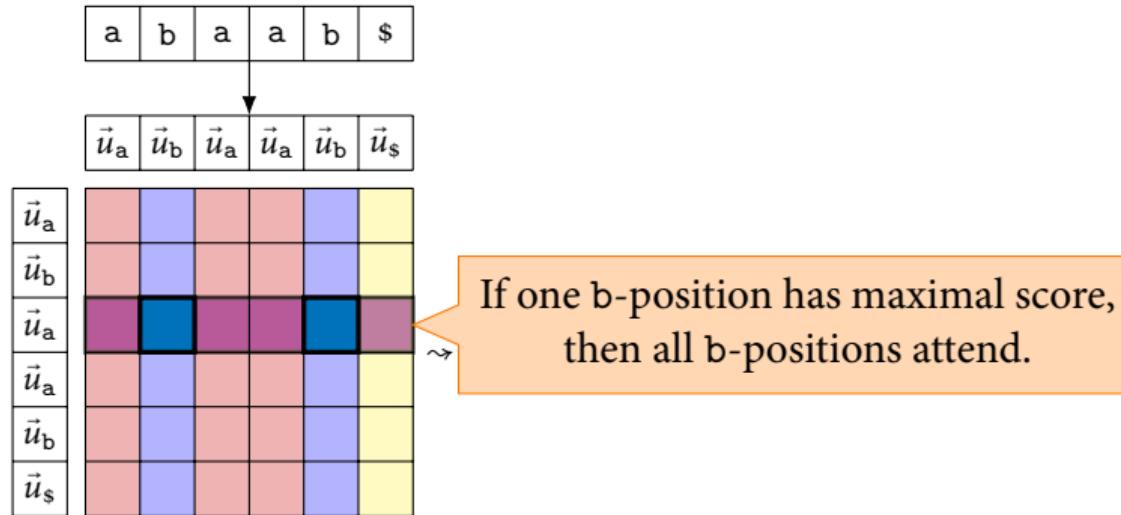
Overview



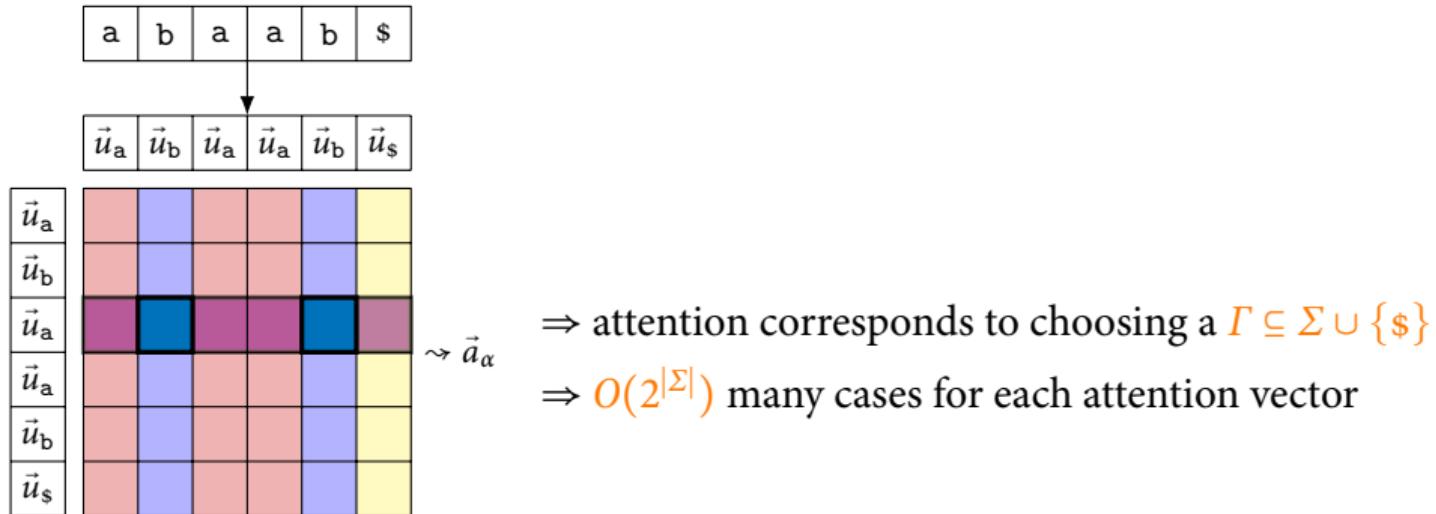
NoPE-AHAT \subseteq SemiAlg (1)



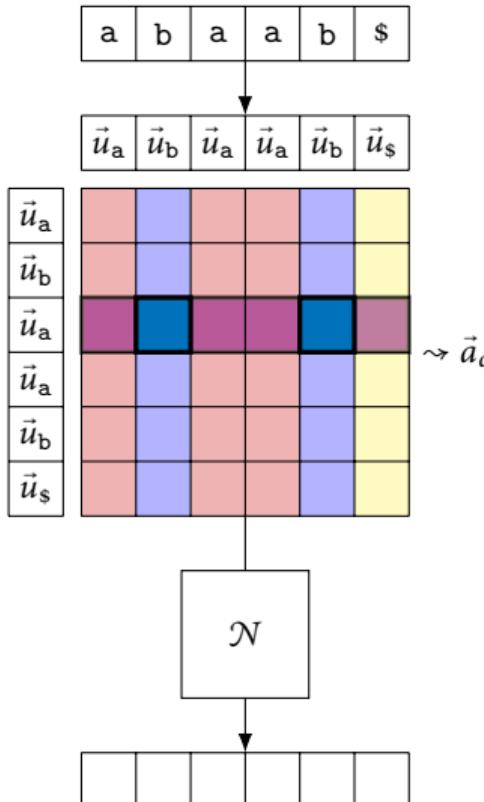
NoPE-AHAT \subseteq SemiAlg (1)



NoPE-AHAT \subseteq SemiAlg (1)



NoPE-AHAT \subseteq SemiAlg (1)



\Rightarrow attention corresponds to choosing a $\Gamma \subseteq \Sigma \cup \{\$\}$

$\Rightarrow O(2^{|\Sigma|})$ many cases for each attention vector

\Rightarrow two cases per letter and ReLU,
 $O(2^{|\Sigma|k})$ many cases if \mathcal{N} contains k ReLUs

$\Rightarrow O(2^{|\Sigma|^2 k})$ many different outcomes

$\Rightarrow O(2^{|\Sigma|^2 k L})$ many different outcomes after L layers

- For each of the $O(2^{|\Sigma|^2 kL})$ many choices construct a conjunction of polynomial inequalities that verify:
 - i the choices maximize the attention scores
 - ii the resulting vector $\vec{v}_\$^{(L)}$ satisfies the accepting condition
- Choice of maximal score can be done via equations of the form

$$\langle K(\vec{v}_\alpha), Q(\vec{v}_\beta) \rangle \geq \langle K(\vec{v}_\gamma), Q(\vec{v}_\beta) \rangle$$

- Attention vectors can be expressed via

$$\vec{a}_\alpha[i] = \frac{\sum_{\beta \in \Gamma} x_\beta \cdot \vec{v}_\beta[i]}{\sum_{\beta \in \Gamma} x_\beta}$$

- Finally, take a disjunction over all $O(2^{|\Sigma|^2 kL})$ conjunctions.

□

- For each of the $O(2^{|\Sigma|^2 kL})$ many choices construct a conjunction of polynomial inequalities that verify:
 - i the choices maximize the attention scores
 - ii the resulting vector $\vec{v}_s^{(L)}$ satisfies the accepting condition
- Choice of maximal score can be done via equations of the form

$$\langle K(\vec{v}_\alpha), Q(\vec{v}_\beta) \rangle \geq \langle K(\vec{v}_\gamma), Q(\vec{v}_\beta) \rangle$$

- Attention vectors K and Q are affine transformations

\Rightarrow polynomials

$$\vec{a}_\alpha[i] = \frac{\rho_{\alpha\Gamma} - \rho_{\alpha\beta}}{\sum_{\beta \in \Gamma} x_\beta}$$

- Finally, take a disjunction over all $O(2^{|\Sigma|^2 kL})$ conjunctions.

□

- For each of the $O(2^{|\Sigma|^2 kL})$ many choices construct a conjunction of polynomial inequalities that verify:
 - i the choices maximize the attention scores
 - ii the resulting vector $\vec{v}_s^{(L)}$ satisfies the accepting condition
- Choice of maximal score can be done via equations of the form

$$\langle K(\vec{v}) \rangle \quad \text{Equivalently:} \\ \vec{v}'_\alpha[i] \cdot \sum_{\beta \in \Gamma} x_\beta = \sum_{\beta \in \Gamma} x_\beta \cdot \vec{v}_\beta[i]$$

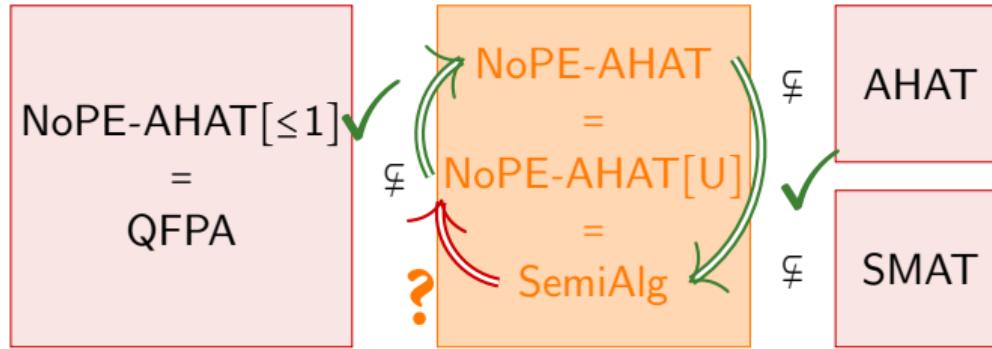
- Attention vectors can be expressed as

$$\vec{a}_\alpha[i] = \frac{\sum_{\beta \in \Gamma} x_\beta \cdot \vec{v}_\beta[i]}{\sum_{\beta \in \Gamma} x_\beta}$$

- Finally, take a disjunction over all $O(2^{|\Sigma|^2 kL})$ conjunctions.

□

Overview



Monomials into AHAT

Lemma

Let $p \in \mathbb{Z}[\vec{x}]$ be a **monomial** of degree L . Then

$$L_p = \{w \in \Sigma^* \mid p(\Psi(w)) > 0\} \in \text{NoPE-AHAT}[\leq L, U].$$

- ① **Word embedding:** map each letter $\alpha \in \Sigma \cup \{\$\}$ to unit vector $\vec{e}_\alpha \in \mathbb{N}^{\Sigma \cup \{\$\}}$
- ② **Compute frequencies:** compute $|\Sigma| + 1$ many new components holding values $\frac{x_\alpha}{|w\$|}$
 - One uniform attention layer \Rightarrow attention vector $\vec{a}_i = \sum_{\alpha \in \Sigma \cup \{\$\}} \frac{|w\$|_\alpha}{|w\$|} \cdot \vec{e}_\alpha = \left(\frac{x_\alpha}{|w\$|} \right)_{\alpha \in \Sigma \cup \{\$\}}$
 - Append the attention vector
- ③ **Multiply:** compute $\frac{x_\alpha}{|w\$|} \cdot y_i$ for some $\alpha \in \Sigma \cup \{\$\}$ and a component y_i
 - Let $u_\alpha \in \{0, 1\}$ be the value in component α .
 - Compute $u_\alpha \cdot y_i$ via a neural network:

$$u_\alpha \cdot y_i = \text{ReLU}(y_i - (1 - u_\alpha)) \in \{0, y_i\}$$

- Accumulate the $u_\alpha \cdot y_i$ to $\frac{x_\alpha \cdot y_i}{|w\$|}$ in a uniform attention layer.
- ④ **Iterate:** compute $\frac{p(\vec{x})}{|w\$|^L}$ by iterated multiplication. □

Monomials into AHAT

Lemma

Let $p \in \mathbb{Z}[\vec{x}]$ be a **monomial** of degree L . Then

$$L_p = \{w \in \Sigma^* \mid p(\Psi(w)) > 0\} \in \text{NoPE-AHAT}[\leq L, U].$$

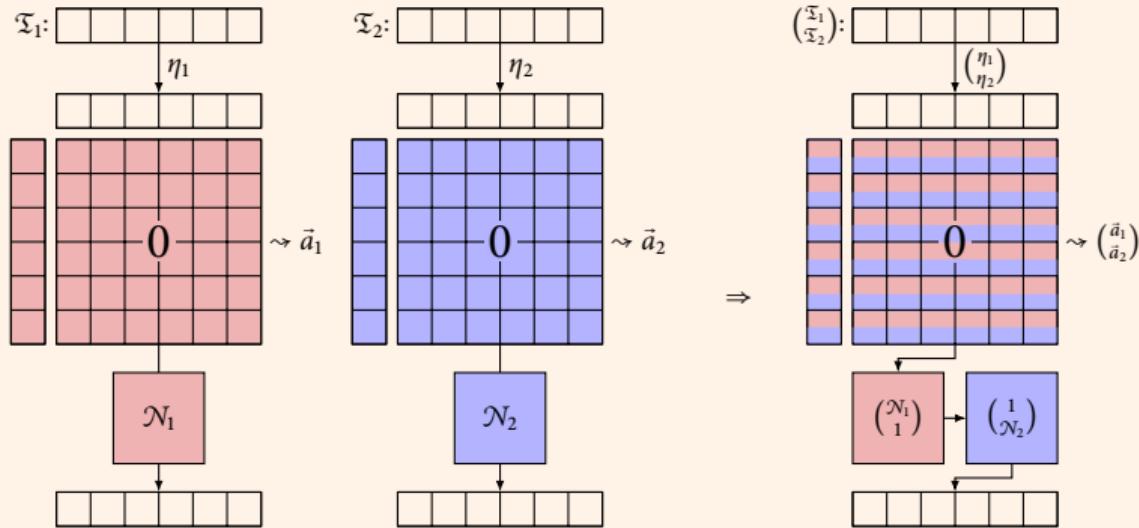
- 0 **Word embedding:** map each letter $\alpha \in \Sigma \cup \{\$\}$ to unit vector $\vec{e}_\alpha \in \mathbb{N}^{\Sigma \cup \{\$\}}$
- 1 **Compute frequencies:** compute $|\Sigma| + 1$ many new components holding values $\frac{x_\alpha}{|w\$|}$
 - One uniform attention layer \Rightarrow attention vector $\vec{a}_i = \sum_{\alpha \in \Sigma \cup \{\$\}} \frac{|w\$|_\alpha}{|w\$|} \cdot \vec{e}_\alpha = \left(\frac{x_\alpha}{|w\$|} \right)_{\alpha \in \Sigma \cup \{\$\}}$
 - Append the attention vector
- 2 **Multiply:** compute $\frac{x_\alpha}{|w\$|} \cdot y_i$ for some $\alpha \in \Sigma \cup \{\$\}$ and a component y_i
 - Let $u_\alpha \in \{0, 1\}$ be the value in component α .
 - Compute $u_\alpha \cdot y_i$ via a neural network:

$$u_\alpha \cdot y_i = \text{ReLU}(y_i - (1 - u_\alpha)) \in \{0, y_i\}$$

- Accumulate the $\frac{x_\alpha \cdot y_i}{|w\$|}$ in the uniform attention layer.
- 3 **Iterate:** compute $\frac{p(\vec{x})}{|w\$|^L} \in \{0, 1\}$ and $\frac{x_\alpha \cdot y_i}{|w\$|} \in [0, 1]$ iteration. □

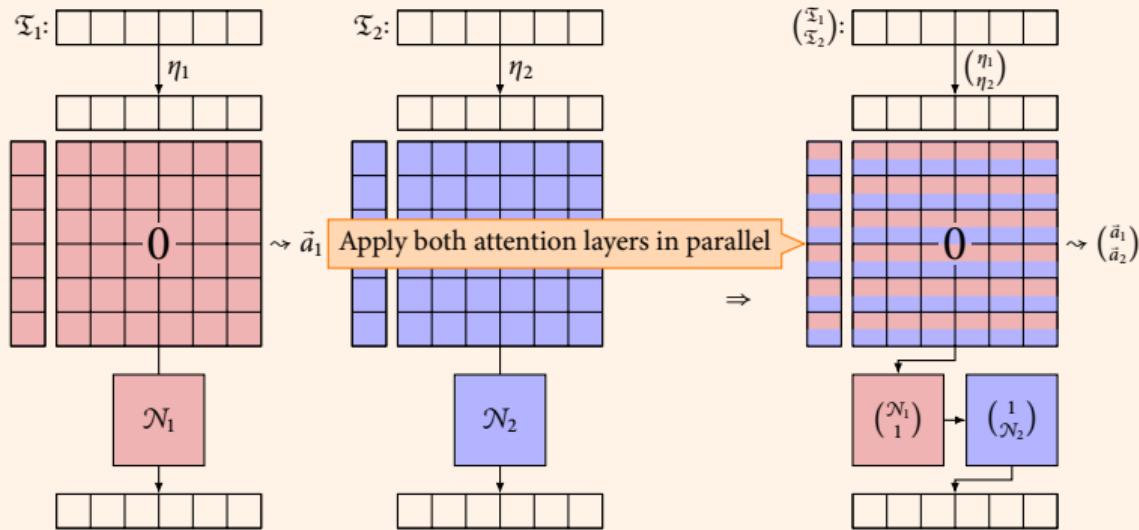
Parallelize Uniform NoPE-AHATs

Observation



Parallelize Uniform NoPE-AHATs

Observation



Polynomials into AHAT

Lemma

Let $p \in \mathbb{Z}[\vec{x}]$ be a **polynomial** of degree L . Then

$$L_p = \{w \in \Sigma^* \mid p(\Psi(w)) > 0\} \in \text{NoPE-AHAT}[\leq L, U].$$

- We may assume that p is **homogeneous**, i.e. all monomials of p have degree L :

$$p'(x_0, \vec{x}) := x_0^L \cdot p\left(\frac{x_1}{x_0}, \dots, \frac{x_m}{x_0}\right)$$

- $p(\vec{x}) > 0 \iff p'(1, \vec{x}) > 0$
- We will represent x_0 via the end marker \$.
- Construct AHATs for all monomials, parallelize them, and compute the sum in the neural network in the last layer. □

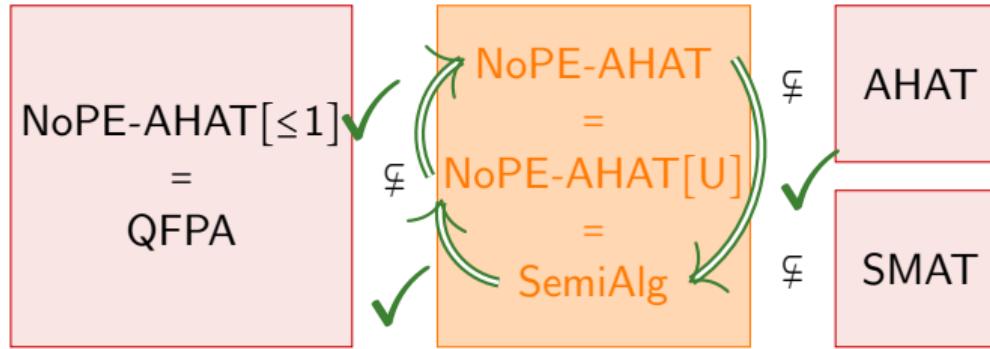
Proposition

$\text{SemiAlg} \subseteq \text{NoPE-AHAT}[U]$

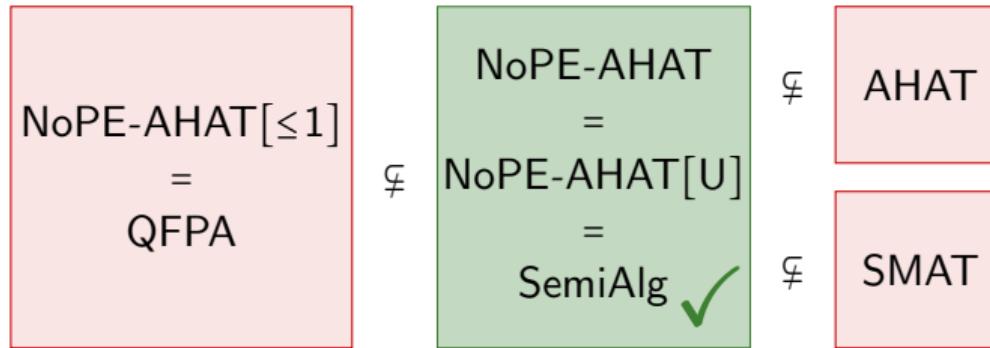
- Let $L \in \text{SemiAlg}$ and $p_1, \dots, p_k \in \mathbb{Z}[\vec{x}]$ be polynomials with $L = L_{p_1} \cup \dots \cup L_{p_k}$.
- Compute AHATs for all L_{p_i} , parallelize them, and compute the maximum of their outputs in the neural network in the last layer.
 - Note: $\max\{x, y\} = \max\{0, x - y\} + y = \text{ReLU}(x - y) + y$

□

Overview



Overview



Related Results

Corollary

- 1 $\text{PARITY} = \{w \in \{a, b\}^*: |w|_a \text{ is even}\} \notin \text{NoPE-AHAT}$
- 2 $\text{Proj}(\text{NoPE-AHAT}) = \text{RE} \cap \text{PI}$

Theorem

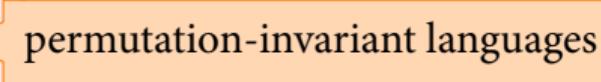
$$\text{Proj}(\text{NoPE-AHAT}[\leq 2, U]) = \text{RE} \cap \text{PI}$$

Corollary

$\text{NoPE-AHAT}[\leq 2, U]$ can express languages not recognized by higher-order recursion schemes (HORS), Petri nets, simplified multi-counter machines, or LTL with counting.

Related Results

Corollary

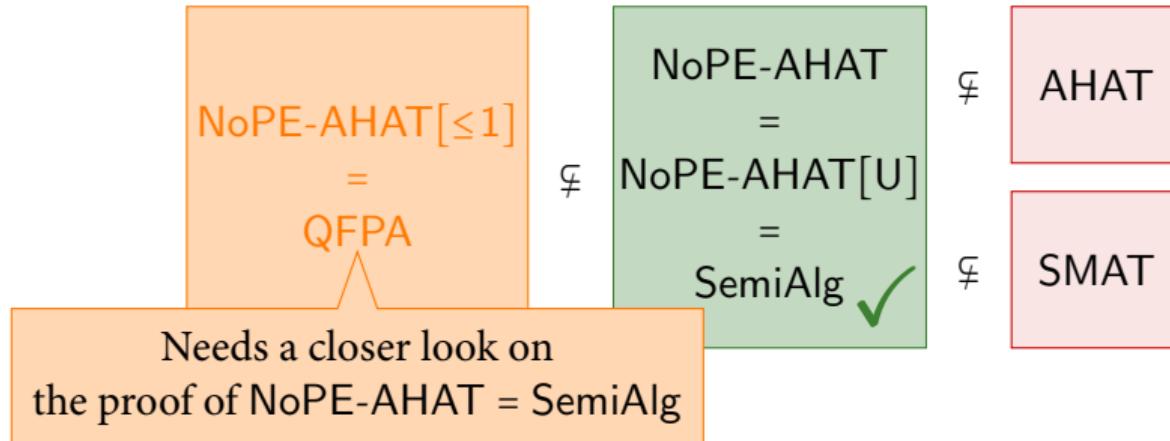
- ① PARITY  recursively enumerable languages
- ② $\text{Proj}(\text{NoPE-AHAT}) = \text{RE} \cap \text{PI}$  permutation-invariant languages

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

Corollary

NoPE-AHAT $[\leq 2, U]$ can express languages not recognized by higher-order recursion schemes (HORS), Petri nets, simplified multi-counter machines, or LTL with counting.

Overview



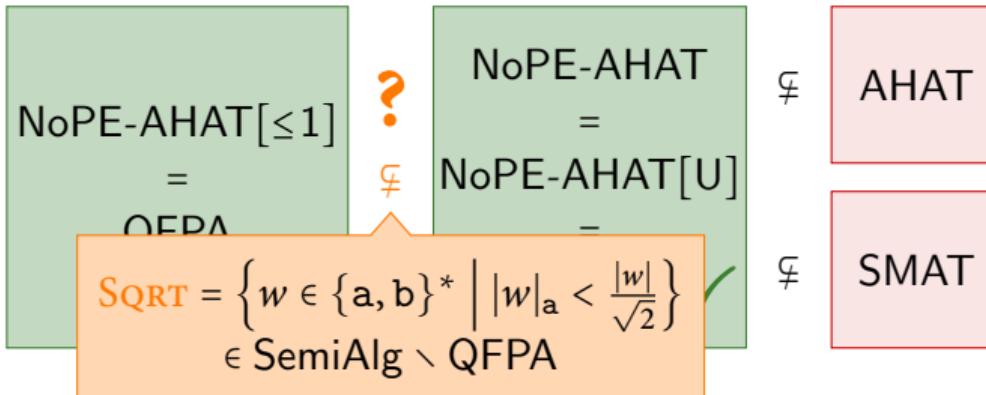
Overview

NoPE-AHAT[≤ 1]
= QFPA

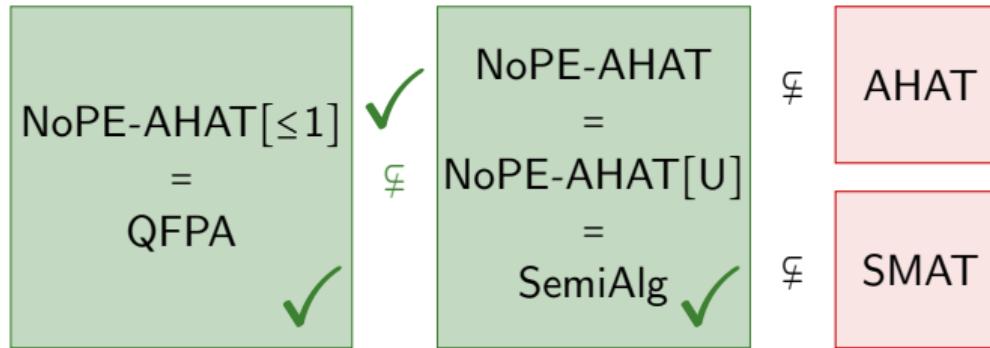
\subsetneq NoPE-AHAT
= NoPE-AHAT[U]
= SemiAlg

\subsetneq AHAT
 \subsetneq SMAT

Overview



Overview



Emptiness Problem

Corollary

The emptiness problem is

- 1 *decidable for NoPE-AHAT[≤1].*
- 2 *undecidable for NoPE-AHAT[≤2].*

Theorem

The following two problems are inter-reducible:

- 1 *The emptiness problem for AHAT without positional encoding and **without end marker***
- 2 *The solvability problem of Diophantine equations over rationals*

Emptiness Problem

Corollary Since Presburger arithmetic is decidable
[Presburger 1929]

The emptiness problem

- 1 decidable for NoPE-AHAT [≤ 1].
- 2 undecidable for NoPE-AHAT [≤ 2].

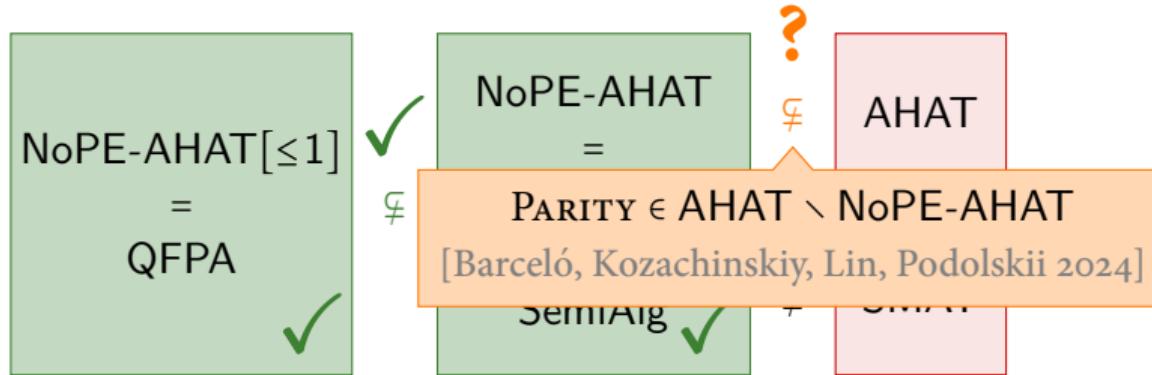
Theorem Since $\text{Proj}(\text{NoPE-AHAT}[\leq 2]) = \text{RE} \cap \text{PI}$

The following two problems are inter-reducible:

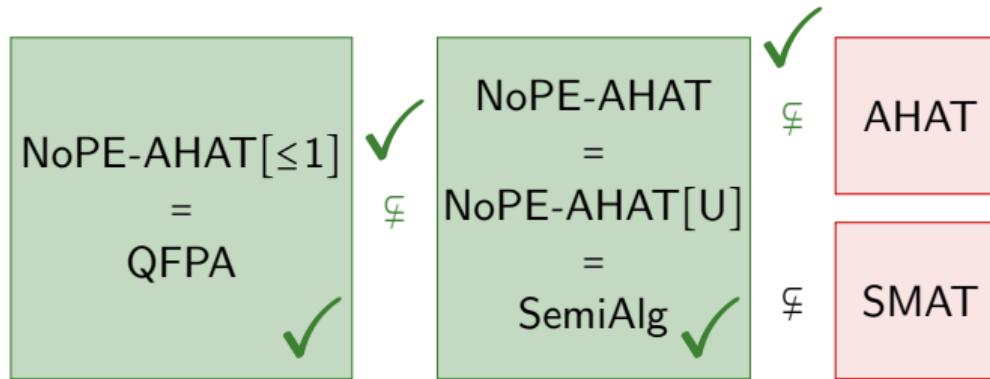
- 1 The emptiness problem for AHAT without positional encoding and *without end marker*
- 2 The solvability problem of Diophantine equations over rationals

(Un-)decidability of this problem is still open!

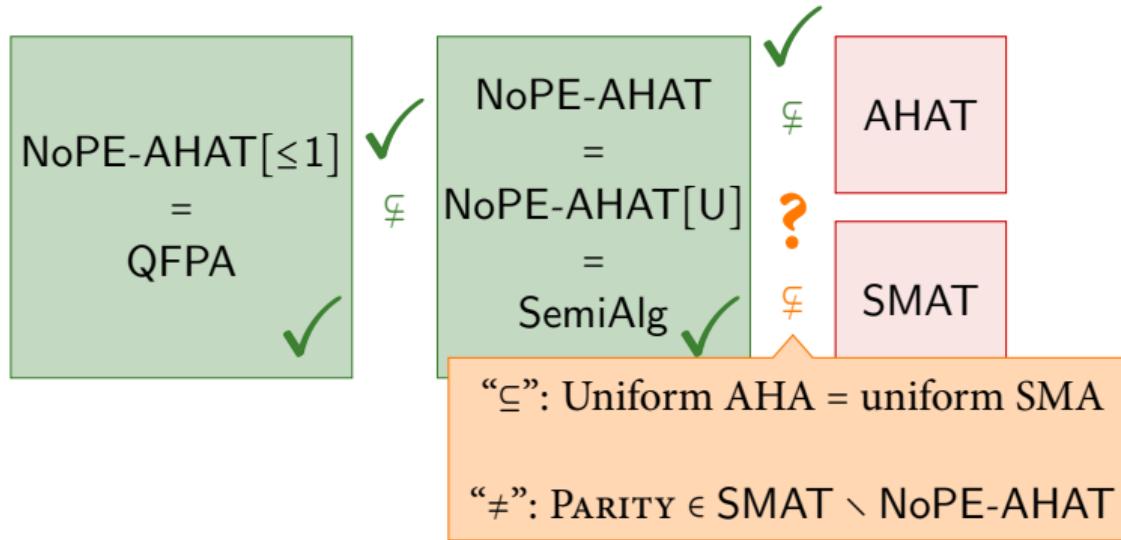
Overview



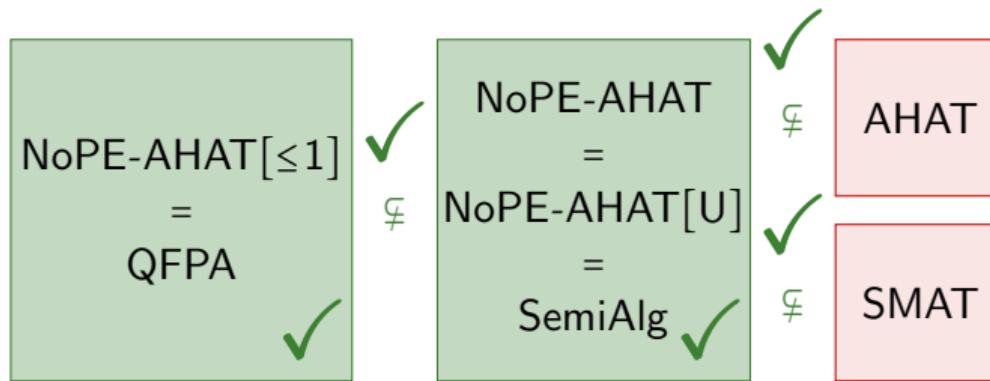
Overview



Overview



Overview



■ Open Problems:

- Do AHAT[U] = AHAT or NoPE-SMAT[U] = NoPE-SMAT hold?
- Is NoPE-AHAT[$\leq L, U$] equal to semi-algebraic sets with polynomials of degree $\leq L$?
- Are intersection, separability, ... decidable?

Thank you!

