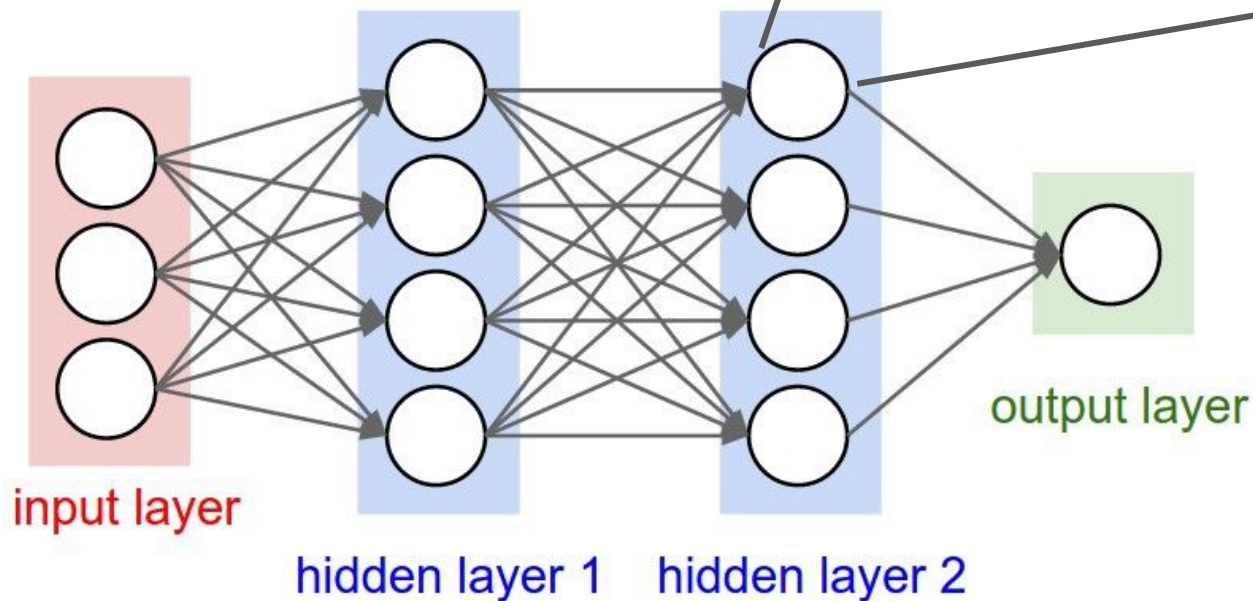
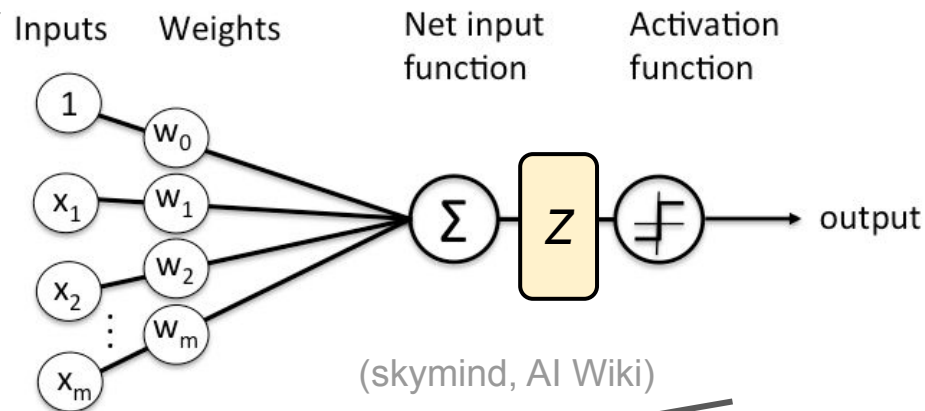


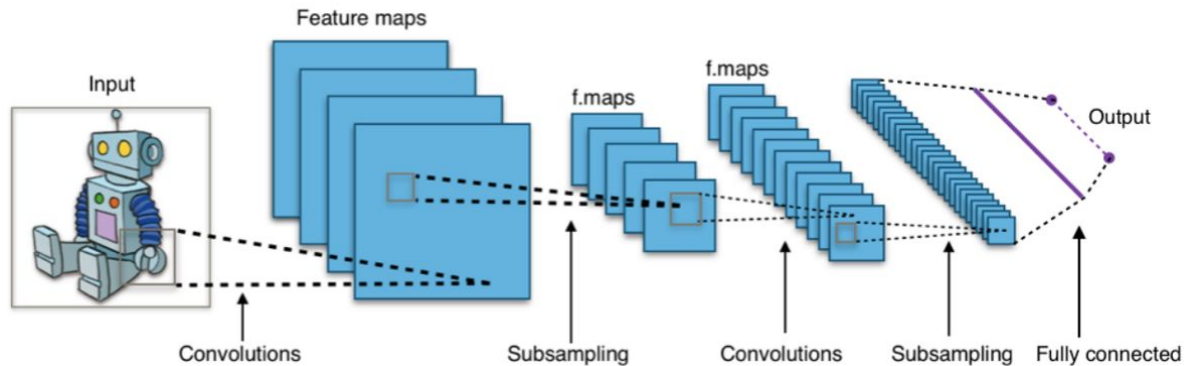
# Transformer Models

CSE545 - Spring 2019

# Review: Feed Forward Network (full-connected)

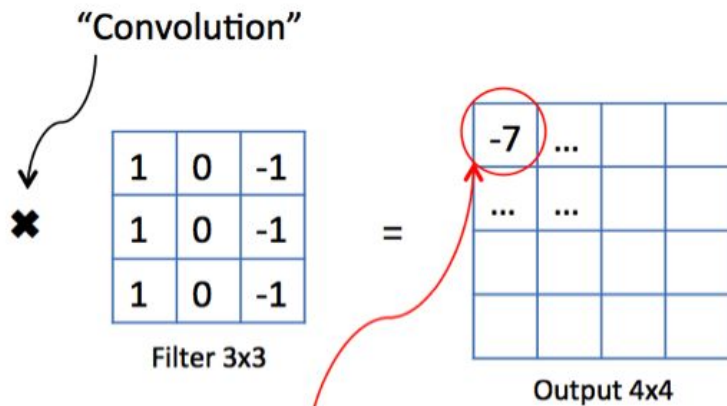


# Review: Convolutional NN



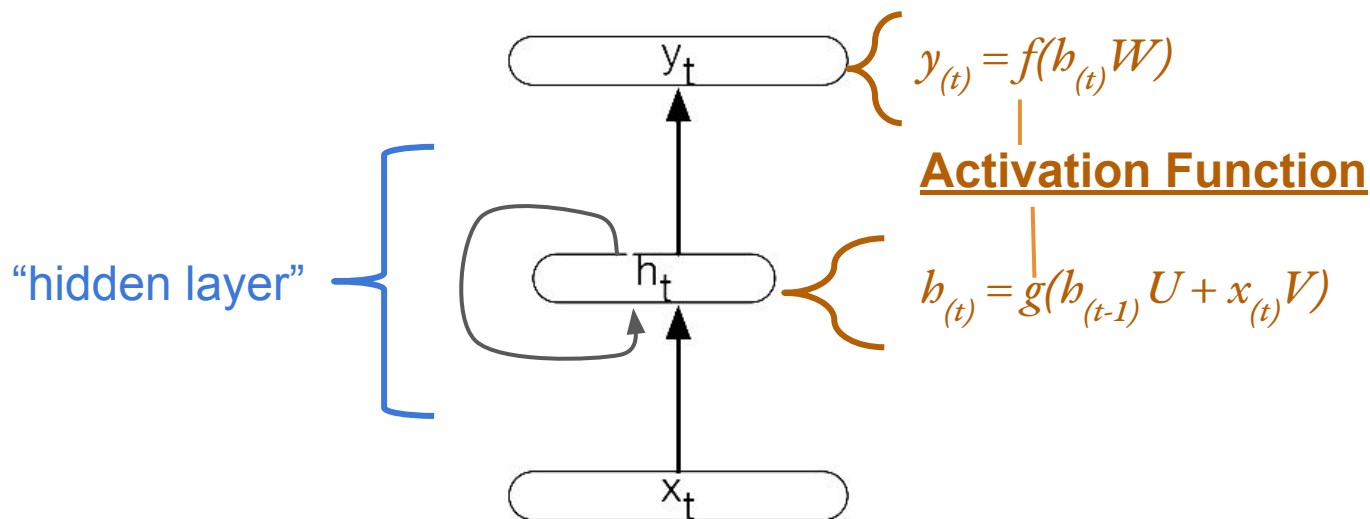
3	1	1	2	8	4
1	0	7	3	2	6
2	3	5	1	1	3
1	4	1	2	6	5
3	2	1	3	7	2
9	2	6	2	5	1

Original image 6x6



Result of the element-wise product and sum of the filter matrix and the original image

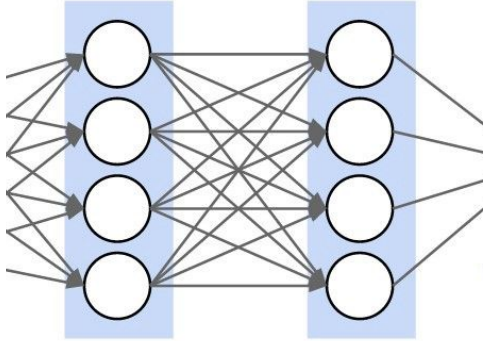
# Review: Recurrent Neural Network



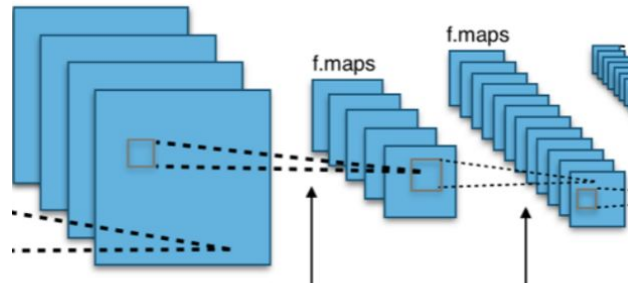
**Figure 9.2** Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep.

(Jurafsky, 2019)

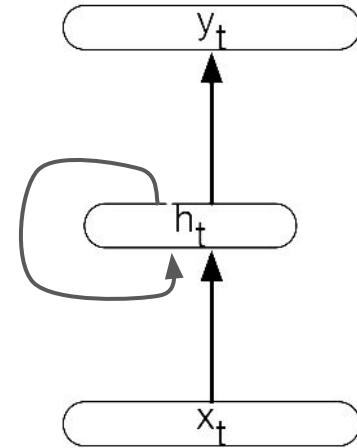
## FFN



## CNN



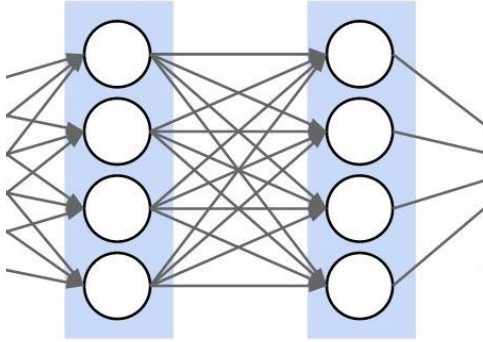
## RNN



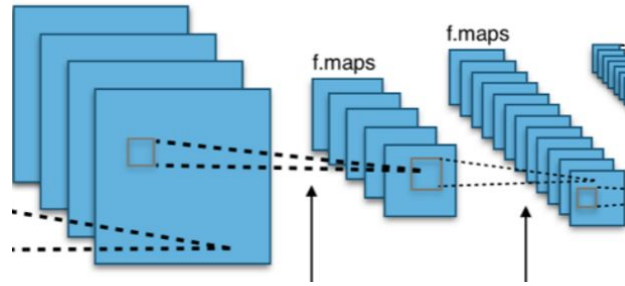
*Can model computation (e.g. matrix operations for a single input) be parallelized?*



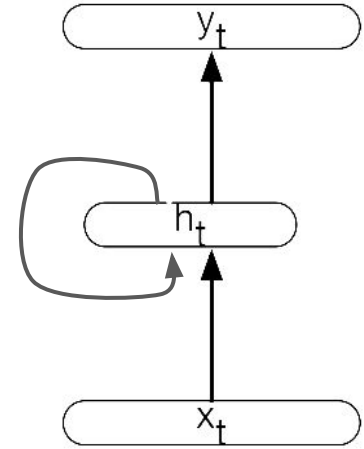
# FFN



# CNN



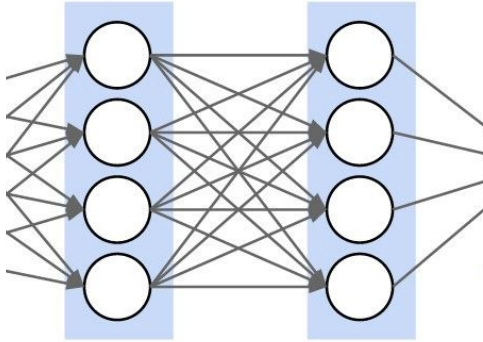
# RNN



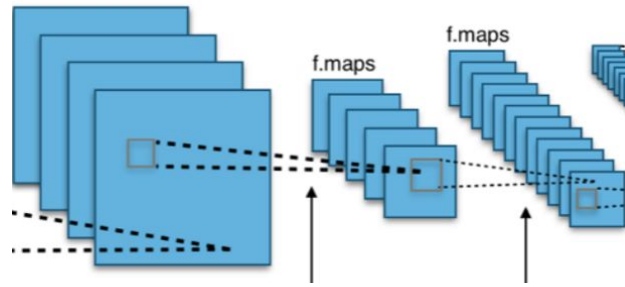
*Can model computation (e.g. matrix operations for a single input) be parallelized?*



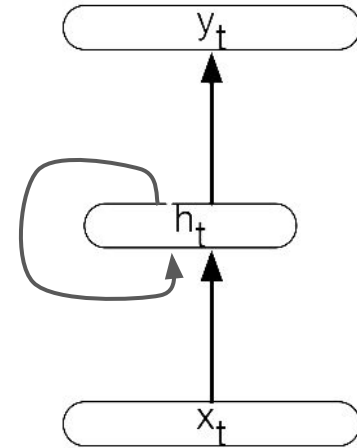
# FFN



# CNN



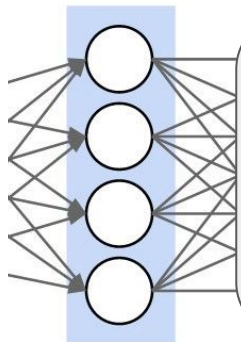
# RNN



*Can model computation (e.g. matrix operations for a single input) be parallelized?*



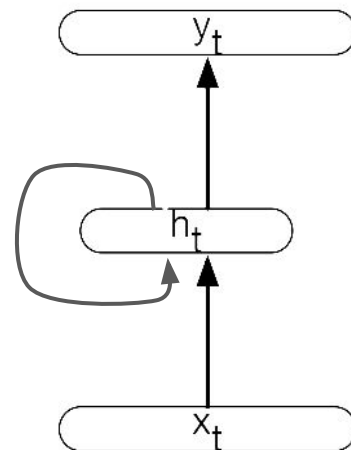
## FFN



## CNN

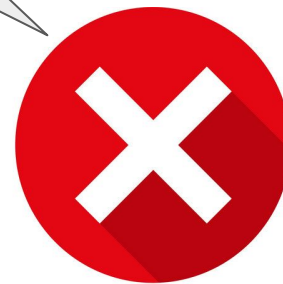


## RNN



*Ultimately limits how complex the model can be (i.e. it's total number of paramers/weights) as compared to a CNN.*

*Can model computation (e.g. matrix operations for a single input) be parallelized?*





# The Transformer: “Attention-only” models

**Can handle sequences and long-distance dependencies, but....**

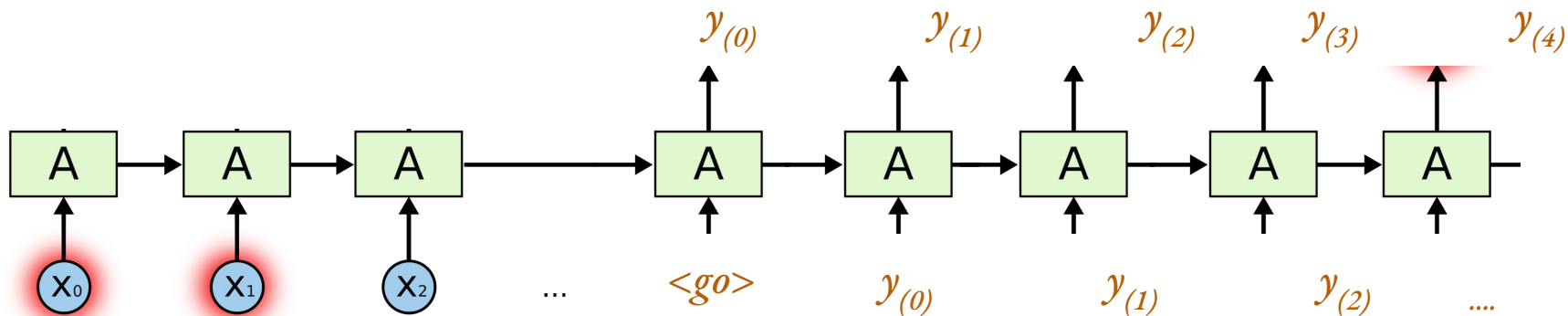
- Don't want complexity of LSTM/GRU cells
- Constant num edges between input steps
- Enables “interactions” (i.e. adaptations) between words
- **Easy to parallelize -- don't need sequential processing.**

# The Transformer: “Attention-only” models

Challenge:

*The ball was kicked by kayla.*

- Long distance dependency when translating:



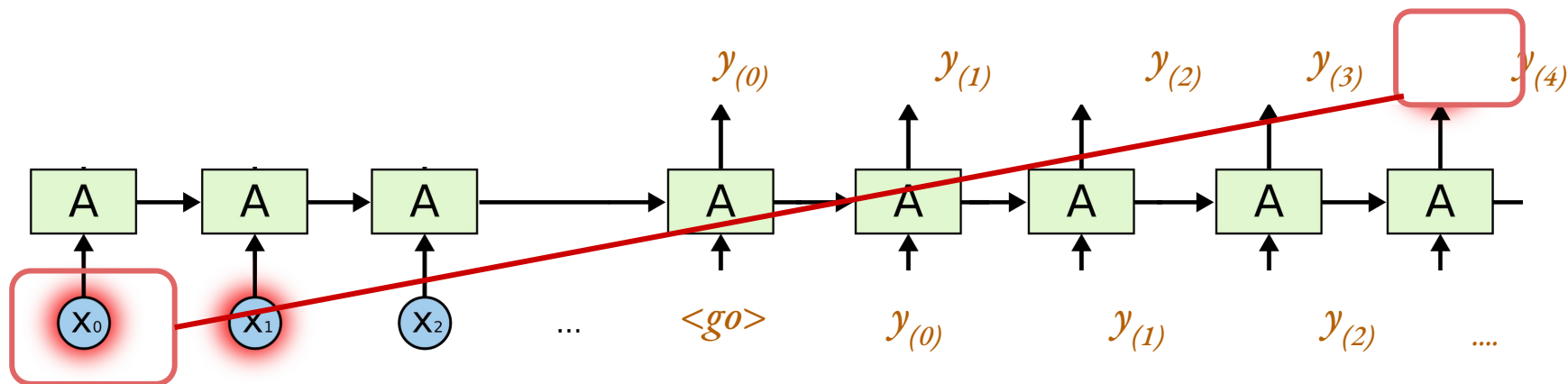
*Kayla kicked the ball.*

# The Transformer: “Attention-only” models

Challenge:

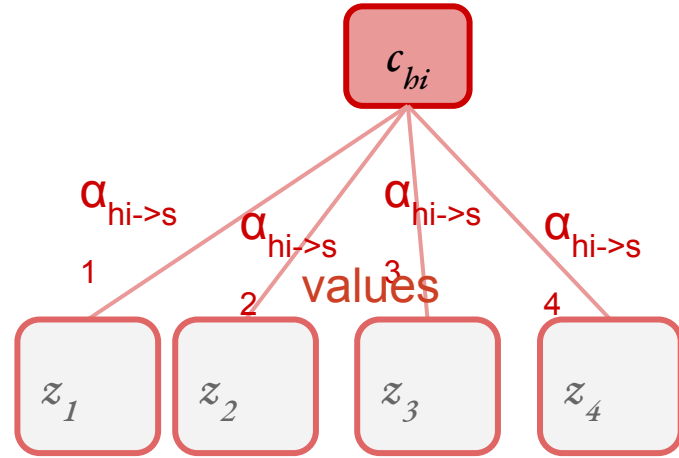
*The ball was kicked by kayla.*

- Long distance dependency when translating:

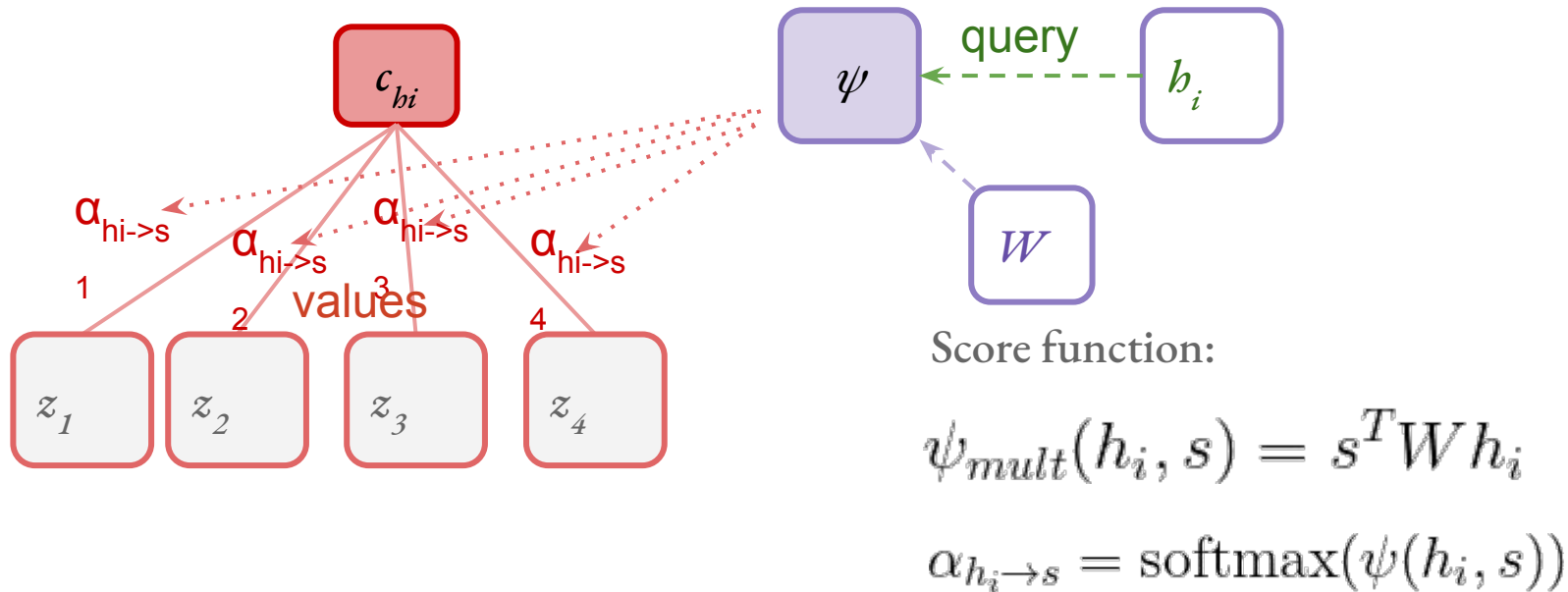


*Kayla kicked the ball.*

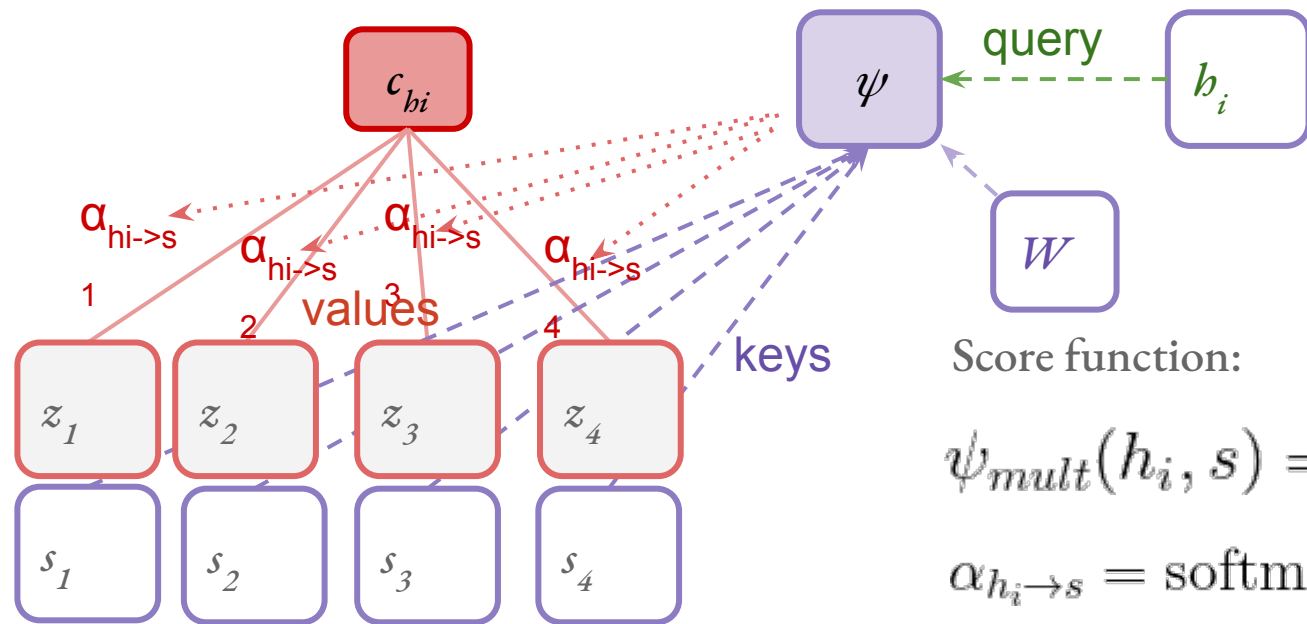
# Attention



# Attention



# Attention



Score function:

$$\psi_{mult}(h_i, s) = s^T W h_i$$

$$\alpha_{h_i \rightarrow s} = \text{softmax}(\psi(h_i, s))$$

$$c_{h_i} = \sum_{n=1}^{|s|} \alpha_{h_i \rightarrow s_n} z_n$$

# The Transformer: “**Attention**-only” models

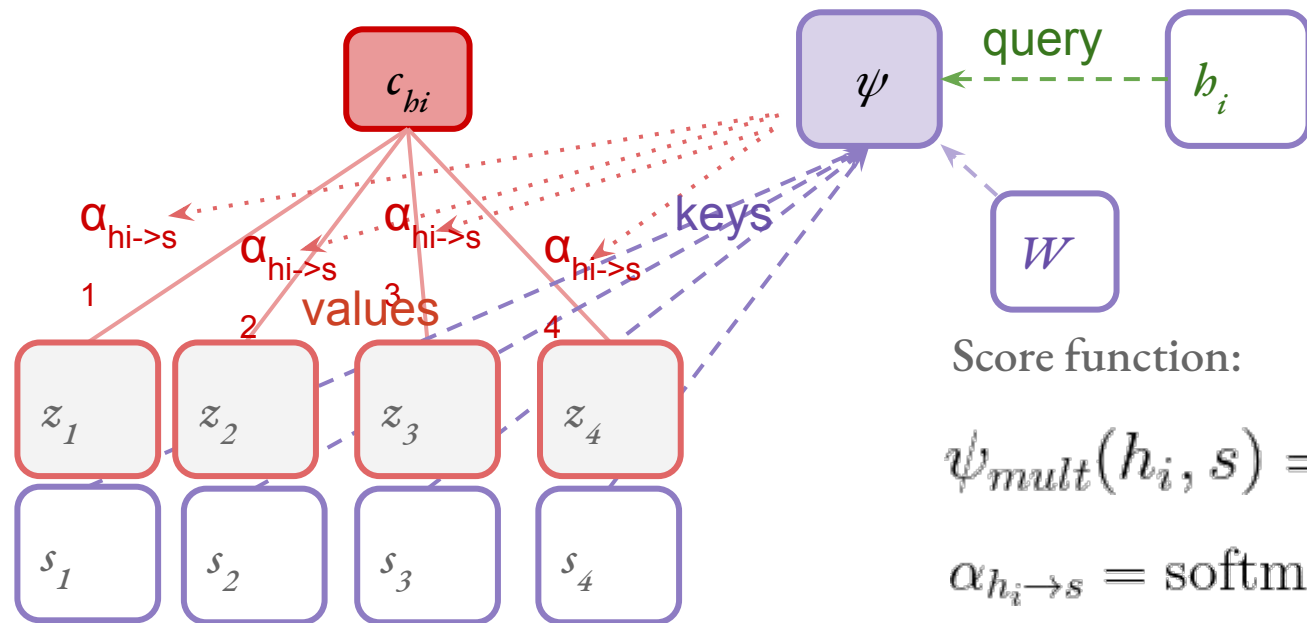
Challenge:

- Long distance dependency when translating:

Attention came about for encoder decoder models.

Then self-attention was introduced:

# Attention



Score function:

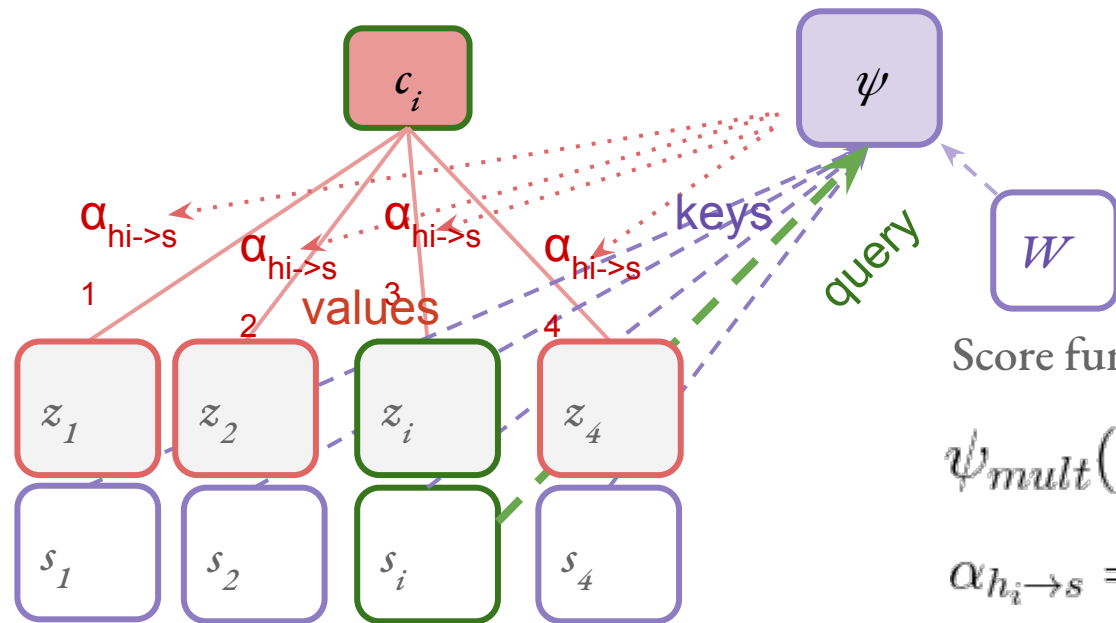
$$\psi_{mult}(h_i, s) = s^T W h_i$$

$$\alpha_{h_i \rightarrow s} = \text{softmax}(\psi(h_i, s))$$

$$c_{h_i} = \sum_{n=1}^{|s|} \alpha_{h_i \rightarrow s_n} z_n$$



# Self-Attention



Score function:

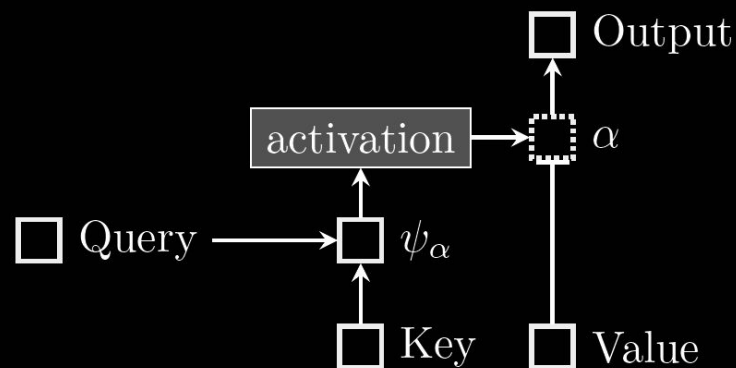
$$\psi_{mult}(h_i, s) = s^T W h_i$$

$$\alpha_{h_i \rightarrow s} = \text{softmax}(\psi(h_i, s))$$

$$c_{h_i} = \sum_{n=1}^{|s|} \alpha_{h_i \rightarrow s_n} z_n$$

# The Transformer: “Attention-only” models

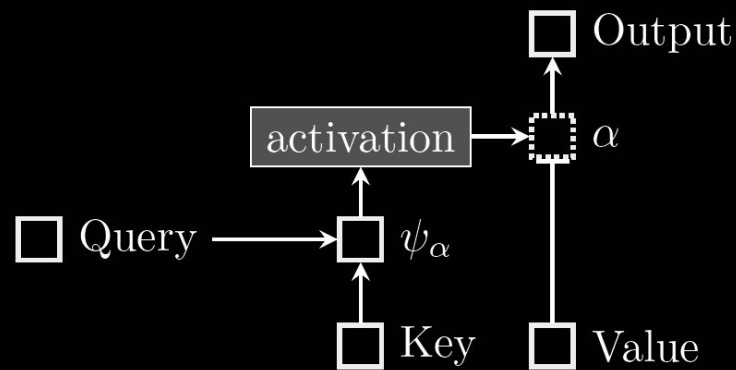
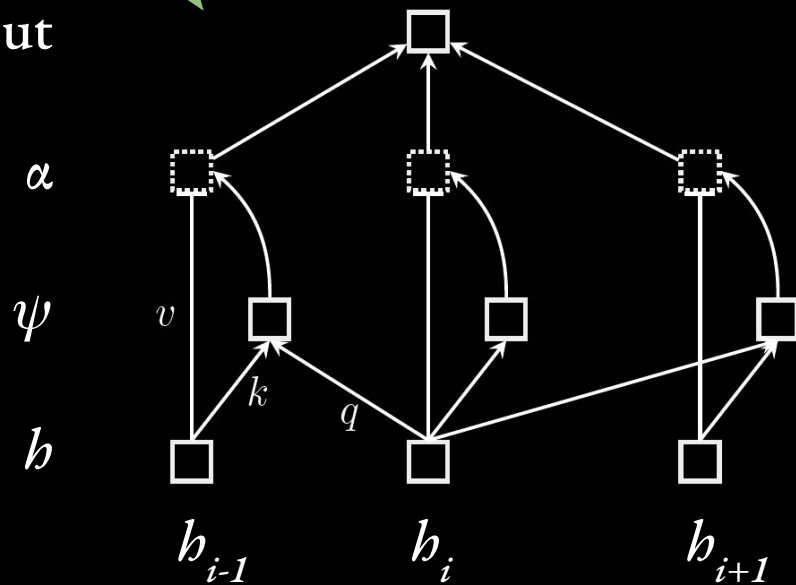
Attention as weighting a value based on a query and key:



(Eisenstein, 2018)

# The Transformer: “Attention-only” models

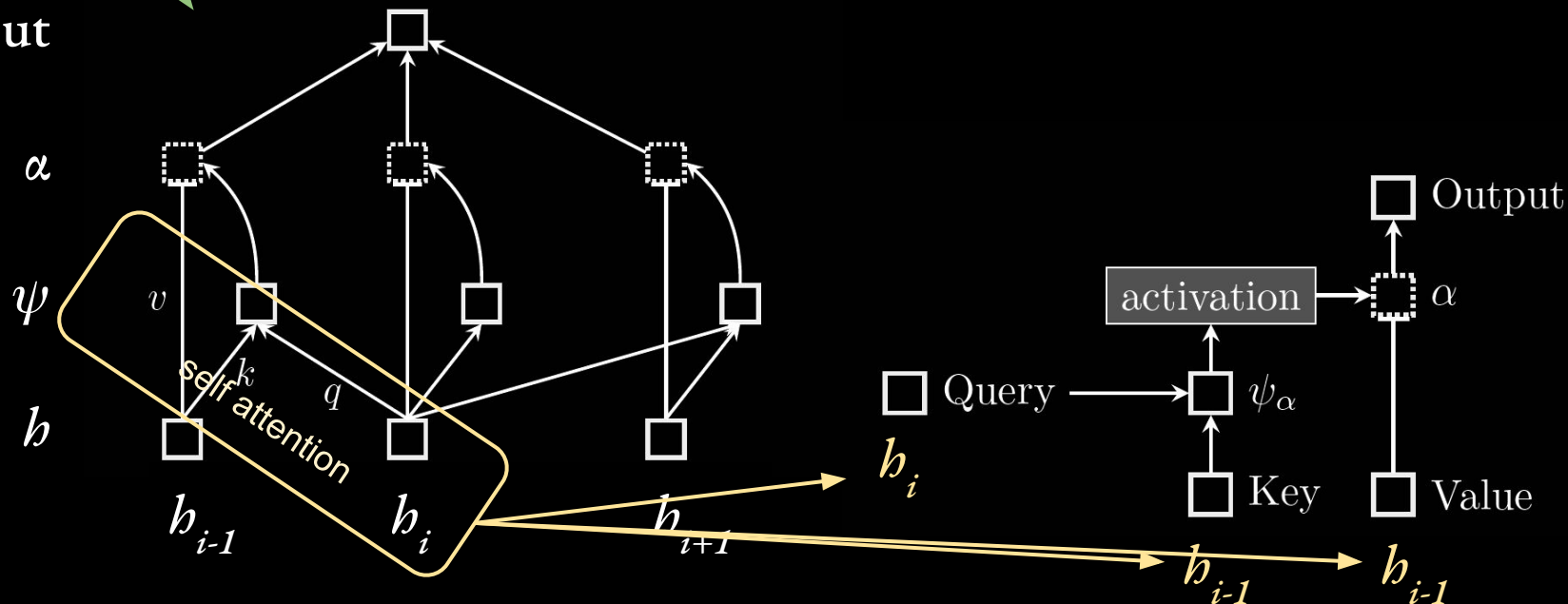
Output



(Eisenstein, 2018)

# The Transformer: “Attention-only” models

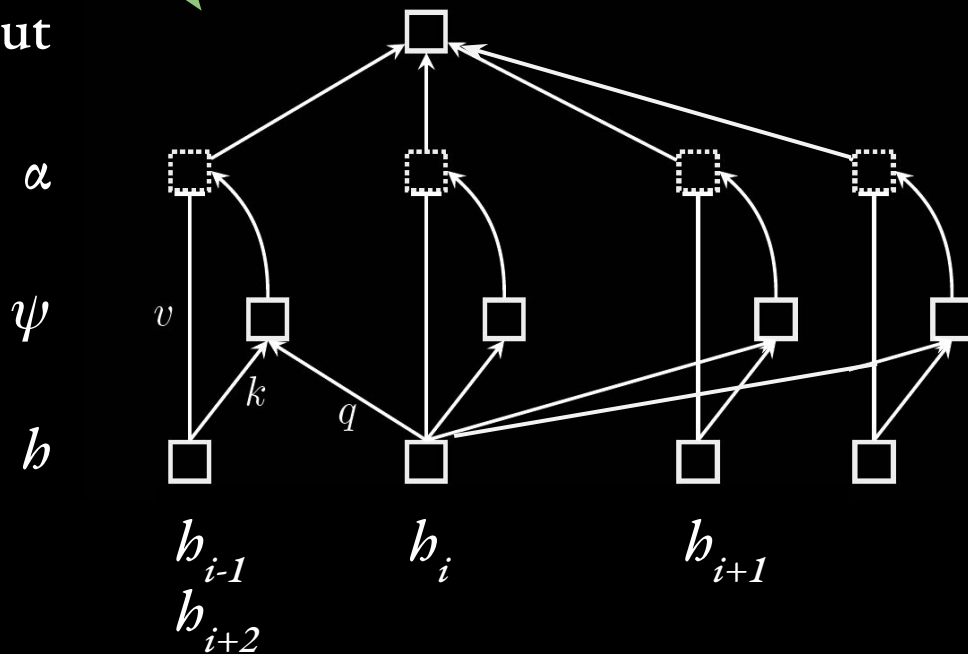
Output



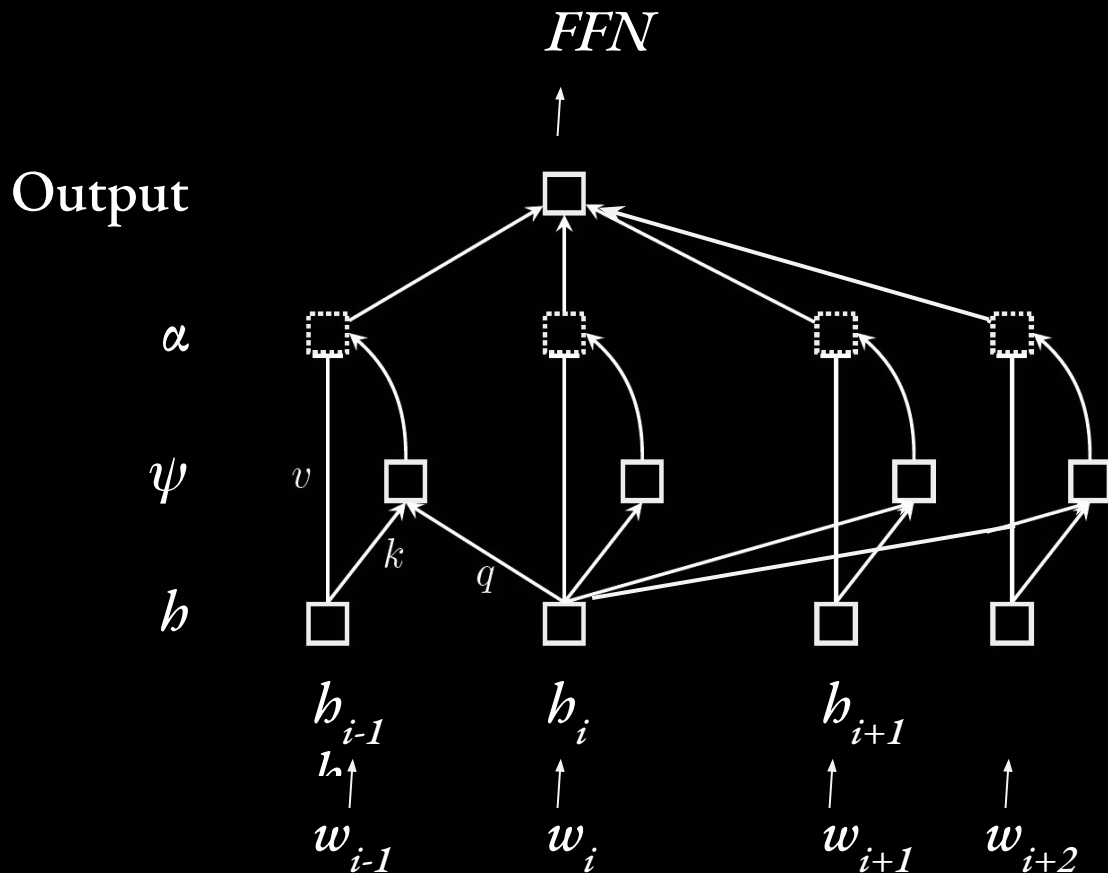
(Eisenstein, 2018)

# The Transformer: “Attention-only” models

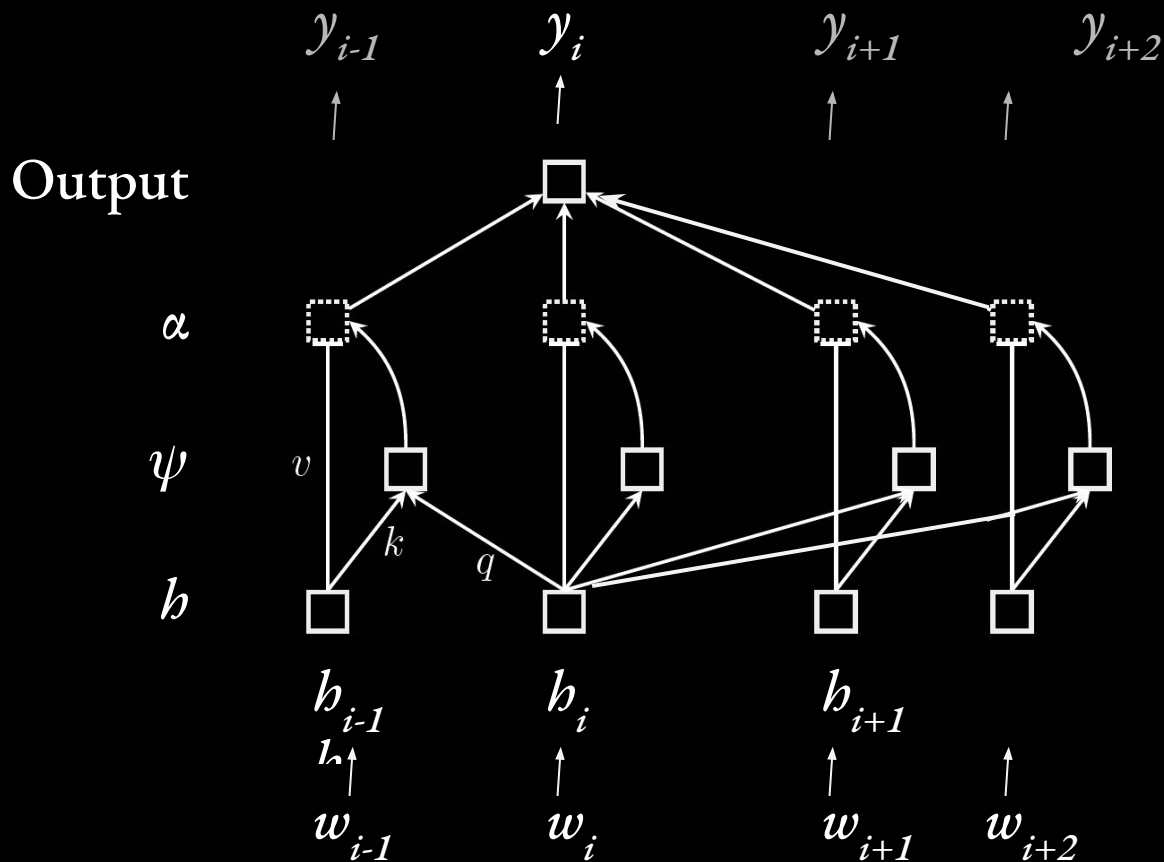
Output



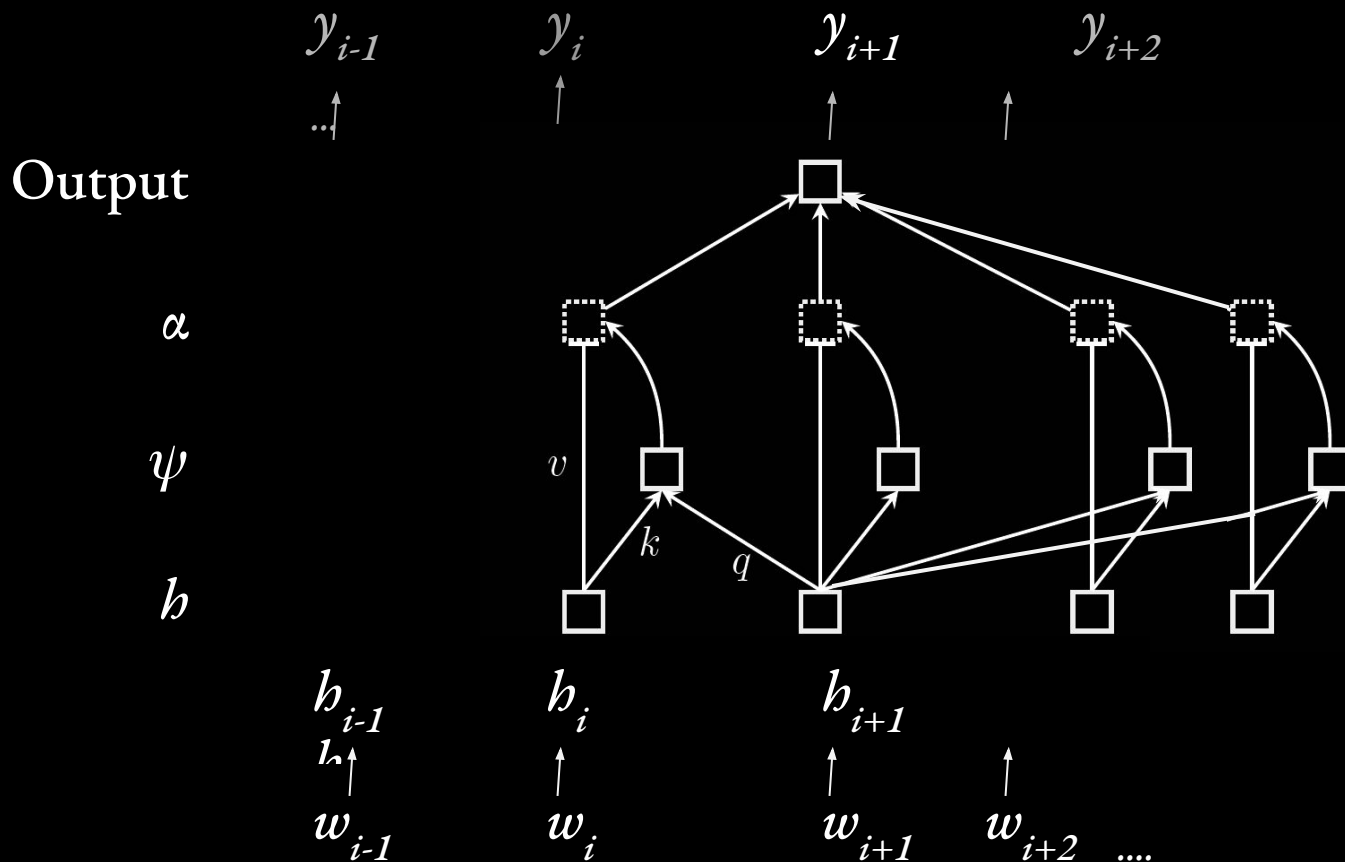
# The Transformer: “Attention-only” models



# The Transformer: “Attention-only” models

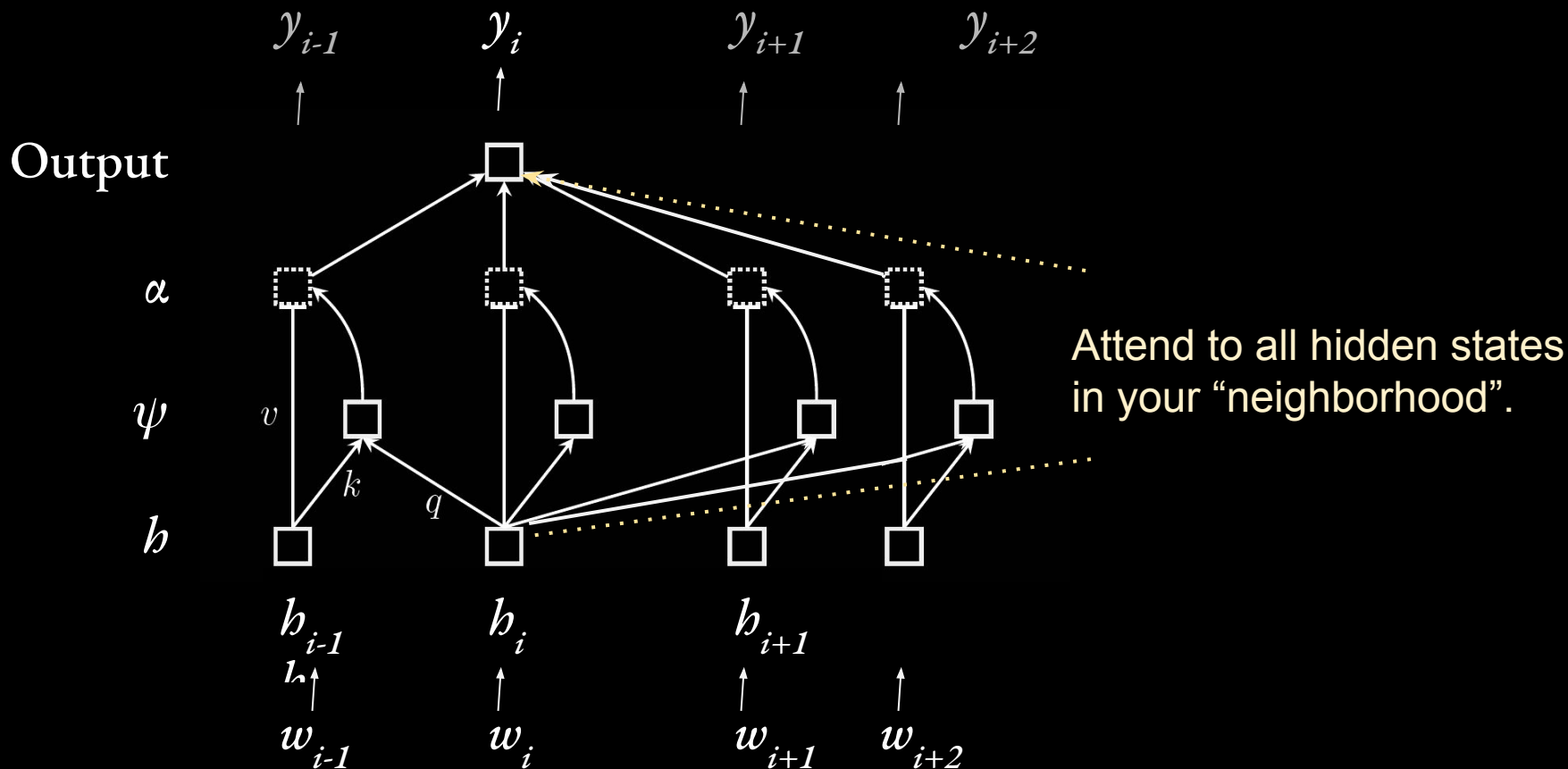


# The Transformer: “Attention-only” models

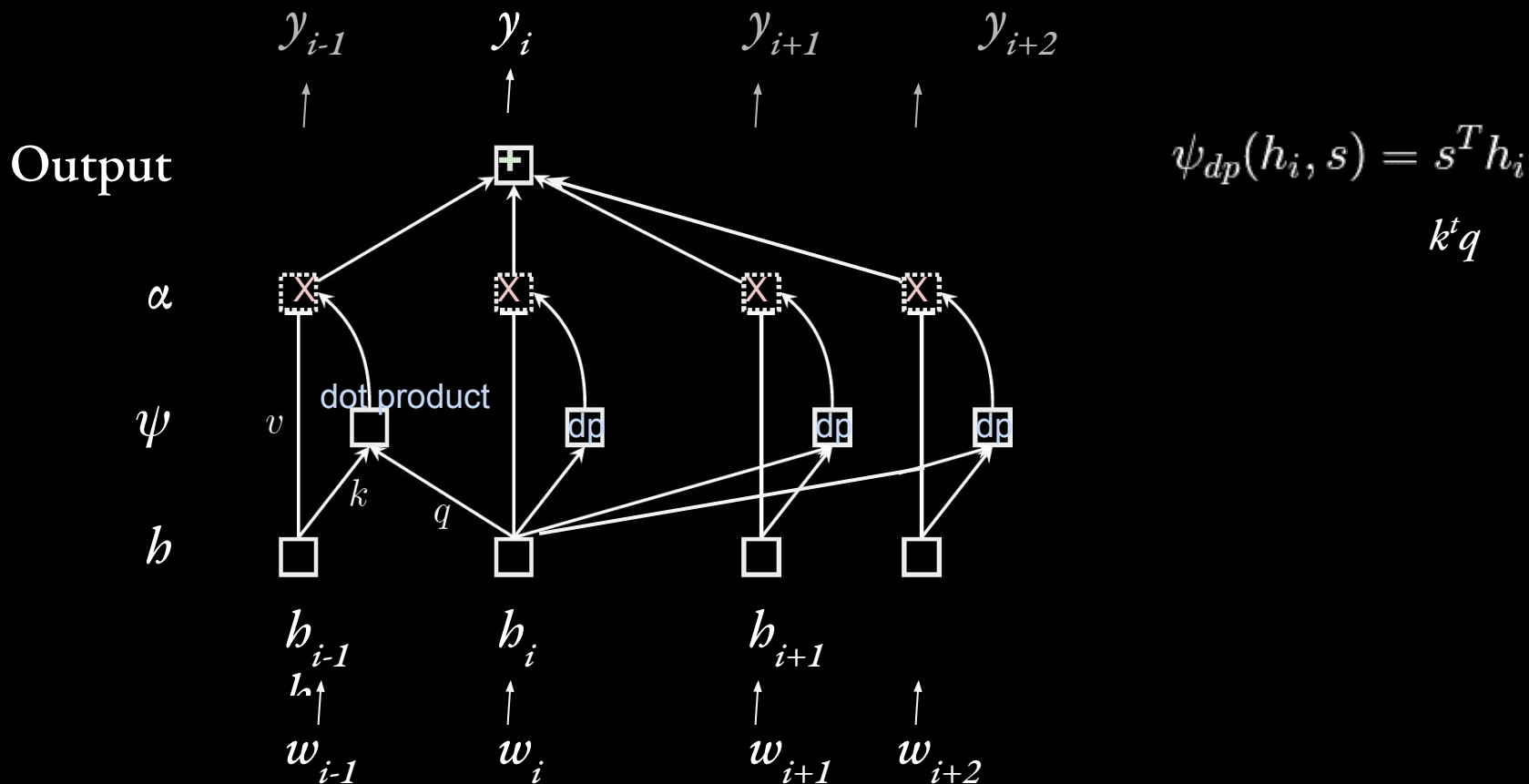




# The Transformer: “Attention-only” models

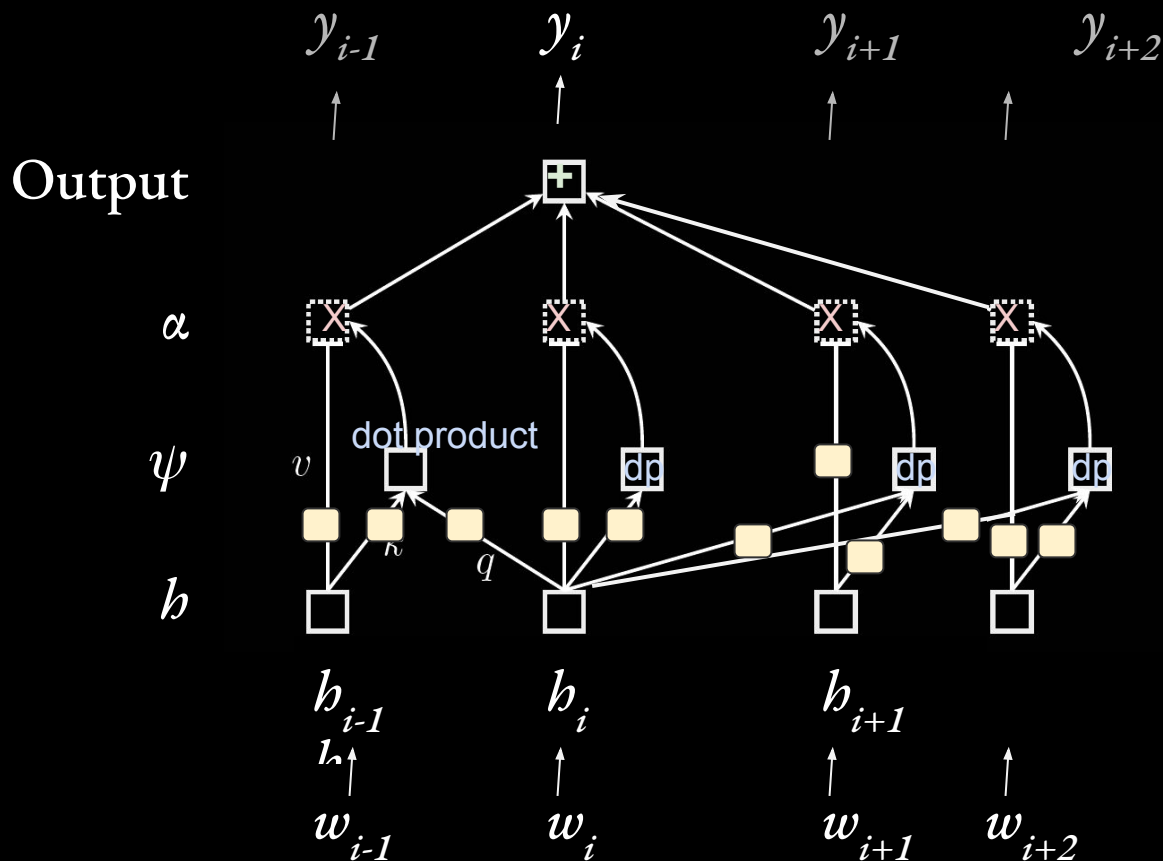


# The Transformer: "Attention-only" models





# The Transformer: "Attention-only" models



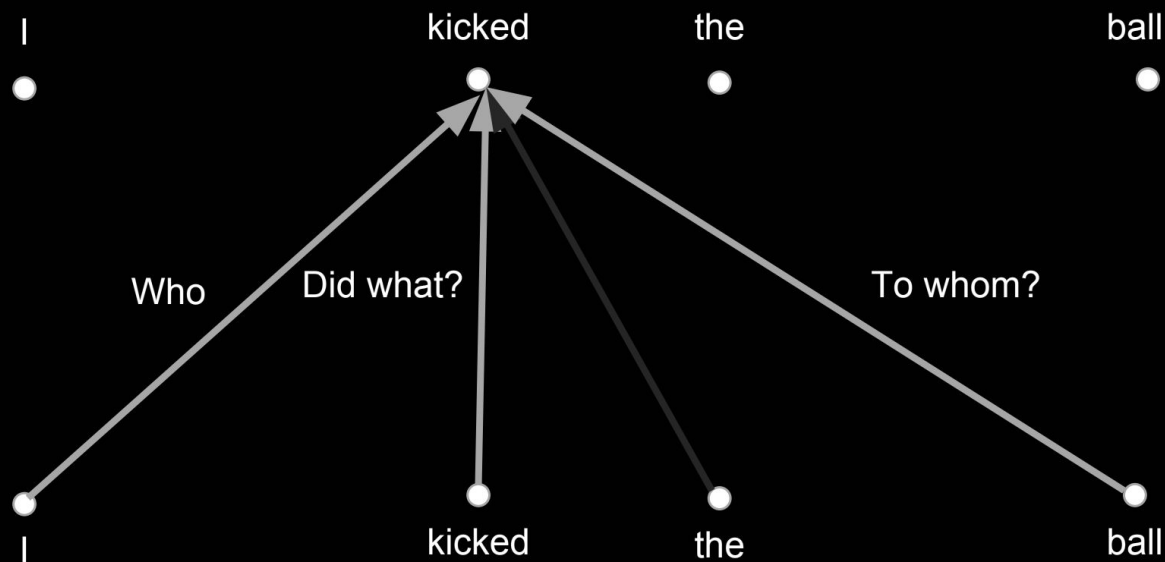
$$\psi_{dp}(k, q) = (k^t q) \sigma$$

Linear layer:  
 $W^T X$

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

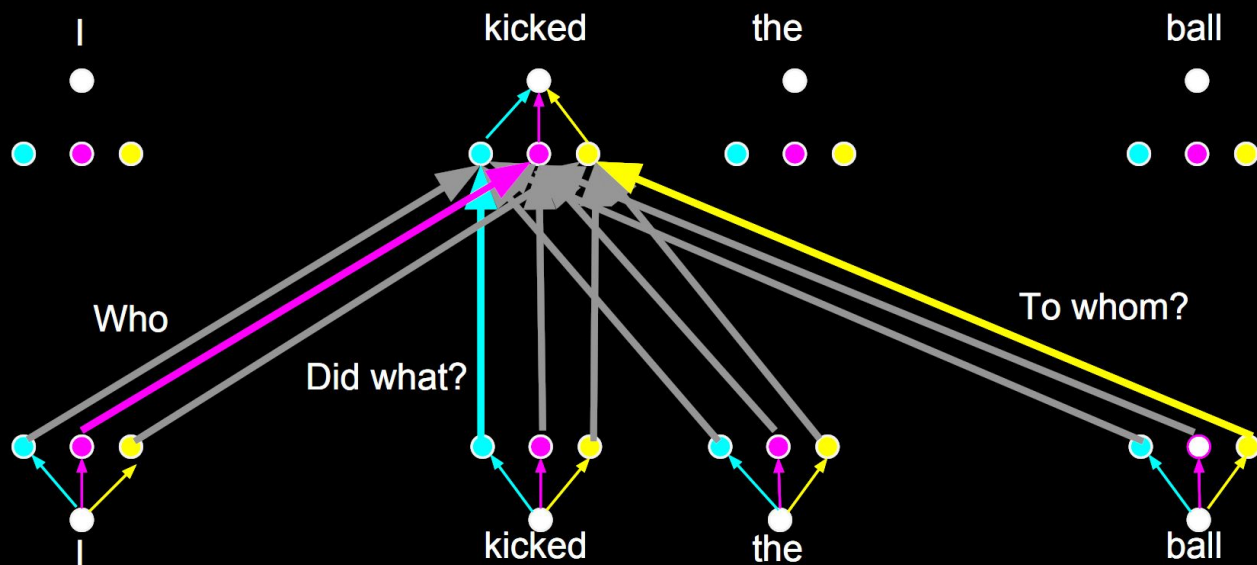
# The Transformer

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

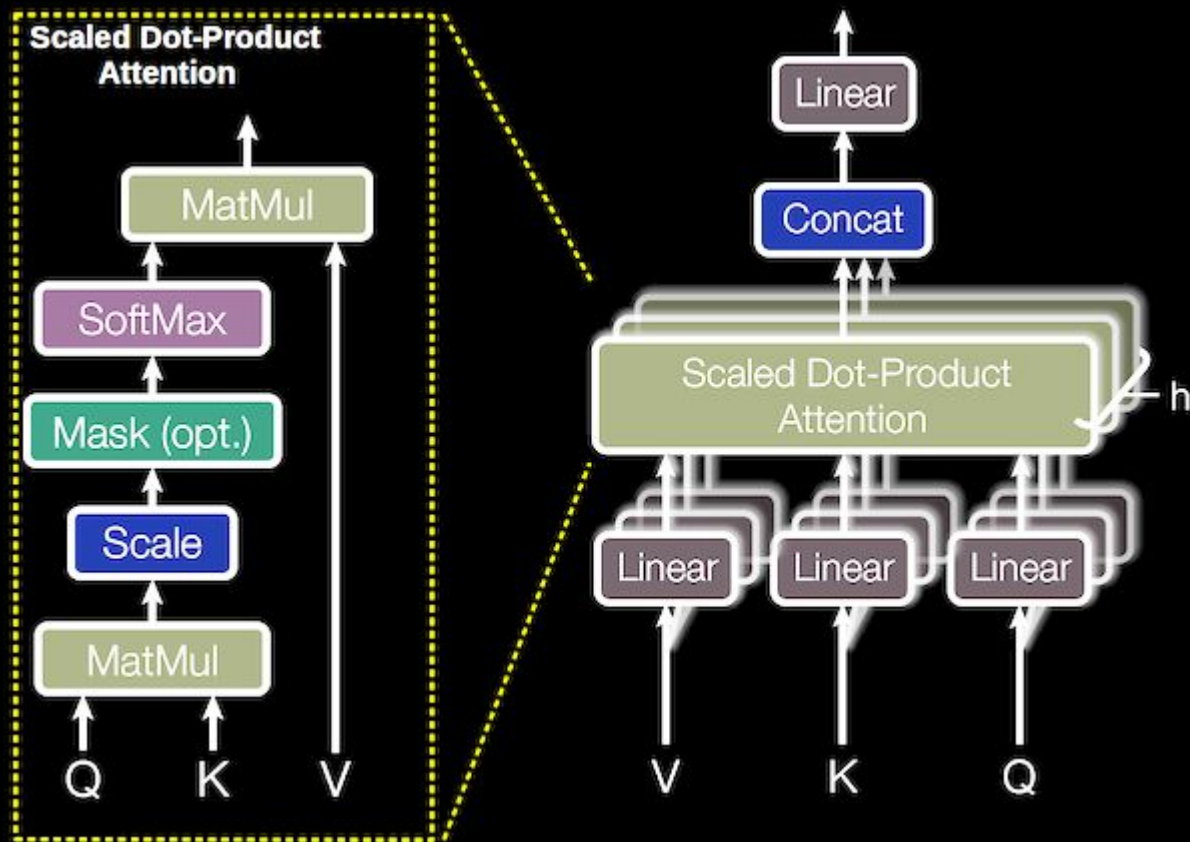


# The Transformer

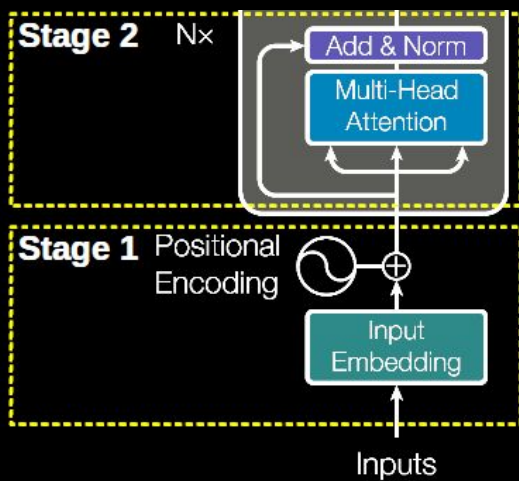
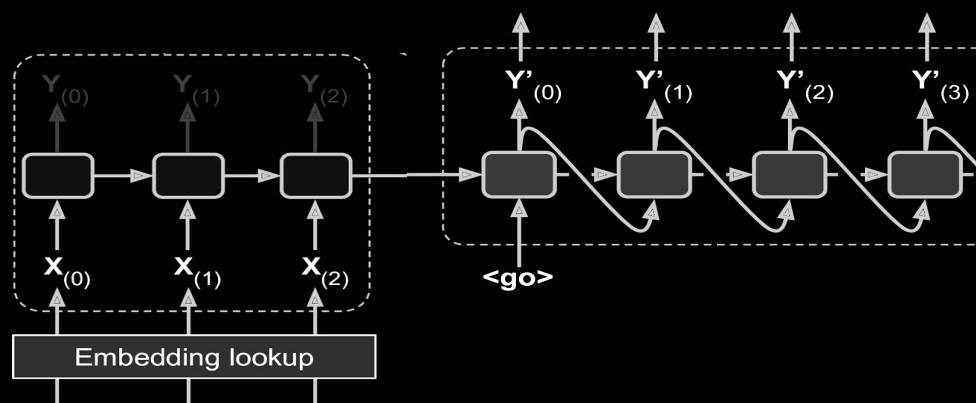
Solution: Multi-head attention



# Multi-head Attention

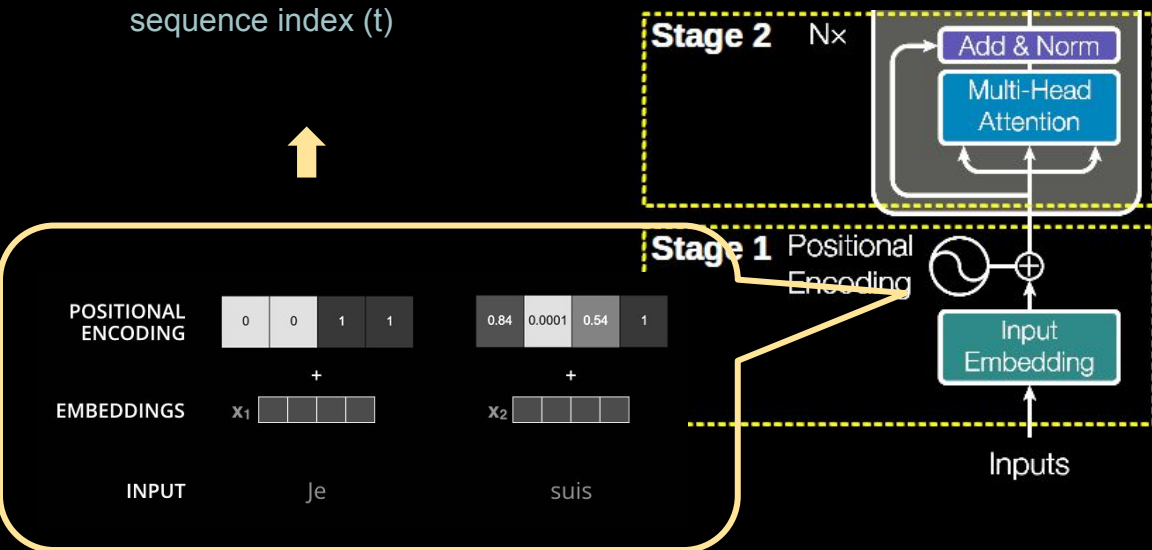
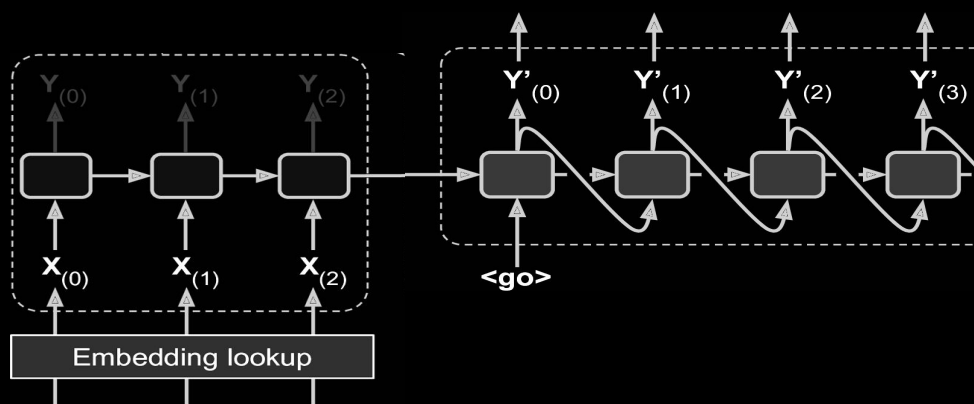
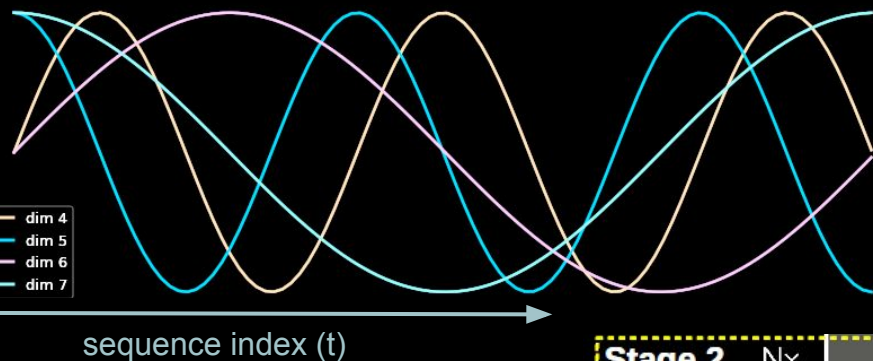


# Transformer for Encoder-Decoder

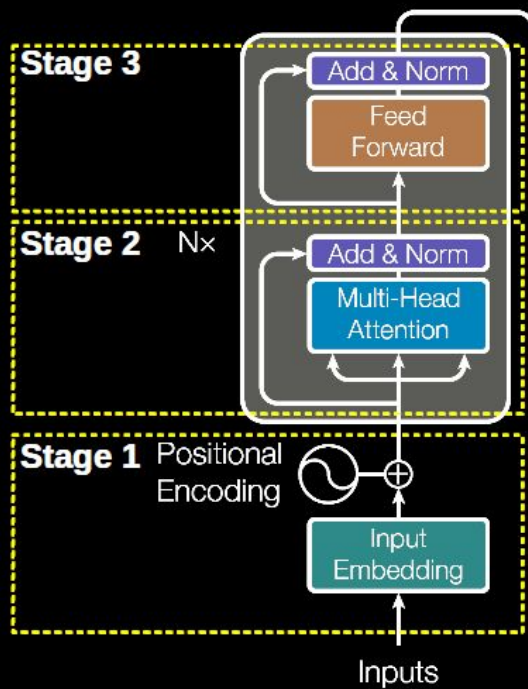
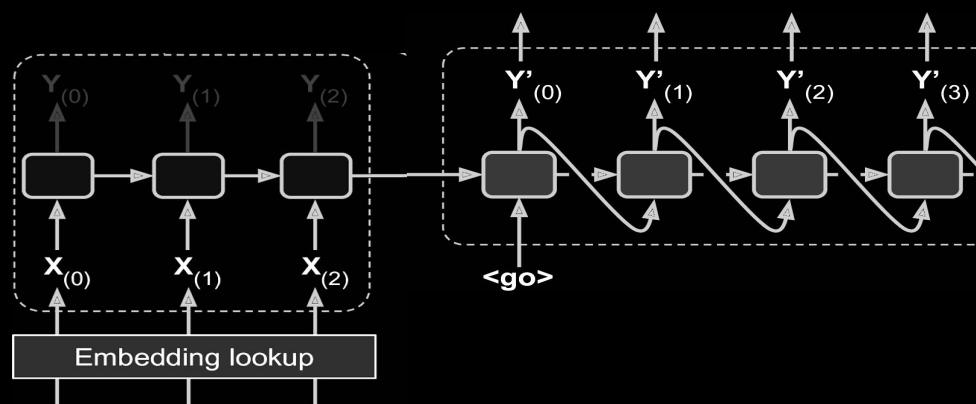




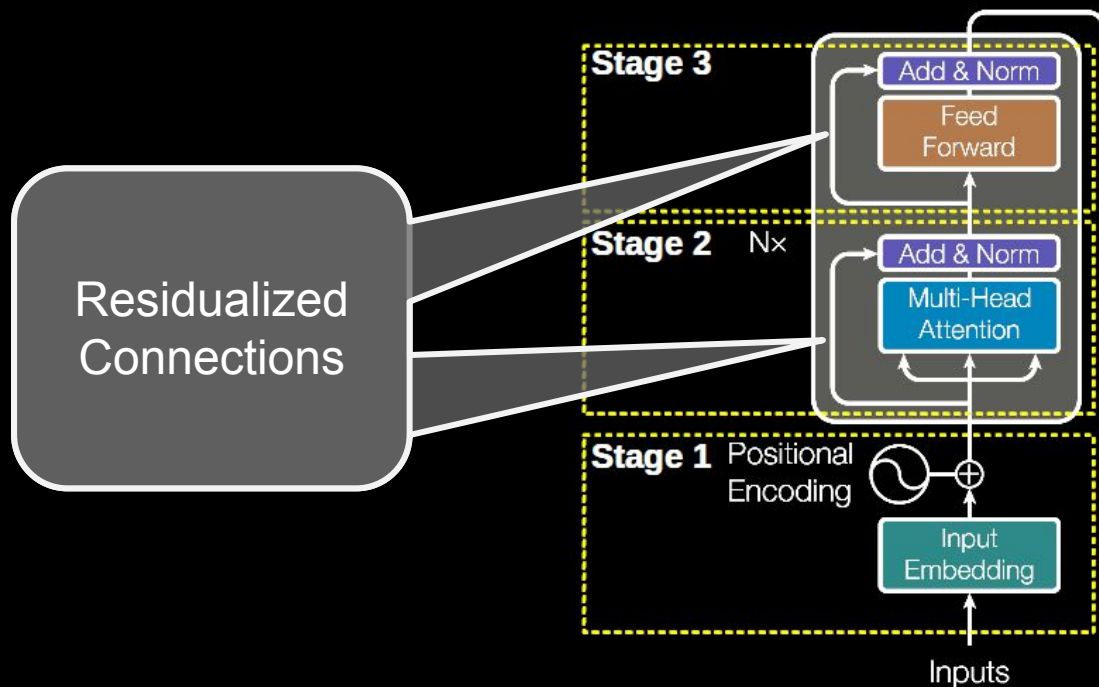
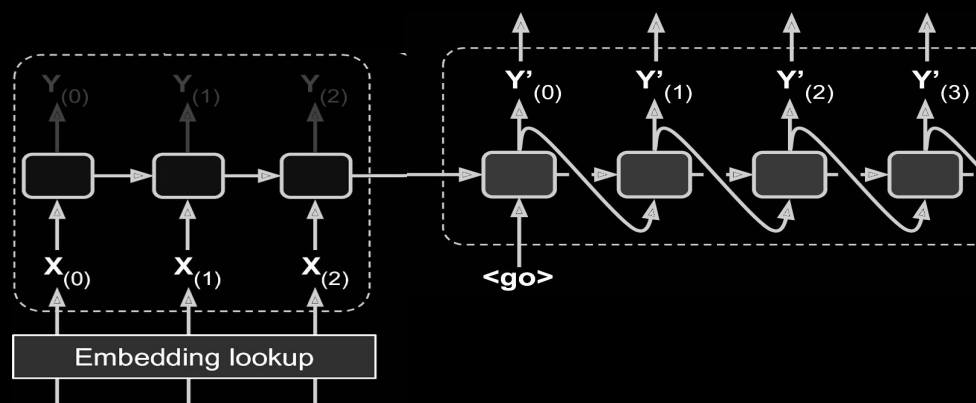
# Transformer for Encoder-Decoder



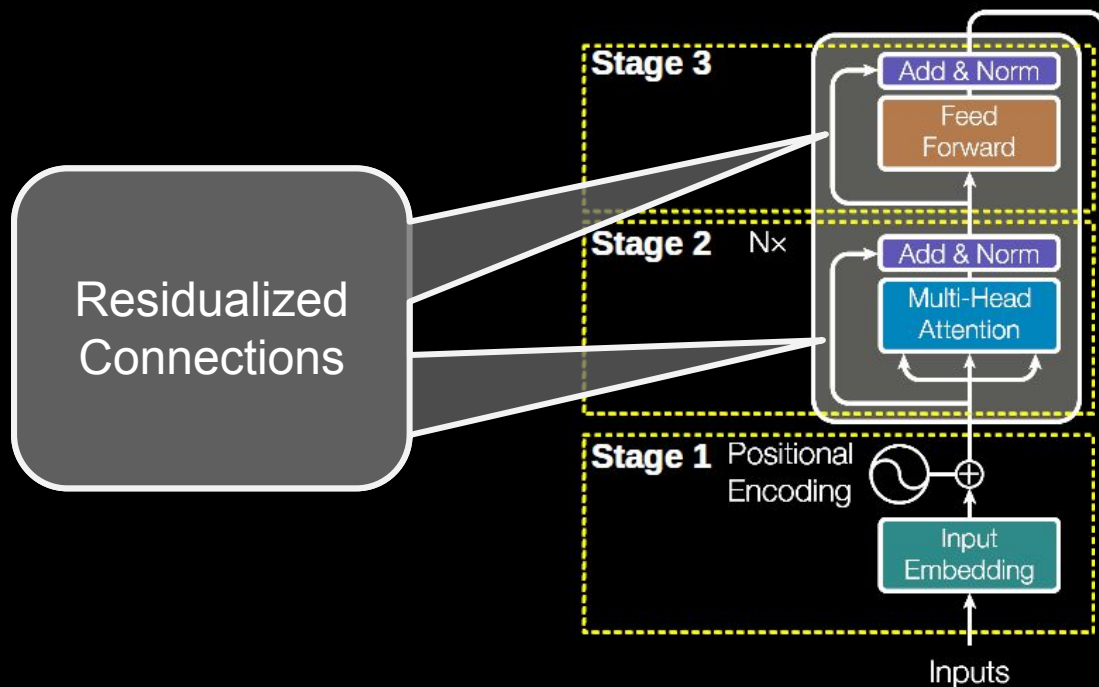
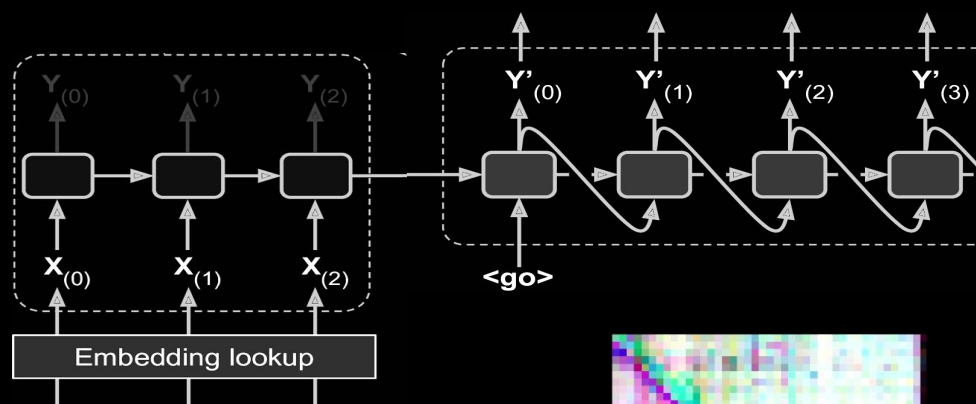
# Transformer for Encoder-Decoder



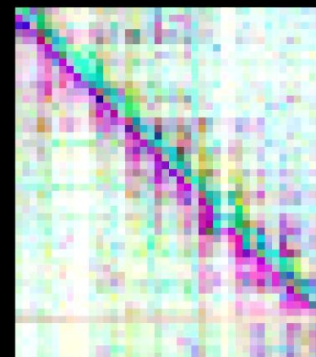
# Transformer for Encoder-Decoder



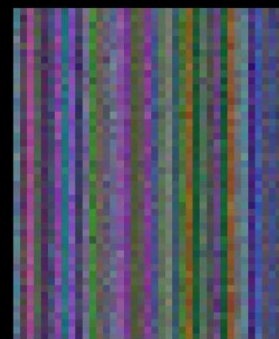
# Transformer for Encoder-Decoder



residuals enable positional information to be passed along

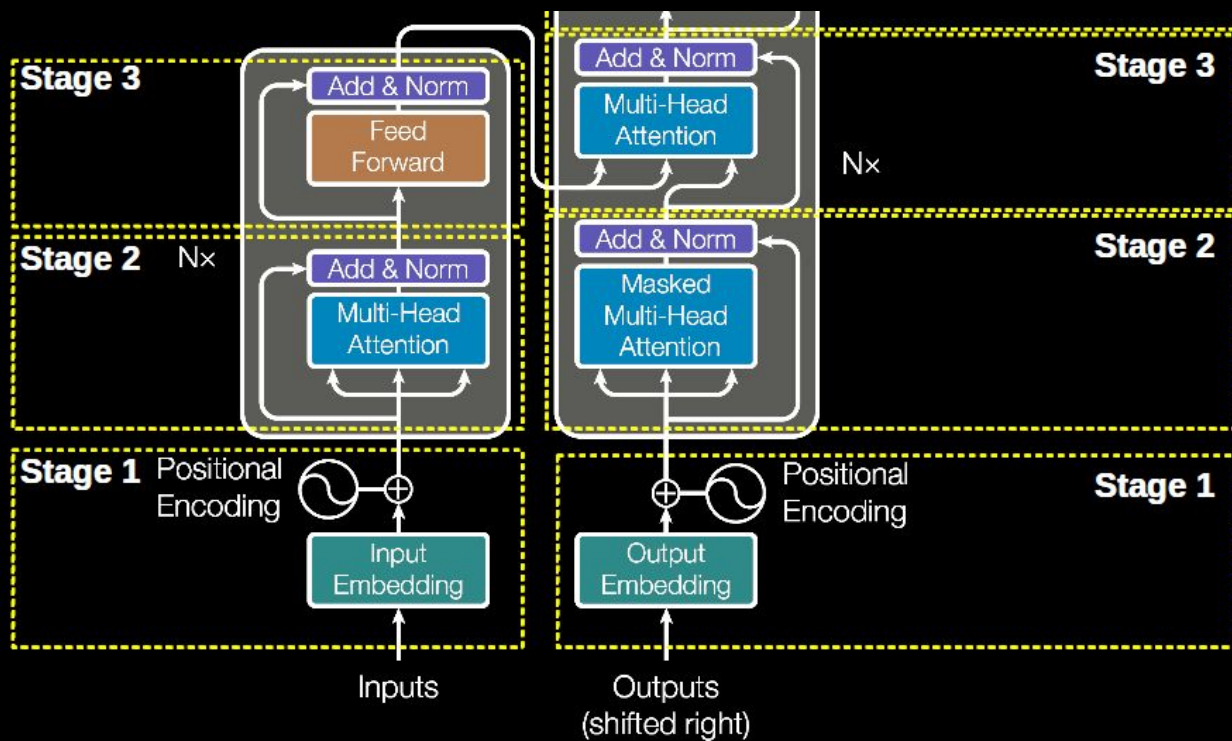
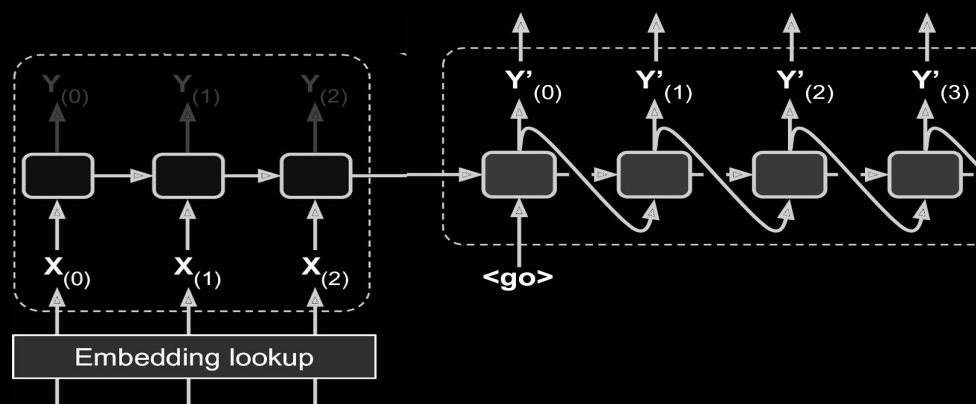


With residuals

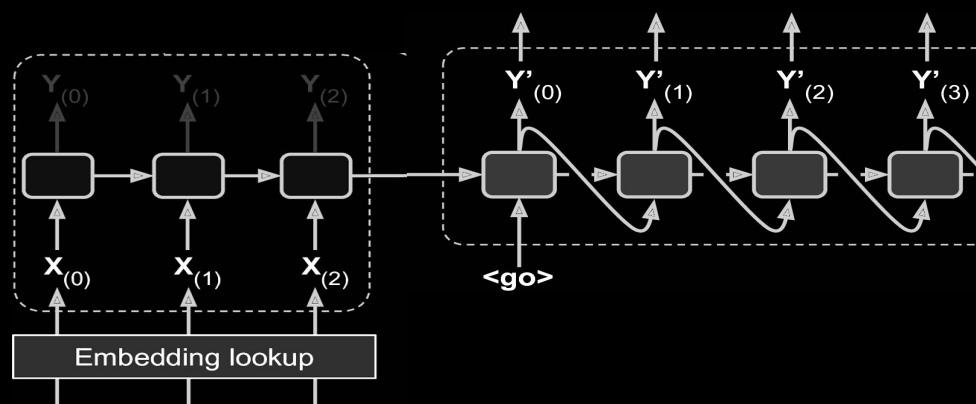


Without residuals

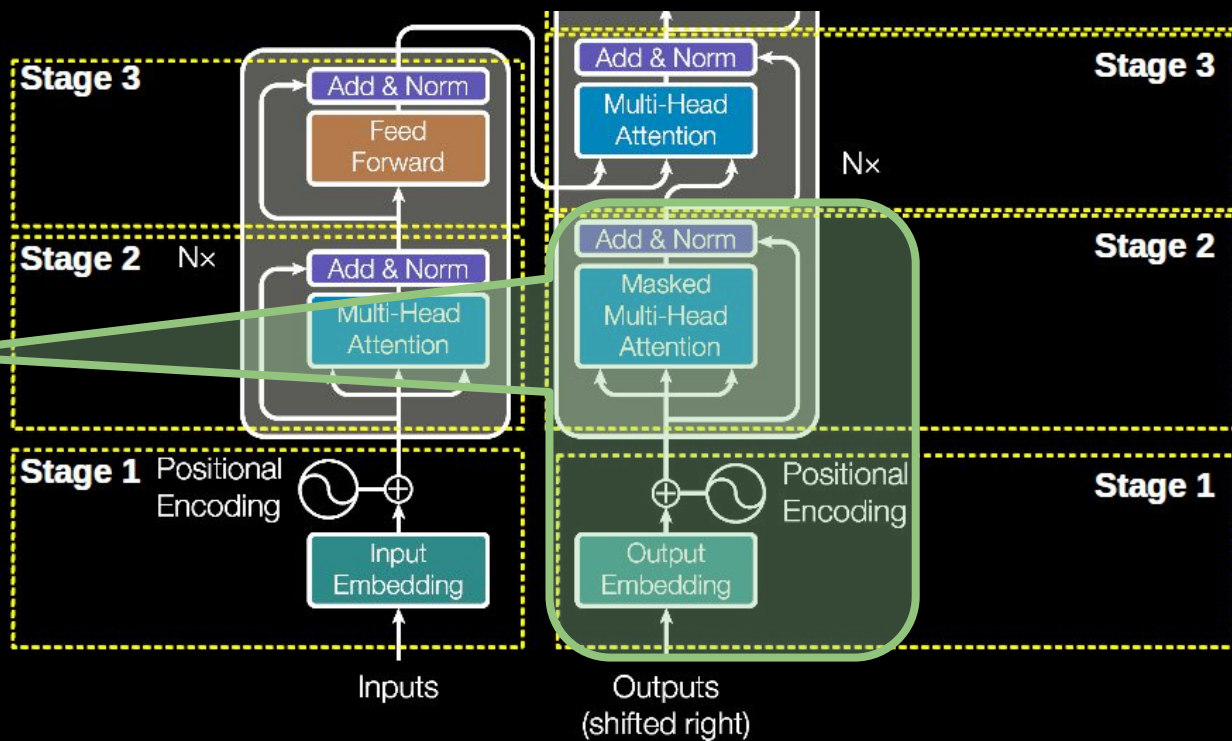
# Transformer for Encoder-Decoder



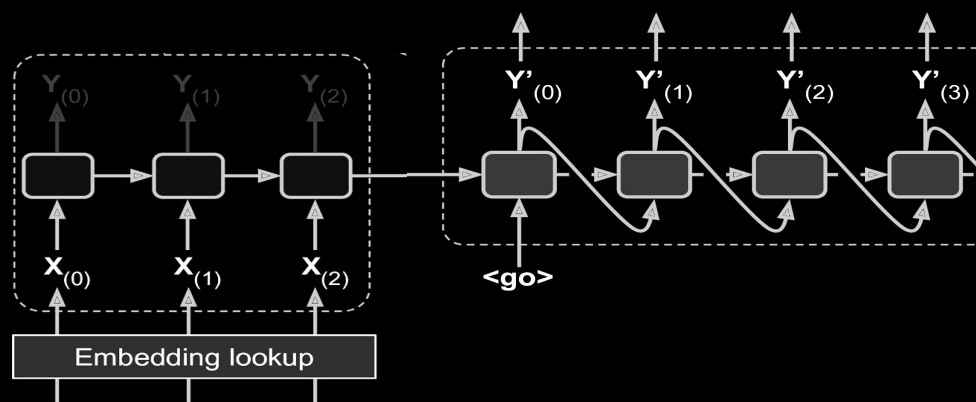
# Transformer for Encoder-Decoder



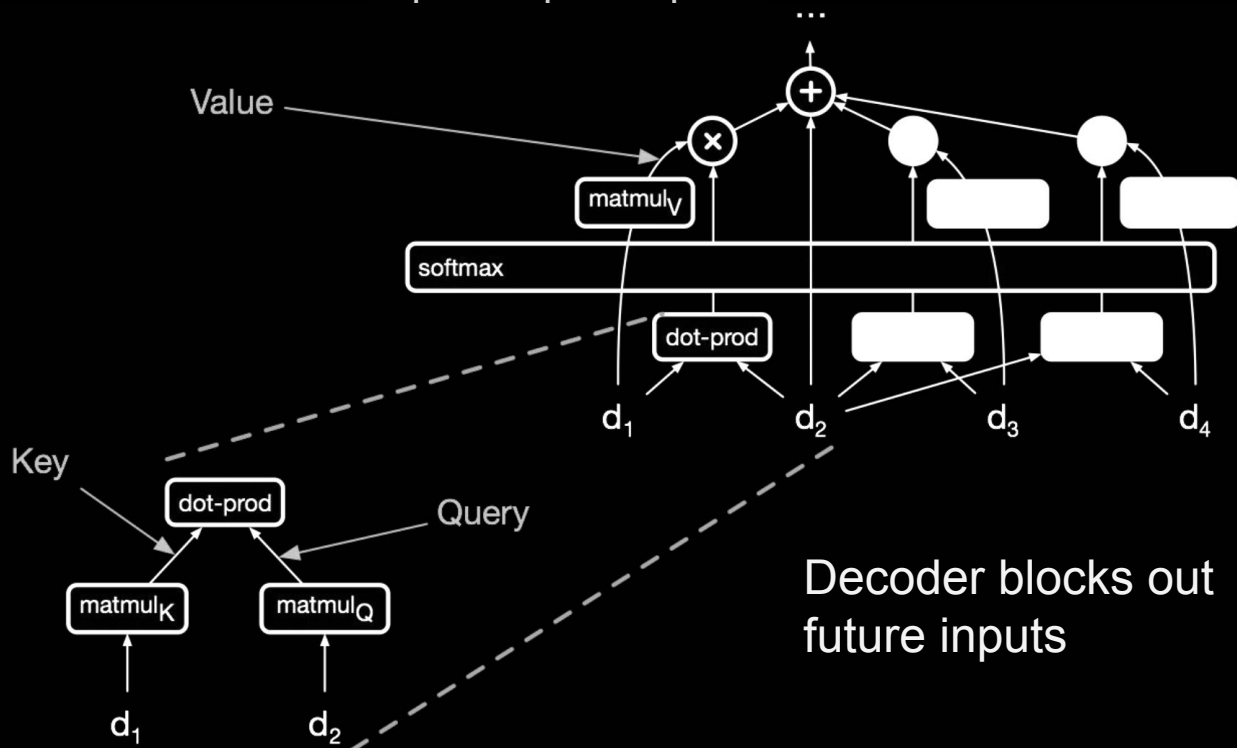
essentially, a language model



# Transformer for Encoder-Decoder

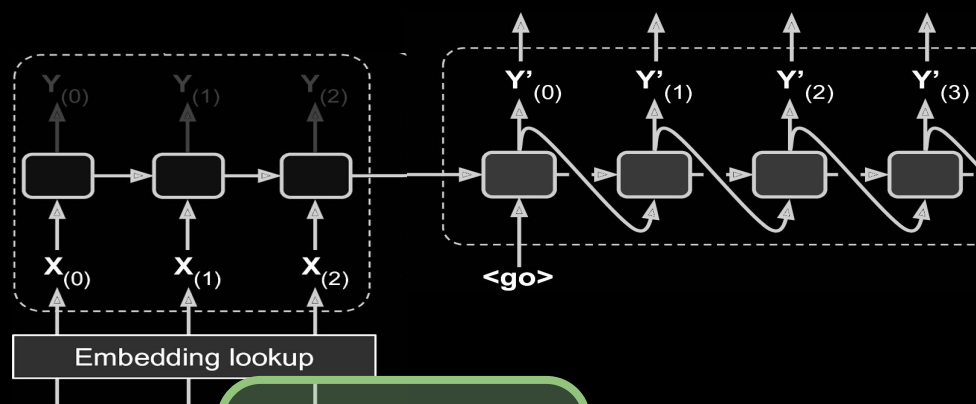


essentially, a language model



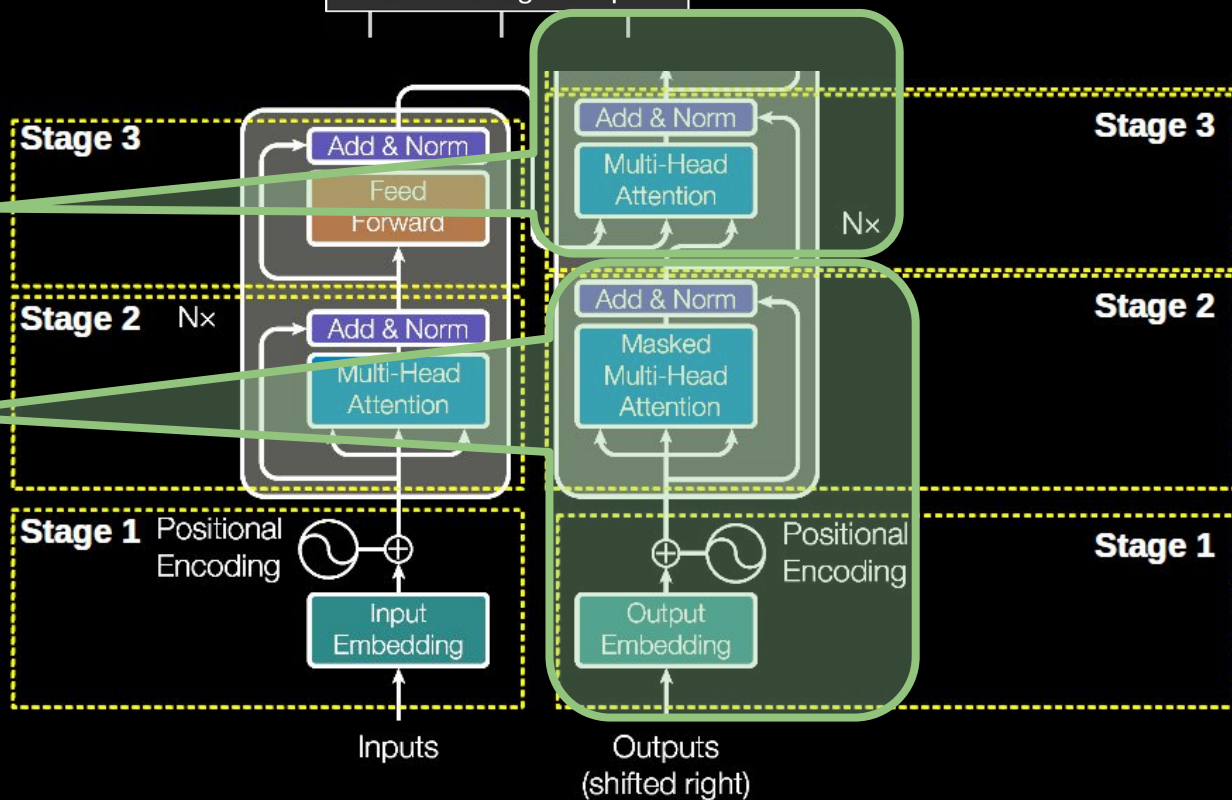
Decoder blocks out future inputs

# Transformer for Encoder-Decoder



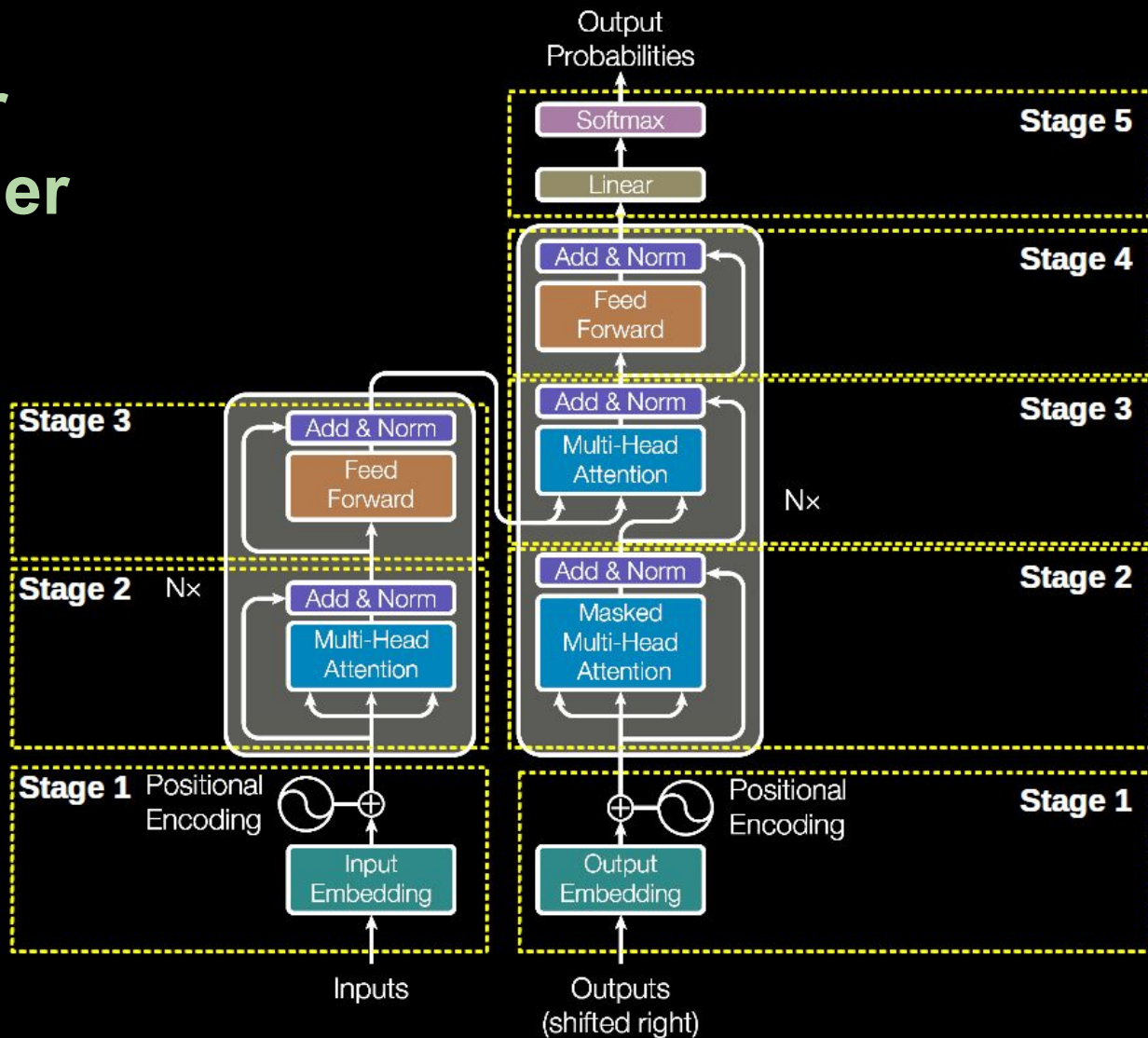
Add conditioning of the LM based on the encoder

essentially, a language model





# Transformer for Encoder-Decoder



# Transformer (as of 2017)

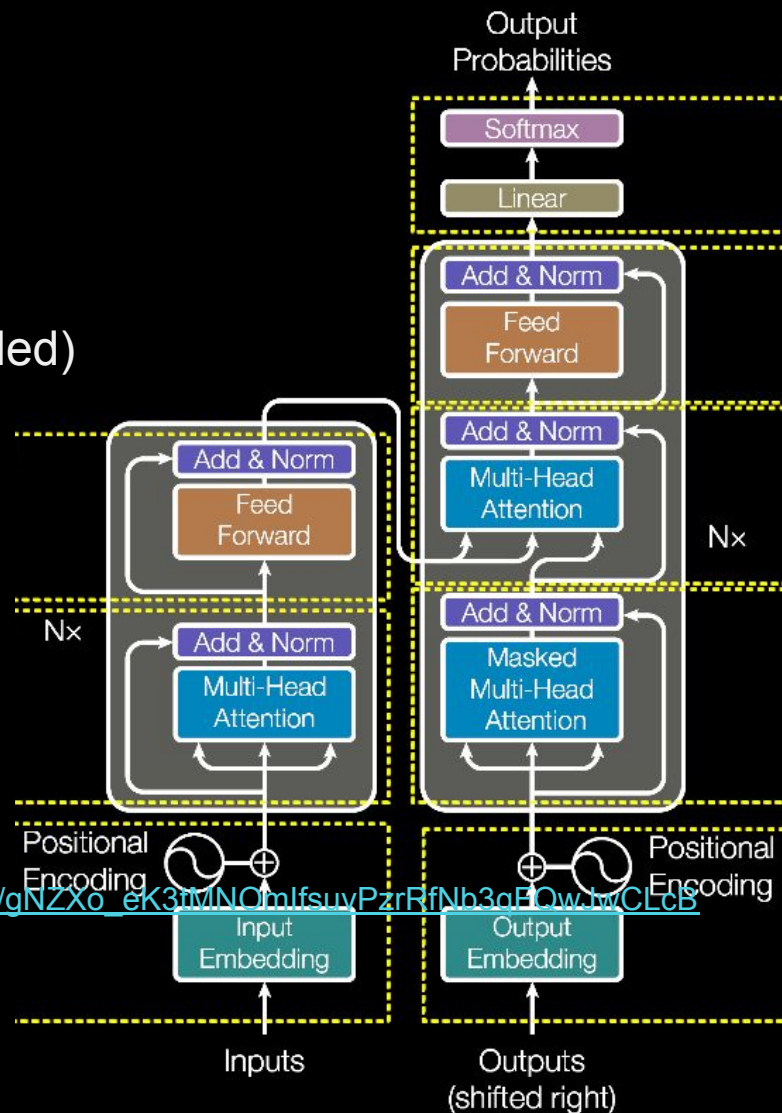
“WMT-2014” Data Set. BLEU scores:

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

# Transformer

- Utilize Self-Attention
- Simple att scoring function (dot product, scaled)
- Added linear layers for Q, K, and V
- Multi-head attention
- Added positional encoding
- Added residual connection
- Simulate decoding by masking

[https://4.bp.blogspot.com/-OlrV-PAtEkQ/W3RkOJCBkaI/AAAAAAAAADOg/gNZXo\\_eK3tMNOmIfsuyPzrRfNb3qEQwIwCLcB/GAs/s640/image1.gif](https://4.bp.blogspot.com/-OlrV-PAtEkQ/W3RkOJCBkaI/AAAAAAAAADOg/gNZXo_eK3tMNOmIfsuyPzrRfNb3qEQwIwCLcB/GAs/s640/image1.gif)



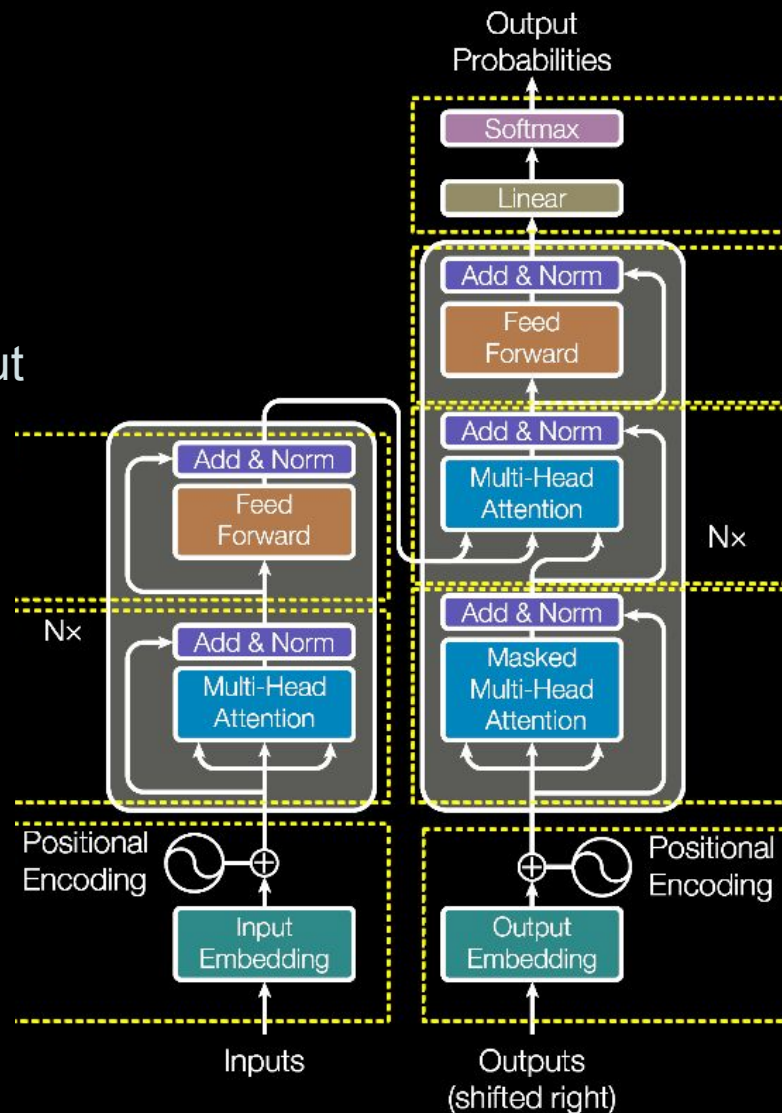
# Transformer

## Why?

- Don't need complexity of LSTM/GRU cells
- Constant num edges between words (or input steps)
- Enables “interactions” (i.e. adaptations) between words
- Easy to parallelize -- don't need sequential processing.

## Drawbacks:

- Only unidirectional by default
- Only a “single-hop” relationship per layer (multiple layers to capture multiple)



# BERT

## Bidirectional Encoder Representations from Transformers

Produces contextualized embeddings  
(or pre-trained contextualized encoder)

### **Drawbacks of Vanilla Transformers:**

- Only unidirectional by default
- Only a “single-hop” relationship per layer  
(multiple layers to capture multiple)

# BERT

## Bidirectional Encoder Representations from Transformers

Produces contextualized embeddings  
(or pre-trained contextualized encoder)

- Bidirectional context by “masking” in the middle
- A lot of layers, hidden states, attention heads.

### **Drawbacks of Vanilla Transformers:**

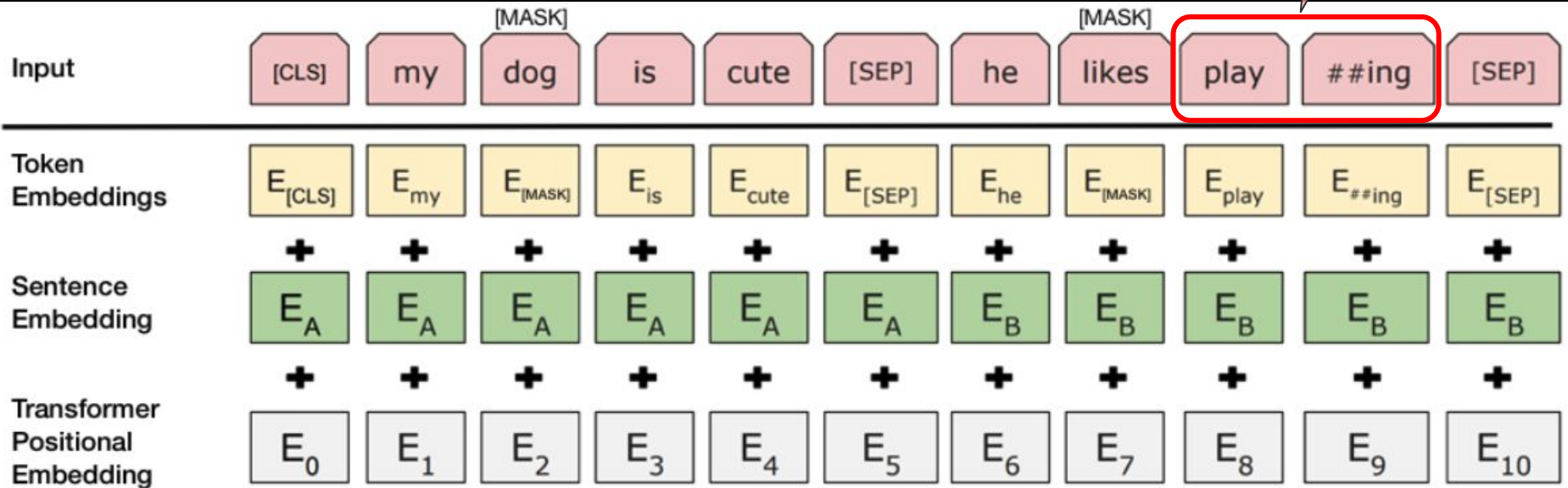
- Only unidirectional by default
- Only a “single-hop” relationship per layer  
(multiple layers to capture multiple)

# BERT

**Sentence A** = The man went to the store.  
**Sentence B** = He bought a gallon of milk.  
**Label** = IsNextSentence

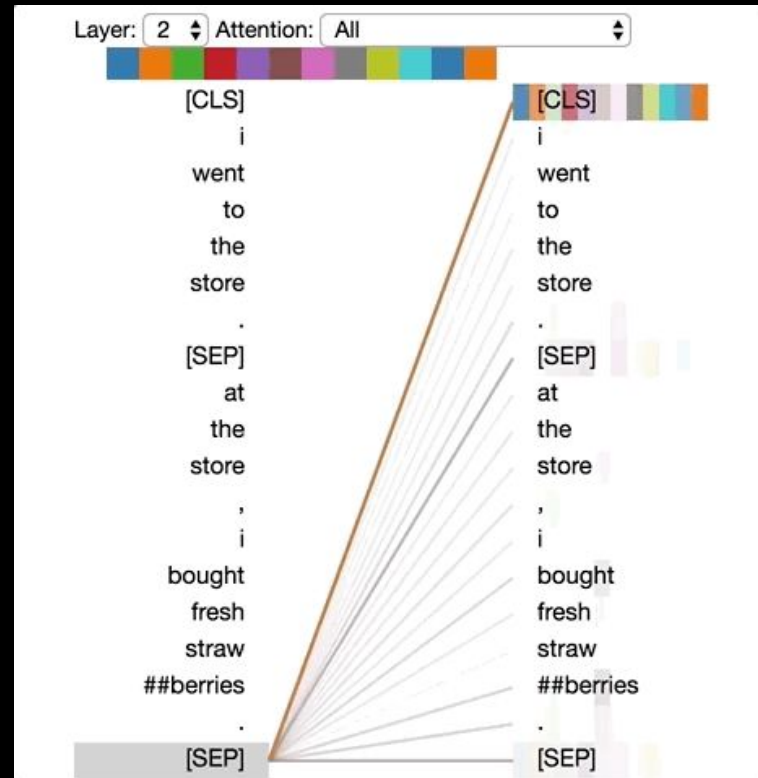
**Sentence A** = The man went to the store.  
**Sentence B** = Penguins are flightless.  
**Label** = NotNextSentence

tokenize into "word pieces"



# Bert: Attention by Layers

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

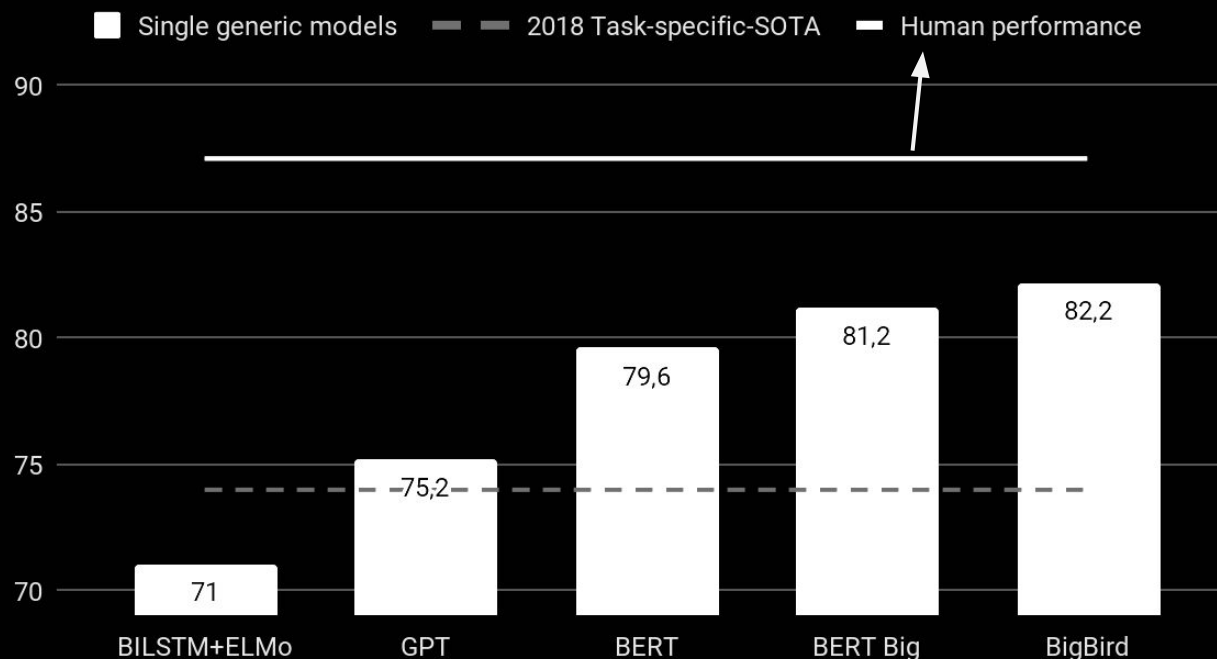


(Vig, 2019)



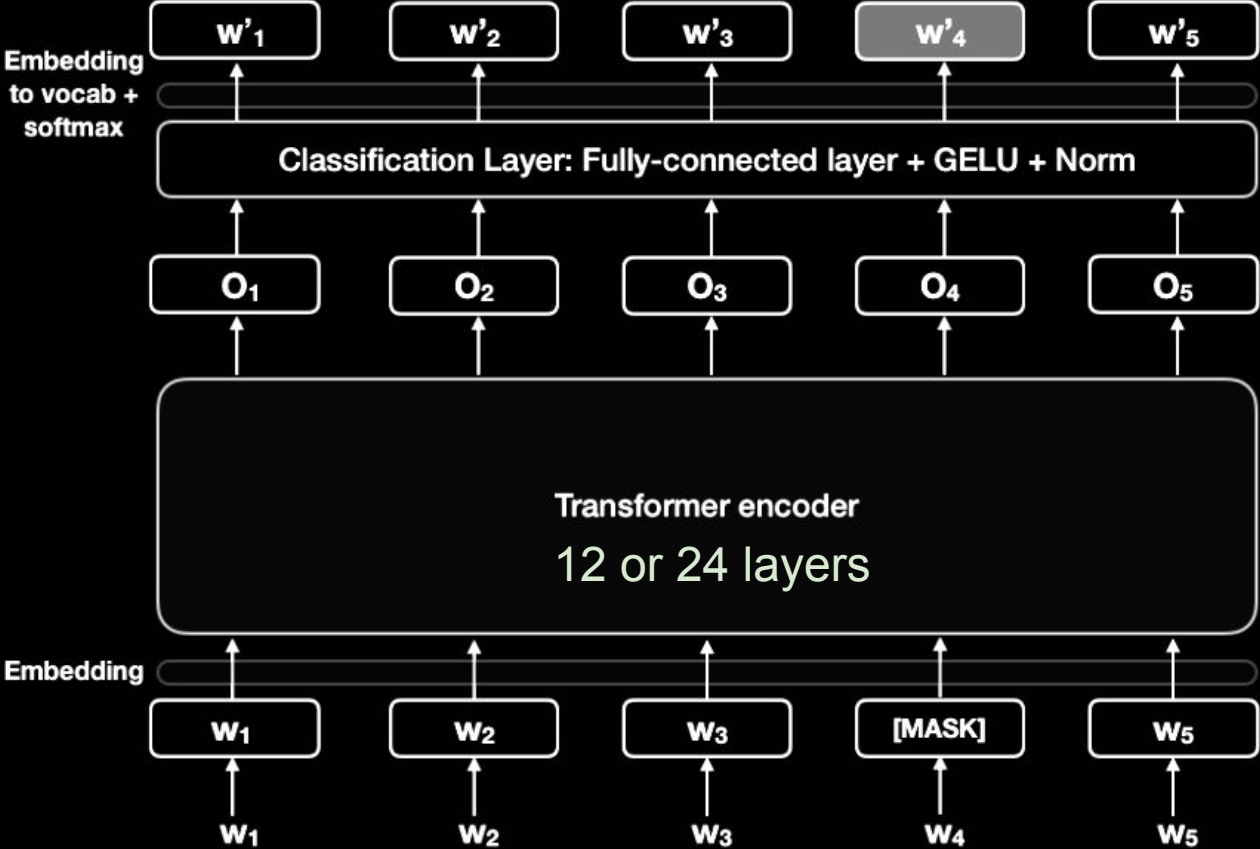
# BERT Performance: e.g. Question Answering

GLUE scores evolution over 2018-2019

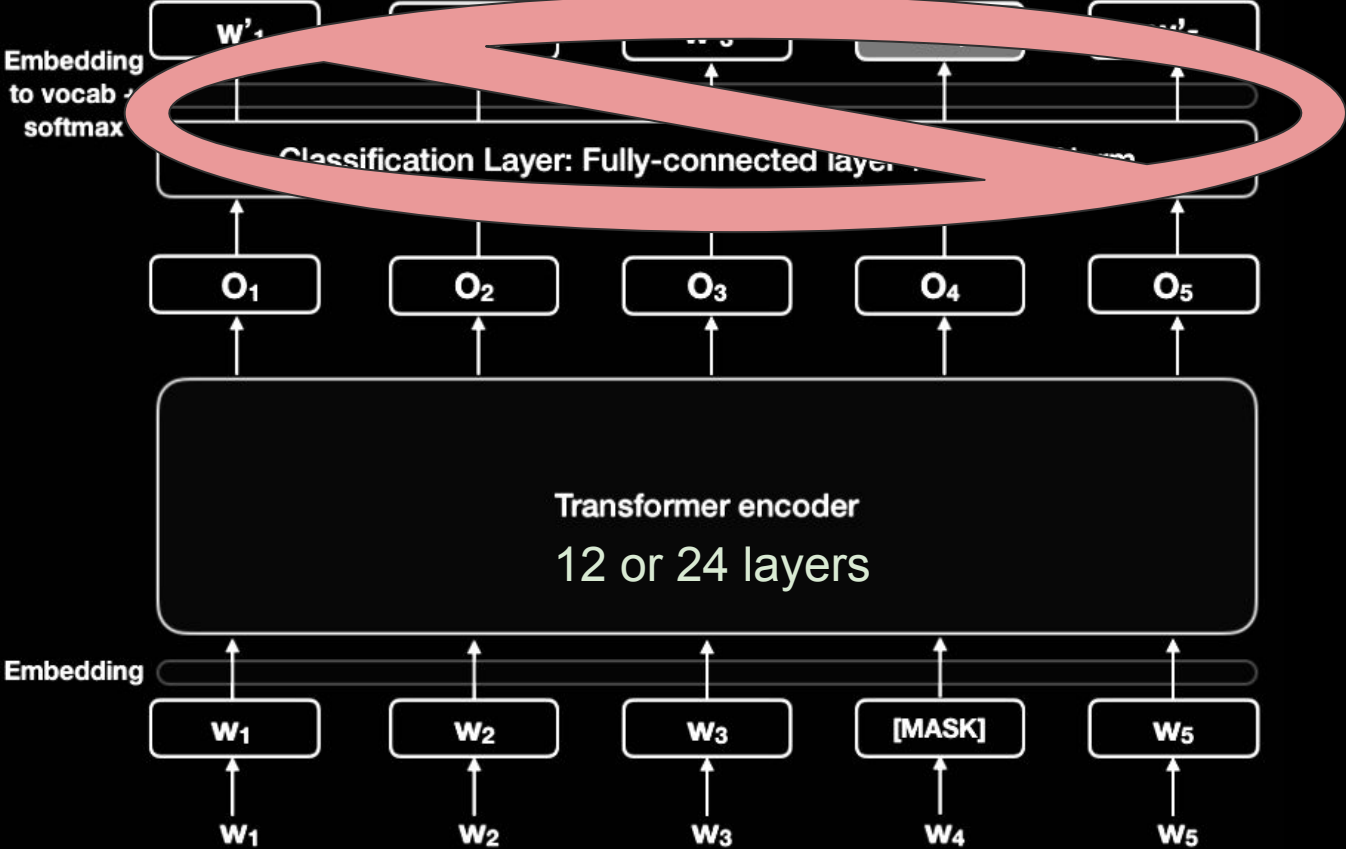


<https://rajpurkar.github.io/SQuAD-explorer/>

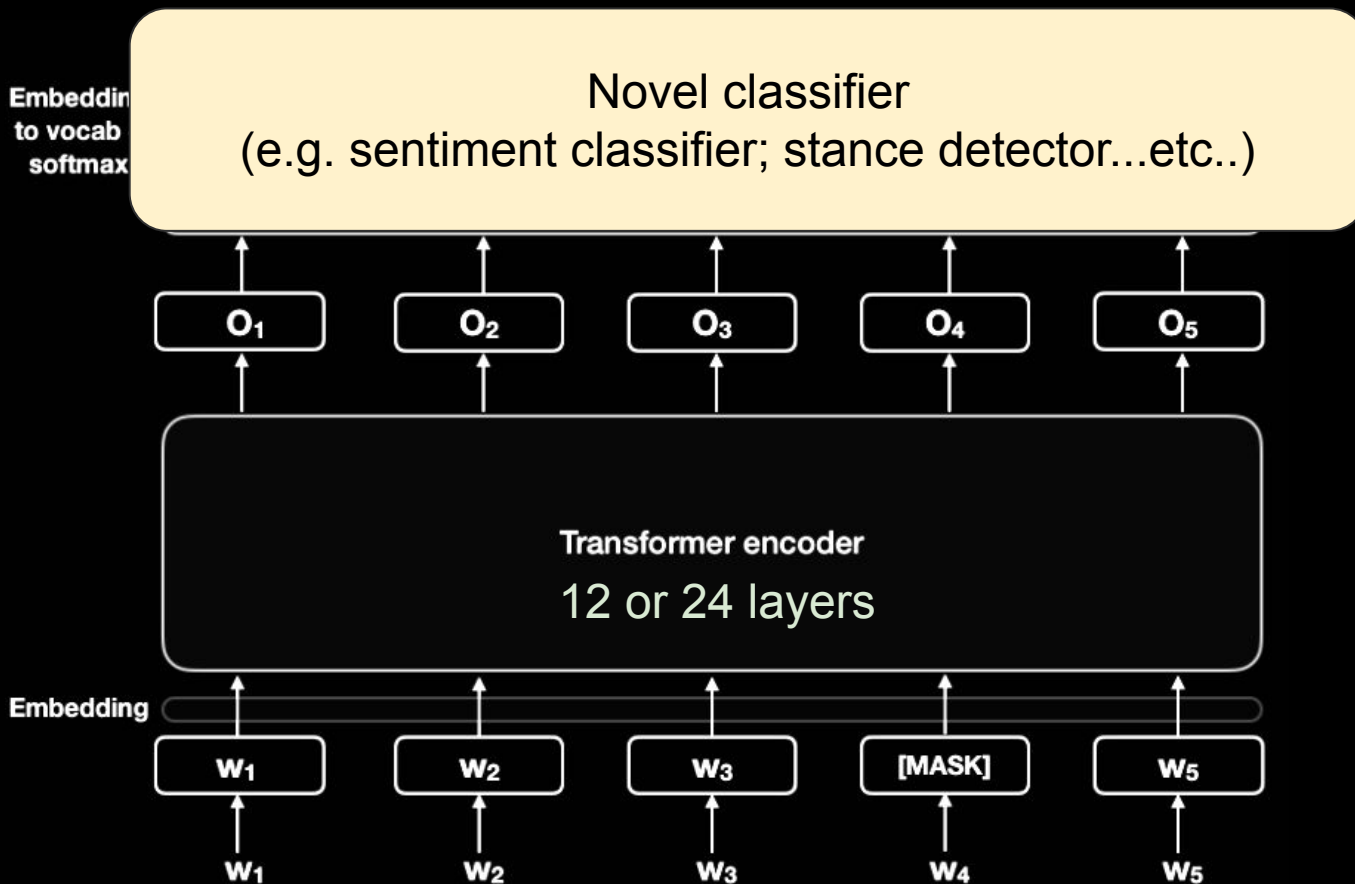
# BERT: Pre-training; Fine-tuning



# BERT: Pre-training; Fine-tuning



# BERT: Pre-training; Fine-tuning



# The Transformer: “Attention-only” models

**Can handle sequences and long-distance dependencies, but....**

- Don't want complexity of LSTM/GRU cells
- Constant num edges between input steps
- Enables “interactions” (i.e. adaptations) between words
- **Easy to parallelize -- don't need sequential processing.**