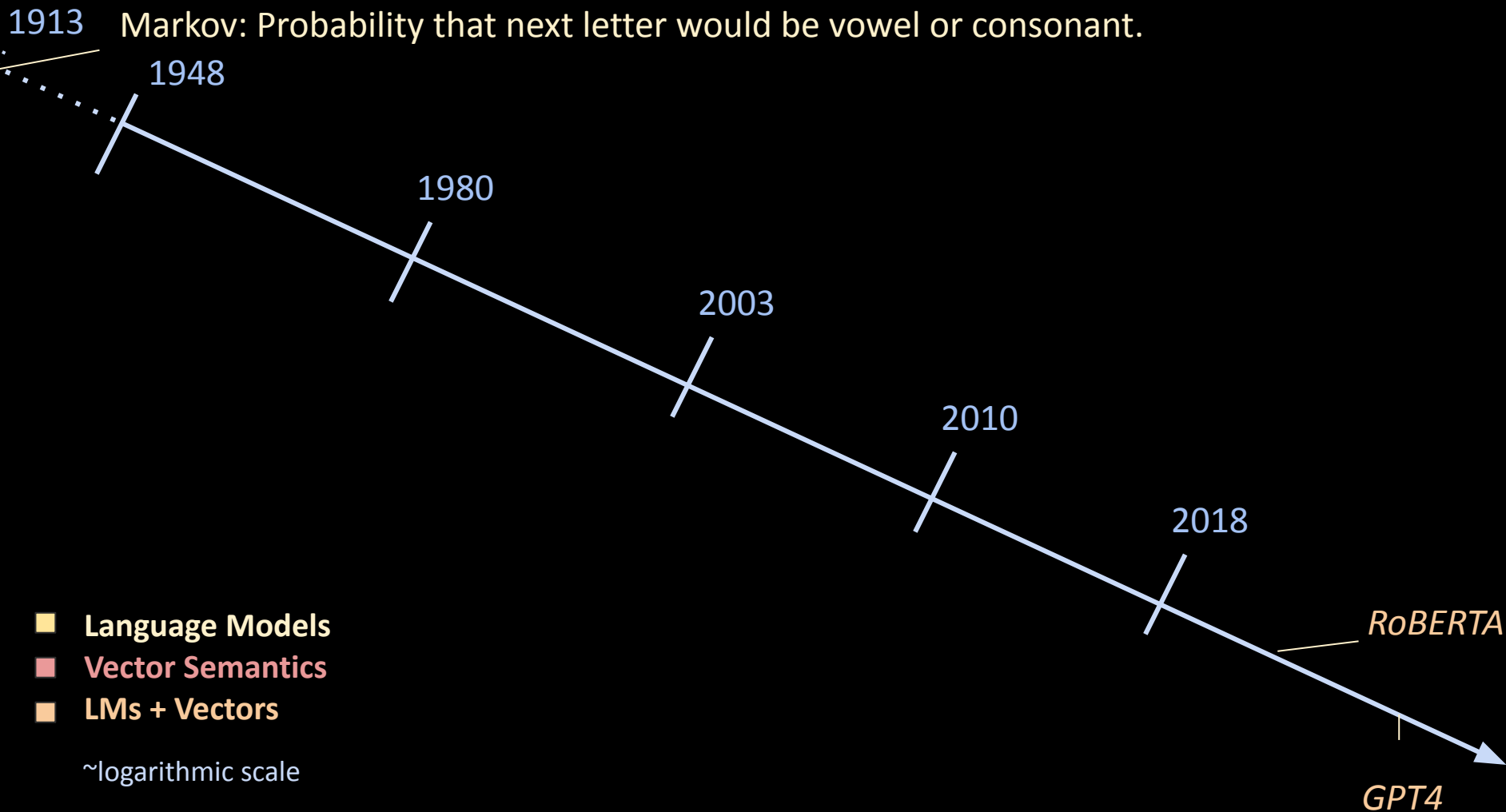


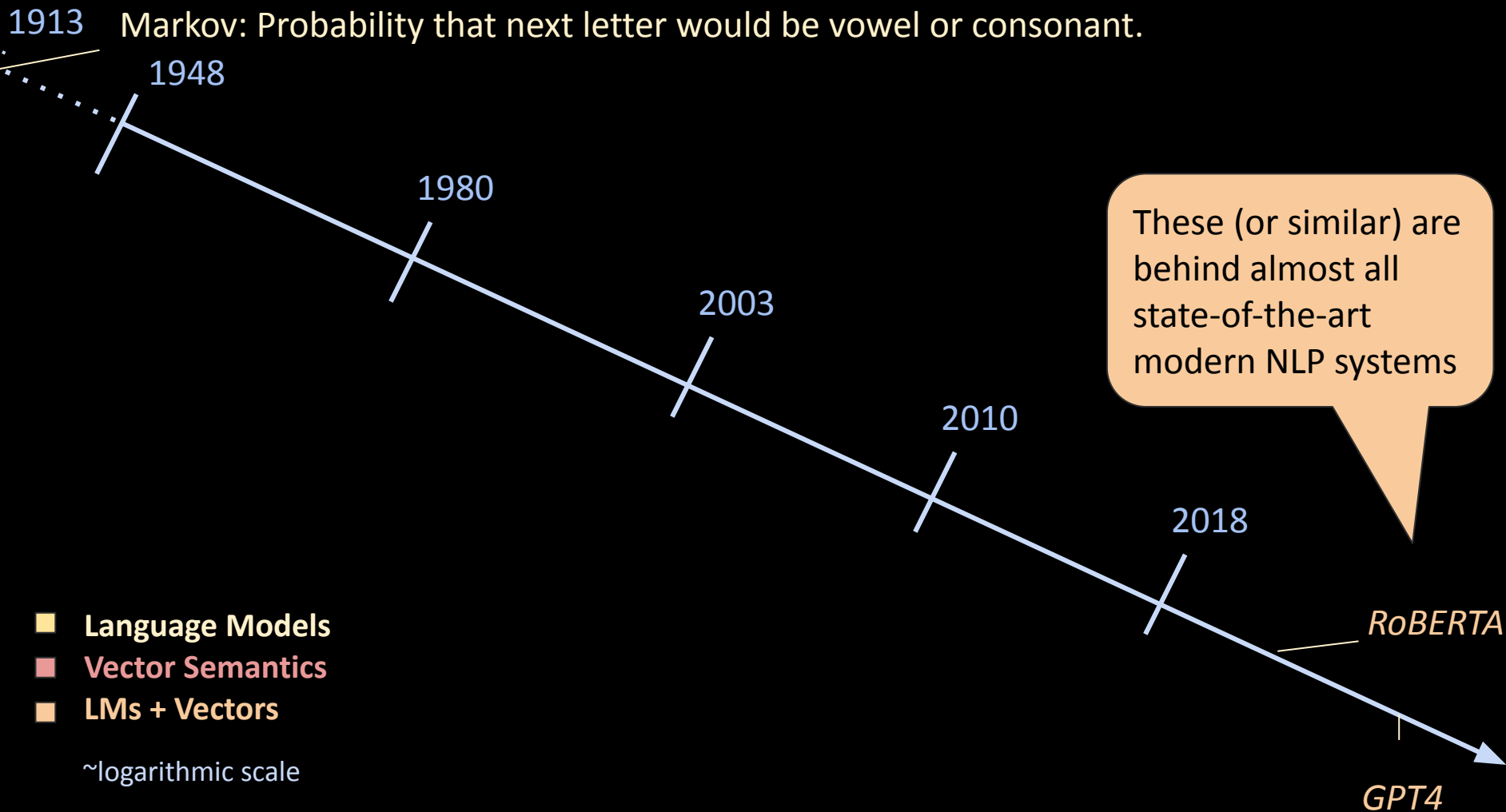
# Deep Learning and Transformers

CSE538 - Spring 2024

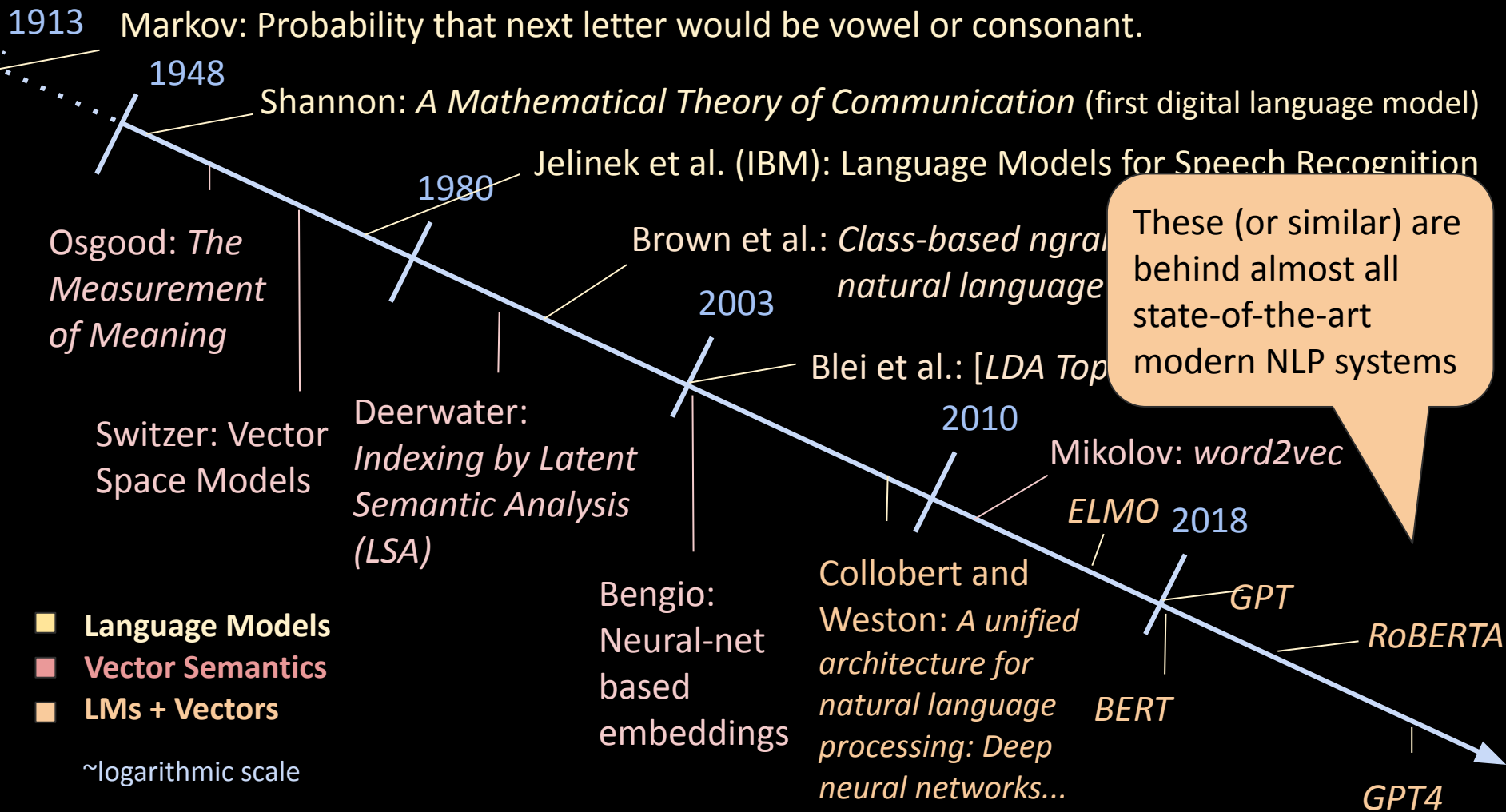
# Timeline: *Language Modeling* and *Vector Semantics*



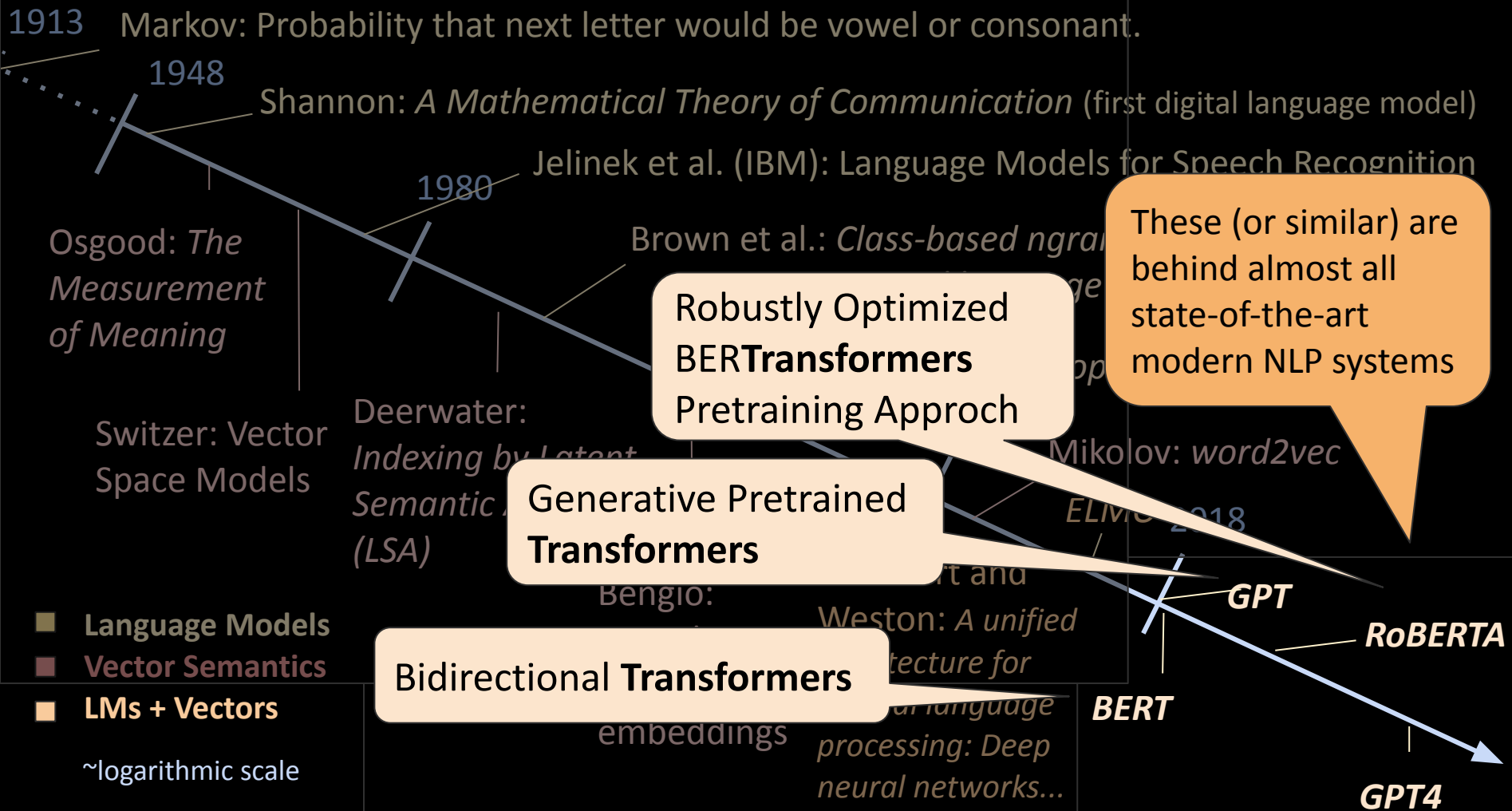
# Timeline: *Language Modeling* and *Vector Semantics*



# Timeline: *Language Modeling* and *Vector Semantics*



# Timeline: *Language Modeling* and *Vector Semantics*



# Transformers

## Attention Is All You Need

Ashish Vaswani\*  
Google Brain  
avaswani@google.com

Noam Shazeer\*  
Google Brain  
noam@google.com

Niki Parmar\*  
Google Research  
nikip@google.com

Jakob Uszkoreit\*  
Google Research  
usz@google.com

Llion Jones\*  
Google Research  
llion@google.com

Aidan N. Gomez\*<sup>†</sup>  
University of Toronto  
aidan@cs.toronto.edu

Lukasz Kaiser\*  
Google Brain  
lukaszkaizer@google.com

Illia Polosukhin\*<sup>‡</sup>  
illia.polosukhin@gmail.com

### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including the state-of-the-art single-model state-of-the-art BLEU score of 41.0 after scaling up over 2 BLEU. On the WMT 2014 English-to-French translation task, our model achieves 40.6 BLEU, improving over the existing best results, including the state-of-the-art single-model state-of-the-art BLEU score of 41.0 after scaling up over 2 BLEU. Our model achieves a small fraction of the training costs of the

Vaswani, A., Shazeer, N.,  
Parmar, N., Uszkoreit, J.,  
Jones, L., Gomez, A. N., ... &  
Polosukhin, I. (2017). Attention  
is all you need. *Advances in neural  
information processing systems*, 30.

# Transformers



## Self-Attention



## Deep Learning



## Neural Networks

# Transformers



## Self-Attention



## Deep Learning



## Neural Networks



- multi-headed attention
  - positional embeddings
  - residual links
- (to be introduced later)



# Timeline: *Language Modeling* and *Vector Semantics*

1913 Markov: Probability that next letter would be vowel or consonant.

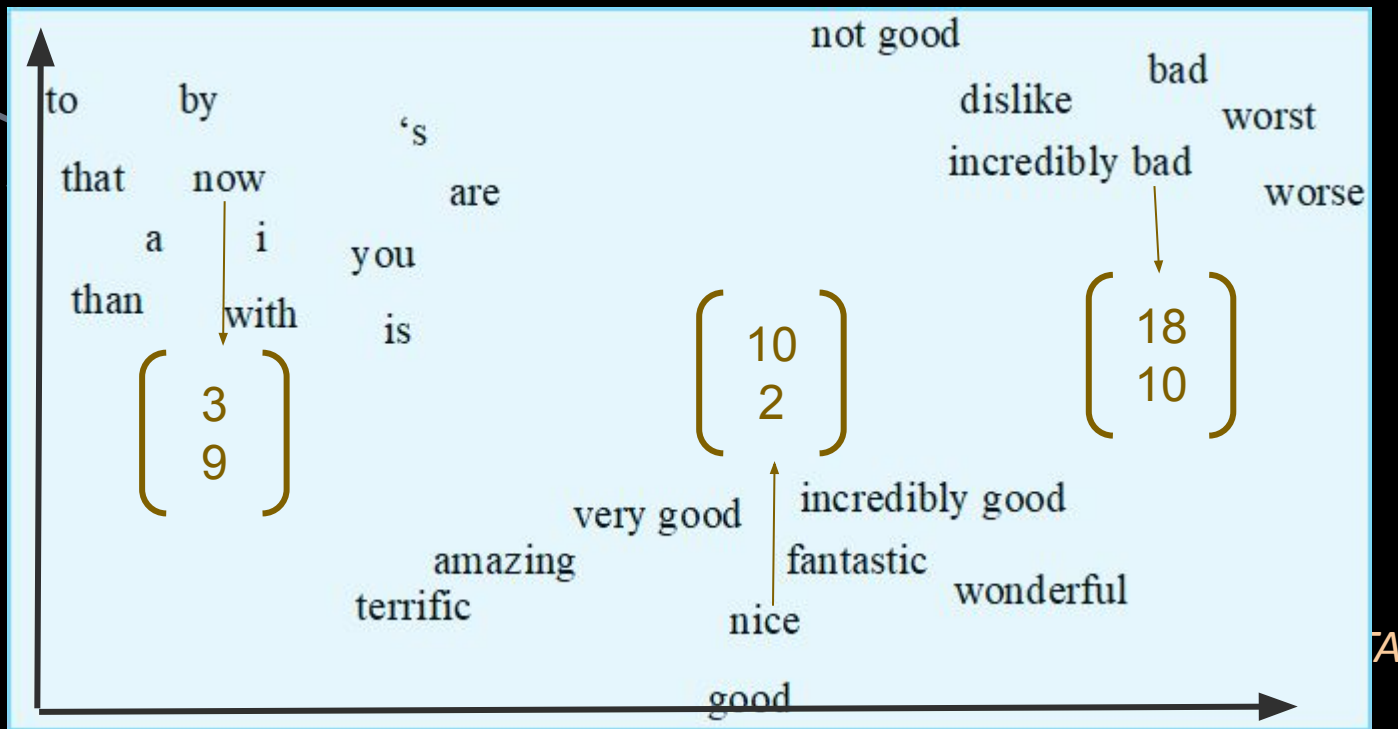
1948

Shannon: *A Mathematical Theory of Communication* (first digital language model)

Osgood: *The Measurement of Meaning*

- Language Models
- Vector Semantics
- LMs + Vectors

~logarithmic scale



(Li et al., 2015; Jurafsky et al., 2019)

GPT4

# Word Vectors

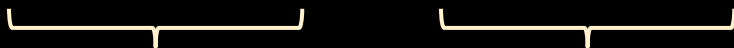
To embed: convert a token (or sequence) to a vector that represents **meaning**.

Wittgenstein, 1945: “*The meaning of a word is its use in the language*”

Distributional hypothesis -- A word's meaning is defined by all the different contexts it appears in (i.e. how it is “distributed” in natural language).

Firth, 1957: “*You shall know a word by the company it keeps*”

*The nail hit the beam behind the wall.*



# Word Vectors

## Person A

*How are you?*

I feel *fine* —even *great*!

*What is going on?*

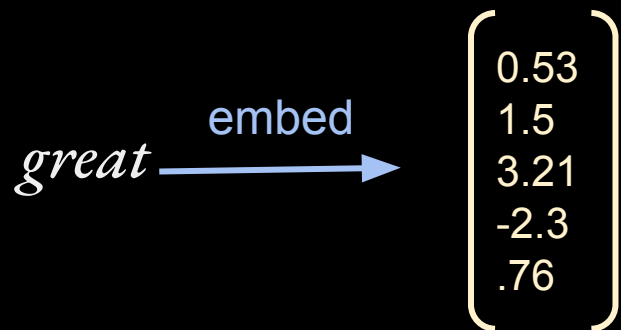
Earlier, I *played* the *game* Yahtzee with my *partner*. I could not get that *die* to roll a 1! Now I'm *lying* on my bed for a *rest*.

## Person B

My life is a *great* mess! I'm having a very hard time being happy.

My business *partner* was *lying* to me. He was trying to *game* the system and *played* me. I think I am going to *die* —he left and now I have to pay the *rest* of his *fine*.

# Objective



# Objective

*great* → embed

$$\begin{bmatrix} 0.53 \\ 1.5 \\ 3.21 \\ -2.3 \\ .76 \end{bmatrix}$$

?

**great.a.1** (relatively large in size or number or extent; larger than others of its kind)

**great.a.2**, outstanding (of major significance or importance)

**great.a.3** (remarkable or out of the ordinary in degree or magnitude or effect)

bang-up, bully, corking, cracking, dandy, **great.a.4**, groovy, keen, neat, nifty, not bad, peachy, slap-up, swell, smashing, old (very good)

capital, **great.a.5**, majuscule (uppercase)

big, enceinte, expectant, gravid, **great.a.6**, large, heavy, with child (in an advanced stage of pregnancy)

# Objective

*great* → embed

$$\begin{bmatrix} 0.53 \\ 1.5 \\ 3.21 \\ -2.3 \\ .76 \end{bmatrix}$$

?

**great.a.1** (relatively large in size or number or extent; larger than others of its kind)

**great.a.2**, outstanding (of major significance or importance)

**great.a.3** (remarkable or out of the ordinary in degree or magnitude or effect)

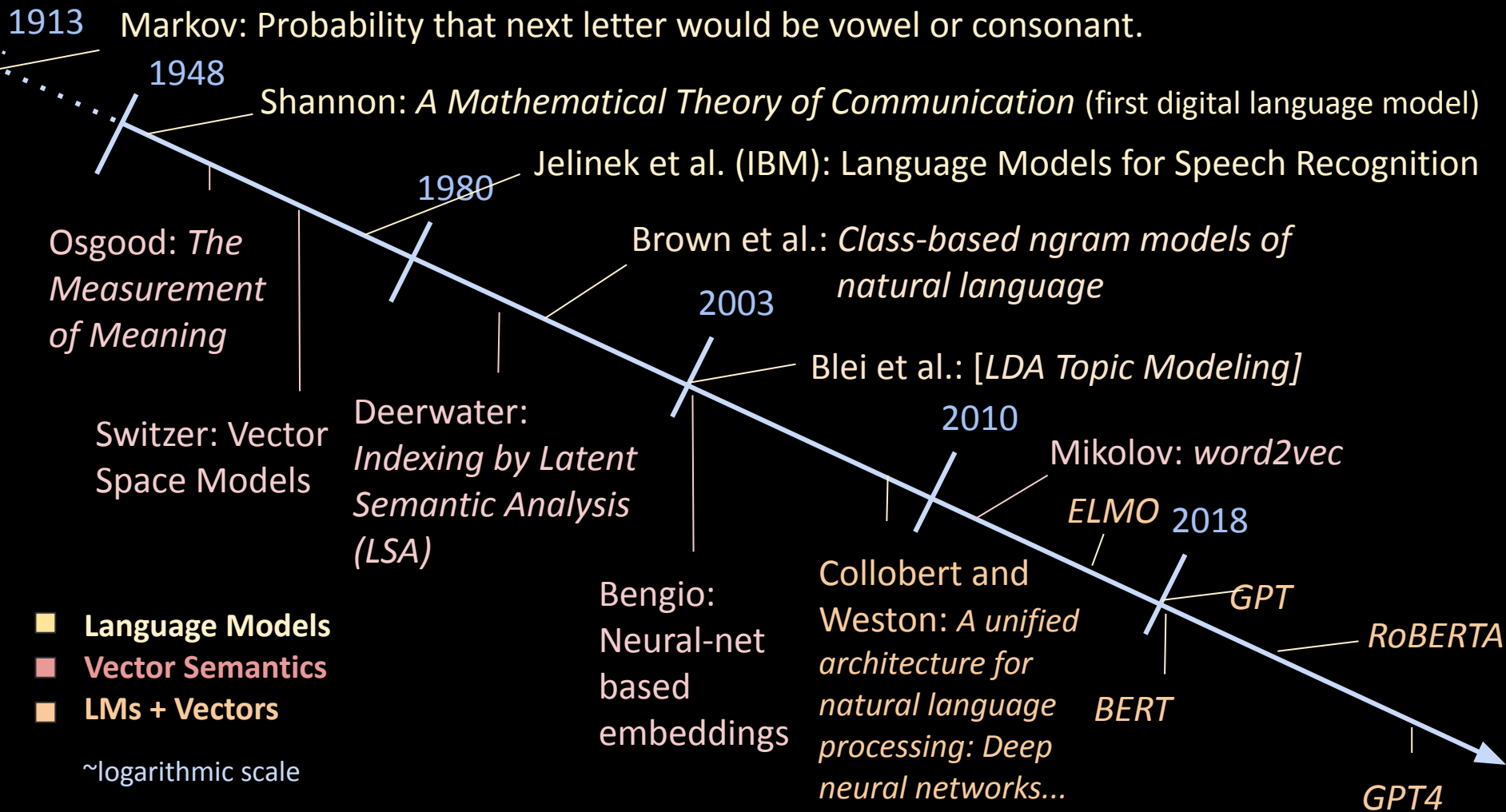
bang-up, bully, corking, cracking, dandy, **great.a.4**, groovy, keen, neat, nifty, not bad, peachy, slap-up, swell, smashing, old (very good)

capital, **great.a.5**, majuscule (uppercase)

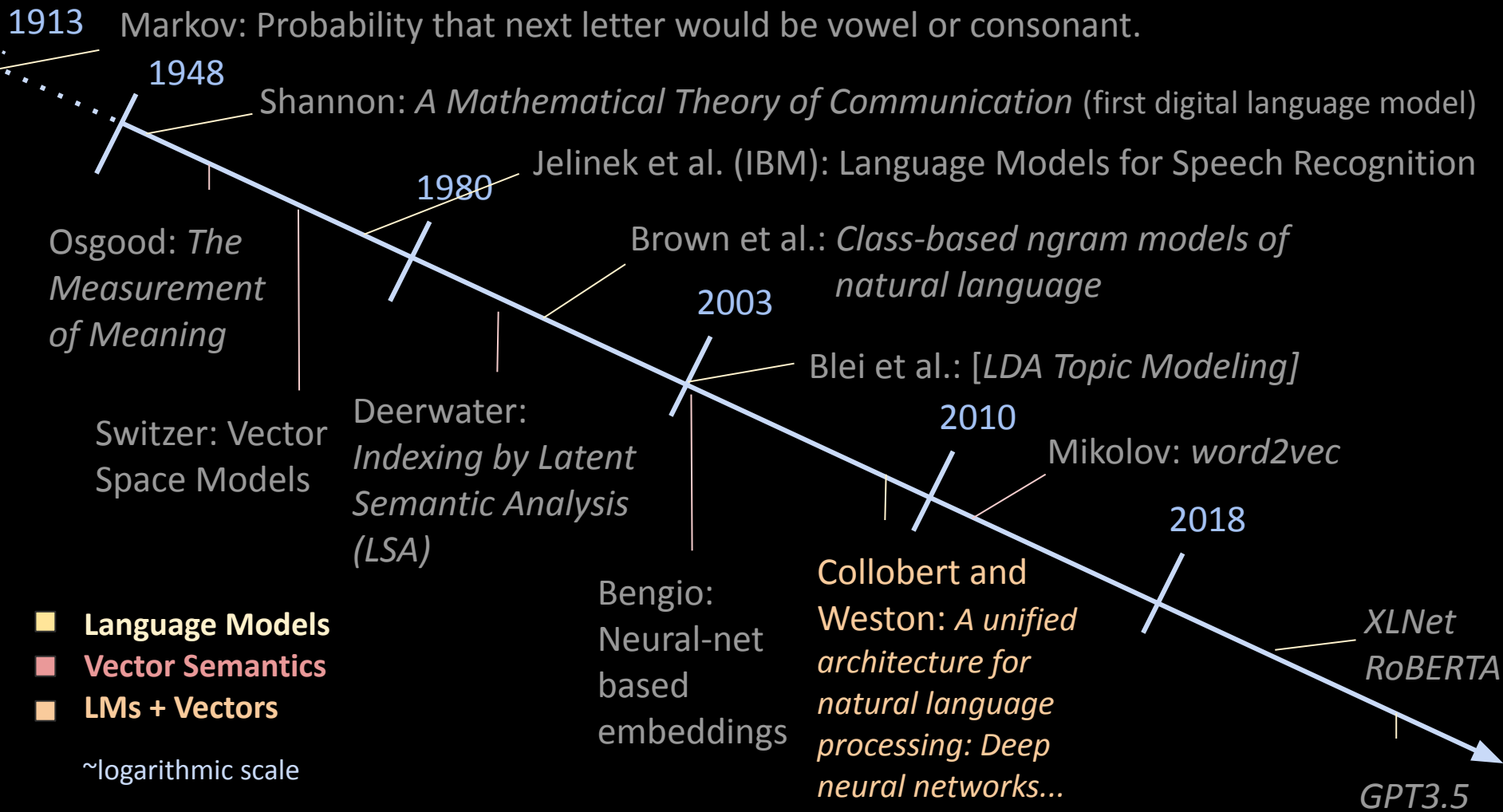
big, enceinte, expectant, gravid, **great.a.6**, large, heavy, with child (in an advanced stage of pregnancy)

**great.n.1** (a person who has achieved distinction and honor in some field)

# Timeline: *Language Modeling* and *Vector Semantics*



# Timeline: *Language Modeling* and *Vector Semantics*





# Modeling and Vector Semantics

... would be vowel or consonant.

... Theory of Communication (first digital language model)

... et al. (IBM): Language Models for Speech Recognition

... Brown et al.: *Class-based ngram models of natural language*

2003

Blei et al.: [LDA Topic Modeling]

2010

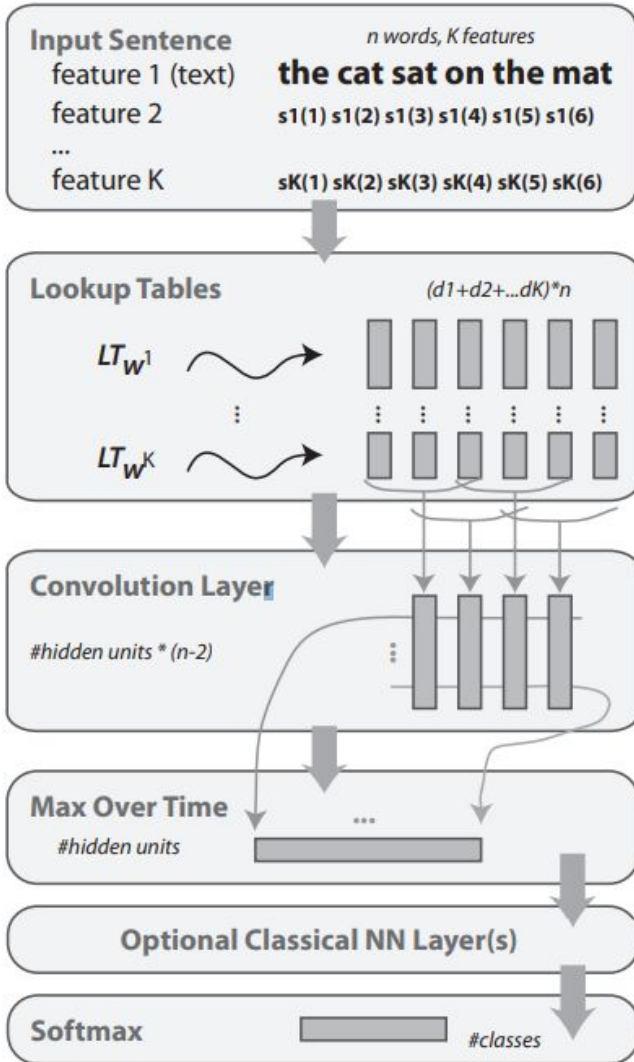
Mikolov: *word2vec*

2018

Collobert and Weston: *A unified architecture for natural language processing: Deep neural networks...*

XLNet  
RoBERTA

GPT3



1913

Os  
M  
of



logarithmic scale

# Modeling and Vector Semantics

r would be vowel or consonant.

l Theory of  
et al. (IBM)  
rown et al  
Speech Recognition

Models of

200

Blei et al.: [LDA Topic Modeling]

2010

Mikolov: word2vec

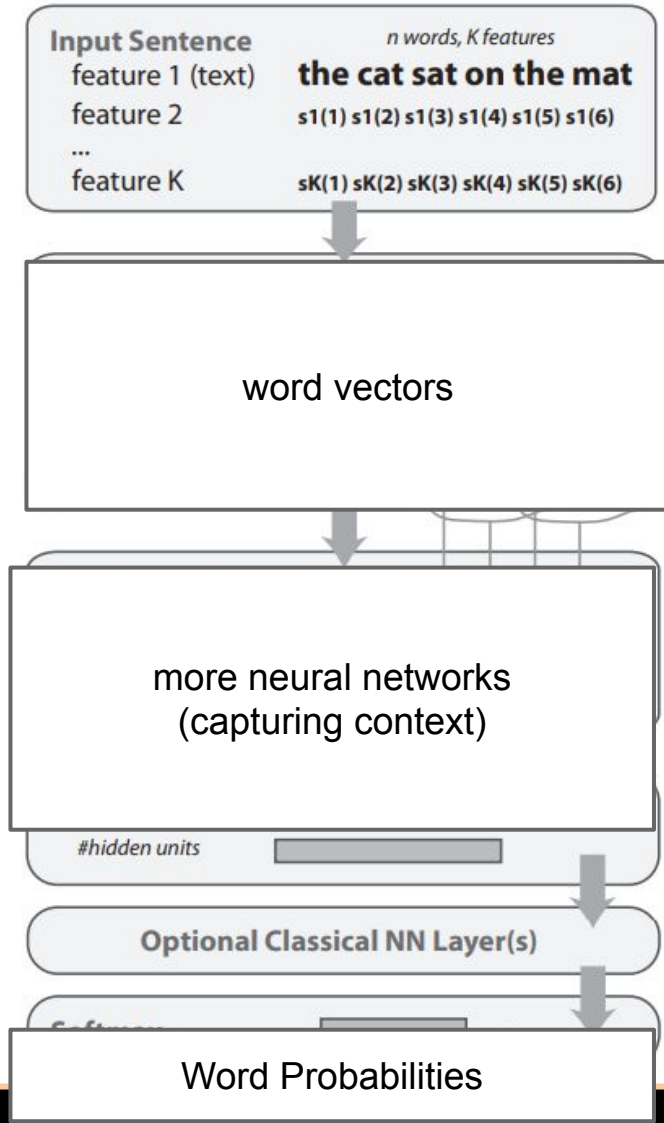
2018

Collobert and Weston: A unified architecture for natural language processing: Deep neural networks...

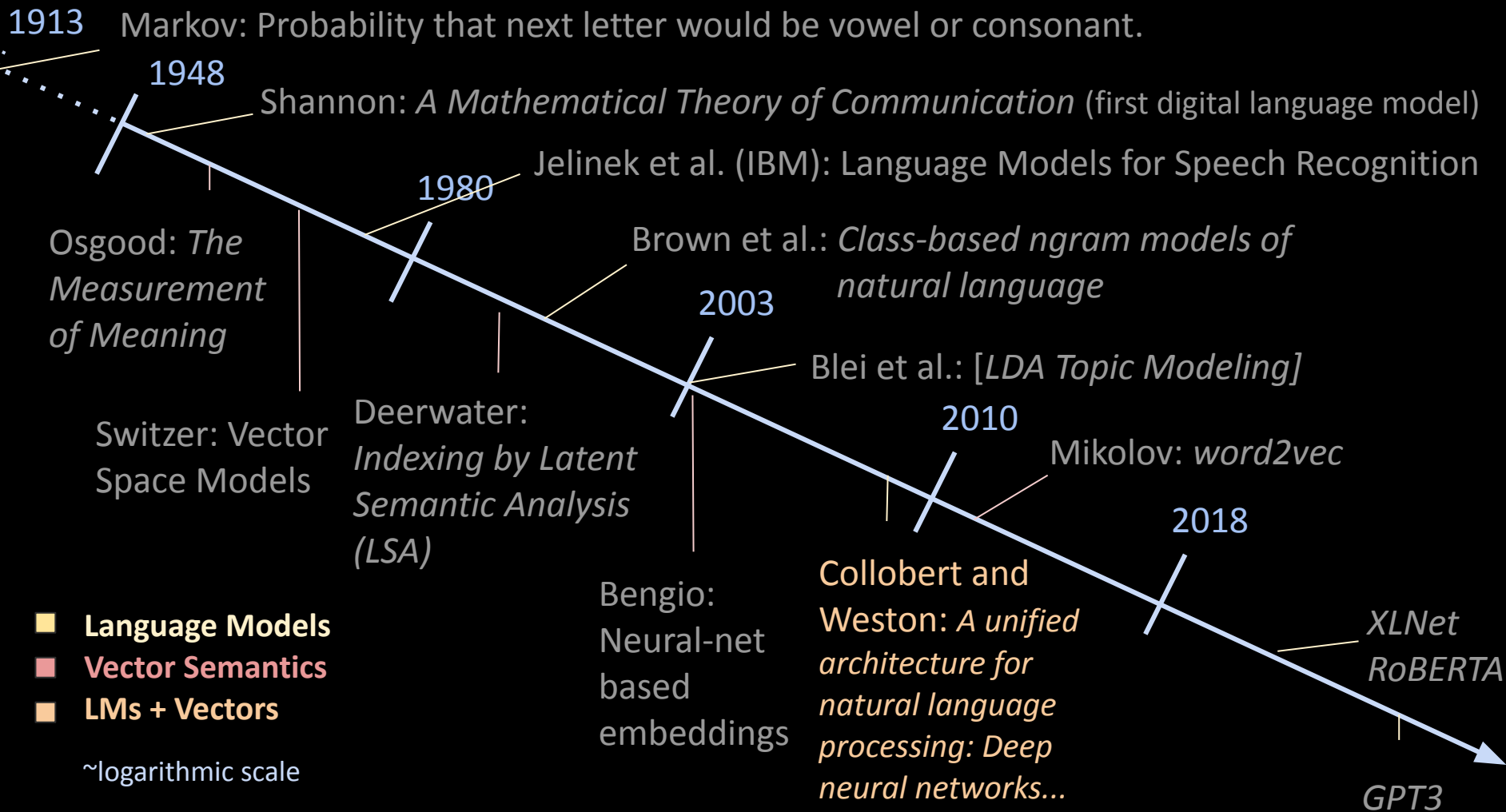
XLNet  
RoBERTA

GPT3

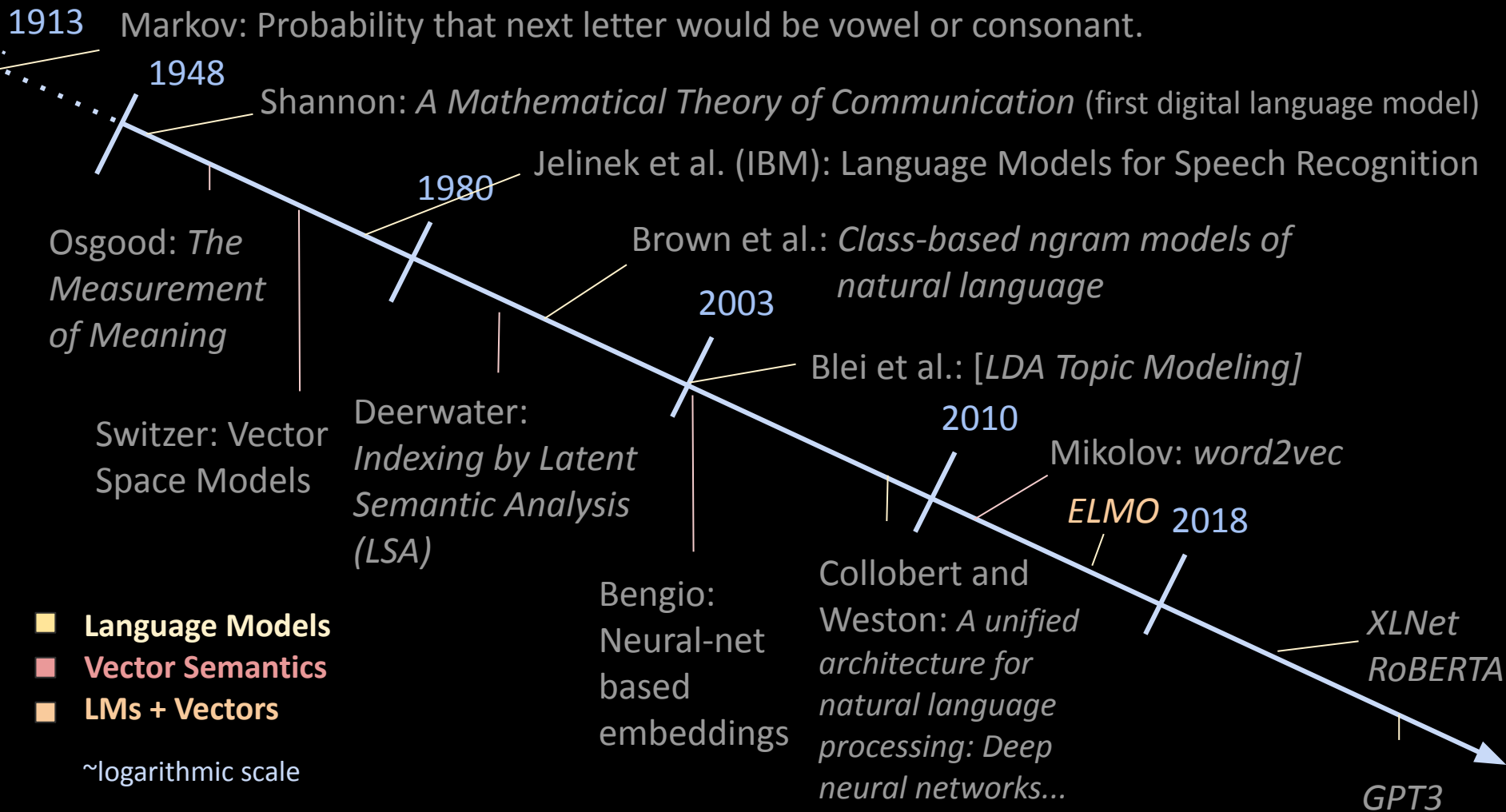
Deep Learning  
for Language  
Modeling!  
(time for break)



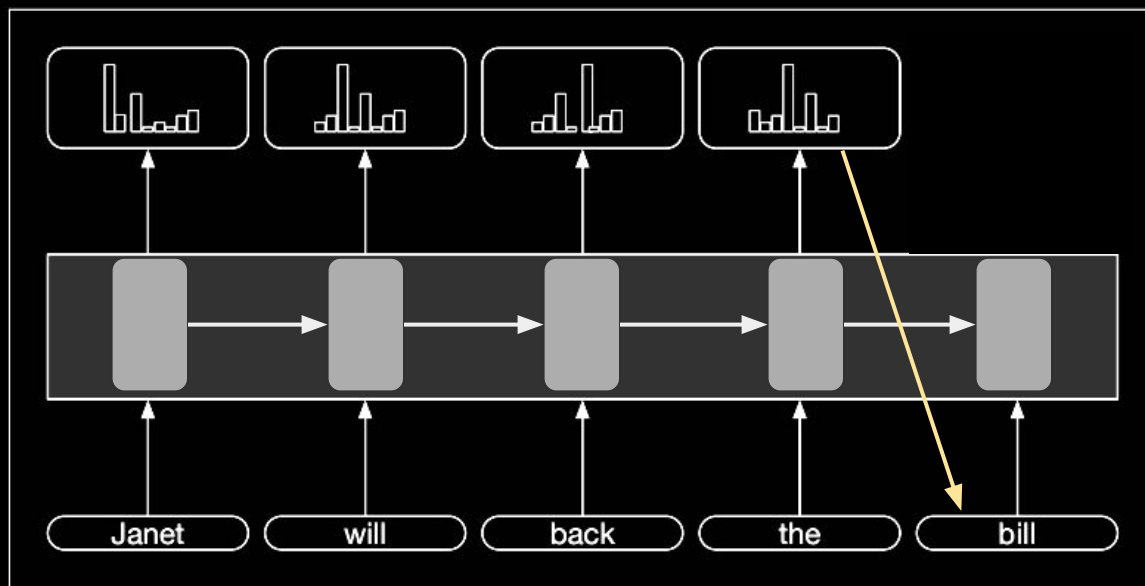
# Timeline: *Language Modeling* and *Vector Semantics*



# Timeline: *Language Modeling* and *Vector Semantics*



# Recurrent Neural Network

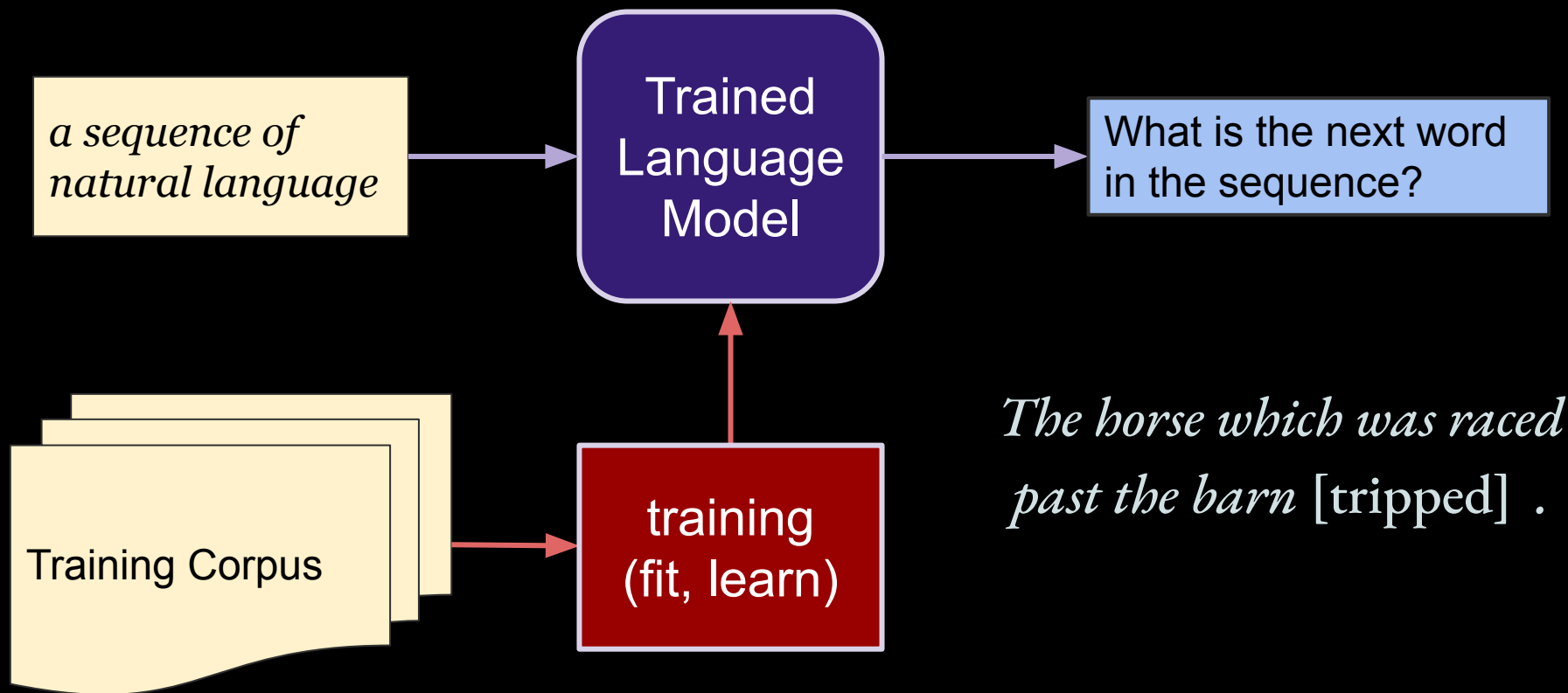


Language modeling  
with an RNN

# Neural Networks

# Language Modeling

Building a model (or system / API) that can answer the following:

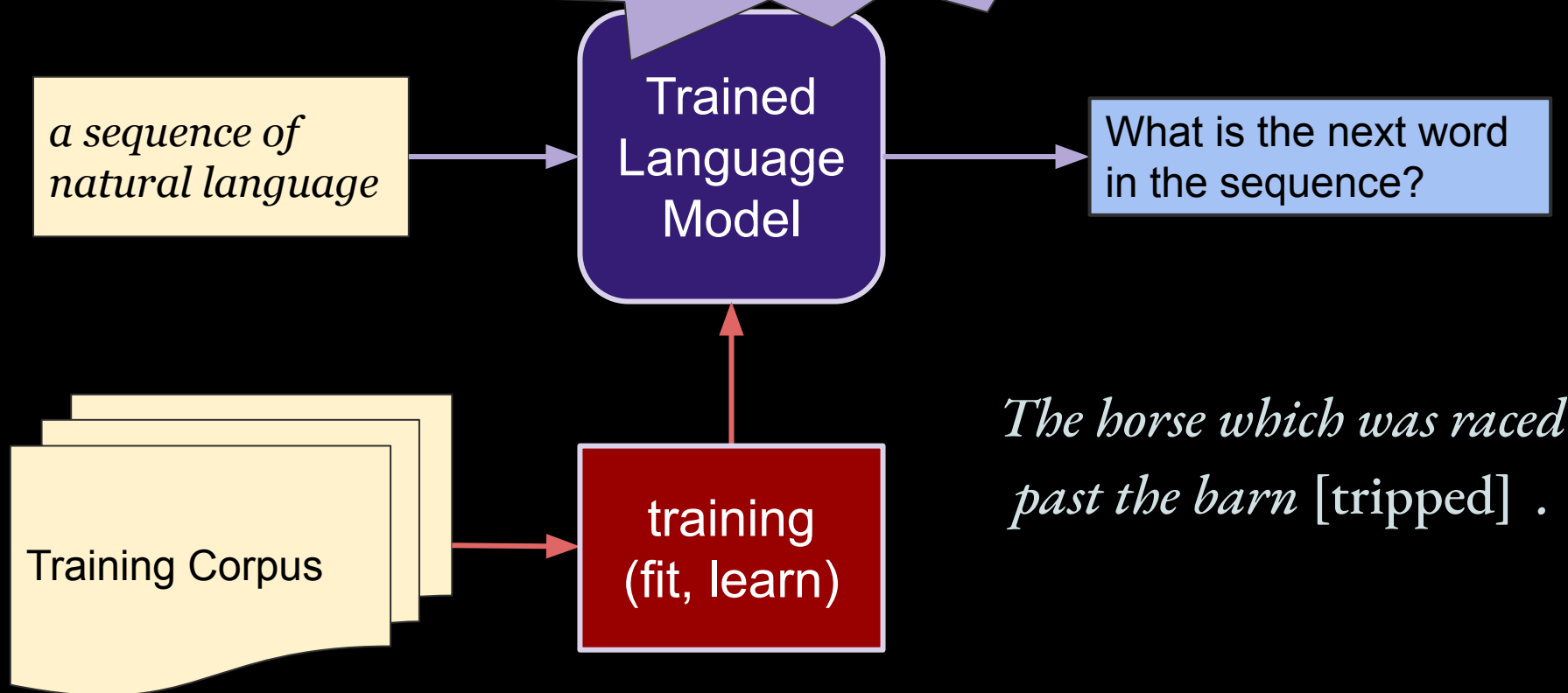


# Language Modeling

Building a model (or a

To fully capture natural language, models get very complex!

er the following:

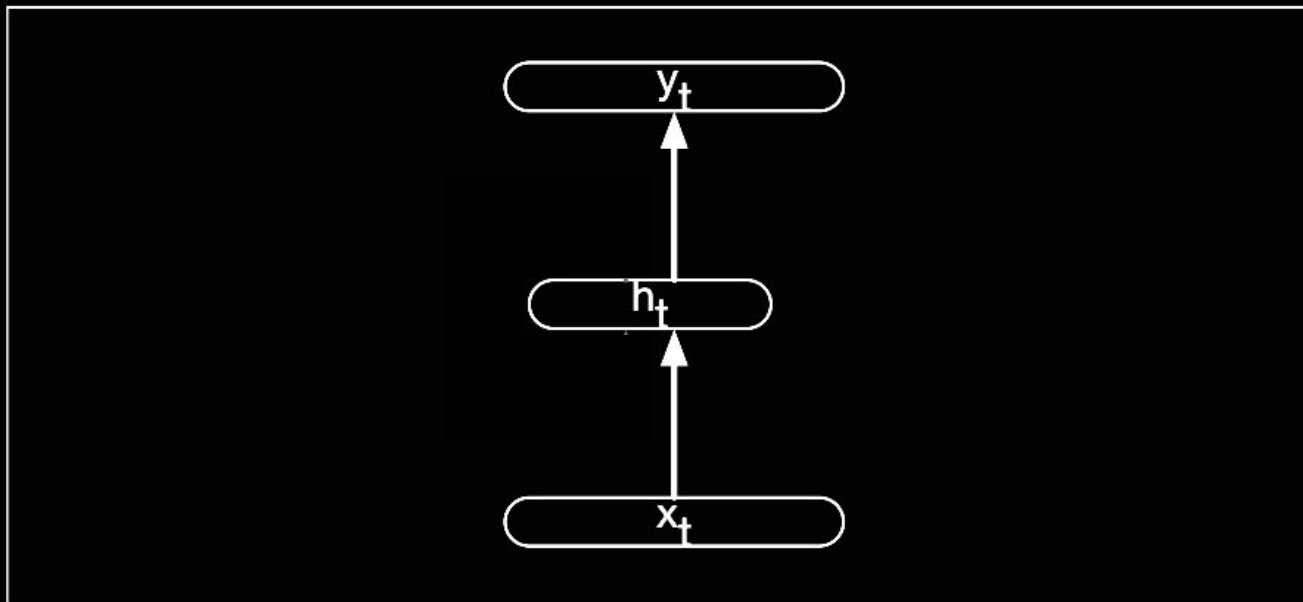


*The horse which was raced  
past the barn [tripped] .*



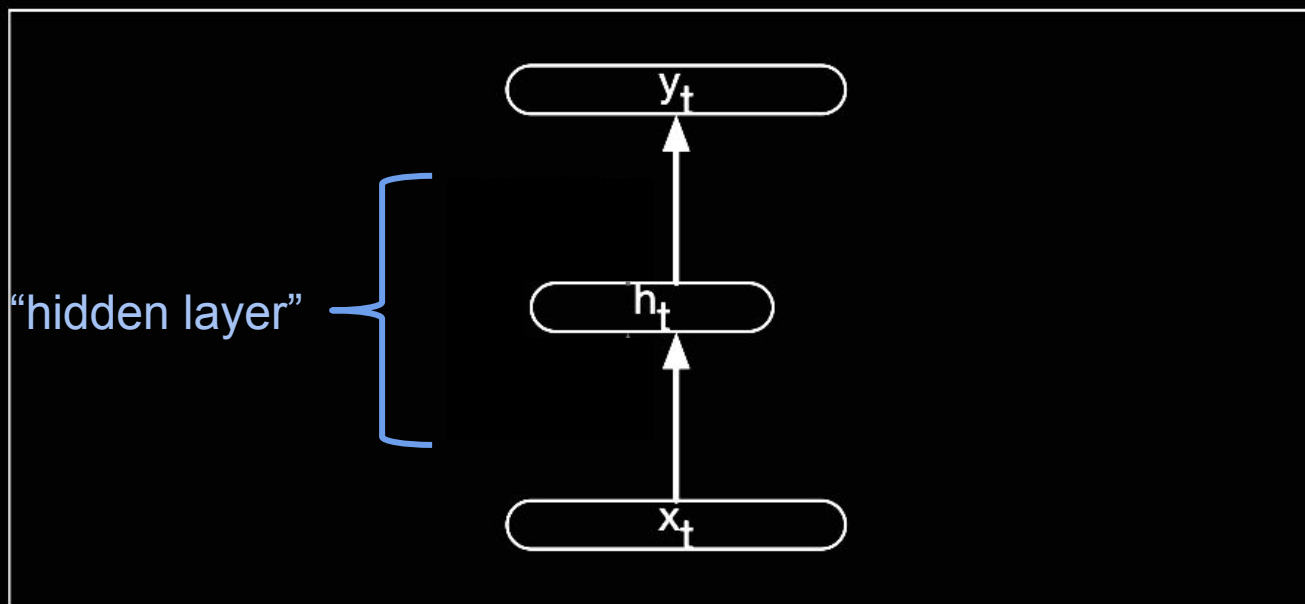
# Neural Networks: Graphs of Operations

# Neural Networks: Graphs of Operations (excluding the optimization nodes)



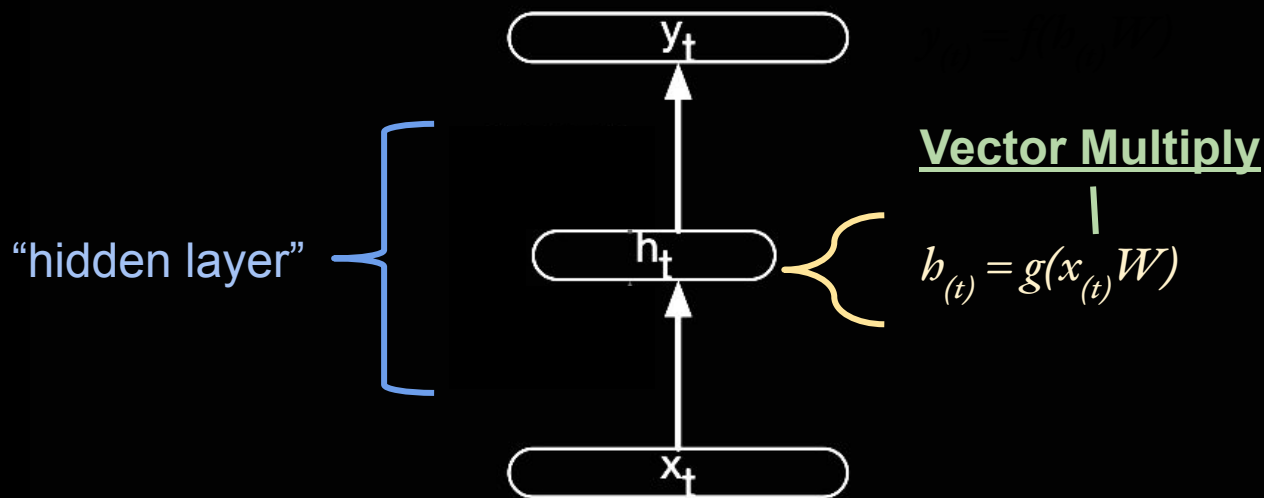
**Figure 9.2** Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)

# Neural Networks: Graphs of Operations (excluding the optimization nodes)



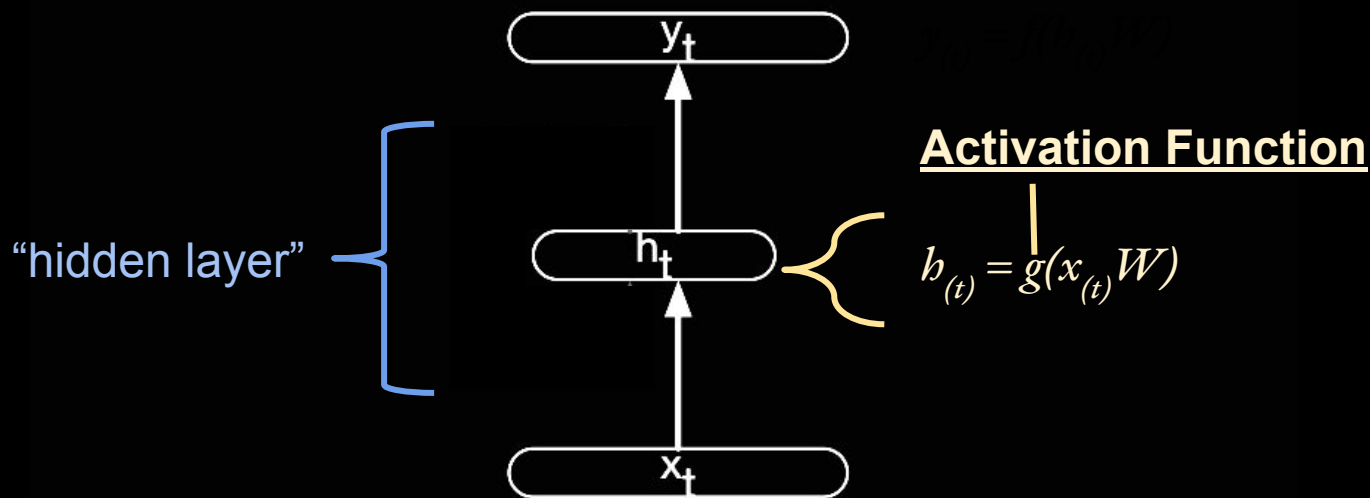
**Figure 9.2** Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)

# Neural Networks: Graphs of Operations (excluding the optimization nodes)



**Figure 9.2** Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)

# Neural Networks: Graphs of Operations (excluding the optimization nodes)



**Figure 9.2** Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)

# Common Activation Functions

$$z = h_{(t)}W$$

Logistic:  $\sigma(z) = 1 / (1 + e^{-z})$

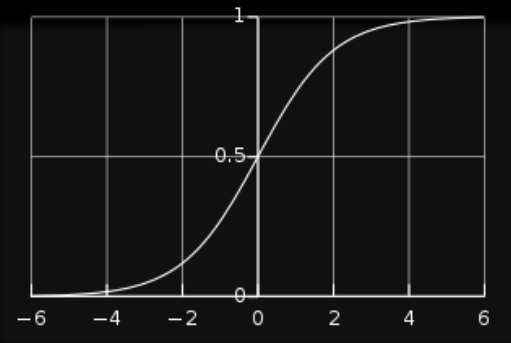
Hyperbolic tangent:  $\tanh(z) = 2\sigma(2z) - 1 = (e^{2z} - 1) / (e^{2z} + 1)$

Rectified linear unit (ReLU):  $ReLU(z) = \max(0, z)$

# Common Activation Functions

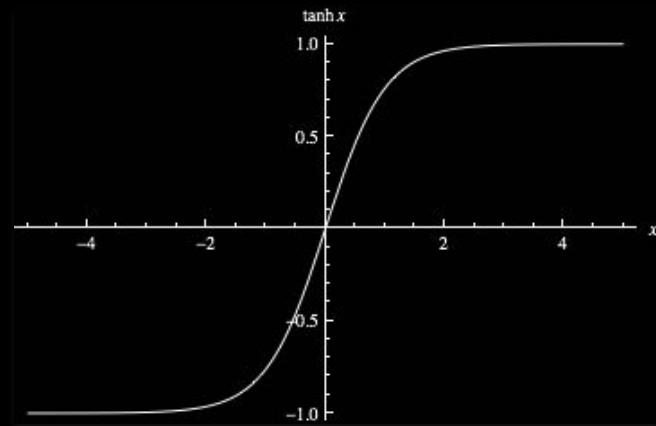
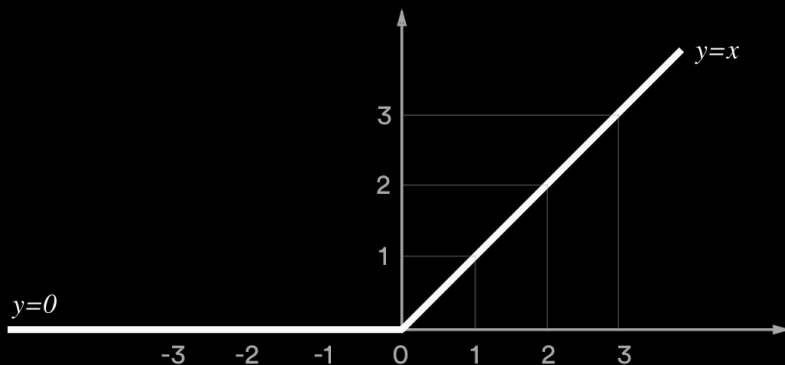
$$z = h_{(t)}W$$

Logistic:  $\sigma(z) = 1 / (1 + e^{-z})$

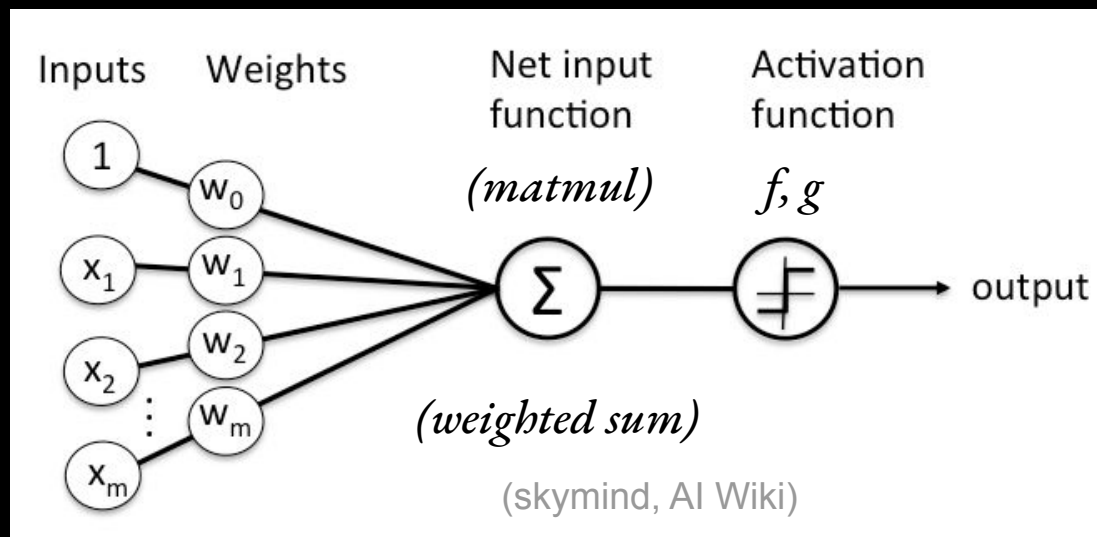


Hyperbolic tangent:  $\tanh(z) = 2\sigma(2z) - 1 = (e^{2z} - 1) / (e^{2z} + 1)$

Rectified linear unit (ReLU):  $ReLU(z) = \max(0, z)$



# Neural Networks: Graphs of Operations (excluding the optimization nodes)



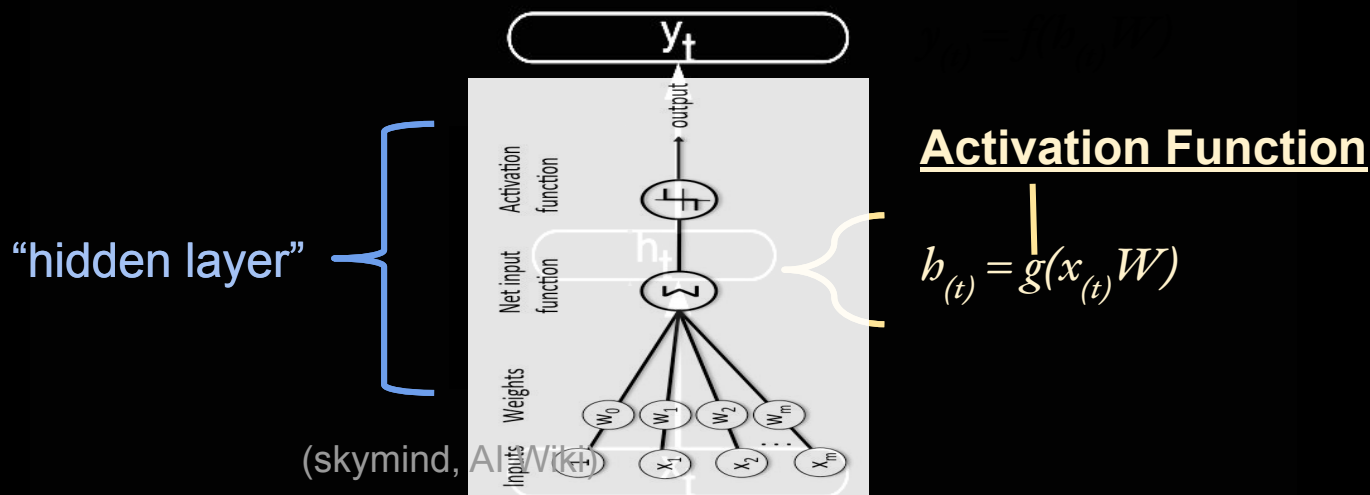
## Activation Function

$$h_{(t)} = g(x_{(t)}W)$$

**Figure 9.2** Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)

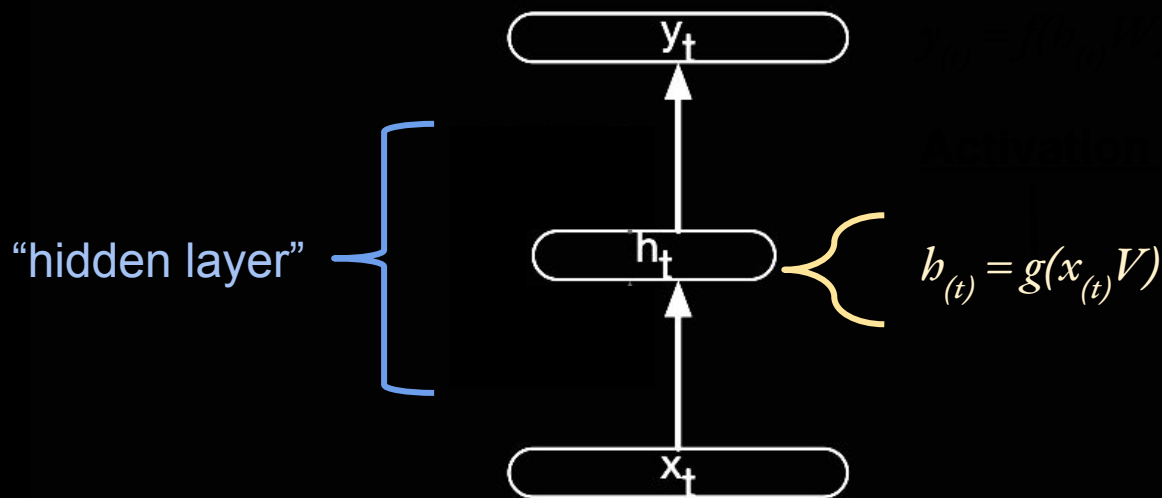


# Neural Networks: Graphs of Operations (excluding the optimization nodes)



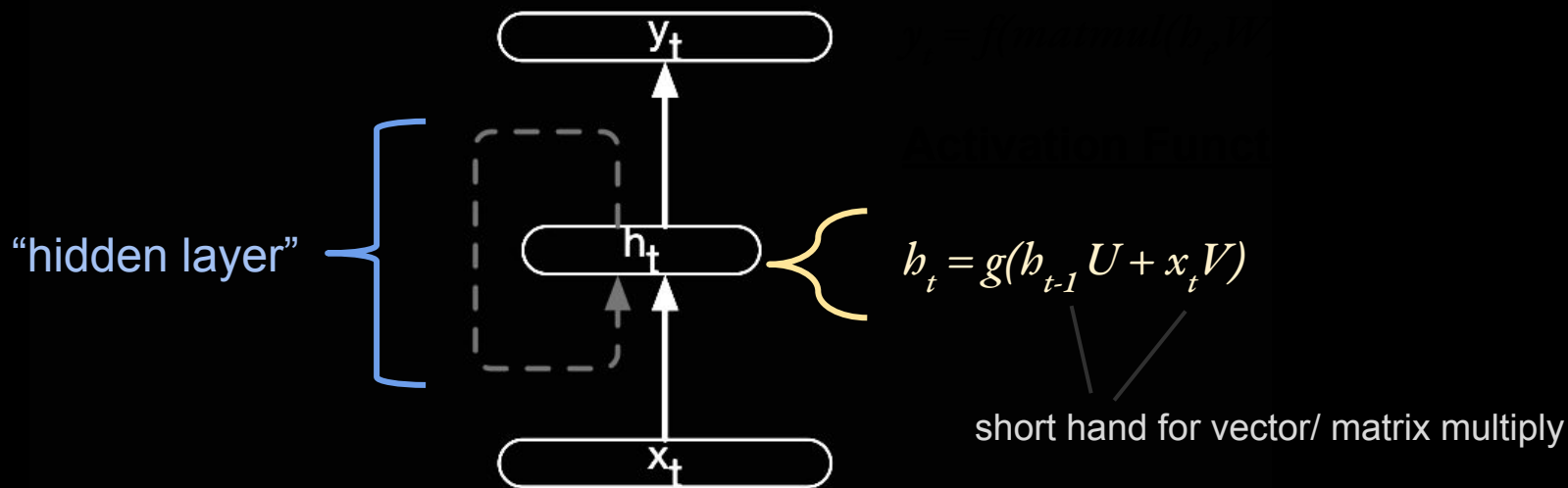
**Figure 9.2** Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)

# Neural Networks: Graphs of Operations (excluding the optimization nodes)



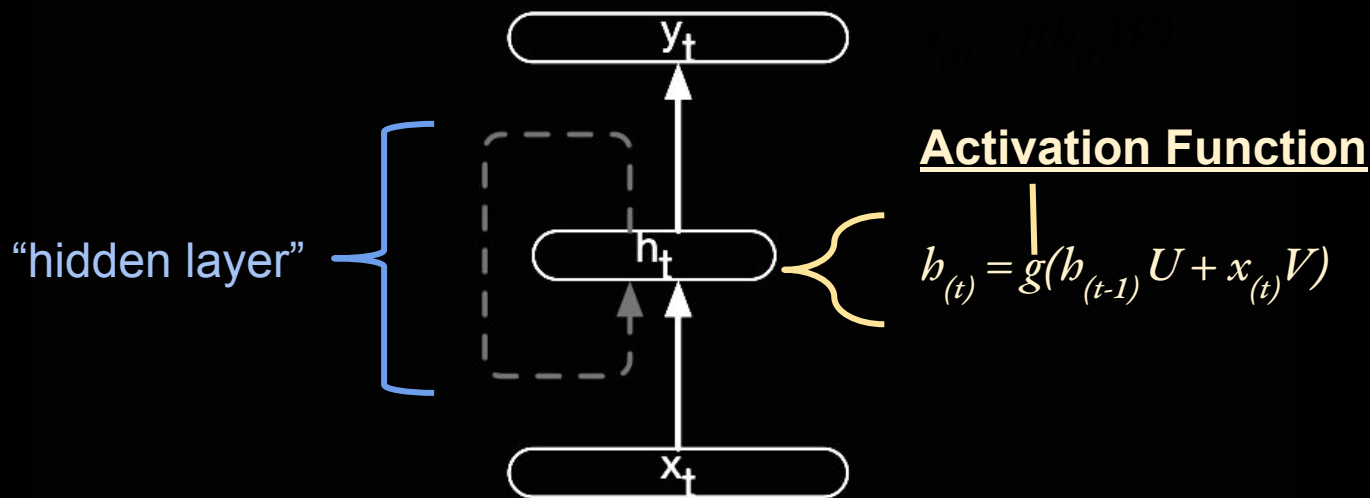
**Figure 9.2** Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)

# Neural Networks: Graphs of Operations (excluding the optimization nodes)



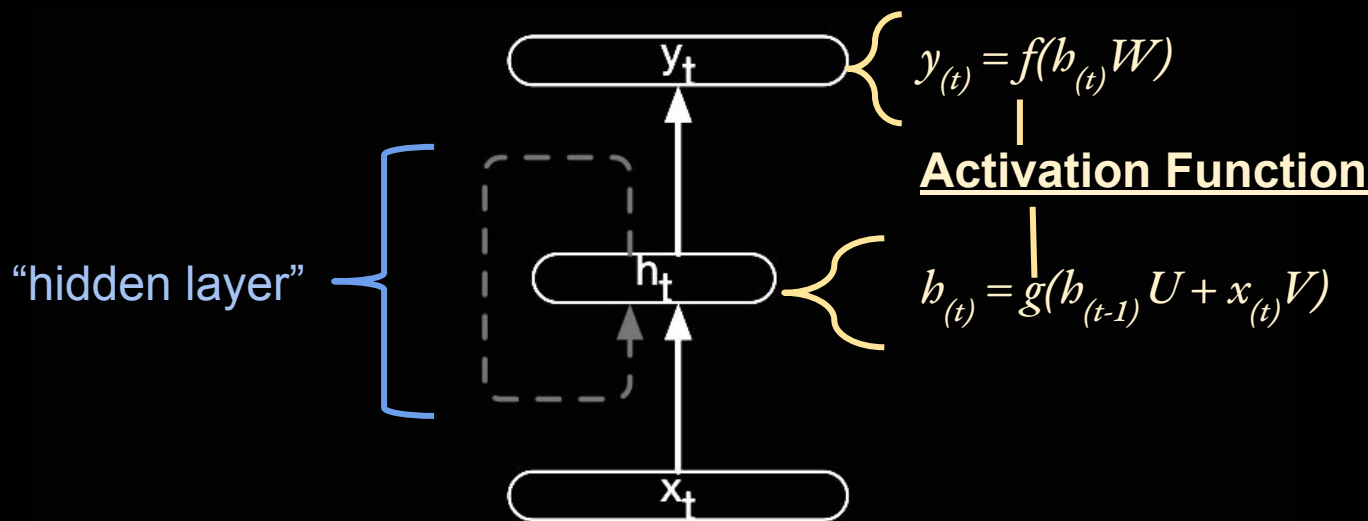
**Figure 9.2** Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)

# Neural Networks: Graphs of Operations (excluding the optimization nodes)



**Figure 9.2** Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)

# Neural Networks: Graphs of Operations (excluding the optimization nodes)

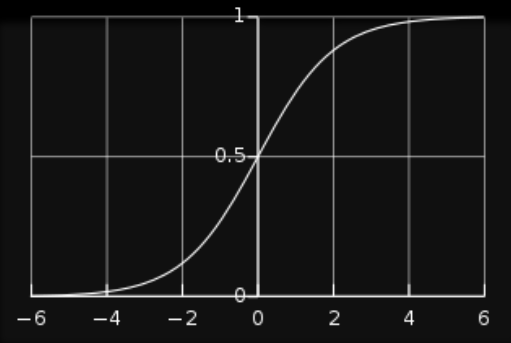


**Figure 9.2** Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)

# Common Activation Functions

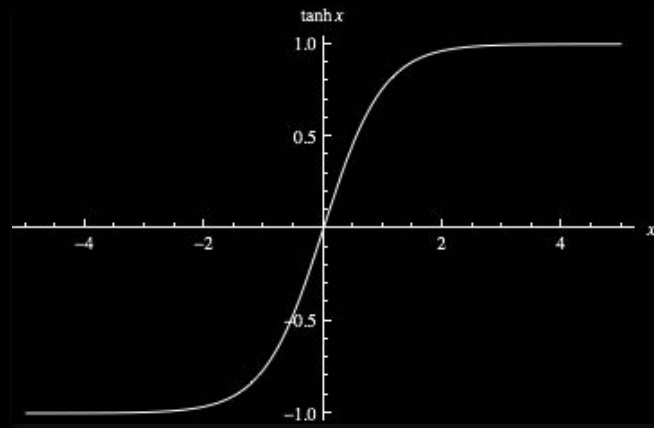
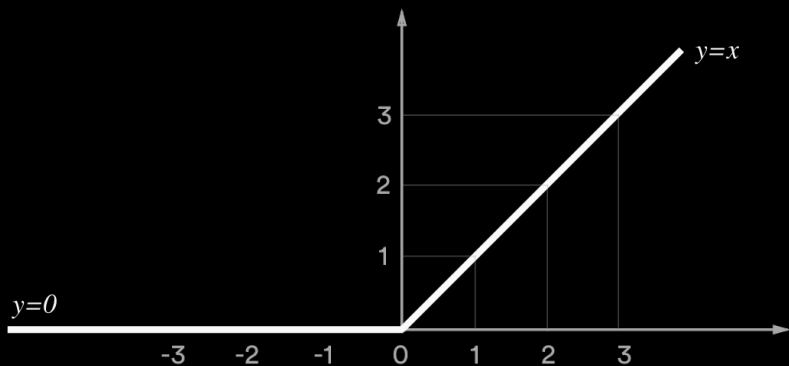
$$z = h_{(t)}W$$

Logistic:  $\sigma(z) = 1 / (1 + e^{-z})$



Hyperbolic tangent:  $\tanh(z) = 2\sigma(2z) - 1 = (e^{2z} - 1) / (e^{2z} + 1)$

Rectified linear unit (ReLU):  $ReLU(z) = \max(0, z)$



# Example: Forward Pass



(Geron, 2017)

```
#define forward pass graph:
```

```
 $h_{(0)} = 0$ 
```

```
for i in range(1, len(x)):
```

```
     $h_{(i)} = g(U h_{(i-1)} + W x_{(i)})$  #update hidden state
```

```
     $y_{(i)} = f(V h_{(i)})$  #update output
```

# Example: Forward Pass



*#define forward pass graph:*

$h_{(0)} = 0$

for  $i$  in range(1, len( $x$ )):

$h_{(i)} = g(U h_{(i-1)} + W x_{(i)})$  *#update hidden state*

$y_{(i)} = f(V h_{(i)})$  *#update output*



# Example: Forward Pass



```
#define forward pass graph:
```

```
 $h_{(0)} = 0$ 
```

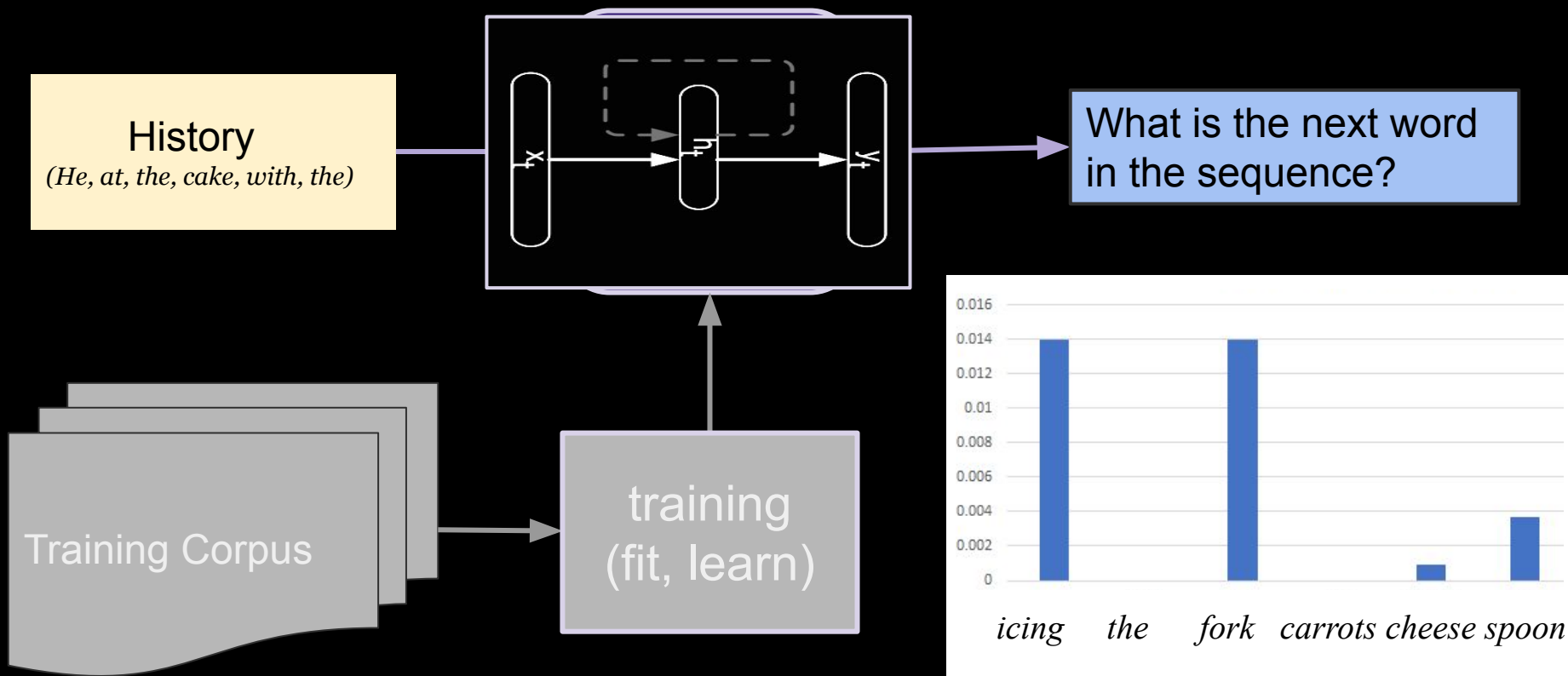
```
for i in range(1, len(x)):
```

```
     $h_{(i)} = \tanh(\text{matmul}(U, h_{(i-1)}) + \text{matmul}(W, x_{(i)}))$  #update hidden state
```

```
     $y_{(i)} = \text{softmax}(\text{matmul}(V, h_{(i)}))$  #update output
```

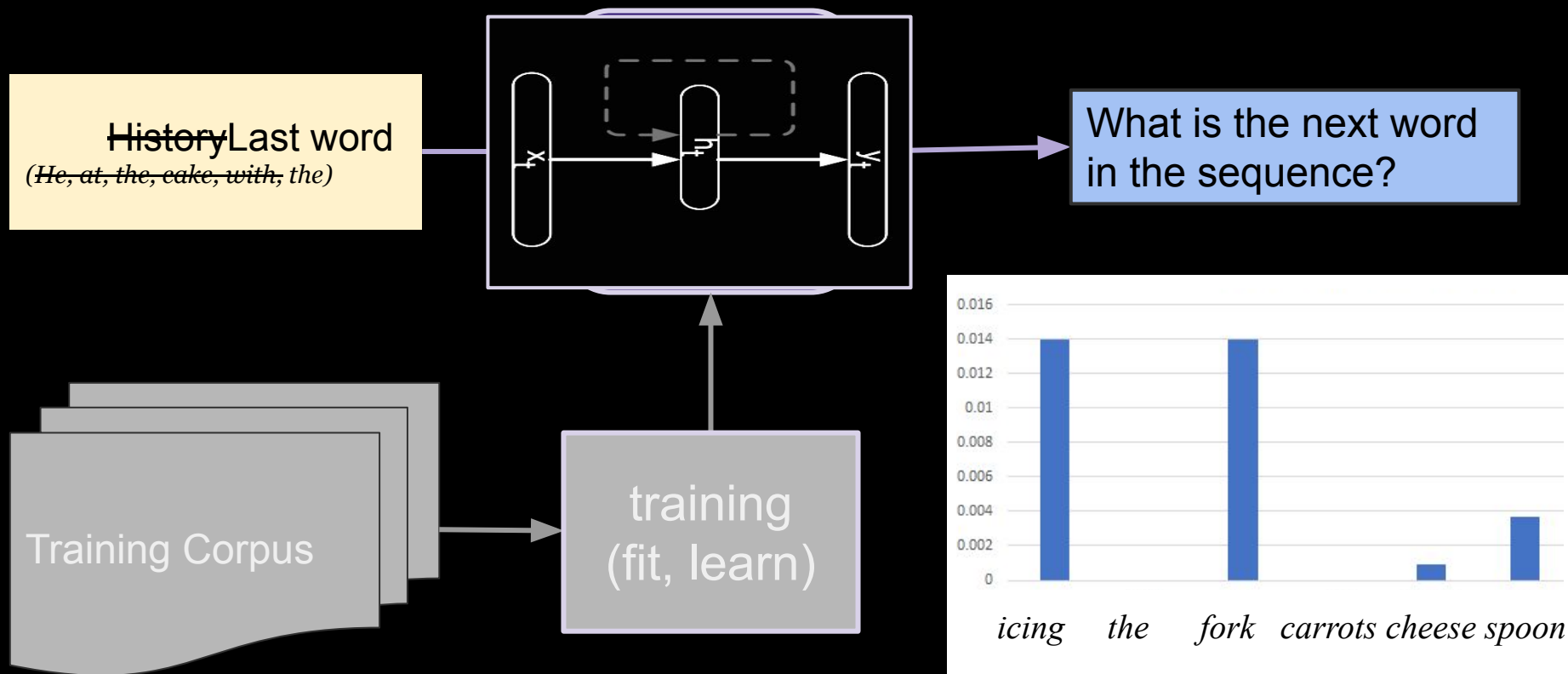
# Language Modeling

*Task: Estimate  $P(w_n | w_1, w_2, \dots, w_{n-1})$*   
:probability of a next word given history  
 *$P(\text{fork} | \text{He ate the cake with the}) = ?$*



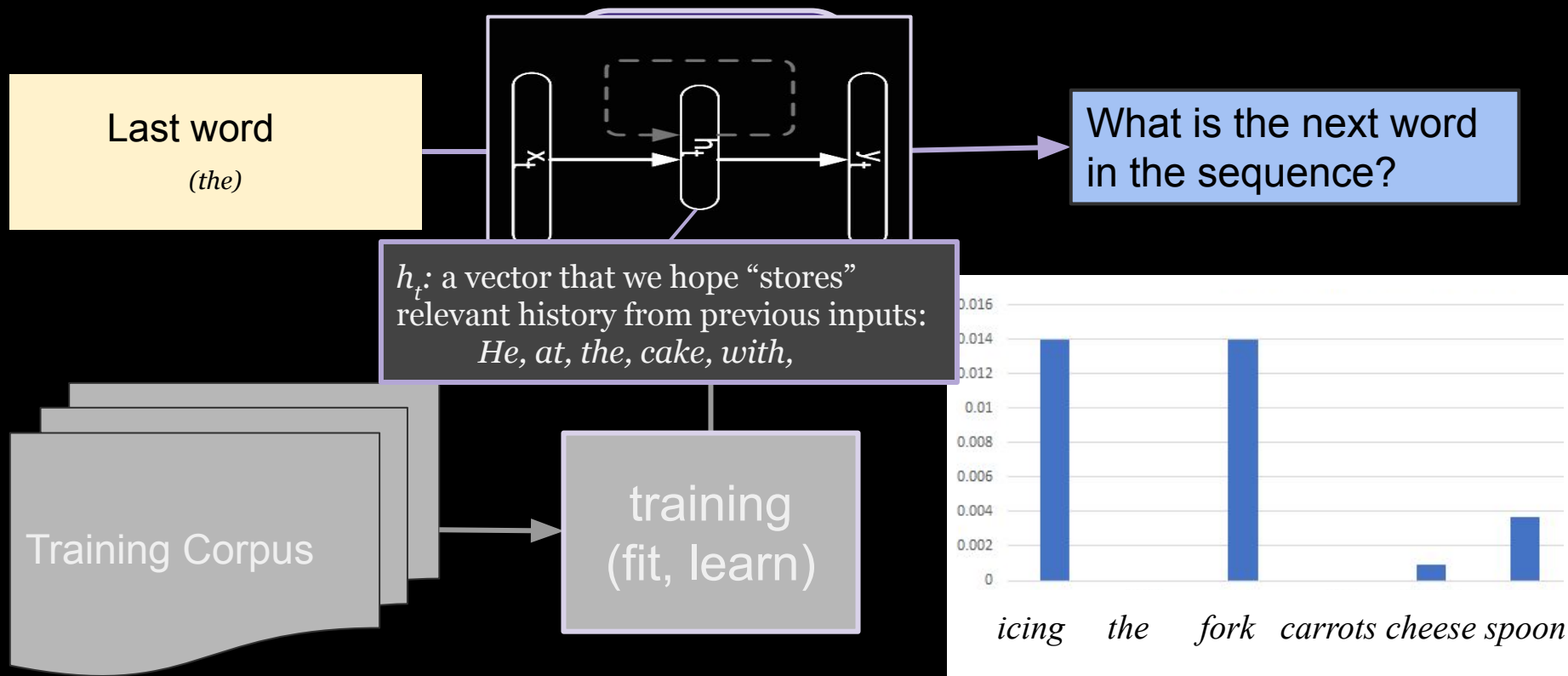
# Language Modeling

*Task: Estimate  $P(w_n | w_1, w_2, \dots, w_{n-1})$*   
:probability of a next word given history  
 *$P(\text{fork} | \text{He ate the cake with the}) = ?$*

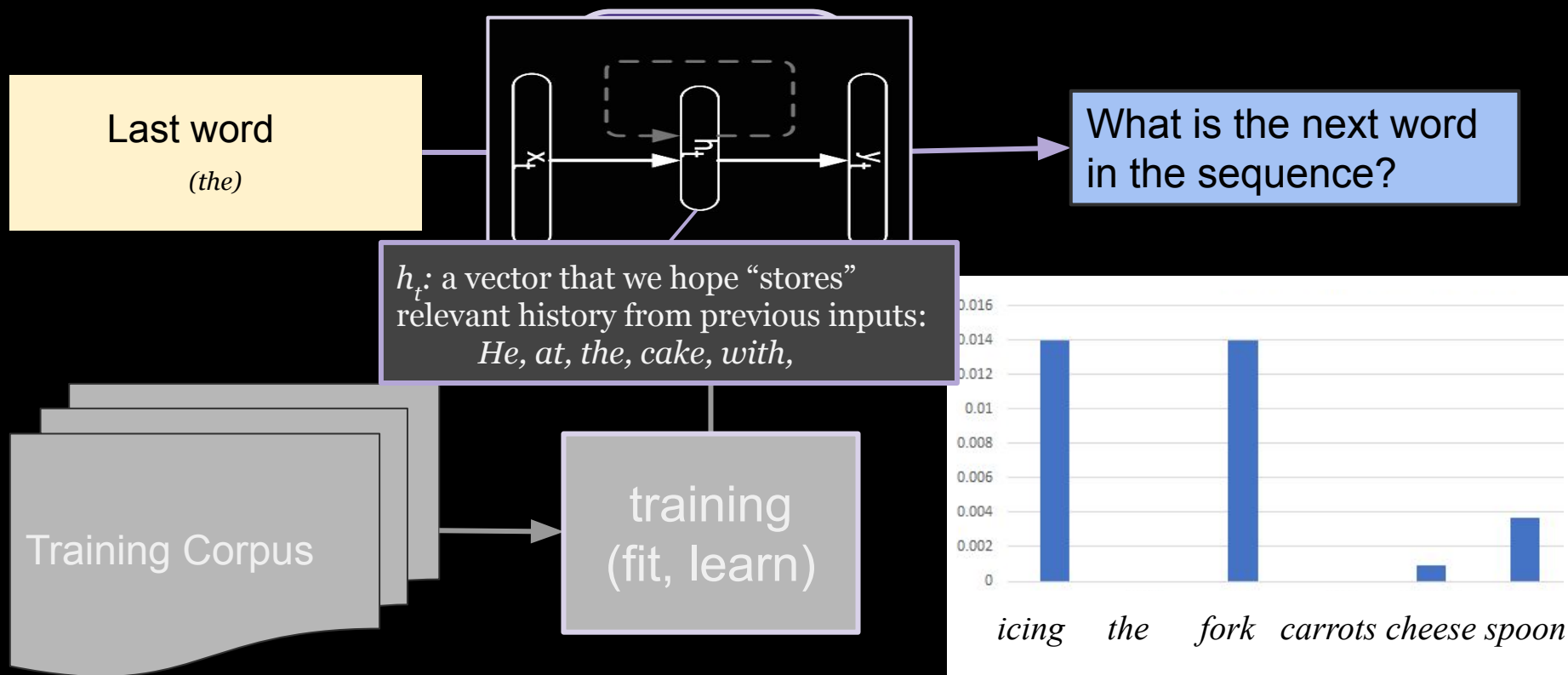


# Language Modeling

*Task: Estimate  $P(w_n | w_1, w_2, \dots, w_{n-1})$*   
:probability of a next word given history  
 *$P(\text{fork} | \text{He ate the cake with the}) = ?$*



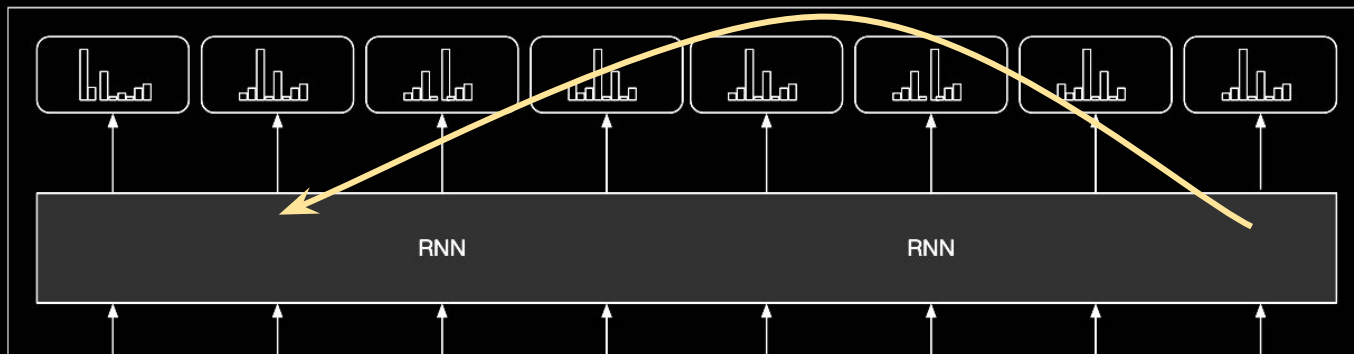
# RNN Limitation: Not para



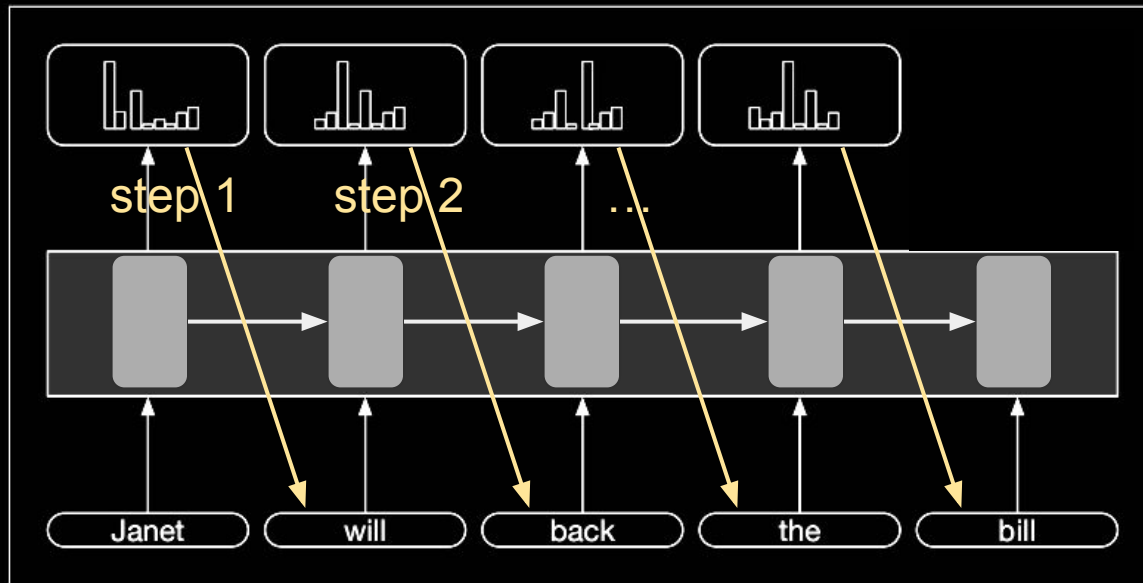
# RNN: Limitation:

## Losing Track of Long Distance Dependencies

*The horse which was raced past the barn tripped.*



# RNN: Limitation: Not parallelizable



Language modeling  
with an RNN

# The Transformer: Motivation

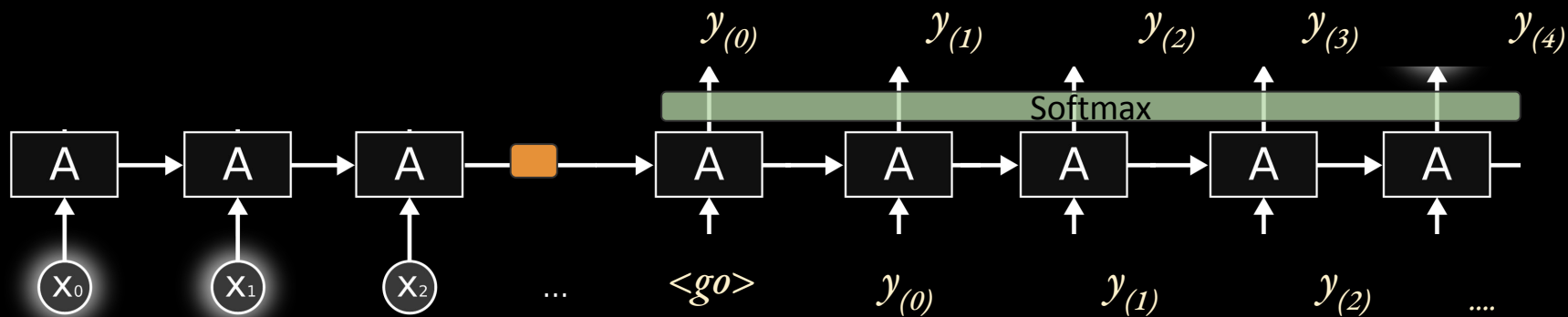
---

## Challenges to sequential representation learning

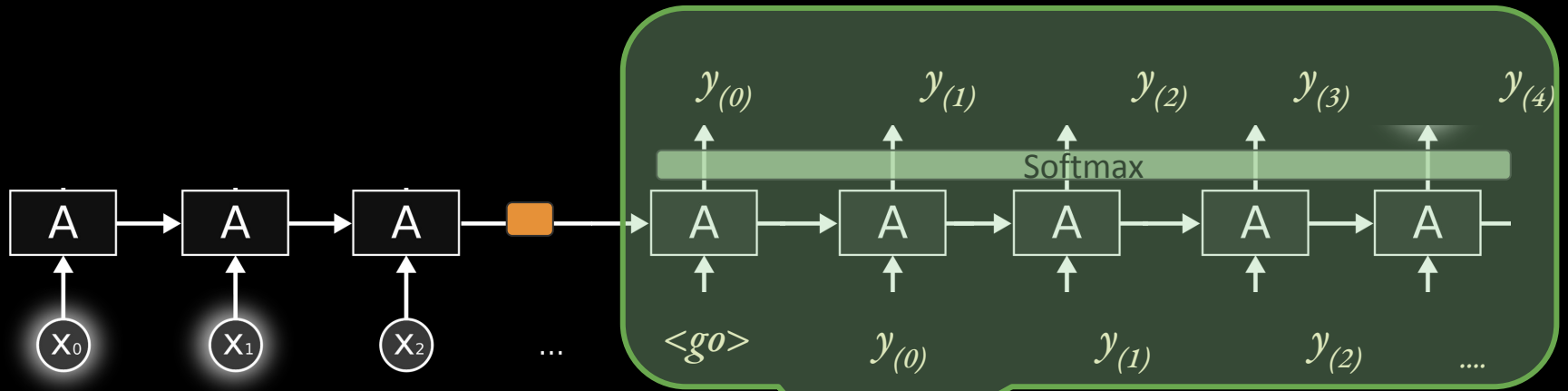
- Capture long-distance dependencies
- Preserving sequential distances / periodicity
- Capture multiple relationships
- Easy to parallelize -- don't need sequential processing.



# Encoder-Decoder (Simpler Representation)

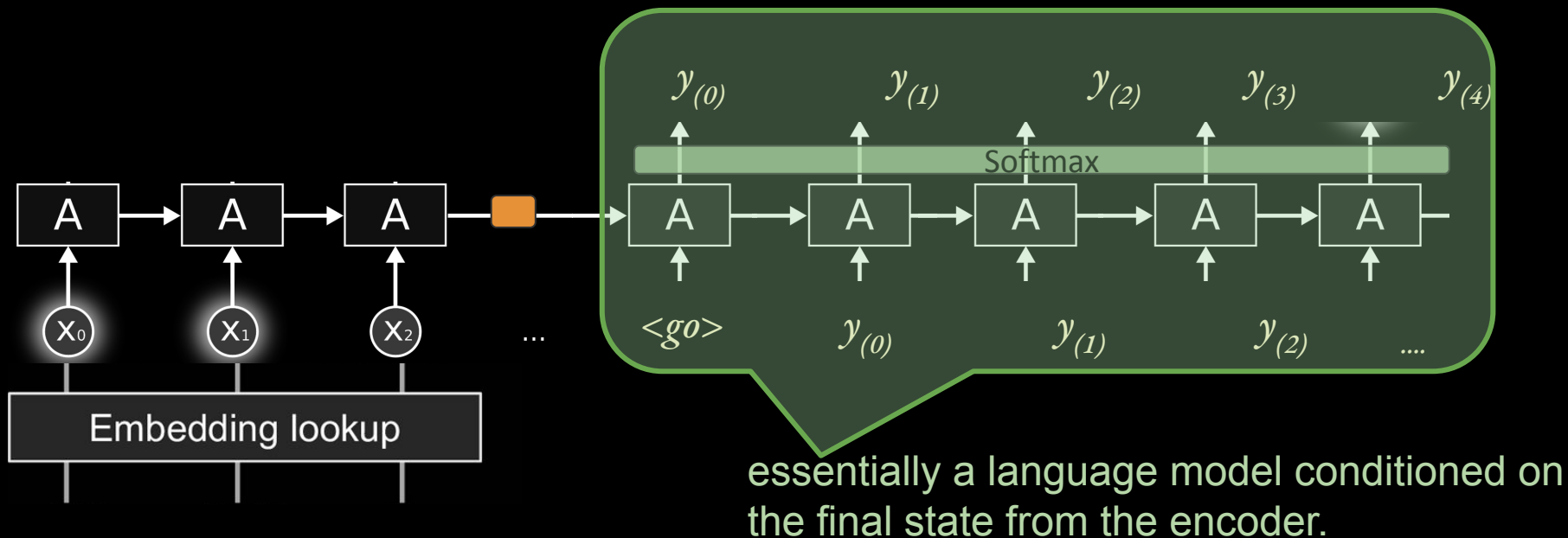


# Encoder-Decoder (Simpler Representation)

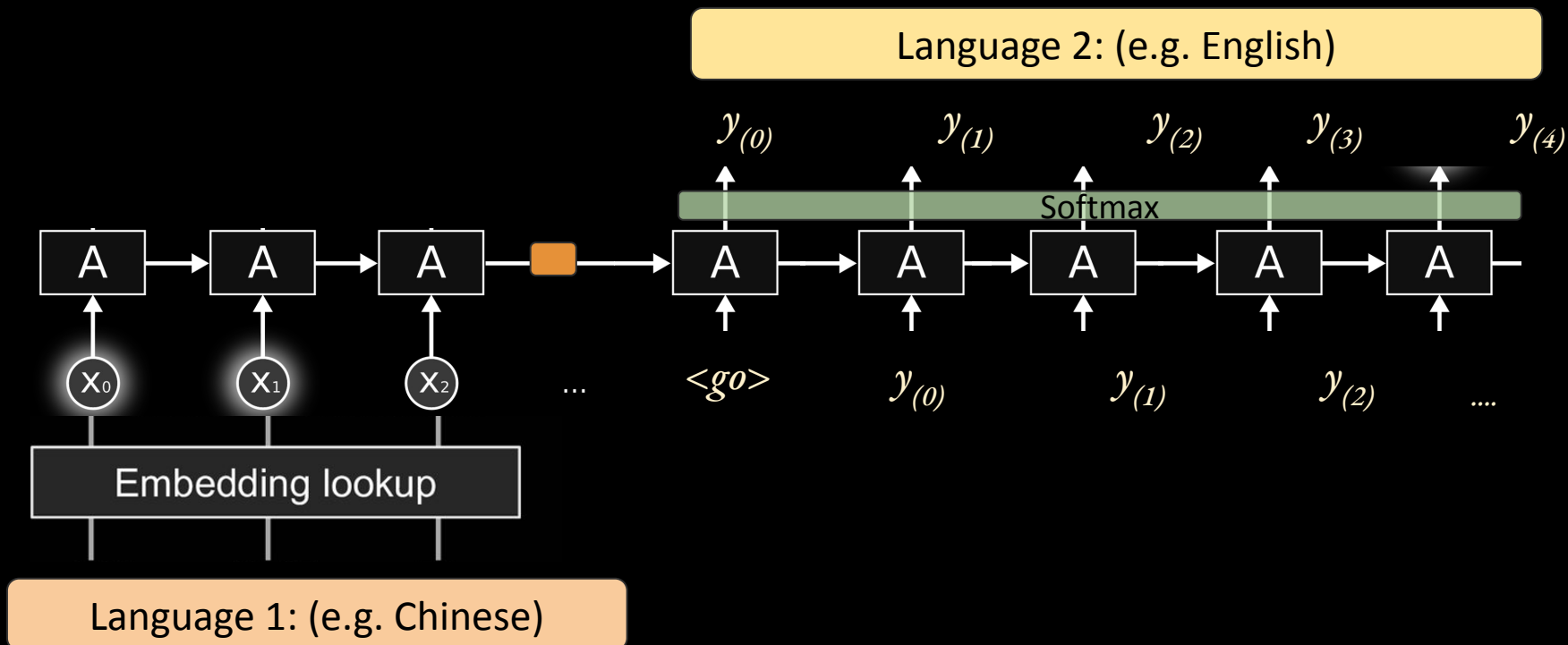


essentially a language model conditioned on the final state from the encoder.

# Encoder-Decoder (Simpler Representation)



# Encoder-Decoder (Simpler Representation)

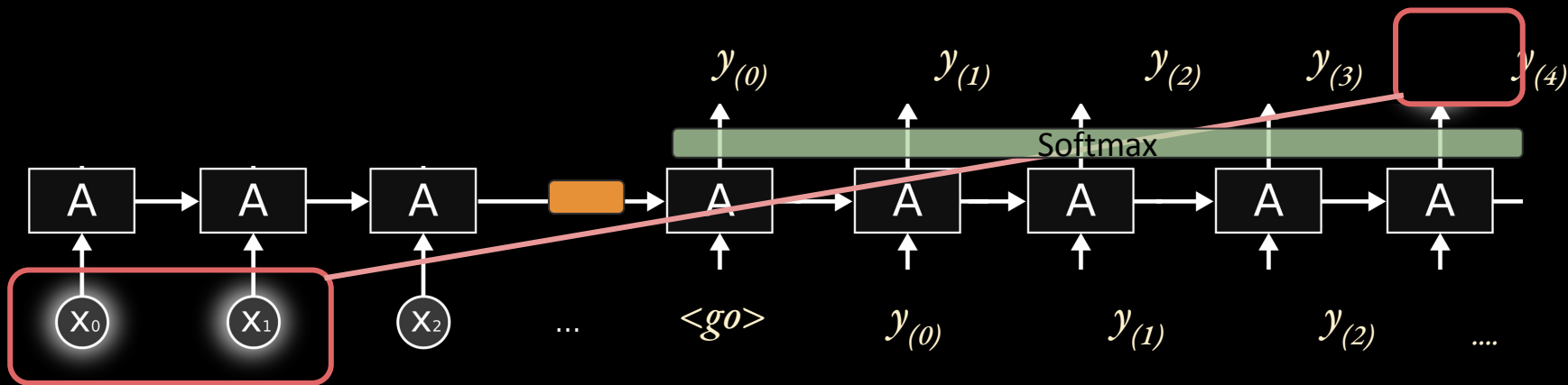


# Encoder-Decoder

Challenge:

*The ball was kicked by kayla.*

- Long distance dependency when translating:



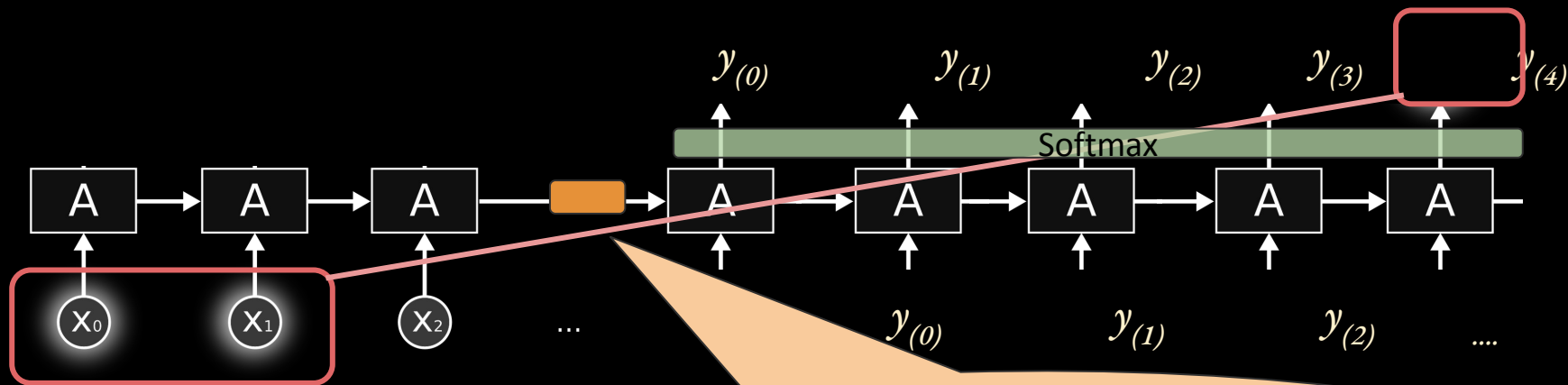
*Kayla kicked the ball.*

# Encoder-Decoder

Challenge:

*The ball was kicked by kayla.*

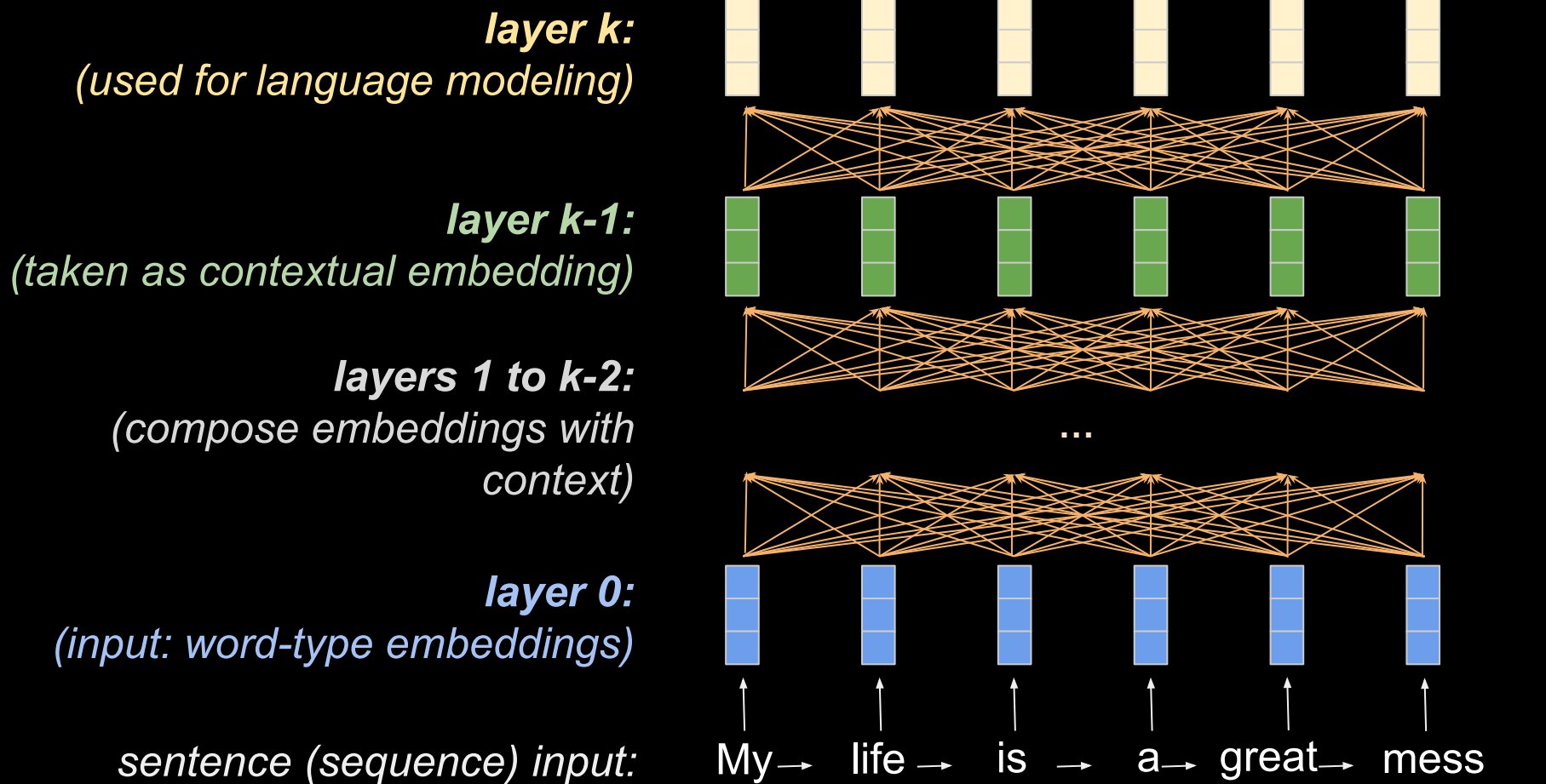
- Long distance dependency when translating:



*Kayla kicked the ball.*

A lot of responsibility put fixed-size hidden state passed from encoder to decoder

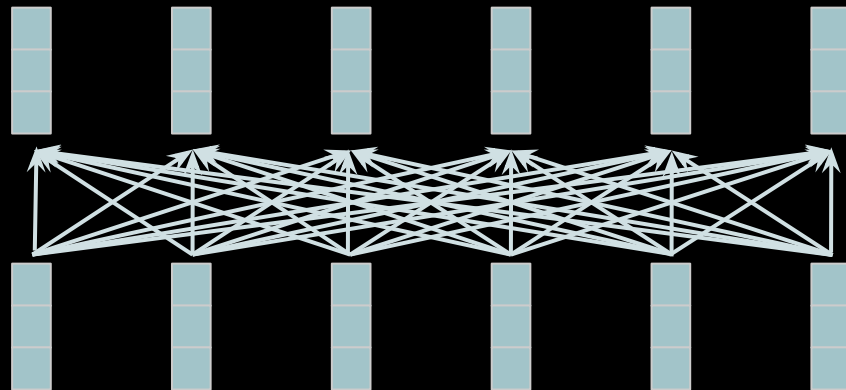
Transformer Language Models: Uses multiple layers of a **transformer**



(Kjell, Kjell, and Schwartz, 2023)

## *auto-encoder:*

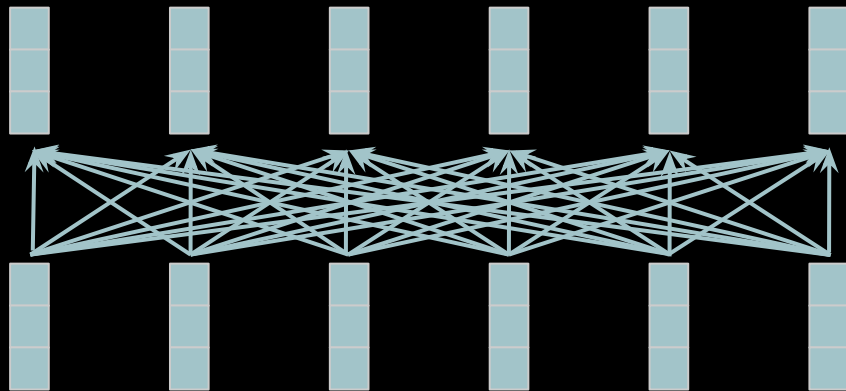
- Connections go both directions.
- Task is predict word in middle:  
 $p(w_i \mid \dots, p_{w_i-2}, w_{i-1}, w_{i+1}, w_{i+2} \dots)$
- Better for:
  - embeddings
  - fine-tuning (transfer learning)





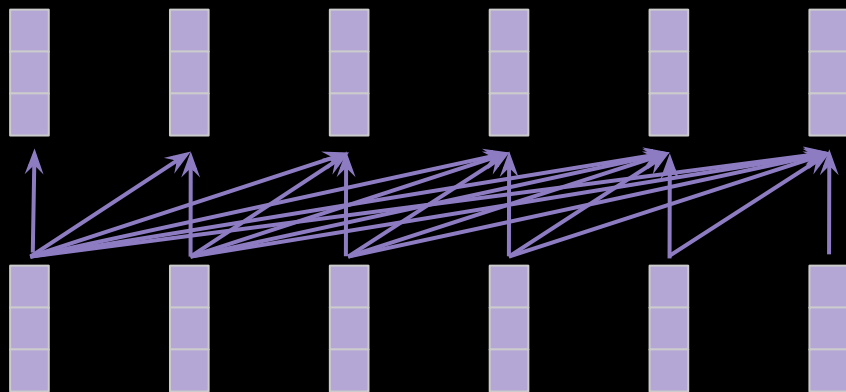
## **auto-encoder:**

- Connections go both directions.
- Task is predict word in middle:  
 $p(w_i | \dots, p_{w_i-2}, w_{i-1}, w_{i+1}, w_{i+2} \dots)$
- Better for:
  - embeddings
  - fine-tuning (transfer learning)



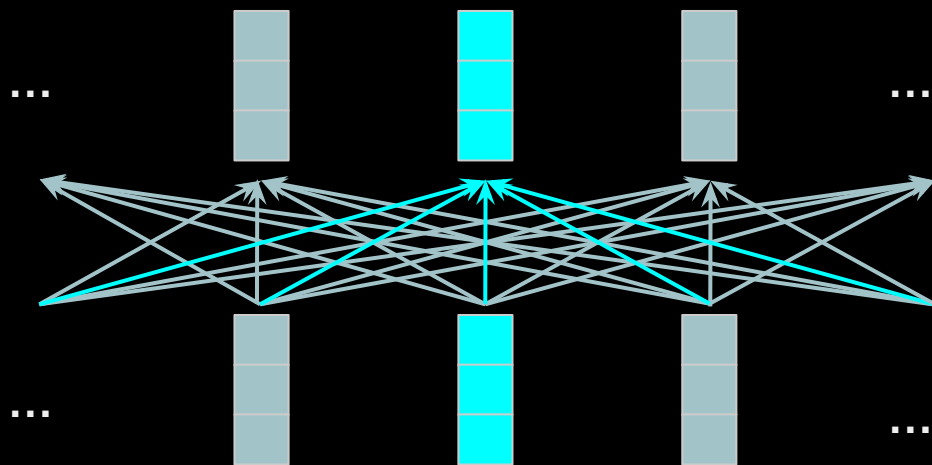
## **auto-regressor (generator):**

- Connections go forward only
- Task is predict word next word:  
 $p(w_i | w_{i-1}, w_{i-2}, \dots)$
- Better for:
  - generating text
  - zero-shot learning



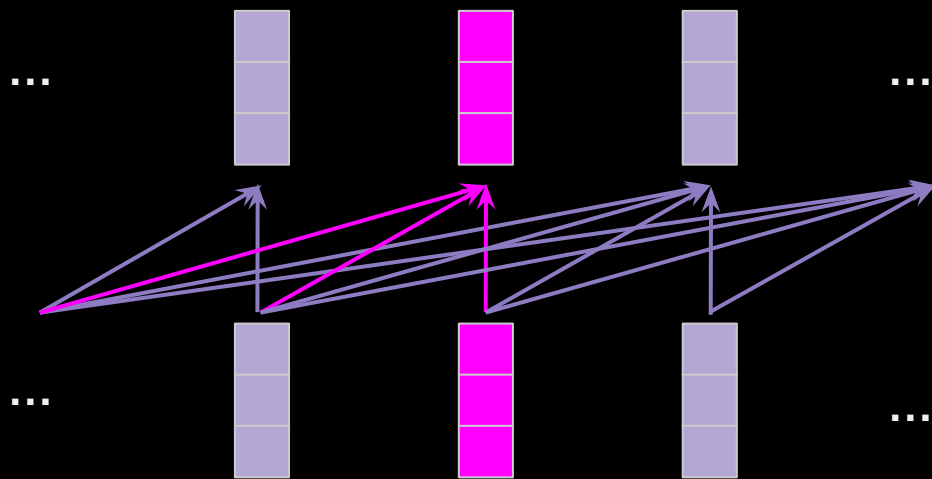
## **auto-encoder:**

- Connections go both directions.
- Task is predict word in middle:  
 $p(w_i \mid \dots, p_{w_{i-2}}, w_{i-1}, w_{i+1}, w_{i+2} \dots)$
- Better for:
  - embeddings
  - fine-tuning (transfer learning)

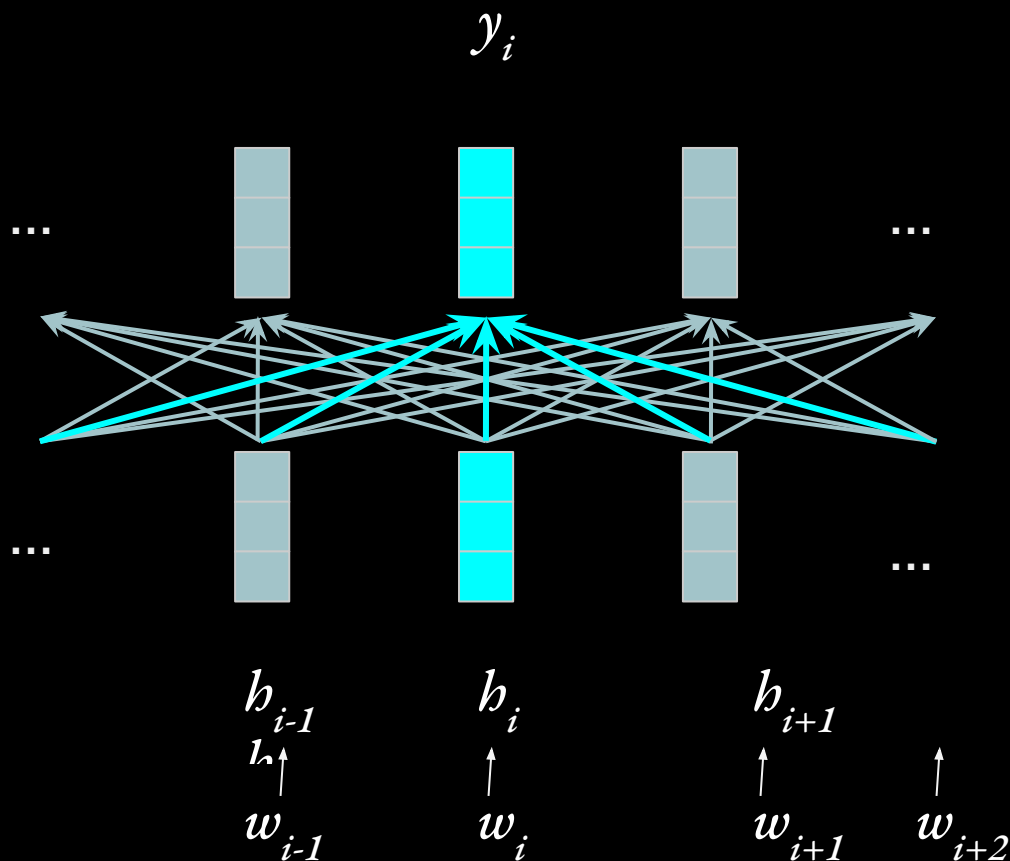


## **auto-regressor (generator):**

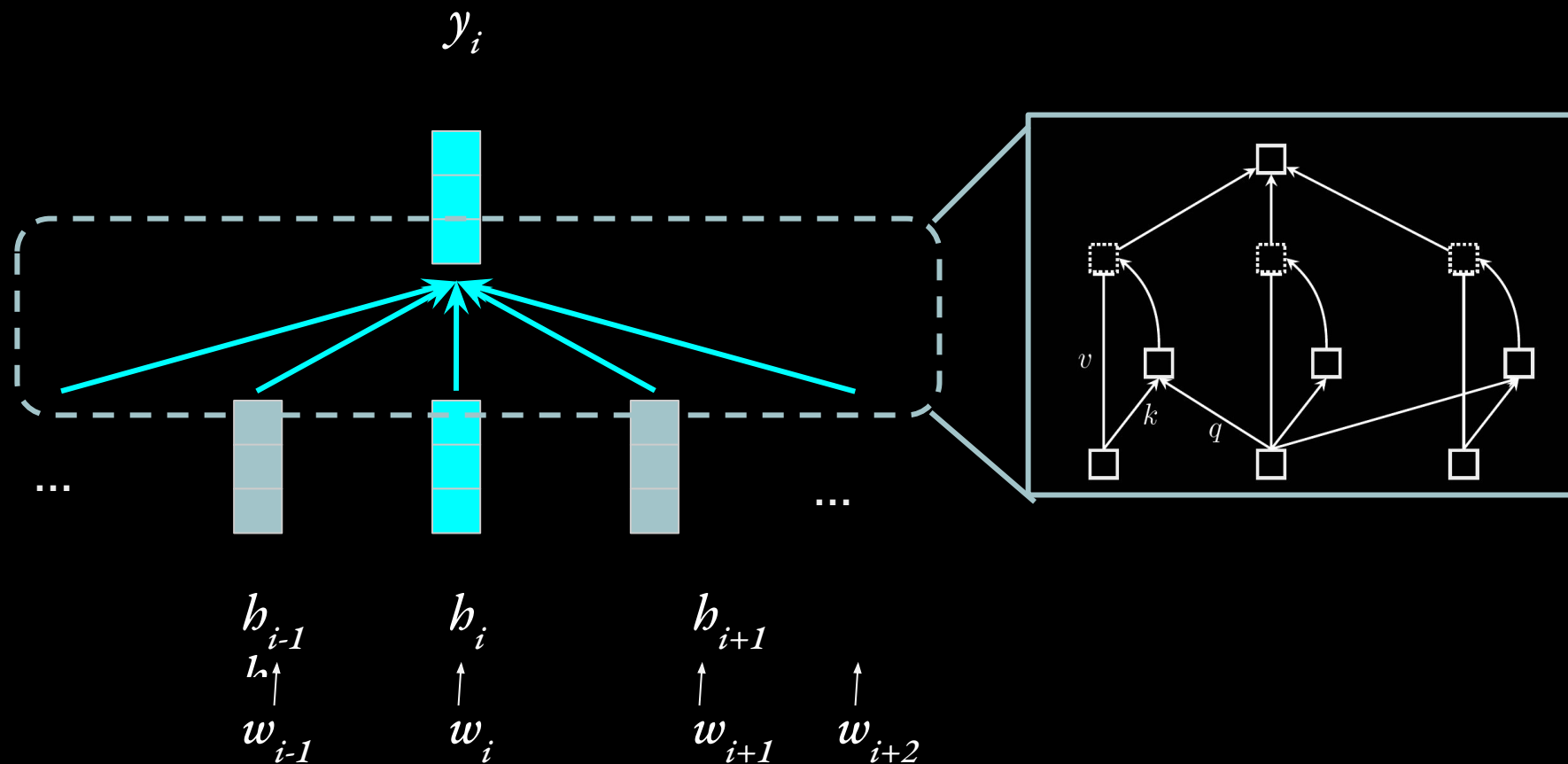
- Connections go forward only
- Task is predict word next word:  
 $p(w_i \mid w_{i-1}, w_{i-2}, \dots)$
- Better for:
  - generating text
  - zero-shot learning



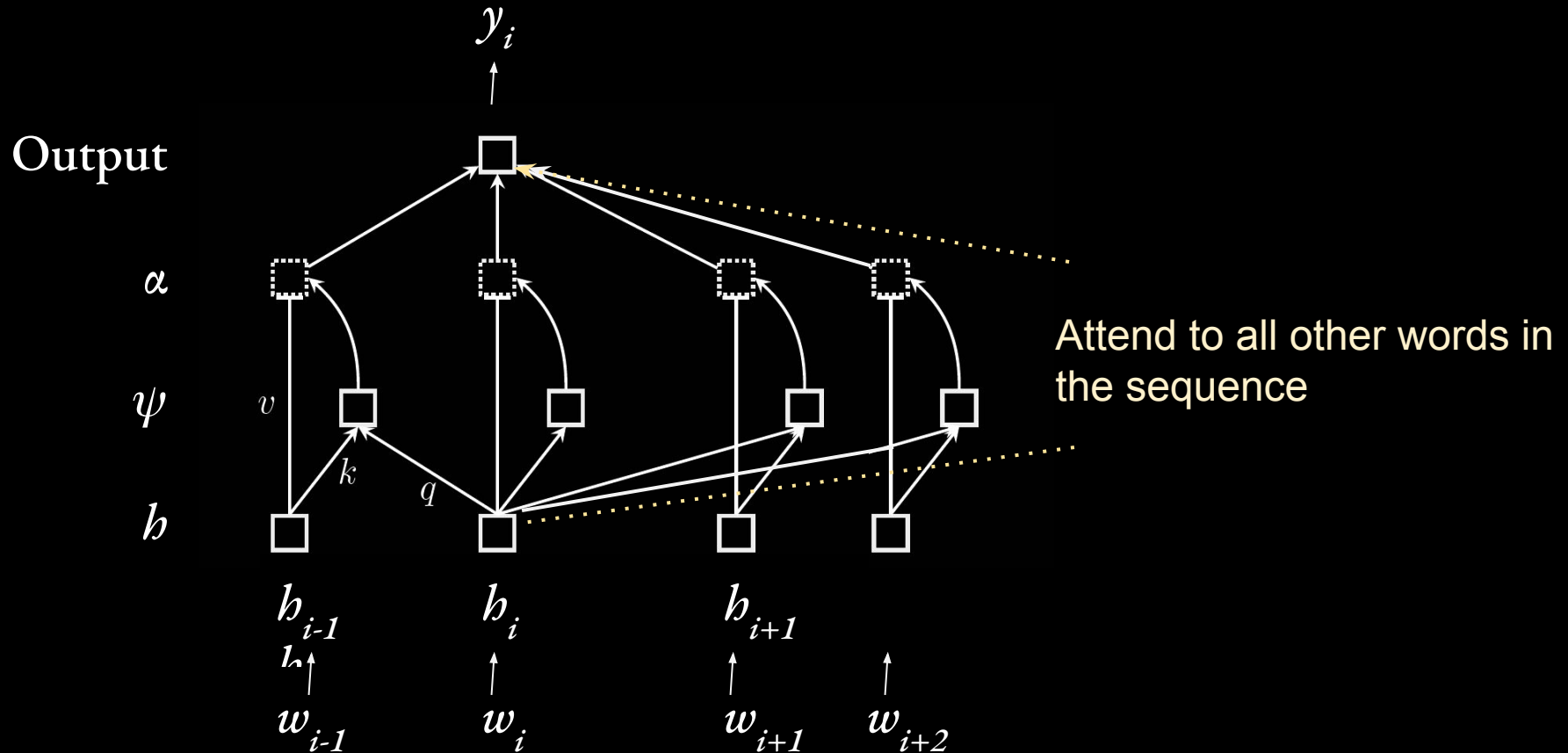
# The Transformer's Heart: Self-Attention



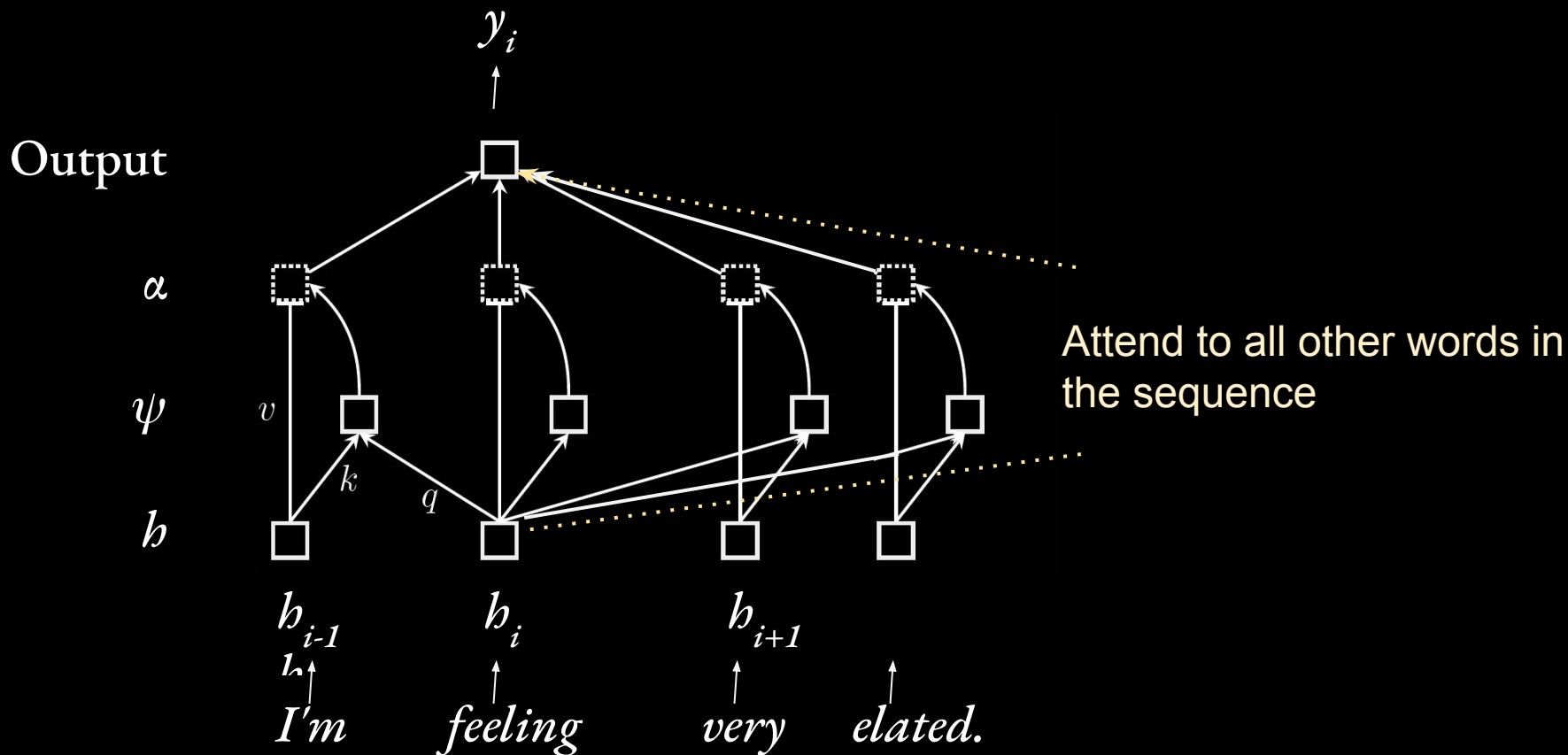
# The Transformer's Heart: Self-Attention



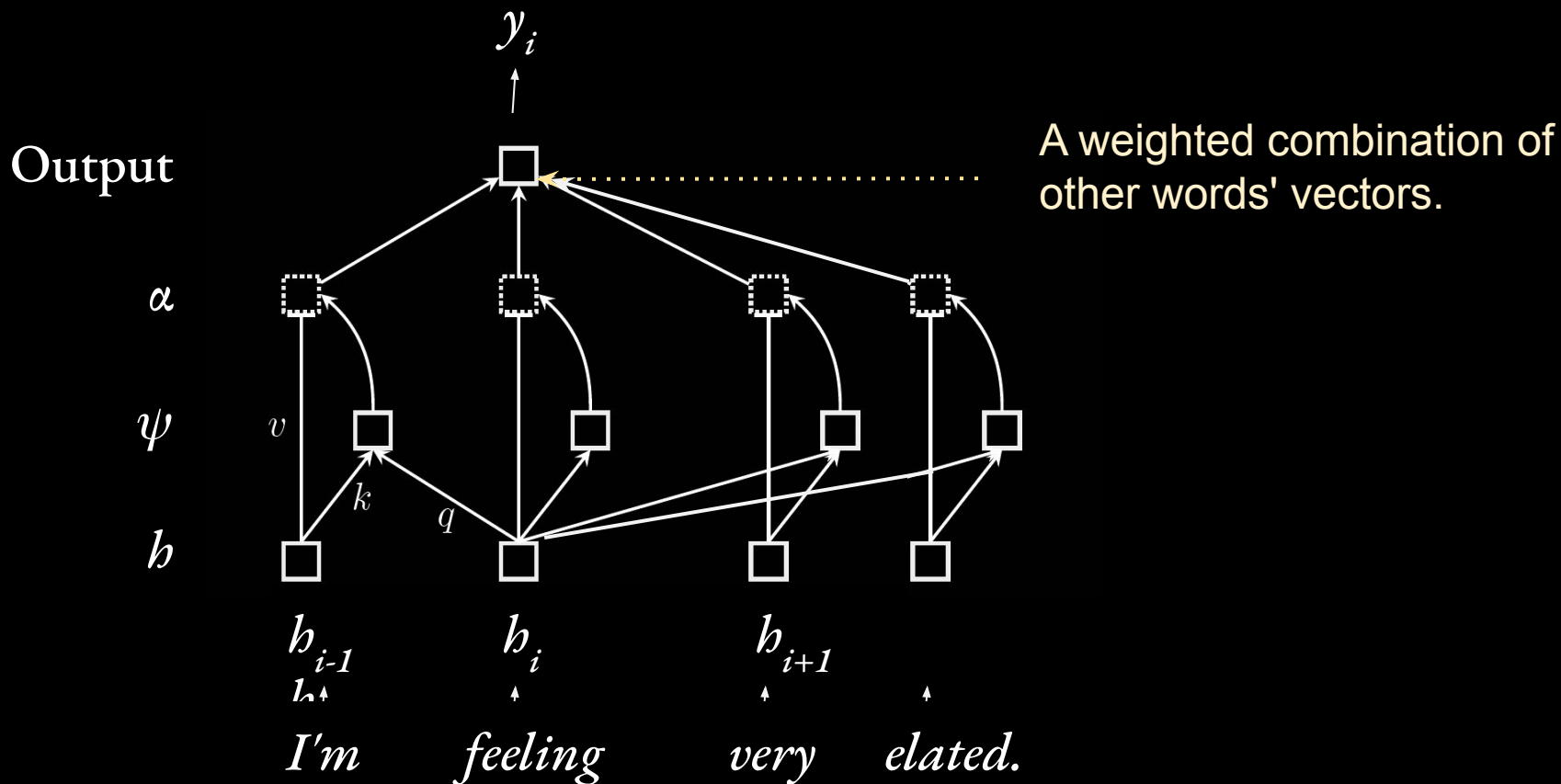
# The Transformer's Heart: Self-Attention



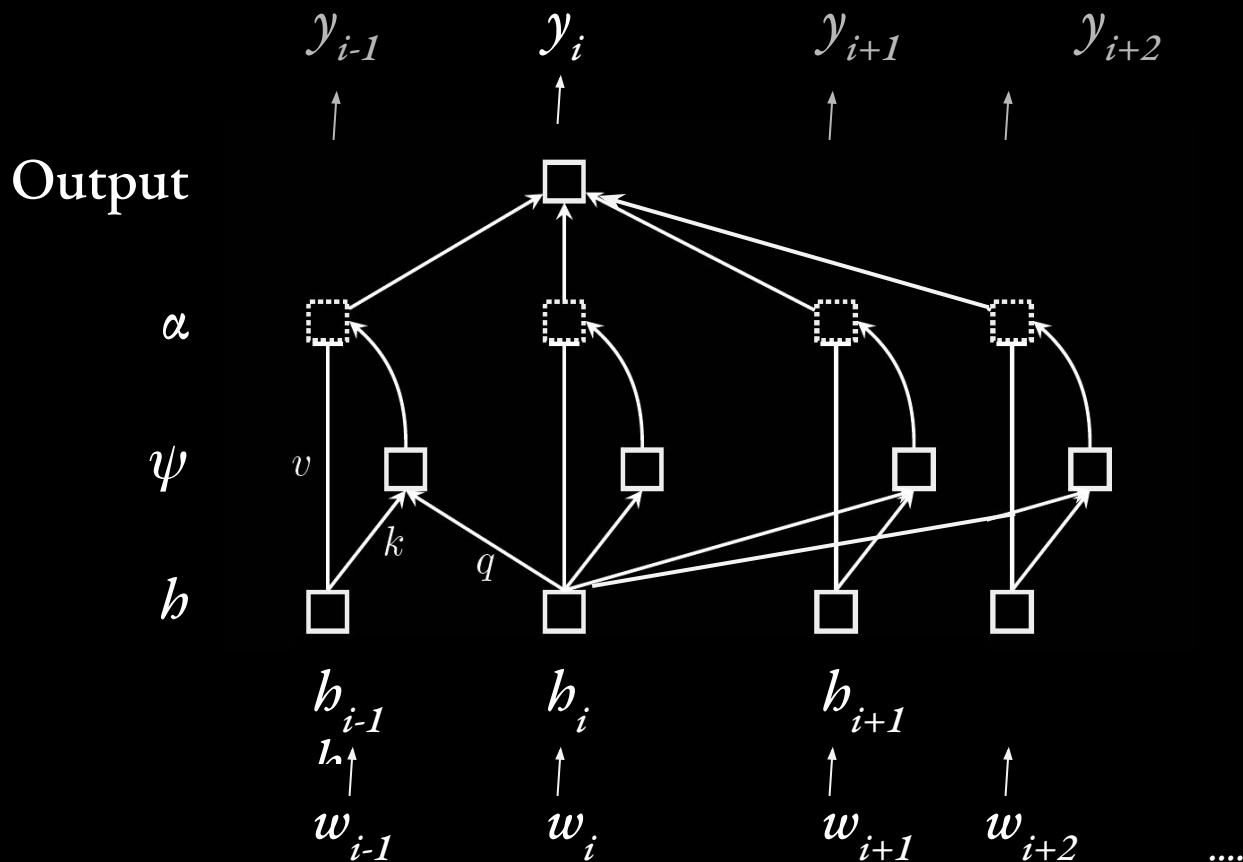
# The Transformer's Heart: Self-Attention



# The Transformer's Heart: Self-Attention

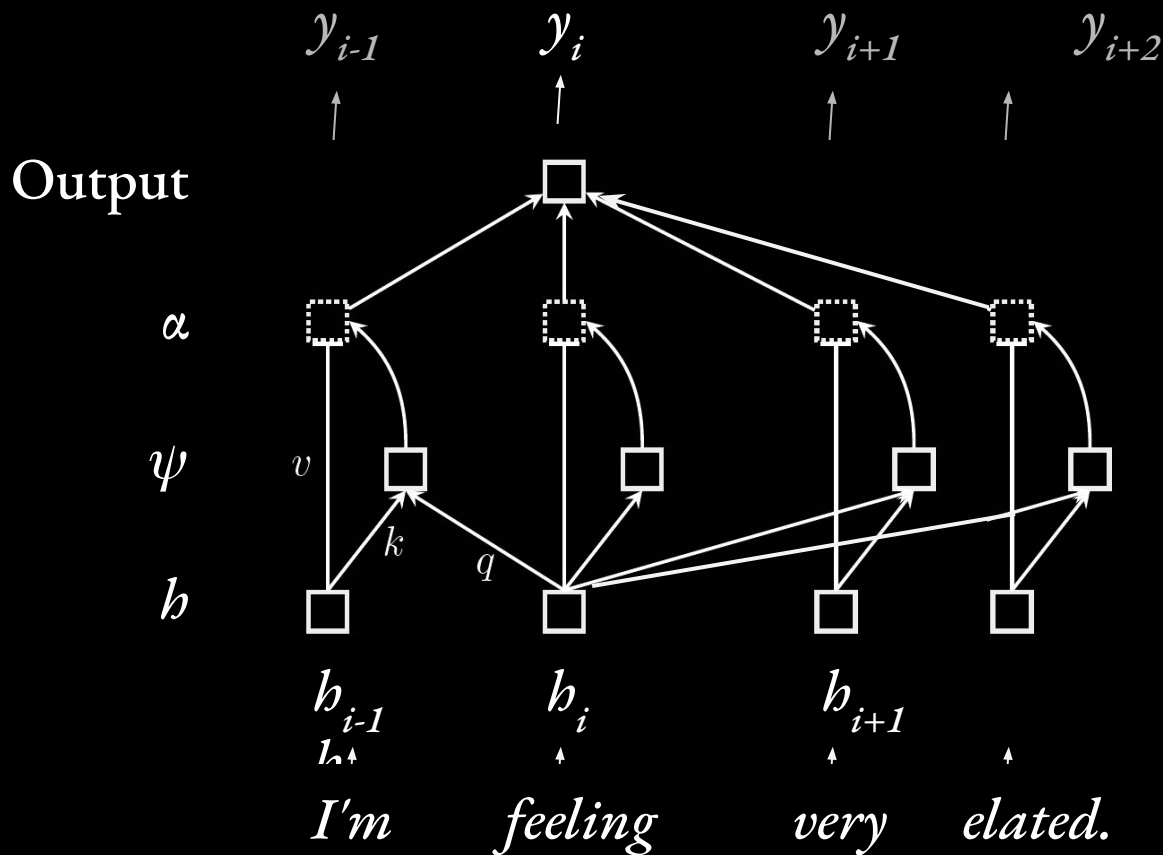


# The Transformer's Heart: Self-Attention

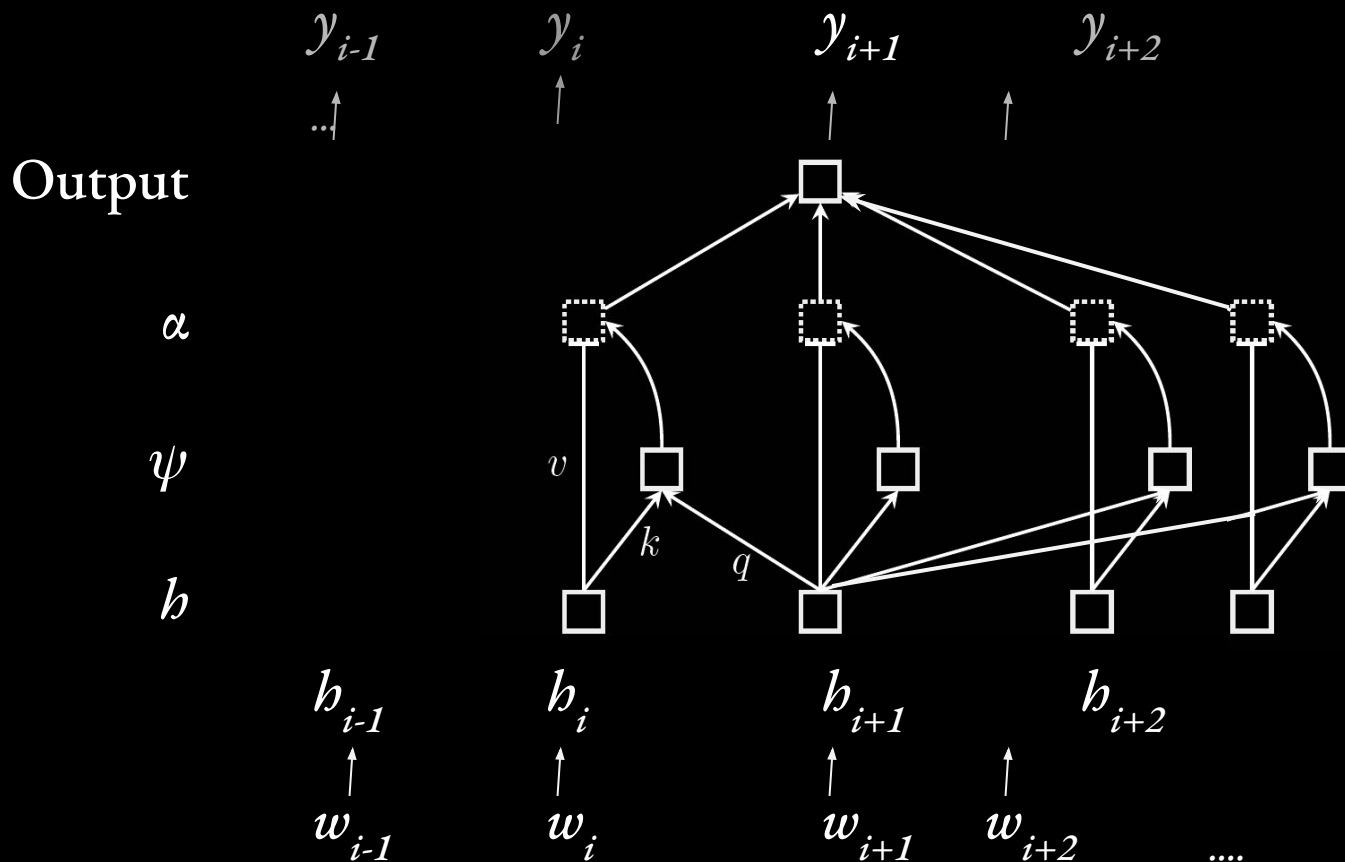




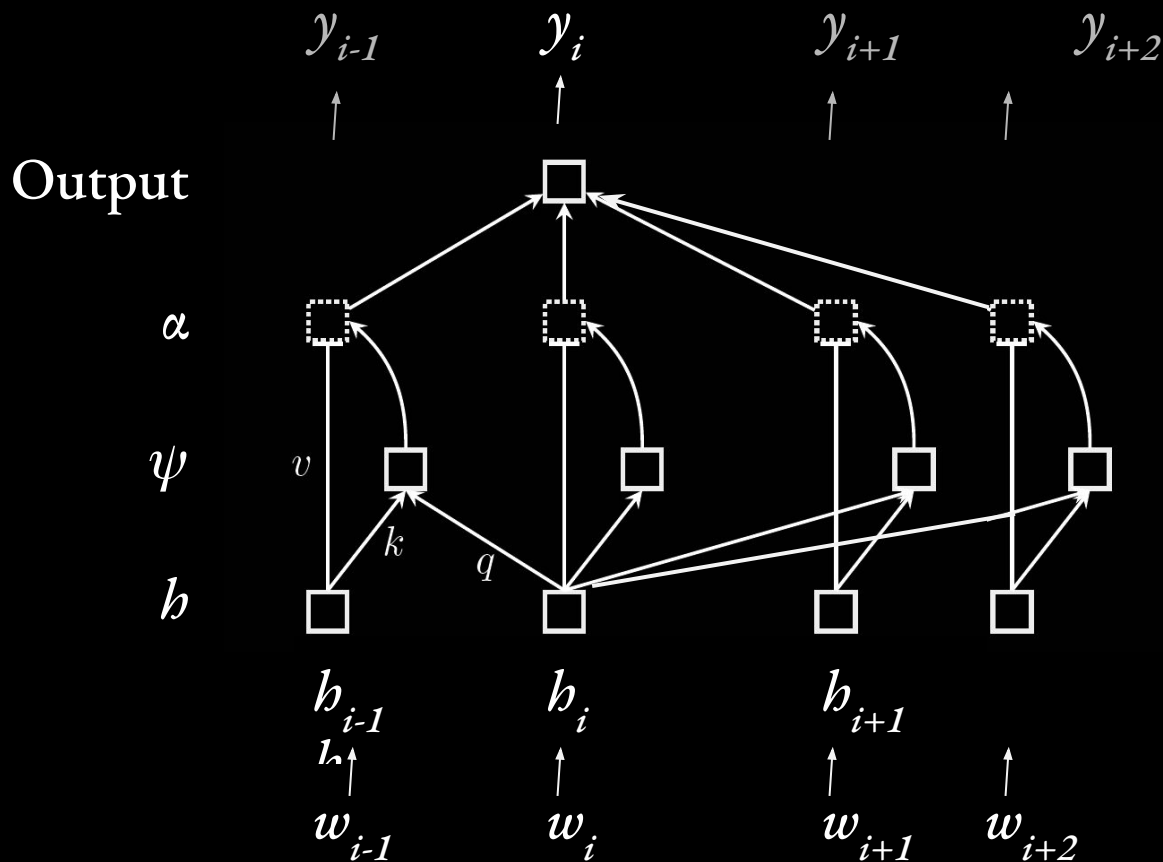
# The Transformer's Heart: Self-Attention



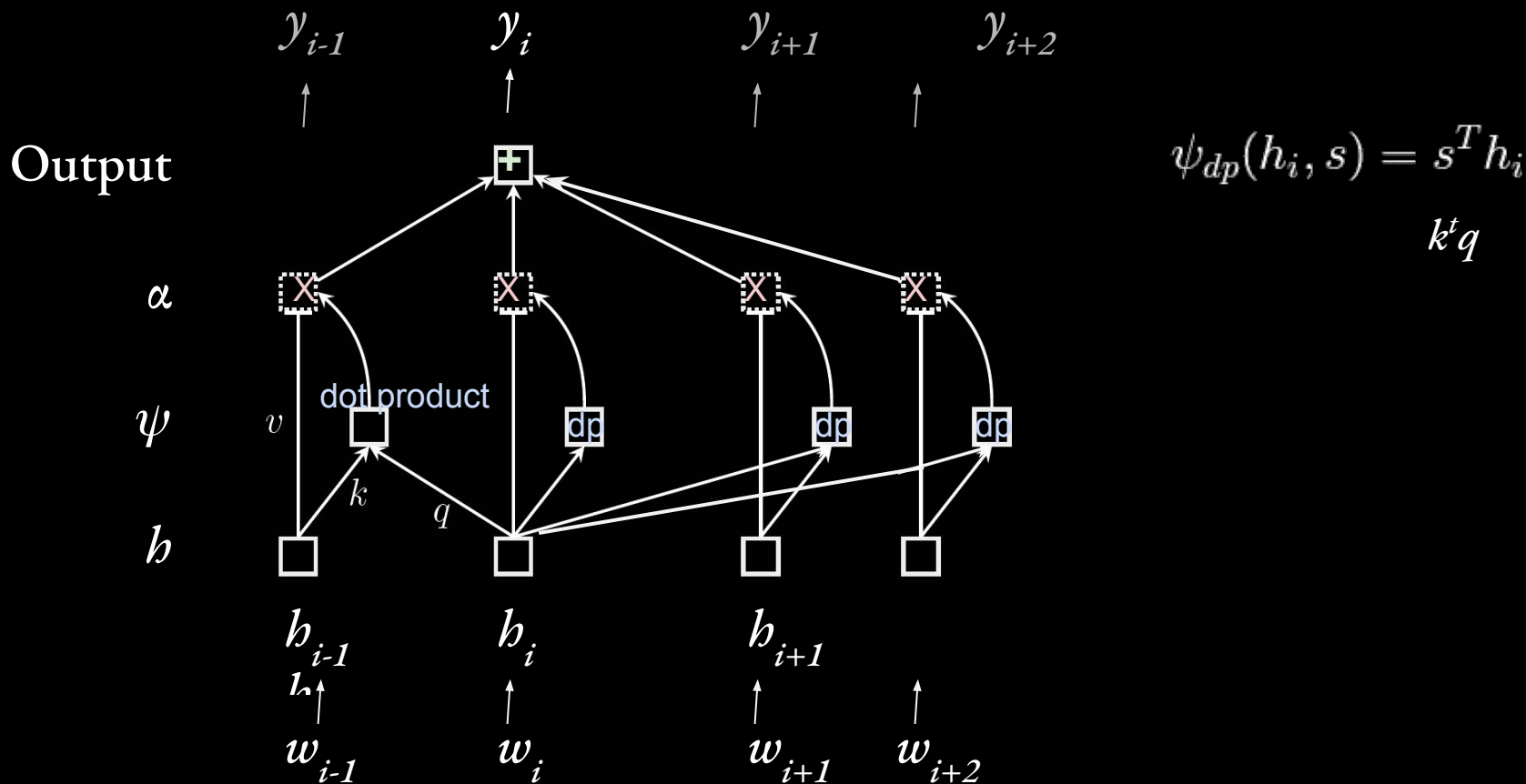
# The Transformer's Heart: Self-Attention



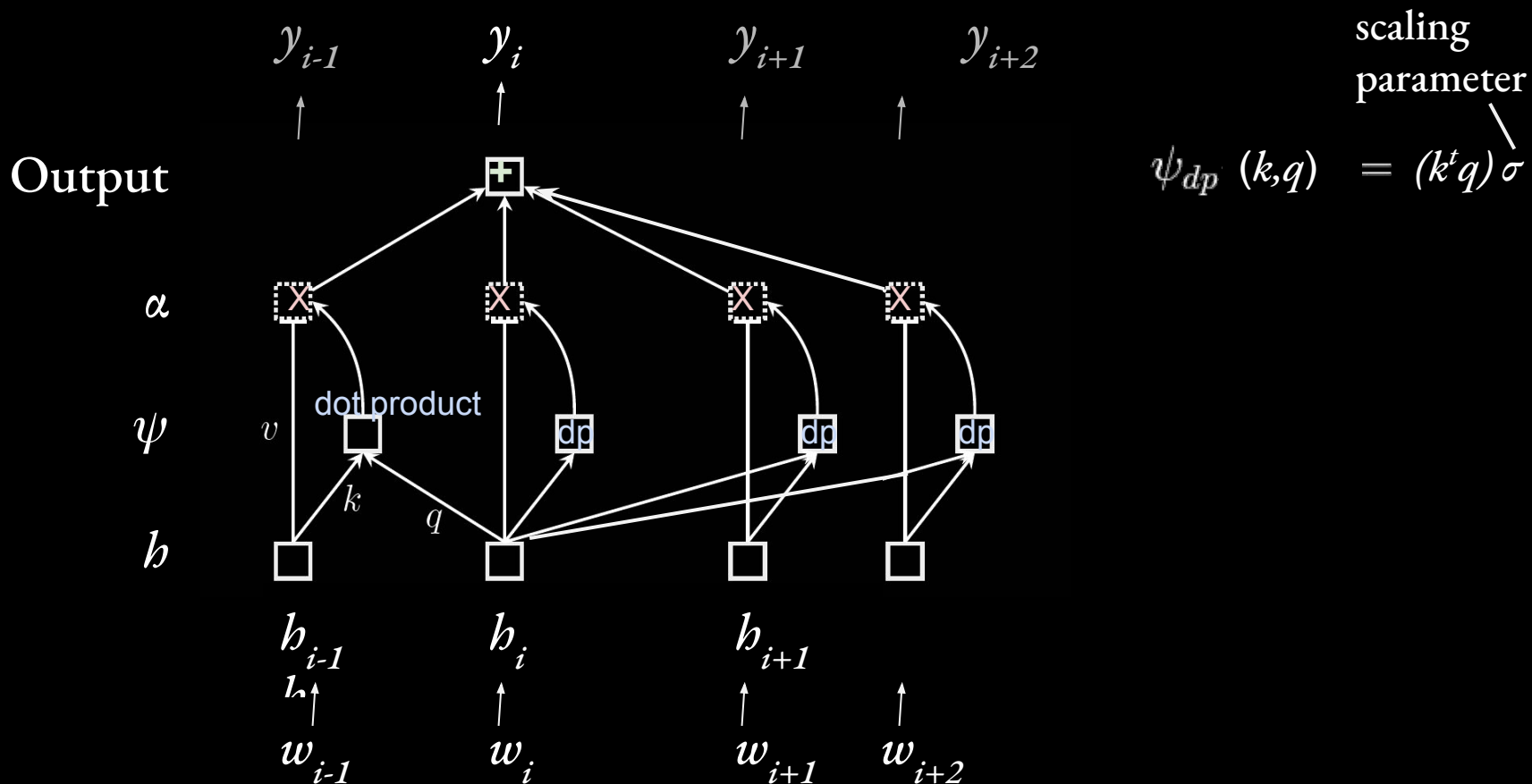
# The Transformer's Heart: Self-Attention



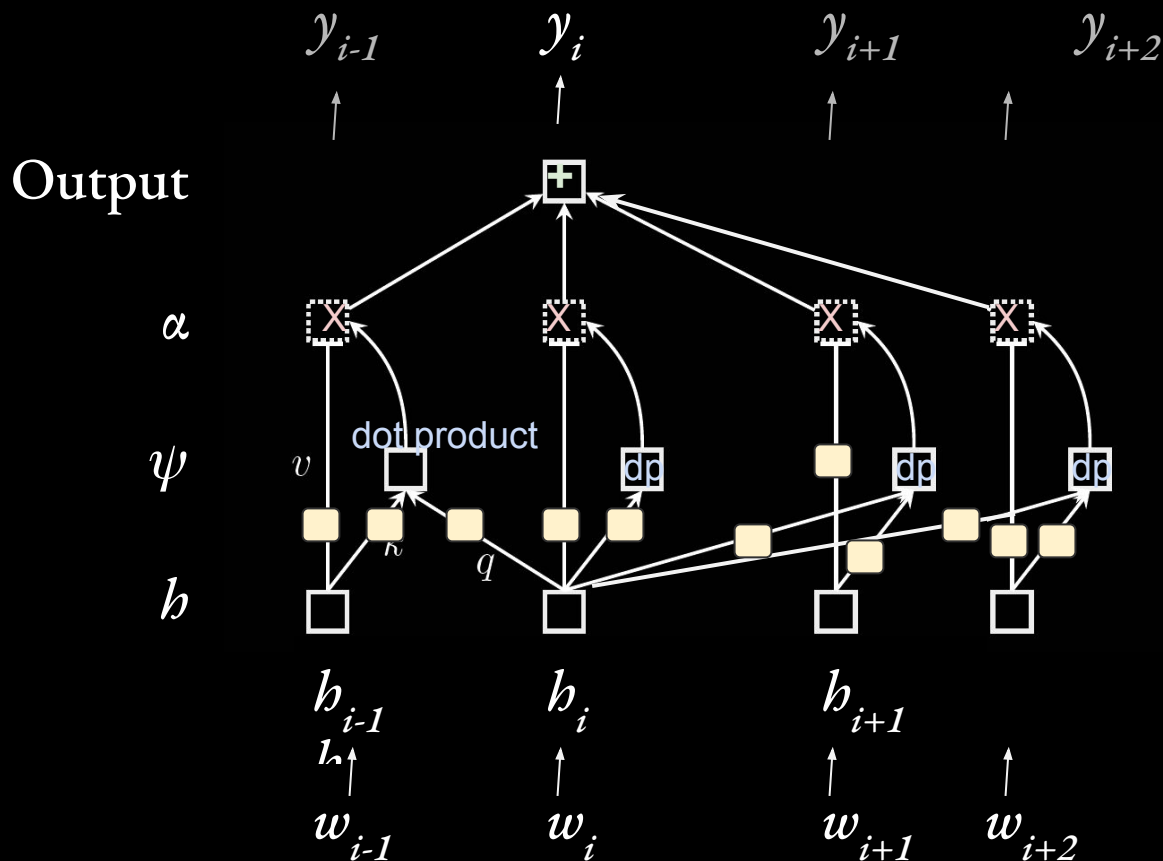
# The Transformer's Heart: Self-Attention



# The Transformer's Heart: Self-Attention



# The Transformer's Heart: Self-Attention



$$\psi_{dp}(k, q) = (k^t q) \sigma$$

Linear layer:  
 $W^T X$

One set of weights for each of for K, Q, and V

# Self-Attention in PyTorch

```
import nn.functional as f

class SelfAttention(nn.Module):
    def __init__(self, h_dim:int):
        self.Q = nn.Linear(h_dim, h_dim) #1 head
        self.K = nn.Linear(h_dim, h_dim)
        self.V = nn.Linear(h_dim, h_dim)

    def forward(hidden_states:torch.Tensor):
        v = self.V(hidden_states)
        k = self.K(hidden_states)
        q = self.Q(hidden_states)
        attn_scores = torch.matmul(q, k.T)
        attn_probs = f.Softmax(attn_scores)

        context = torch.matmul(attn_probs, v)
        return context
```

$$\psi_{dp}(k, q) = (k^t q) \sigma$$

Linear layer:

$$W^T X$$

One set of weights  
for each of for K,  
Q, and V

# Self-Attention in PyTorch

```
import nn.functional as f

class SelfAttention(nn.Module):
    def __init__(self, h_dim:int):
        self.Q = nn.Linear(h_dim, h_dim) #1 head
        self.K = nn.Linear(h_dim, h_dim)
        self.V = nn.Linear(h_dim, h_dim)
        self.dropout = nn.dropout(p=0.1)

    def forward(hidden_states:torch.Tensor):
        v = self.V(hidden_states)
        k = self.K(hidden_states)
        q = self.Q(hidden_states)
        attn_scores = torch.matmul(q, k.T)
        attn_probs = f.Softmax(attn_scores)
        attn_probs = self.dropout(attn_probs)
        context = torch.matmul(attn_probs, v)
        return context
```

$$\psi_{dp}(k, q) = (k^t q) \sigma$$

Linear layer:

$$W^T X$$

One set of weights  
for each of for K,  
Q, and V



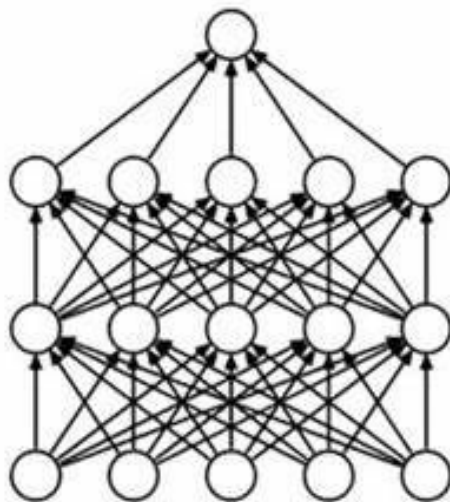
# Self-Attention in PyTorch

```
import nn.functional as f
class SelfAttention(nn.Module):
```

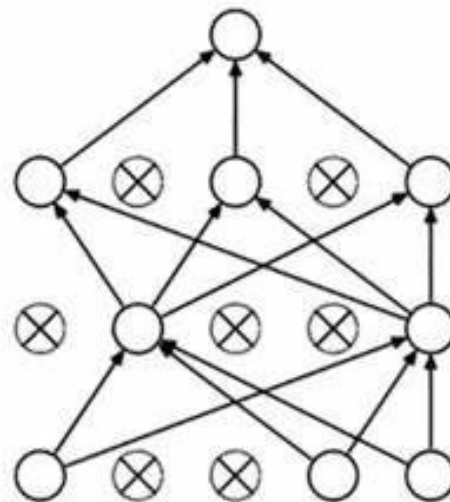
```
    def __init__(
        self.Q =
        self.K =
        self.V =
        self.drop
```

```
    def forward(h
        v = self.
        k = self.
        q = self.
        attn_scor
        attn_prob
```

```
        attn_probs = self.dropout(attn_probs)
        context = torch.matmul(attn_probs, v)
        return context
```



(a) Standard Neural Net



(b) After applying dropout.

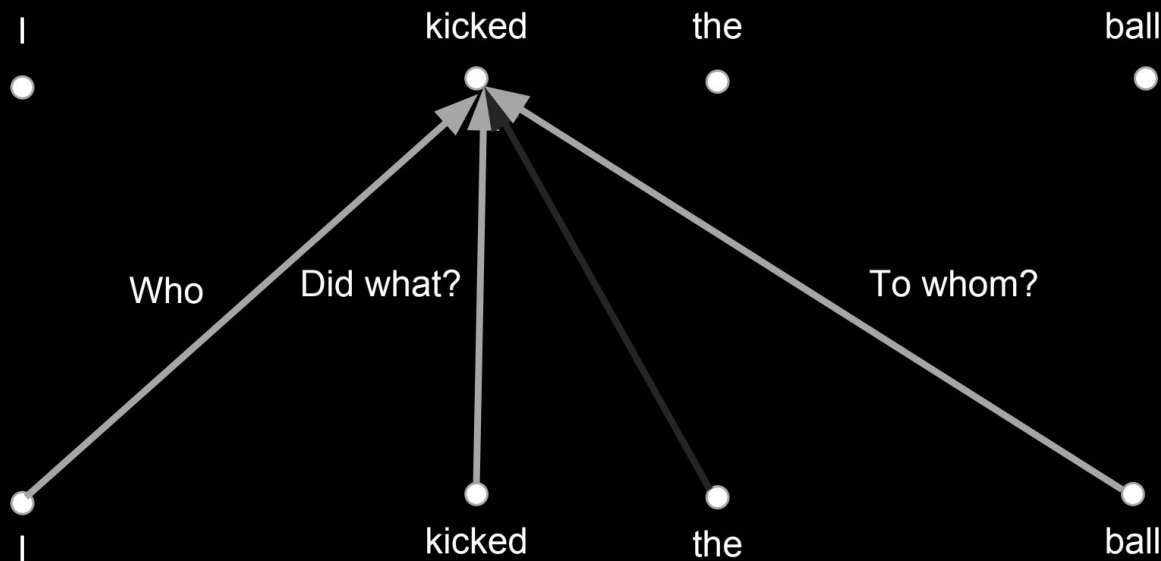
$$= (k^t q) \sigma$$

ayer:

of weights  
of for K,

# The Transformer: Beyond Self-Attention

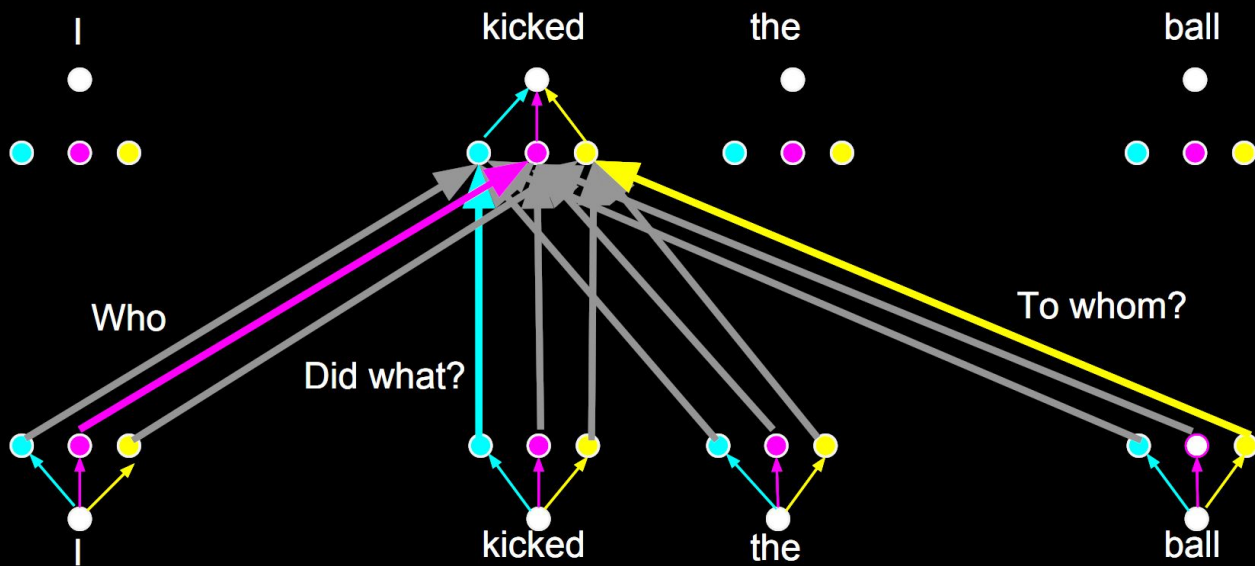
Limitation (thus far): Can't capture multiple types of dependencies between words.



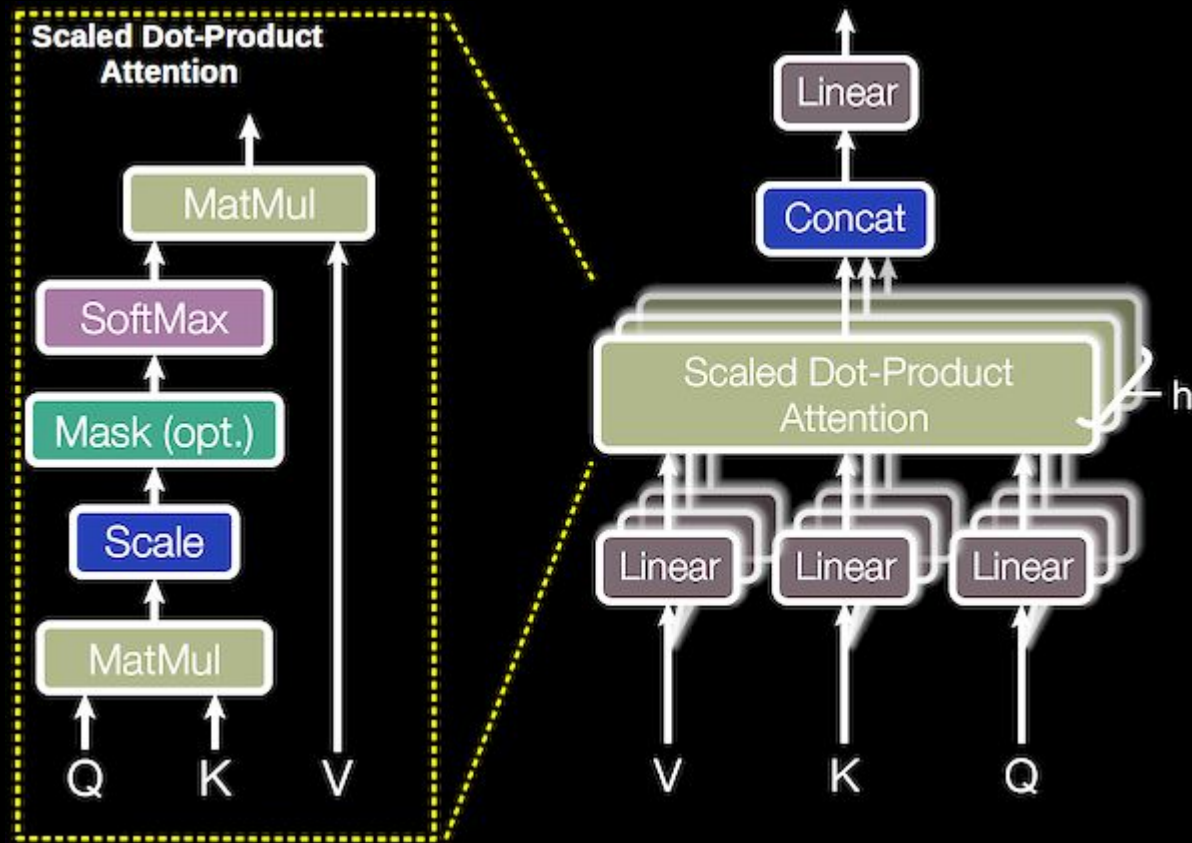
# The Transformer: Beyond Self-Attention

Limitation (thus far): Can't capture multiple types of dependencies between words.

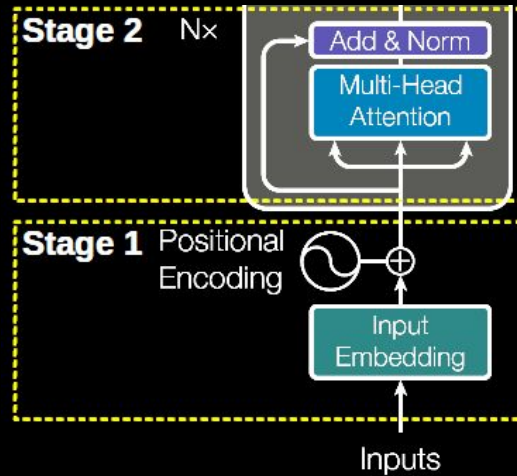
Solution: Multi-head attention



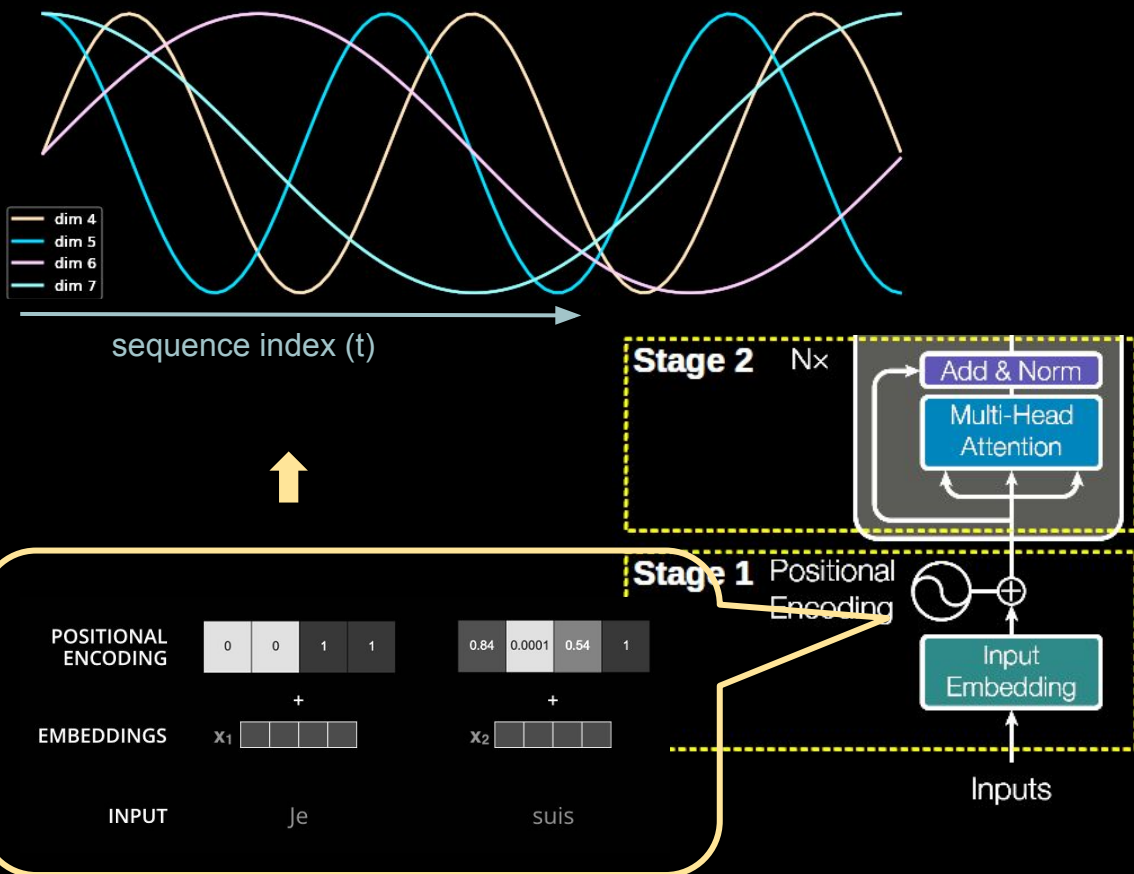
# The Transformer: Multi-headed Attention



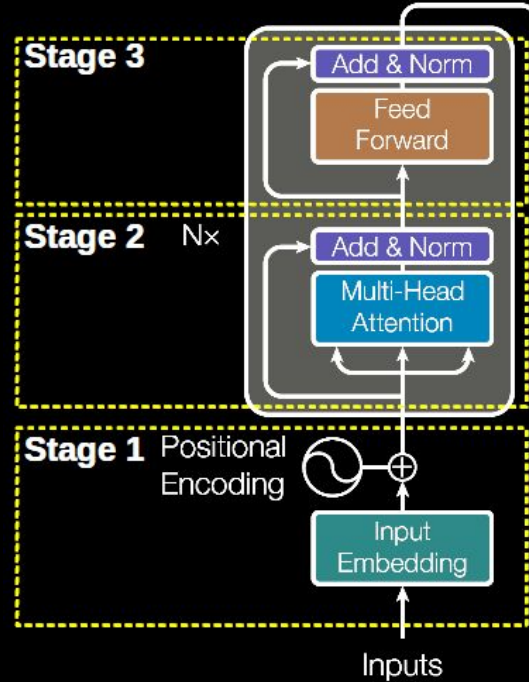
# The Transformer



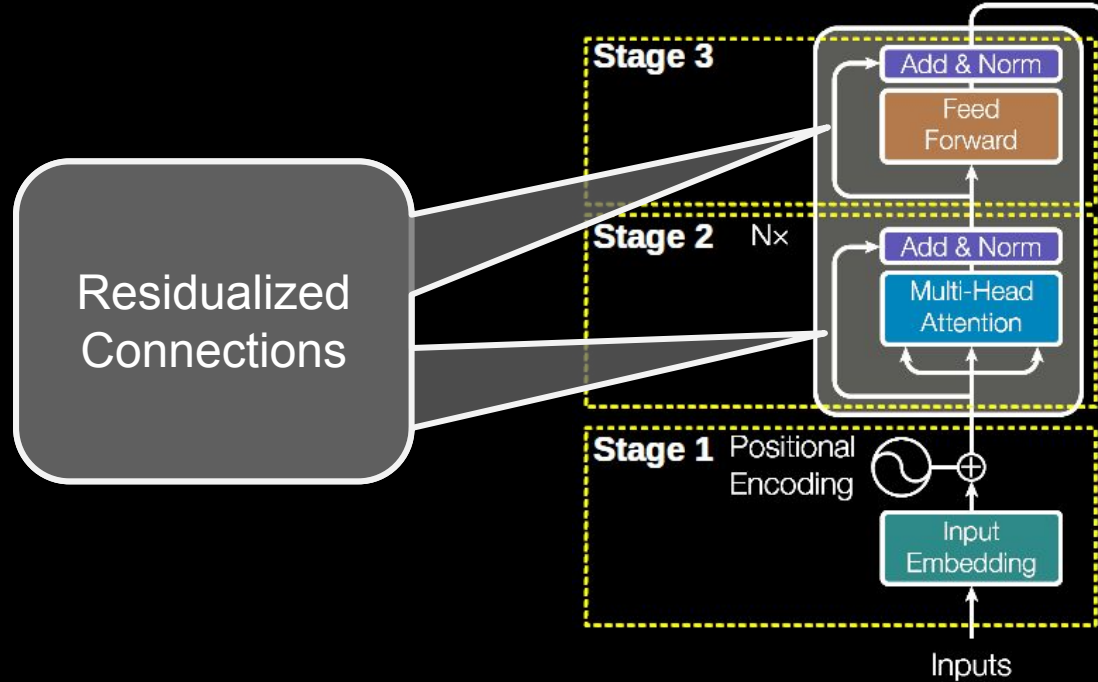
# The Transformer



# The Transformer

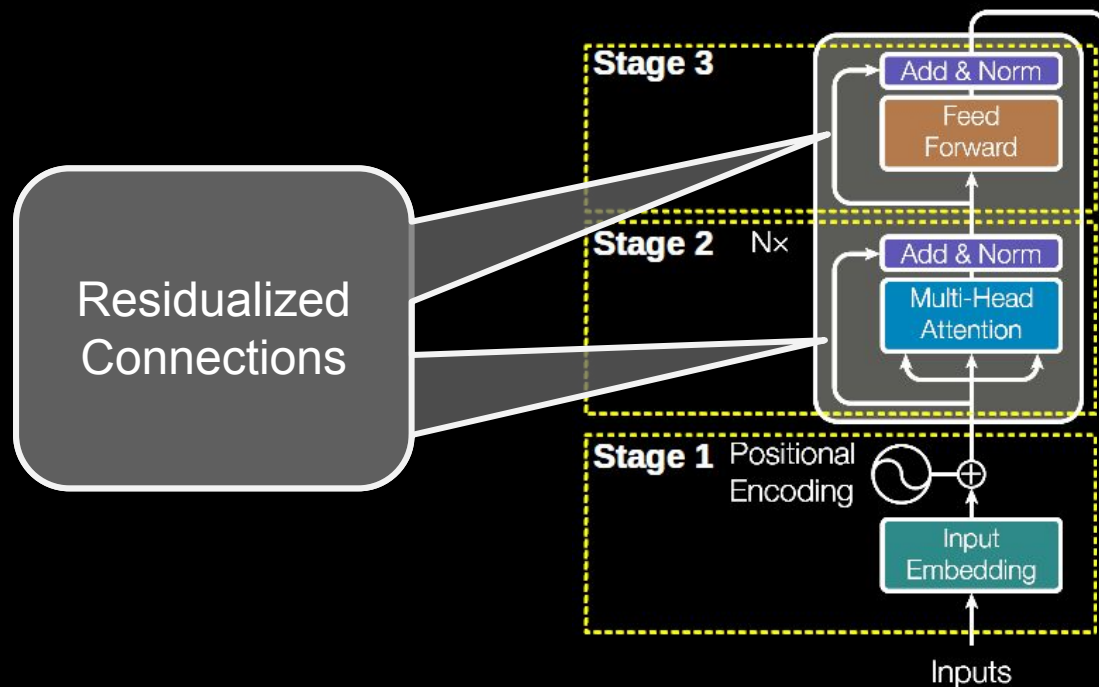


# The Transformer

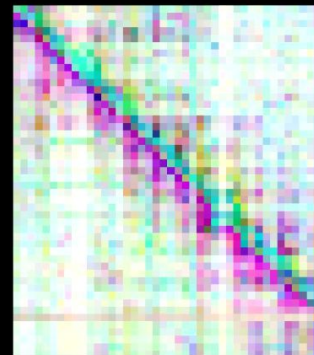




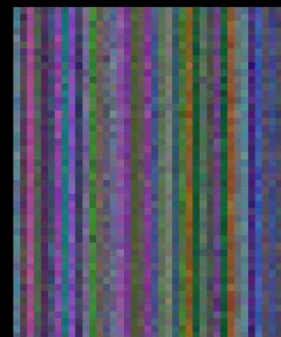
# The Transformer



residuals enable positional information to be passed along



With residuals



Without residuals

# The Transformer: Motivation

## Challenges to sequential representation learning

- Capture long-distance dependencies
- Preserving sequential distances / periodicity
- Capture multiple relationships
- Easy to parallelize -- don't need sequential processing.

# The Transformer: Motivation

## Challenges to sequential representation learning

- Capture long-distance dependencies  
***Self-attention** treats far away words similar to those close.*
- Preserving sequential distances / periodicity  
***Positional embeddings** encode distances/periods.*
- Capture multiple relationships  
***Multi-headed attention** enables multiple compositions.*
- Easy to parallelize -- don't need sequential processing.  
*Entire layer can be computed at once. Is only matrix multiplications + standardizing.*

# Transformer (as of 2017)

“WMT-2014” Data Set. BLEU scores:

	EN-DE	EN-FR
GNMT (orig)	24.6	39.9
ConvSeq2Seq	25.2	40.5
Transformer*	<b>28.4</b>	<b>41.8</b>

# Transformers as of 2023

General Language Understanding Evaluations:

<https://gluebenchmark.com/leaderboard>

<https://super.gluebenchmark.com/leaderboard/>

## ChatGPT



ChatGPT is an artificial intelligence chatbot developed by OpenAI and launched in November 2022. It is built on top of OpenAI's GPT-3.5 and GPT-4 families of large language models and has been fine-tu...



# Transformers as of 2023

Machine  
Translation

Web  
Search

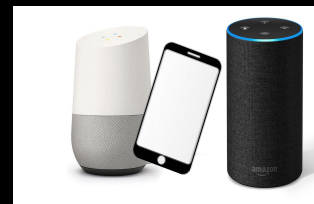
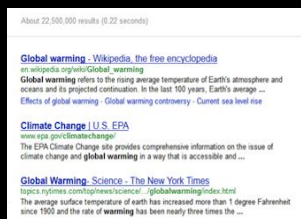
Sentiment  
Analysis

Document  
Classification

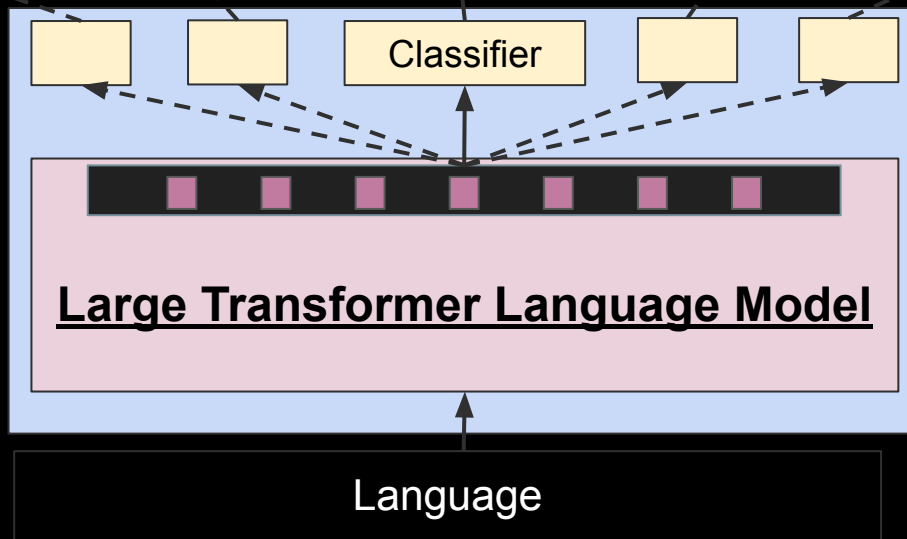
Assistant,  
QA

...

absolutamente  
me gustaría ir  
de excursión



Soni, N., Matero, M.,  
Balasubramanian, N., &  
Schwartz, H. (2022, May).  
Human Language Modeling. In  
*Findings of the Association for  
Computational Linguistics: ACL  
2022* (pp. 622-636).

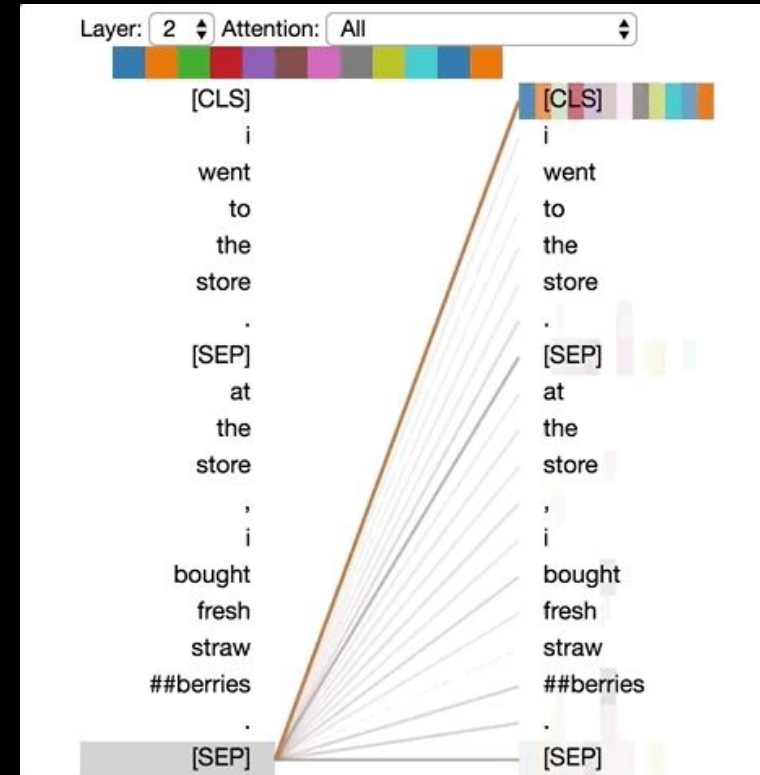


(NLP System)



# Bert: Attention by Layers

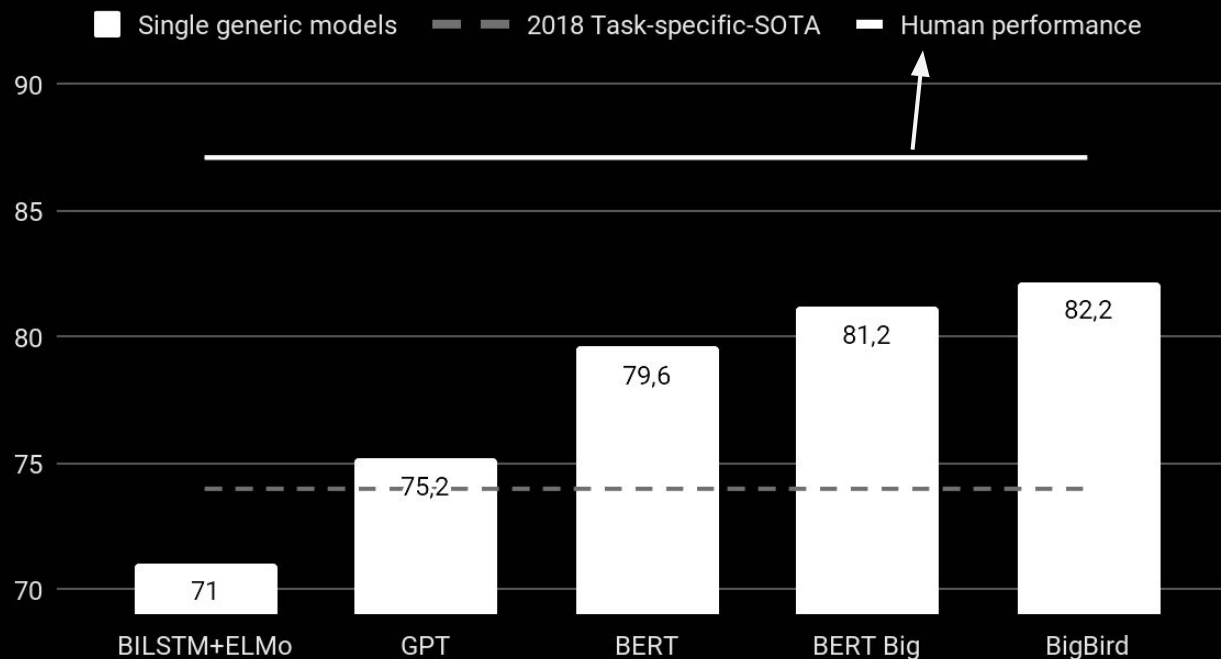
<https://colab.research.google.com/drive/1vIOJ1lhdujVjfH857hvYKIdKPTD9Kid8>



(Vig, 2019)

# BERT Performance: e.g. Question Answering

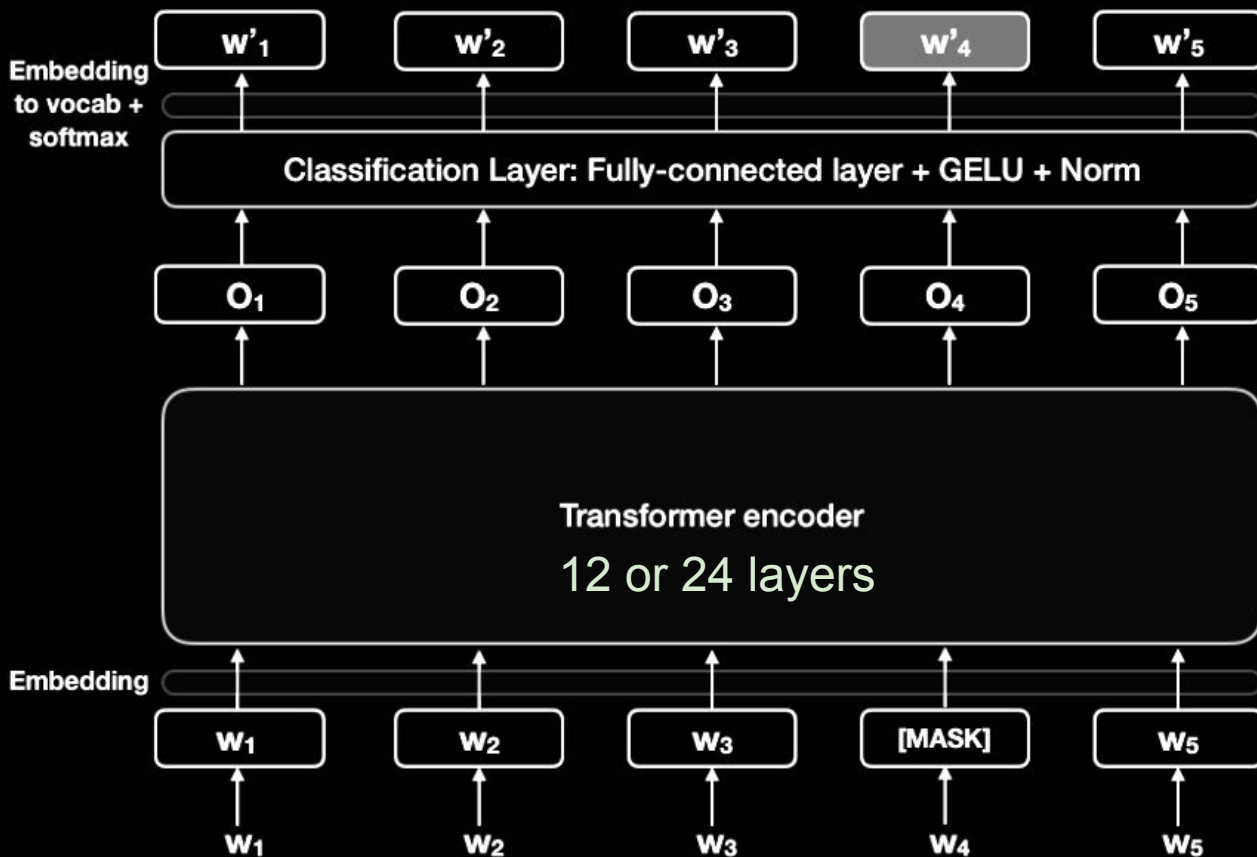
GLUE scores evolution over 2018-2019



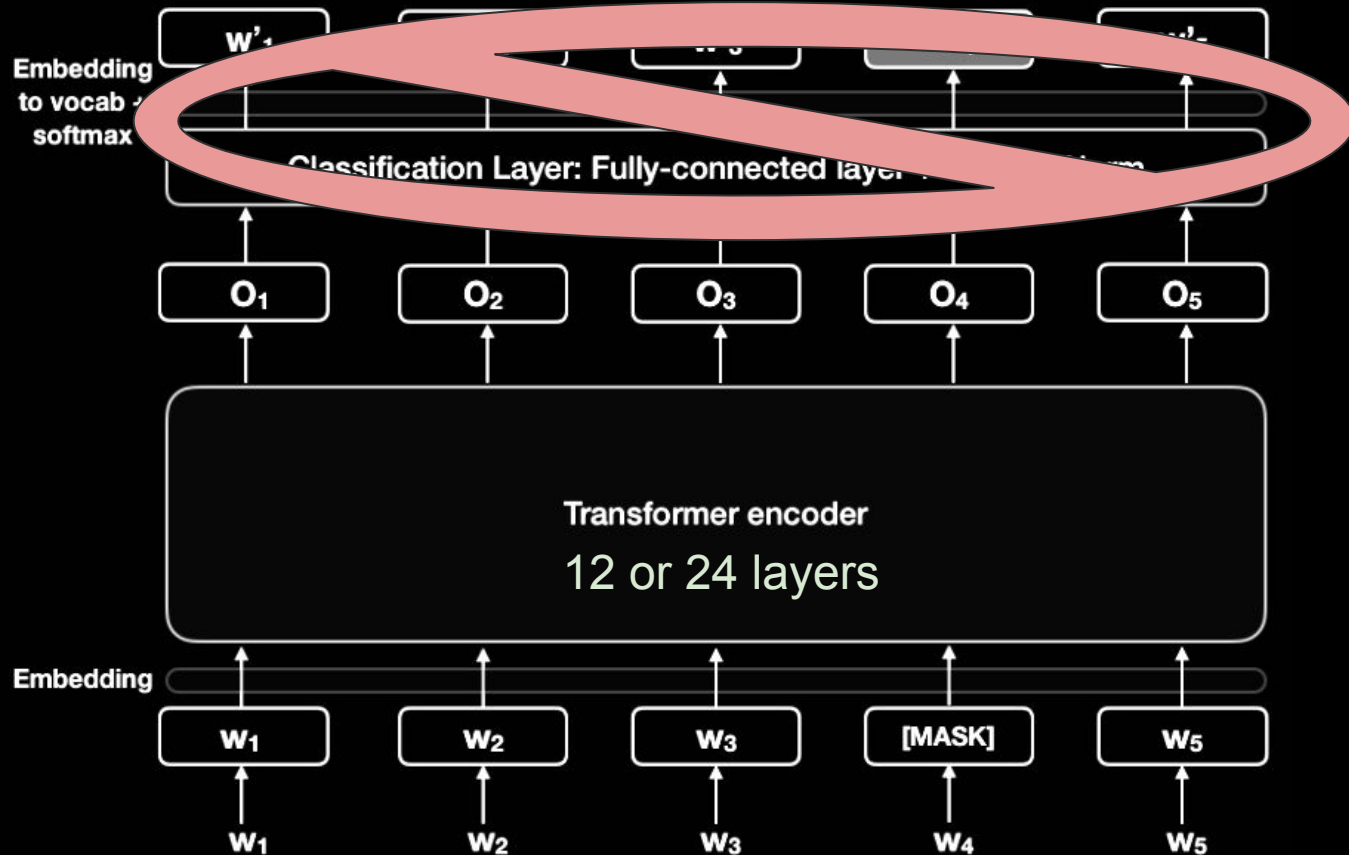
<https://rajpurkar.github.io/SQuAD-explorer/>



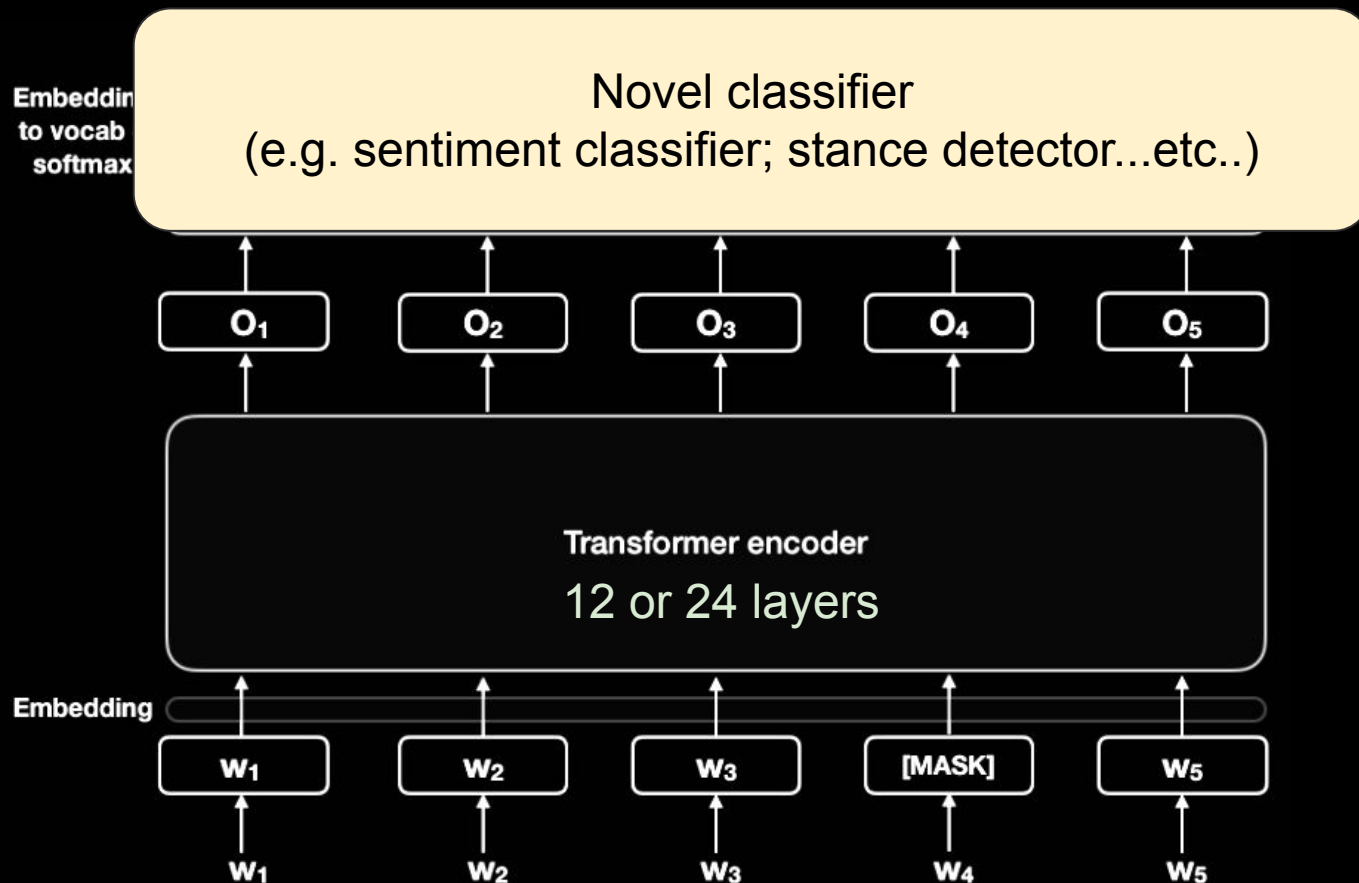
# BERT: Pre-training; Fine-tuning



# BERT: Pre-training; Fine-tuning



# BERT: Pre-training; Fine-tuning

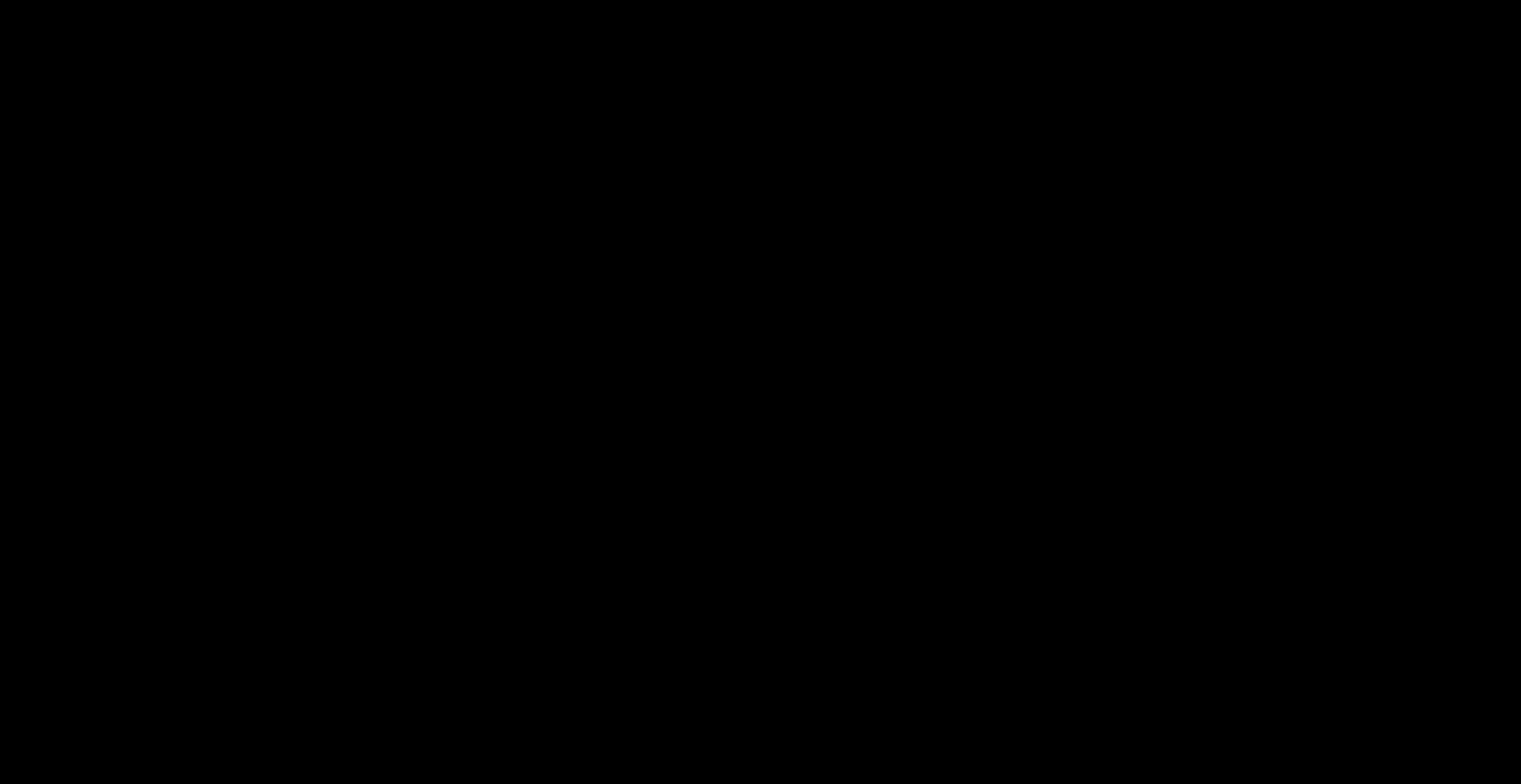


# The Transformer: Motivation

## Challenges to sequential representation learning

- Capture long-distance dependencies  
*Self-attention treats far away words similar to those close.*
- Preserving sequential distances / periodicity  
*Positional embeddings encode distances/periods.*
- Capture multiple relationships  
*Multi-headed attention enables multiple compositions.*
- Easy to parallelize -- don't need sequential processing.  
*Entire layer can be computed at once. Is only matrix multiplications + standardizing.*

# Extra Material



# When L2, L1 work well?

*Depends on data, but generally we find...*

*k: number of features; N: number of observations*

	$k \ll N$	$k < N$	$k == N$	$k > N$	$k \gg N$
<i>None</i>	Often	Sometimes	Almost Never	Almost Never	Almost Never
<i>L2</i>	Sometimes	Often	Often	Sometimes	Almost Never
<i>L1</i>	Almost Never	Sometimes	Often	Sometimes	Almost Never