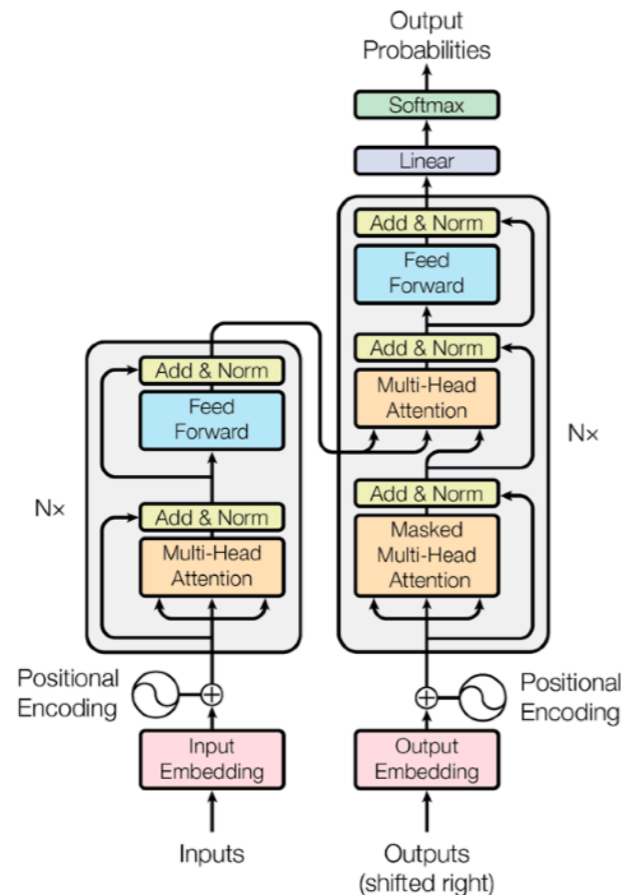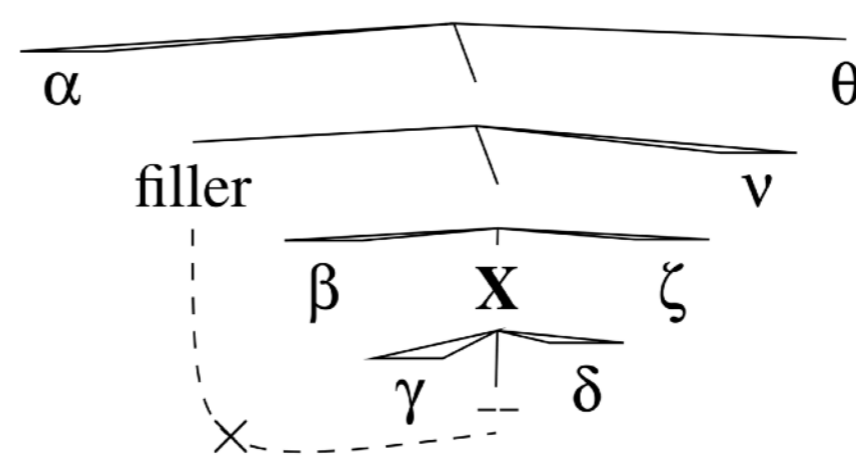# Transformer language models, targeted syntactic evaluation, and learnability



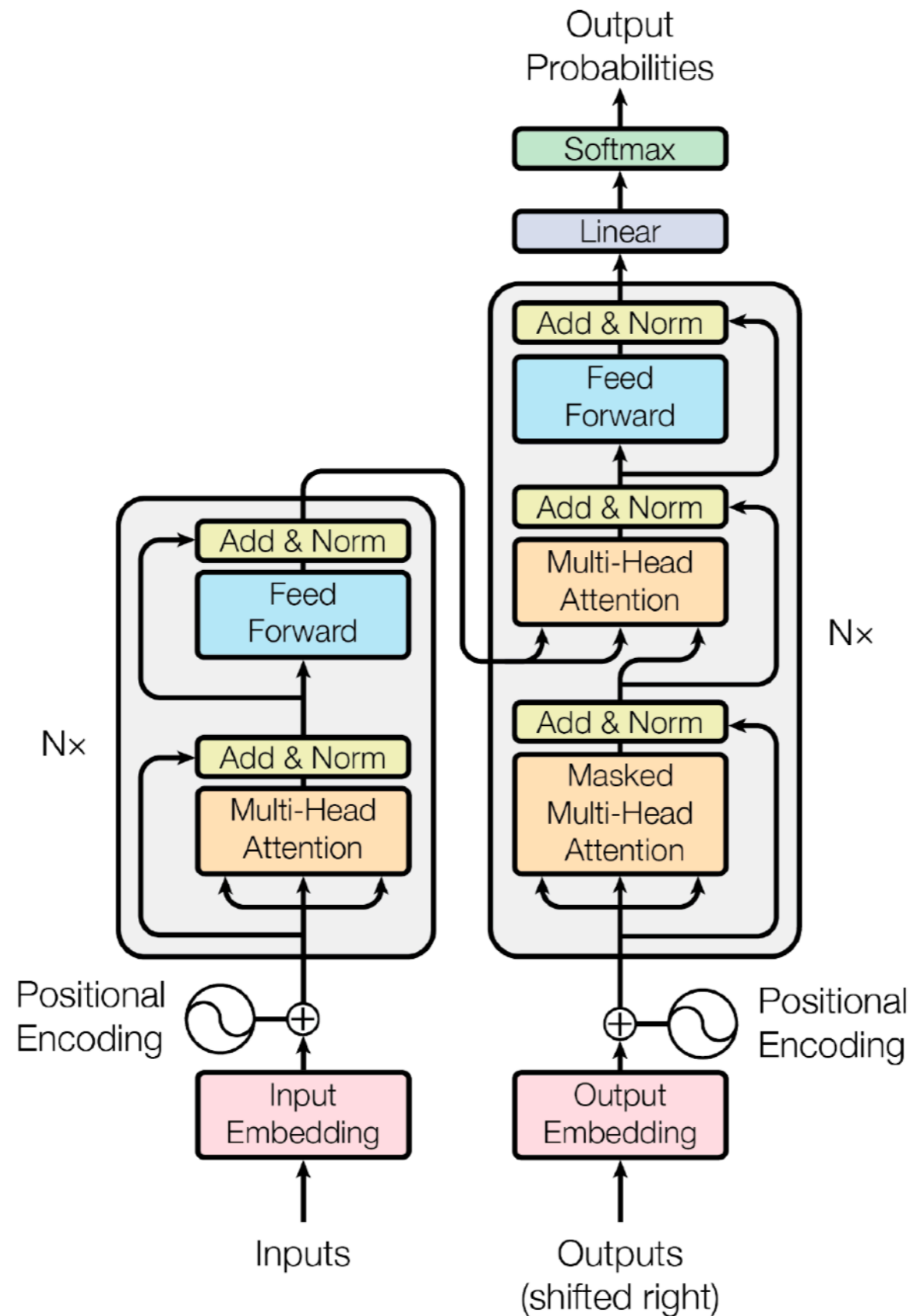*(Vaswani et al., 2017)*

*(Wilcox et al., 2019)*

Roger Levy

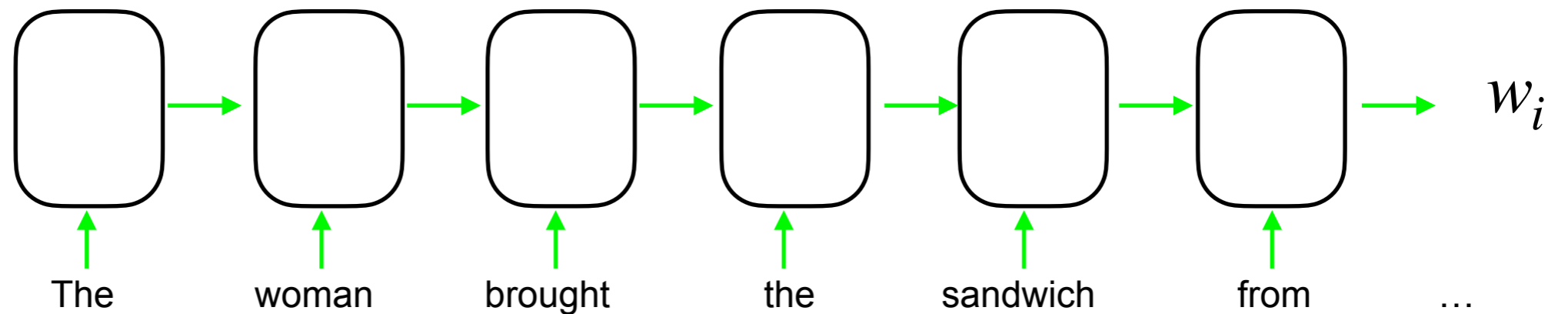9.19: Computational Psycholinguistics

6 November 2023

# Agenda for today

- The Transformer
- Targeted syntactic testing: filler–gap dependencies
- Learnability: syntactic **islands**

# The Transformer model



*(Vaswani et al., 2017)*

# Motivating the Transformer model



The    woman    brought    the    sandwich    from    …

$w_i$

# Motivating the Transformer model

- With RNNs, a fixed-dimension model could propagate information indefinitely into the future...but it's hard!

$w_i$

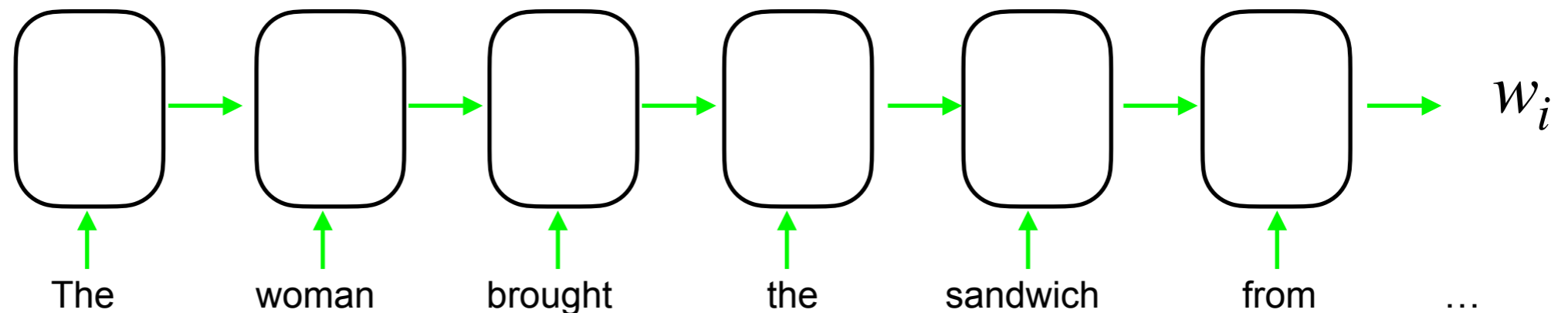The    woman    brought    the    sandwich    from    …

# Motivating the Transformer model

- With RNNs, a fixed-dimension model could propagate information indefinitely into the future...but it's hard!

- We can make RNNs **deep** by stacking them...

$w_i$
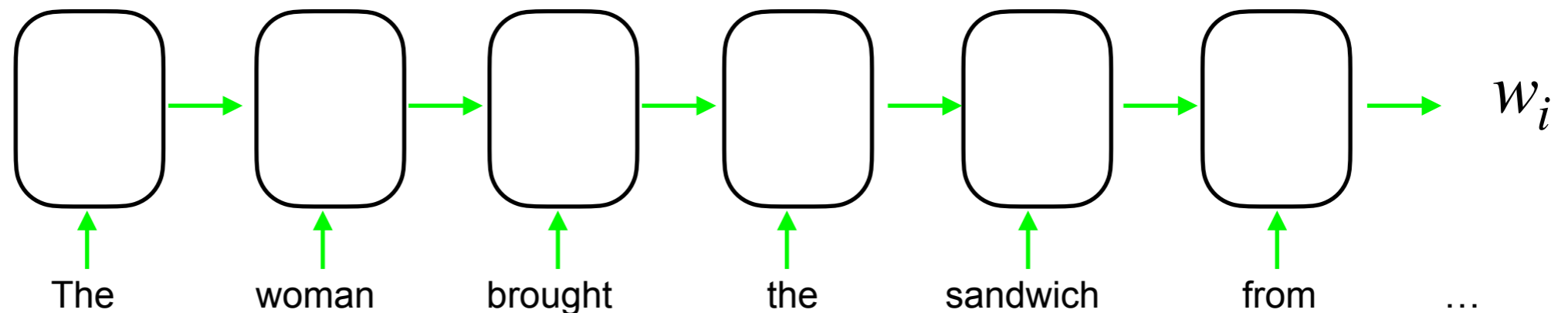
The    woman    brought    the    sandwich    from    ...

# Motivating the Transformer model

- With RNNs, a fixed-dimension model could propagate information indefinitely into the future...but it's hard!

- We can make RNNs *deep* by stacking them...
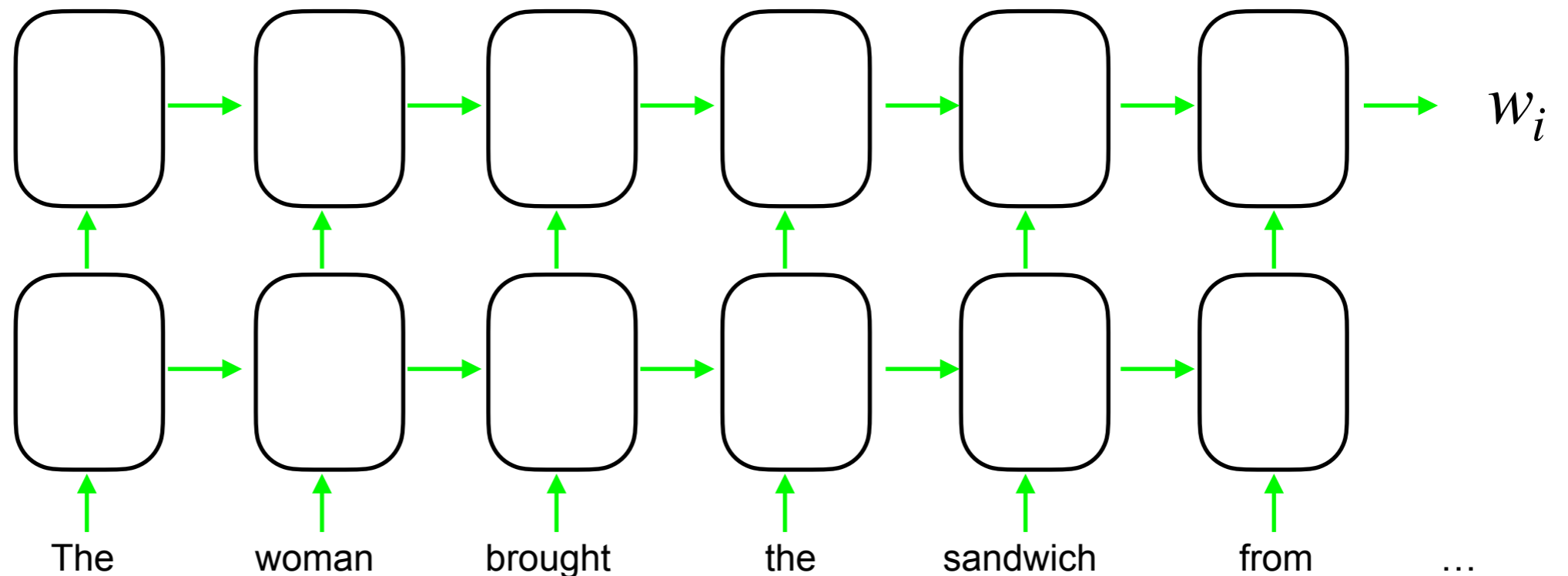
# Motivating the Transformer model

- With RNNs, a fixed-dimension model could propagate information indefinitely into the future...but it's hard!

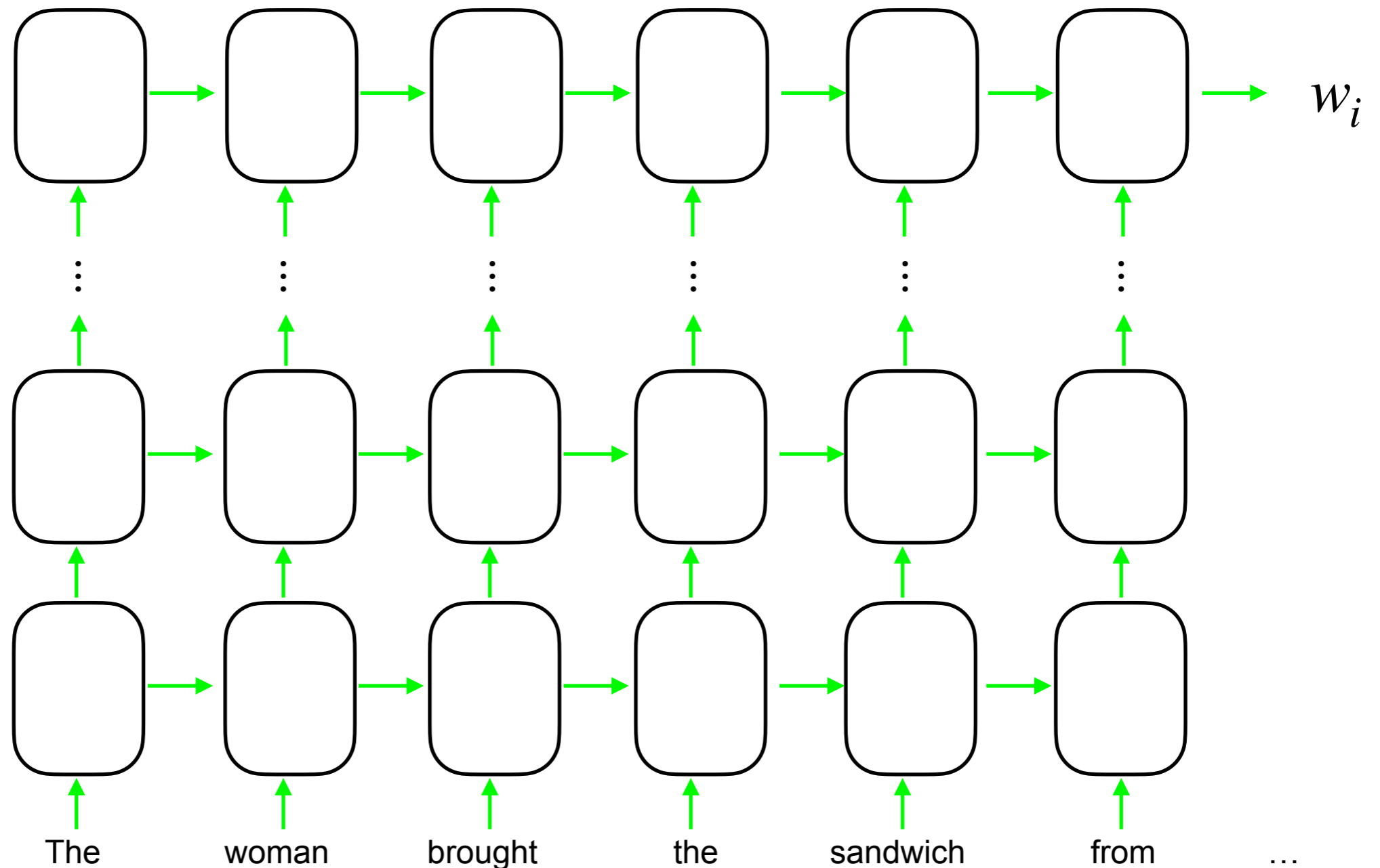- We can make RNNs **deep** by stacking them...
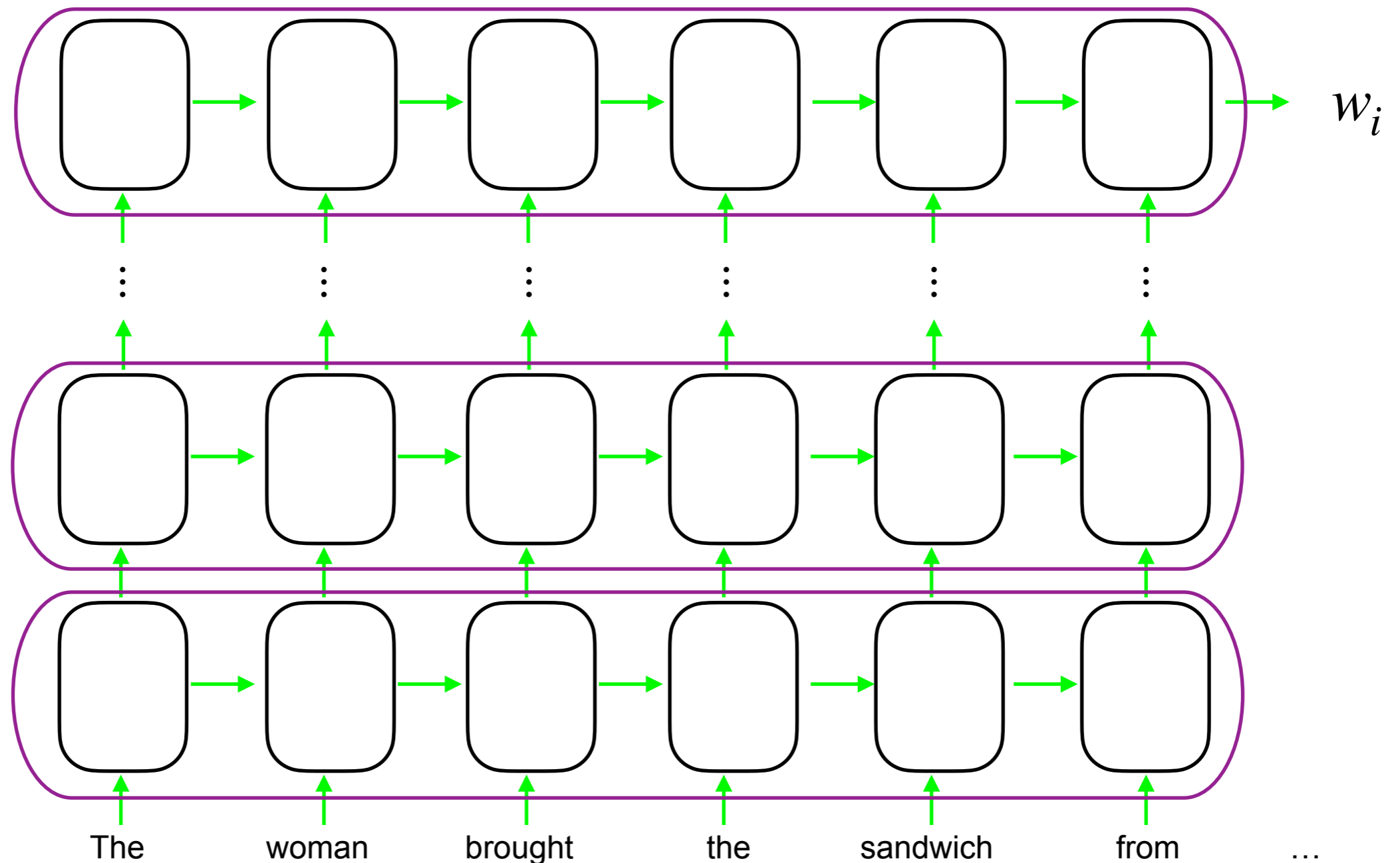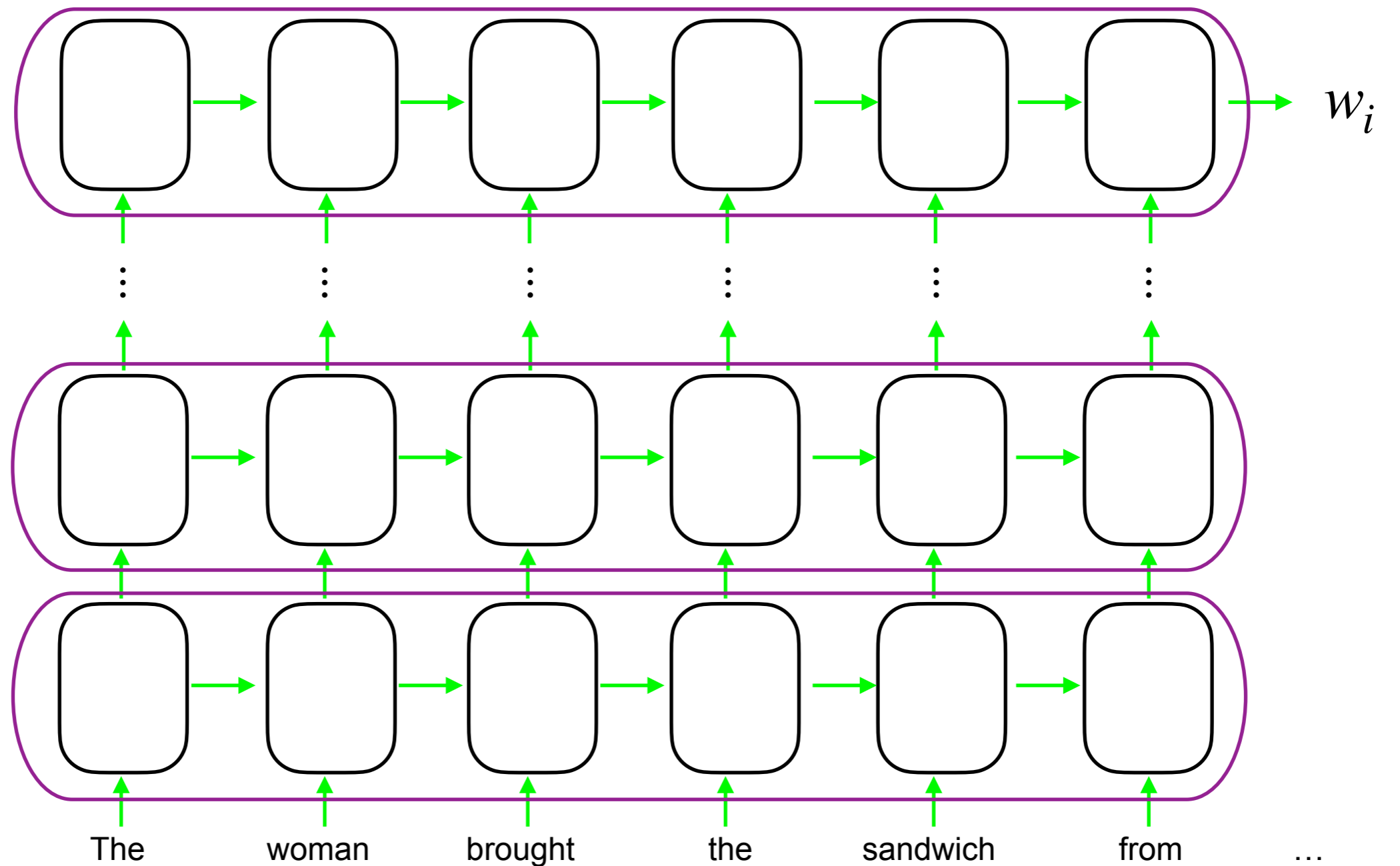
# Motivating the Transformer model

- With RNNs, a fixed-dimension model could propagate information indefinitely into the future...but it's hard!

- We can make RNNs **deep** by stacking them...

# Motivating the Transformer model



The    woman    brought    the    sandwich    from    ...

$w_i$

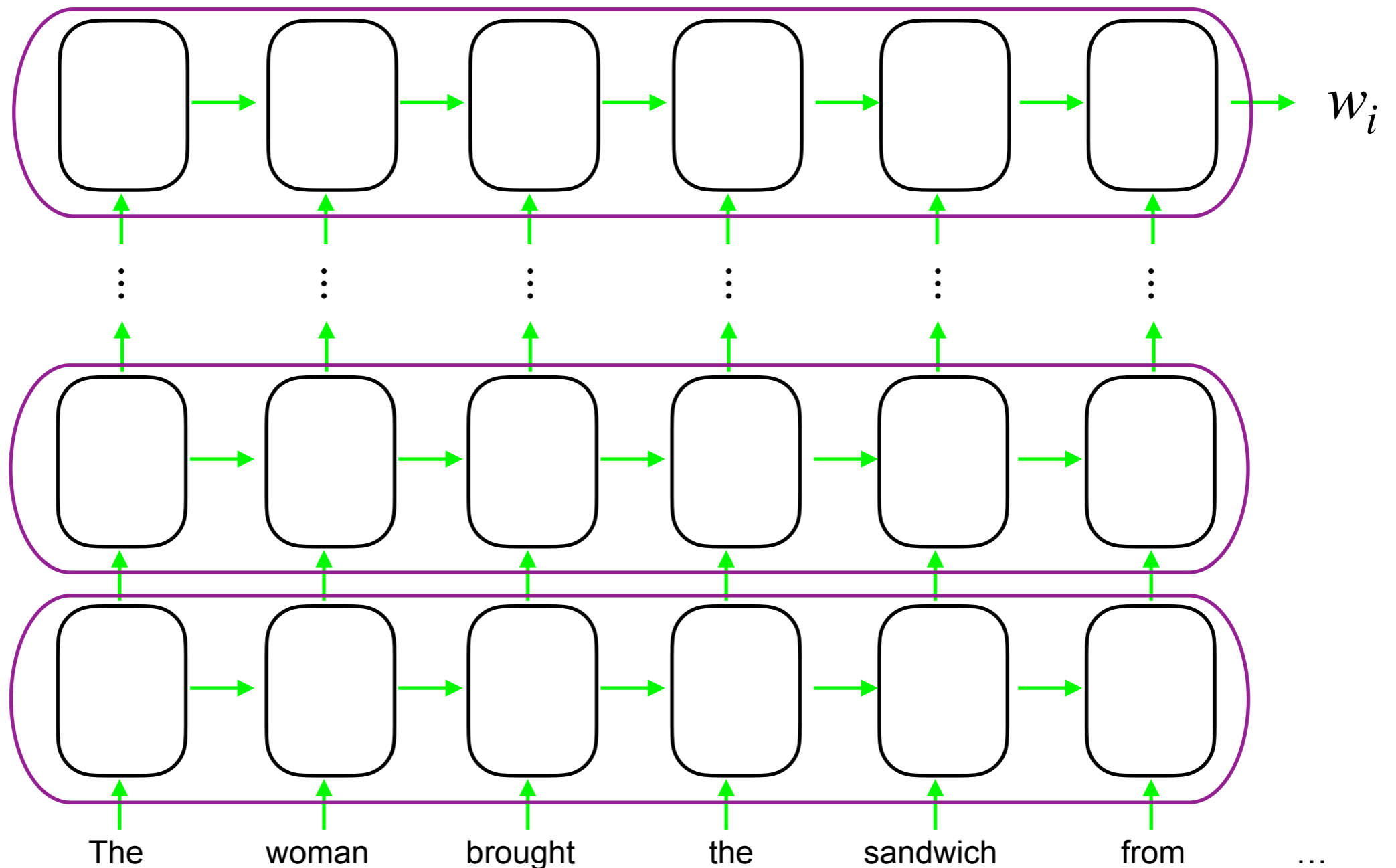# Motivating the Transformer model

- ...but input distant in the context is still far away.

# Motivating the Transformer model

- ...but input distant in the context is still far away.

- Solution: make all context words equally distant from $w_i$!

# Motivating the Transformer model

- ...but input distant in the context is still far away.
- Solution: make all context words equally distant from $w_i$!

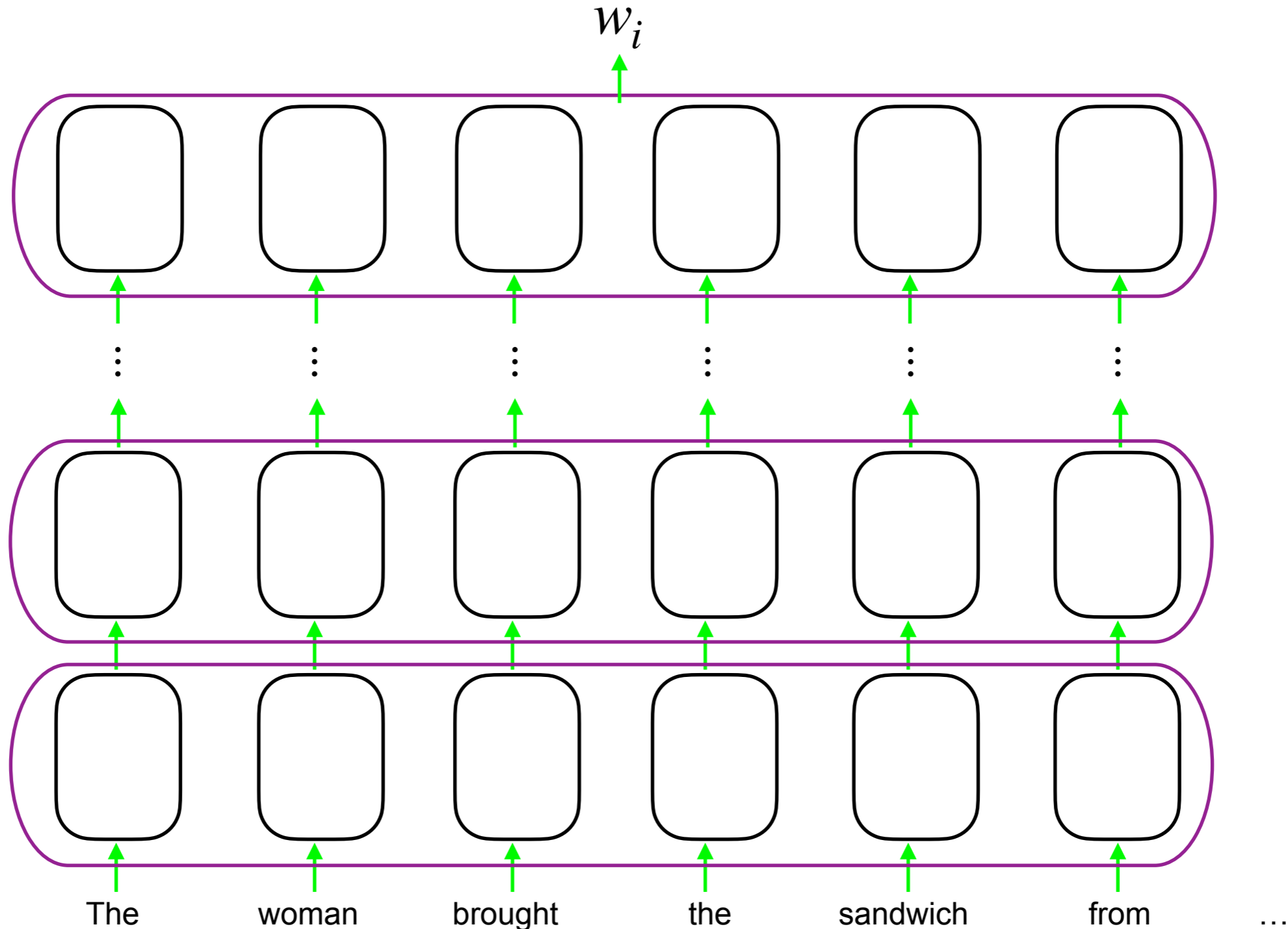# Motivating the Transformer model

- ...but input distant in the context is still far away.

- Solution: make all context words equally distant from $w_i$!

$w_i$

*...but now how do words interact with each other in context?*

The    woman    brought    the    sandwich    from    …

# Input + Positional Embedding

the                    dog                    ate                    the                    ...

# Input + Positional Embedding

**Word embedding matrix:** $d$

$|V|$



the        dog        ate        the        ...

# Input + Positional Embedding

**Word embedding matrix:** $d$

$|V|$

the            dog            ate            the            ...

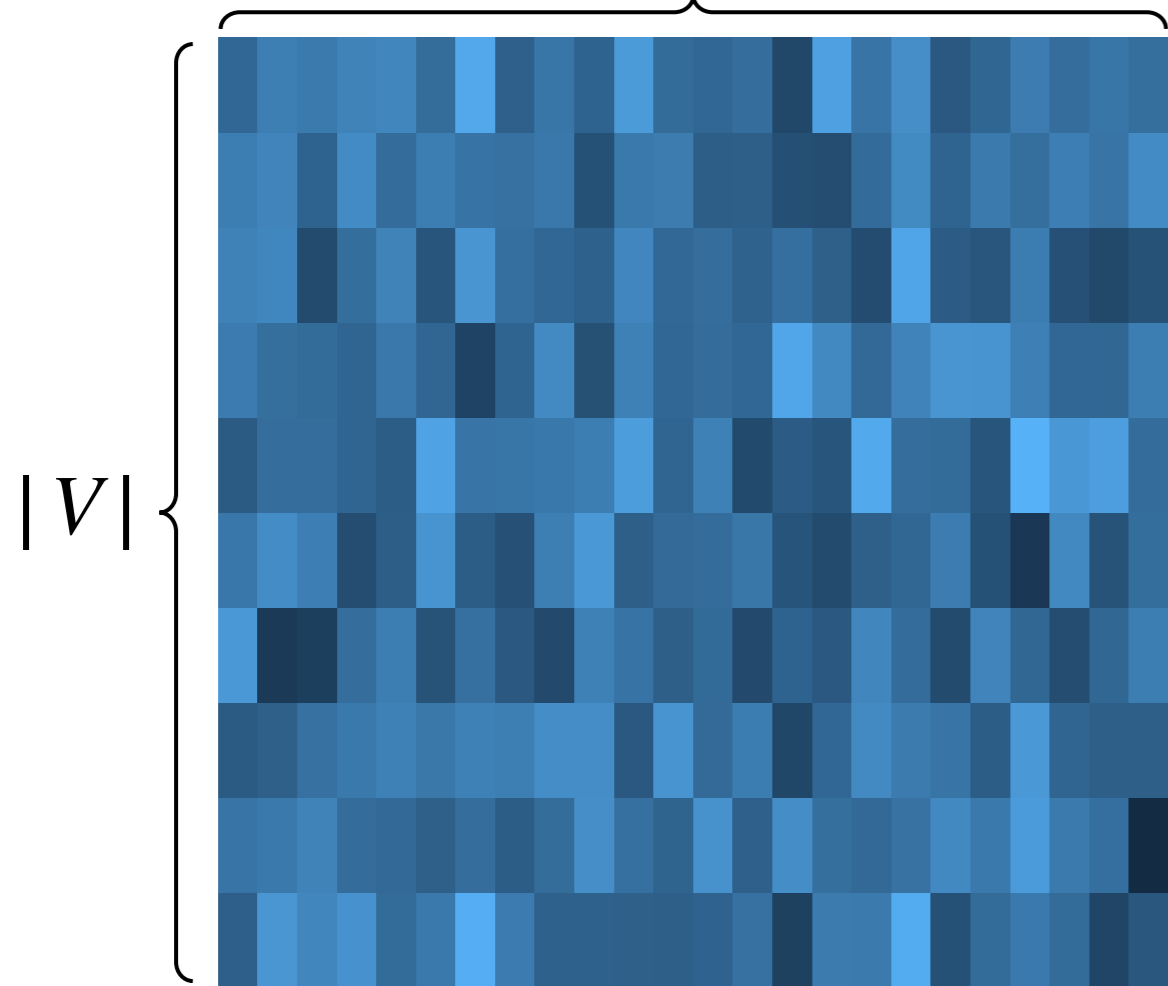# Input + Positional Embedding

**Word embedding matrix:**



$d$

$|V|$

the           dog           ate           the           ...

# Input + Positional Embedding

**Word embedding matrix:** $d$

**Position embedding matrix:** $d$

$|V|$

*word position*

1

20

40

⋮

the        dog        ate        the        ...

# Input + Positional Embedding

**Word embedding matrix:** $d$

**Position embedding matrix:** $d$

$|V|$

*word position*

1

20

40

$\vdots$

$\oplus$

the      dog      ate      the      ...

# Input + Positional Embedding

**Word embedding matrix:** $d$

**Position embedding matrix:** $d$



$|V|$

*word position*

1

20

40

⋮

‖

$\oplus$

the        dog        ate        the        ...

# Input + Positional Embedding

**Word embedding matrix:** $d$

**Position embedding matrix:** $d$



$|V|$

*word position*

1

20

40

$\vdots$

||

$\oplus$

$\oplus$

the                    dog                    ate                    the                    ...
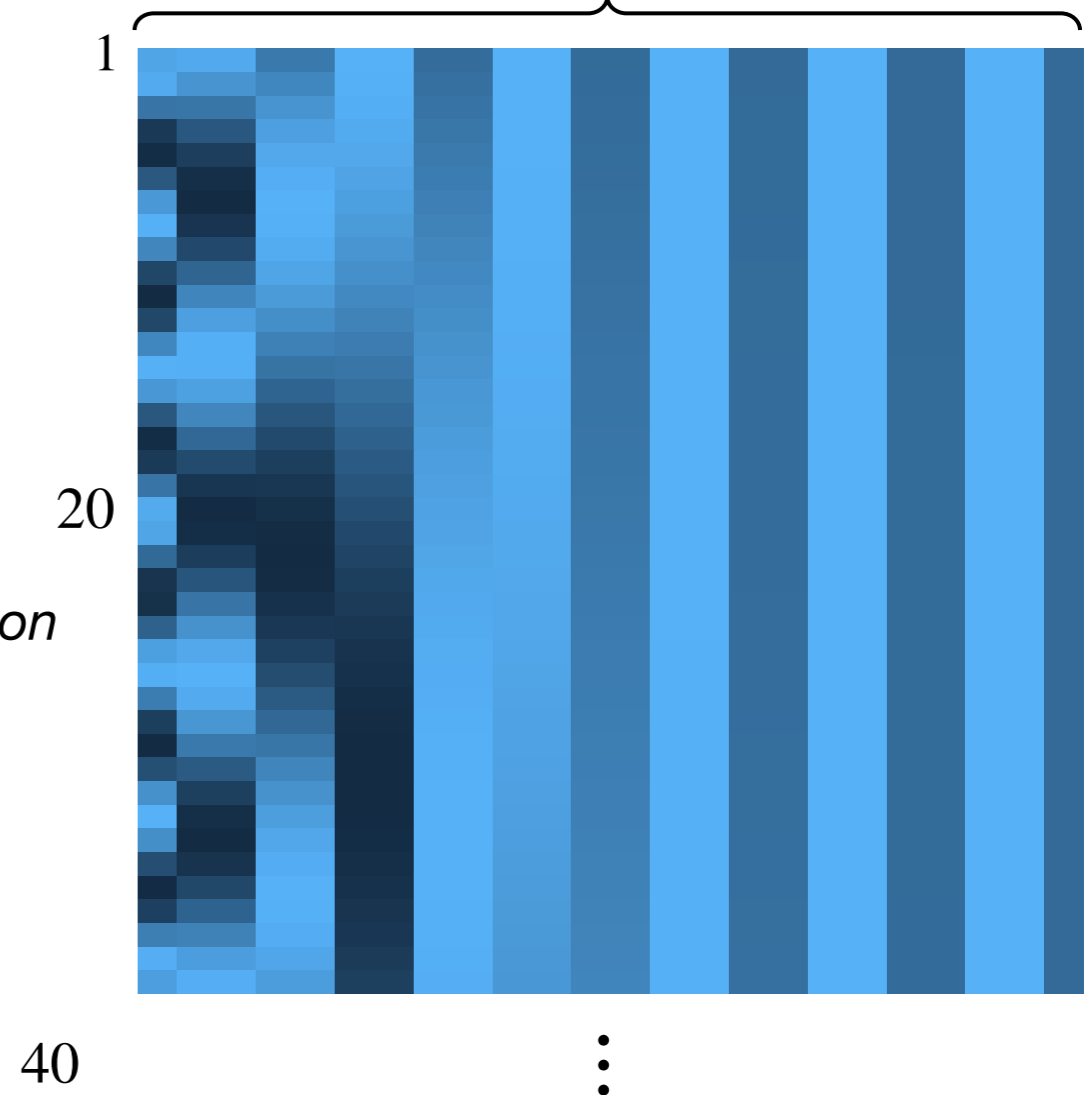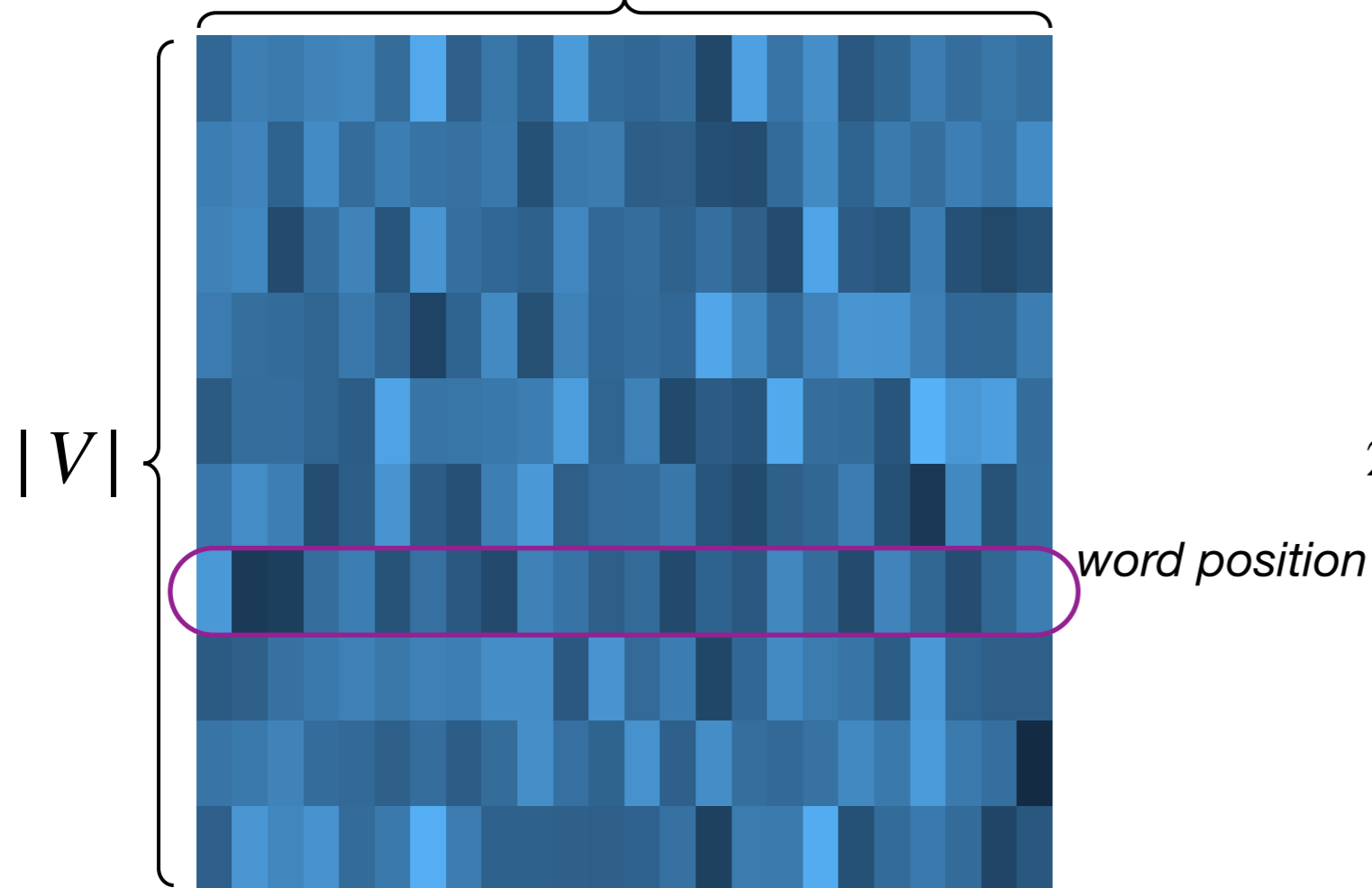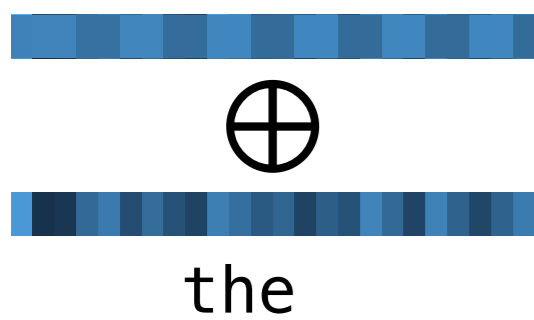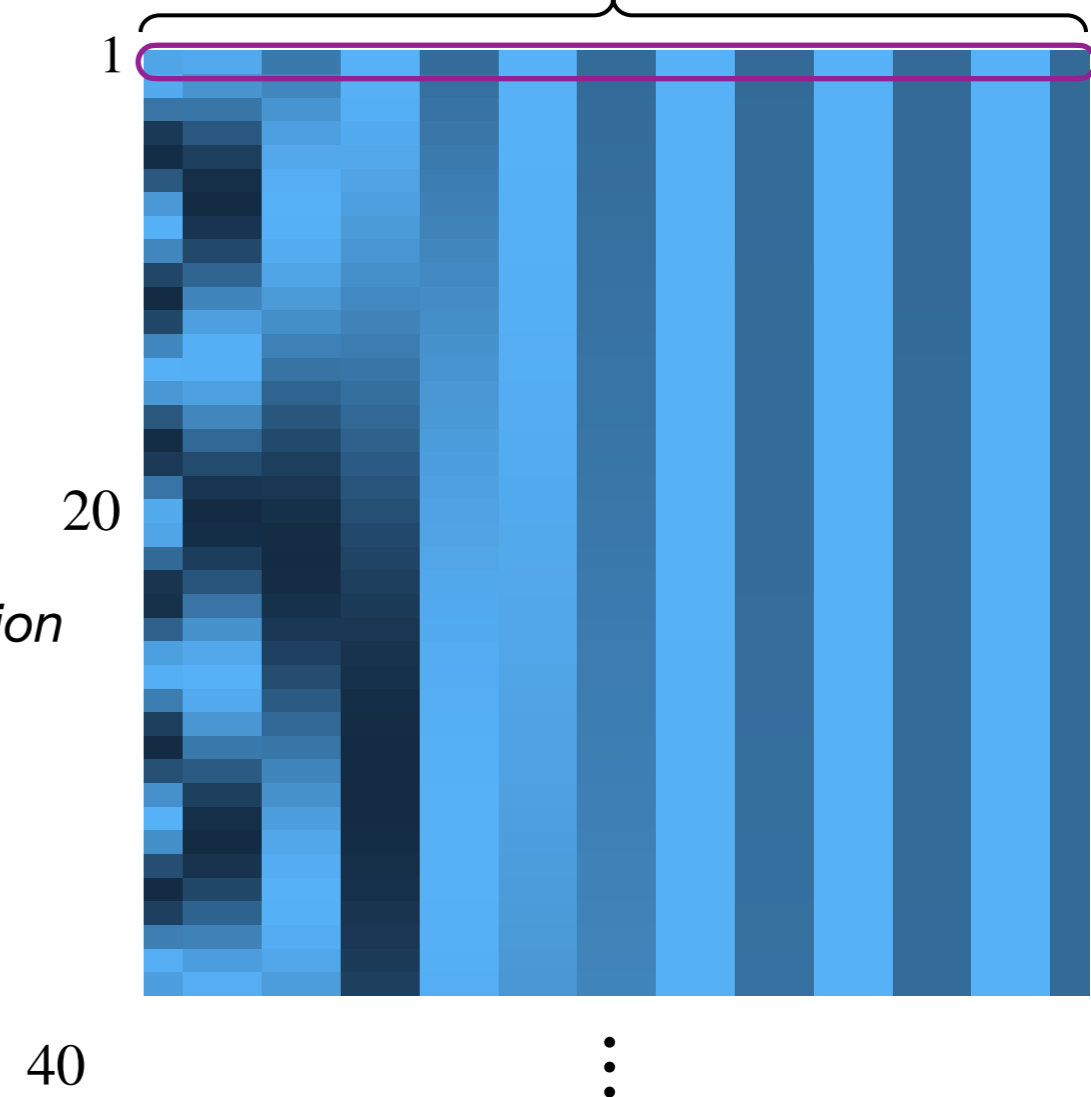
# Input + Positional Embedding

**Word embedding matrix:** $d$

**Position embedding matrix:** $d$



$|V|$

*word position*

1

20

40

$\vdots$

$\|$

$\oplus$

the   dog   ate   the   ...

# The positional embedding function

$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right) \qquad PE(pos, 2i+1) = \cos\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$

$$d = 512$$



word position

# The Transformer unit



$w_i$

The woman brought the sandwich from

*(Vaswani et al., 2017)*

# The Transformer unit



Output $o$

Layer Norm

Feed Forward

$\oplus$

Layer Norm

$\oplus$

Masked Multi Self Attention

Input $x$

$w_i$

The woman brought the sandwich from

*(Vaswani et al., 2017)*

*(Figure from Radford et al., 2018)*

# The Transformer unit



Output $o$

Layer Norm

Feed Forward

Layer Norm

Masked Multi Self Attention

Input $x$

$w_i$

The woman brought the sandwich from

*(Vaswani et al., 2017)*

*(Figure from Radford et al., 2018)*

# Neural Attention

$$V_1 \quad V_2 \quad \cdots \quad V_{i-1}$$

$$K_1 \quad K_2 \quad \cdots \quad K_{i-1}$$

$$Q$$

# Neural Attention

**Q**uery, **K**ey, and **V**alue

# Neural Attention

Attention function options:   $e_j = \begin{cases} v \tanh \left[ W_Q Q + W_K K_j \right] & \textbf{(Bahdanau et al., 2014)} \\[2em] Q^T W K_j & \textbf{(Luong et al., 2015)} \\[2em] \dfrac{Q^T K}{\sqrt{|K|}} & \textbf{(Vaswani et al., 2017)} \end{cases}$

# Neural Attention

*Q**uery, **K**ey, and **V**alue*

$$\alpha_1, \ldots, \alpha_{i-1} = \textbf{softmax}(e_1, \ldots, e_{i-1})$$



Attention function options: $e_j = \begin{cases} v \tanh \left[ W_Q Q + W_K K_j \right] & \textbf{(Bahdanau et al., 2014)} \\ Q^T W K_j & \textbf{(Luong et al., 2015)} \\ \dfrac{Q^T K}{\sqrt{|K|}} & \textbf{(Vaswani et al., 2017)} \end{cases}$

# Neural Attention

$$\alpha_1, \ldots, \alpha_{i-1} = \textbf{softmax}(e_1, \ldots, e_{i-1})$$

$$\alpha_j = \frac{e^{e_j}}{\sum_{j'=1}^{i-1} e^{e_{j'}}}$$

$V_1$ $V_2$ $\cdots$ $V_{i-1}$

$K_1$ $K_2$ $\cdots$ $K_{i-1}$

$e_2$

$e_{i-1}$

$\cdots$

$e_1$

$Q$

Attention function options: $e_j = \begin{cases} v \tanh \left[ W_Q Q + W_K K_j \right] & \textbf{(Bahdanau et al., 2014)} \\ Q^T W K_j & \textbf{(Luong et al., 2015)} \\ \dfrac{Q^T K}{\sqrt{|K|}} & \textbf{(Vaswani et al., 2017)} \end{cases}$

# Neural Attention

$$\alpha_1, \dots, \alpha_{i-1} = \textbf{softmax}(e_1, \dots, e_{i-1})$$

$$\alpha_j = \frac{e^{e_j}}{\sum_{j'=1}^{i-1} e^{e_{j'}}}$$

Attention function options:
$$e_j = \begin{cases} v \tanh \left[ W_Q Q + W_K K_j \right] & \textbf{(Bahdanau et al., 2014)} \\[2ex] Q^T W K_j & \textbf{(Luong et al., 2015)} \\[2ex] \dfrac{Q^T K}{\sqrt{|K|}} & \textbf{(Vaswani et al., 2017)} \end{cases}$$

# Neural Attention

$$\textbf{Output} = \sum_{j=1}^{i-1} \alpha_j V_j$$

$$\alpha_1, \ldots, \alpha_{i-1} = \textbf{softmax}(e_1, \ldots, e_{i-1})$$

$$\alpha_j = \frac{e^{e_j}}{\sum_{j'=1}^{i-1} e^{e_{j'}}}$$

$V_1$ $V_2$ $\cdots$ $V_{i-1}$

$K_1$ $K_2$ $\cdots$ $K_{i-1}$

$\alpha_1$ $\alpha_2$ $\cdots$ $\alpha_{i-1}$

$e_2$ $e_{i-1}$ $e_1$ $\cdots$

$Q$

Attention function options: $e_j = \begin{cases} v \tanh \left[ W_Q Q + W_K K_j \right] & \textbf{(Bahdanau et al., 2014)} \\ Q^T W K_j & \textbf{(Luong et al., 2015)} \\ \dfrac{Q^T K}{\sqrt{|K|}} & \textbf{(Vaswani et al., 2017)} \end{cases}$
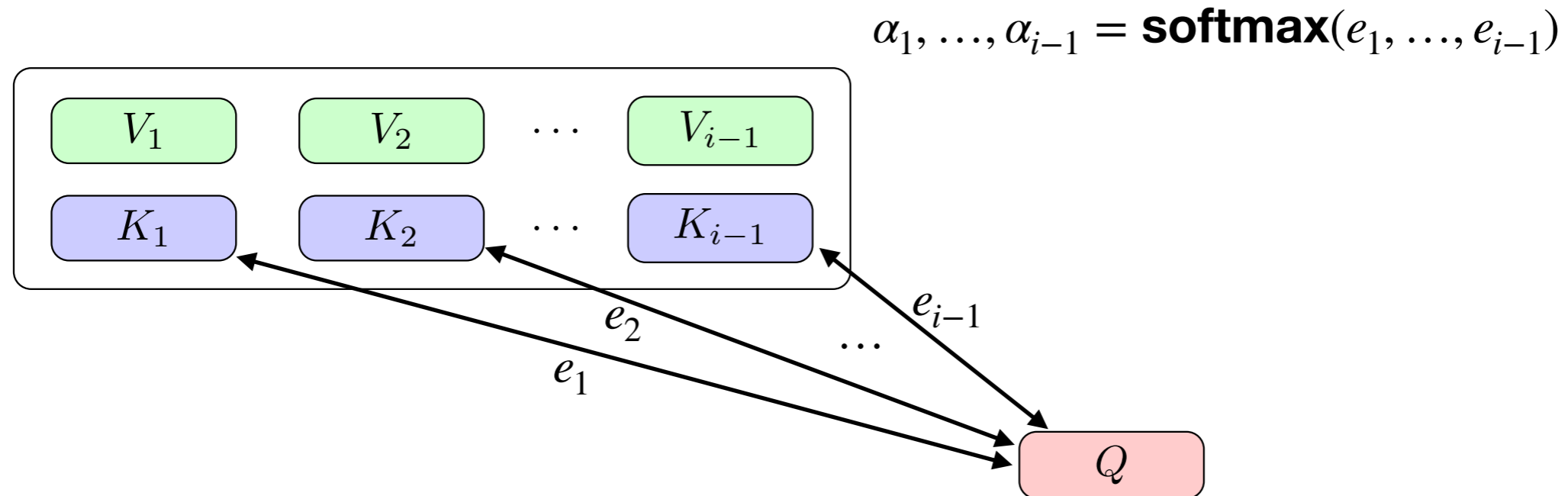
# A single masked attention "head"

$x_1$

$x_2$

$x_3$

$w_1$

$w_2$

$w_3$

...

$W_K$, $W_V$, and $W_Q$ are all learned during training

(Vaswani et al., 2017)

# A single masked attention "head"



$K_1$

$= W_K x_1$

$x_1$

$K_2$

$= W_K x_2$

$x_2$

$K_3$

$= W_K x_3$

$x_3$

$w_1$ $w_2$ $w_3$ $\ldots$

$W_K$, $W_V$, and $W_Q$ are all learned during training

*(Vaswani et al., 2017)*

# A single masked attention "head"

$V_1$

$= W_V x_1$

$V_2$

$= W_V x_2$

$V_3$

$= W_V x_3$

$K_1$

$= W_K x_1$

$K_2$

$= W_K x_2$

$K_3$

$= W_K x_3$

$x_1$

$x_2$

$x_3$

$w_1$  $w_2$  $w_3$  $\ldots$

$W_K$, $W_V$, and $W_Q$ are all learned during training

*(Vaswani et al., 2017)*

# A single masked attention "head"

$V_1$    $= W_V x_1$

$V_2$    $= W_V x_2$

$V_3$    $= W_V x_3$

$Q_3$

$= W_Q x_3$

$K_1$

$= W_K x_1$

$K_2$

$= W_K x_2$

$K_3$

$= W_K x_3$

$x_1$

$x_2$

$x_3$

$w_1$      $w_2$      $w_3$      $\dots$

$W_K$, $W_V$, and $W_Q$ are all learned during training

*(Vaswani et al., 2017)*

# A single masked attention "head"



$V_1 = W_V x_1$

$V_2 = W_V x_2$

$V_3 = W_V x_3$

$Q_3 = W_Q x_3$

$K_1 = W_K x_1$

$K_2 = W_K x_2$

$K_3 = W_K x_3$

$\alpha_1$

$\alpha_2$

$\alpha_3$

$w_1$

$w_2$

$w_3$

$\cdots$

$W_K$, $W_V$, and $W_Q$ are all learned during training

*(Vaswani et al., 2017)*

# A single masked attention "head"



$$o_3 = \sum_{j=1}^{3} \alpha_j V_j$$

$V_1 = W_V x_1$

$V_2 = W_V x_2$

$V_3 = W_V x_3$

$Q_3 = W_Q x_3$

$\alpha_1$

$\alpha_2$

$\alpha_3$

$K_1 = W_K x_1$

$K_2 = W_K x_2$

$K_3 = W_K x_3$

$w_1$

$w_2$

$w_3$

$\cdots$

$W_K$, $W_V$, and $W_Q$ are all learned during training

*(Vaswani et al., 2017)*

# A single masked attention "head"



$$o_3 = \sum_{j=1}^{3} \alpha_j V_j$$

$o_2$  **?**  $o_3$

$V_1 = W_V x_1$   $V_2 = W_V x_2$   $V_3 = W_V x_3$   $Q_3 = W_Q x_3$

$\alpha_1$   $\alpha_2$   $\alpha_3$

$K_1 = W_K x_1$   $K_2 = W_K x_2$   $K_3 = W_K x_3$

$x_1$   $x_2$   $x_3$

$w_1$   $w_2$   $w_3$   $\ldots$

*(Vaswani et al., 2017)*

# A single masked attention "head"

Subsequent context words are **masked** from attention

**?**

$o_2$

$o_3$

$o_3 = \sum\limits_{j=1}^{3} \alpha_j V_j$

$V_1$    $= W_V x_1$

$V_2$    $= W_V x_2$

$V_3$    $= W_V x_3$

$Q_3$    $= W_Q x_3$

$\alpha_1$

$\alpha_2$

$\alpha_3$

$K_1$   $= W_K x_1$

$K_2$   $= W_K x_2$

$K_3$   $= W_K x_3$

$x_1$

$x_2$

$x_3$

$w_1$

$w_2$

$w_3$

$\cdots$

*(Vaswani et al., 2017)*

# A single masked attention "head"

Subsequent context words are *masked* from attention

**?** $o_2$

$V_1$

$= W_V x_1$

$V_2$

$= W_V x_2$

$K_1$

$= W_K x_1$

$K_2$

$= W_K x_2$

$x_1$

$x_2$

$x_3$

$w_1$

$w_2$

$w_3$

$\ldots$

*(Vaswani et al., 2017)*

# A single masked attention "head"

$o_2$

$V_1$

$= W_V x_1$

$V_2$

$= W_V x_2$

$K_1$

$= W_K x_1$

$K_2$

$= W_K x_2$

$x_1$

$x_2$

$x_3$

$w_1$

$w_2$

$w_3$

$\cdots$

*(Vaswani et al., 2017)*

# A single masked attention "head"



*(Vaswani et al., 2017)*

# A single masked attention "head"



$o_2$

$V_1$    $= W_V x_1$      $V_2$    $= W_V x_2$      $Q_2$

$\alpha_1$

$\alpha_2$

$K_1$

$= W_K x_1$      $K_2$

$= W_K x_2$

$x_1$          $x_2$          $x_3$

$w_1$          $w_2$          $w_3$      $\ldots$

*(Vaswani et al., 2017)*

# A single masked attention "head"



$$o_2 = \sum_{j=1}^{2} \alpha_j V_j$$

$o_2$

$V_1$        $V_2$        $Q_2$

$= W_V x_1$        $= W_V x_2$

$\alpha_1$        $\alpha_2$

$K_1$        $K_2$

$= W_K x_1$        $= W_K x_2$

$x_1$        $x_2$        $x_3$

$w_1$        $w_2$        $w_3$        $\dots$

*(Vaswani et al., 2017)*

# Multi-headed attention

The $h$ output vectors are concatenated and then linearly transformed with learned matrix $W_O$

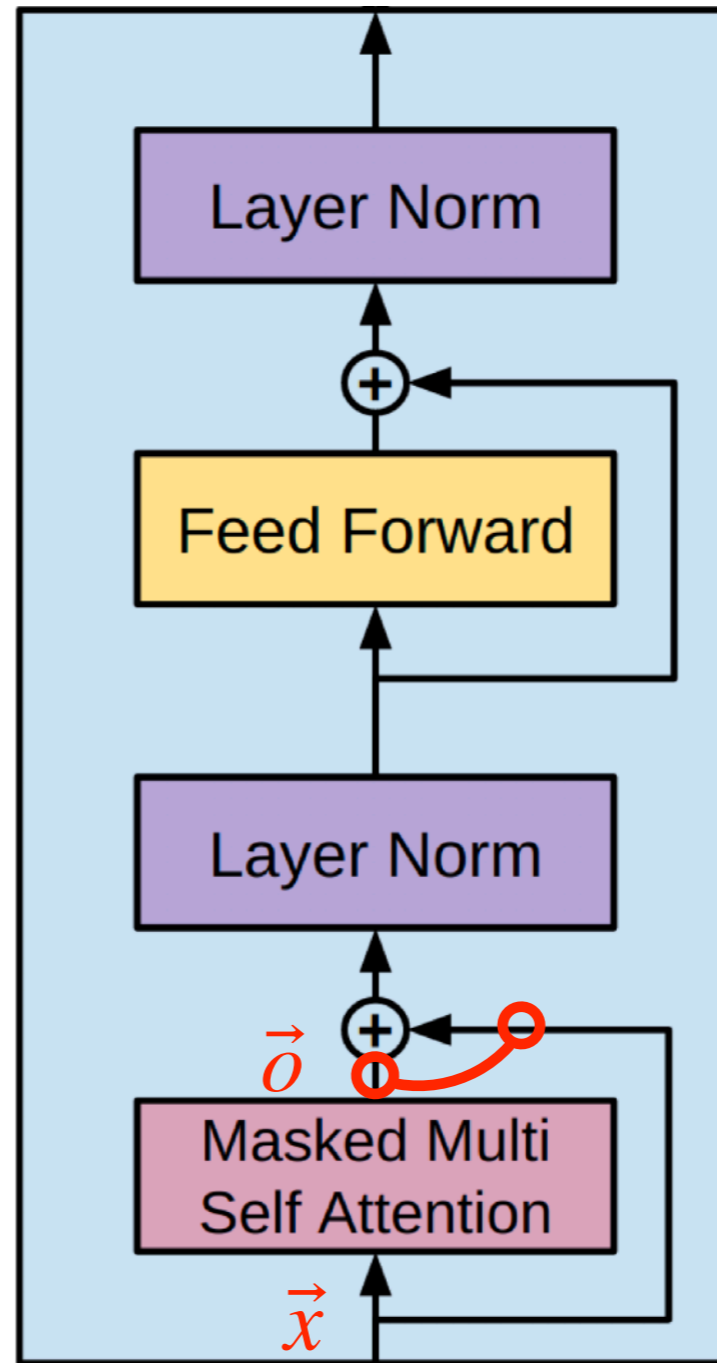$h$ heads, each with its own $W_K$, $W_V$, $W_Q$ matrices

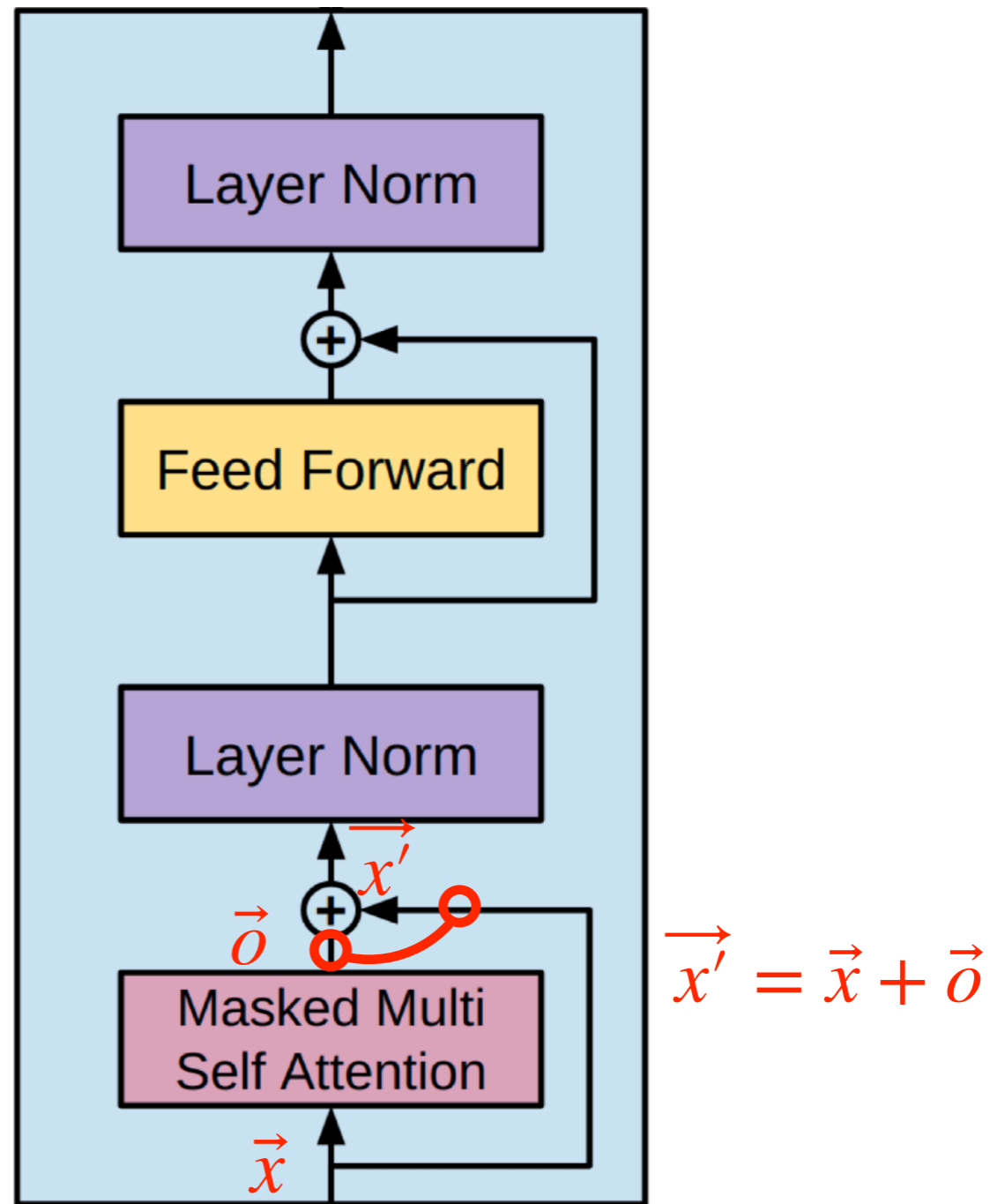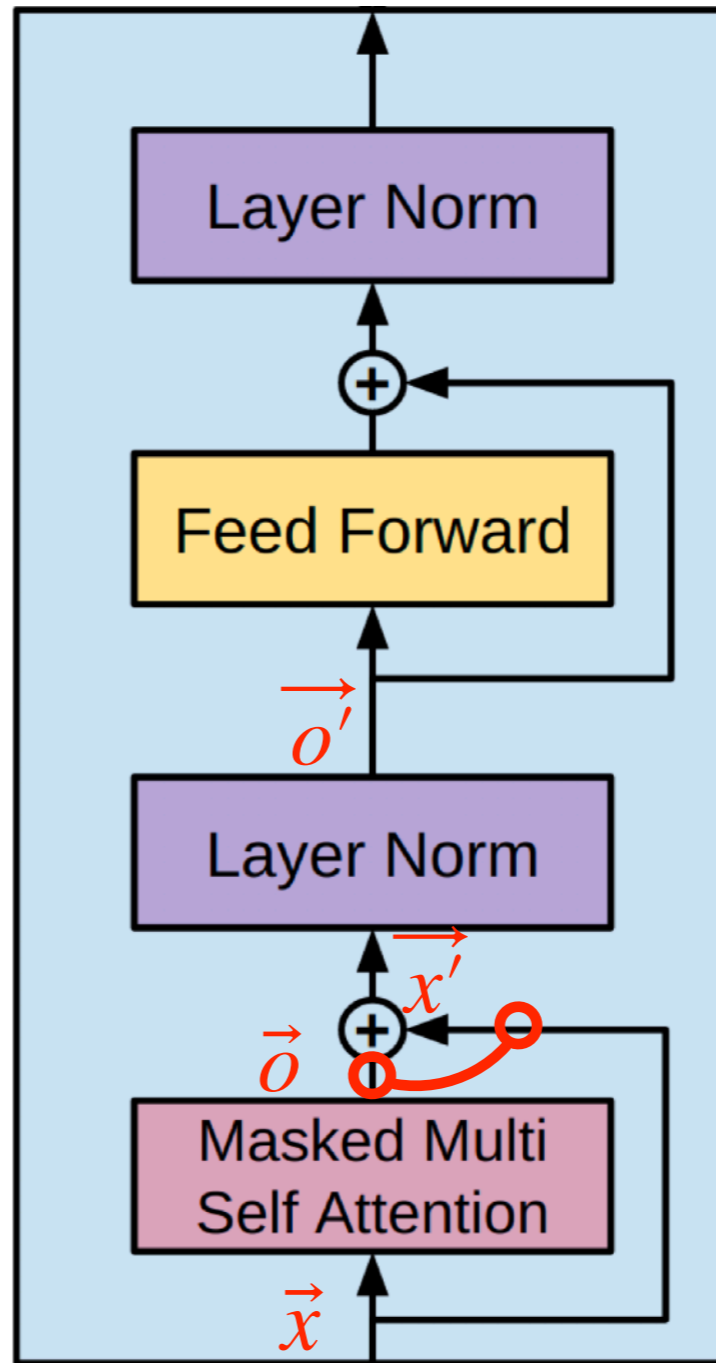# Residual connection & layer normalization

# Residual connection & layer normalization

# Residual connection & layer normalization



$$\vec{x'} = \vec{x} + \vec{o}$$

# Residual connection & layer normalization



$$\overrightarrow{o'} = \frac{\overrightarrow{x'} - \textbf{Mean}(\overrightarrow{x'})}{\textbf{StdDev}(\overrightarrow{x'})}$$

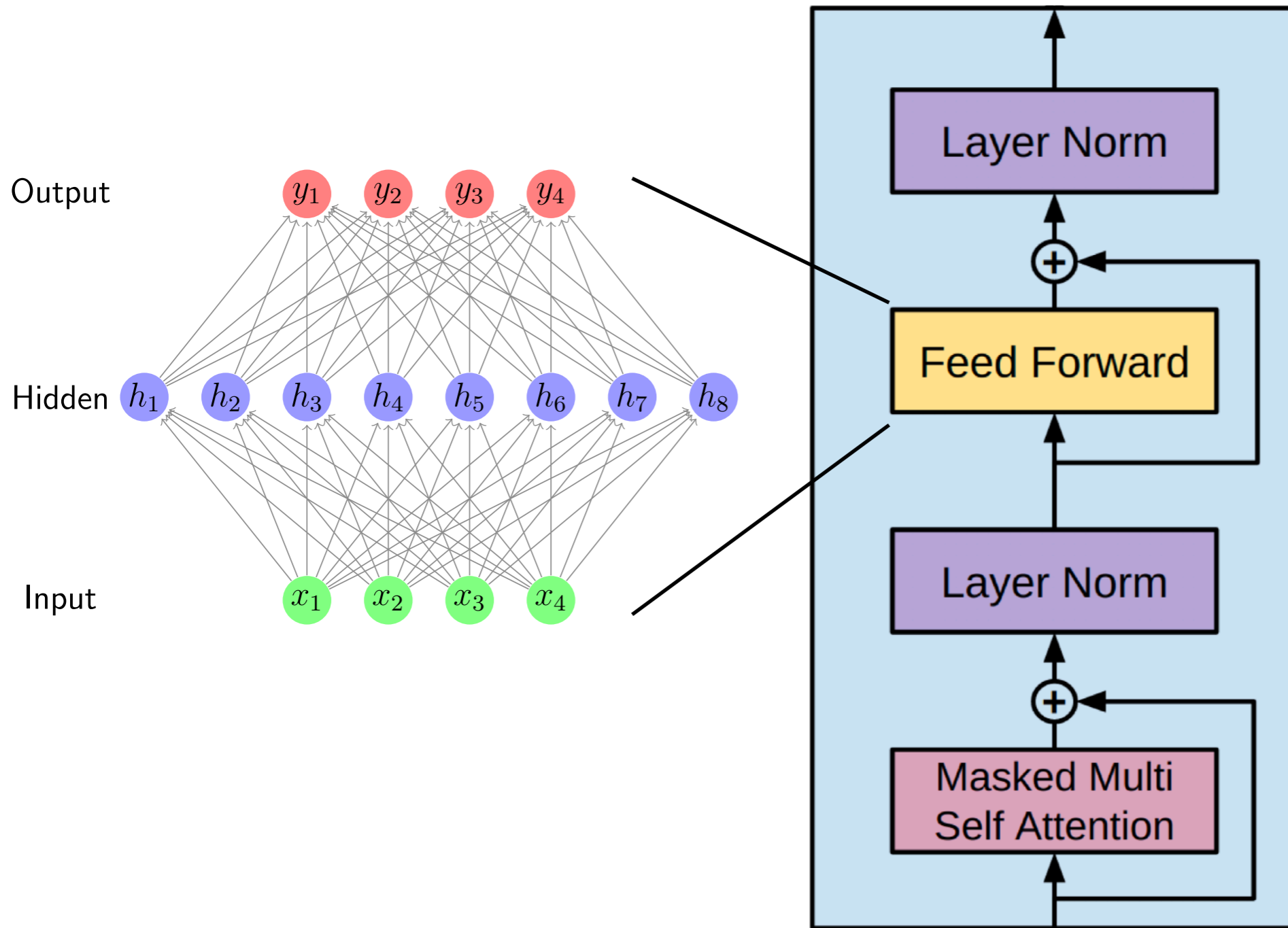$$\overrightarrow{x'} = \vec{x} + \vec{o}$$
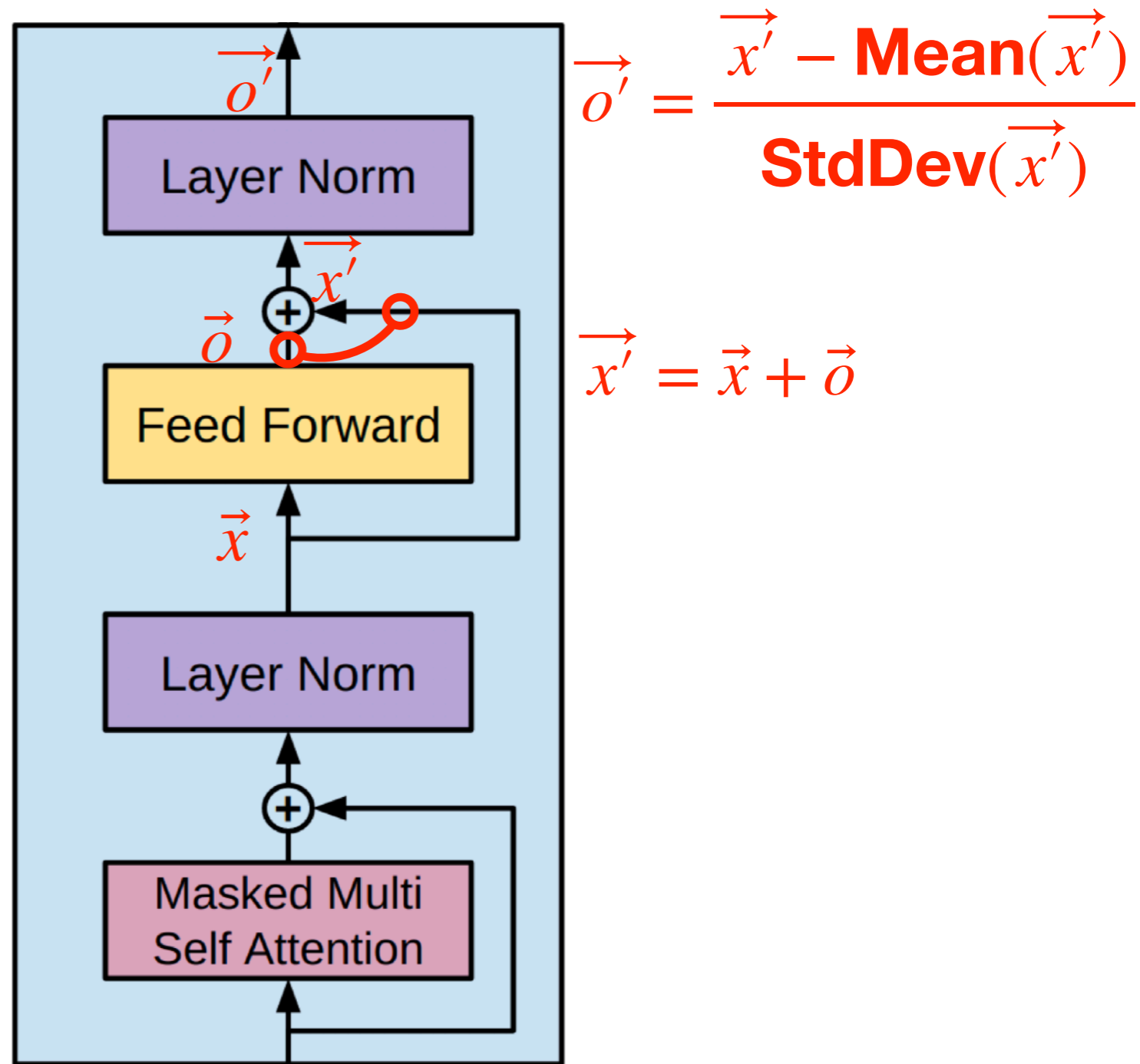
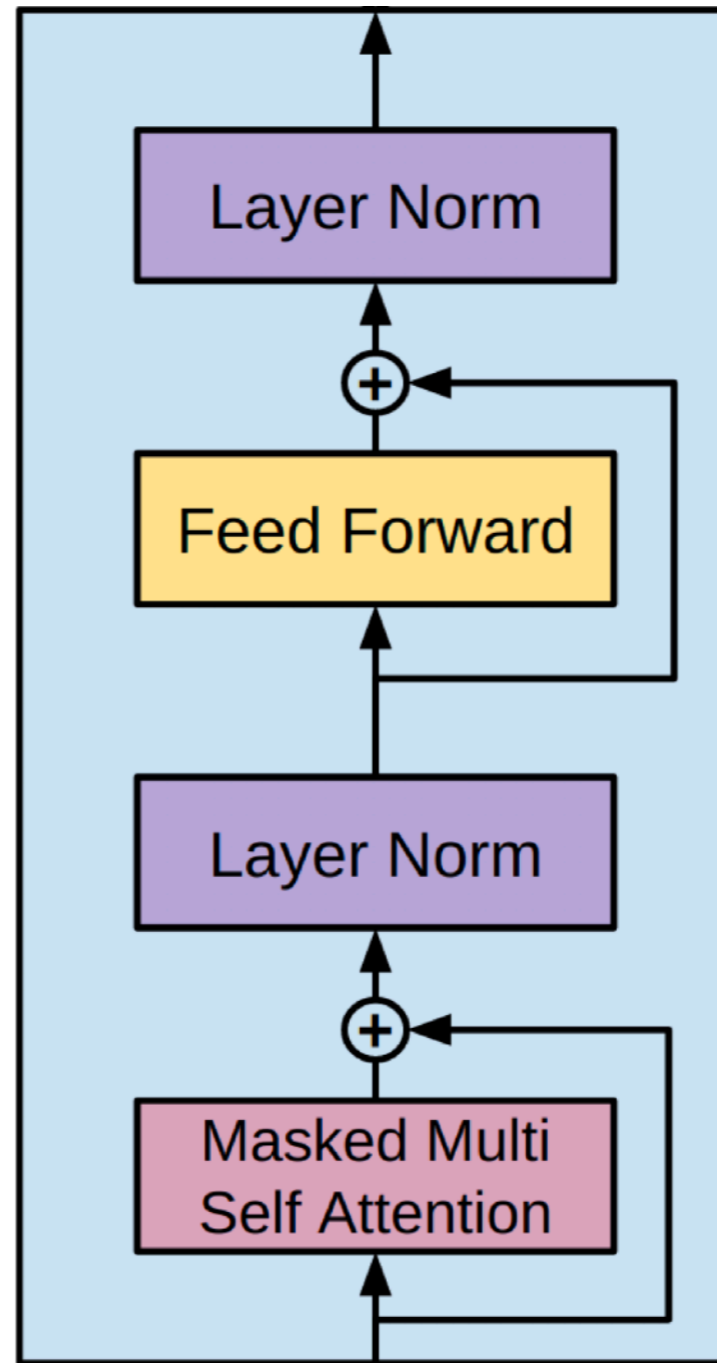# Feed-forward layer

# Res. connection & layer norm. (again)



$$\vec{o'} = \frac{\vec{x'} - \textbf{Mean}(\vec{x'})}{\textbf{StdDev}(\vec{x'})}$$

$$\vec{x'} = \vec{x} + \vec{o}$$

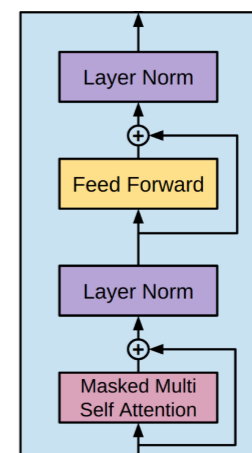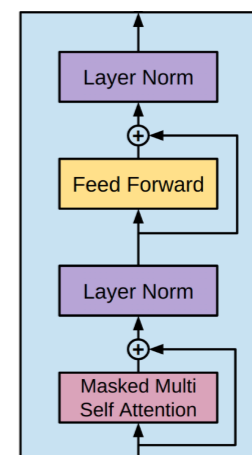# Res. connection & layer norm. (again)

# Res. connection & layer norm. (again)

# Res. connection & layer norm. (again)

# Res. connection & layer norm. (again)

# Res. connection & layer norm. (again)

# Transformer + a huge corpus = ...?

**New AI fake text generator may be too dangerous to release, say creators**

▶ The Guardian

- **OpenAI text-generating tool GPT2 won't be released for fear of misuse**
  ▶ Business Insider

📰 View Full Coverage        Feb 14, 2019 ⌄

▽ The Verge

**OpenAI has published the text-generating AI it said was too dangerous to share**

GPT-2 is part of a new breed of text-generation systems that have impressed experts with their ability to generate coherent text from minimal ...

Nov 7, 2019

# Write With Transformer |

transformer.huggingface.co

*Giant language model testing room:* `http://gltr.io/dist/index.html`

# Papers to read to understand GPT-2

- Radford et al. (2019): the GPT-2 paper itself
- Radford et al. (2018): the GPT architecture, mostly shared by GPT-2
- Liu et al. (2018): the Transformer decoder
- Vaswani et al. (2017): the original Transformer paper
- Ba et al. (2016): layer normalization

# The full Transformer model



*(Vaswani et al., 2017)*

# The full Transformer model

- In ML/NLP, the model we just studied is called the ***Transformer decoder***



*(Vaswani et al., 2017)*

# The full Transformer model

- In ML/NLP, the model we just studied is called the *Transformer decoder*

- Sometimes, the Transformer is conditioned on a string that doesn't itself get predicted—this is called the *encoder*

*(Vaswani et al., 2017)*

# The full Transformer model

- In ML/NLP, the model we just studied is called the ***Transformer decoder***

- Sometimes, the Transformer is conditioned on a string that doesn't itself get predicted—this is called the ***encoder***

- Only difference: in encoder, attention is over the ***entire string***, not just words to the left



*(Vaswani et al., 2017)*

# The full Transformer model

- In ML/NLP, the model we just studied is called the **_Transformer decoder_**

- Sometimes, the Transformer is conditioned on a string that doesn't itself get predicted—this is called the **_encoder_**

- Only difference: in encoder, attention is over the **_entire string_**, not just words to the left

- BERT = Transformer encoder!

Google has updated its search algorithm: Say hello to BERT

SmartCompany.com.au · Nov 4

*(Devlin et al., 2018)*

*(Vaswani et al., 2017)*

# GPT-2 on targeted syntax testing

## syntaxgym.org



*(Gauthier et al., 2020; Hu et al., 2020)*

# Filler—gap dependencies

✓      *I know that the lion devoured the gazelle at sunrise.*

*(Wilcox et al. 2018, Blackbox NLP)*

# Filler—gap dependencies

✓      *I know* **that** *the lion devoured* **the gazelle** *at sunrise.*

*(Wilcox et al. 2018, Blackbox NLP)*

# Filler—gap dependencies



✓ *I know **that** the lion devoured **the gazelle** at sunrise.*

*(Wilcox et al. 2018, Blackbox NLP)*

# Filler—gap dependencies

✓ *I know **that** the lion devoured **the gazelle** at sunrise.*

*(Wilcox et al. 2018, Blackbox NLP)*

# Filler—gap dependencies



✓ *I know **that** the lion devoured **the gazelle** at sunrise.*

**-FILLER
-GAP**

*(Wilcox et al. 2018, Blackbox NLP)*

# Filler—gap dependencies



✓  *I know that the lion devoured the gazelle at sunrise.*

✓  *I know* `what` *the lion devoured* ___ *at sunrise.*

**+FILLER
+GAP**

*(Wilcox et al. 2018, Blackbox NLP)*

# Filler—gap dependencies



✓ *I know that the lion devoured the gazelle at sunrise.*

\* *I know what the lion devoured the gazelle at sunrise.*

✓ *I know what the lion devoured ___ at sunrise.*

**+FILLER—GAP**

*(Wilcox et al. 2018, Blackbox NLP)*

# Filler—gap dependencies

## Approach: *Wh-Licensing Interaction*



✓ *I know that the lion devoured the gazelle at sunrise.*

✗ *I know what the lion devoured the gazelle at sunrise.*

✗ *I know that the lion devoured ▮▮▮ at sunrise.*

✓ *I know what the lion devoured ___ at sunrise.*

**−FILLER**
**+GAP**

*(Wilcox et al. 2018, Blackbox NLP)*

✓ *I know that my brother said our aunt devoured the cake at the party.*

�direct *I know what my brother said our aunt devoured the cake at the party.*

✓   *I know* <span style="color:teal">*that*</span> *my brother said our aunt devoured* <span style="color:green">*the cake*</span> *at the party.*

✳   *I know* <span style="color:salmon">*what*</span> *my brother said our aunt devoured* <span style="color:green">*the cake*</span> *at the party.*

✓ *I know **that** my brother said our aunt devoured **the cake** at the party.*

✱ *I know **what** my brother said our aunt devoured **the cake** at the party.*

✱ *I know **that** my brother said our aunt devoured _____ at the party.*

✓ *I know **what** my brother said our aunt devoured _____ at the party.*

✓  *I know **that** my brother said our aunt devoured **the cake** at the party.*

✱  *I know **what** my brother said our aunt devoured **the cake** at the party.*

✱  *I know **that** my brother said our aunt devoured _____ at the party.*

✓  *I know **what** my brother said our aunt devoured _____ at the party.*

✓    *I know* *that* *my brother said our aunt devoured* *the cake* *at the party.*

✳    *I know* *what* *my brother said our aunt devoured* *the cake* *at the party.*

✳    *I know* *that* *my brother said our aunt devoured* ＿＿＿ *at the party.*

✓    *I know* *what* *my brother said our aunt devoured* ＿＿＿ *at the party.*

✓    *I know* *that* *my brother said our aunt devoured* *the cake* *at the party.*

✳    *I know* *what* *my brother said our aunt devoured* *the cake* *at the party.*

✳    *I know* *that* *my brother said our aunt devoured* \_\_\_\_\_ *at the party.*

✓    *I know* *what* *my brother said our aunt devoured* \_\_\_\_\_ *at the party.*



27

# Unboundedness of *wh*-dependencies

*I know what our mother gave __ to Mary last weekend.*

# Unboundedness of *wh*-dependencies

*I know* *what* *our mother gave __ to Mary last weekend.*

*I know* *what* *our mother said that your friend gave _ to Mary last weekend.*

# Unboundedness of *wh*-dependencies

**0** *I know what our mother gave __ to Mary last weekend.*

**1** *I know what our mother said that your friend gave _ to Mary last weekend.*

# Unboundedness of *wh*-dependencies

**0** *I know what our mother gave __ to Mary last weekend.*

**1** *I know what our mother said that your friend gave __ to Mary last weekend.*

**2** *I know what our mother said that her friend remarked that your friend gave __ to Mary last weekend.*

# Unboundedness of *wh*-dependencies

**0** *I know what our mother gave __ to Mary last weekend.*

**1** *I know what our mother said that your friend gave __ to Mary last weekend.*

**2** *I know what our mother said that her friend remarked that your friend gave __ to Mary last weekend.*

**3** *I know what our mother said that her friend remarked that the park attendant wondered that your friend gave __ to Mary last weekend.*

# Unboundedness of *wh*-dependencies

**0** *I know what our mother gave __ to Mary last weekend.*

**1** *I know what our mother said that your friend gave __ to Mary last weekend.*

**2** *I know what our mother said that her friend remarked that your friend gave __ to Mary last weekend.*

**3** *I know what our mother said that her friend remarked that the park attendant wondered that your friend gave __ to Mary last weekend.*

**4** *I know what our mother said that her friend remarked that the park attendant wondered that the people stated that your friend gave __ to Mary last weekend.*

Unboundedness: Object Gap

# Potential concern #1

Couldn't the models be learning a **linear** dependency between filler and gap, not a **hierarchical** dependency?

# Syntactic Hierarchy

- A filler must be appropriately "above" its gap

32

S

NP — VP

the fact — SBAR

that — S

NP — VP

the mayor — knows — RC

who — S

NP — VP

the criminal — shot — NP — the teller

VP

shocked — NP — __

32

# Potential concern #1

Couldn't the models be learning a ***linear*** dependency between filler and gap, not a ***hierarchical*** dependency?

# Potential concern #1 — *addressed*

Couldn't the models be learning a ***linear*** dependency between filler and gap, not a ***hierarchical*** dependency?

Our results suggest that RNN models trained on enough data are sensitive to syntactic hierarchy for *wh*-dependency

# Does syntactic supervision help?



**Grammar-based model (RNNG)**

# Does syntactic supervision help?



**Grammar-based model (RNNG)**

# Does syntactic supervision help?



**Grammar-based model (RNNG)**

*stripped away*

# Syntactic supervision helps a lot!

- With small-dataset training (1m words):



Syntactic Hierarchy

**Grammar-based model (RNNG)**

**Sequence model (LSTM)**

# Syntactic island constraints

# Syntactic island constraints

- Some types of phrases are *islands*: filler–gap dependencies cannot link from outside to inside of them

# Syntactic island constraints

- Some types of phrases are ***islands***: filler–gap dependencies cannot link from outside to inside of them



- Islands are prominent in learnability debates: they'd require learning from negative evidence, and are rare structures

# Syntactic island constraints

- Some types of phrases are ***islands***: filler–gap dependencies cannot link from outside to inside of them



- Islands are prominent in learnability debates: they'd require learning from negative evidence, and are rare structures

- We take a language model to have learned an island constraint if it *fails* to propagate filler-generated expectations for gaps into phrases that should be islands

# Syntactic islands

**Wh-complementizers** block filler—gap dependencies:

*I know* what *Alex said…*

*…your friend devoured* __ *at the party.*
[null complementizer]

'

# Syntactic islands

**Wh-complementizers** block filler—gap dependencies:

*I know* what *Alex said…*

✓         *…your friend devoured* __ *at the party.*
                    [null complementizer]

,

# Syntactic islands

**Wh*-complementizers** block filler—gap dependencies:

*I know* what *Alex said…*

✓          *…your friend devoured* __ *at the party.*
                              [null complementizer]

*…**that** your friend devoured* __ *at the party.*
                              [*that* complementizer]

'

# Syntactic islands

**Wh-complementizers** block filler—gap dependencies:

*I know* what *Alex said…*

✓       *…your friend devoured* __ *at the party.*
[null complementizer]

✓     *…**that** your friend devoured* __ *at the party.*
[*that* complementizer]

# Syntactic islands

**Wh-complementizers** block filler—gap dependencies:

*I know* what *Alex said…*

✓      *…your friend devoured __ at the party.*
[null complementizer]

✓    ***…that** your friend devoured __ at the party.*
[*that* complementizer]

   ***…whether** your friend devoured __ at the party.*
[*wh*–complementizer]

# Syntactic islands

**Wh-complementizers** block filler—gap dependencies:

*I know* what *Alex said…*

✓        *…your friend devoured* __ *at the party.*
                              [null complementizer]

✓   *…**that** your friend devoured* __ *at the party.*
                              [*that* complementizer]

✳   *…**whether** your friend devoured* __ *at the party.*
                              [*wh*–complementizer]

# Syntactic islands

**Wh-complementizers** block filler—gap dependencies:

*I know* what *Alex said…*

✓     *…your friend devoured* __ *at the party.*
[null complementizer]

✓  *…that your friend devoured* __ *at the party.*
[*that* complementizer]

✳ *…whether your friend devoured* __ *at the party.*
[*wh*-complementizer]

Do the RNNs learn this?

✓ *I know that my brother said our aunt devoured the cake at the party.*

✱ *I know* *what* *my brother said our aunt devoured the cake at the party.*

✱ *I know that my brother said our aunt devoured* _____ *at the party.*

✓ *I know* *what* *my brother said our aunt devoured* _____ *at the party.*

∗ *I know that my brother said **whether** our aunt devoured the cake at the party.*

∗ *I know what my brother said **whether** our aunt devoured the cake at the party.*

∗ *I know that my brother said **whether** our aunt devoured _____ at the party.*

∗ *I know what my brother said **whether** our aunt devoured _____ at the party.*

# Potential concern #2

Could RNNs have difficulty threading *any* type of expectation into a syntactic island?

# Gendered-pronoun Expectation Control

# Gendered-pronoun Expectation Control

- Worry: Can the models thread **any** expectation into islands?

# Gendered-pronoun Expectation Control

- Worry: Can the models thread **any** expectation into islands?

- Test with expectation for **gendered pronouns** set up by **culturally or morphologically gendered subjects**.

*(Wilcox et al., 2019, CogSci)*

# Gendered-pronoun Expectation Control

- Worry: Can the models thread **any** expectation into islands?

- Test with expectation for **gendered pronouns** set up by **culturally or morphologically gendered subjects**.

    ✓ *The actress said that they insulted* **her** *friends.*
    **[CONTROL, MATCH]**

*(Wilcox et al., 2019, CogSci)*

# Gendered-pronoun Expectation Control

- Worry: Can the models thread **any** expectation into islands?
- Test with expectation for **gendered pronouns** set up by **culturally or morphologically gendered subjects**.

✓ *The actress said that they insulted* **her** *friends.*
**[CONTROL, MATCH]**

\# *The actress said that they insulted* **his** *friends.*
**[CONTROL, MISMATCH]**

*(Wilcox et al., 2019, CogSci)*  47

# Gendered-pronoun Expectation Control

- Worry: Can the models thread **any** expectation into islands?

- Test with expectation for **gendered pronouns** set up by **culturally or morphologically gendered subjects**.

**Gender Expectation Effect (#-✓ should be *positive*)**

✓ *The actress said that they insulted **her** friends.* **[CONTROL, MATCH]**

\# *The actress said that they insulted **his** friends.* **[CONTROL, MISMATCH]**

*(Wilcox et al., 2019, CogSci)*

# Gendered-pronoun Expectation Control

- Worry: Can the models thread **any** expectation into islands?

- Test with expectation for **gendered pronouns** set up by **culturally or morphologically gendered subjects**.

**Gender Expectation Effect (#-✓ should be *positive*)**

✓ *The actress said that they insulted* **her** *friends.* **[CONTROL, MATCH]**

# *The actress said that they insulted* **his** *friends.* **[CONTROL, MISMATCH]**

✓ *The actress said whether they insulted* **her** *friends.* **[ISLAND, MATCH]**

# Gendered-pronoun Expectation Control

- Worry: Can the models thread **any** expectation into islands?

- Test with expectation for **gendered pronouns** set up by **culturally or morphologically gendered subjects**.

**Gender Expectation Effect (#-✓ should be *positive*)**

✓ *The actress said that they insulted* **her** *friends.* **[CONTROL, MATCH]**

\# *The actress said that they insulted* **his** *friends.* **[CONTROL, MISMATCH]**

✓ *The actress said whether they insulted* **her** *friends.* **[ISLAND, MATCH]**

\# *The actress said whether they insulted* **his** *friends.* **[ISLAND, MISMATCH]**

*(Wilcox et al., 2019, CogSci)*

# Gendered-pronoun Expectation Control

- Worry: Can the models thread **any** expectation into islands?
- Test with expectation for **gendered pronouns** set up by **culturally or morphologically gendered subjects**.

**Gender Expectation Effect (#-✓ should be _positive_)**

✓ *The actress said that they insulted* **her** *friends.*
**[CONTROL, MATCH]**

\# *The actress said that they insulted* **his** *friends.*
**[CONTROL, MISMATCH]**

**?**

✓ *The actress said whether they insulted* **her** *friends.*
**[ISLAND, MATCH]**

\# *The actress said whether they insulted* **his** *friends.*
**[ISLAND, MISMATCH]**

*(Wilcox et al., 2019, CogSci)*

# Gendered-pronoun Expectation Control

- Worry: Can the models thread **any** expectation into islands?
- Test with expectation for **gendered pronouns** set up by **culturally or morphologically gendered subjects**.

**Gender Expectation Effect (#-✓ should be *positive*)**

✓ *The actress said that they insulted* **her** *friends.* **[CONTROL, MATCH]**

\# *The actress said that they insulted* **his** *friends.* **[CONTROL, MISMATCH]**

**?**

✓ *The actress said whether they insulted* **her** *friends.* **[ISLAND, MATCH]**

\# *The actress said whether they insulted* **his** *friends.* **[ISLAND, MISMATCH]**

S
NP — VP
the actress   said   SBAR
whether   S
NP   VP
they   insulted   NP
her friends

*(Wilcox et al., 2019, CogSci)*

# Gendered-pronoun Expectation Control

- Worry: Can the models thread **any** expectation into islands?
- Test with expectation for **gendered pronouns** set up by **culturally or morphologically gendered subjects**.

**Gender Expectation Effect (#-✓ should be *positive*)**

✓ *The actress said that they insulted* **her** *friends.*
**[CONTROL, MATCH]**

# *The actress said that they insulted* **his** *friends.*
**[CONTROL, MISMATCH]**

✓ *The actress said whether they insulted* **her** *friends.*
**[ISLAND, MATCH]**

# *The actress said whether they insulted* **his** *friends.*
**[ISLAND, MISMATCH]**

**?**

S
— NP — VP
the actress — said
SBAR
whether — S
NP — VP
they — insulted — NP
**her** friends

*(Wilcox et al., 2019, CogSci)*

# Gendered-pronoun Expectation Control

- Worry: Can the models thread **any** expectation into islands?
- Test with expectation for **gendered pronouns** set up by **culturally or morphologically gendered subjects**.

**Gender Expectation Effect (#-✓ should be *positive*)**

✓ *The actress said that they insulted **her** friends.* **[CONTROL, MATCH]**

\# *The actress said that they insulted **his** friends.* **[CONTROL, MISMATCH]**

**?**

✓ *The actress said whether they insulted **her** friends.* **[ISLAND, MATCH]**

\# *The actress said whether they insulted **his** friends.* **[ISLAND, MISMATCH]**

**???**

*(Wilcox et al., 2019, CogSci)*

# Gendered-pronoun Expectation Control

- Worry: Can the models thread **any** expectation into islands?
- Test with expectation for **gendered pronouns** set up by **culturally or morphologically gendered subjects**.

**Gender Expectation Effect (#-✓ should be *positive*)**

✓ *The actress said that they insulted **her** friends.* **[CONTROL, MATCH]**

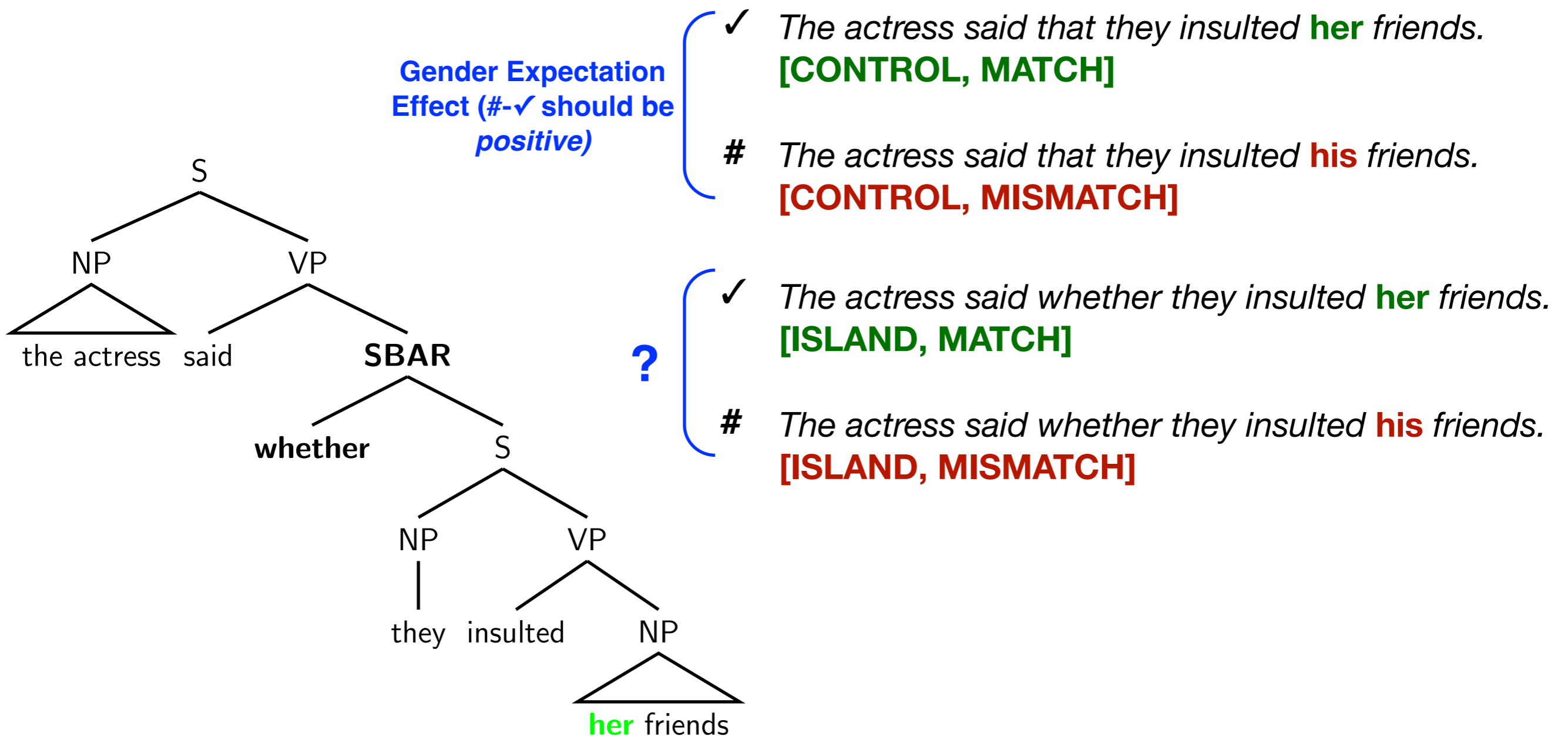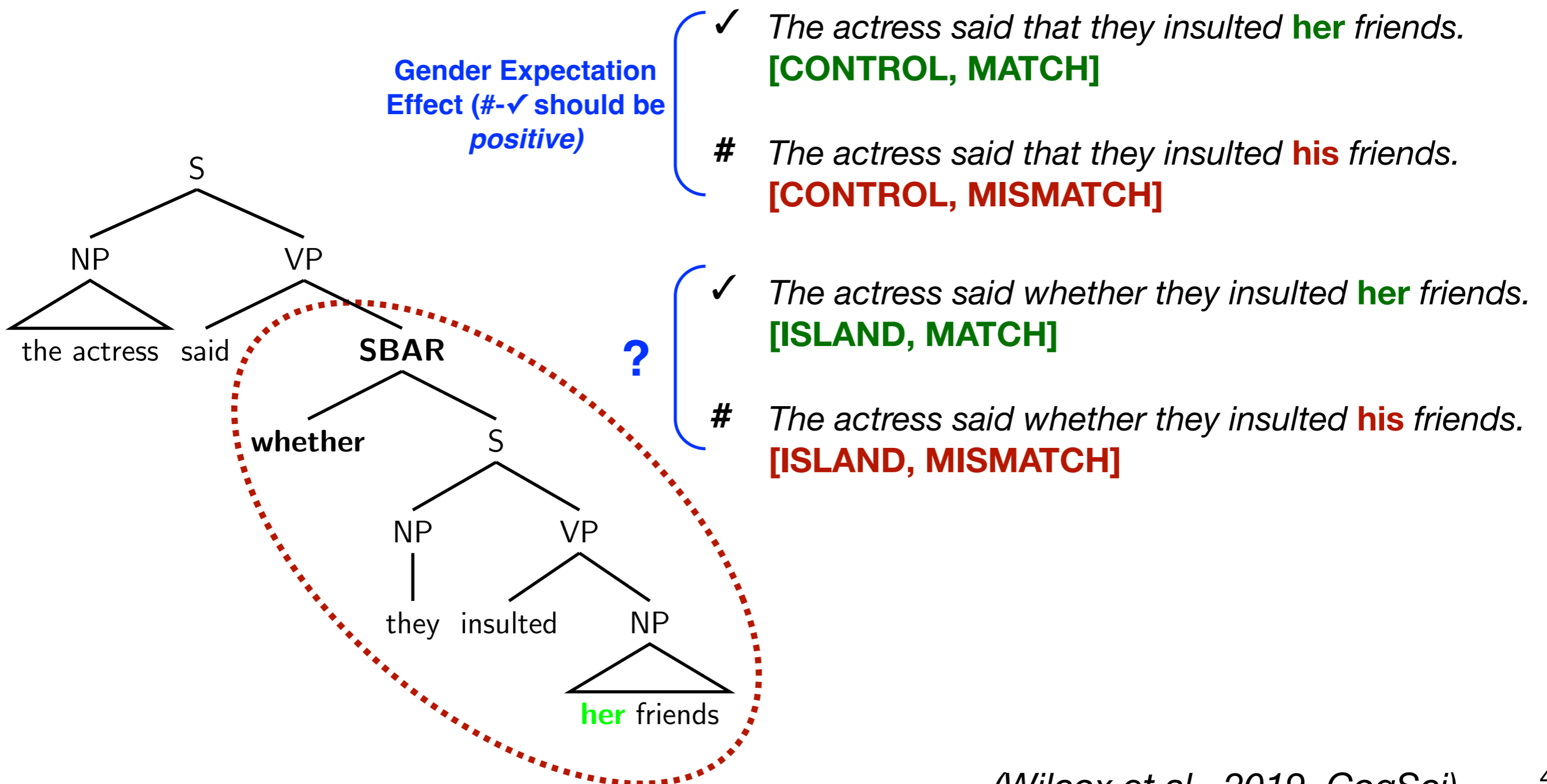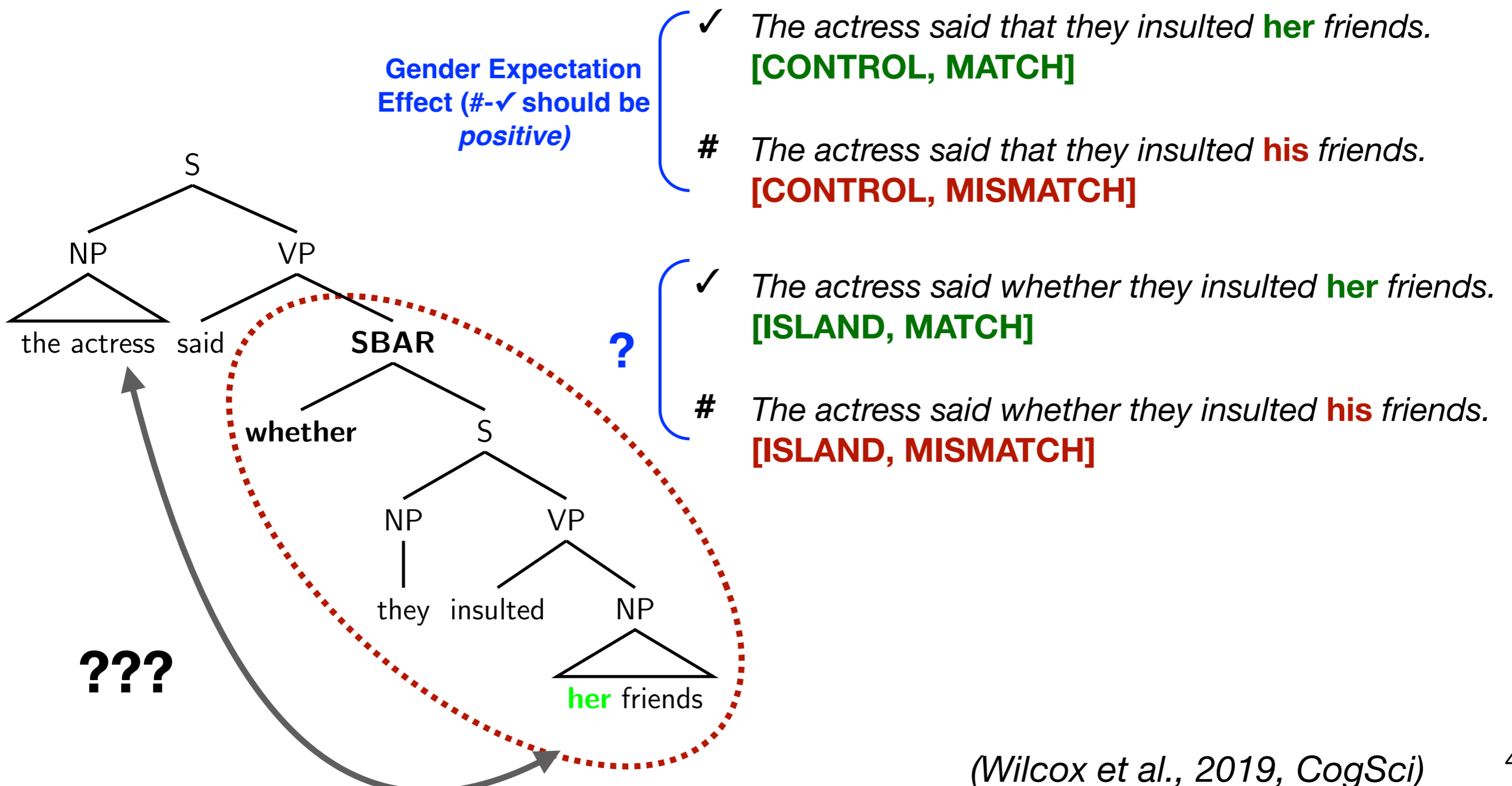\# *The actress said that they insulted **his** friends.* **[CONTROL, MISMATCH]**

✓ *The actress said whether they insulted **her** friends.* **[ISLAND, MATCH]**

**?**

\# *The actress said whether they insulted **his** friends.* **[ISLAND, MISMATCH]**

S
NP        VP
the actress   said        SBAR
whether        S
NP        VP
they   insulted      NP
her friends

**???**

If models can thread gender expectation into islands, the gender expectation effect should **look the same in islands as in the control conditions.**

*(Wilcox et al., 2019, CogSci)*

*The actress said that they insulted **her** friends.*

*The actress said that they insulted **his** friends.*

*The actress said that they insulted* **her** *friends.*

*The actress said that they insulted* **his** *friends.*



48

*The actress said that they insulted **her** friends.*

*The actress said that they insulted **his** friends.*

*The actress said that they insulted **her** friends.*

*The actress said that they insulted **his** friends.*

*The actress said whether they insulted **her** friends.*

*The actress said whether they insulted **his** friends.*

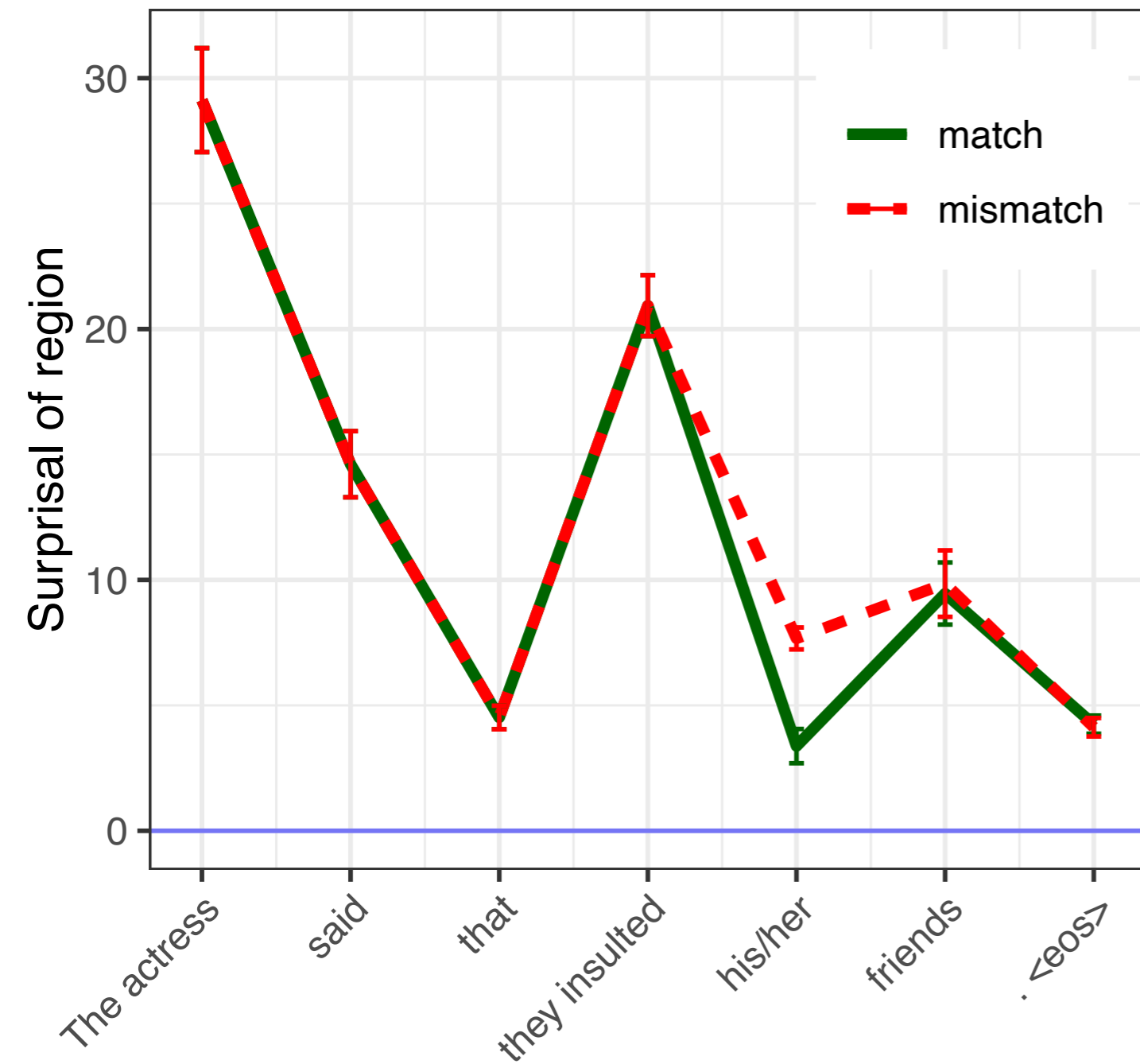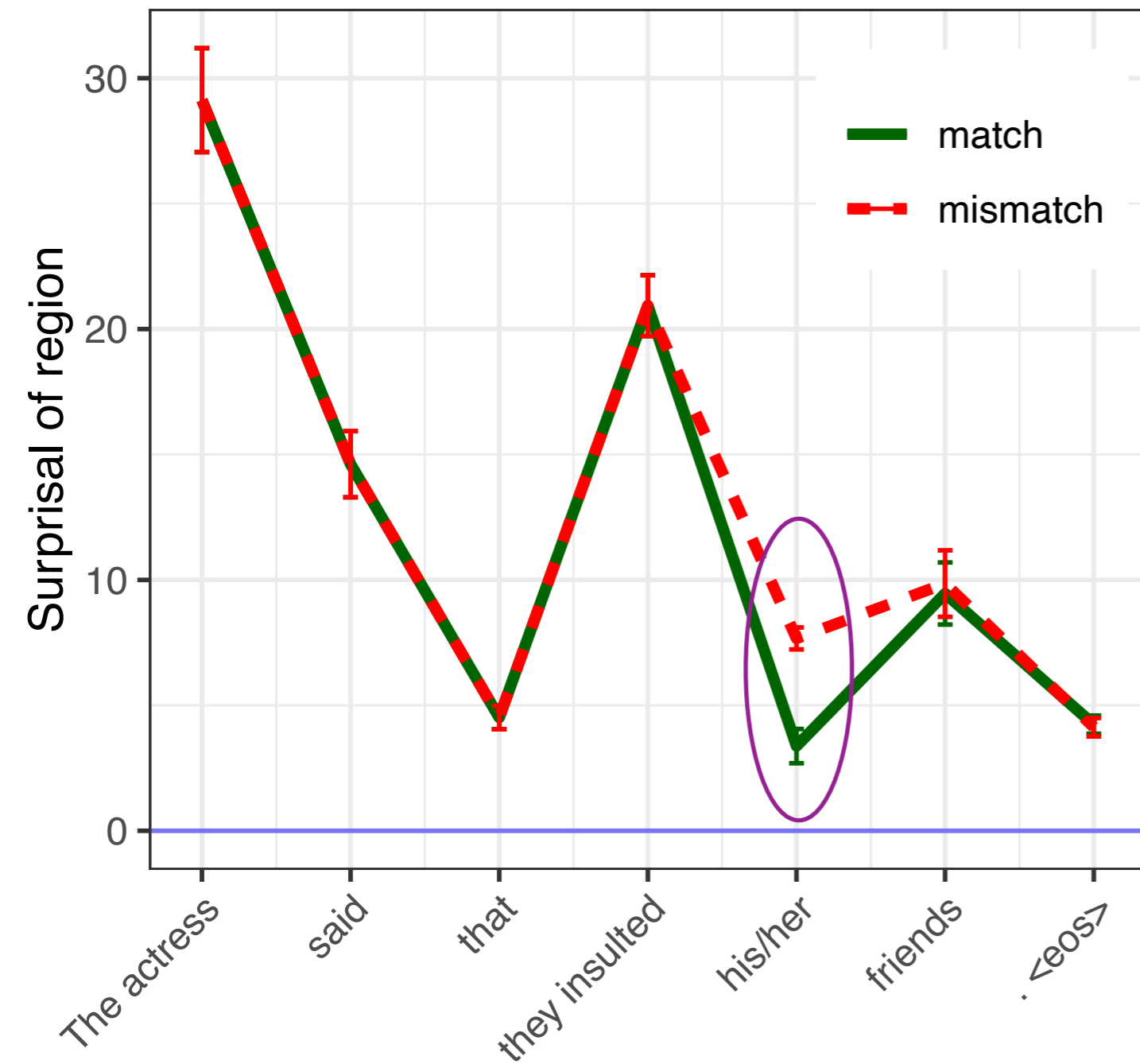The actress said that they insulted *her* friends.

The actress said that they insulted *his* friends.

The actress said whether they insulted *her* friends.
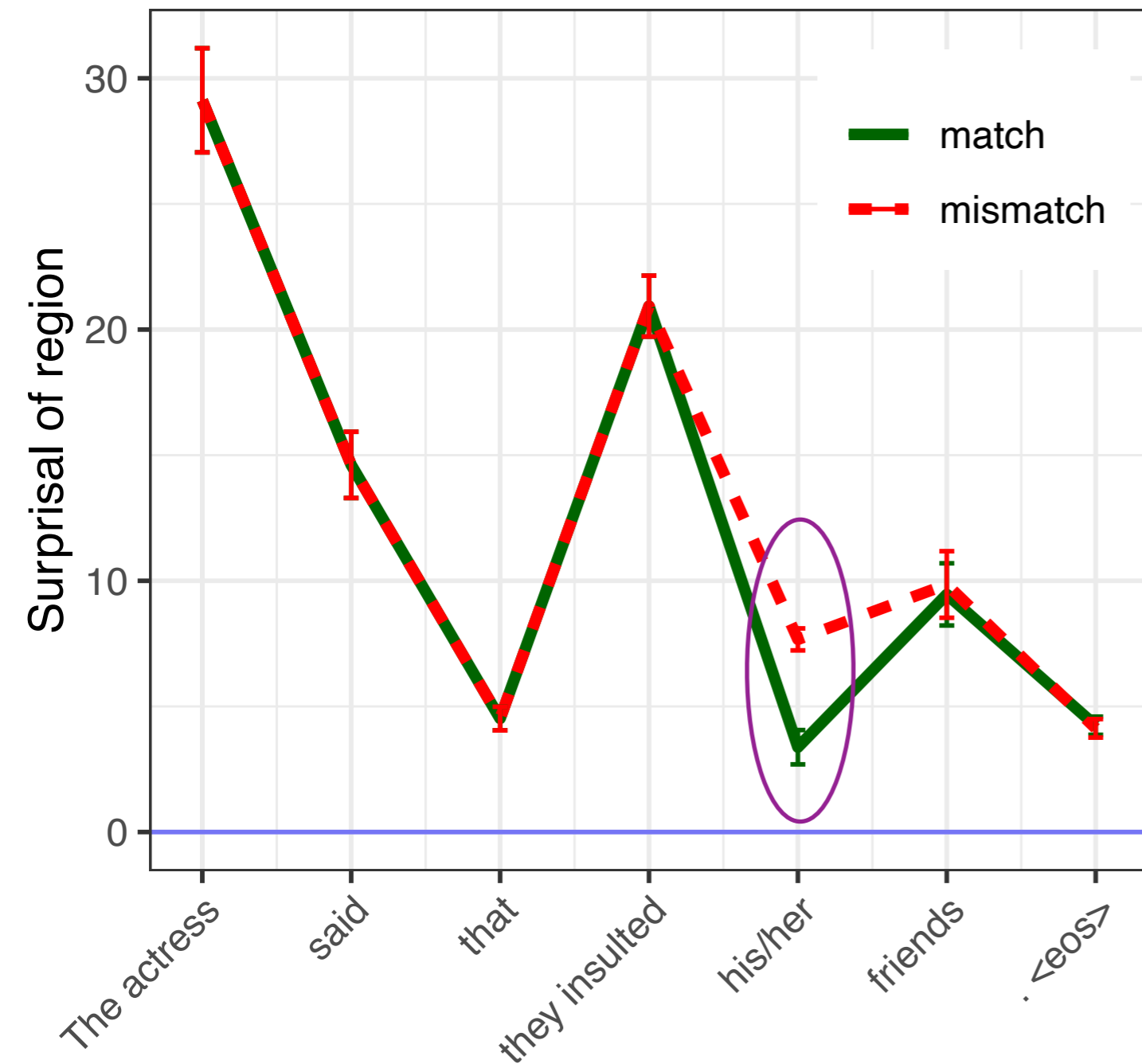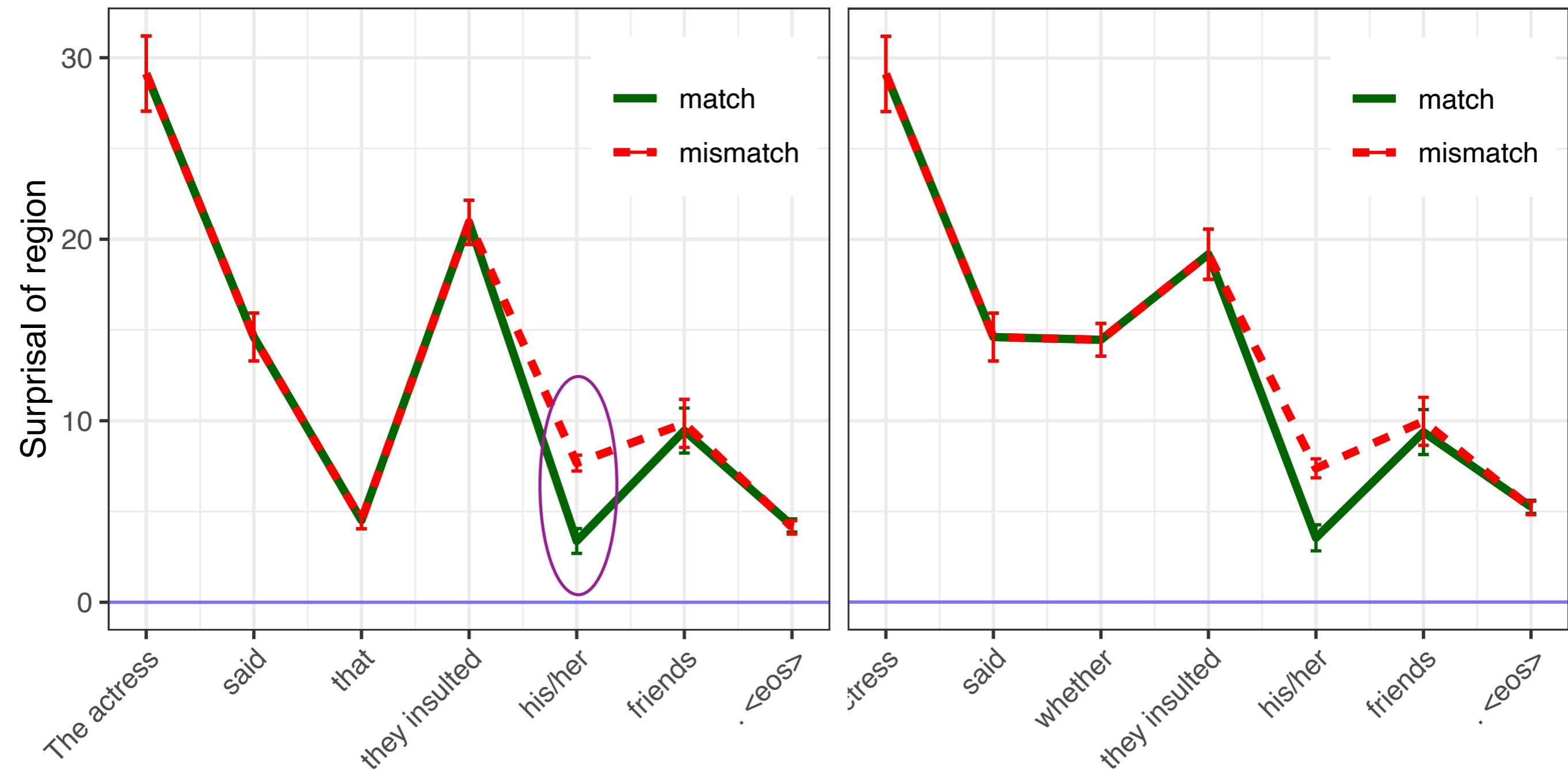
The actress said whether they insulted *his* friends.

The actress said that they insulted *her* friends.

The actress said that they insulted *his* friends.

The actress said whether they insulted *her* friends.

The actress said whether they insulted *his* friends.

# Potential concern #2

Could RNNs have difficulty threading **any** type of expectation into a syntactic island?

# Potential concern #2 — *addressed*

Could RNNs have difficulty u~~n~~ading *any* type of expectation into a sy~~n~~tic island?

RNN models that learn island constraints still propagate pronoun gender expectations into islands

# References

Ba, J. L., Kiros, J. R., & Hinton, G. E. (2016). Layer normalization. arXiv preprint arXiv:1607.06450.

Bahdanau, D., Cho, K., & Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. CoRR, abs/1409.0473. arXiv: 1409.0473

Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of deep bidirectional transformers for language understanding. Proceedings of the 17th Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies.

Gauthier, J., Hu, J., Wilcox, E., Qian, P., & Levy, R. P. (2020). SyntaxGym: An online platform for targeted evaluation of language models. In Proceedings of the 58th annual meeting of the Association for Computational Linguistics.

Hu, J., Gauthier, J., Qian, P., Wilcox, E., & Levy, R. P. (2020). A systematic assessment of syntactic generalization in neural language models. In Proceedings of the 58th annual meeting of the Association for Computational Linguistics.

Liu, P. J., Saleh, M., Pot, E., Goodrich, B., Sepassi, R., Kaiser, L., & Shazeer, N. (2018). Generating wikipedia by summarizing long sequences. In Proceedings of ICLR.

Luong, T., Pham, H., & Manning, C. D. (2015). Effective approaches to attention-based neural machine translation. In Proceedings of the 2015 conference on empirical methods in natural language processing (pp. 1412–1421).

Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding by generative pre-training.

Vaswani, A., Shazeer, N., Parmar, N.,Uszkoreit, J.,Jones,L.,Gomez, A.N.,Kaiser,L[?].,& Polosukhin, I.(2017). Attention is all you need. In Proceedings of Neural Information Processing Systems (pp. 5998–6008).

Wilcox, E., Levy, R. P., & Futrell, R. (2019). What syntactic structures block dependencies in RNN language models? In Proceedings of the 41st annual meeting of the Cognitive Science Society (pp. 1199–1205).

Wilcox, E., Levy, R. P., Morita, T., & Futrell, R. (2018). What do RNN language models learn about filler–gap dependencies? In Proceedings of the workshop on analyzing and interpreting neural networks for NLP.

Wilcox, E., Qian, P., Futrell, R., Ballesteros, M., & Levy, R. (2019). Structural supervision improves learning of non-local grammatical dependencies. In Proceedings of the 18th Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (pp. 3302–3312).