# Neural Sequence Models

Hi, how are you? → Salut, comment ça va?

Enjoy your meal! → Bon appétit!

That was great! → POSITIVE

Can you [MASK] me [MASK] salt? → Can you pass me the salt?

...

# Neural Sequence Models

Hi, how are you? → ❓ → Salut, comment ça va?

Enjoy your meal! → Bon appétit!

That was great! → ❓ → POSITIVE

Can you [MASK] me [MASK] salt? → ❓ → Can you pass me the salt?

...

# Neural Sequence Models

Hi, how are you? → Salut, comment ça va?

Enjoy your meal! → Bon appétit!

**Understanding the Black Box**

That was great! → POSITIVE

Can you [MASK] me [MASK] salt? → Can you pass me the salt?

...

# Neural Sequence Models

**Hi, how are you?** → ⬛ **?** → **Salut, comment ça va?**

**Enjoy your meal!** → **Bon appétit!**

## Understanding the Black Box

- Reliability (Verifiability)
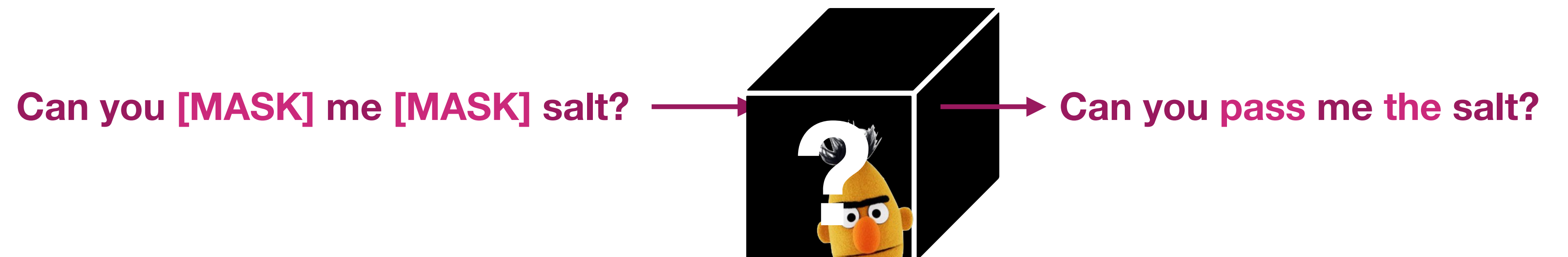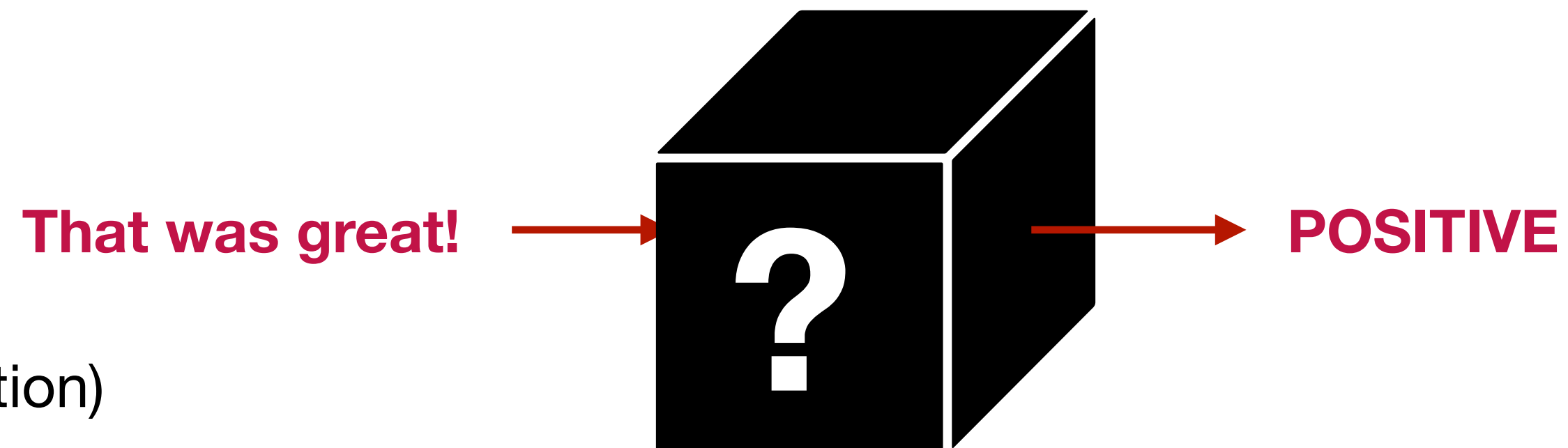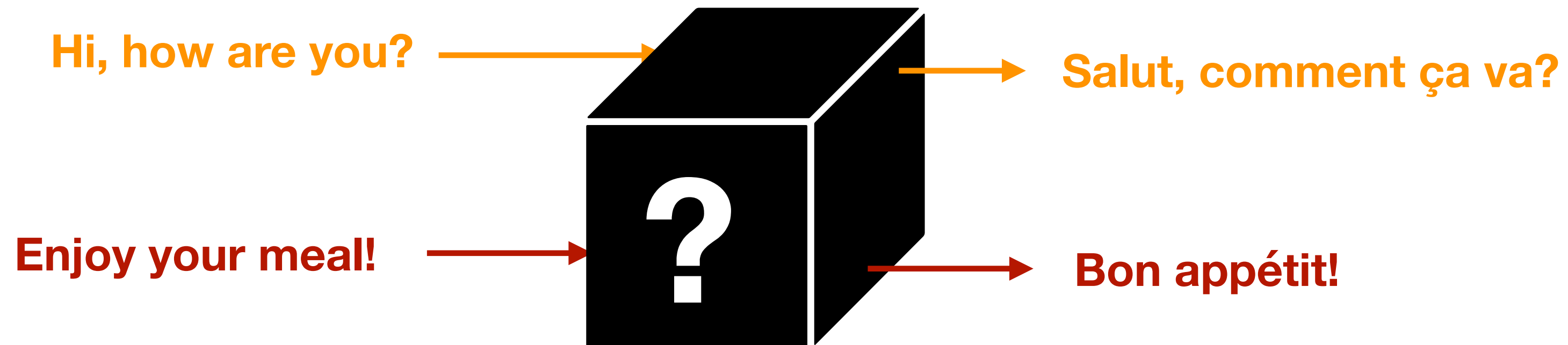- Intuition (biases, model selection)
- Knowledge extraction

**That was great!** → ⬛ **?** → **POSITIVE**

**Can you [MASK] me [MASK] salt?** → ⬛ **?** → **Can you pass me the salt?**

...

# Neural Sequence Models

Hi, how are you? → Salut, comment ça va?

Enjoy your meal! → **?** → Bon appétit!

## Understanding the Black Box

That was great! → **?** → POSITIVE

- Reliability (Verifiability)
- Intuition (biases, model selection)
- Knowledge extraction

- Model design

Can you [MASK] me [MASK] salt? → **?** → Can you pass me the salt?

...

# Neural Sequence Models

**Hi, how are you?** → ■ → **Salut, comment ça va?**

**Enjoy your meal!** → ? → **Bon appétit!**

## Understanding the Black Box

**That was great!** → ? → **POSITIVE**

- Reliability (Verifiability)
- Intuition (biases, model selection)
- Knowledge extraction

- Model design

**Can you [MASK] me [MASK] salt?** → ? → **Can you pass me the salt?**

- **Just kinda cool**

…

# Neural Sequence Models

Hi, how are you? → ? → Salut, comment ça va?

Enjoy your meal! → ? → Bon appétit!

## Understanding the Black Box

That was great! → ? → POSITIVE

- Reliability (Verifiability)
- Intuition (biases, model selection)
- Knowledge extraction

- Model design

Can you [MASK] me [MASK] salt? → ? → Can you pass me the salt?

- **Just kinda cool**

**Natural language is hard…**

…

# Neural Sequence Models:
# A Formal Lens



()(()))(  →  **PTPPPTFF**

)()()  →  **FFFFF**

**Understanding the Black Box**

- Reliability (Verifiability)
- Intuition (biases, model selection)
- Knowledge extraction

- Model design

- **Just kinda cool**
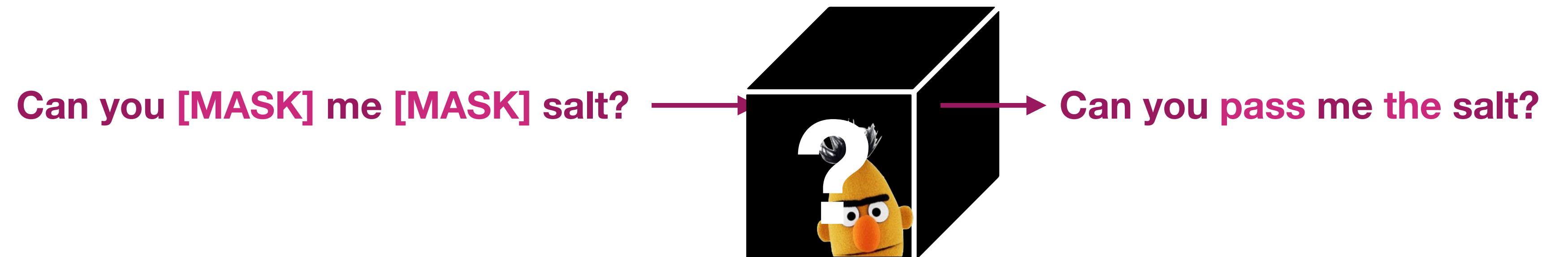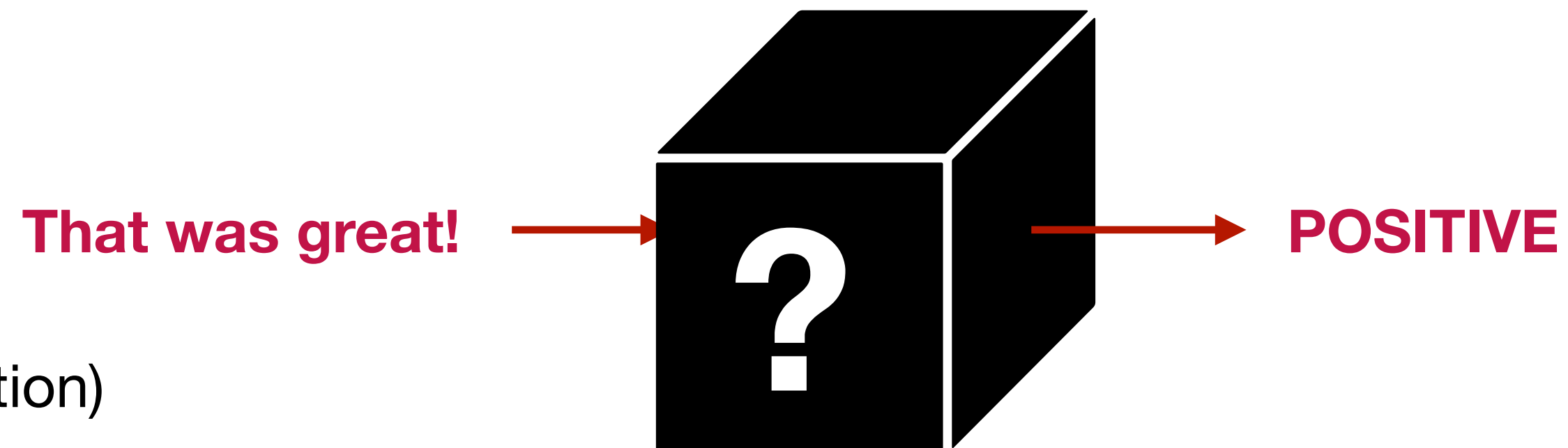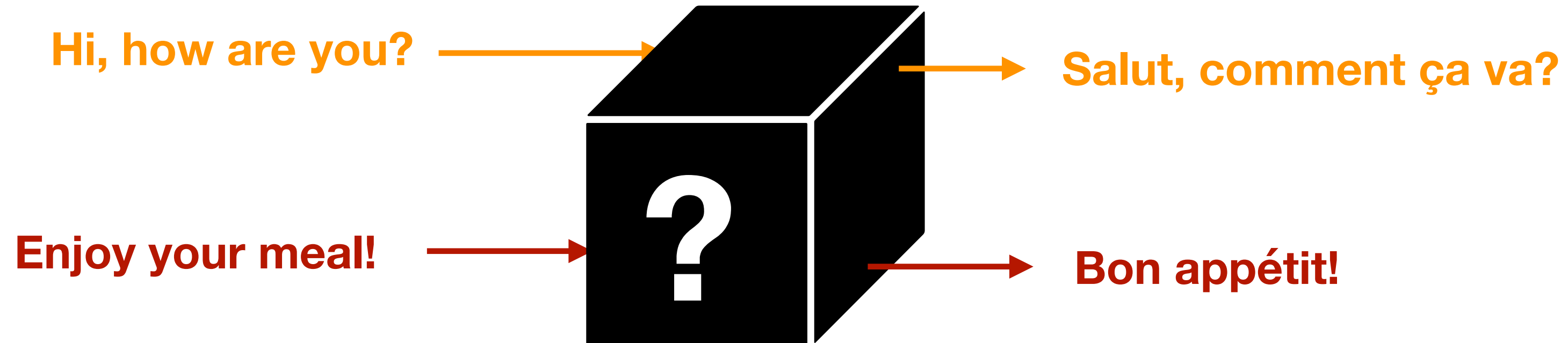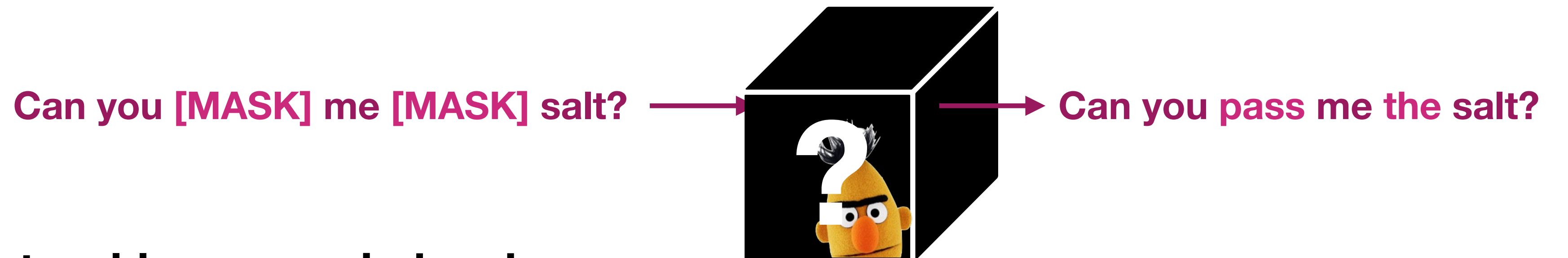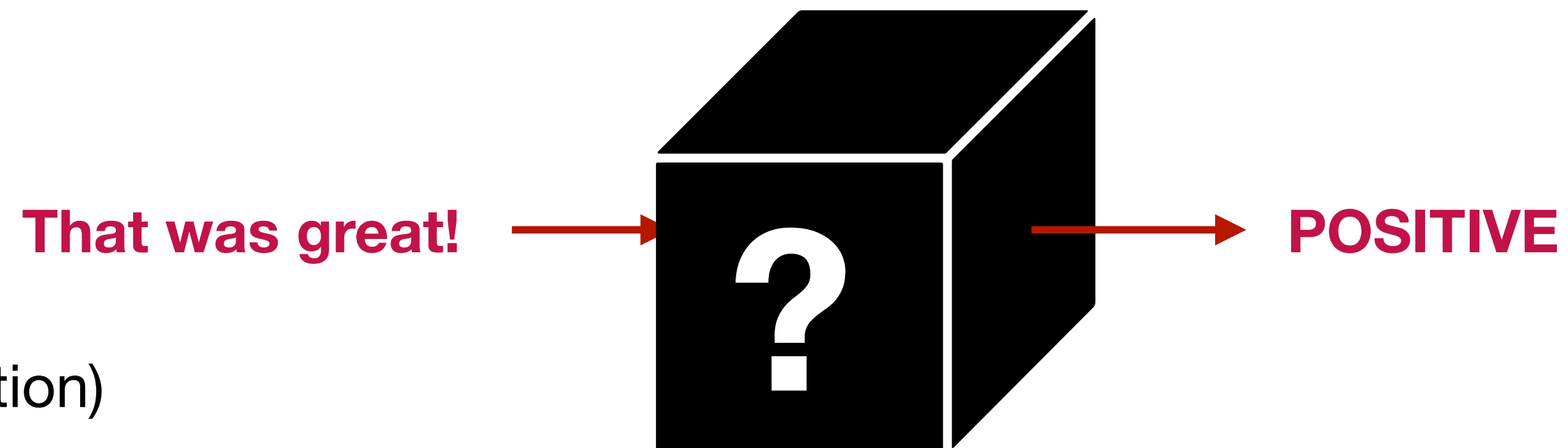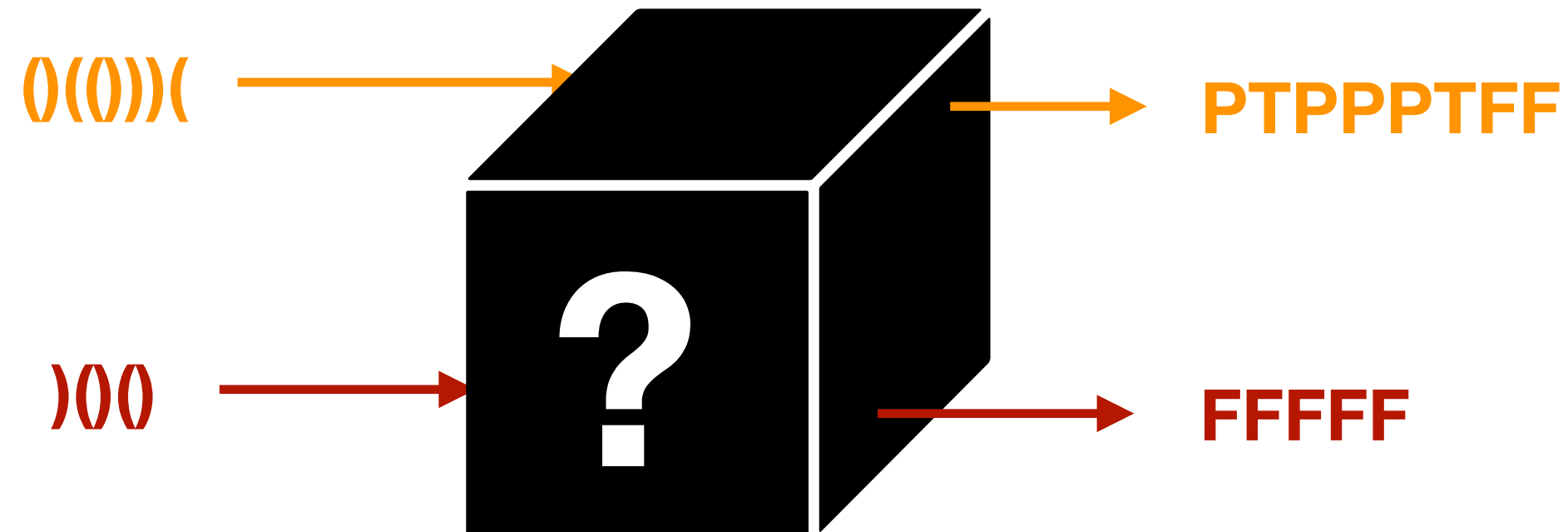
**aaabbb**  →  **Accept**

**abcdefg**  →  **gfedcba**

...

# Neural Sequence Models:
# a Formal Lens

**Counting**
LSTMs are counter machines, GRUs aren't **(ACL 2018)**

**RASP**
Finding a formalism to describe transformers **(ICML 2021)**

**DFAs from RNNs**
Applying L* to learn DFAs from RNNs **(ICML 2018)**

    + using the result for CFGs **(TACAS 2021)**

# Neural Sequence Models:
# a Formal Lens

## Counting
LSTMs are counter machines, GRUs aren't **(ACL 2018)**

## RASP
Finding a formalism to describe transformers **(ICML 2021)**

## DFAs from RNNs
Applying L* to learn DFAs from RNNs **(ICML 2018)**
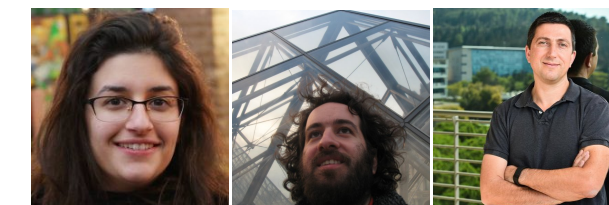
+ using the result for CFGs **(TACAS 2021)**

# RNNs

(Elman, **1990**)
Introduction of RNNs



**General RNN concept:** $\qquad h_t = f(x_t, h_{t-1})$

**Elman RNN:** $\qquad h_t = \sigma_h(W_h x_t + U_h h_{t-1} + b_h)$

On the Practical Computational Power of Finite Precision RNNs for Language Recognition (Weiss, Goldberg, Yahav, ACL 2018)

# RNNs

(Elman, **1990**)
Introduction of RNNs

General RNN concept:

$$h_t = f(x_t, h_{t-1})$$

Elman RNN:

$$h_t = \sigma_h(W_h x_t + U_h h_{t-1} + b_h)$$

On the Practical Computational Power of Finite Precision RNNs for Language Recognition (Weiss, Goldberg, Yahav, ACL 2018)

# RNNs

⭐ (Elman, **1990**)
Introduction of RNNs

⭐ (Siegelmann and Sonntag, **1993**)
RNNs are Turing Complete

**Theoretical Power**



| 1 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | ... |

On the Practical Computational Power of Finite Precision RNNs for Language Recognition (Weiss, Goldberg, Yahav, ACL 2018)

# RNNs

★ (Elman, **1990**)
Introduction of RNNs

★ (Siegelmann and Sonntag, **1993**)
RNNs are Turing Complete

**Theoretical Power**

★ (Hochreiter and Schmidhuber, **1997**)
LSTMs

**Practical Modifications**

★ (Cho et al, **2014**)
GRUs

On the Practical Computational Power of Finite Precision RNNs for Language Recognition (Weiss, Goldberg, Yahav, ACL 2018)

# RNNs

(Elman, **1990**)
Introduction of RNNs

(Siegelmann and Sonntag, **1993**)
RNNs are Turing Complete

**Theoretical Power**

(Hochreiter and Schmidhuber, **1997**)
LSTMs

**Practical
Modifications**

(Cho et al, **2014**)
GRUs

On the Practical Computational Power of Finite Precision RNNs for Language Recognition (Weiss, Goldberg, Yahav, ACL 2018)

# RNNs

(Elman, **1990**)
Introduction of RNNs

(Siegelmann and Sonntag, **1993**)
RNNs are Turing Complete

**Theoretical Power**

**RNN Turing Completeness Proof (1993):**

1. **Requires Infinite Precision:**
   **Uses stack(s), with zeros pushed using division:** $g = g/4 + 1/4$
   *In 32 bits, this reaches the limit after **15** pushes*

2. **Requires Infinite Time:**
   And specifically, allows processing beyond reading input
   **(Non standard use case!)**

# RNNs

⭐ (Elman, **1990**)
Introduction of RNNs

⭐ (Siegelmann and Sonntag, **1993**)
RNNs are Turing Complete

**Theoretical Power**

⭐ (Hochreiter and Schmidhuber, **1997**)
LSTMs

**?** }

**Practical
Modifications**

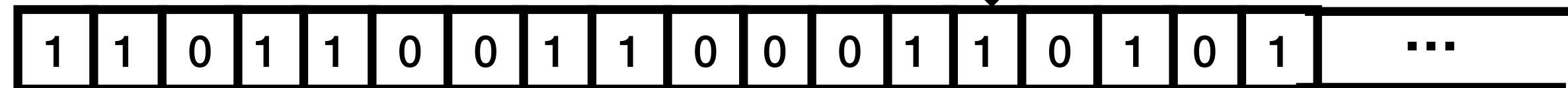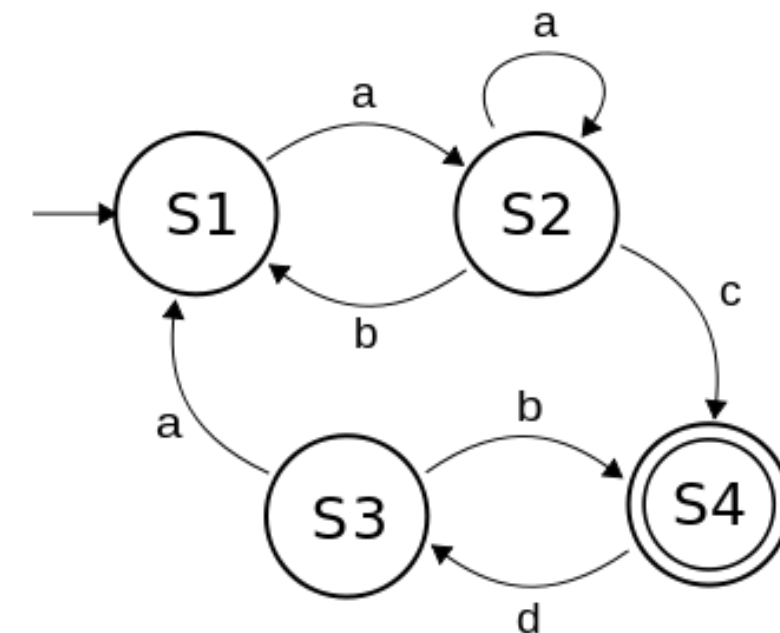⭐ (Cho et al, **2014**)
GRUs

On the Practical Computational Power of Finite Precision RNNs for Language Recognition (Weiss, Goldberg, Yahav, ACL 2018)

$$h_t = f(x_t, h_{t-1})$$

# Practical RNNs

## GRU

$$z_t = \sigma(W^z x_t + U^z h_{t-1} + b^z)$$
$$r_t = \sigma(W^r x_t + U^r h_{t-1} + b^r)$$
$$\tilde{h}_t = \tanh(W^h x_t + U^h(r_t \circ h_{t-1}) + b^h)$$
$$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$$

## LSTM

$$f_t = \sigma(W^f x_t + U^f h_{t-1} + b^f)$$
$$i_t = \sigma(W^i x_t + U^i h_{t-1} + b^i)$$
$$o_t = \sigma(W^o x_t + U^o h_{t-1} + b^o)$$
$$\tilde{c}_t = \tanh(W^c x_t + U^c h_{t-1} + b^c)$$
$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$
$$h_t = o_t \circ g(c_t)$$

On the Practical Computational Power of Finite Precision RNNs for Language Recognition (Weiss, Goldberg, Yahav, ACL 2018)

# Practical RNNs

**GRU**

$$z_t = \sigma(W^z x_t + U^z h_{t-1} + b^z)$$
$$r_t = \sigma(W^r x_t + U^r h_{t-1} + b^r)$$
$$\tilde{h}_t = \tanh(W^h x_t + U^h(r_t \circ h_{t-1}) + b^h)$$
$$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$$

**gates**

**candidate vectors**

**update functions**

**LSTM**

$$f_t = \sigma(W^f x_t + U^f h_{t-1} + b^f)$$
$$i_t = \sigma(W^i x_t + U^i h_{t-1} + b^i)$$
$$o_t = \sigma(W^o x_t + U^o h_{t-1} + b^o)$$
$$\tilde{c}_t = \tanh(W^c x_t + U^c h_{t-1} + b^c)$$
$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$
$$h_t = o_t \circ g(c_t)$$

On the Practical Computational Power of Finite Precision RNNs for Language Recognition (Weiss, Goldberg, Yahav, ACL 2018)

# Practical RNNs

**GRU**

**LSTM**

$z_t \in (0,1)$

$r_t \in (0,1)$

**gates**

$f_t \in (0,1)$

$i_t \in (0,1)$

$o_t \in (0,1)$

$\tilde{h}_t = \tanh(W^h x_t + U^h(r_t \circ h_{t-1}) + b^h)$

$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$

$\tilde{c}_t = \tanh(W^c x_t + U^c h_{t-1} + b^c)$

$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$

**candidate vectors**

$h_t = o_t \circ g(c_t)$

**update functions**

# Practical RNNs

**GRU**

**LSTM**

$z_t \in (0,1)$

$r_t \in (0,1)$

**gates**

$f_t \in (0,1)$

$i_t \in (0,1)$

$o_t \in (0,1)$

$\tilde{h}_t \in (-1,1)$

$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$

$\tilde{c}_t \in (-1,1)$

$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$

$h_t = o_t \circ g(c_t)$

**candidate vectors**

**update functions**

On the Practical Computational Power of Finite Precision RNNs for Language Recognition (Weiss, Goldberg, Yahav, ACL 2018)

# Practical RNNs

**GRU**

$$z_t \in (0,1)$$
$$r_t \in (0,1)$$

**Bounded!**

$$\tilde{h}_t \in (-1,1)$$

$$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$$

**Interpolation**

**LSTM**

$$f_t \in (0,1)$$
$$i_t \in (0,1)$$
$$o_t \in (0,1)$$
$$\tilde{c}_t \in (-1,1)$$
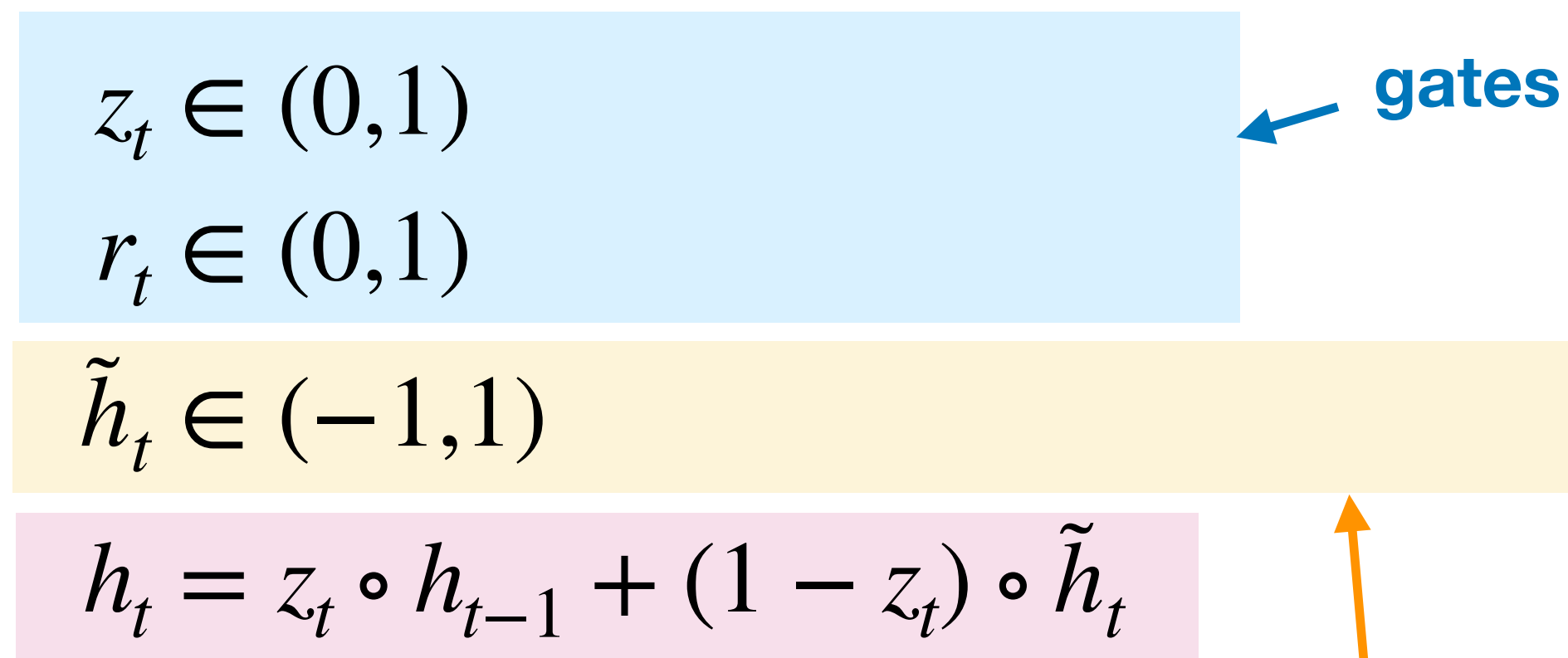$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$
$$h_t = o_t \circ g(c_t)$$

23

On the Practical Computational Power of Finite Precision RNNs for Language Recognition (Weiss, Goldberg, Yahav, ACL 2018)

# Practical RNNs

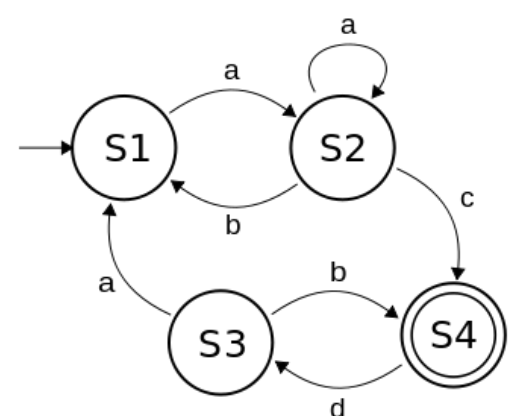## GRU

$$z_t \in (0,1)$$

$$r_t \in (0,1)$$

$$\tilde{h}_t \in (-1,1)$$

$$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$$

## LSTM

$$f_t \in (0,1) \quad \textbf{reset/keep, then -}$$

$$i_t \in (0,1) \quad \textbf{stay/step, by -}$$

$$o_t \in (0,1)$$

$$\tilde{c}_t \in (-1,1) \quad \textbf{subtract/add}$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

$$h_t = o_t \circ g(c_t)$$

On the Practical Computational Power of Finite Precision RNNs for Language Recognition (Weiss, Goldberg, Yahav, ACL 2018)
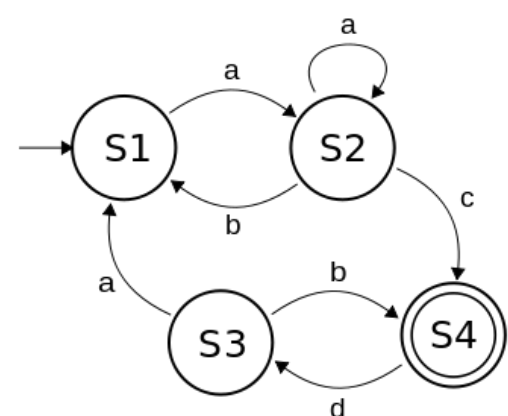
# Practical RNNs

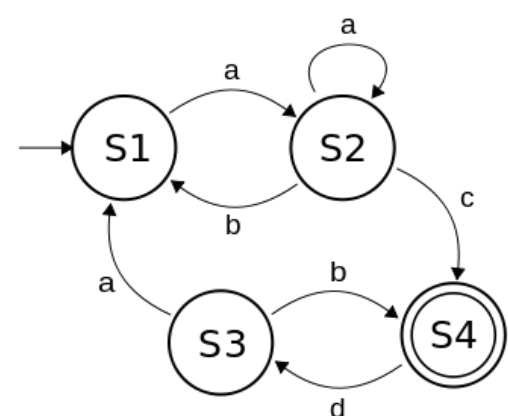## GRU

$$z_t \in (0,1)$$

$$r_t \in (0,1)$$

$$\tilde{h}_t \in (-1,1)$$

$$h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$$



## LSTM

$$f_t \in (0,1)$$  **reset/keep, then -**

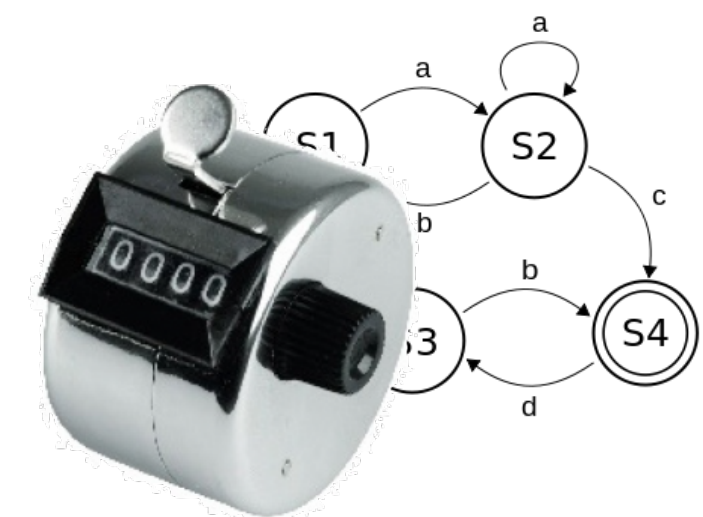$$i_t \in (0,1)$$  **stay/step, by -**

$$o_t \in (0,1)$$

$$\tilde{c}_t \in (-1,1)$$  **subtract/add**

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

$$h_t = o_t \circ g(c_t)$$

On the Practical Computational Power of Finite Precision RNNs for Language Recognition (Weiss, Goldberg, Yahav, ACL 2018)

# Counting



$a^n b^n$

$a^n b^n c^n$

?

**Regular Languages (RL)**

*Palindromes* ✗

**Context Free Languages (CFL)**

**Context Sensitive Languages (CSL)**

**Recursively Enumerable Languages (RE)**

# Practical RNNs

**GRU**

**LSTM**



**Activations on** $a^{1000}b^{1000}$

**Trained** $a^n b^n$ **,** (on positive examples up to length 100)

GRU begins failing at length 39

On the Practical Computational Power of Finite Precision RNNs for Language Recognition (Weiss, Goldberg, Yahav, ACL 2018)
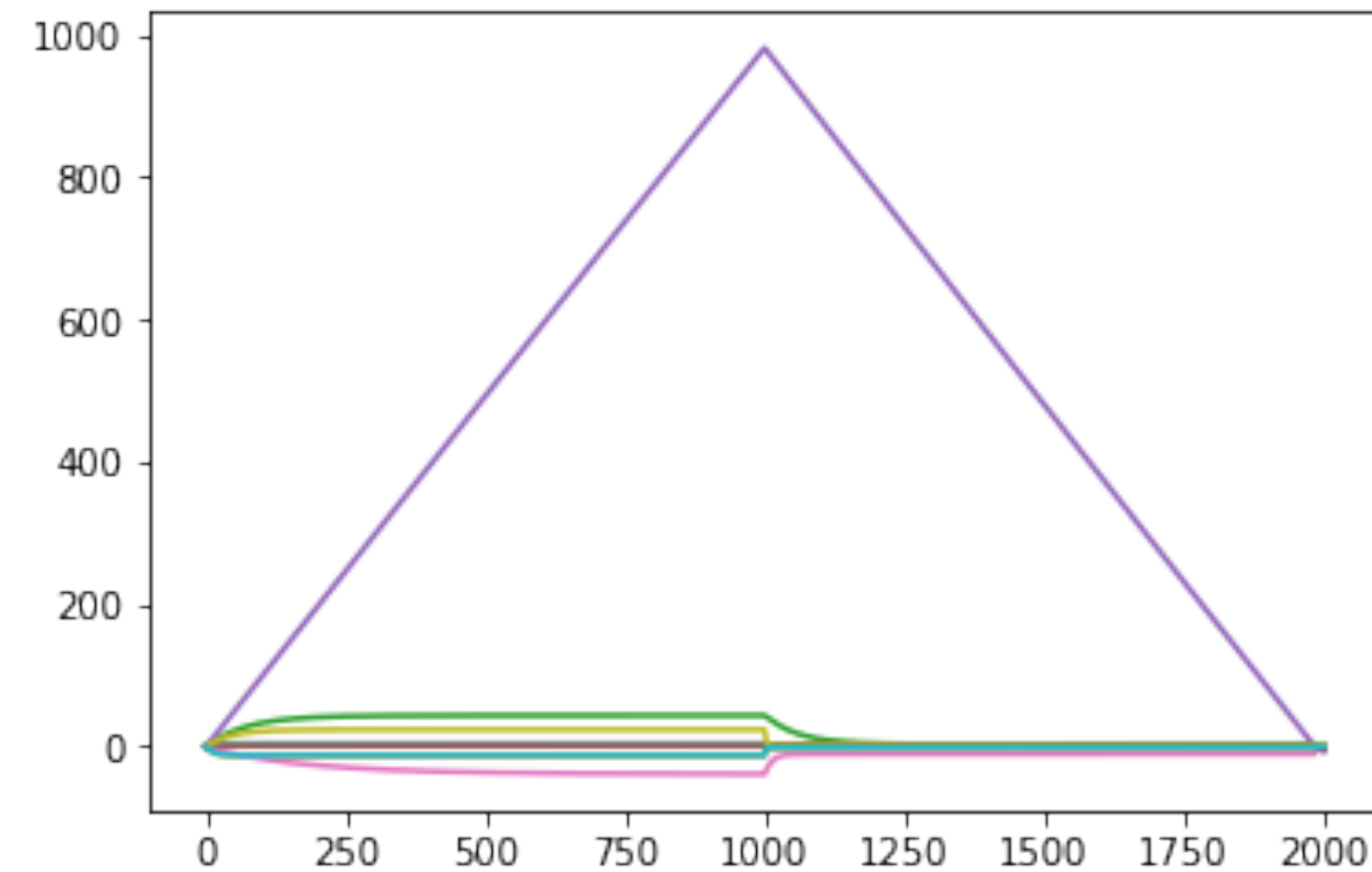
# Practical RNNs

**GRU**

**LSTM**



**Activations on** $a^{100}b^{100}c^{100}$

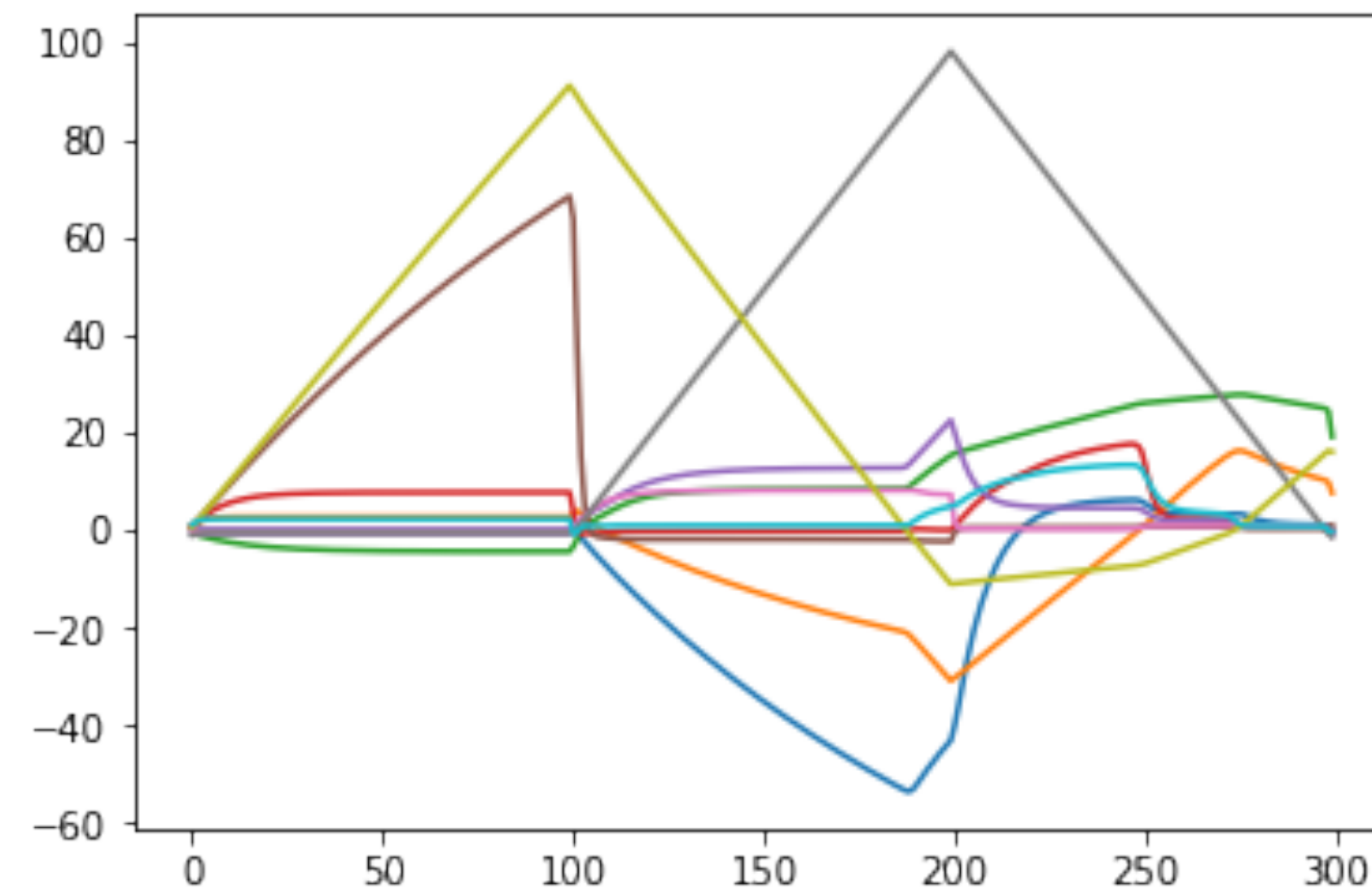**Trained** $a^n b^n c^n$,    (on positive examples up to length 100)

GRU begins failing at length 9

# Neural Sequence Models:
# a Formal Lens

## Counting
LSTMs are counter machines, GRUs aren't **(ACL 2018)**

## RASP
Finding a formalism to describe transformers **(ICML 2021)**

## DFAs from RNNs
Applying L* to learn DFAs from RNNs **(ICML 2018)**

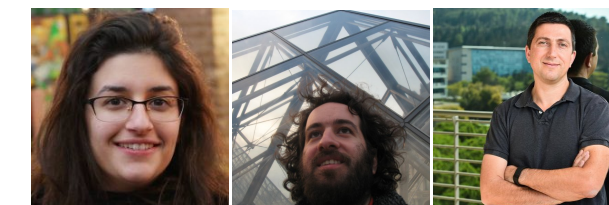  + using the result for CFGs **(TACAS 2021)**

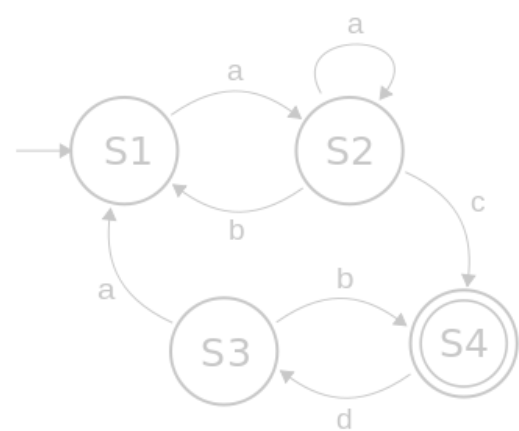# Neural Sequence Models:
# a Formal Lens
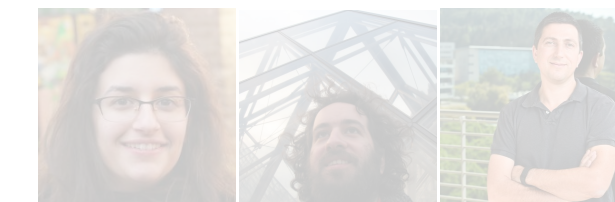
**Counting**
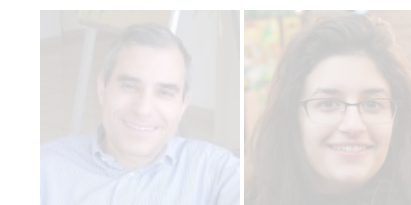LSTMs are counter machines, GRUs aren't **(ACL 2018)**

**RASP**
Finding a formalism to describe transformers **(ICML 2021)**

**DFAs from RNNs**
Applying L* to learn DFAs from RNNs **(ICML 2018)**
    + using the result for CFGs **(TACAS 2021)**

# RASP

# RASP



How is the transformer… doing things?

(How) does it count?

(How) does it reason?

(What) does attention explain?

# Transformers

**Attention Is All You Need**

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit,

Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Transformers

**Attention Is All You Need**

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit,

Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin

**Decoder**

**Encoder**

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

N×

Positional Encoding

Output Embedding

Outputs (shifted right)

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Positional Encoding

Input Embedding

Inputs

35

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Transformers

**Attention Is All You Need**

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit,

Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin

## Encoder

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Transformers



N×

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Positional
Encoding

Input
Embedding

Inputs

- Receive their entire input 'at once', processing all tokens in parallel

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Transformers



Positional Encoding

Inputs

- Receive their entire input 'at once', processing all tokens in parallel

- Have a fixed number of layers, such that the output of one is the input of the next

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Transformers



- Receive their entire input 'at once', processing all tokens in parallel

- Have a fixed number of layers, such that the output of one is the input of the next

Computation "progresses" along network depth… not input length

# Transformers

| I | Like | Dogs |
|---|---|---|

$e(\text{I})$ | $e(\text{Like})$ | $e(\text{dogs})$     $p(0)$ | $p(1)$ | $p(2)$

$\oplus$

| $x_1$ | $x_2$ | $x_3$ |
|---|---|---|

**Encoder Layer 1**

| $y_1^1$ | $y_2^1$ | $y_3^1$ |
|---|---|---|

…

**Encoder Layer L**

| $y_1^L$ | $y_2^L$ | $y_3^L$ |
|---|---|---|

tokens = positionwise_embeddings(input)

indices = positionwise_indices(input)

$x$ = tokens+indices

$$y^1 = L_1(x)$$

$$y^2 = L_2(y^1)$$

…

$$y = y^L = L_L\,(y^{L-1}\,)$$

Layer input/outputs are "variables" of a transformer "program"
The layers themselves are "operations"

# RASP (Restricted Access Sequence Processing)



Layer input/outputs are "variables" of a transformer "program"
The layers themselves are "operations"

# RASP (Restricted Access Sequence Processing)

- A transformer-encoder is a sequence to sequence function ("sequence operator", or, "**s-op**")

- Its **layers apply operations** to the sequences

- **RASP builds s-ops**, constrained to a transformer's inputs and possible operations

  - (The s-ops are the transformer abstractions!)



Layer input/outputs are "variables" of a transformer "program"
The layers themselves are "operations"

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# RASP base s-ops

**Word Embedding**

| |
|:---:|
| $e(\text{I})$ |
| $e(\text{Like})$ |
| $e(\text{dogs})$ |

$$d_x$$

I

Like

Dogs

**Positional Embedding**

| |
|:---:|
| $p(0)$ |
| $p(1)$ |
| $p(2)$ |

$$d_x$$

The information before a transformer has done anything ("0 layer transformer")

*tokens* and *indices* are RASP built-ins:

```
>> tokens;
    s-op: tokens

>> indices;
    s-op: indices
```

# RASP base s-ops

**Word Embedding**

$e(\text{I})$

$e(\text{Like})$

$e(\text{dogs})$

$d_x$

I

Like

Dogs

**Positional Embedding**

$p(0)$

$p(1)$

$p(2)$

$d_x$

The information before a
transformer has done anything
("0 layer transformer")

*tokens* and *indices* are RASP built-ins:

```
>> tokens;
    s-op: tokens
        Example: tokens("hello") = [h, e, l, l, o] (strings)
>> indices;
    s-op: indices
        Example: indices("hello") = [0, 1, 2, 3, 4] (ints)
```

The RASP REPL gives you
examples (until you ask it not to)

44

# Okay, now what?

```
>> tokens;
    s-op: tokens
        Example: tokens("hello") = [h, e, l, l, o] (strings)
>> indices;
    s-op: indices
        Example: indices("hello") = [0, 1, 2, 3, 4] (ints)
```

To know what operations RASP may have, we must
inspect the transformer-encoder layers!

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Transformer-Encoder Layer

**Input**

$x_1$

$x_2$

$x_3$

$d_x$

Multi-Head Attention

$o_1$

$o_2$

$o_3$

$d_x$

Linear Transformation

$A$

Residual ("Skip") Connection

Layer Norm 1

$d_x$

Feed-Forward Sublayer

$W_1^{ff}$

$ff_1'$

$ff_2'$

$ff_3'$

$d_{ff}$

ReLU

$W_2^{ff}$

$o_1''$

$o_2''$

$o_3''$

$d_x$

Residual ("Skip") Connection

Layer Norm 2

**There's a lot in here…**

**Output**

$out_1$

$out_2$

$out_3$

$d_x$

46

# Feed-Forward Sublayer

**Input**



$x_1$

$x_2$

$x_3$

$d_x$

Multi-Head Attention

$o_1$

$o_2$

$o_3$

$d_x$

Linear Transformation

$A$

Residual ("Skip") Connection

Layer Norm 1

$d_x$

**Feed-Forward Sublayer**

$\otimes$

$W^{ff}_1$

$ff_1$

$ff_2$

$ff_3$

$d_{ff}$

ReLU

$\otimes$

$W^{ff}_2$

$o''_1$

$o''_3$

Layer Norm 2

Multilayer Feedforward Networks are Universal Approximators (Hornik et al, 1989)

**Output**

$out_1$

$out_2$

$out_3$

$d_x$

# Feed-Forward Sublayer



**Feed-Forward Sublayer**

$$W^{ff}_1 \qquad ff_1 \quad ff_2 \quad ff_3 \qquad \text{ReLU} \qquad W^{ff}_2$$

$d_{ff}$

### Multilayer Feedforward Networks are Universal Approximators

KURT HORNIK

Technische Universität Wien

MAXWELL STINCHCOMBE AND HALBERT WHITE

University of California, San Diego

**Abstract**—*This paper rigorously establishes that standard multilayer feedforward networks with as few as one hidden layer using arbitrary squashing functions are capable of approximating any Borel measurable function from one finite dimensional space to another to any desired degree of accuracy, provided sufficiently many hidden units are available. In this sense, multilayer feedforward networks are a class of universal approximators.*

```
|>> indices+1;
    s-op: out
        Example: out("hello") = [1, 2, 3, 4, 5] (ints)
|>> tokens=="e" or tokens=="o";
    s-op: out
        Example: out("hello") = [F, T, F, F, T] (bools)
```

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)
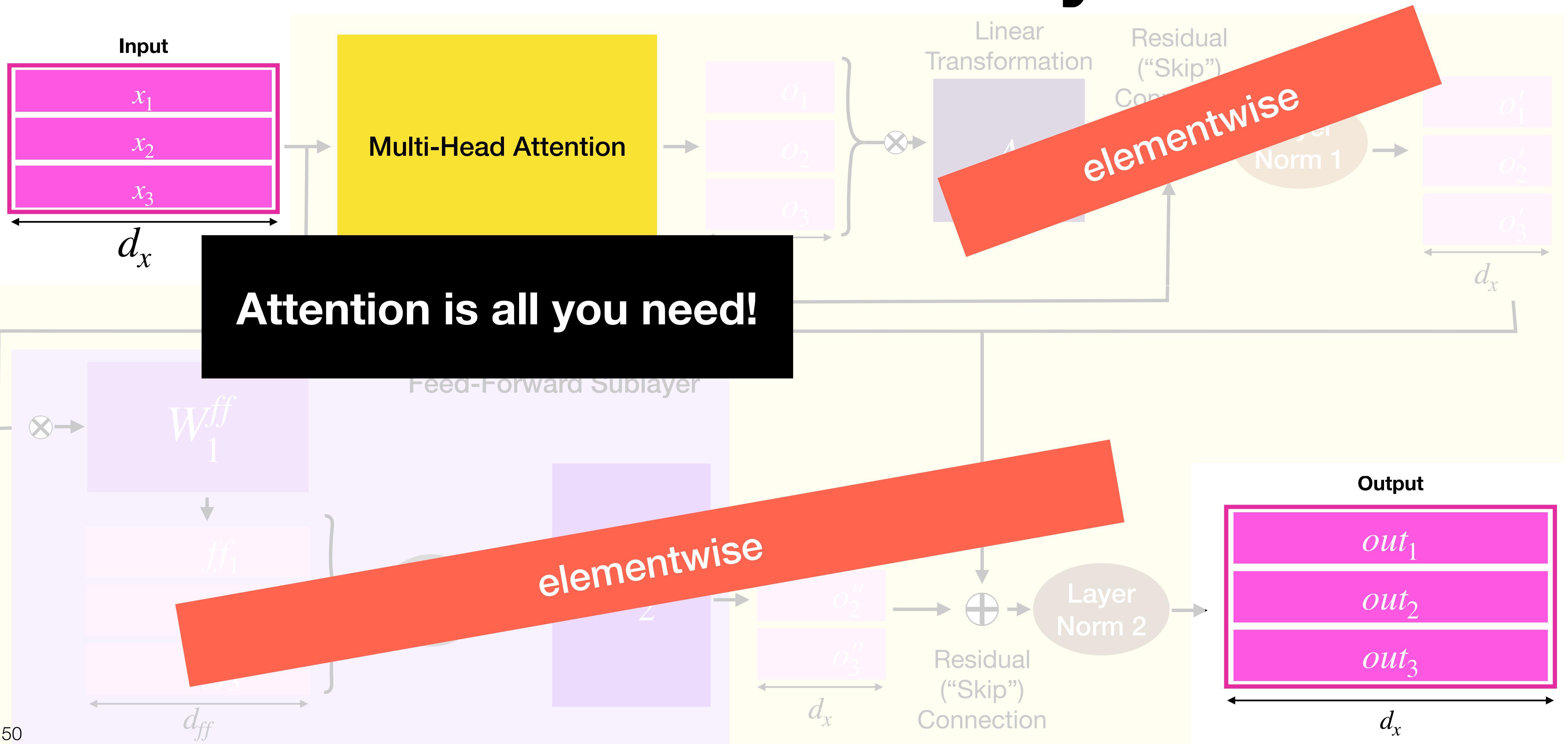
# So far

```
>> tokens;
    s-op: tokens
        Example: tokens("hello") = [h, e, l, l, o] (strings)
>> indices;
    s-op: indices
        Example: indices("hello") = [0, 1, 2, 3, 4] (ints)


>> indices+1;
    s-op: out
        Example: out("hello") = [1, 2, 3, 4, 5] (ints)
>> tokens=="e" or tokens=="o";
    s-op: out
        Example: out("hello") = [F, T, F, F, T] (bools)
```

**Are we all-powerful
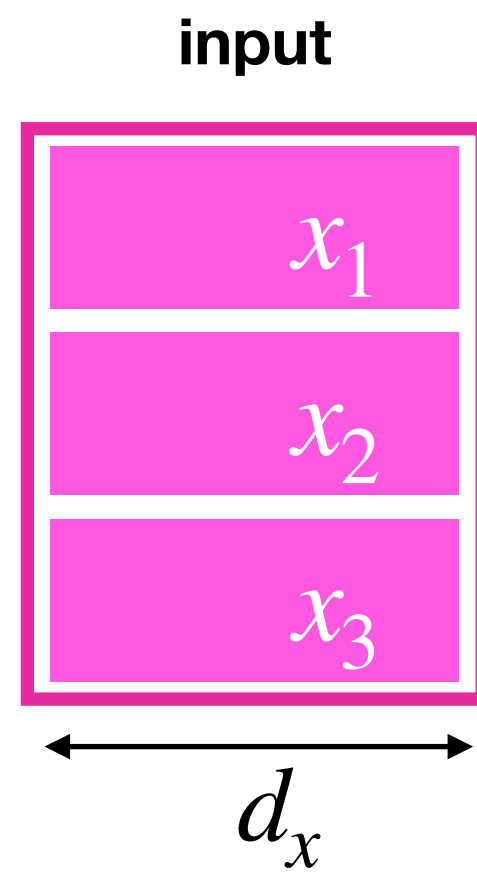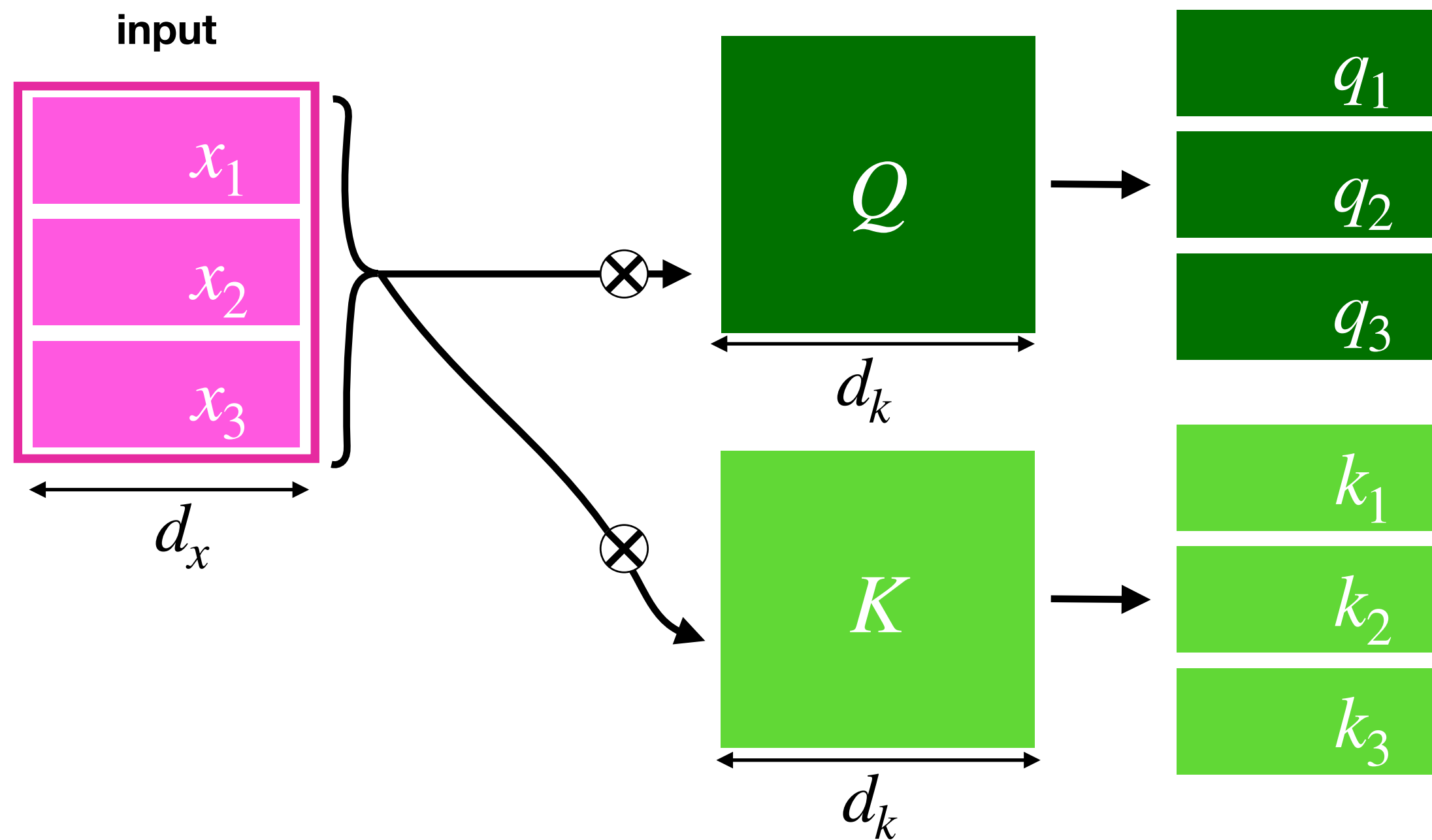(well, transformer-powerful) yet?**

# Attention Sublayer

**Input**

| |
|---|
| $x_1$ |
| $x_2$ |
| $x_3$ |

$d_x$

Multi-Head Attention

$o_1$

$o_2$

$o_3$

Linear
Transformation

Residual
("Skip")
Con...

**elementwise**

Layer
Norm 1

$d_x$

**Attention is all you need!**

Feed-Forward Sublayer

$W_1^{ff}$

$ff_1$

$d_{ff}$

**elementwise**

$z$

$o_3''$

$o_3''$

$d_x$

Residual
("Skip")
Connection

Layer
Norm 2

**Output**

| |
|---|
| $out_1$ |
| $out_2$ |
| $out_3$ |

$d_x$

# Background - Multi Head Attention

Starting from single-head attention…

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Background - Self Attention (Single Head)

**input**

$$x_1$$

$$x_2$$

$$x_3$$

$$d_x$$

# Background - Self Attention (Single Head)

input

$x_1$

$x_2$

$x_3$

$d_x$

$Q$

$d_k$

$K$

$d_k$

$q_1$

$q_2$

$q_3$

$k_1$

$k_2$

$k_3$

# Background - Self Attention (Single Head)

input

$x_1$

$x_2$

$x_3$

$d_x$

$\otimes$

$Q$

$d_k$

$\otimes$

$K$

$d_k$

$q_1$

$q_2$

$q_3$

$k_1$

$k_2$

$k_3$

# Background - Self Attention (Single Head)

**scores**

**input**

$x_1$

$x_2$

$x_3$

$d_x$

$\otimes$

$\otimes$

$Q$

$d_k$

$K$

$d_k$

$q_1$

$q_2$

$q_3$

$k_1$

$k_2$

$k_3$

$q_1 \cdot k_1$

# Background - Self Attention (Single Head)



scores

input

$x_1$

$x_2$

$x_3$

$d_x$

$Q$

$d_k$

$K$

$d_k$

$q_1$

$q_2$

$q_3$

$k_1$

$k_2$

$k_3$

$q_1 \cdot k_1$   $q_1 \cdot k_2$

# Background - Self Attention (Single Head)

input

$x_1$

$x_2$

$x_3$

$d_x$

$Q$

$d_k$

$K$

$d_k$

$q_1$

$q_2$

$q_3$

$k_1$

$k_2$

$k_3$

$q_1 \cdot k_1$  $q_1 \cdot k_2$  $q_1 \cdot k_3$

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Background - Self Attention (Single Head)

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Background - Self Attention (Single Head)



**scores**

**input**

$x_1$

$x_2$

$x_3$

$d_x$

$Q$

$d_k$

$q_1$

$q_2$

$q_3$

$q_1 \cdot k_1$  $q_1 \cdot k_2$  $q_1 \cdot k_3$

normalise (i.e. $\times 1/\sqrt{d_k}$)

softmax

$K$

$d_k$

$k_1$

$k_2$

$k_3$

$w_{1,1}$  $w_{1,2}$  $w_{1,3}$

**weights**

$V$

$d_v$

$v_1$

$v_2$

$v_3$

# Background - Self Attention (Single Head)

**input**

$x_1$

$x_2$

$x_3$

$d_x$

$Q$

$q_1$

$q_2$

$q_3$

$d_k$

$K$

$k_1$

$k_2$

$k_3$

$d_k$

$q_1 \cdot k_1$  $q_1 \cdot k_2$  $q_1 \cdot k_3$

normalise (i.e. $\times 1/\sqrt{d_k}$ )

softmax

$w_{1,1}$ $w_{1,2}$ $w_{1,3}$

**weights**

$V$

$d_v$

$v_1$

$v_2$

$v_3$

$out_1$

$d_v$

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Background - Self Attention (Single Head)



scores

input

$x_1$

$x_2$

$x_3$

$d_x$

$Q$

$d_k$

$q_1$

$q_2$

$q_3$

$k_1$

$k_2$

$k_3$

$K$

$d_k$

$V$

$d_v$

$q_2 \cdot k_1 \quad q_2 \cdot k_2 \quad q_2 \cdot k_3$

normalise (i.e. $\times 1/\sqrt{d_k}$)

softmax

$w_{2,1} \quad w_{2,2} \quad w_{2,3}$

weights

$v_1$

$v_2$

$v_3$

$out_1$

$out_2$

$d_v$

# Background - Self Attention (Single Head)

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Background - Self Attention (Single Head)



**Attention Head**

input

$x_1$

$x_2$

$x_3$

$d_x$

$Q$

$d_k$

$K$

$d_k$

$V$

$d_v$

$q_1$

$q_2$

$q_3$

$k_1$

$k_2$

$k_3$

**scores**

$q_3 \cdot k_1$ $q_3 \cdot k_2$ $q_3 \cdot k_3$

normalise (i.e. $\times 1/\sqrt{d_k}$)

softmax

$w_{3,1}$ $w_{3,2}$ $w_{3,3}$

**weights**

$v_1$

$v_2$

$v_3$

$out_1$

$out_2$

$out_3$

$d_v$

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# So, how do we present an attention head?

# Self Attention (Single Head)



**Attention Head**

**input**

**scores**

$x_1$
$x_2$
$x_3$

$d_x$

$Q$

$d_k$

$K$

$d_k$

$V$

$d_x$

$q_1$
$q_2$
$q_3$

$k_1$
$k_2$
$k_3$

$q_3 \cdot k_1$  $q_3 \cdot k_2$  $q_3 \cdot k_3$

normalise (i.e. $\times 1/\sqrt{d_k}$ )

softmax

$w_{3,1}$ $w_{3,2}$ $w_{3,3}$

**weights**

$v_1$
$v_2$
$v_3$

$out_1$
$out_2$
$out_3$

$d_x$

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Self Attention (Single Head)



input

$x_1$

$x_2$

$x_3$

$d_x$

Attention Head

$Q$

$d_k$

$K$

$d_k$

$q_1$

$q_2$

$q_3$

$k_1$

$k_2$

$k_3$

scores

$q_3 \cdot k_1$ $q_3 \cdot k_2$ $q_3 \cdot k_3$

normalise (i.e. $\times 1/\sqrt{d_k}$ )

softmax

$w_{3,1}$ $w_{3,2}$ $w_{3,3}$

weights

**Pairwise!**

$V$

$d_x$

$v_1$

$v_2$

$v_3$

$out_1$

$out_2$

$out_3$

$d_x$

# Self Attention (Single Head)

**Attention Head**

scores

**input**

$x_1$

$x_2$

$x_3$

$d_x$

$Q$

$d_k$

$q_1$

$q_2$

$q_3$

$q_3 \cdot k_1$   $q_3 \cdot k_2$   $q_3 \cdot k_3$

normalise (i.e. $\times 1/\sqrt{d_k}$)

$K$

$d_k$

$k_1$

$k_2$

$k_3$

softmax

$w_{3,1}$ $w_{3,2}$ $w_{3,3}$

**weights**

$V$

$d_x$

$v_1$

$v_2$

$v_3$

$out_1$

$out_2$

$out_3$

$d_x$

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Single Head: Scoring $\leftrightarrow$ Selecting

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Single Head: Scoring ↔ Selecting

Decision: RASP abstracts to binary
select/don't select decisions

**sel = select([2,0,0],[0,1,2],==)**

|   | 2 | 0 | 0 |
|---|---|---|---|
| **0** | F | **T** | **T** |
| **1** | F | F | F |
| **2** | **T** | F | F |

# Single Head: Scoring ↔ Selecting

Decision: RASP abstracts to binary
select/don't select decisions

**sel = select([2,0,0],[0,1,2],==)**

|   | 2 | 0 | 0 |
|---|---|---|---|
| 0 | F | T | T |
| 1 | F | F | F |
| 2 | T | F | F |

Decision: RASP abstracts to binary
select/don't select decisions

**sel = select([2,0,0], [0,1,2], ==)**

```
        2   0   0
   0    F   T   T
   1    F   F   F
   2    T   F   F
```

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Single Head: Scoring ↔ Selecting

Decision: RASP abstracts to binary
select/don't select decisions

**sel = select([2,0,0],[0,1,2],==)**



|   | 2 | 0 | 0 |
|---|---|---|---|
| 0 | F | T | T |
| 1 | F | F | F |
| 2 | T | F | F |

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Single Head: Scoring ↔ Selecting

Decision: RASP abstracts to binary
select/don't select decisions

**sel = select([2,0,0],[0,1,2],==)**

$$
\begin{array}{c c c c}
 & 2 & 0 & 0 \\
0 & F & T & T \\
1 & F & F & F \\
2 & T & F & F \\
\end{array}
$$

# Single Head: Scoring ↔ Selecting

Decision: RASP abstracts to binary
select/don't select decisions

**sel = select([2,0,0],[0,1,2],==)**

|       | 2 | 0 | 0 |
|-------|---|---|---|
| **0** | F | T | T |
| **1** | F | F | F |
| **2** | T | F | F |

# Single Head: Scoring ↔ Selecting

Decision: RASP abstracts to binary select/don't select decisions

**sel = select([2,0,0],[0,1,2],==)**

# Single Head: Scoring ↔ Selecting

Decision: RASP abstracts to binary select/don't select decisions

**sel** = **select([2,0,0],[0,1,2],==)**



|   | 2 | 0 | 0 |
|---|---|---|---|
| 0 | F | T | T |
| 1 | F | F | F |
| 2 | T | F | F |

Another example:

**sel2 = select([2,0,0],[0,1,2] >=)**

|   | 2 | 0 | 0 |
|---|---|---|---|
| 0 | T | T | T |
| 1 | T | F | F |
| 2 | T | F | F |

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Single Head: Scoring ↔ Selecting

**prevs** = **select([0,1,2],[0,1,2],<=)**

|   | 0 | 1 | 0 |
|---|---|---|---|
| 0 | T | F | F |
| 1 | T | T | F |
| 2 | T | T | T |

# Single Head: Scoring ↔ Selecting

**prevs** = **select([0,1,2],[0,1,2],<=)**

$(1, 0, 0, \ldots)\ k_1$

$(0, 1, 0, \ldots)\ k_2$

$(0, 0, 1, \ldots)\ k_3$

|   | 0 | 1 | 2 |
|---|---|---|---|
| 0 | T | F | F |
| 1 | T | T | F |
| 2 | T | T | T |

# Single Head: Scoring ↔ Selecting

**prevs = select([0,1,2],[0,1,2],<=)**

|   | 0 | 1 | 2 |
|---|---|---|---|
| **0** | T | F | F |
| **1** | T | T | F |
| **2** | T | T | T |

$(1, 0, 0, \ldots)\ k_1$

$(0, 1, 0, \ldots)\ k_2$

$(0, 0, 1, \ldots)\ k_3$

$(1, 0, 0, \ldots)\ q_1$

$(1, 1, 0, \ldots)\ q_2$

$(1, 1, 1, \ldots)\ q_3$

# Single Head: Weighted Average ↔ Aggregation



**Attention Head**

scores

input

$x_1$

$x_2$

$x_3$

$d_x$

$Q$

$d_k$

$q_3$

$q_3 \cdot k_1$  $q_3 \cdot k_2$  $q_3 \cdot k_3$

normalise (i.e. $\times 1/\sqrt{d_k}$)

$K$

$d_k$

$k_1$

$k_2$

$k_3$

softmax

$w_{3,1}$ $w_{3,2}$ $w_{3,3}$

**weights**

$V$

$d_x$

$v_1$

$v_2$

$v_3$

$out_1$

$out_2$

$out_3$

$d_x$

# Single Head: Weighted Average ↔ Aggregation

**new=aggregate(sel, [1,2,4])**

```
           1 2 4
F  T  T    1 2 4  =>  3
F  F  F    1 2 4  =>  0  =>  [3,0,1]
T  F  F    1 2 4  =>  1
```

# Single Head: Weighted Average ↔ Aggregation

**new**=**aggregate(sel, [1,2,4])**

$$
\begin{array}{ccc}
 & \mathbf{1\ 2\ 4} & \\
F\ \mathbf{T}\ \mathbf{T} & \boxed{\mathbf{1\ 2\ 4}} & => \quad 3 \\
F\ F\ F & 1\ 2\ 4 & => \quad 0 \quad => \quad \mathbf{[3,0,1]} \\
\mathbf{T}\ F\ F & \mathbf{1}\ 2\ 4 & => \quad 1 \\
\end{array}
$$

# Single Head: Weighted Average ↔ Aggregation

**new=aggregate(sel, [1,2,4])**



```
      1 2 4
F  T  T   1 2 4  =>  3
F  F  F   1 2 4  =>  0   =>  [3,0,1]
T  F  F   1 2 4  =>  1
```

# Single Head: Weighted Average ↔ Aggregation

**new=aggregate(sel, [1,2,4])**

```
        1 2 4
F T T   1 2 4  =>   3
F F F   1 2 4  =>   0   =>  [3,0,1]
T F F   1 2 4  =>   1
```

# Single Head: Weighted Average ↔ Aggregation

**new**=**aggregate**(**sel**, **[1,2,4]**)

```
        1 2 4
  F  T  T   1 2 4  =>    3
  F  F  F   1 2 4  =>    0    =>   [3,0,1]
  T  F  F   1 2 4  =>    1
```

# Single Head: Weighted Average ↔ Aggregation

**new**=**aggregate**(**sel**, **[1,2,4]**)

```
         1 2 4
F  T  T  1 2 4  =>   3
F  F  F  1 2 4  =>   0   =>   [3,0,1]
T  F  F  1 2 4  =>   1
```

# Single Head: Weighted Average ↔ Aggregation

**new**=**aggregate**(**sel**, **[1,2,4]**)

```
        1 2 4
F  T  T  1 2 4  =>  3
F  F  F  1 2 4  =>  0   =>  [3,0,1]
T  F  F  1 2 4  =>  1
```



Symbolic language + no averaging when only
one position selected allows (for example):

**reverse**=**aggregate**(**flip**, **[A,B,C]**)

```
        A B C
F  F  T  A B C  =>  C
F  T  F  A B C  =>  B  =>  [C,B,A]
T  F  F  A B C  =>  A
```

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Great!
# Now do multi-headed attention

# Background - Multi-Headed Self Attention

**Input**

$x_1$

$x_2$

$x_3$

$d_x$

$x_1$

$x_2$

$x_3$

$d_h$  $d_h$  $\cdots$  $d_h$

$$d_k = d_v = d_h = \frac{d_x}{H}$$

Head 1

Head 2

$\bullet\ \bullet\ \bullet$

Head H

$out_1^1$

$out_2^1$

$out_3^1$

$d_h$

$out_1^2$

$out_2^2$

$out_3^2$

$d_h$

$\cdots$

$\cdots$

$\cdots$

$out_1^H$

$out_2^H$

$out_3^H$

$d_h$

**Concatenate**

**Output**

$out_1$

$out_2$

$out_3$

$d_x$

89

# The multi-headed attention lets one layer do multiple single head operations

We do not need 'new' RASP operations to describe it!

(We will just let the RASP compiler know it can place multiple heads on the same layer)

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Transformer-Encoder Layer

**Input**

$x_1$

$x_2$

$x_3$

$d_x$

Multi-Head Attention

$o_1$

$o_2$

$o_3$

$d_x$

Linear Transformation

$A$

Residual ("Skip") Connection

Layer Norm 1

$d_x$

Feed-Forward Sublayer

$W^{ff}_1$

$ff_1$

$ff_2$

$ff_3$

$d_{ff}$

ReLU

$W^{ff}_2$

$o''_1$

$o''_2$

$o''_3$

$d_x$

Residual ("Skip") Connection

Layer Norm 2

**Output**

$out_1$

$out_2$

$out_3$

$d_x$

# Transformer-Encoder Layer

# RASP (Restricted Access Sequence Processing)

**Initial Sequences**

```
>> tokens;
    s-op: tokens
        Example: tokens("hello") = [h, e, l, l, o] (strings)
>> indices;
    s-op: indices
        Example: indices("hello") = [0, 1, 2, 3, 4] (ints)
```

**Elementwise application of atomic operations**

```
>> indices+1;
    s-op: out
        Example: out("hello") = [1, 2, 3, 4, 5] (ints)
>> tokens=="e" or tokens=="o";
    s-op: out
        Example: out("hello") = [F, T, F, F, T] (bools)
```

**Selectors, and aggregate**

sel = select([2,0,0],[0,1,2],==)

```
      2 0 0
  0   F T T
  1   F F F
  2   T F F
```

new=aggregate(sel, [1,2,4])

```
              1 2 4
      F T T   1 2 4  =>  3
      F F F   1 2 4  =>  0  =>  [3,0,1]
      T F F   1 2 4  =>  1
```

```
>> flip = select(length-indices-1,indices,==);
    selector: flip
        Example:
                          h e l l o
                      h |         1
                      e |       1
                      l |     1
                      l |   1
                      o | 1
>> reverse = aggregate(flip,tokens);
    s-op: reverse
        Example: reverse("hello") = [o, l, l, e, h]
```
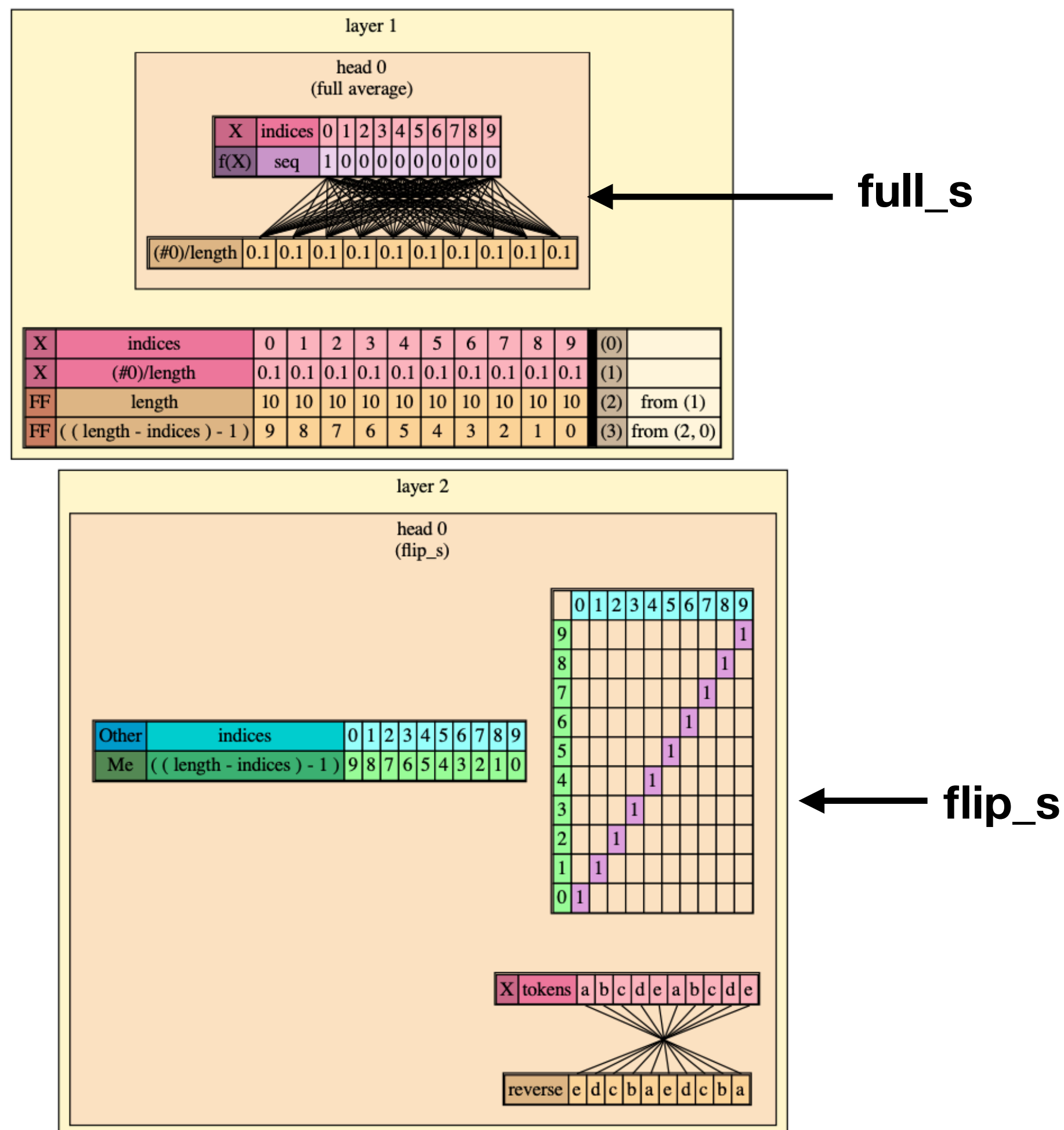
# RASP Extras

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# RASP Extras

**Extra Sequences**

```
>> length;
    s-op: length
        Example: length("hello") = [5]*5 (ints)
```

# RASP Extras

**Extra Sequences**

```
>> length;
    s-op: length
        Example: length("hello") = [5]*5 (ints)
```

**Selector Compositions**

```
>> select(indices,3,==) or select(indices,indices,<=);
    selector: out
        Example:
                    h e l l o
                h | 1     1
                e | 1 1   1
                l | 1 1 1 1
                l | 1 1 1 1
                o | 1 1 1 1 1
```

# RASP Extras

**Extra Sequences**

```
>> length;
    s-op: length
        Example: length("hello") = [5]*5 (ints)
```

**Functions**

```
>> def in_range(min,val,max) {
..          return (min<=val) and (val<=max);
..      }
    console function: in_range(min, val, max)


>> in_range(1,indices,3);
    s-op: out
        Example: out("hello") = [F, T, T, T, F]
```

**Selector Compositions**

```
>> select(indices,3,==) or select(indices,indices,<=);
    selector: out
        Example:

                  h e l l o
              h | 1       1
              e | 1 1     1
              l | 1 1 1 1
              l | 1 1 1 1
              o | 1 1 1 1 1
```

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# RASP Extras

**Extra Sequences**

```
>> length;
    s-op: length
        Example: length("hello") = [5]*5 (ints)
```

**Selector Compositions**

```
>> select(indices,3,==) or select(indices,indices,<=);
    selector: out
        Example:
                    h e l l o
                h | 1       1
                e | 1 1     1
                l | 1 1 1 1
                l | 1 1 1 1
                o | 1 1 1 1 1
```

**Functions**

```
>> def in_range(min,val,max) {
..          return (min<=val) and (val<=max);
..     }
    console function: in_range(min, val, max)


>> in_range(1,indices,3);
    s-op: out
        Example: out("hello") = [F, T, T, T, F]
```

**Library Functions**

```
>> selector_width(select(tokens,tokens,==));
    s-op: out
        Example: out("hello") = [1, 1, 2, 2, 1] (ints)

>> count(tokens,"l");
    s-op: out
        Example: out("hello") = [2]*5 (ints)
```

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# RASP Extras

**Extra Sequences**

```
>> length;
    s-op: length
        Example: length("hello") = [5]*5 (ints)
```

**Functions**

```
>> def in_range(min,val,max) {
..          return (min<=val) and (val<=max);
..     }
    console      tion: in_range(mi        max)
```

**Selector Compositions**

```
>> select(indices,3,==) or select(indices,i  ices,<=);
    selector: out
        Example:
                  e l l o
              h       1
              e       1
              l       1
              l       1
              o | 1 1 1 1 1
```

**Library Functions**

```
>> selector_width select(tokens,tokens,==));
    s-op: out
        Example: out("hello") = [1, 1, 2, 2, 1] (ints)

>> count tokens,"l");
    s-op: out
        Example: out("hello") = [2]*5 (ints)
```

# Small Example

**Computing *length*:**

```
>> full_s = select(1,1,==);
    selector: full_s
        Example:

                    h e l l o
             h | 1 1 1 1 1
             e | 1 1 1 1 1
             l | 1 1 1 1 1
             l | 1 1 1 1 1
             o | 1 1 1 1 1
```

# Small Example

**Computing *length*:**

```
>> full_s = select(1,1,==);
    selector: full_s
        Example:

                        h e l l o
                h |  1 1 1 1 1
                e |  1 1 1 1 1
                l |  1 1 1 1 1
                l |  1 1 1 1 1
                o |  1 1 1 1 1
                    indicator(indices==0)
```

```
>> indicator(indices==0);
    s-op: out
        Example: out("hello") = [1, 0, 0, 0, 0] (ints)
```

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Small Example

**Computing *length*:**

```
>> full_s = select(1,1,==);
    selector: full_s
        Example:

                        h e l l o
                    h | 1 1 1 1 1
                    e | 1 1 1 1 1
                    l | 1 1 1 1 1
                    l | 1 1 1 1 1
                    o | 1 1 1 1 1
>> frac_0=aggregate(full_s,indicator(indices==0));
    s-op: frac_0
        Example: frac_0("hello") = [0.2]*5 (floats)
```

```
>> indicator(indices==0);
    s-op: out
        Example: out("hello") = [1, 0, 0, 0, 0] (ints)
```

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Small Example

**Computing _length_:**

```
>> full_s = select(1,1,==);
    selector: full_s
        Example:

                        h e l l o
                    h | 1 1 1 1 1
                    e | 1 1 1 1 1
                    l | 1 1 1 1 1
                    l | 1 1 1 1 1
                    o | 1 1 1 1 1
>> frac_0=aggregate(full_s,indicator(indices==0));
    s-op: frac_0
        Example: frac_0("hello") = [0.2]*5 (floats)
>> round(1/frac_0);
    s-op: out
        Example: out("hello") = [5]*5 (ints)
```

```
>> indicator(indices==0);
    s-op: out
        Example: out("hello") = [1, 0, 0, 0, 0] (ints)
```

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Small Example

**Computing *length*:**

```
>> full_s = select(1,1,==);
    selector: full_s
        Example:

                        h e l l o
                    h | 1 1 1 1 1
                    e | 1 1 1 1 1
                    l | 1 1 1 1 1
                    l | 1 1 1 1 1
                    o | 1 1 1 1 1
>> frac_0=aggregate(full_s,indicator(indices==0));
    s-op: frac_0
        Example: frac_0("hello") = [0.2]*5 (floats)
>> round(1/frac_0);
    s-op: out
        Example: out("hello") = [5]*5 (ints)
```

```
>> indicator(indices==0);
    s-op: out
        Example: out("hello") = [1, 0, 0, 0, 0] (ints)
```

*Can you see how to use
this trick for
selector_width?*

# Connection to Reality?

RASP expects 2 layers for arbitrary-length reverse

```
>> flip = select(length-indices-1,indices,==);
    selector: flip
        Example:
                        h e l l o
                    h |         1
                    e |       1
                    l |     1
                    l |   1
                    o | 1
>> reverse = aggregate(flip,tokens);
    s-op: reverse
        Example: reverse("hello") = [o, l, l, e, h] (strings)
```

# Connection to Reality?

```
>> draw(reverse,"abcdeabcde")
```



RASP expects 2 layers for arbitrary-length reverse

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Connection to Reality?

```
[>> draw(reverse,"abcdeabcde")
```



full_s



flip_s

RASP expects 2 layers for arbitrary-length reverse

**Test**:
Training small transformers on lengths 0-100:

2 layers: **99.6**% accuracy after 20 epochs
1 layer: **39.6**% accuracy after 50 epochs

Even with compensation for number of heads and parameters!

# Connection to Reality?

```
>> draw(reverse,"abcdeabcde")
```



full_s

flip_s

RASP expects 2 layers for arbitrary-length reverse

**Test**:
Training small transformers on lengths 0-100:

2 layers: **99.6**% accuracy after 20 epochs
1 layer: **39.6**% accuracy after 50 epochs

Even with compensation for number of heads and parameters!

**Bonus**: the 2 layer transformer's attention patterns:

Layer 1 (*full_s*)

Layer 2 (*flip_s*)

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Connection to Reality?

**Example 2:** *histogram*  (assuming BOS)

**in place histogram,
with BOS - examples:**

[§,a,a,a,b] -> [0,3,3,3,1]
[§,a,b,a,c] -> [0,2,1,2,1]
[§,a,b,c,c] -> [0,1,1,2,2]

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Connection to Reality?

**Example 2: *histogram*** (assuming BOS)

```
>> examples off
>> same_or_0 = select(tokens,tokens,==) or select(indices,0,==);
    selector: same_or_0
>> frac_with_0 = aggregate(same_or_0,indicator(indices==0));
    s-op: frac_with_0
>> histogram_assuming_bos = round(1/frac_with_0)-1;
    s-op: histogram_assuming_bos
>> histogram_assuming_bos("§hello");
        =  [0, 1, 1, 2, 2, 1] (ints)
```

**in place histogram,
with BOS - examples:**

[§,a,a,a,b] -> [0,3,3,3,1]
[§,a,b,a,c] -> [0,2,1,2,1]
[§,a,b,c,c] -> [0,1,1,2,2]

# Connection to Reality?

**Example 2: *histogram*** (assuming BOS)

```
[>> examples off
[>> same_or_0 = select(tokens,tokens,==) or select(indices,0,==);
    selector: same_or_0
[>> frac_with_0 = aggregate(same_or_0,indicator(indices==0));
    s-op: frac_with_0
[>> histogram_assuming_bos = round(1/frac_with_0)-1;
    s-op: histogram_assuming_bos
[>> histogram_assuming_bos("§hello");
        =  [0, 1, 1, 2, 2, 1] (ints)
```

**RASP analysis:**

- Just one attention head
- It focuses on:
  1. All positions with same token, and:
  2. Position 0 (regardless of content)

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Connection to Reality?

**Example 2: *histogram*** (assuming BOS)

```
>> examples off
>> same_or_0 = select(tokens,tokens,==) or select(indices,0,==);
   selector: same_or_0
>> frac_with_0 = aggregate(same_or_0,indicator(indices==0));
   s-op: frac_with_0
>> histogram_assuming_bos = round(1/frac_with_0)-1;
   s-op: histogram_assuming_bos
>> histogram_assuming_bos("§hello");
       =  [0, 1, 1, 2, 2, 1] (ints)
```

**RASP analysis:**

- Just one attention head
- It focuses on:
  1. All positions with same token, and:
  2. Position 0 (regardless of content)

Selector pattern vs trained
transformer's attention for same
input sequence:

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Connection to Reality?

**Example 2: *histogram*** (assuming BOS)

```
>> examples off
>> same_or_0 = select(tokens,tokens,==) or select(indices,0,==);
    selector: same_or_0
>> frac_with_0 = aggregate(same_or_0,indicator(indices==0));
    s-op: frac_with_0
>> histogram_assuming_bos = round(1/frac_with_0)-1;
    s-op: histogram_assuming_bos
>> histogram_assuming_bos("§hello");
        =  [0, 1, 1, 2, 2, 1] (ints)
```

### RASP analysis:

- Just one attention head
- It focuses on:
  1. All positions with same token, and:
  2. Position 0 (regardless of content)

Selector pattern vs trained
transformer's attention for same
input sequence:



Try it out!
🌟 github.com/tech-srl/RASP 🌟

Thinking Like Transformers (Weiss, Goldberg, Yahav, ICML 2021)

# Neural Sequence Models: a Formal Lens

## Counting
LSTMs are counter machines, GRUs aren't (ACL 2018)

## RASP
Finding a formalism to describe transformers (ICML 2021)

## **DFAs from RNNs**
Applying L* to learn DFAs from RNNs (ICML 2018)

+ using the result for CFGs (TACAS 2021)

# DFAs from RNNs



**Goal:**
Concise (Meaningful) Model
from Trained RNN

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# Previous Approaches

1. **Partition** RNN state space

2. Explore using **pruned BFS** or transition sampling



e.g.: Omlin and Giles (1996), Cechin et al. (2003)

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# Previous Approaches

1. **Too coarse:** not representative

2. **Too fine:** very large: slow & memory consuming extraction

## Impractical!



e.g.: Omlin and Giles (1996), Cechin et al. (2003)

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# L* (Angluin, 1987)

An exact learning algorithm for DFAs

Learns using:

- **Membership Queries** (request to label input sequence) and
- **Equivalence Queries** (request to accept/reject DFA)

*Creates hypothesis DFA and improves it until accepted by teacher*



Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# Iterative Approach

Apply L* to RNN:

**Membership queries** are trivial

(**Equivalence queries** are hard)

Use **equivalence queries** to induce the **partitioning** of the RNN state space

Use the **partitioning** to answer the **equivalence queries**

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# Iterative Approach

Partitioning

L*

$\varepsilon$? ✔  $bb$? ✔
$a$? ✔  $bab$? ✔
       $baa$? ✘
$b$? ✔  $bba$? ✘
       ...

?

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# Iterative Approach

Partitioning

L*

$\varepsilon$? ✔     $bb$? ✔
             $bab$? ✔
$a$? ✔     $baa$? ✘
$b$? ✔     $bba$? ✘
             ...

✔$aba$✘

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# Iterative Approach

Partitioning

L*

$\varepsilon$? ✔

$a$? ✔

$b$? ✔

$bb$? ✔

$bab$? ✔

$baa$? ✘

$bba$? ✘

...

start

1

2

✔$aba$✘

RNN says:

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)
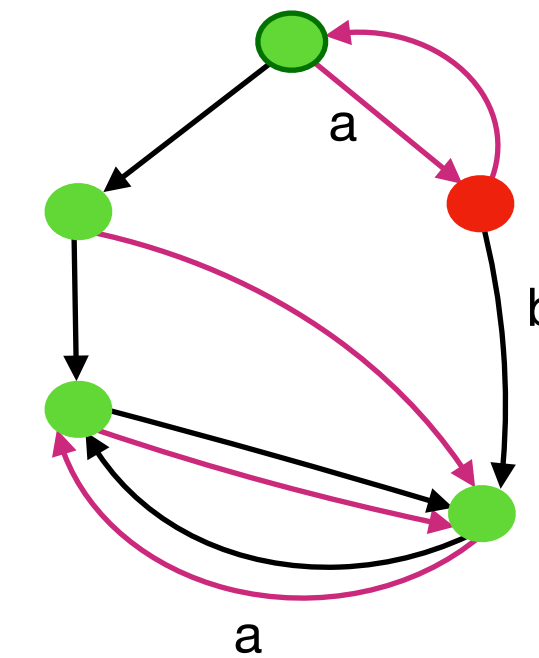
# Iterative Approach

Partitioning

L*

$\varepsilon$? ✔  $bb$? ✔
$a$? ✔  $bab$? ✔
$b$? ✔  $baa$? ✗
       $bba$? ✗
       ...

start  a
  b
  1  b
  a b
  2  a

✔*aba*✗

RNN says: ✗

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# Iterative Approach

Partitioning

L*

$\varepsilon$? ✔  $bb$? ✔
$a$? ✔  $bab$? ✔
$b$? ✔  $baa$? ✘
        $bba$? ✘
        ...

✔$aba$✘

RNN says: ✘

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# Iterative Approach

Partitioning

L*

$\varepsilon$? ✔   $bb$? ✔
$a$? ✔   $bab$? ✔
$b$? ✔   $baa$? ✘
         $bba$? ✘
         ...

start — b — 1 — a,b — 2

✔$aba$✘

RNN says: ✘

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# Iterative Approach



Partitioning

L*

$\varepsilon$? ✔   $bb$? ✔
         $bab$? ✔
$a$? ✔   $baa$? ✘
$b$? ✔   $bba$? ✘
         ...

RNN says:

✔$aba$✘

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# Iterative Approach

Partitioning

L*

$\varepsilon$? ✔   $bb$? ✔
$a$? ✔   $bab$? ✔
$b$? ✔   $baa$? ✘
      $bba$? ✘
      ...

✔$aba$✘

RNN says: ✔

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# Iterative Approach



Partitioning

L*

$\varepsilon$?  ✔  $bb$? ✔
$a$?  ✔  $bab$? ✔
     $baa$? ✘
$b$?  ✔  $bba$? ✘
     ...

✔$aba$✘

RNN says: ✔

$aba$✔

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# Iterative Approach

Partitioning

L*

$\varepsilon$? ✔    $bb$? ✔
          $bab$? ✔
$a$? ✔    $baa$? ✗
$b$? ✔    $bba$? ✗
          ...

✔$aba$✗

RNN says: ✔

$aba$✔

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# Iterative Approach

Partitioning

L*

$\varepsilon$? ✔  $bb$? ✔

$a$? $aa$? $baa$? ✘  $bab$? ✔

$b$? $abab$? $bba$? ✘

$daa$? ✔ ...

$abba$? ✔

$abb$? ✘

...

✔$aba$✘

RNN says: ✔

$aba$✔

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# Iterative Approach



Partitioning

L*

$\varepsilon$? ✔  $bb$? ✔
$bab$? ✔
$a$? $aa$? $baa$? ✘
$abab$? ✘ $bba$? ✘
$b$?
$daa$? ✔ ...
$abba$? ✔
$abb$? ✘
...

**?**

✔$aba$✘

RNN says: ✔

$aba$✔

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# Results

1. Concise, Exact Models in Short Time:

```python
def target(w):
    if len(w)==0:
        return True
    return w[0]==w[-1]
alphabet = "abcd"
```

Training
   (4,400 samples to 100% accuracy)

**RNN**

Extraction
0.2s

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# Results

## 1. Concise, Exact Models in Short Time:

```python
def target(w):
    if len(w)==0:
        return True
    return w[0]==w[-1]
alphabet = "abcd"
```

Training
(4,400 samples to 100% accuracy)

**RNN**

Extraction
0.2s



## 2. **Adversarial Examples** (finding flaws)

Balanced Parentheses GRU
100% train set accuracy
BP up to depth 11, over alphabet: ()a-z

**Counterexamples:**

))                    (1.1s)
(()                   (1.2s)
((()))                (2.1s)
((()))                (3.1s)
(((())))              (3.8s)
((((()))))            (4.4s)
(((((())))))          (6.6s)
((((((()))))))        (9.2s)
(((((((v())))))))     (10.7s)
((((((((a()z))))))))  (8.3s)

**Comparison:
Random sampling
counterexamples:**

))                    (0.4s)
(()i)ma               (32.6s)

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# DFAs from RNNs



**Goal:**
Concise (Meaningful) **Model**
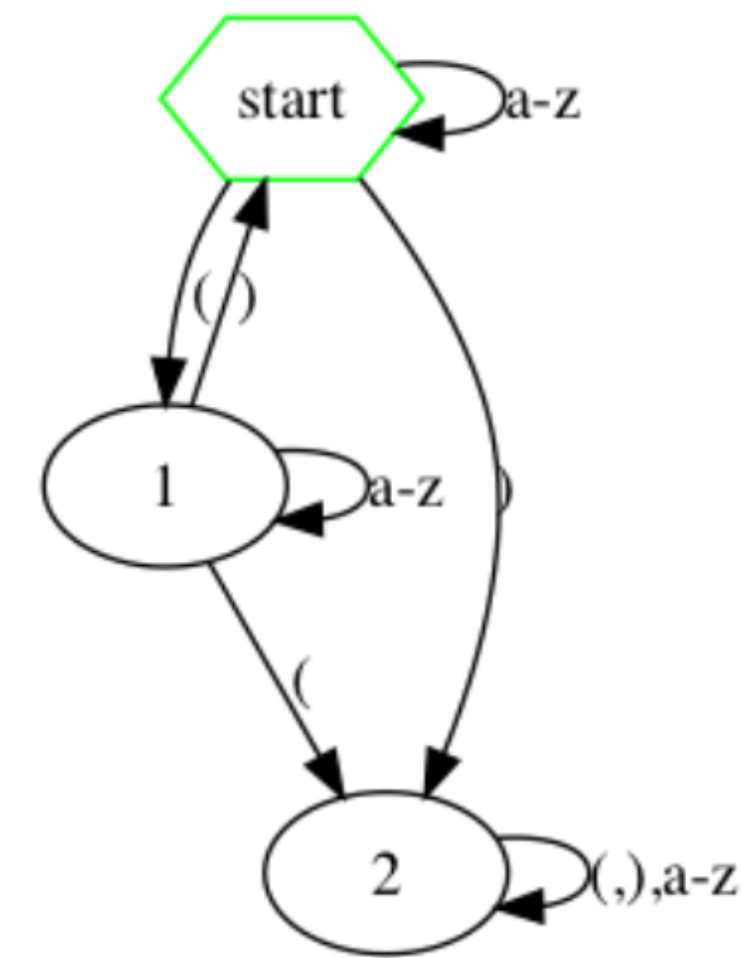from Trained RNN

**DFA**
PDFA
CFG

Extracting Automata from Recurrent Neural Networks Using Queries and Counterexamples (Weiss, Goldberg, Yahav, ICML 2018)

# DFAs from RNNs



**Goal:**
Concise (Meaningful) **Model**
from Trained RNN

DFA
PDFA
**CFG**

Synthesising Context Free Grammars from Recurrent Neural Networks (Yellin, Weiss, TACAS 2021)

# CFGs from RNNs

Observation: L-star learning a CFG seems to have **structured increases** (example on BP)



etc…

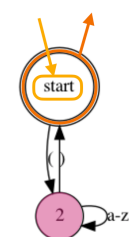Synthesising Context Free Grammars from Recurrent Neural Networks (Yellin, Weiss, TACAS 2021)

# CFGs from RNNs

## Patterns

- Structure
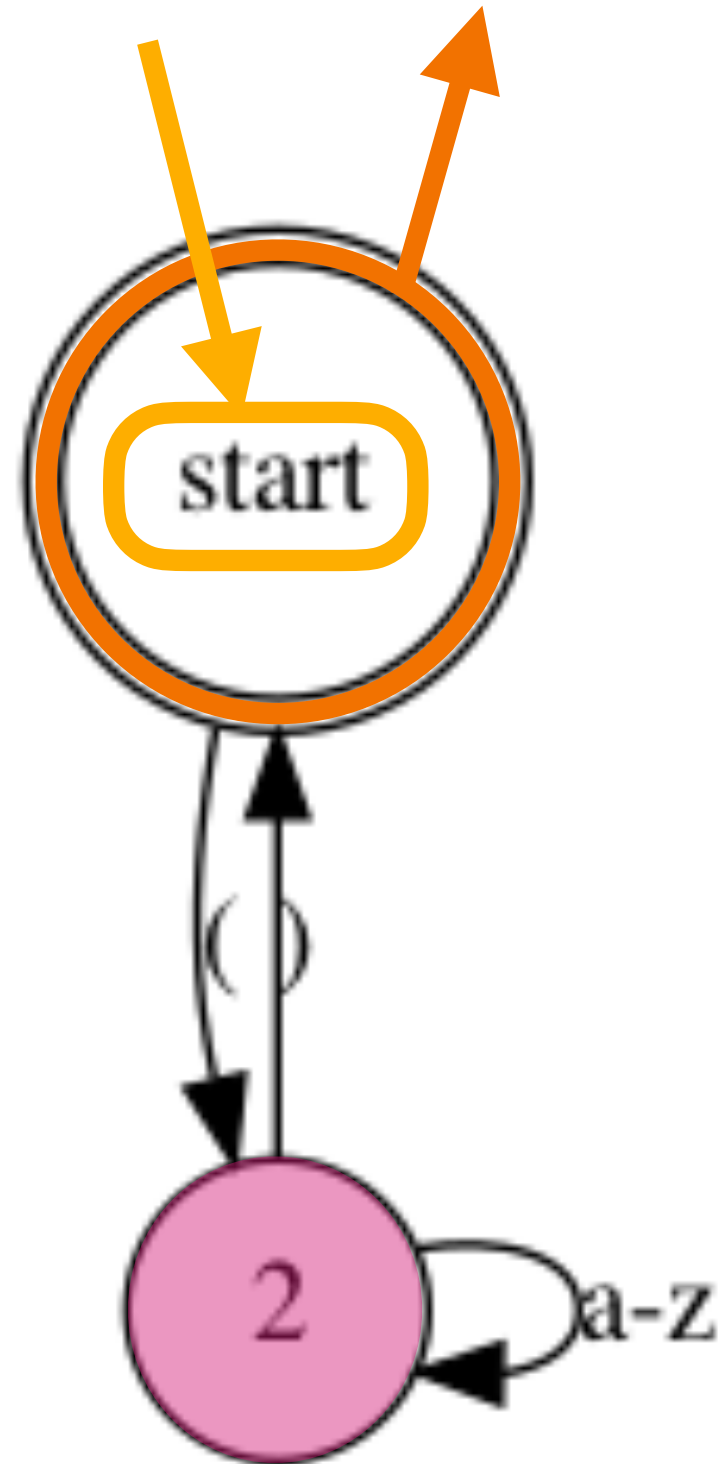  - Entry
  - Exit
- Connection Point(s)
- Composable

Synthesising Context Free Grammars from Recurrent Neural Networks (Yellin, Weiss, TACAS 2021)

# CFGs from RNNs



## Patterns
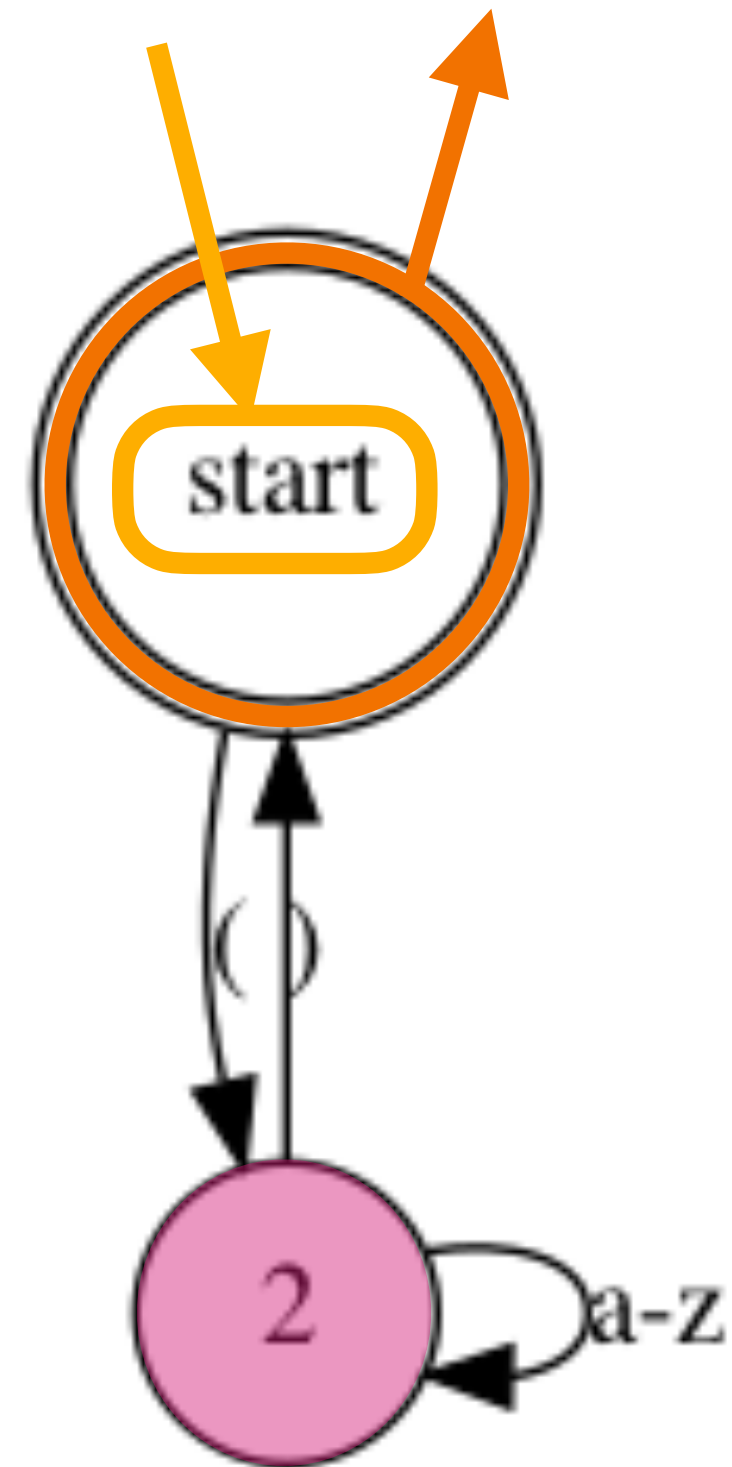
- Structure
  - Entry
  - Exit
- Connection Point(s)
- Composable

## Rules

- Describe legal compositions
  - Legal sequences of DFAs

Synthesising Context Free Grammars from Recurrent Neural Networks (Yellin, Weiss, TACAS 2021)
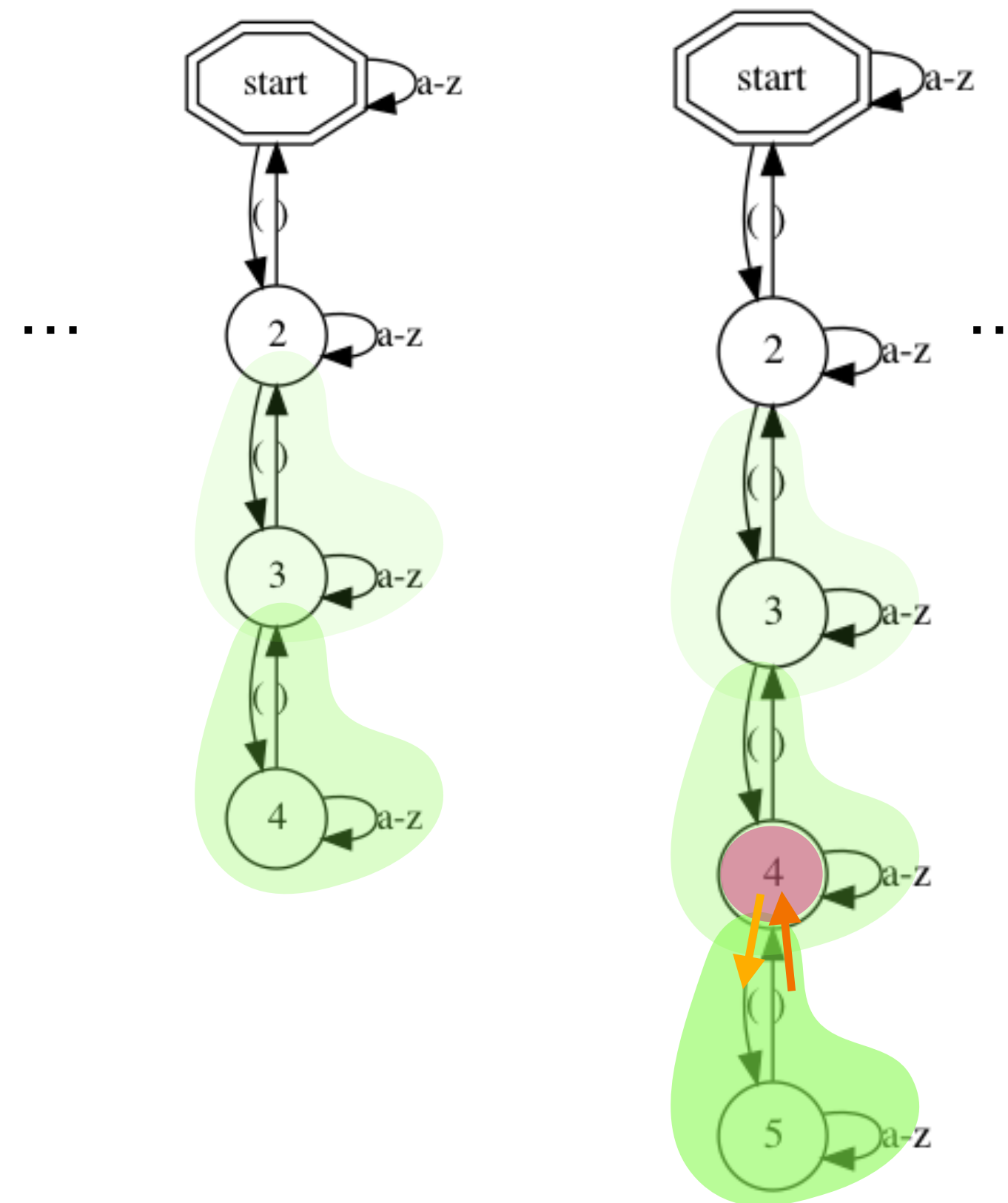
# CFGs from RNNs



## Patterns

- Structure
  - Entry
  - Exit
- Connection Point(s)
- Composable

## Rules

- Describe legal compositions
  - Legal sequences of DFAs

## Result:

Algorithm to recover
*Pattern Rule Sets* from a
sequence of DFAs

Sequence can be obtained
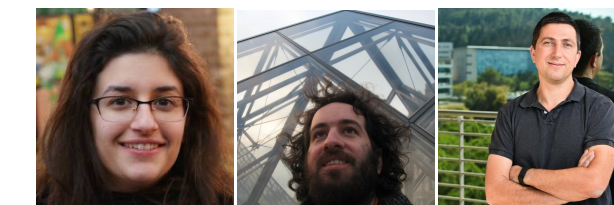from L-star extraction

Some tolerance to noise!
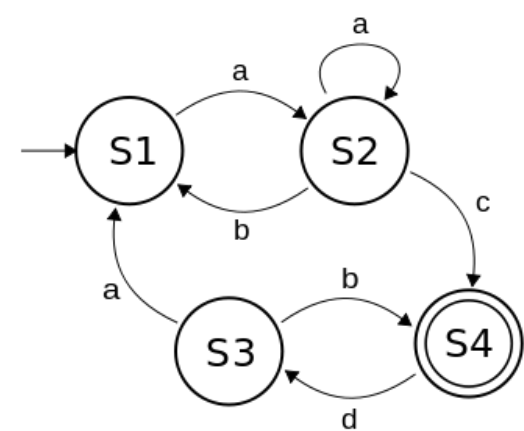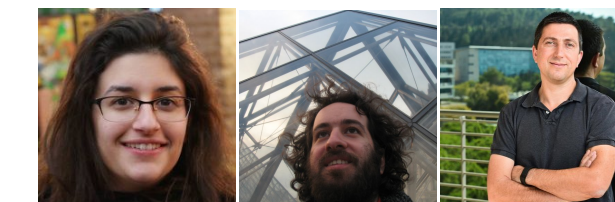
# Neural Sequence Models:
# a Formal Lens

## Counting
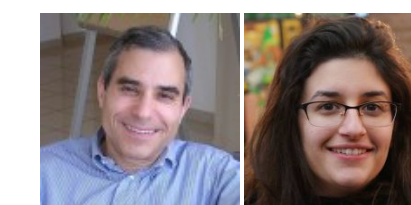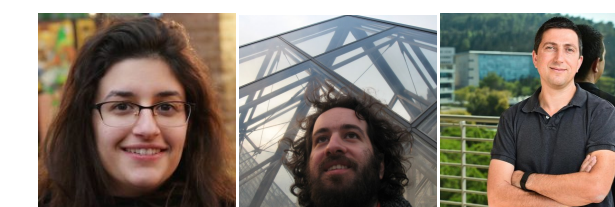LSTMs are counter machines, GRUs aren't **(ACL 2018)**

## RASP
Finding a formalism to describe transformers **(ICML 2021)**

## DFAs from RNNs
Applying L* to learn DFAs from RNNs **(ICML 2018)**

+ using the result for CFGs **(TACAS 2021)**

# Neural Sequence Models: a Formal Lens

## WDFAs from RNNs

Adapting L* to the (noisy!) weighted case (Neurips 2019)

## A Hierarchy of RNNs

Comparing more RNN architectures, with different angles (ACL 2020)

# Thanks!