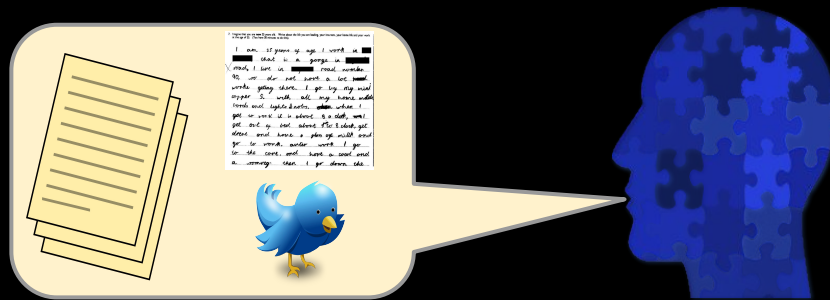


Attention and Transformer Sequence Models

CSE354 - Spring 2021
Natural Language Processing



Most NLP Tasks. E.g.

- Sequence Tasks
 - Language Modeling
 - Machine Translation
 - Speech Recognition
 - Named Entity
- Document Classification

• Transformer Networks

- Transformers
- BERT

Evolution of Sequence Modeling

RNNs

LSTMs

LSTMs with Attention

Attention without LSTMs

Evolution of Sequence Modeling

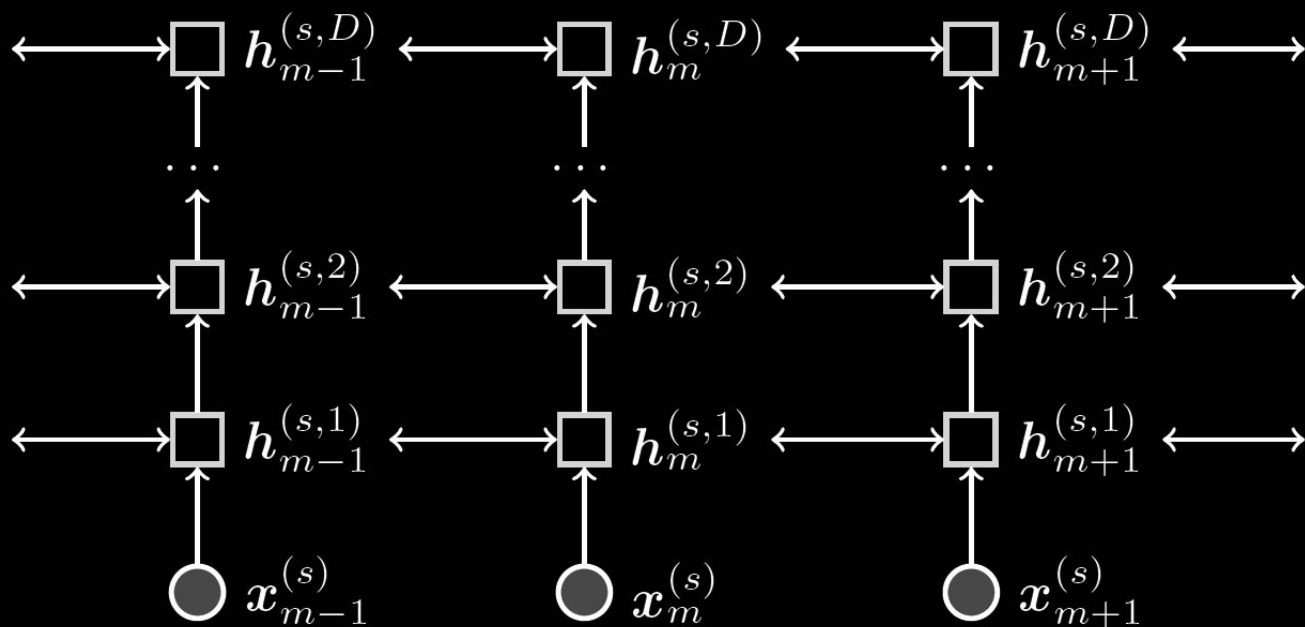
RNNs

LSTMs

LSTMs with Attention

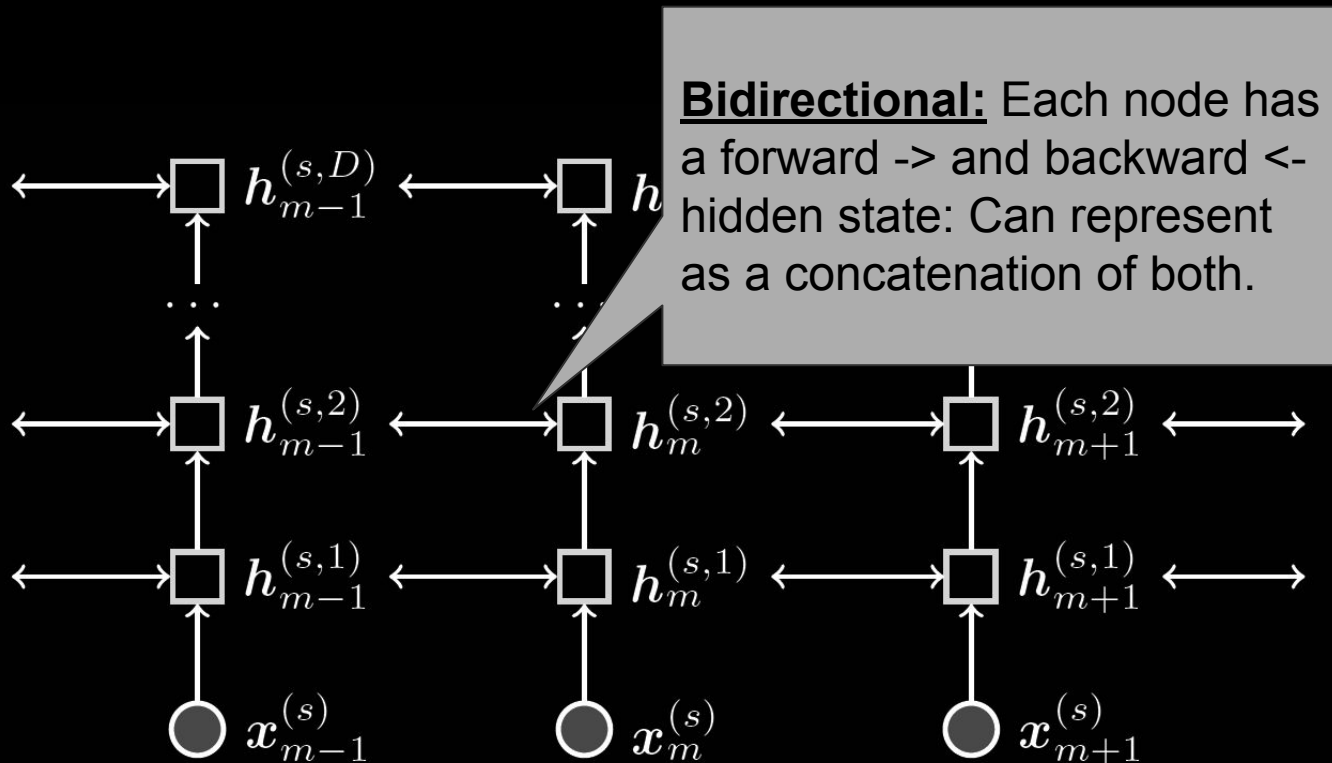
Attention without LSTMs

Multi-level bidirectional RNN (LSTM or GRU)



(Eisenstein, 2018)

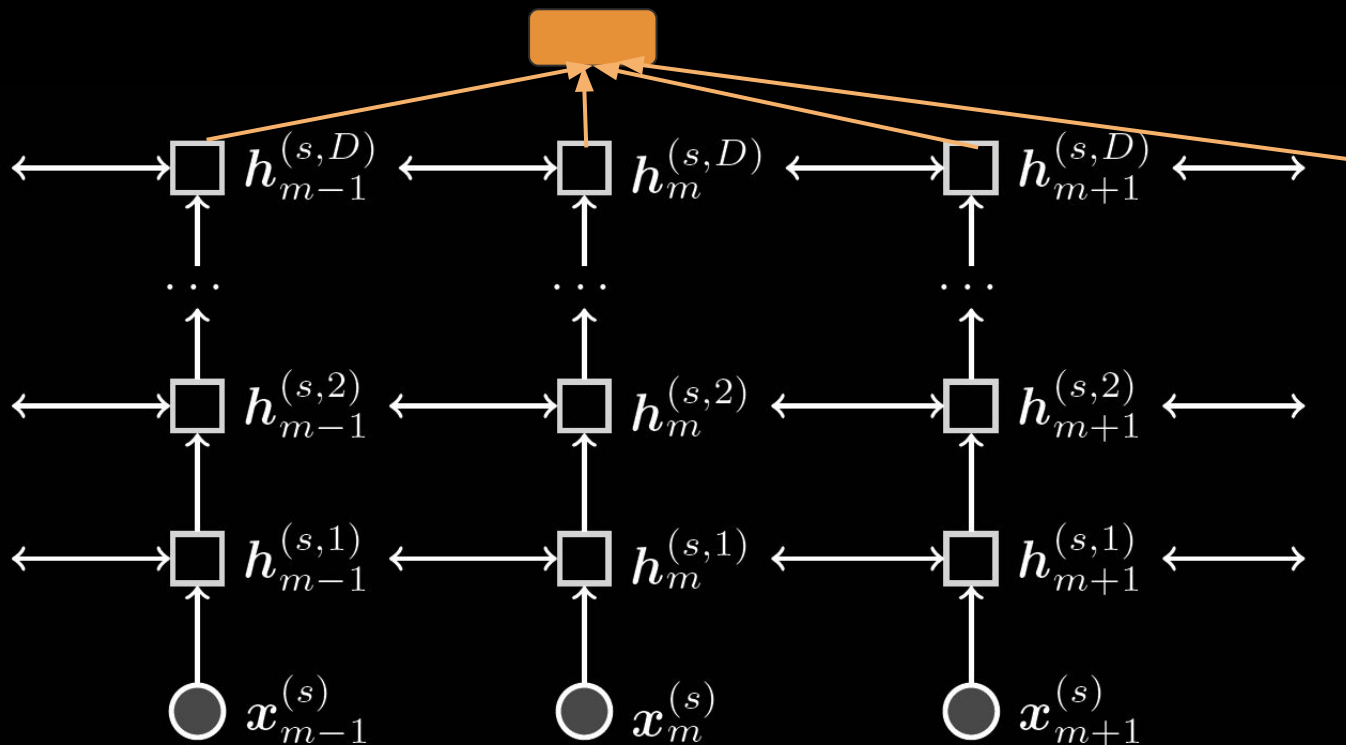
Multi-level bidirectional RNN (LSTM or GRU)



(Eisenstein, 2018)

Multi-level bidirectional RNN (LSTM or GRU)

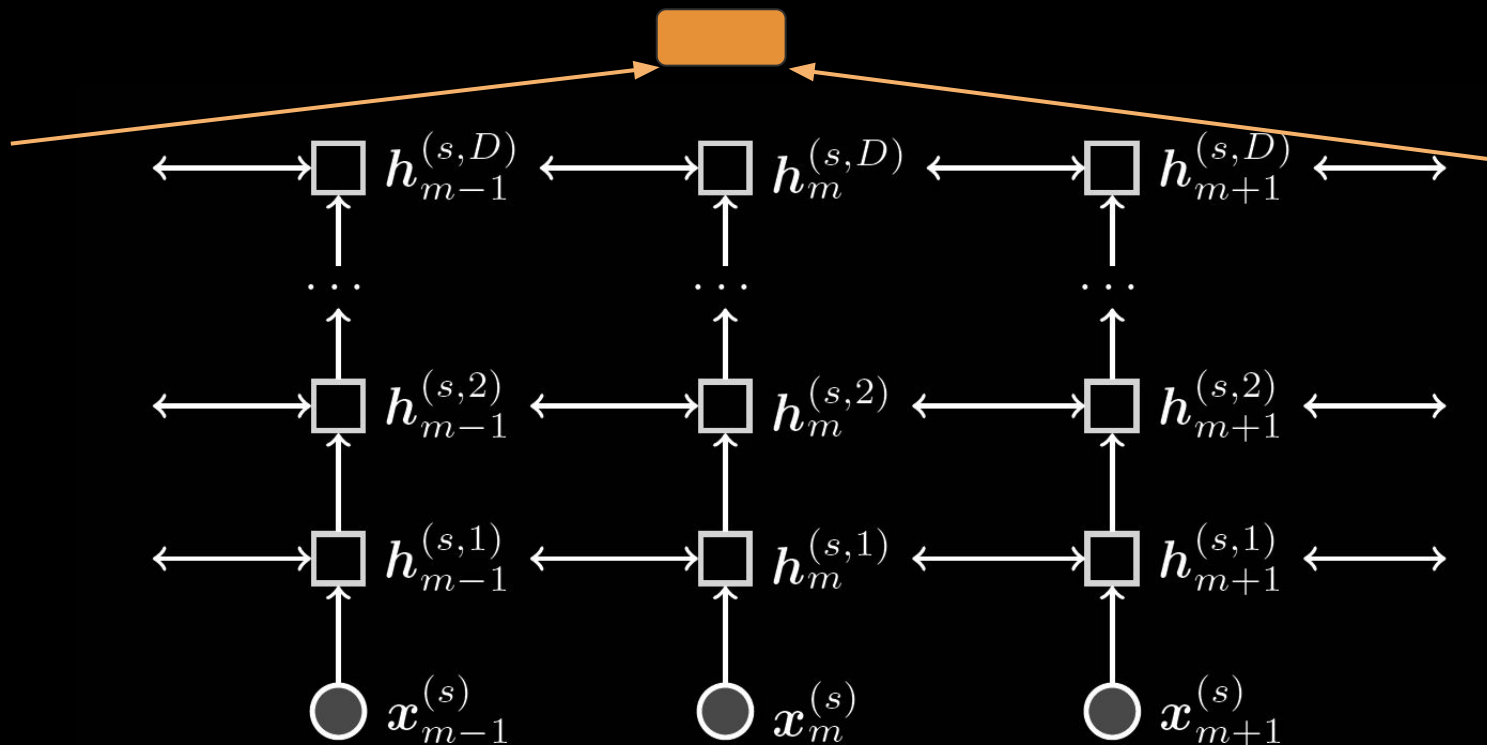
Average of top layer is an embedding (average of concatenated vectors)



(Eisenstein, 2018)

Multi-level bidirectional RNN (LSTM or GRU)

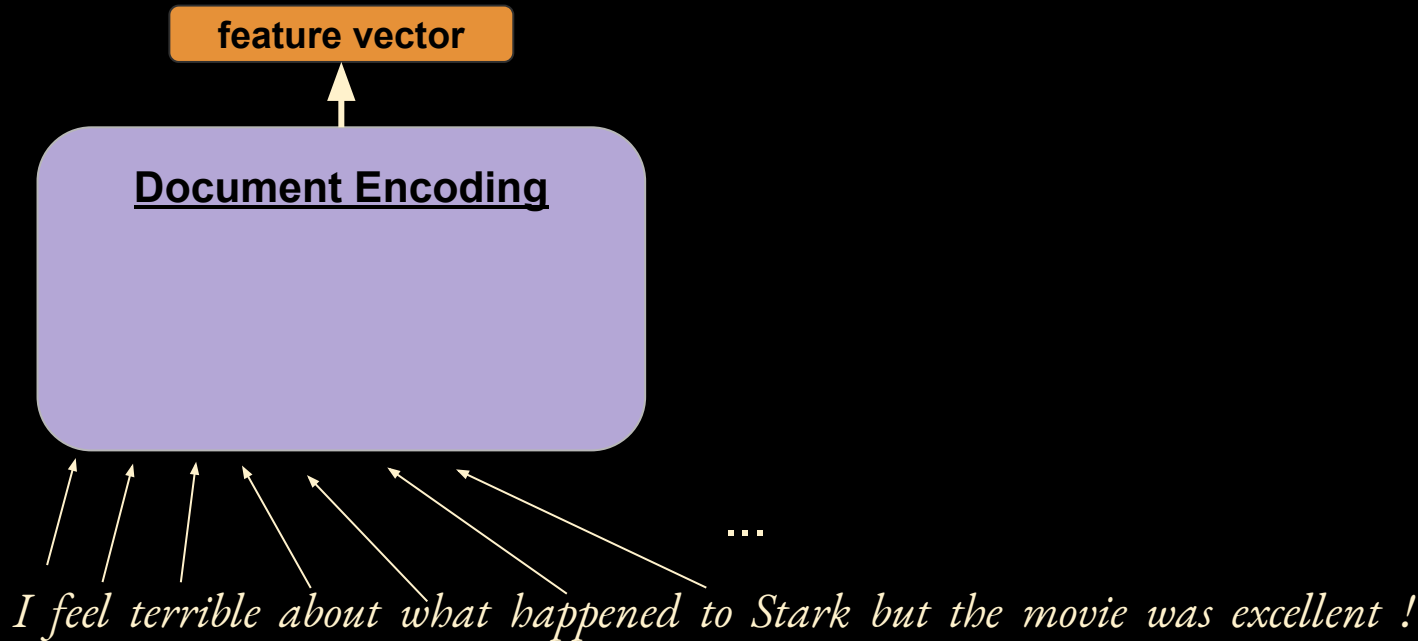
Sometimes just use left-most and right-most hidden state instead



(Eisenstein, 2018)

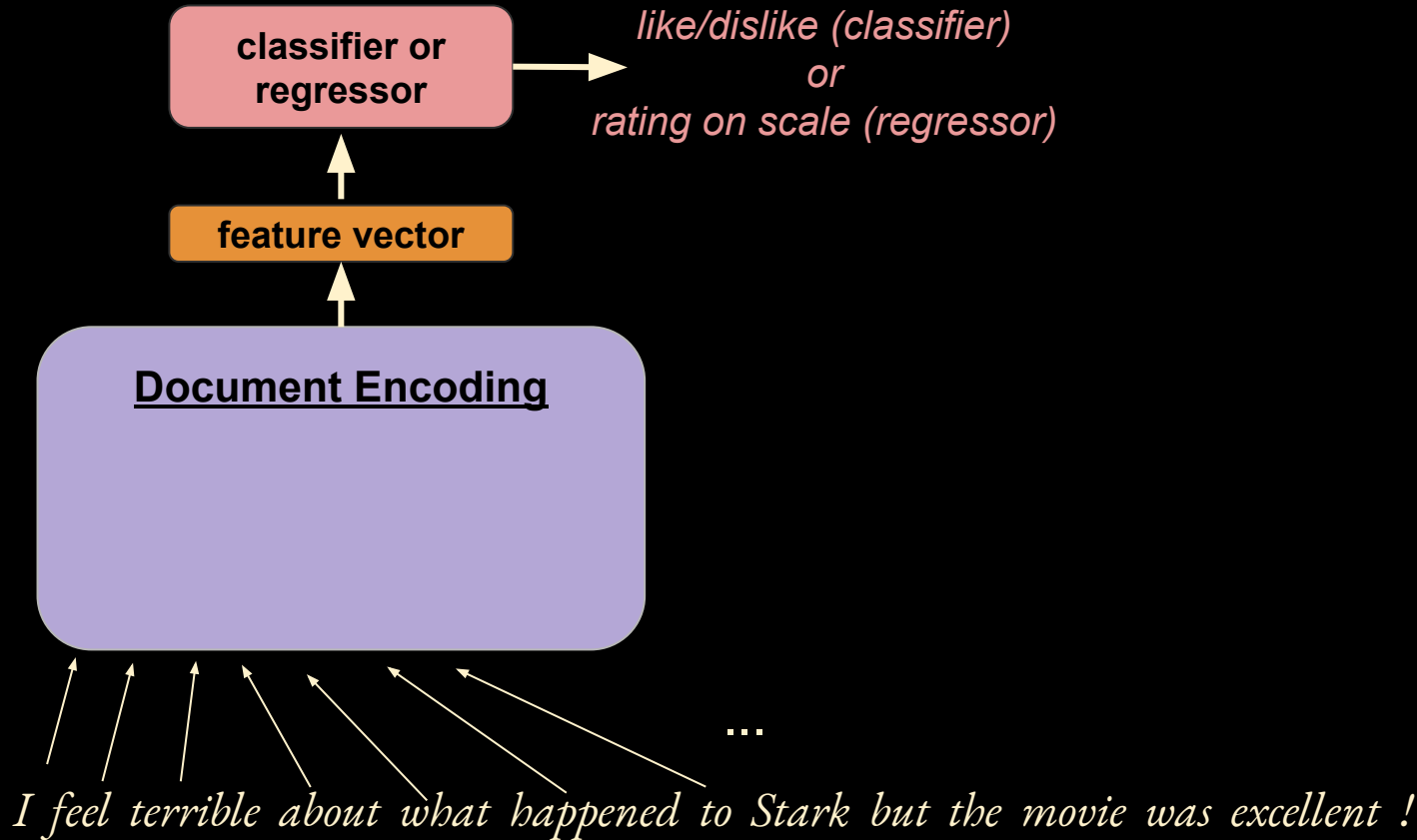
Sentiment Analysis:

Example Application of Single Representation of document



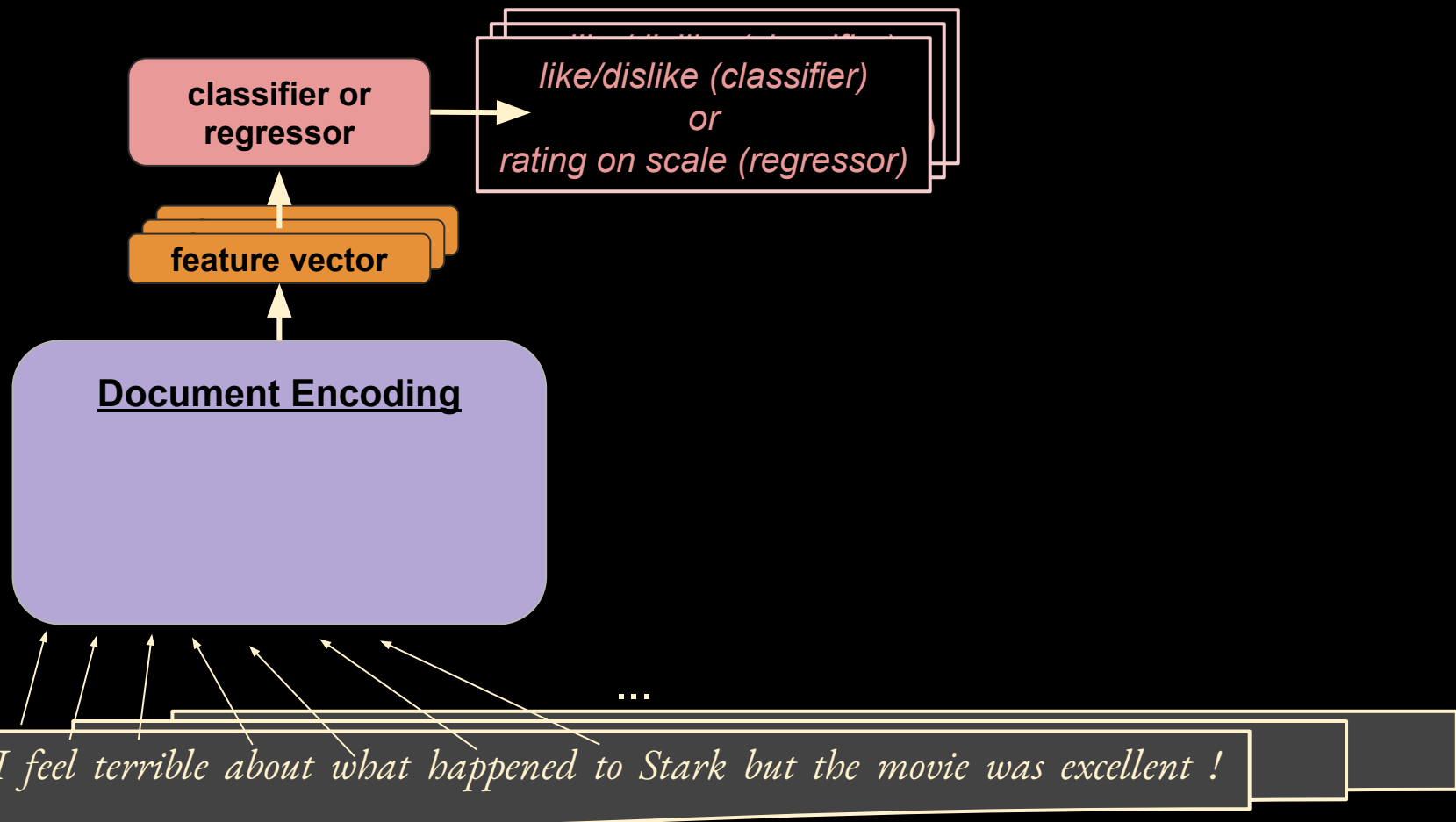
Sentiment Analysis:

Example Application of Single Representation of document



Sentiment Analysis:

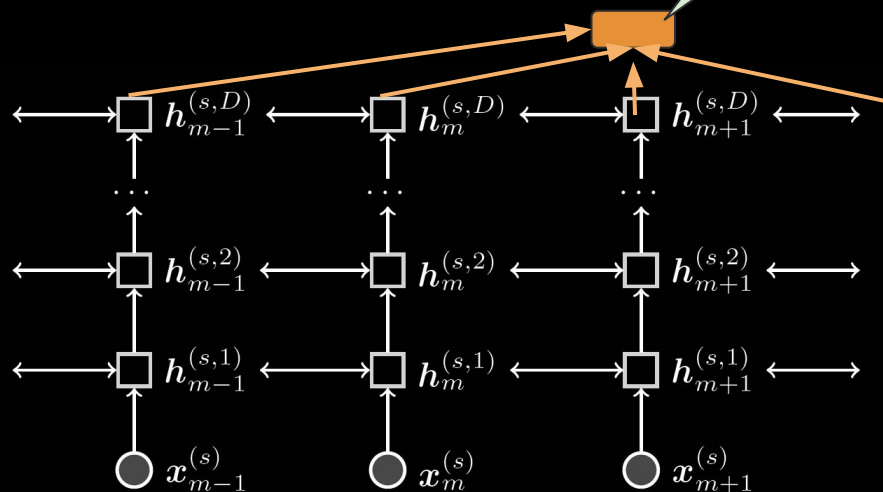
Example Application of Single Representation of document



Encoder

A representation of input.

Document classification only needs an **encoder**. Goes from input into a single representation.

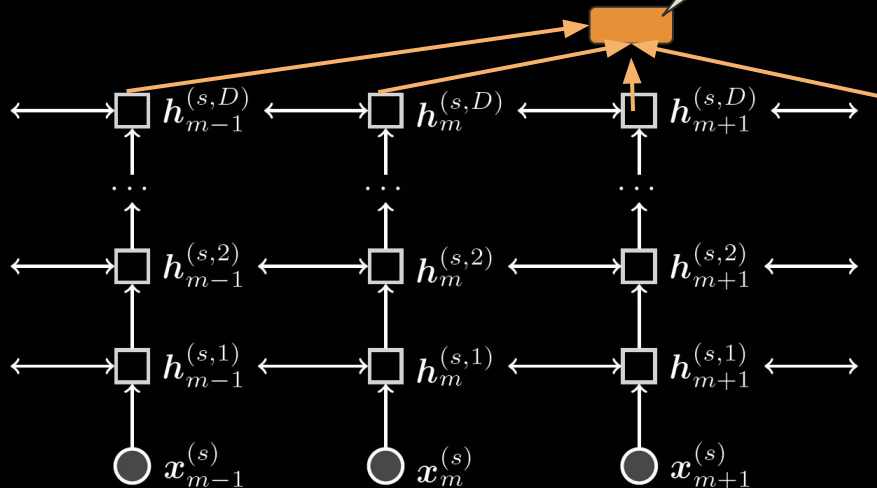


(Eisenstein, 2018)

Encoder

A representation of input.

Document classification only needs an **encoder**. Goes from input into a single representation.



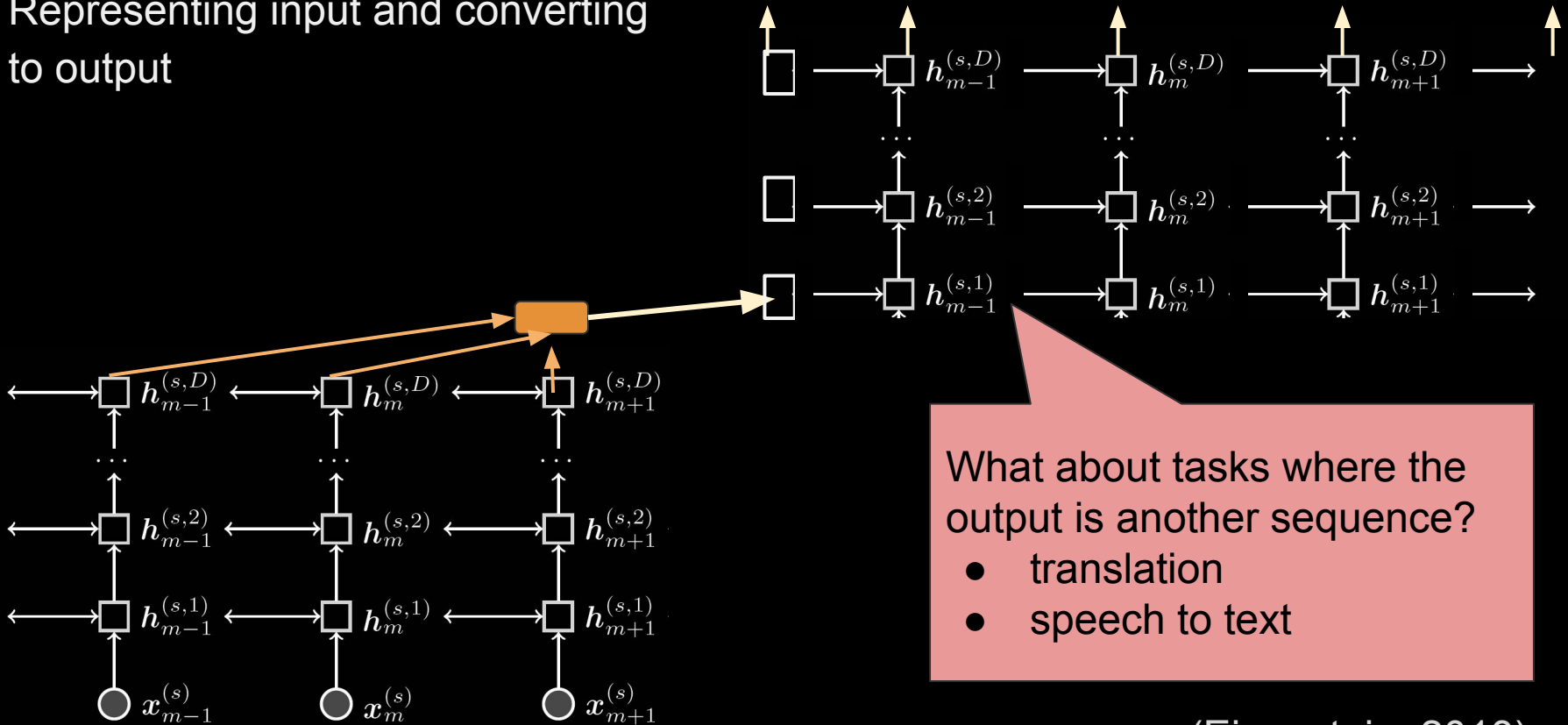
What about tasks where the output is another sequence?

- translation
- speech to text

(Eisenstein, 2018)

Encoder-Decoder (seq to seq models)

Representing input and converting to output

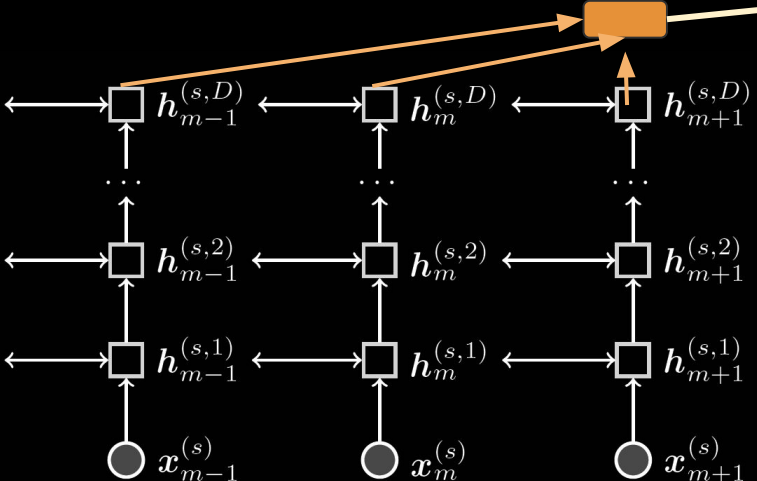
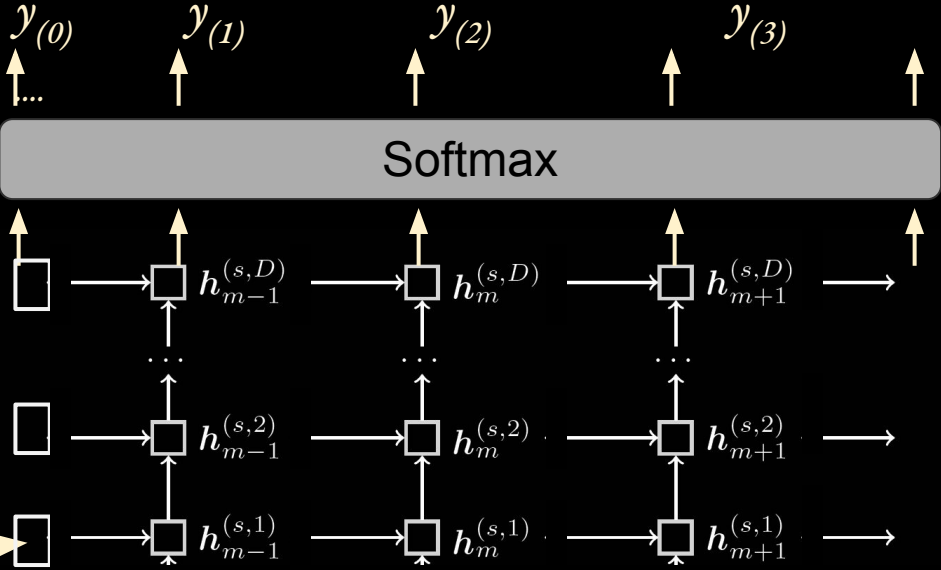


What about tasks where the output is another sequence?

- translation
- speech to text

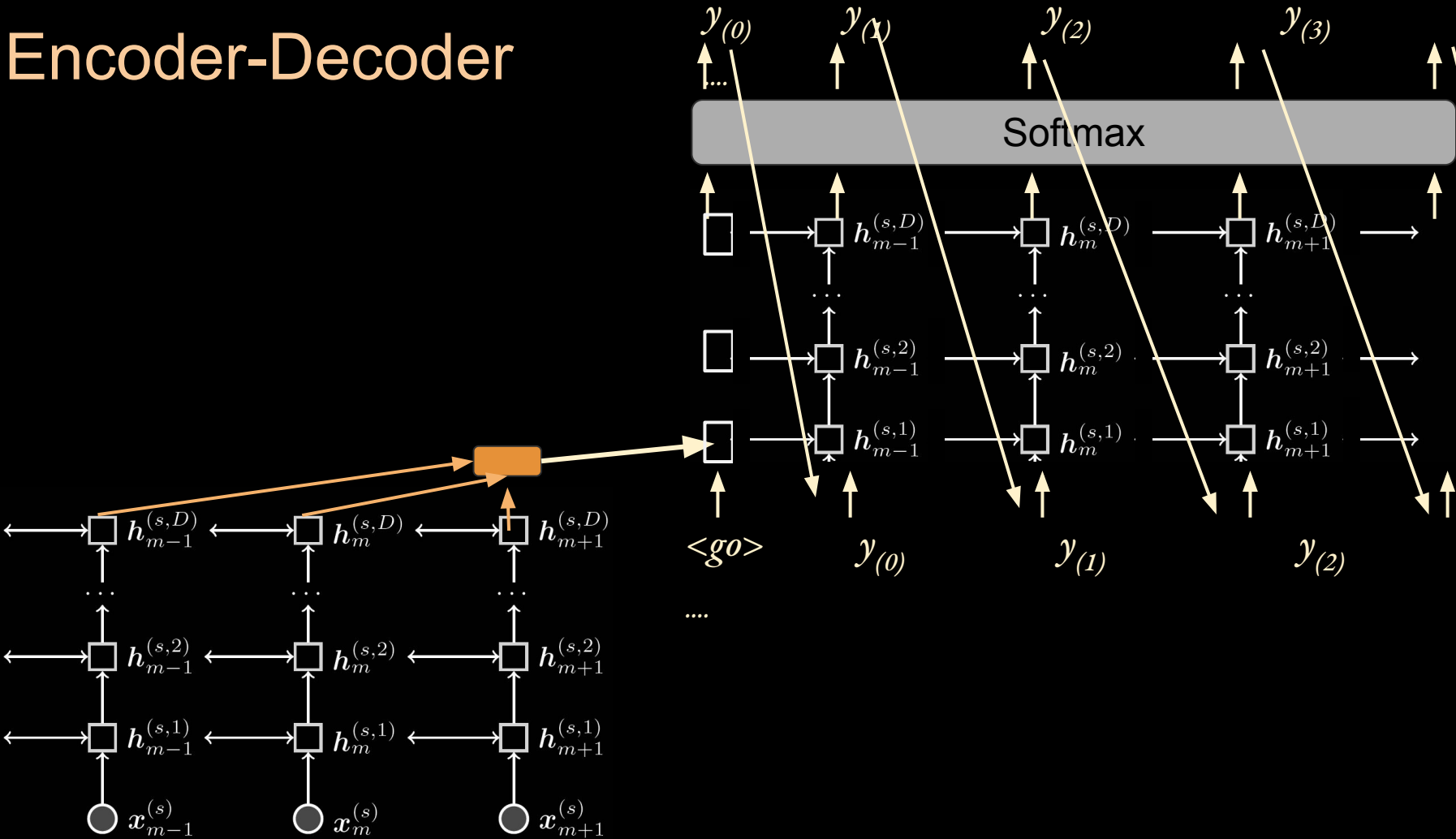
(Eisenstein, 2018)

Encoder-Decoder

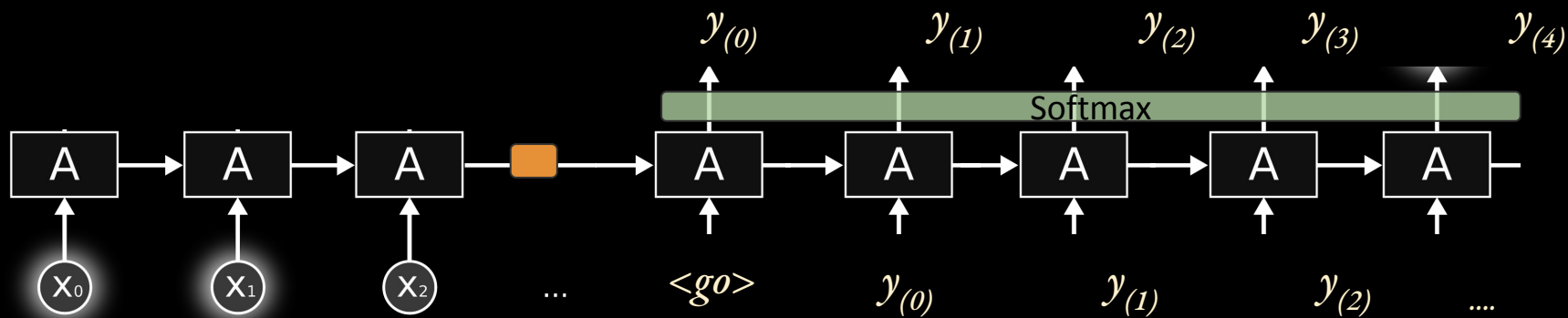


(Eisenstein, 2018)

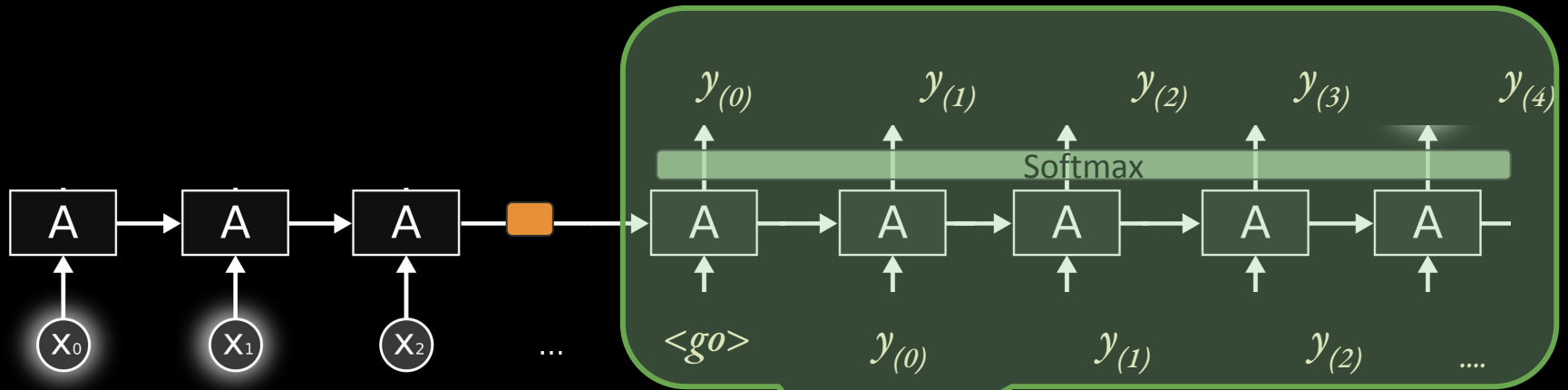
Encoder-Decoder



Encoder-Decoder (Simpler Representation)

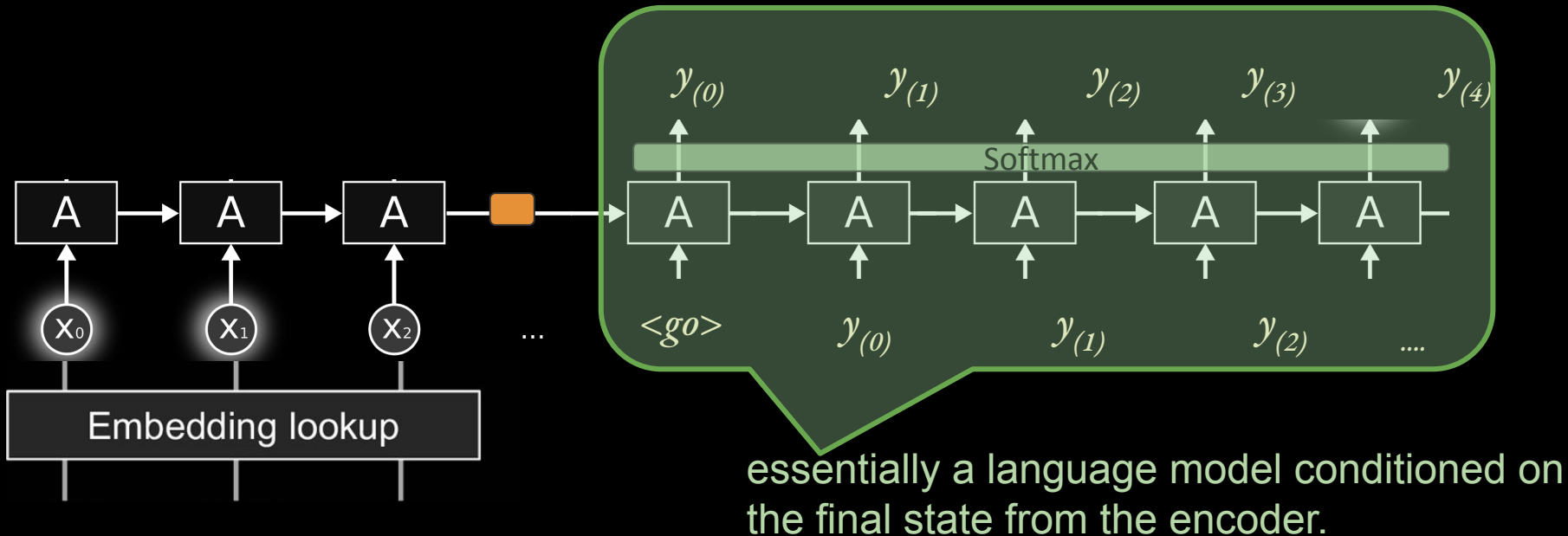


Encoder-Decoder (Simpler Representation)

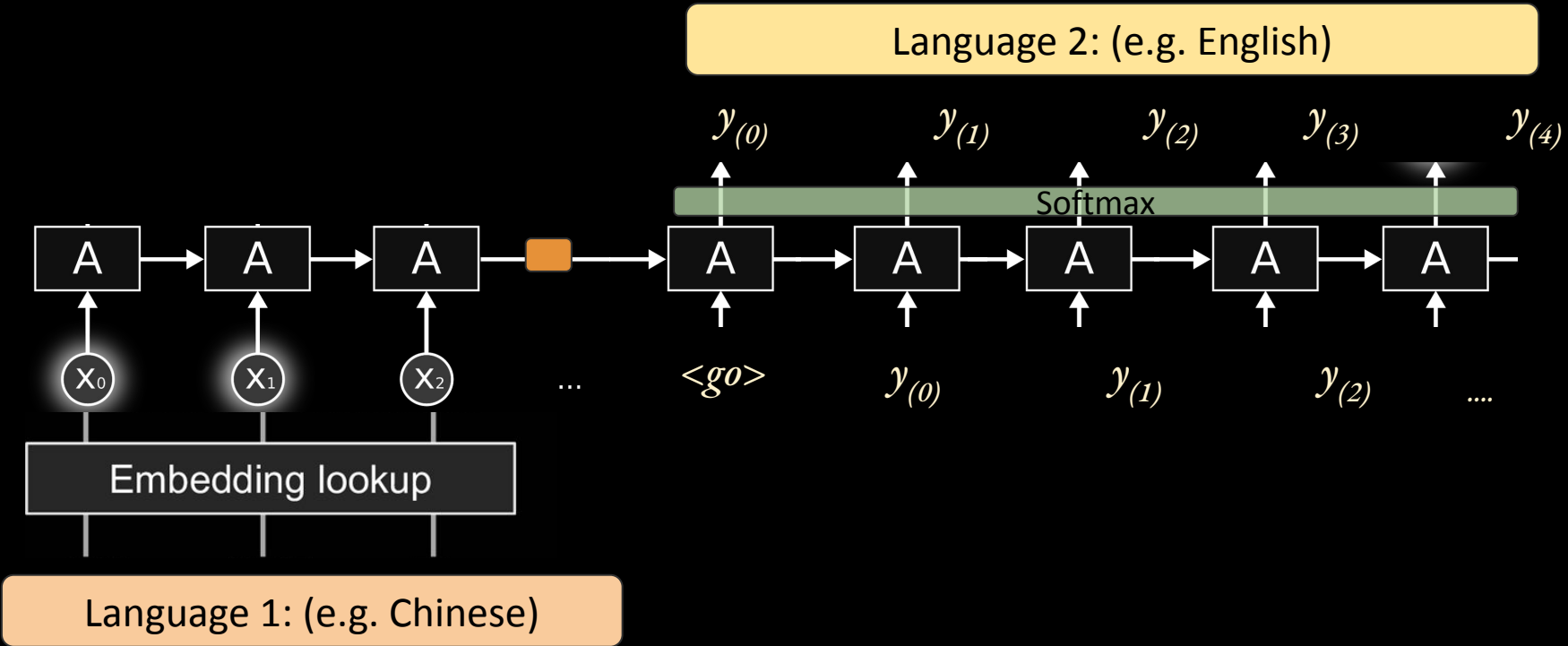


essentially a language model conditioned on the final state from the encoder.

Encoder-Decoder (Simpler Representation)



Encoder-Decoder (Simpler Representation)

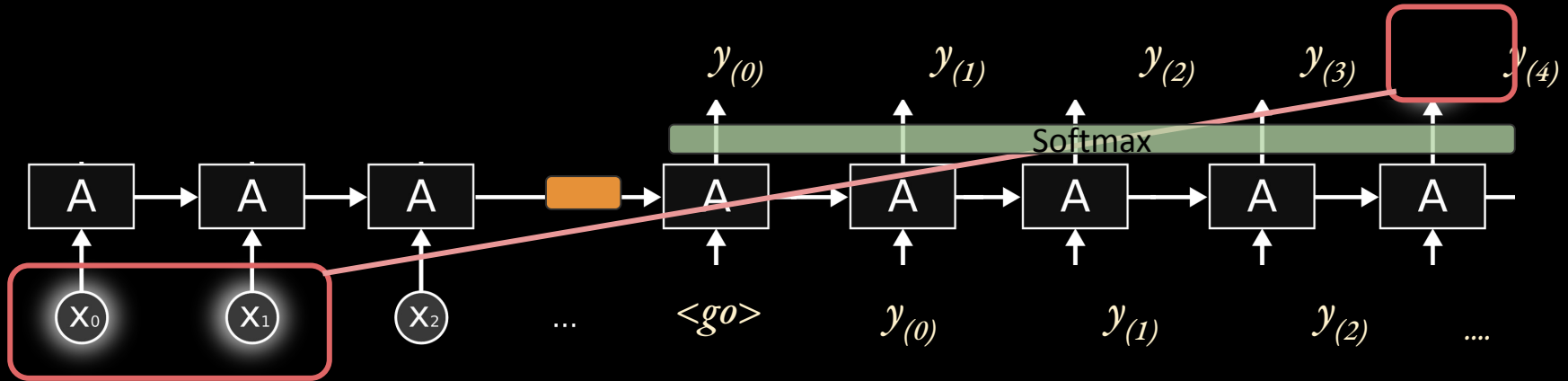


Encoder-Decoder

Challenge:

The ball was kicked by kayla.

- Long distance dependency when translating:



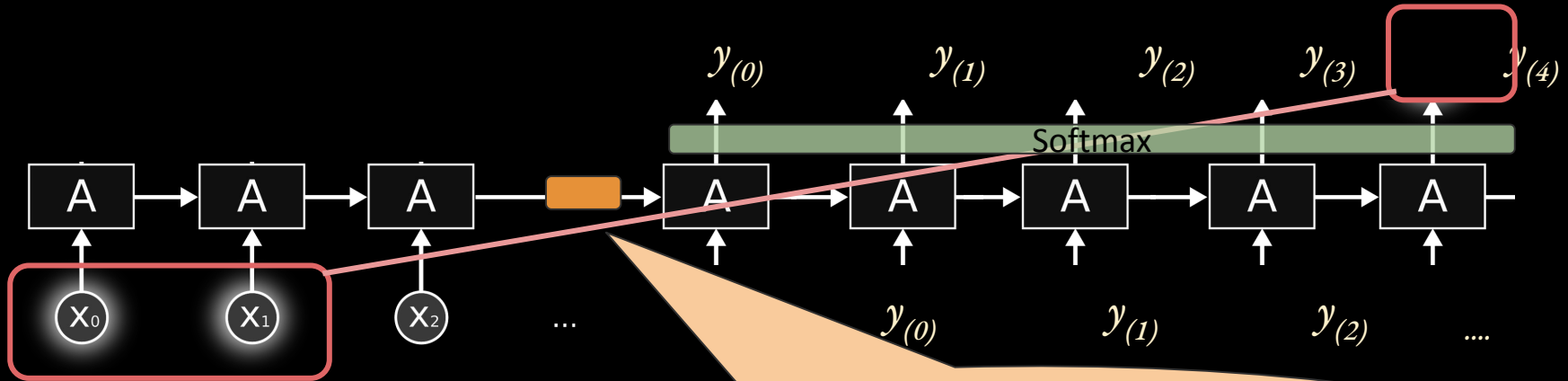
Kayla kicked the ball.

Encoder-Decoder

Challenge:

The ball was kicked by kayla.

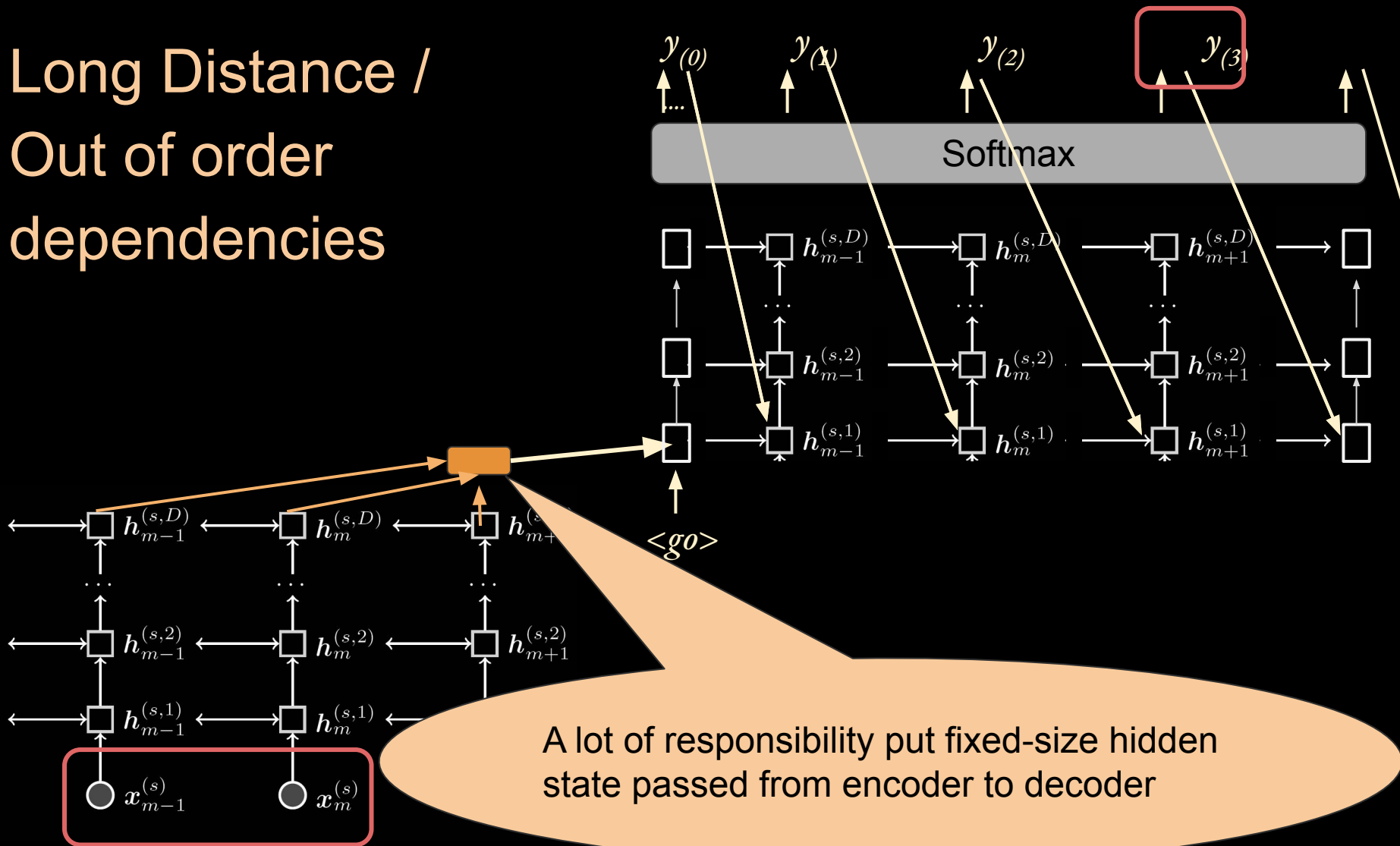
- Long distance dependency when translating:



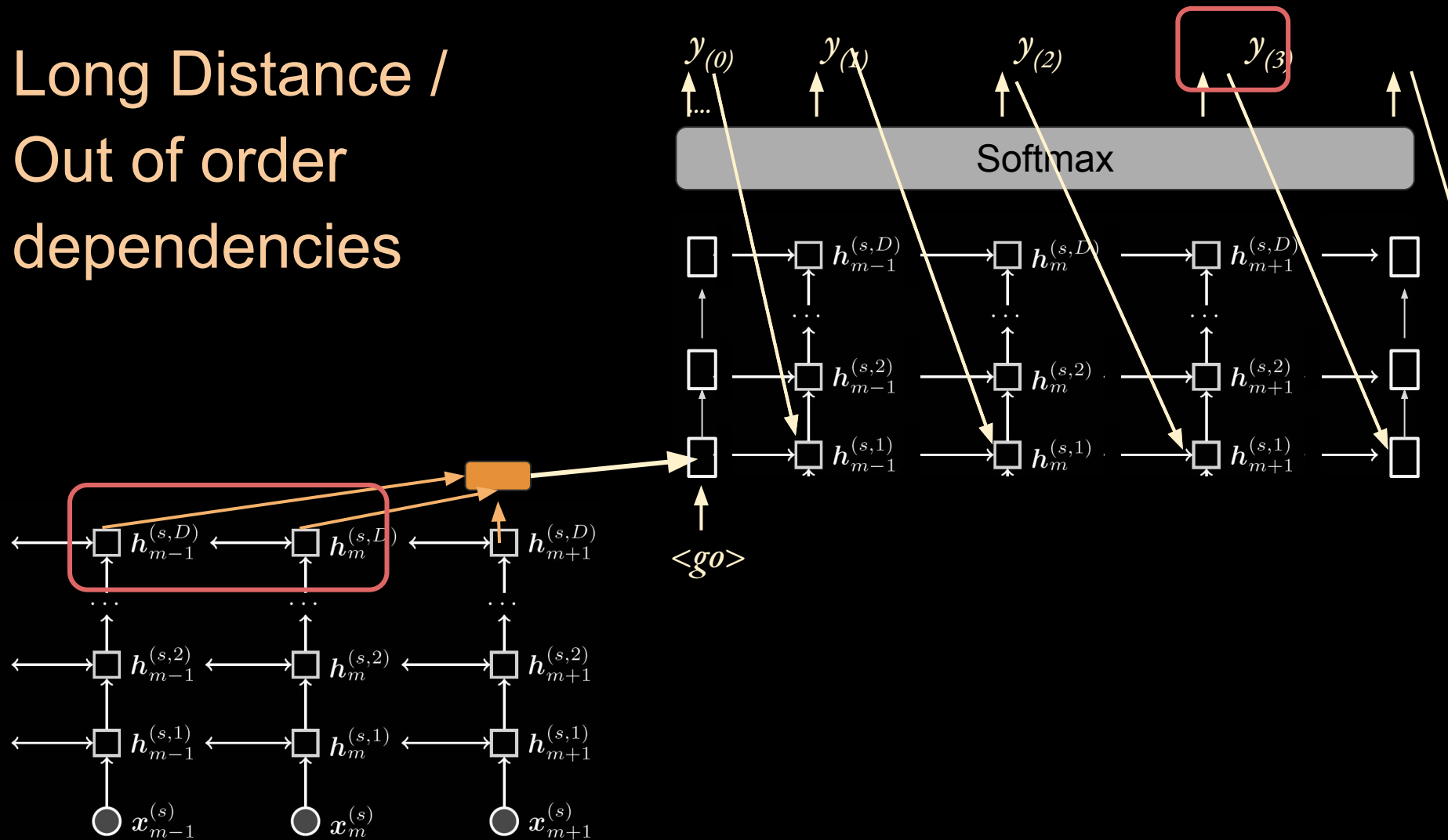
Kayla kicked the ball.

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

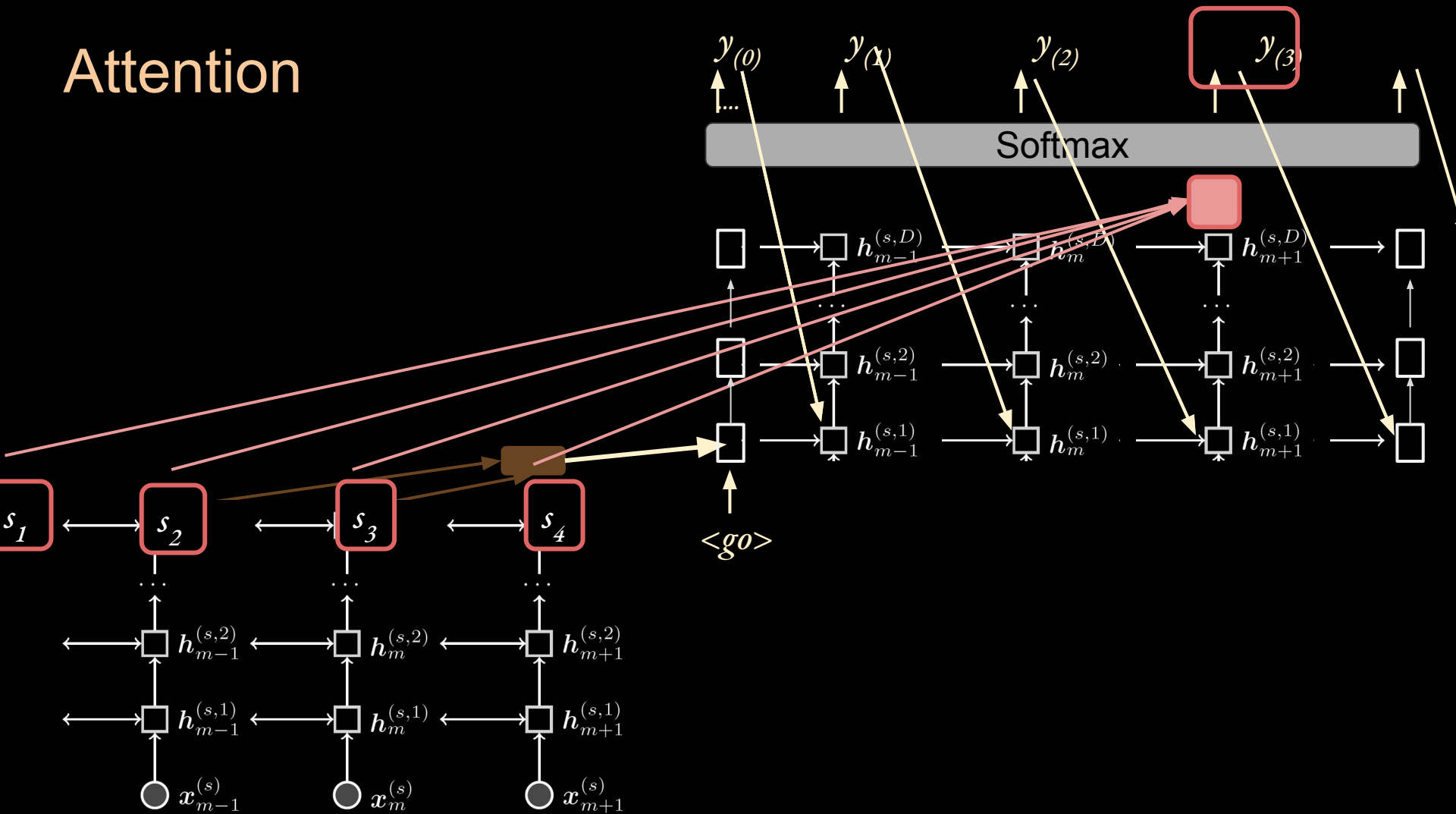
Long Distance /
Out of order
dependencies



Long Distance /
Out of order
dependencies

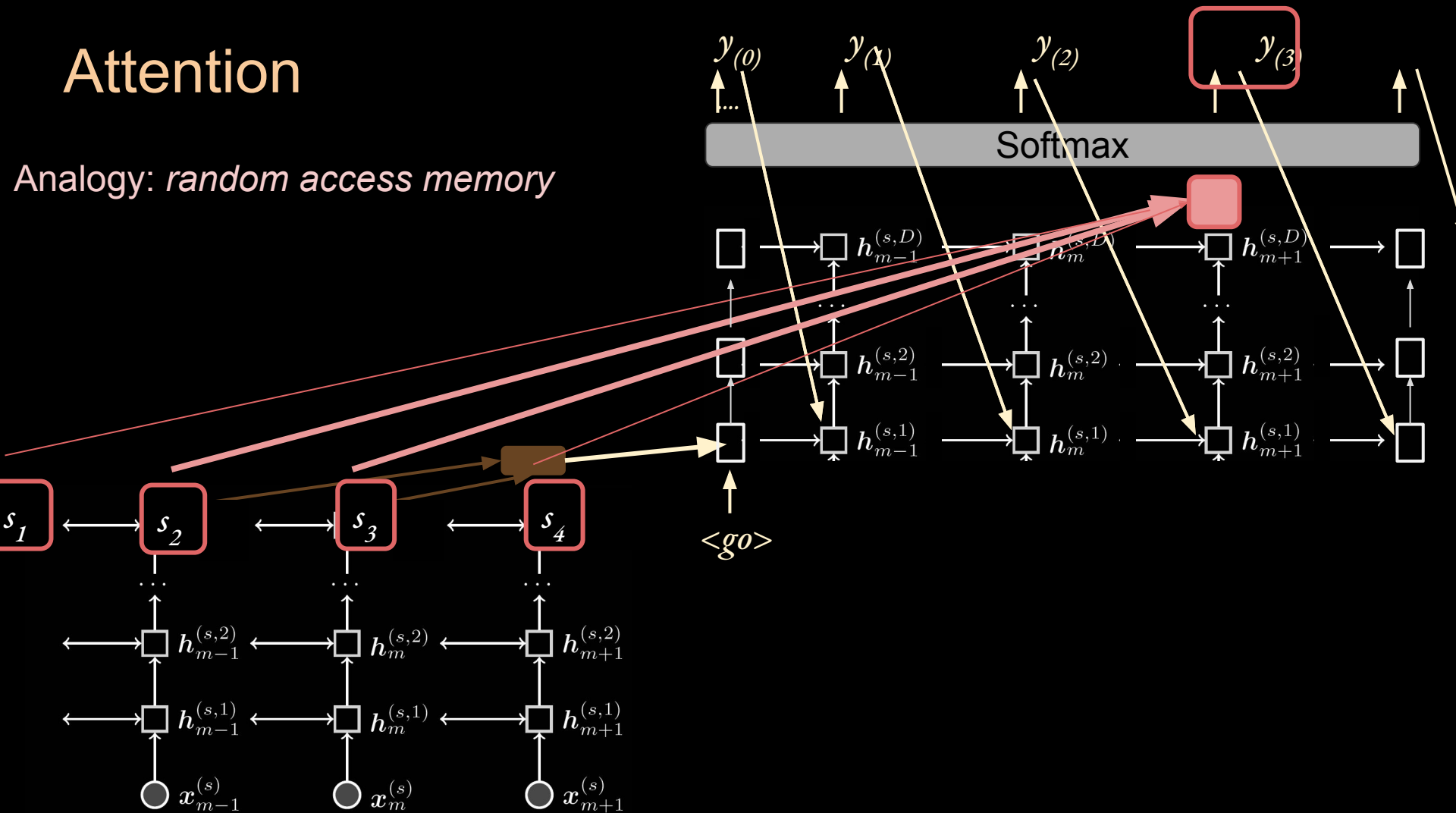


Attention

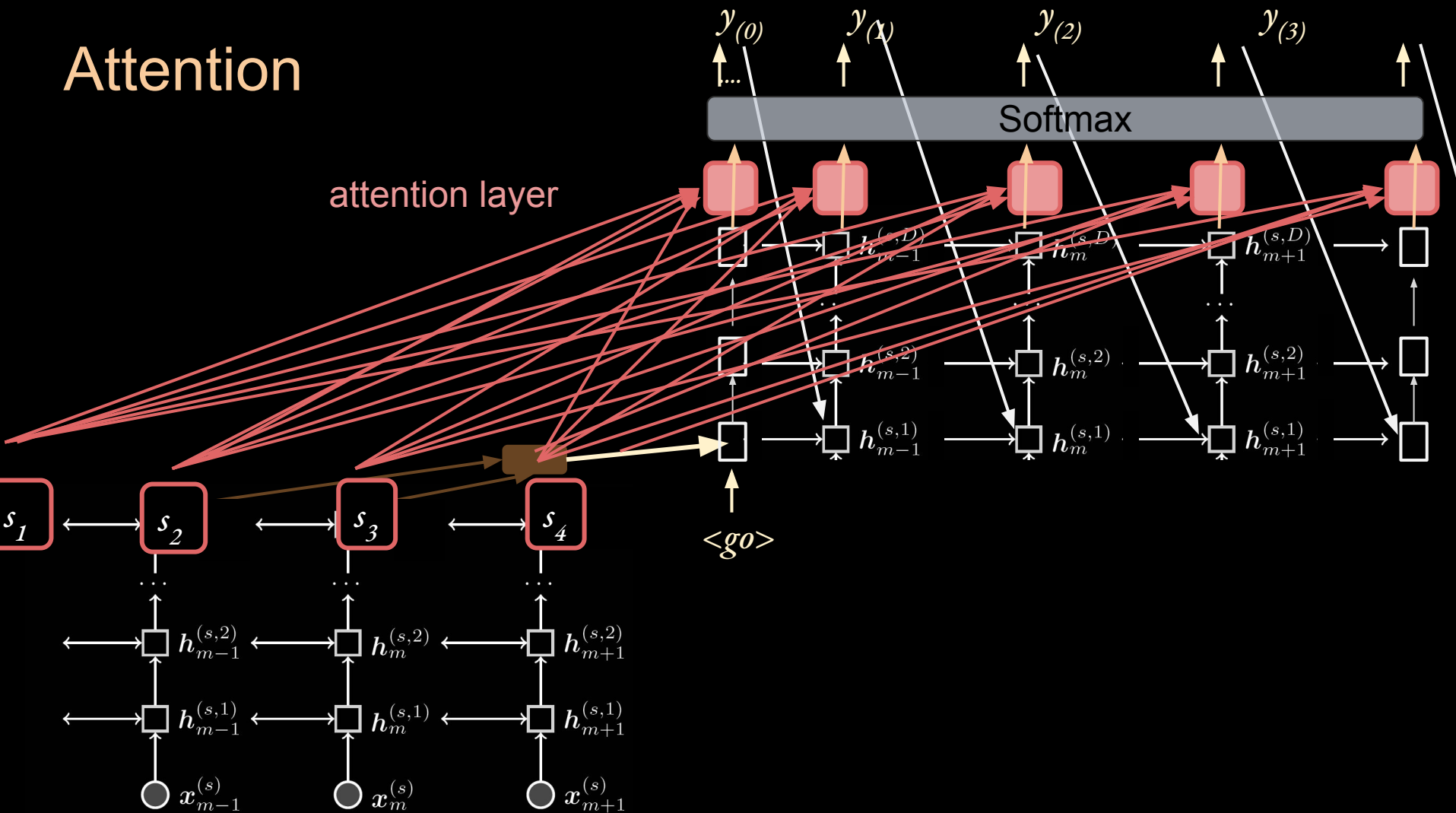


Attention

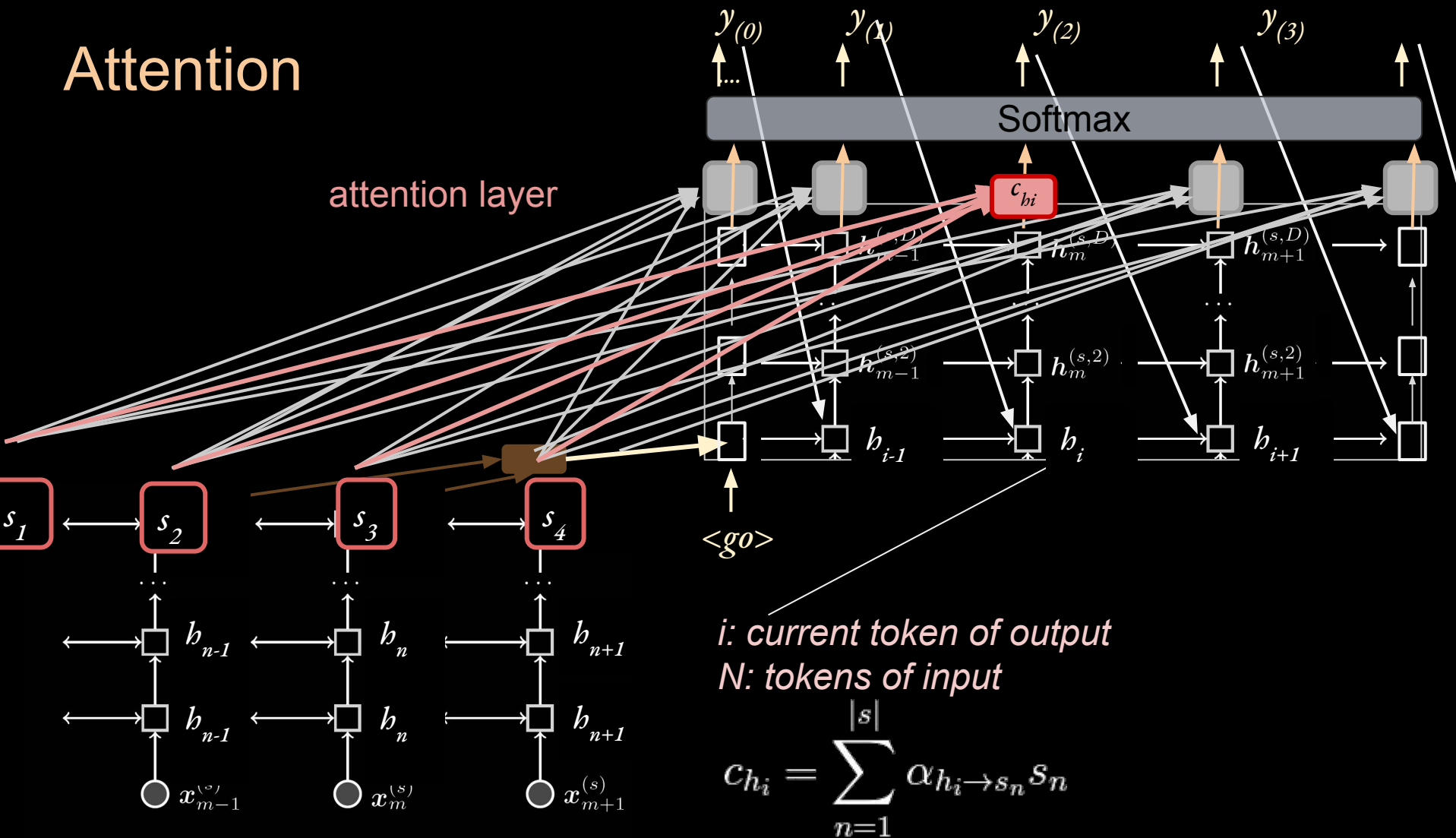
Analogy: *random access memory*



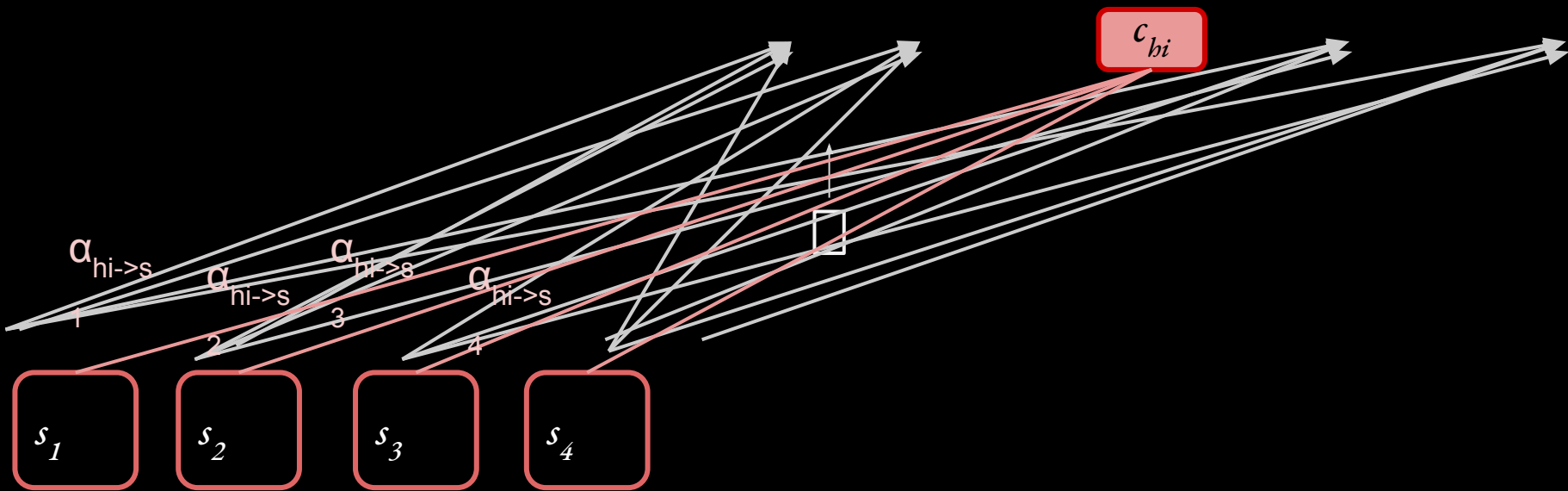
Attention



Attention

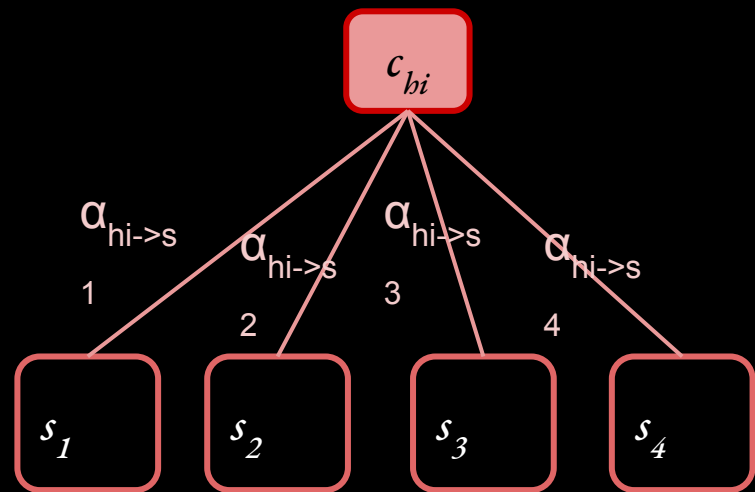


Attention



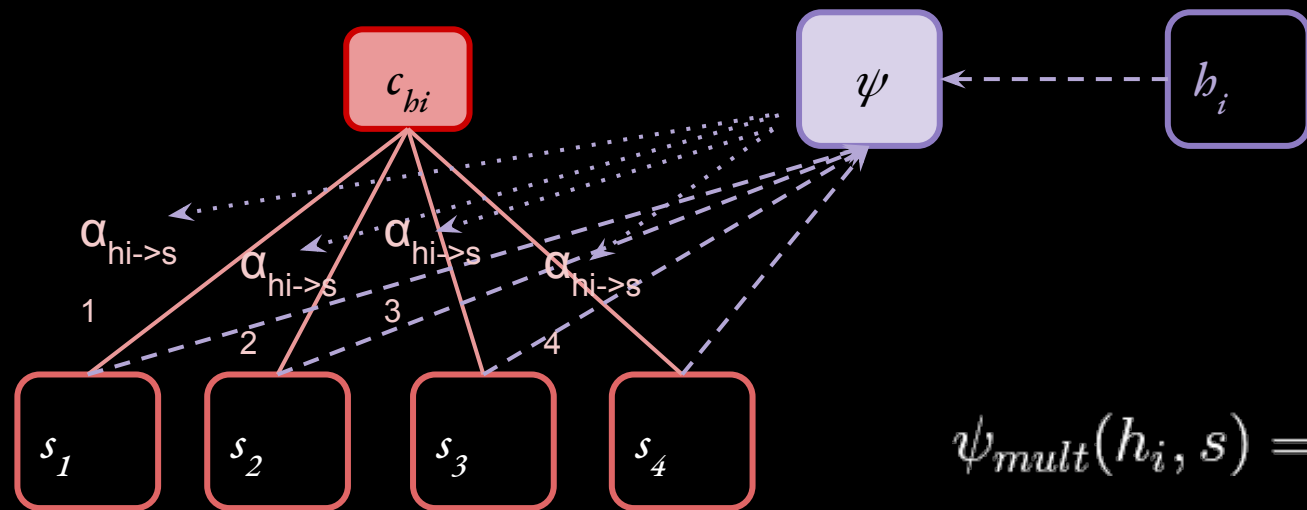
$$c_{h_i} = \sum_{n=1}^{|s|} \alpha_{h_i \rightarrow s_n} s_n$$

Attention



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

Attention

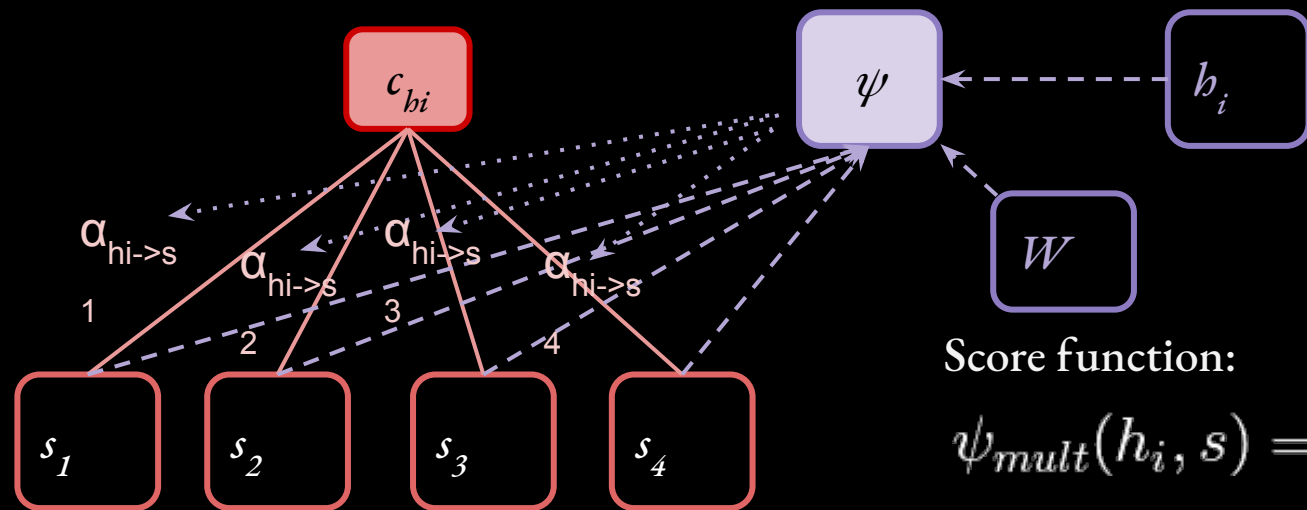


$$\psi_{mult}(h_i, s) = s^T W h_i$$

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

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

Attention



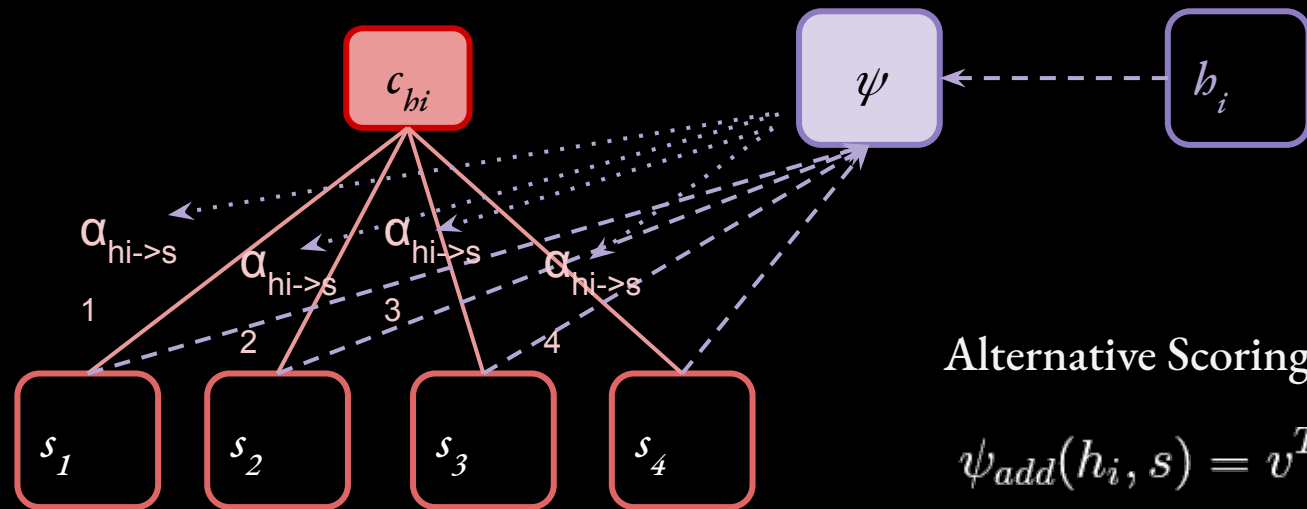
Score function:

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

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

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

Attention



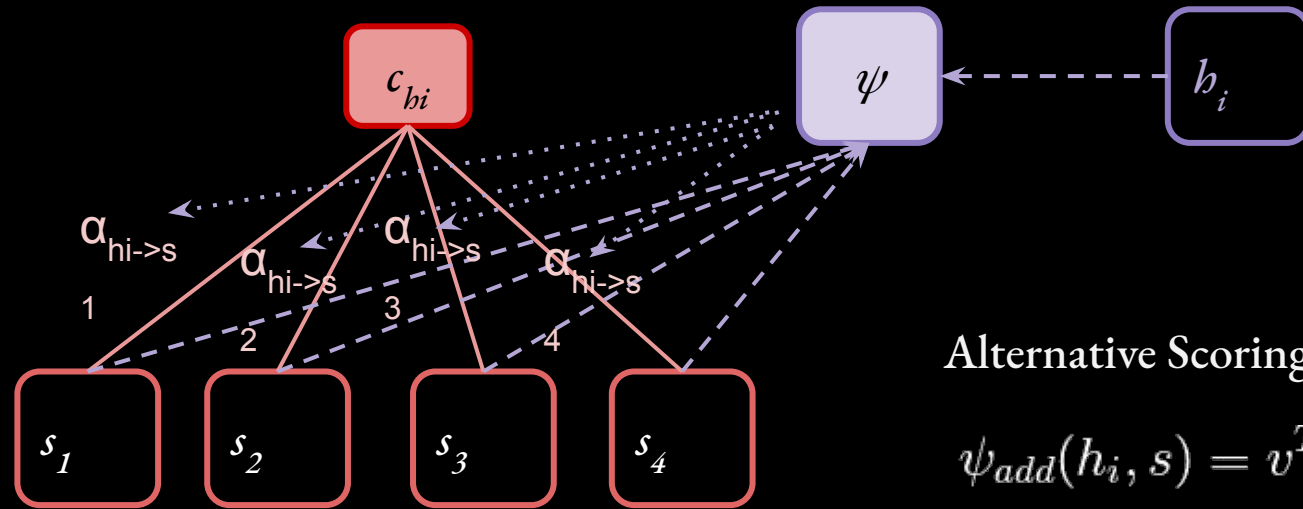
Alternative Scoring Functions

$$\psi_{add}(h_i, s) = v^T \tanh(W_h h_i + W_s s)$$

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

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

Attention



If variables are standardized, matrix multiply produces a similarity score.

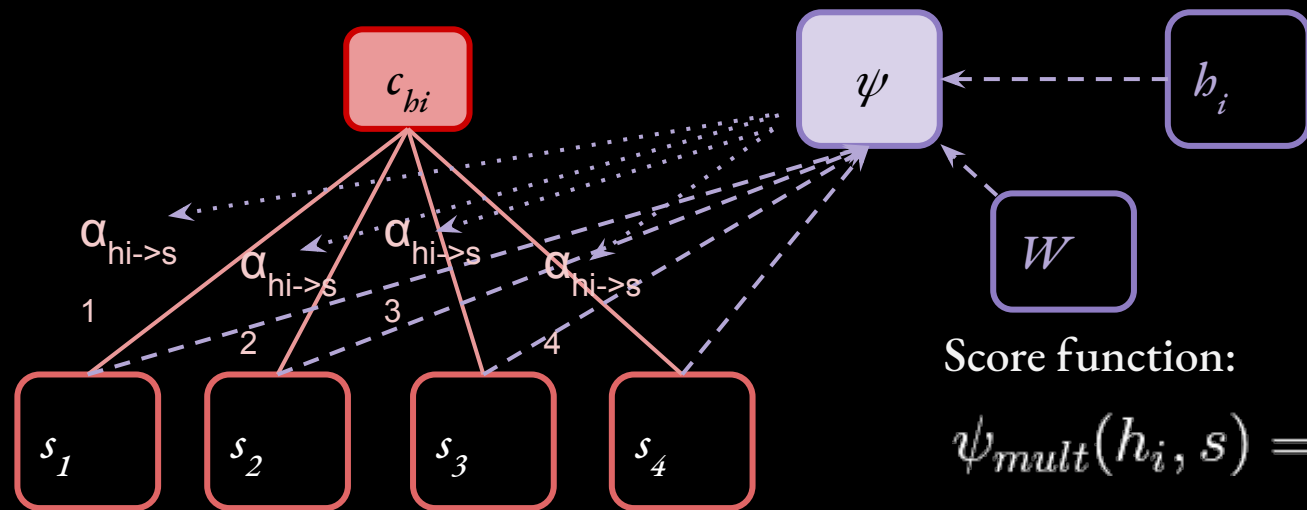
Alternative Scoring Functions

$$\psi_{add}(h_i, s) = v^T \tanh(W_h h_i + W_s s)$$

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

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

Attention



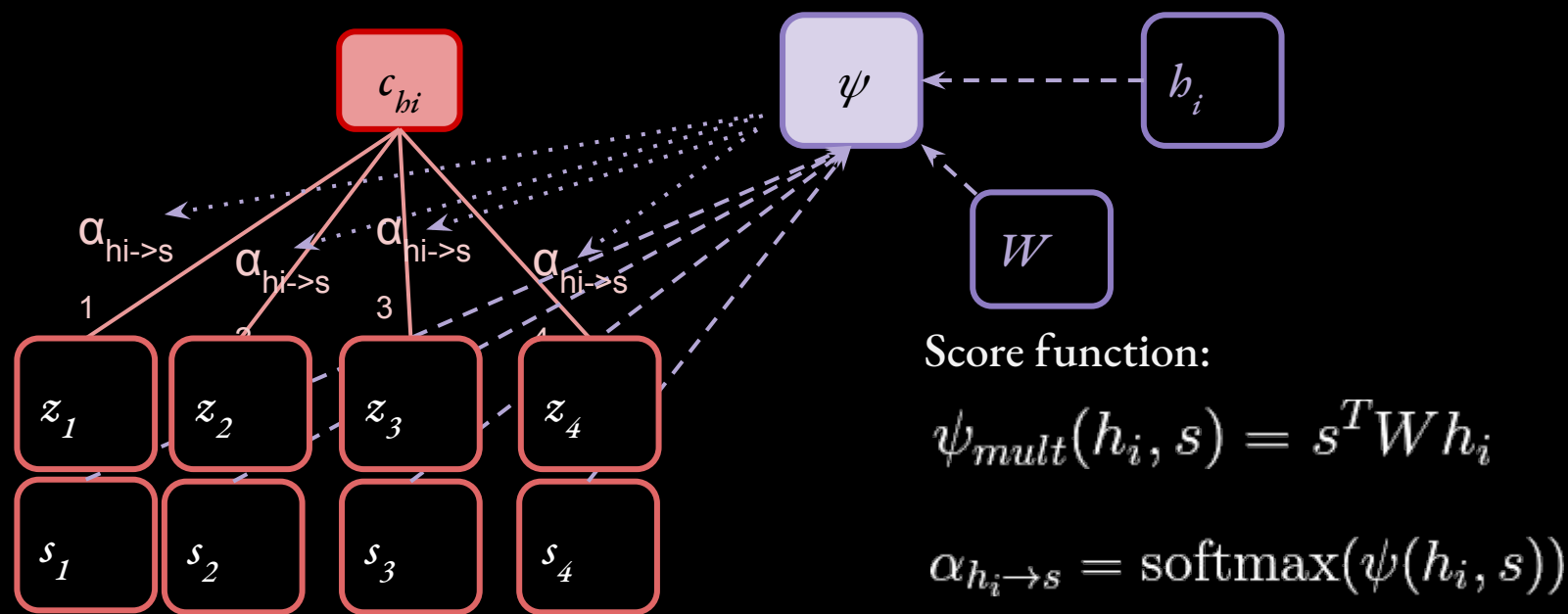
Score function:

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

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

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

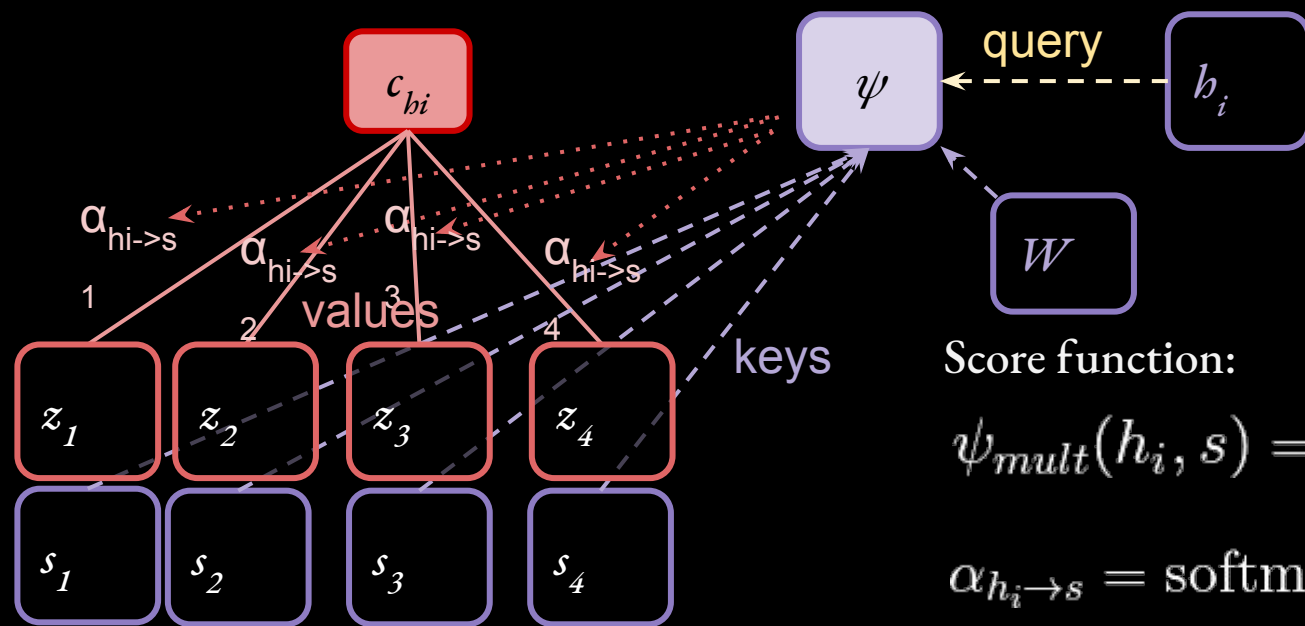
Attention



A useful abstraction is to make the vector attended to (the “value vector”, Z) separate than the “key vector” (s).

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

Attention



Score function:

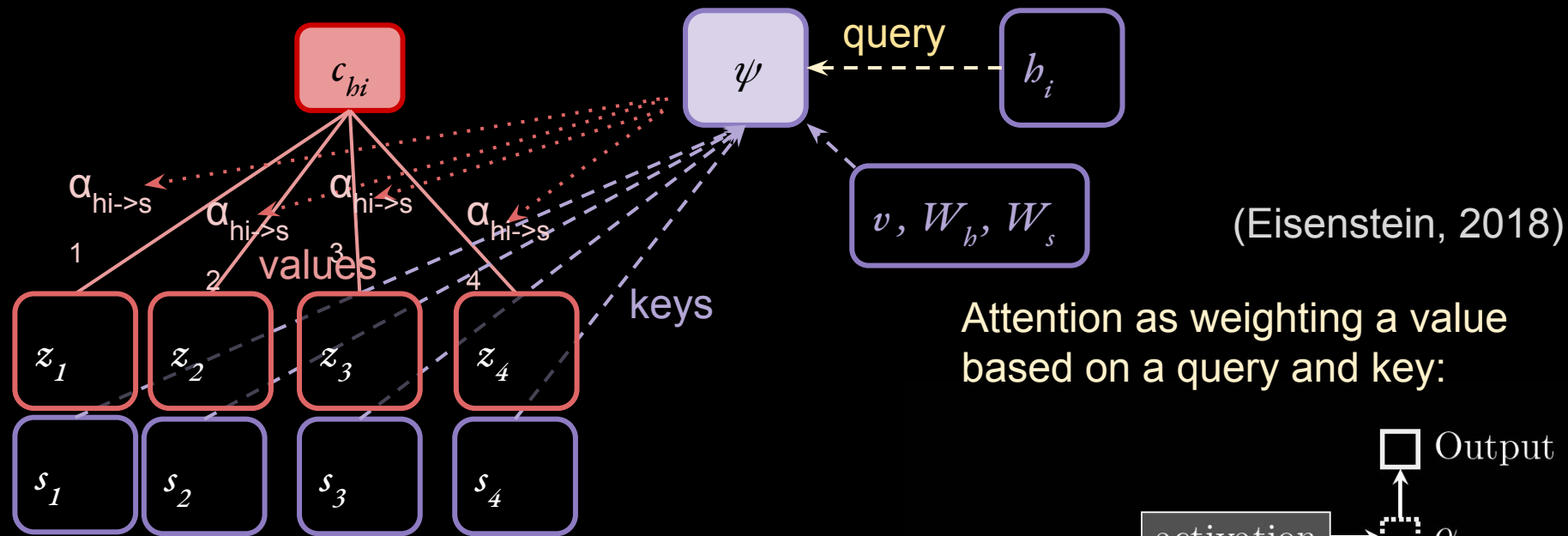
$$\psi_{mult}(h_i, s) = s^T W h_i$$

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

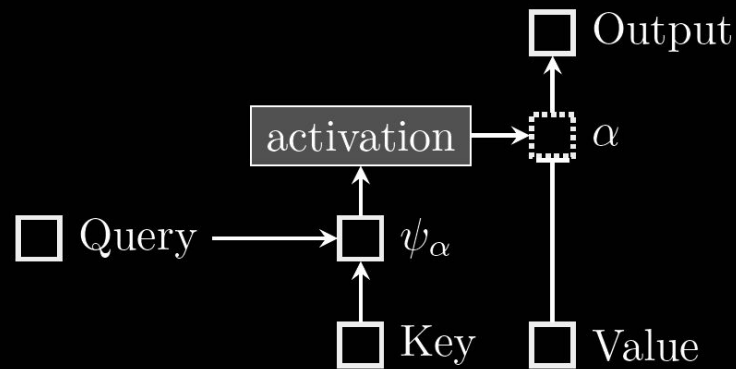
$$c_{h_i} = \sum_{n=1}^{|s|} \alpha_{h_i \rightarrow s_n} z_n$$

A useful abstraction is to make the vector attended to (the “value vector”, Z) separate than the “key vector” (s).

Attention



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



Machine Translation

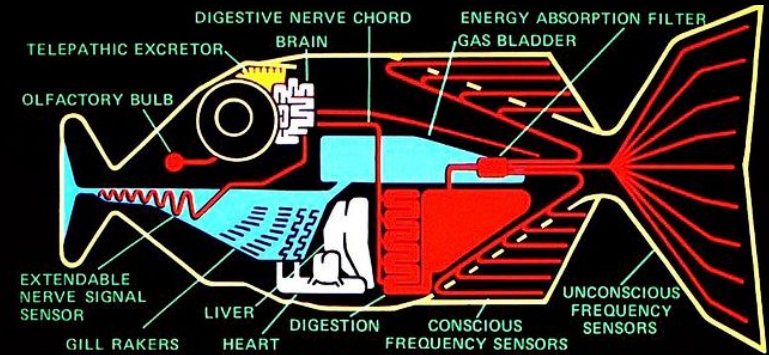
As an optimization problem (Eisenstein, 2018):

$$\hat{w}^{(t)} = \operatorname{argmax}_{w^{(t)}} \Psi(w^{(s)}, w^{(t)})$$

Machine Translation

Why?

- \$40billion/year industry
- A center piece of many genres of science fiction
- A fairly “universal” problem:
 - Language understanding
 - Language generation
- Societal benefits of inter-cultural communication

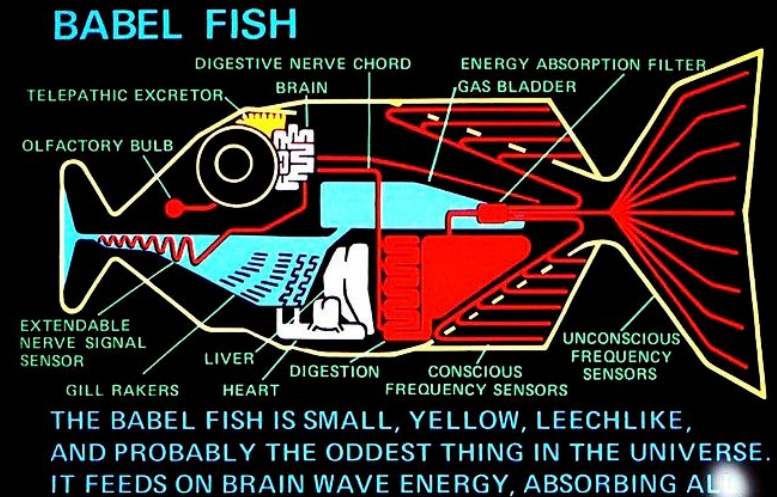


THE BABEL FISH IS SMALL, YELLOW, LEECHLIKE, AND PROBABLY THE ODDEST THING IN THE UNIVERSE. IT FEEDS ON BRAIN WAVE ENERGY, ABSORBING ALL

Machine Translation

Why?

- \$40billion/year industry
- A center piece of many genres of science fiction
- A fairly “universal” problem:
 - Language understanding
 - Language generation
- Societal benefits of inter-cultural communication



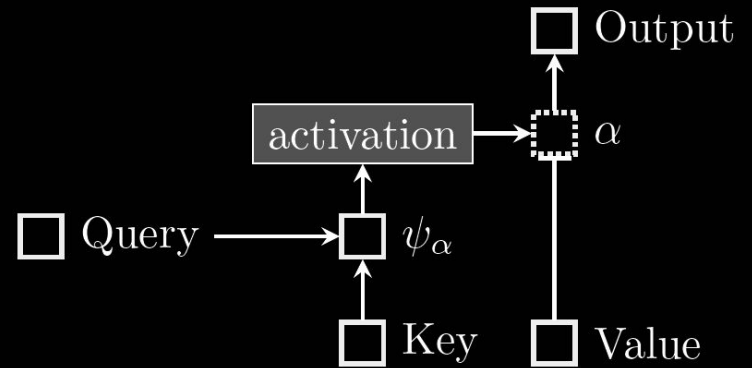
(Douglas Adams)

Machine Translation

Why Neural Network Approach works? (Manning, 2018)

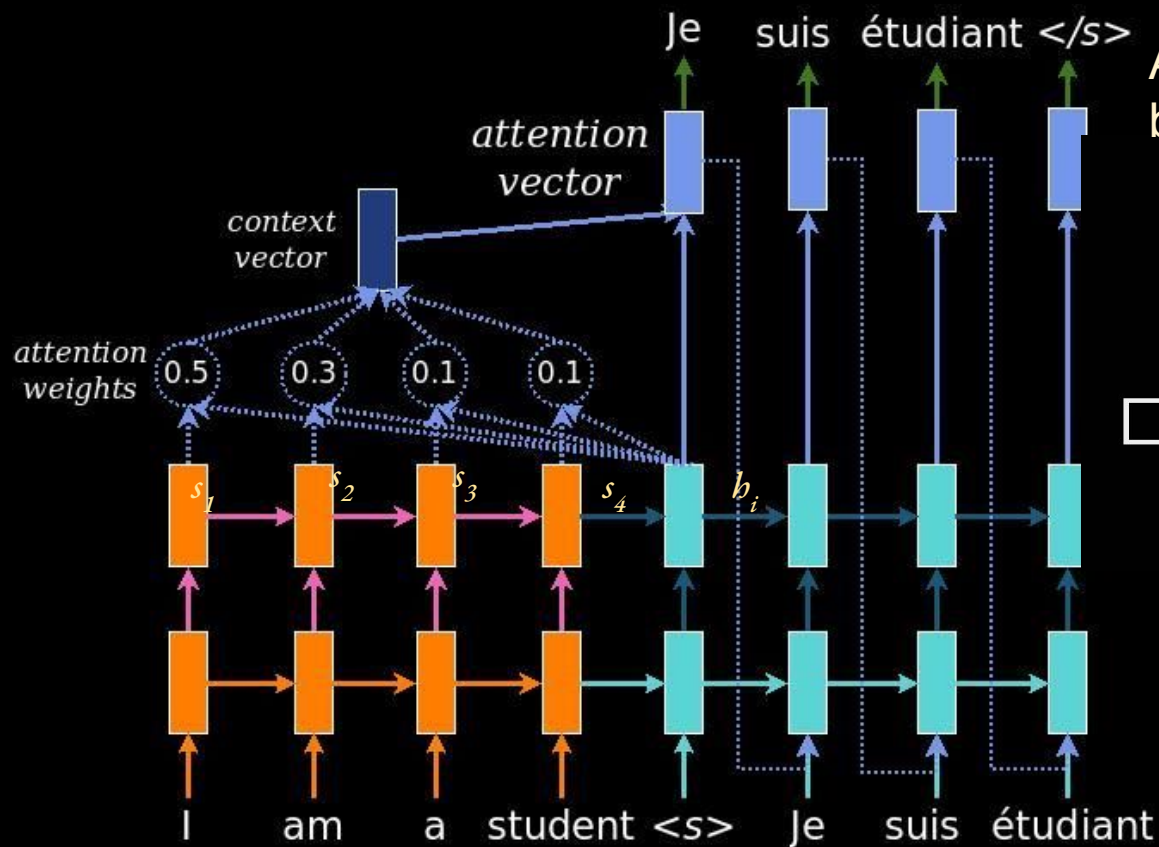
- Joint end-to-end training: learning all parameters at once.
- Exploiting distributed representations (embeddings)
- Exploiting variable-length context
- High quality generation from deep decoders - stronger language models (even when wrong, make sense)

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

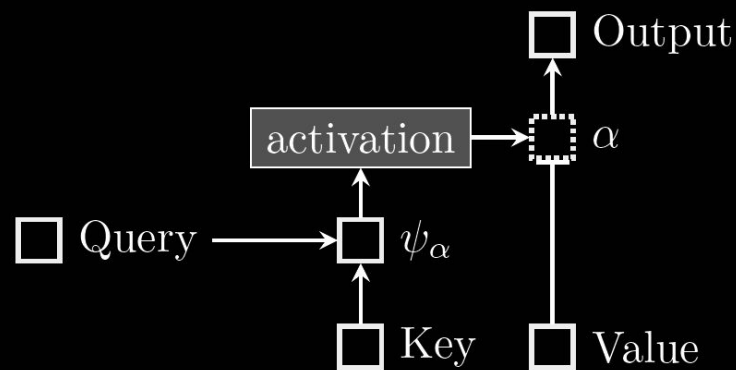


(Eisenstein, 2018)

Attention

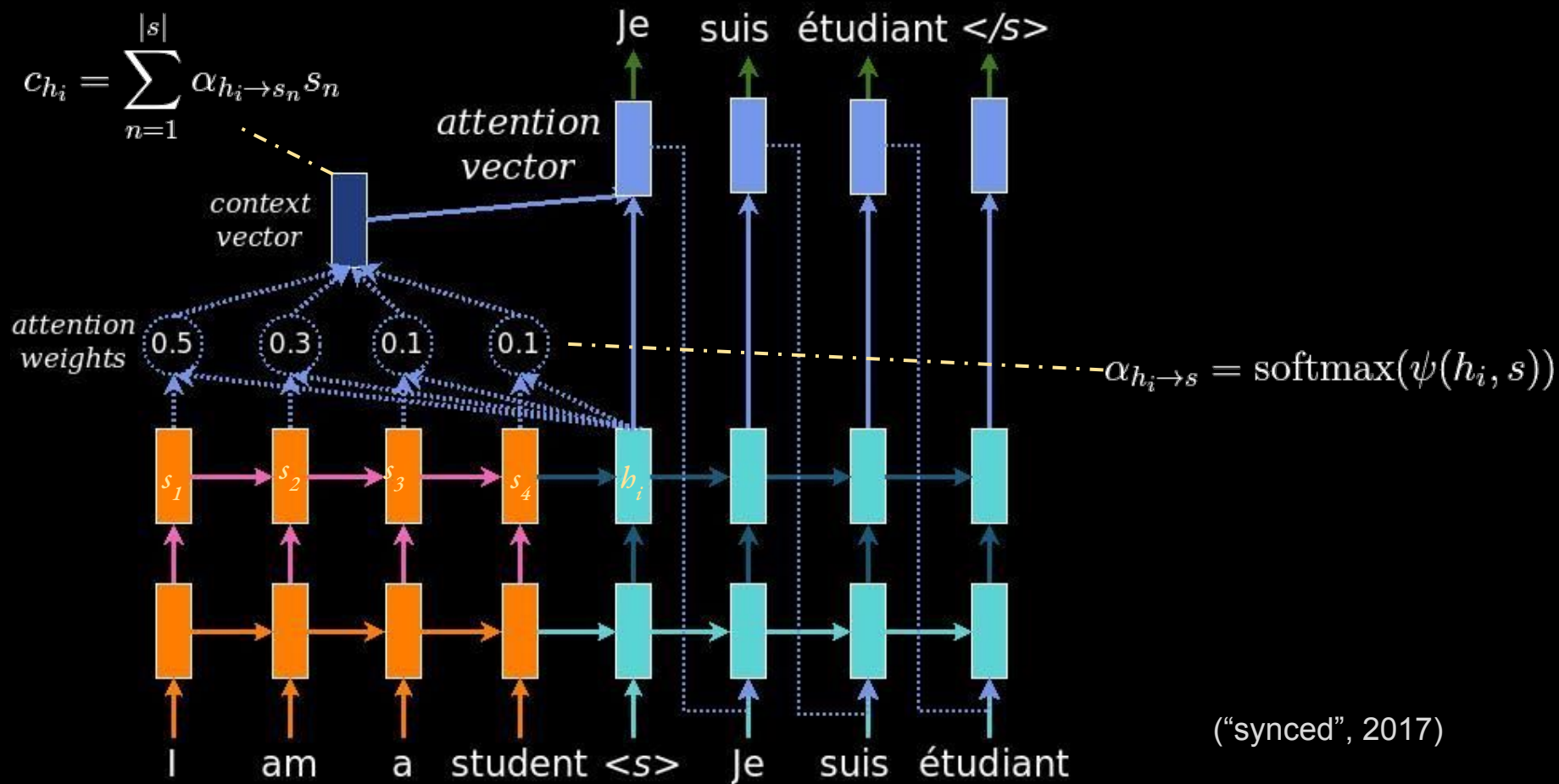


Attention as weighting a value based on a query and key:



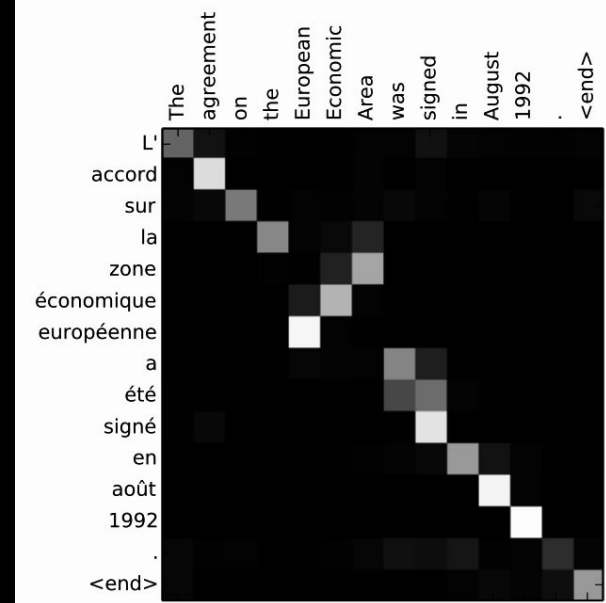
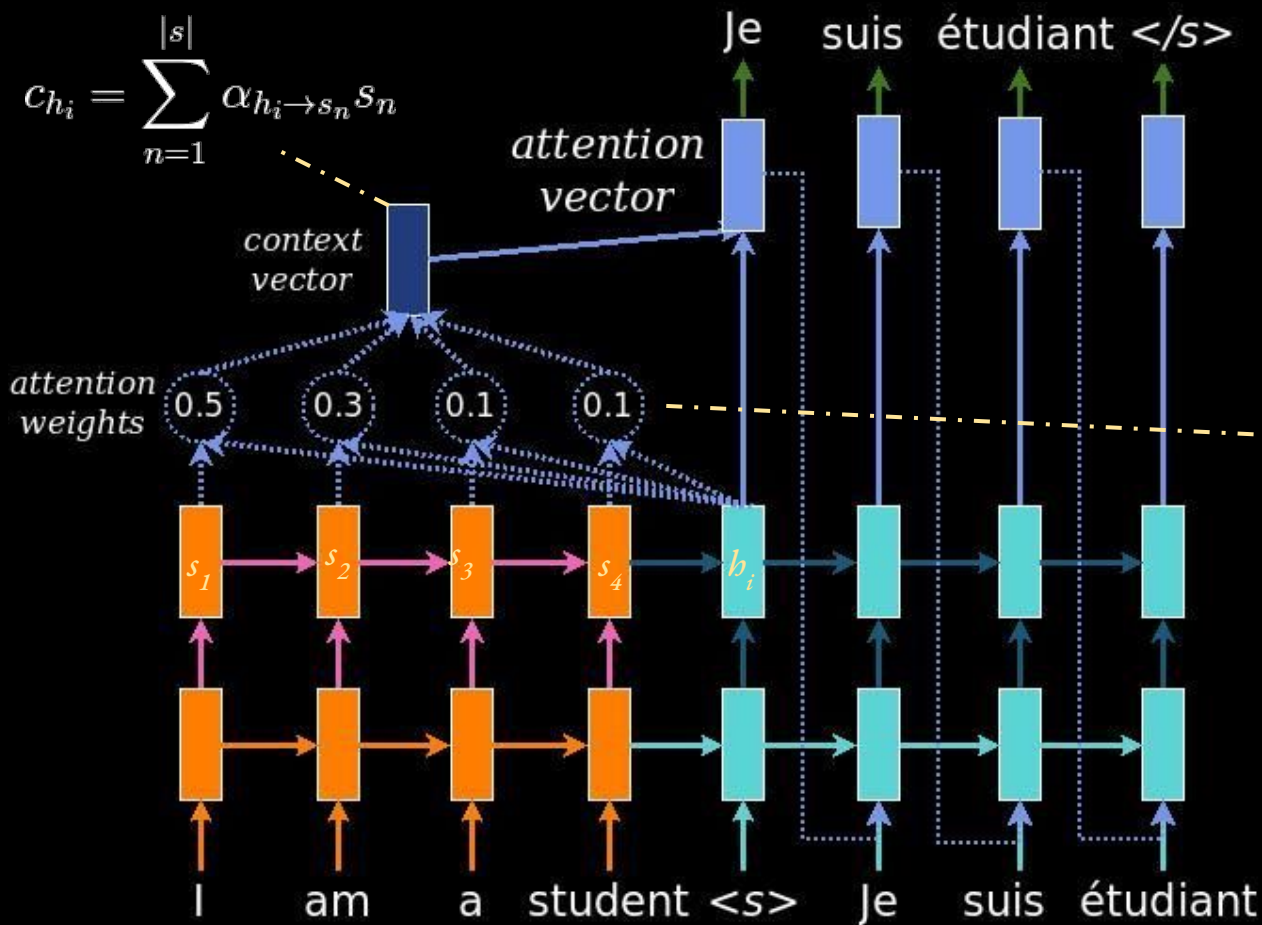
("synced", 2017)

Attention



Attention

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

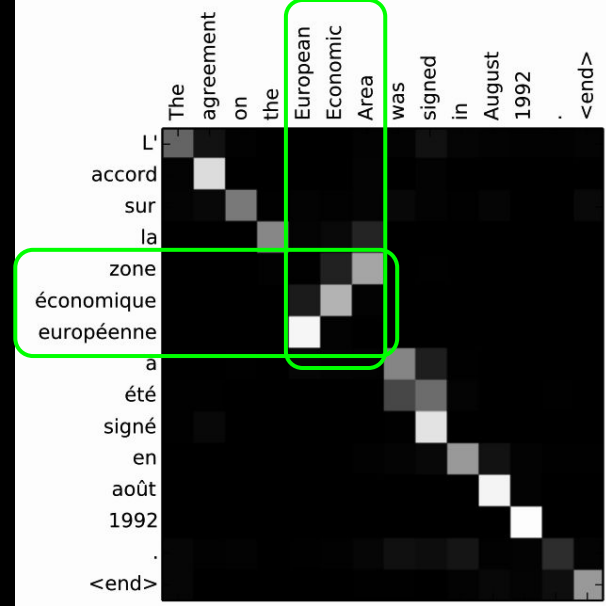
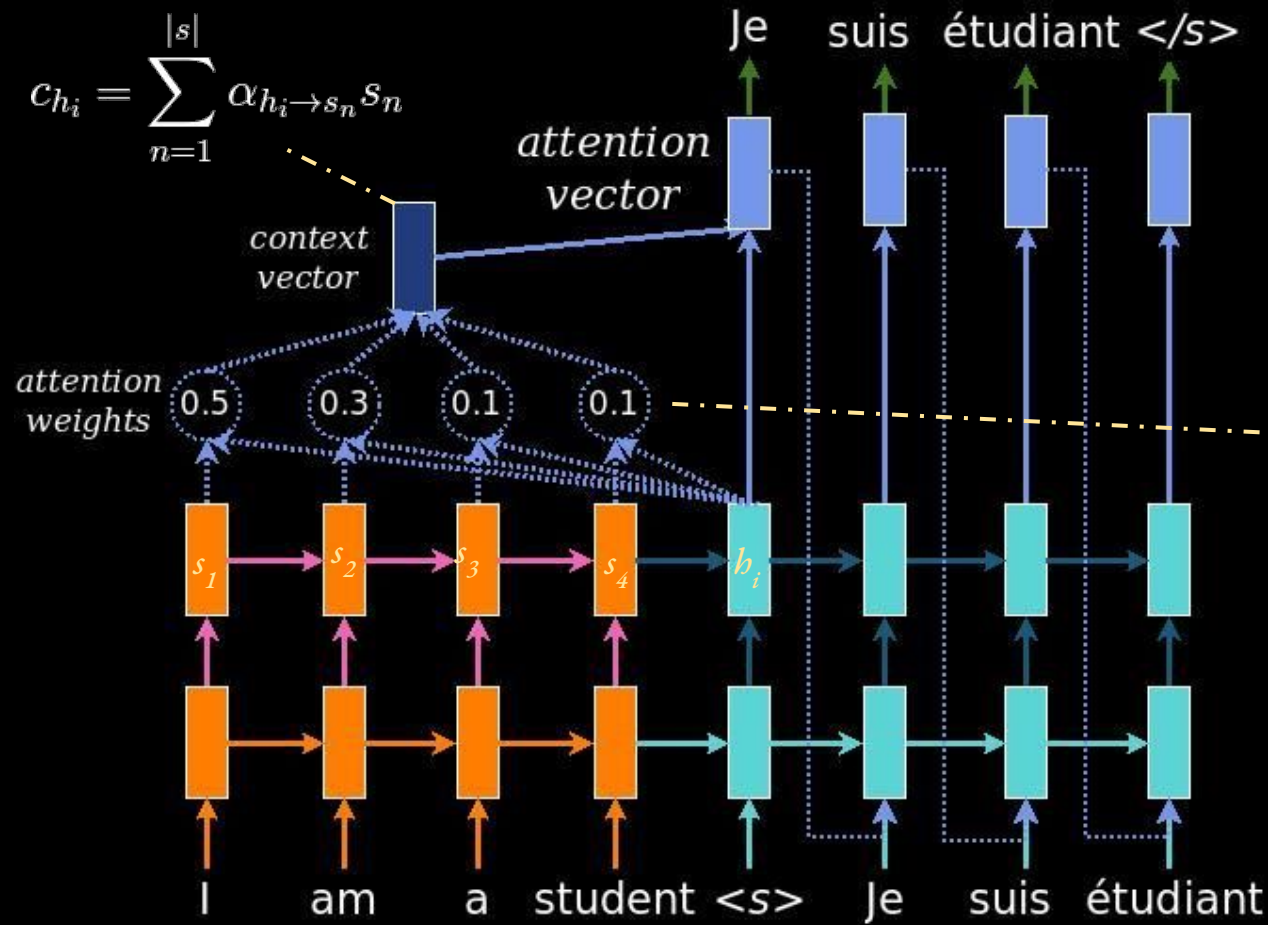


$$\alpha_{h_i \rightarrow s} = \text{softmax}(\psi(h_i, s))$$

(“synced”, 2017)

Attention

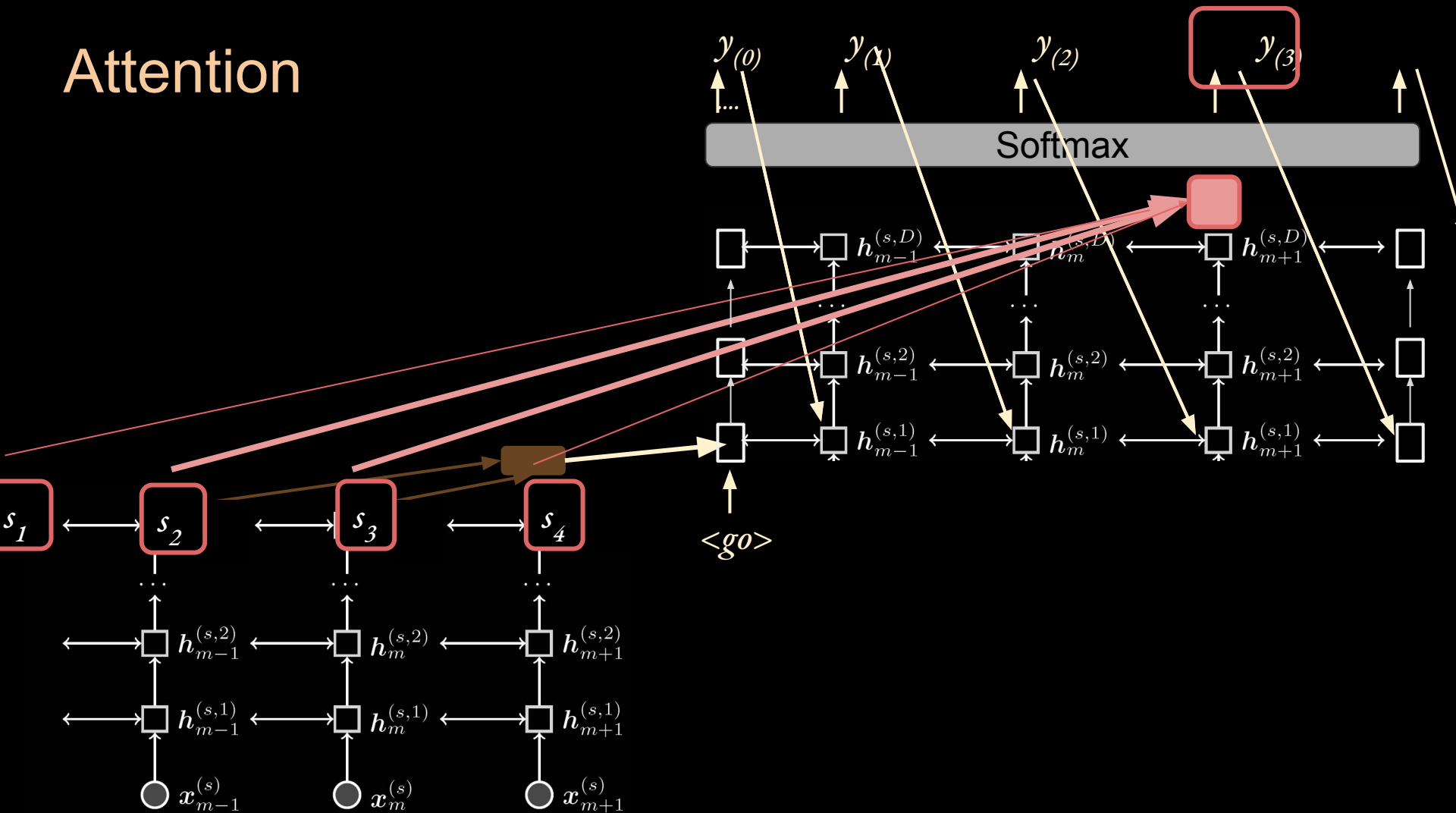
$$c_{h_i} = \sum_{n=1}^{|s|} \alpha_{h_i \rightarrow s_n} s_n$$



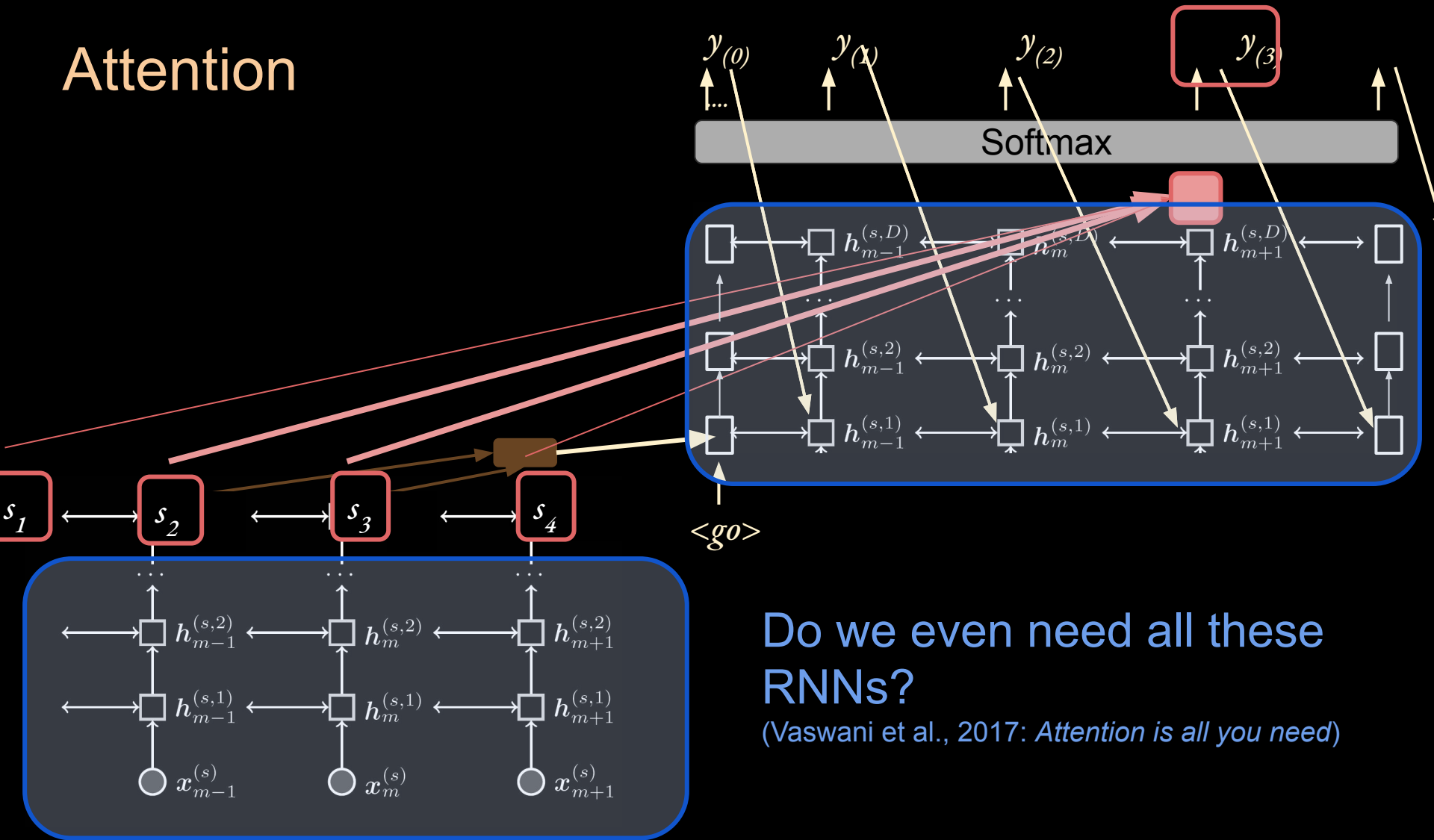
$$\alpha_{h_i \rightarrow s} = \text{softmax}(\psi(h_i, s))$$

(“synced”, 2017)

Attention



Attention



Do we even need all these RNNs?
(Vaswani et al., 2017: *Attention is all you need*)

Evolution of Sequence Modeling

RNNs

LSTMs

LSTMs with Attention

Attention without LSTMs

Evolution of Sequence Modeling

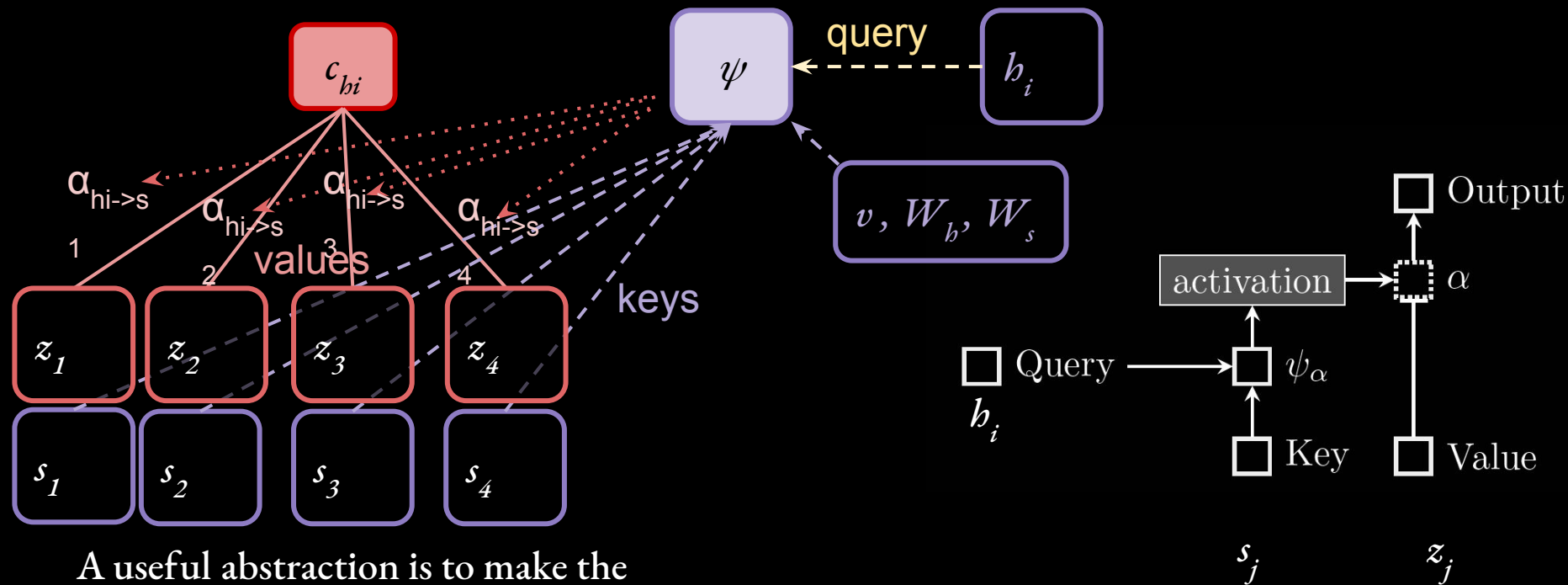
RNNs

LSTMs

LSTMs with Attention

Attention without LSTMs

Attention

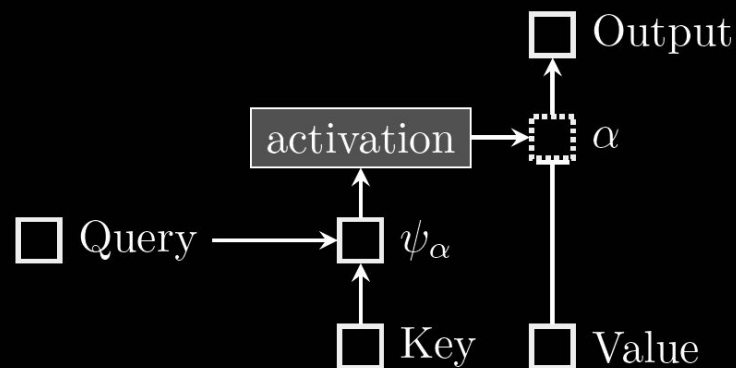


A useful abstraction is to make the vector attended to (the “value vector”, Z) separate than the “key vector” (s).

(Eisenstein, 2018)

The Transformer: “Attention-only” models

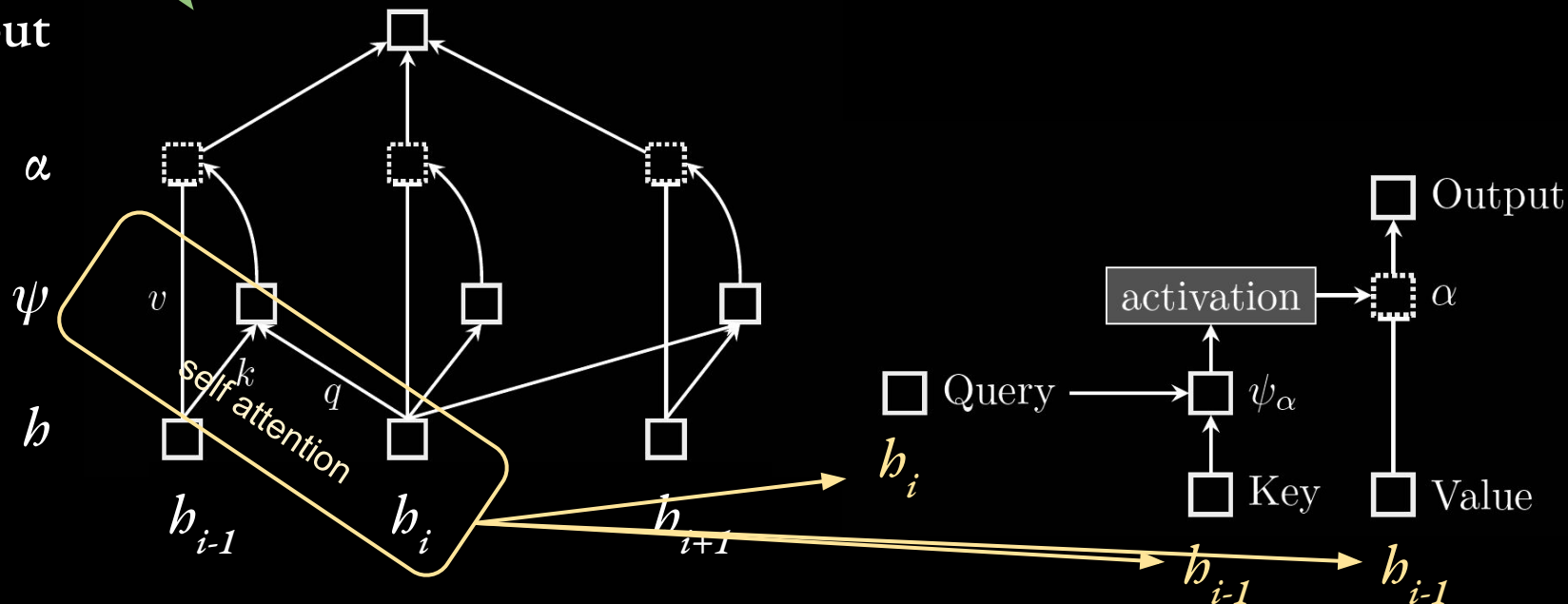
Attention as weighting a value based on a query and key:



(Eisenstein, 2018)

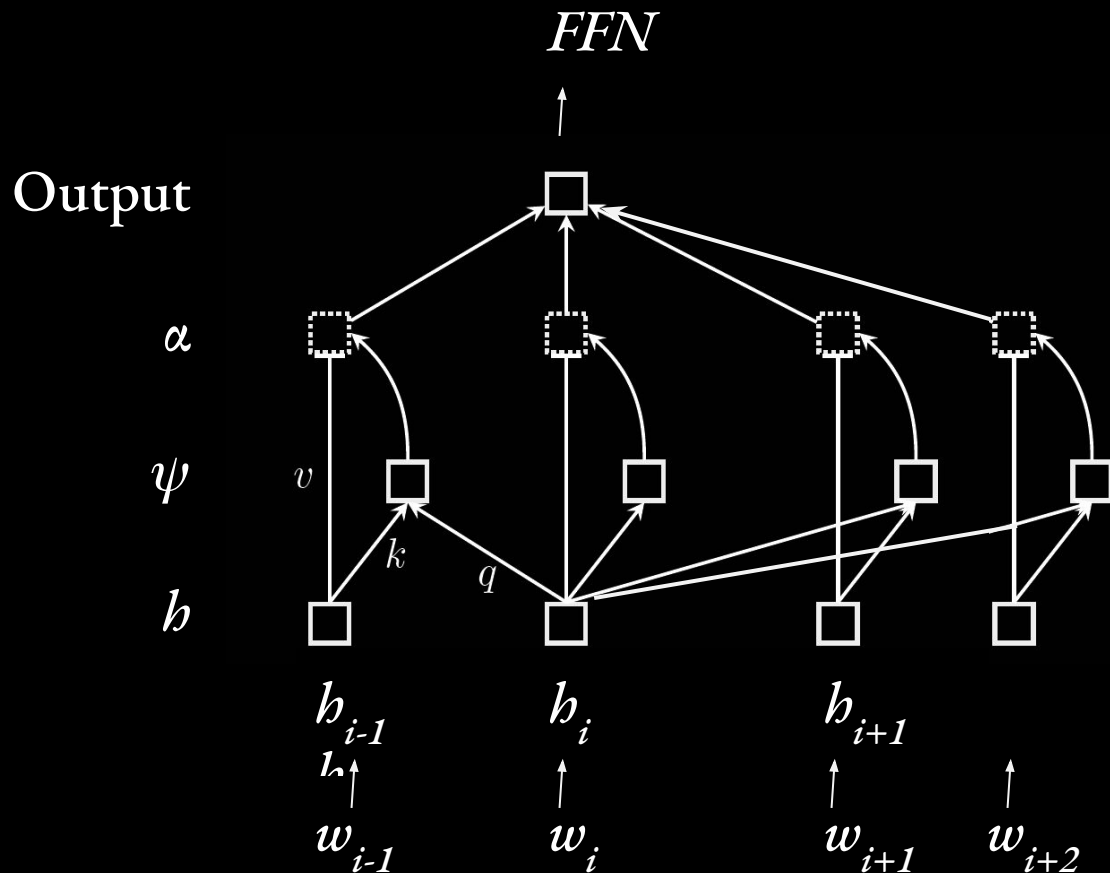
The Transformer: “Attention-only” models

Output

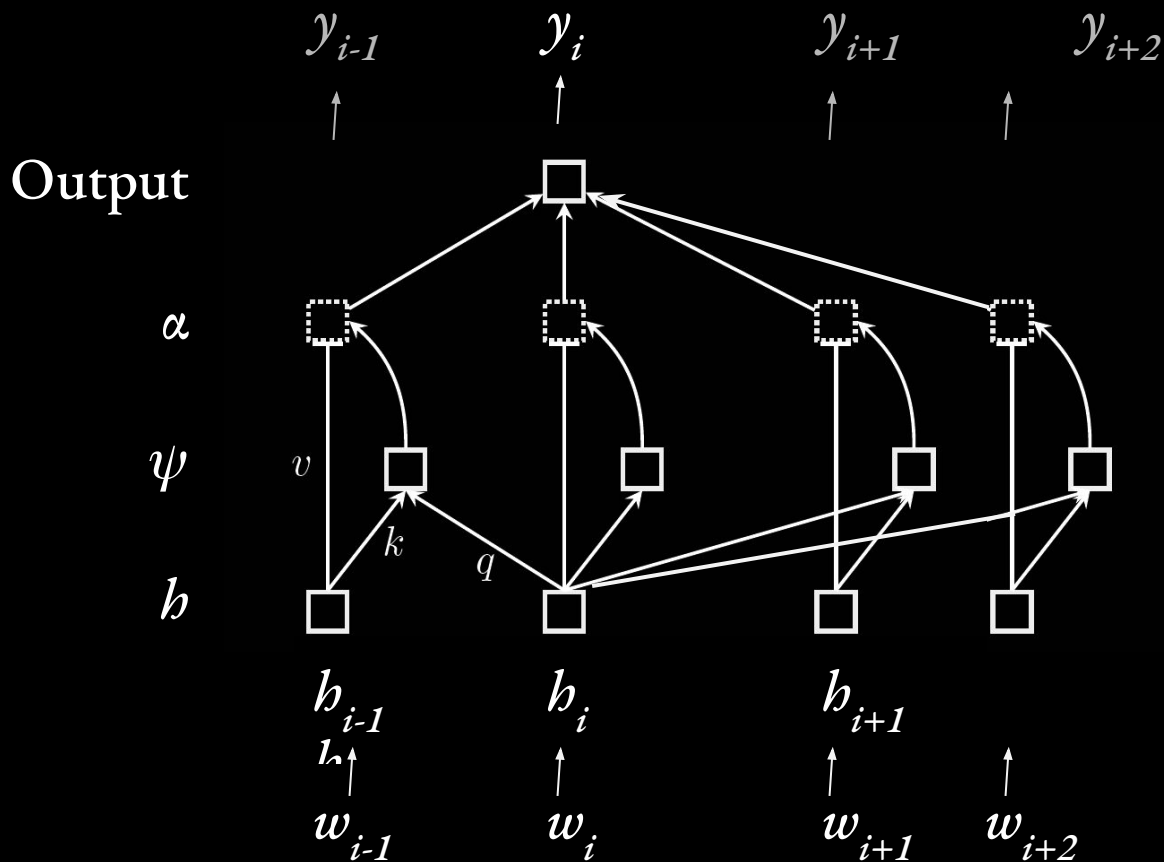


(Eisenstein, 2018)

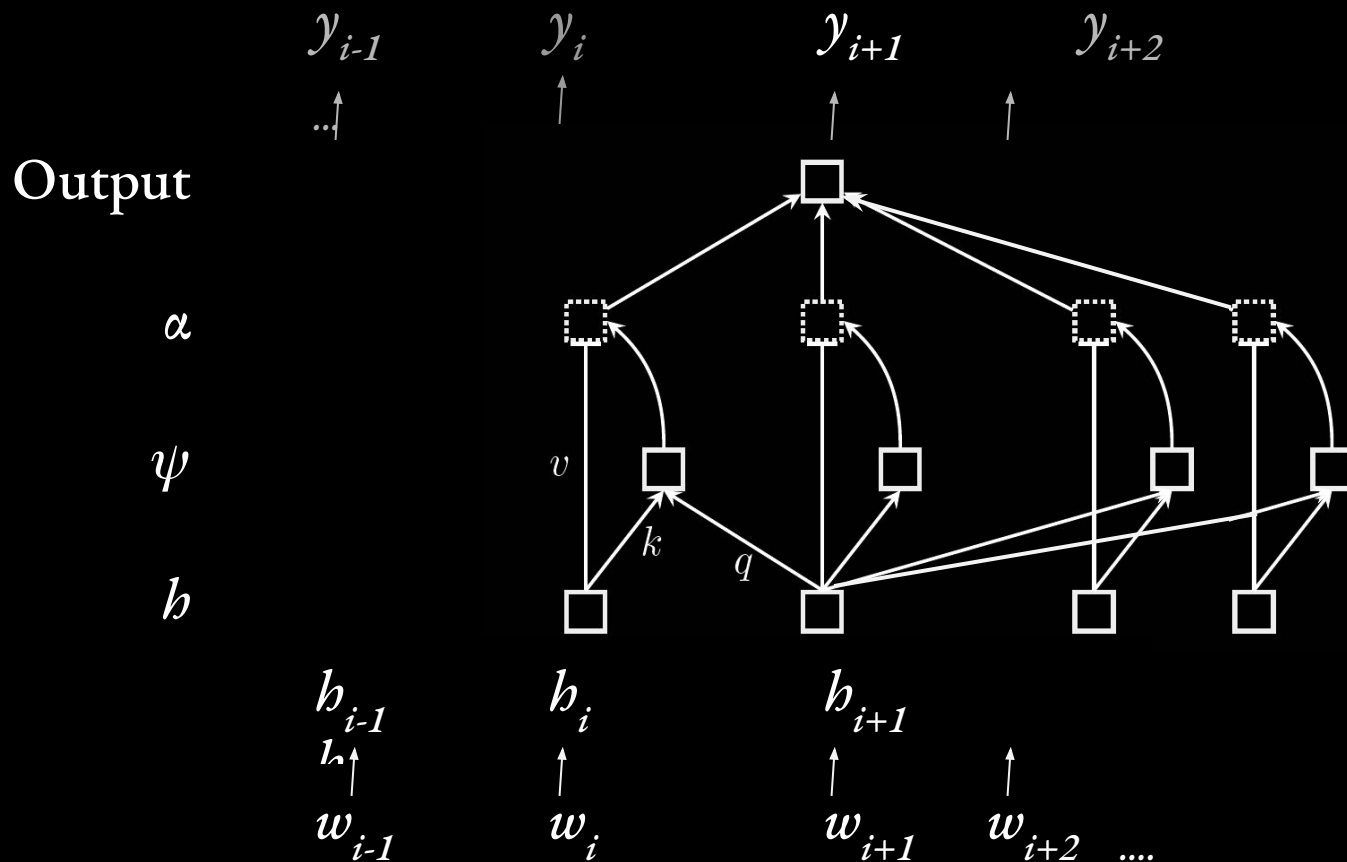
The Transformer: “Attention-only” models



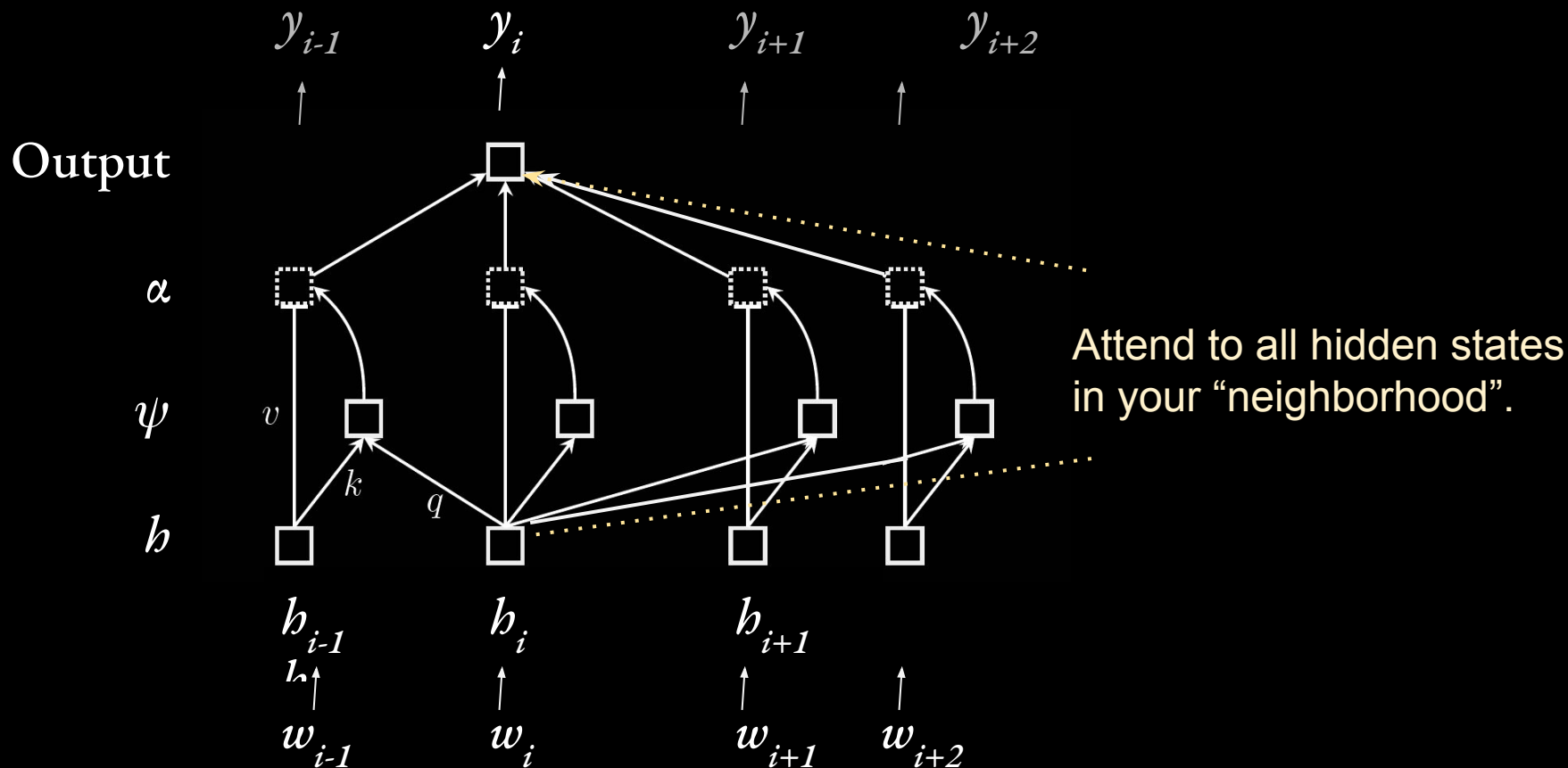
The Transformer: “Attention-only” models



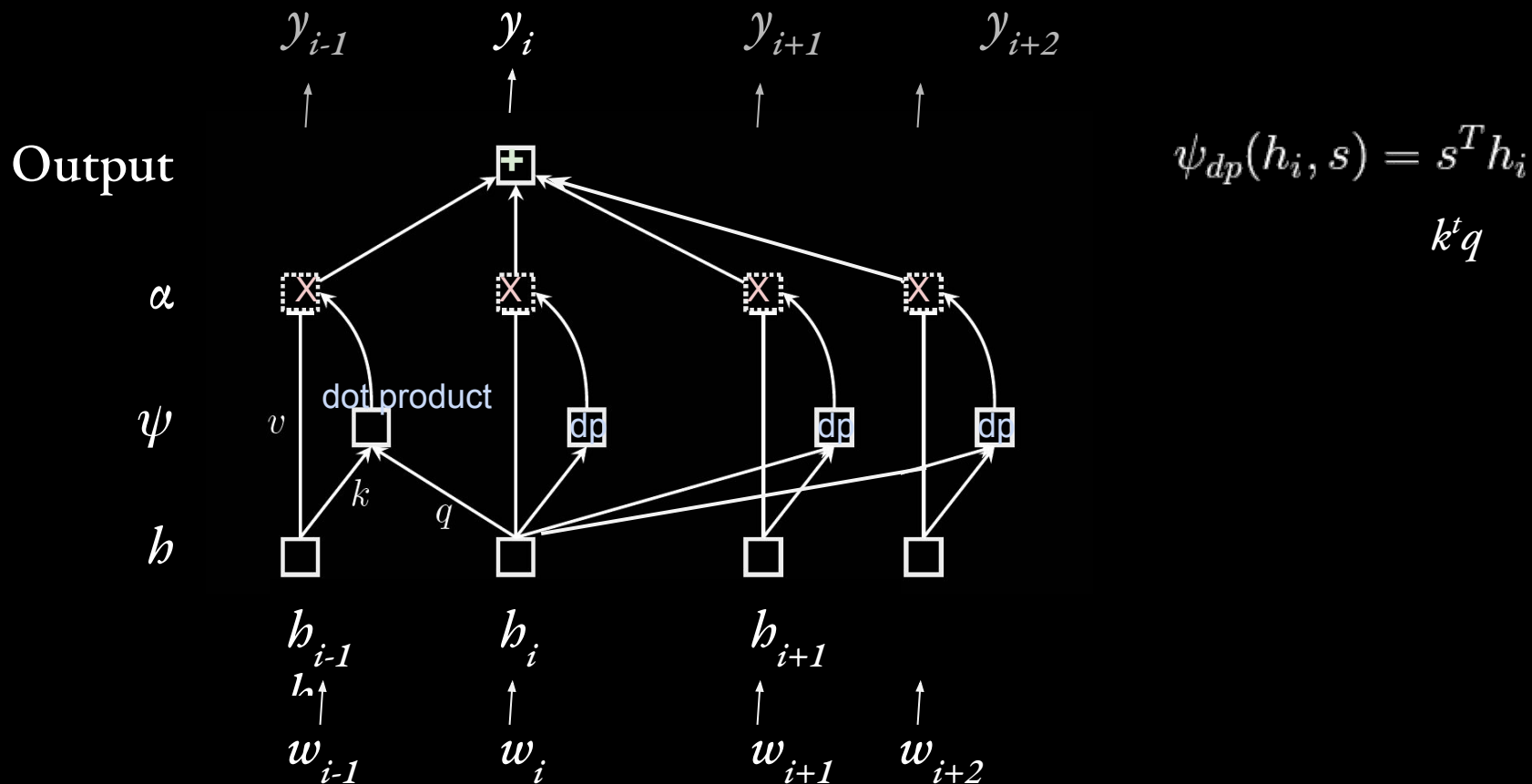
The Transformer: “Attention-only” models



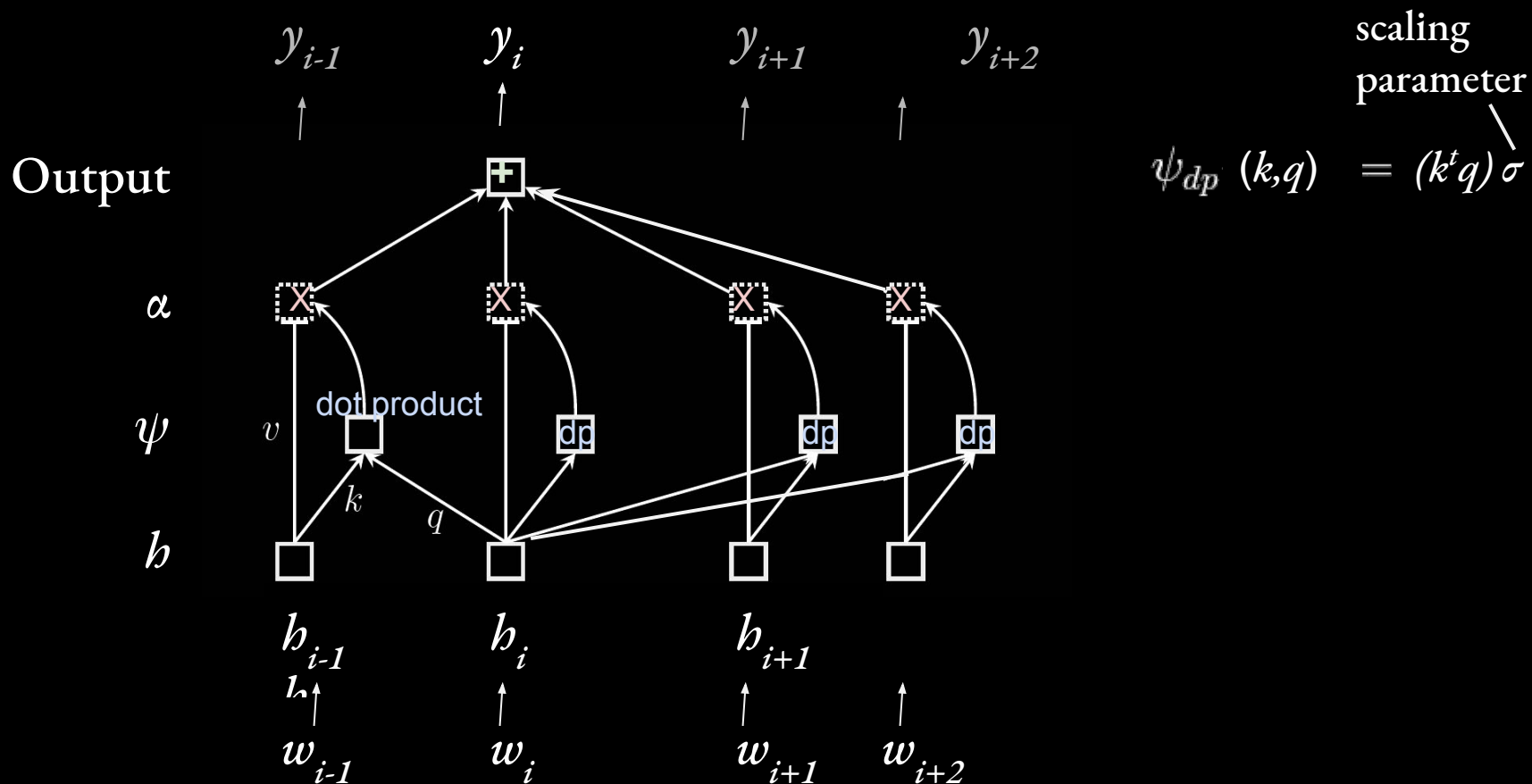
The Transformer: “Attention-only” models



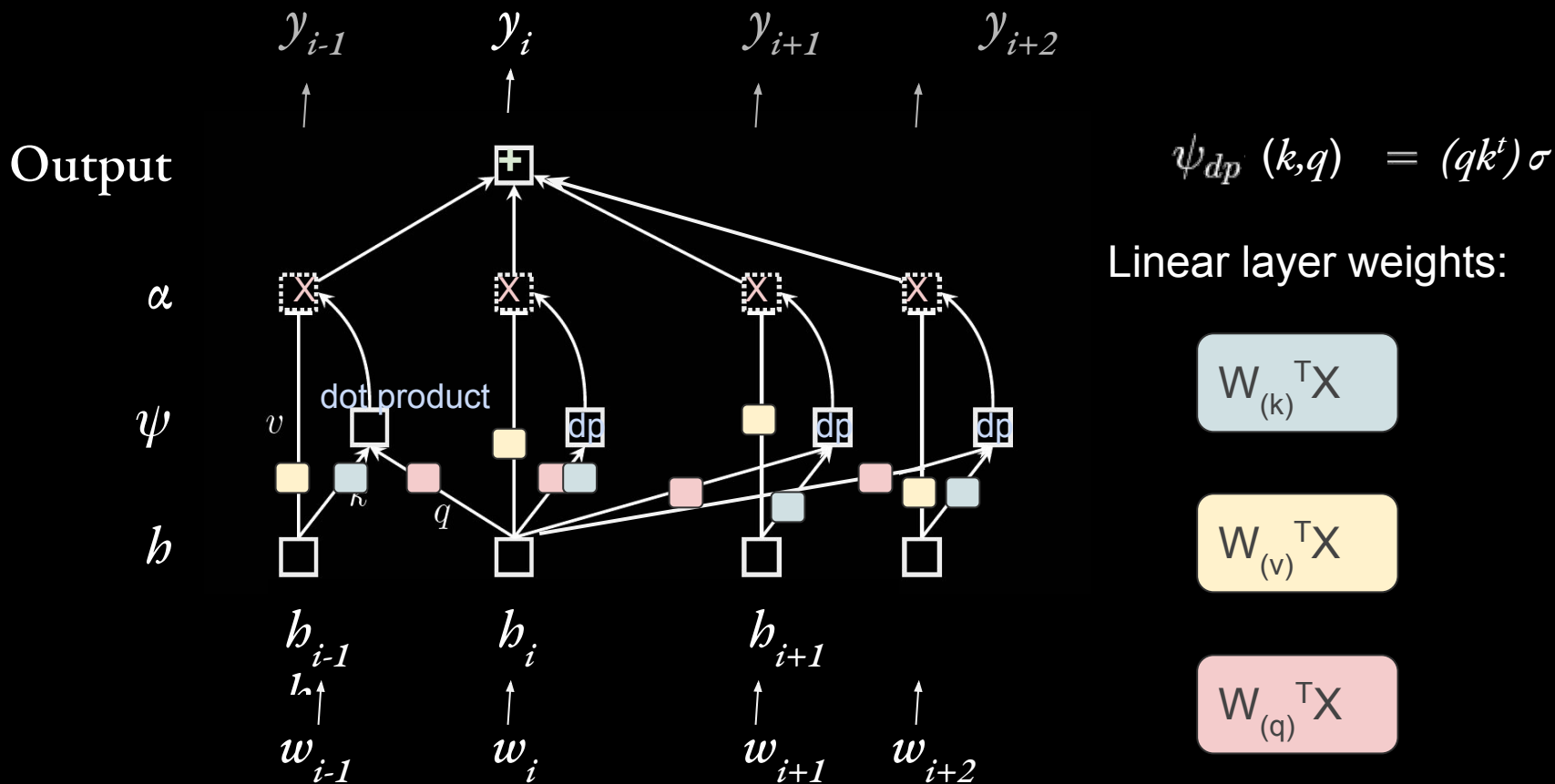
The Transformer: "Attention-only" models



The Transformer: “Attention-only” models



The Transformer: "Attention-only" models



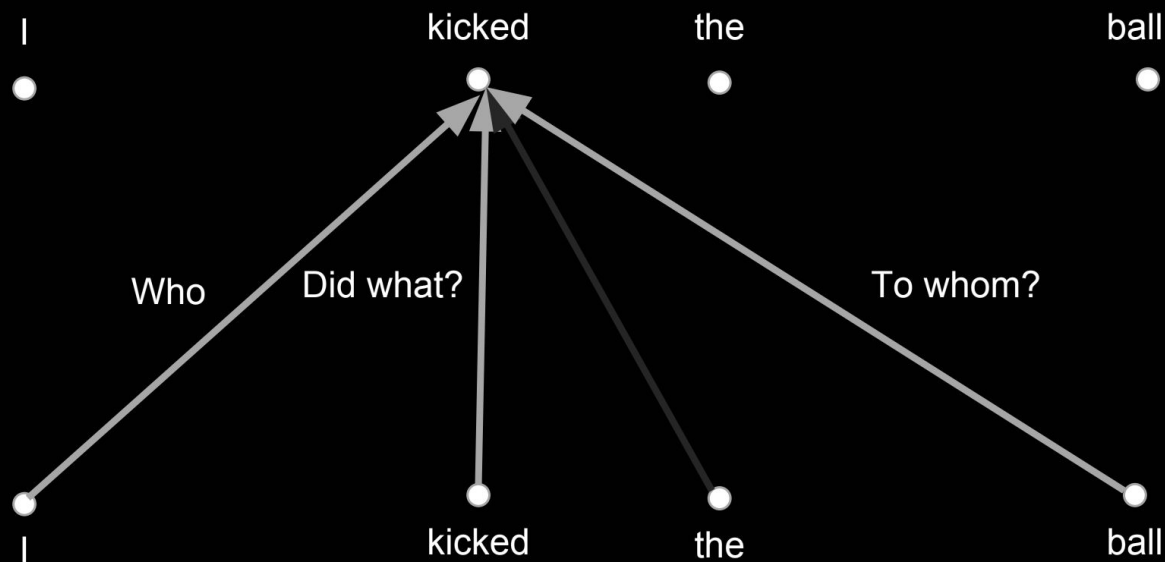
The Transformer: “Attention-only” models

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.

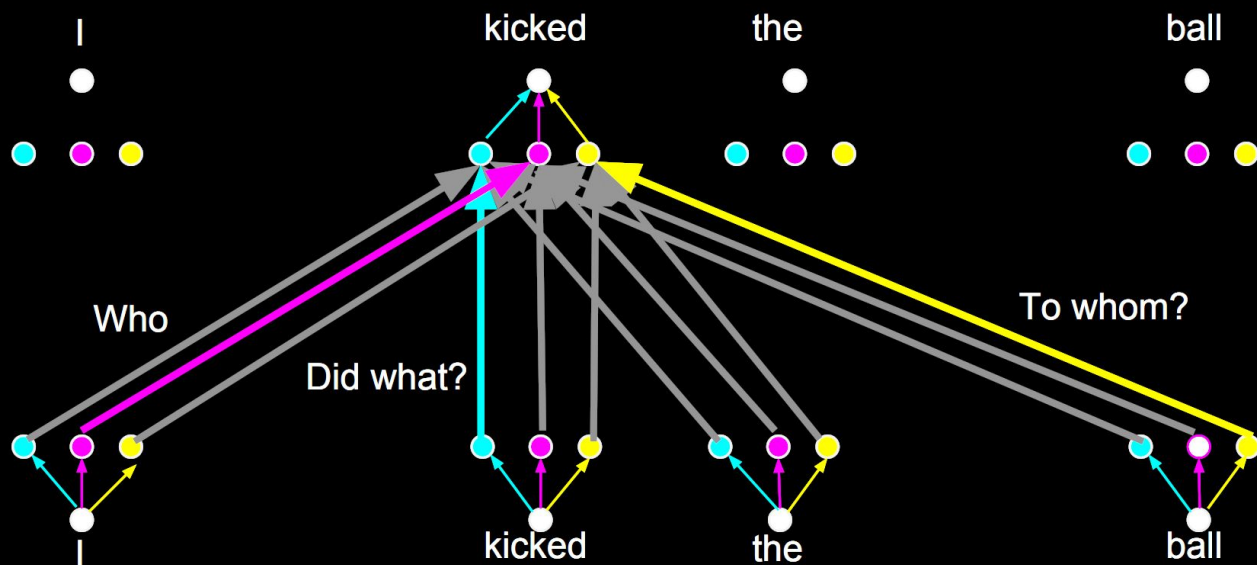
The Transformer

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

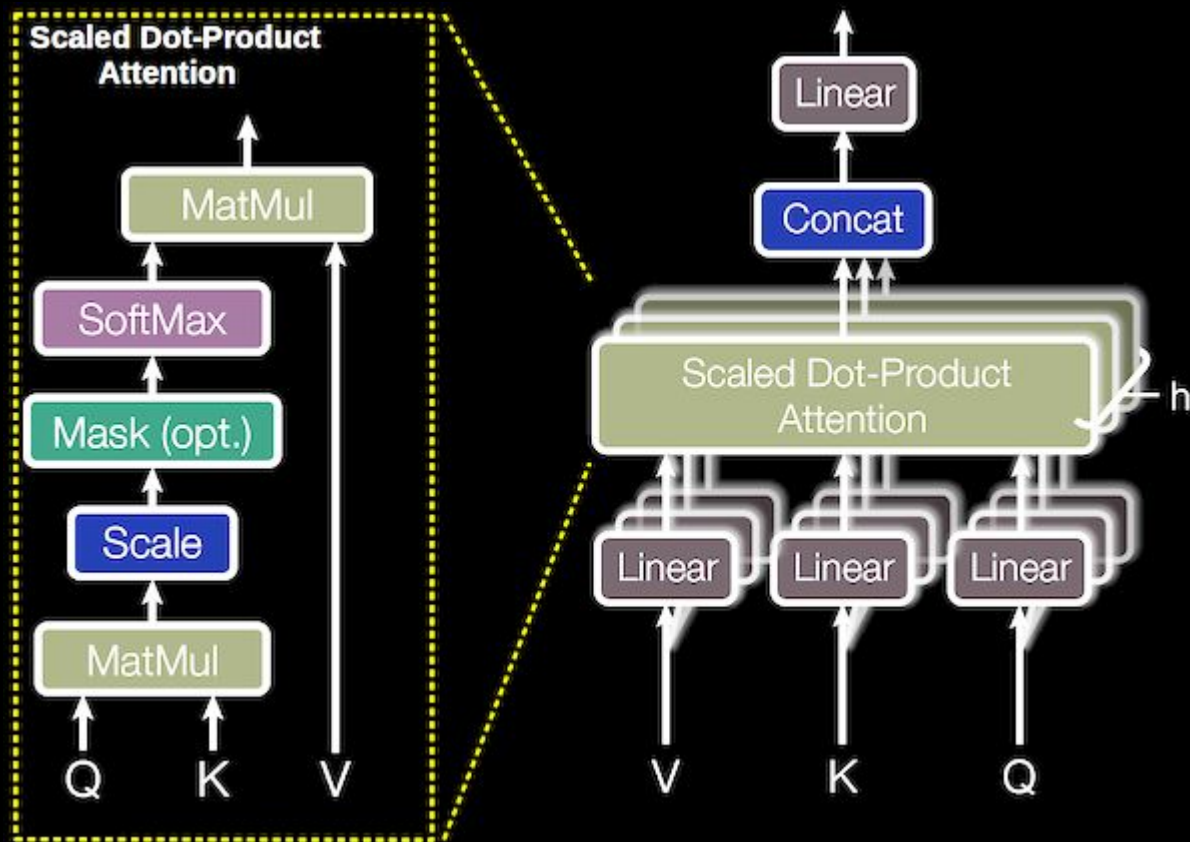


The Transformer

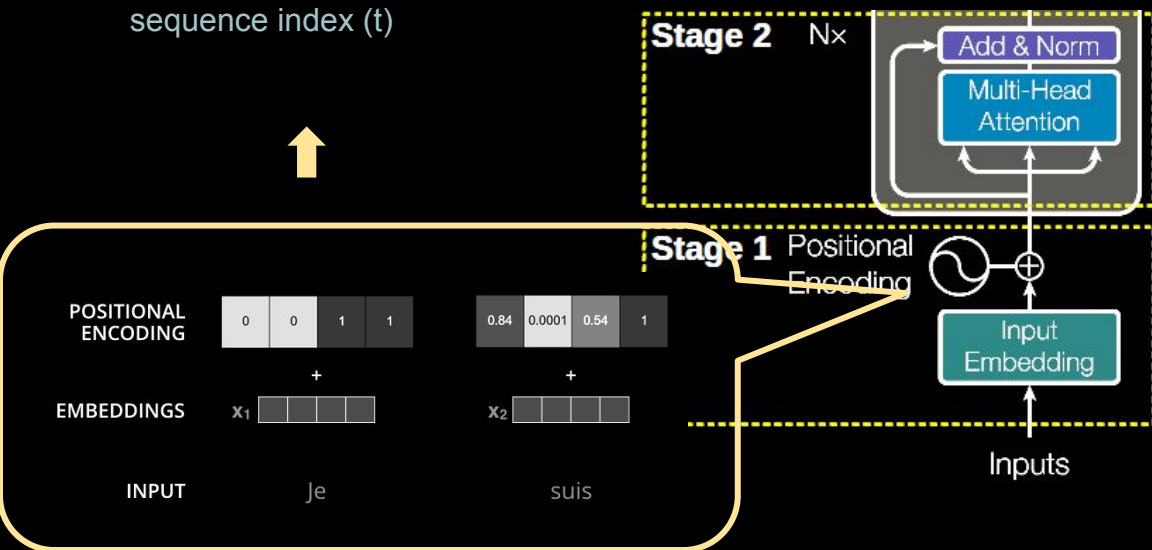
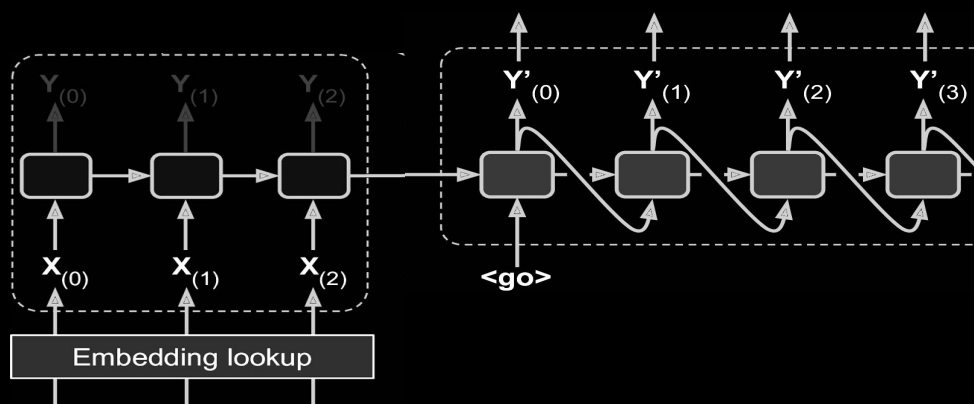
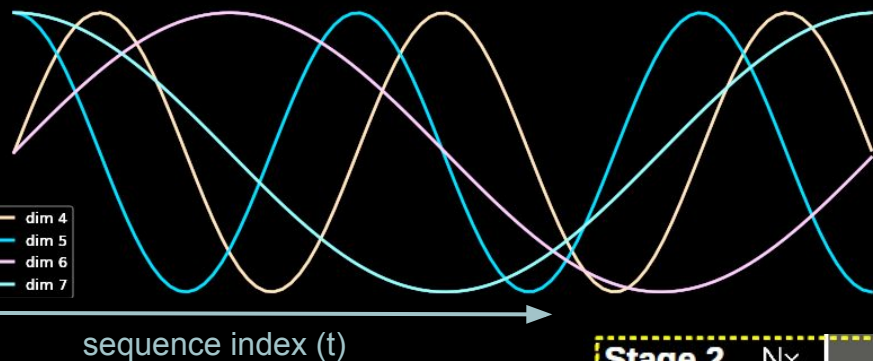
Solution: Multi-head attention



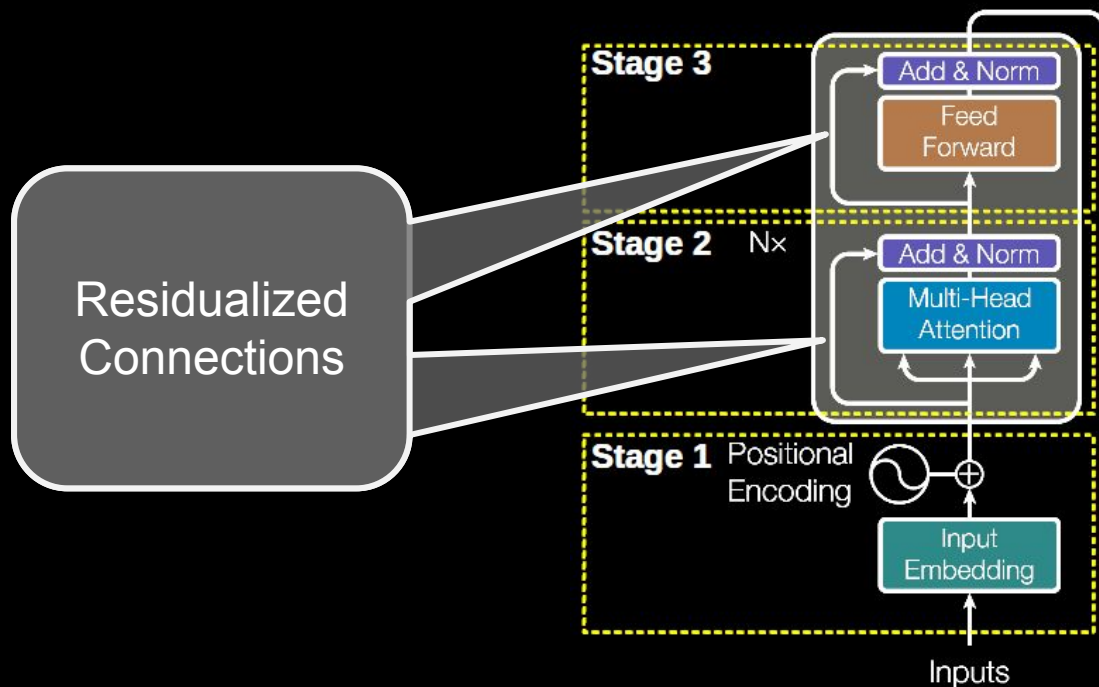
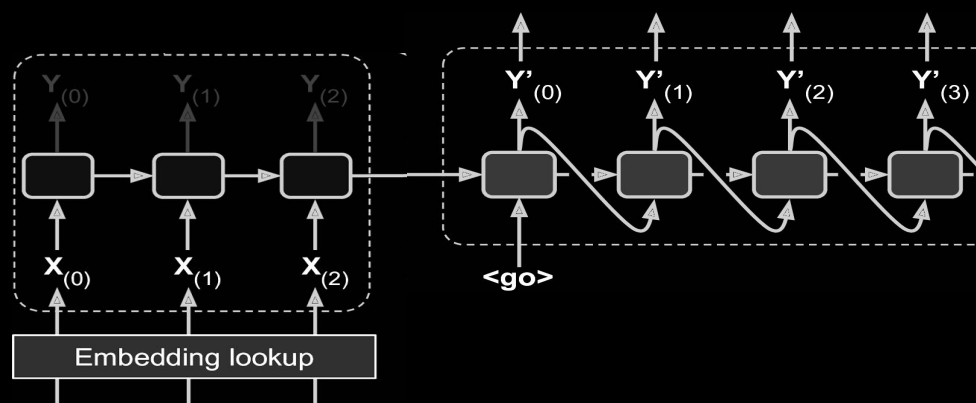
Multi-head Attention



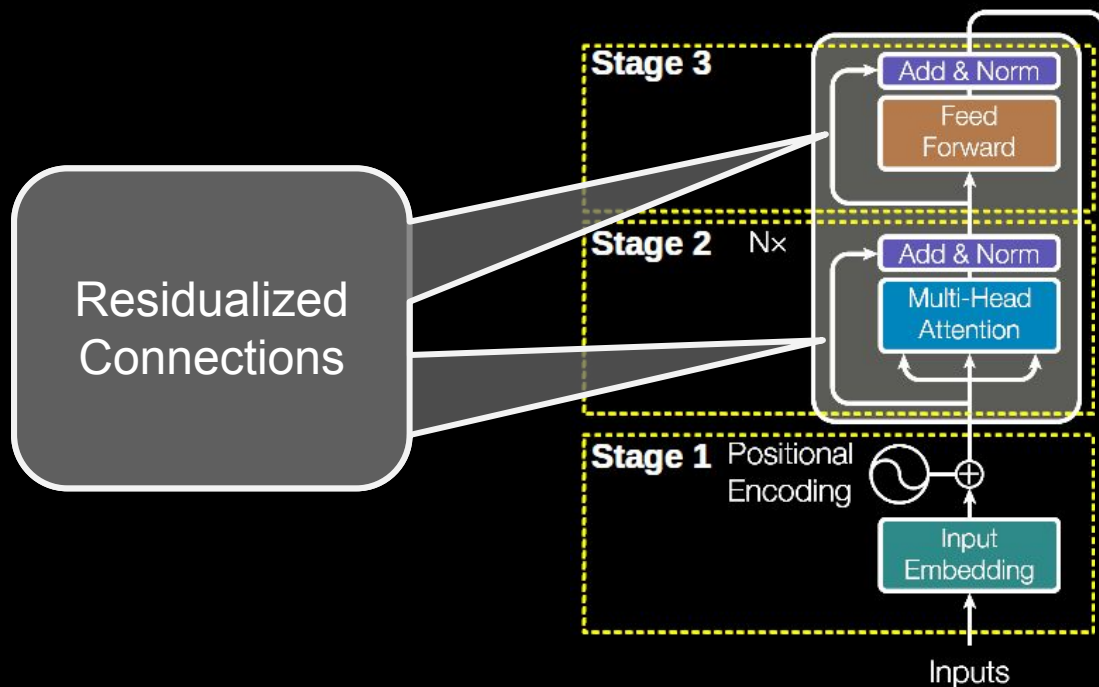
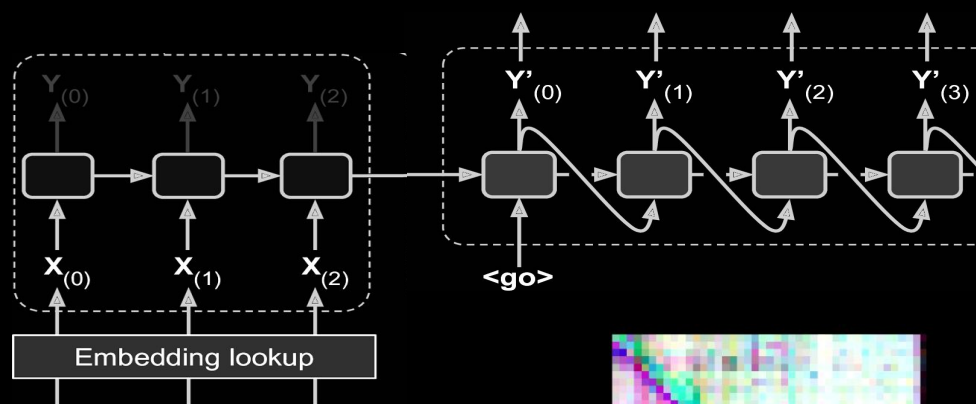
Transformer for Encoder-Decoder



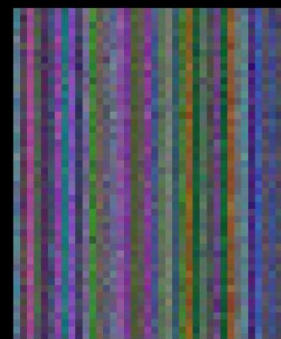
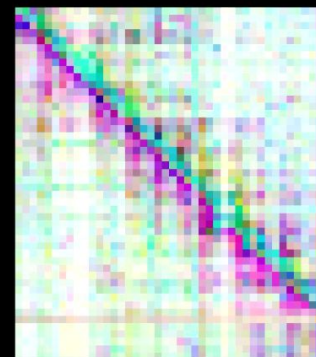
Transformer for Encoder-Decoder



Transformer for Encoder-Decoder

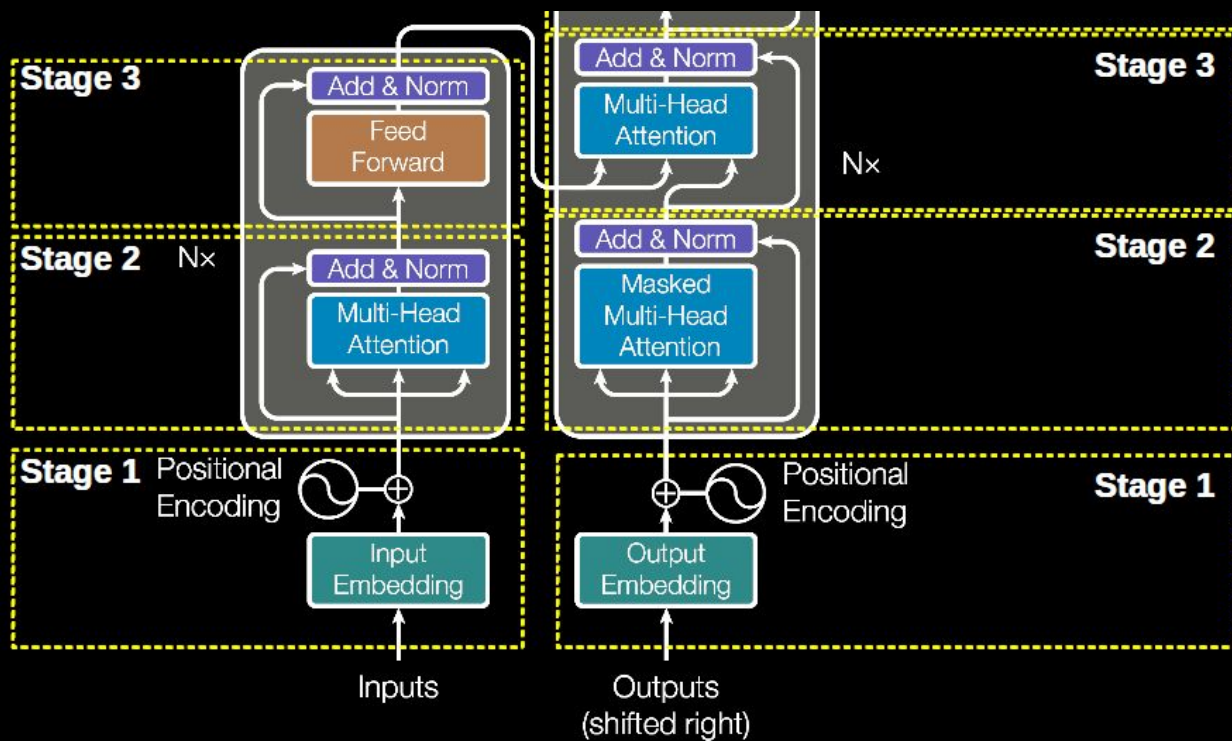
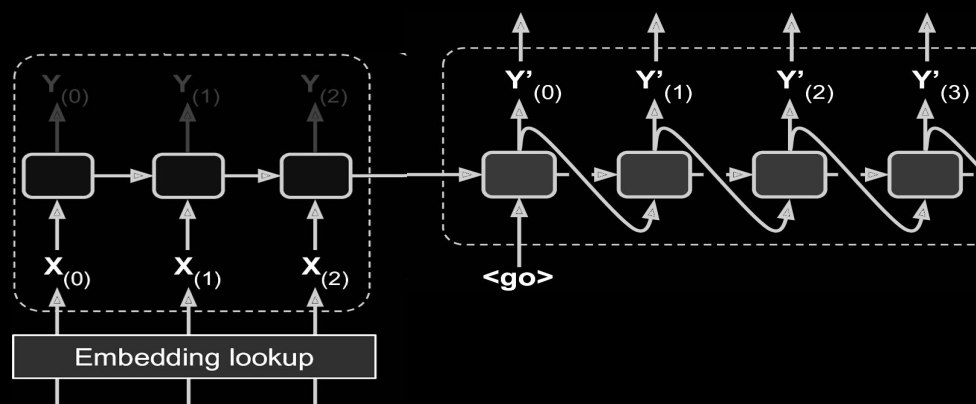


residuals enable positional information to be passed along

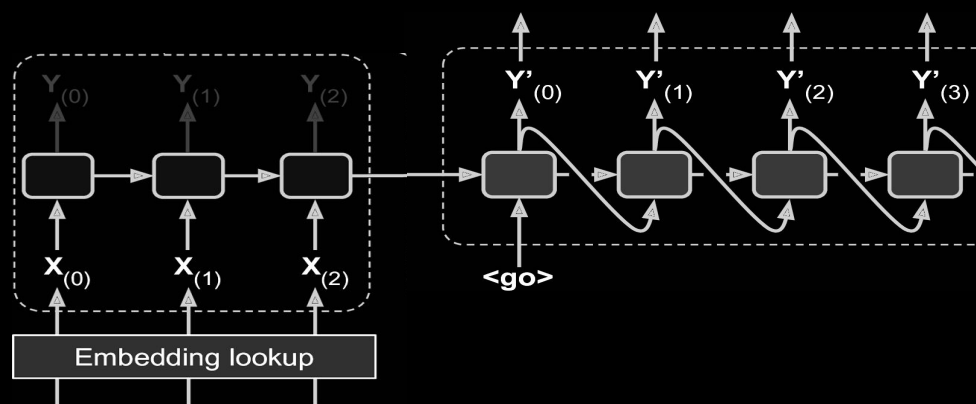


Without residuals

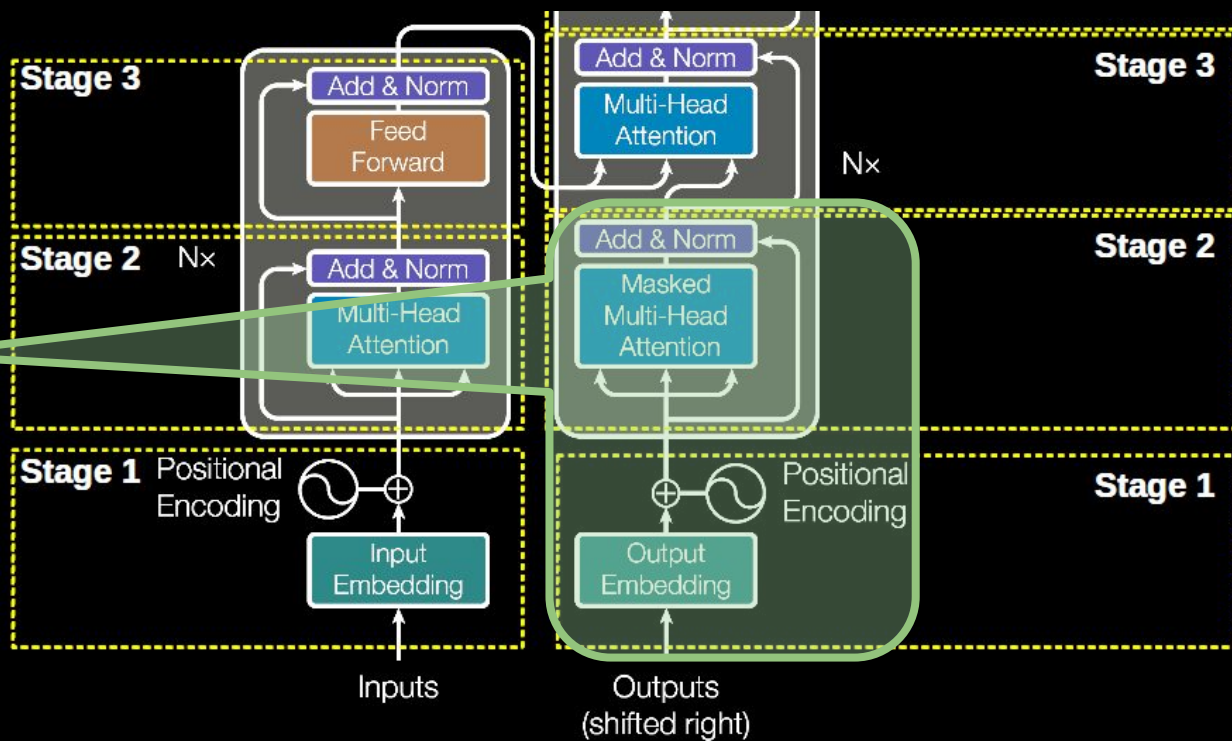
Transformer for Encoder-Decoder



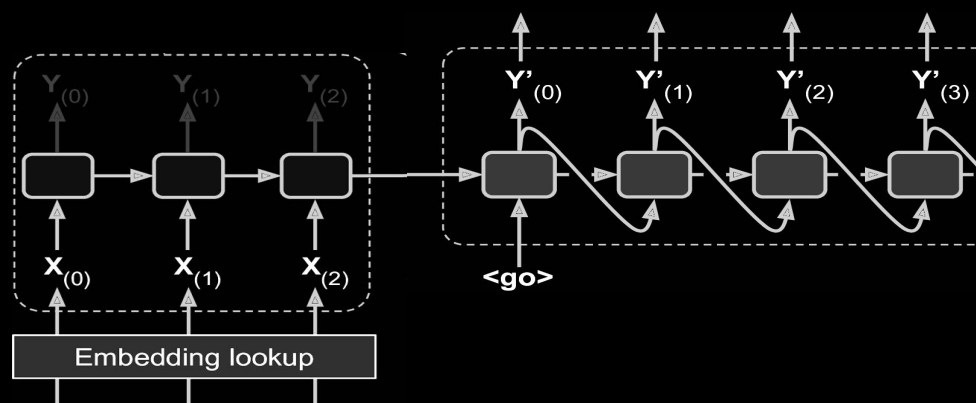
Transformer for Encoder-Decoder



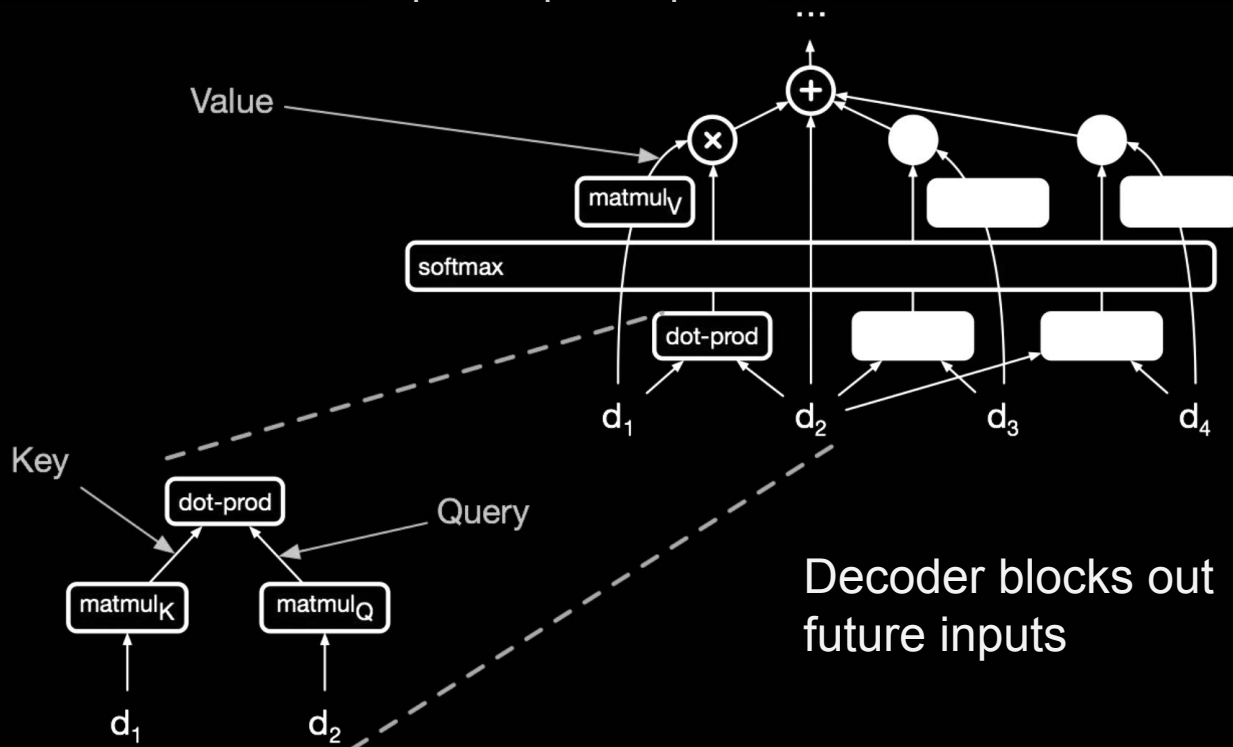
essentially, a language model



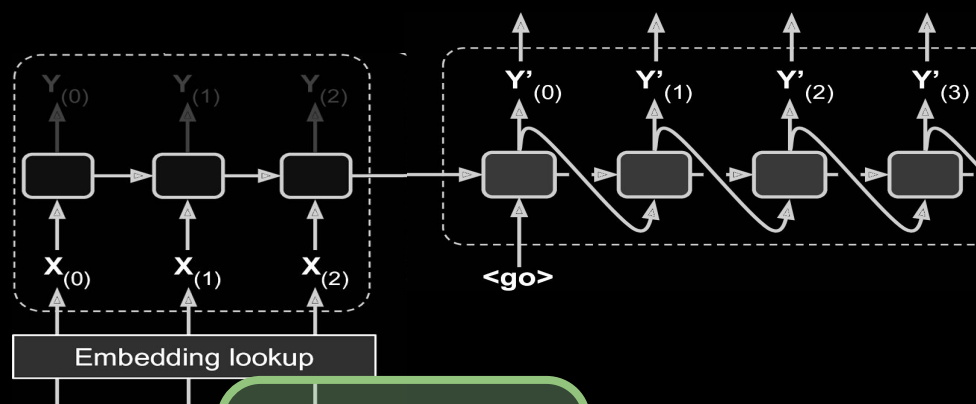
Transformer for Encoder-Decoder



essentially, a language model

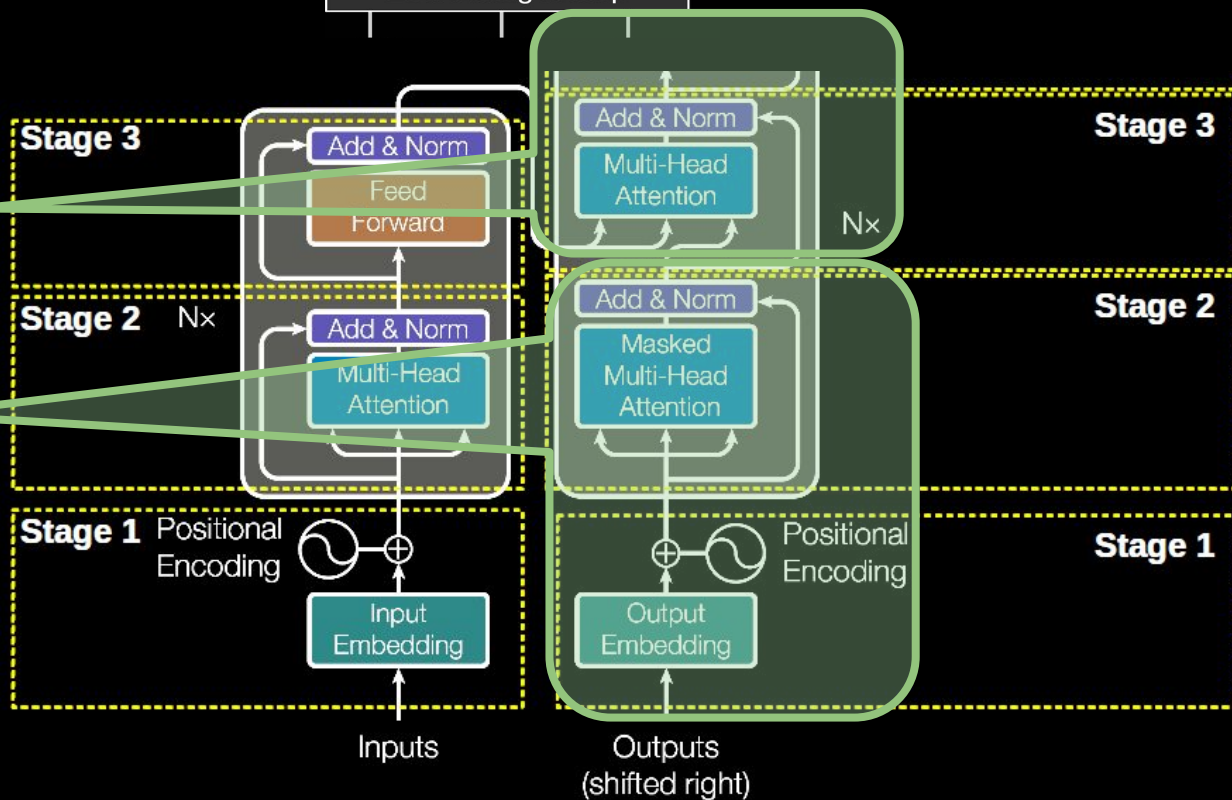


Transformer for Encoder-Decoder

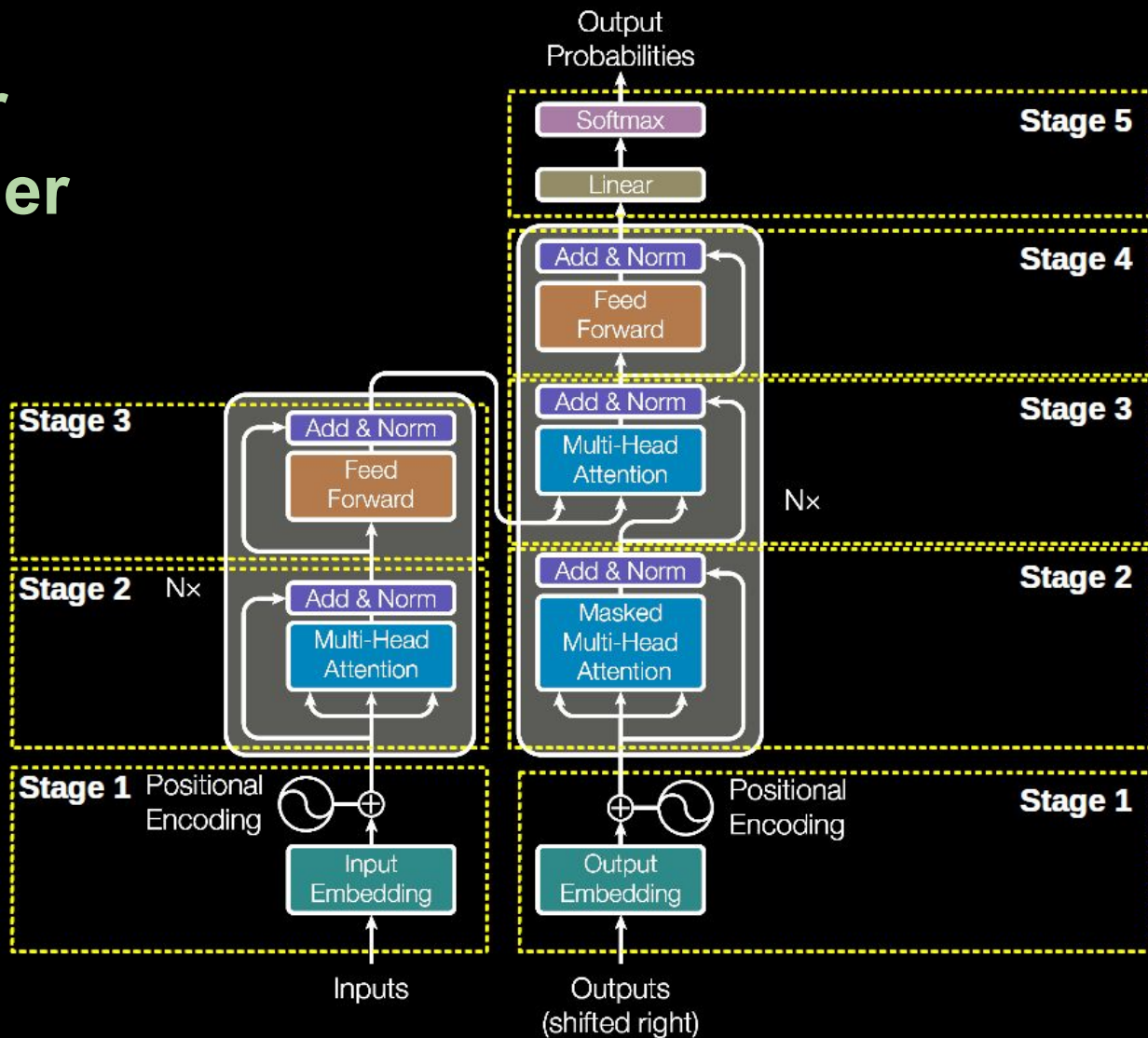


Add conditioning of the LM based on the encoder

essentially, a language model



Transformer for Encoder-Decoder



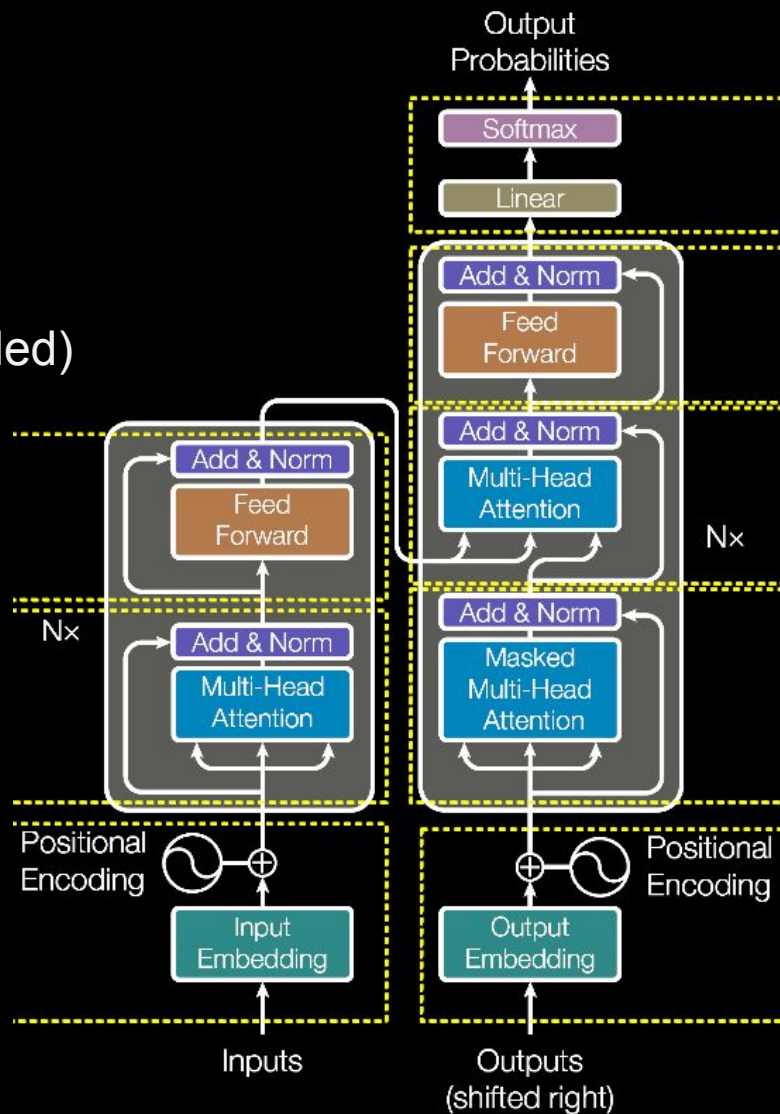
Transformer (as of 2017)

“WMT-2014” Data Set. BLEU scores:

| | EN-DE | EN-FR |
|--------------|-------------|-------------|
| GNMT (orig) | 24.6 | 39.9 |
| ConvSeq2Seq | 25.2 | 40.5 |
| Transformer* | 28.4 | 41.8 |

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



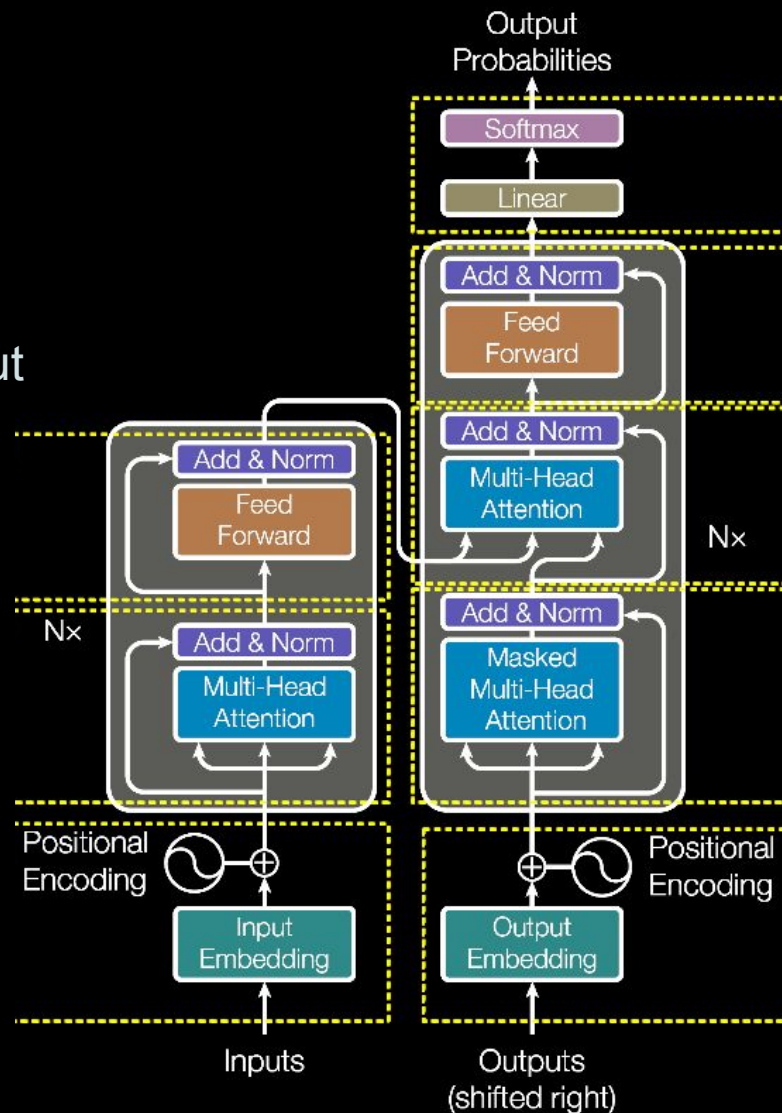
Transformer

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)

- 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

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.

She saw the man on the hill with the telescope.

She [mask] the man on the hill [mask] the telescope.

Mask 1 in 7 words:

- Too few: expensive, less robust
- Too many: not enough context

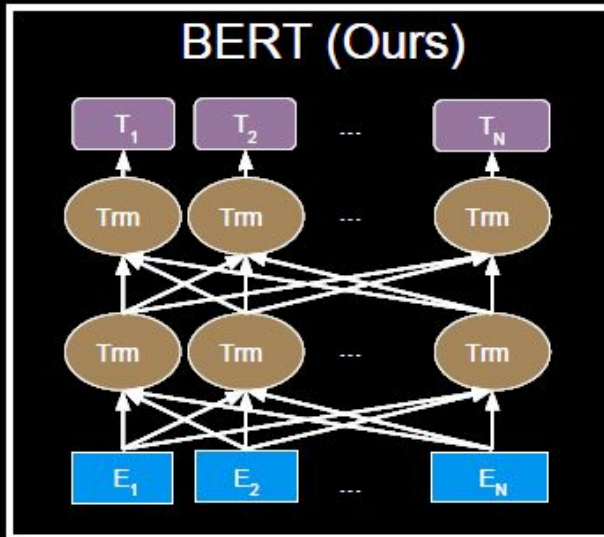
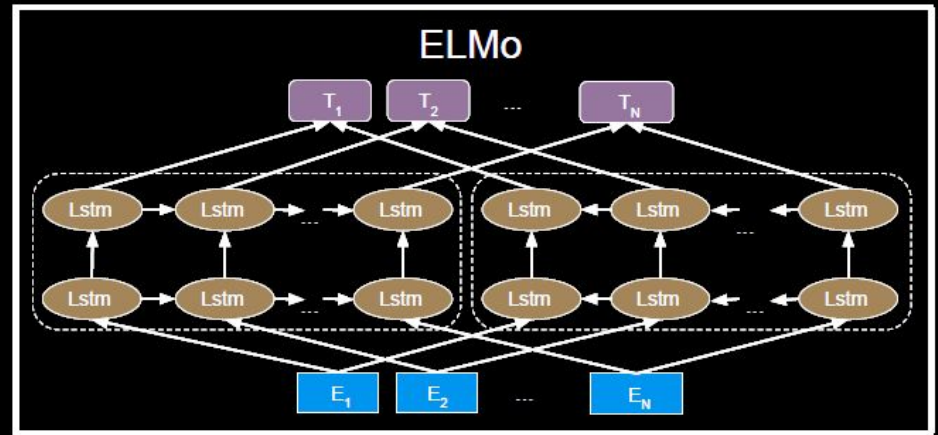
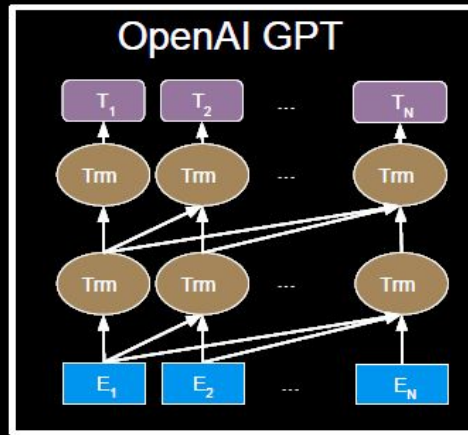
BERT

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.**
 - BERT-Base, Cased:
12-layer, 768-hidden, 12-heads , 110M parameters
 - BERT-Large, Cased:
24-layer, 1024-hidden, 16-heads, 340M parameters
 - BERT-Base, Multilingual Cased:
104 languages, 12-layer, 768-hidden, 12-heads, 110M parameters

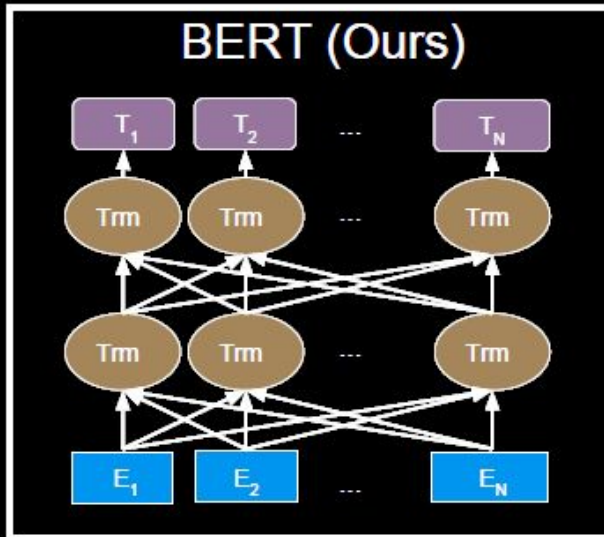
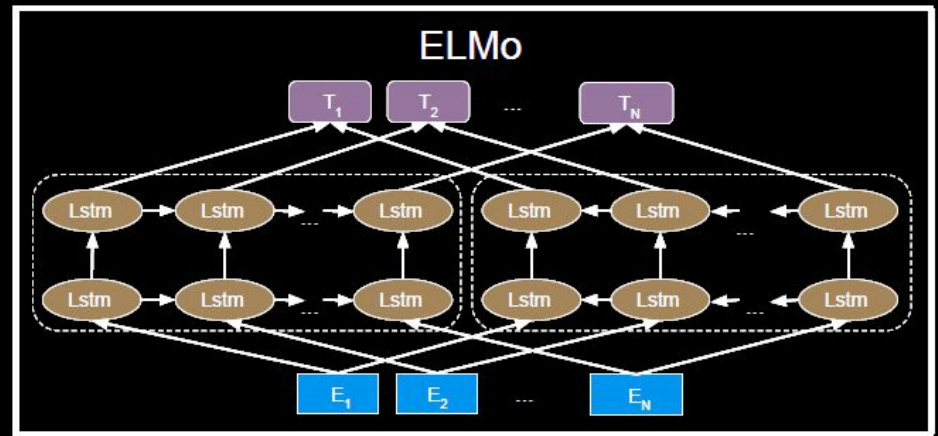
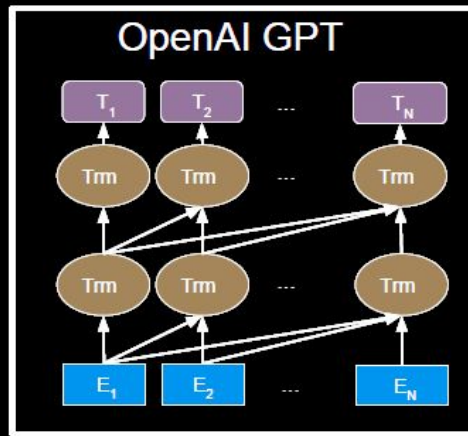
BERT



Differences from previous state of the art:

- Bidirectional transformer (through masking)
- Directions jointly trained at once.

BERT



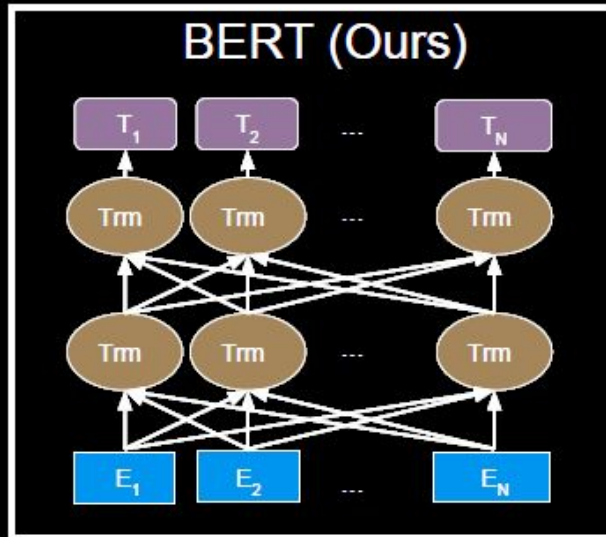
Differences from previous state of the art:

- Bidirectional transformer (through masking)
- Directions jointly trained at once.
- **Capture sentence-level relations**

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



Differences from previous state of the art:

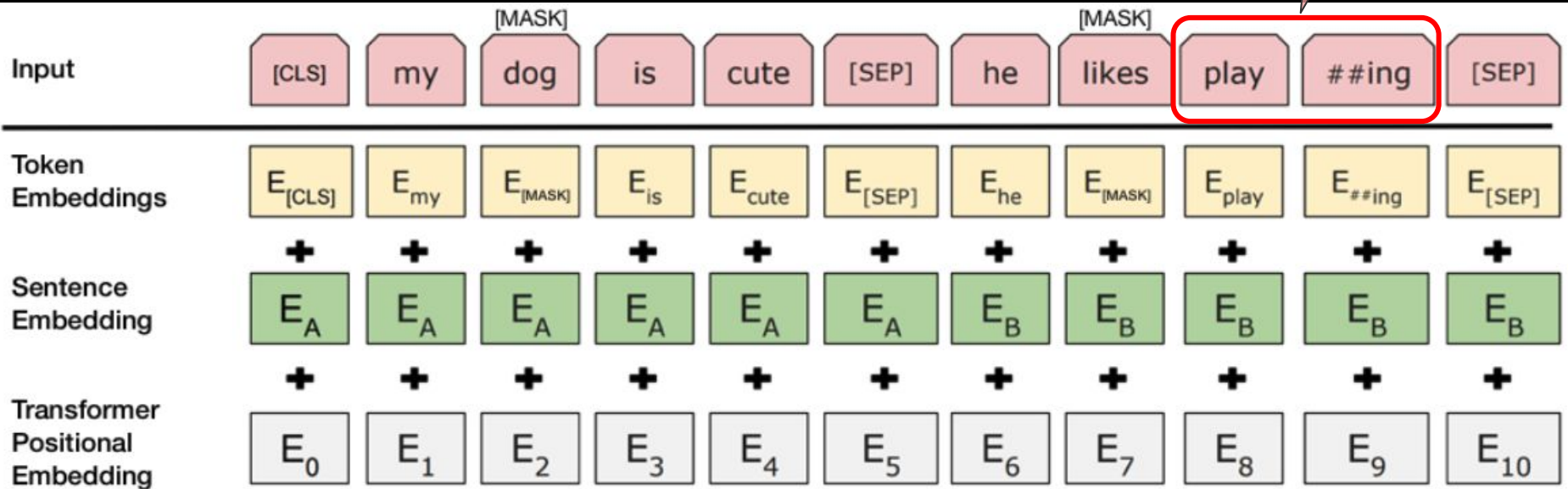
- Bidirectional transformer (through masking)
- Directions jointly trained at once.
- **Capture sentence-level relations**

BERT

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

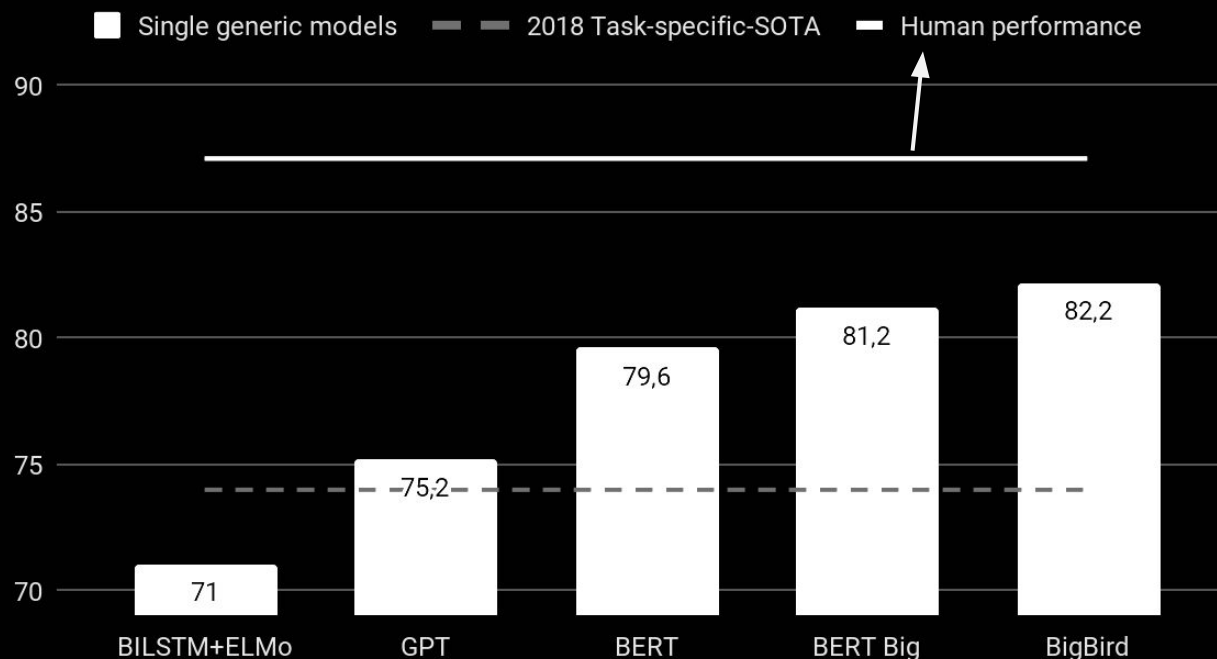
Sentence A = The man went to the store.
Sentence B = Penguins are flightless.
Label = NotNextSentence

tokenize into "word pieces"



BERT Performance: e.g. Question Answering

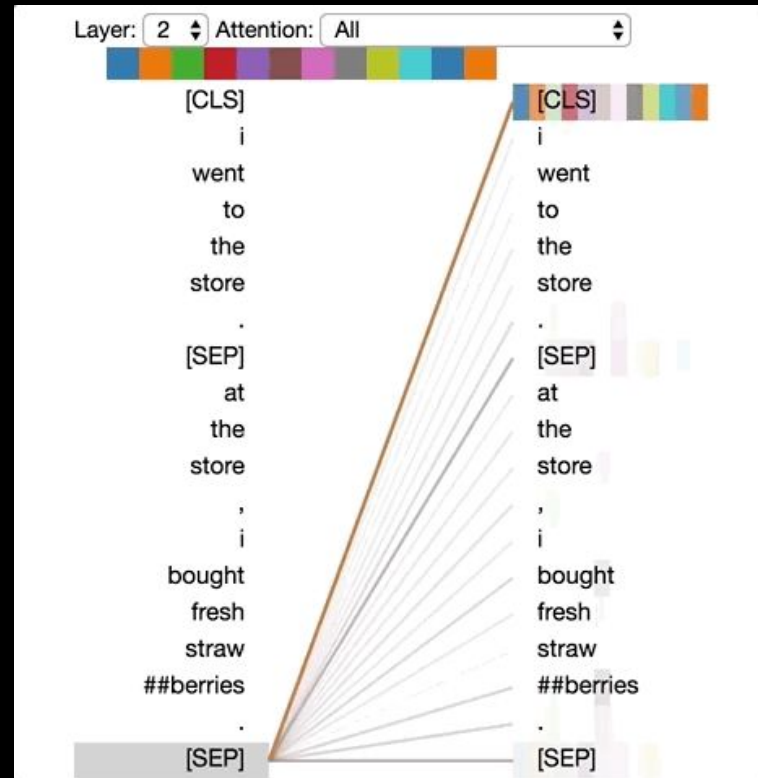
GLUE scores evolution over 2018-2019



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

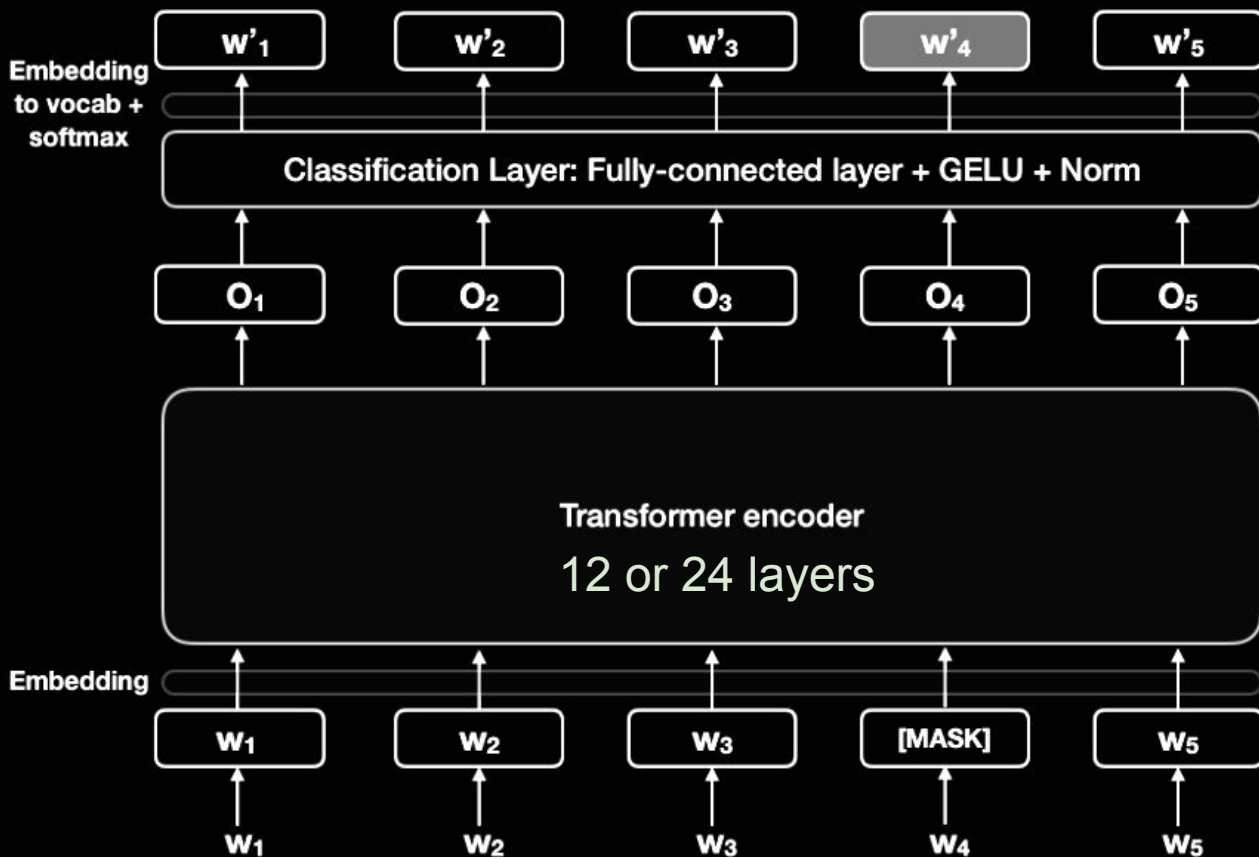
Bert: Attention by Layers

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

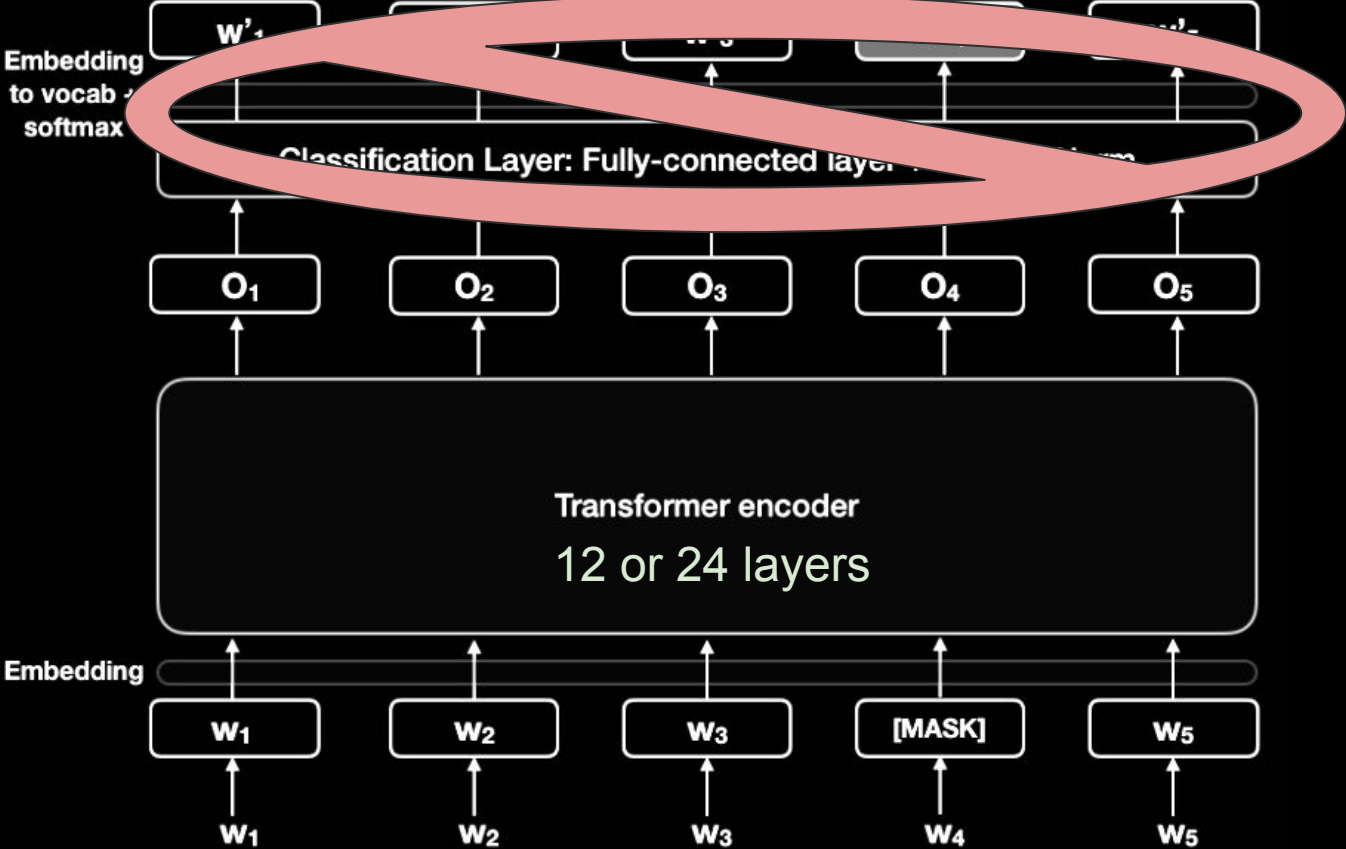


(Vig, 2019)

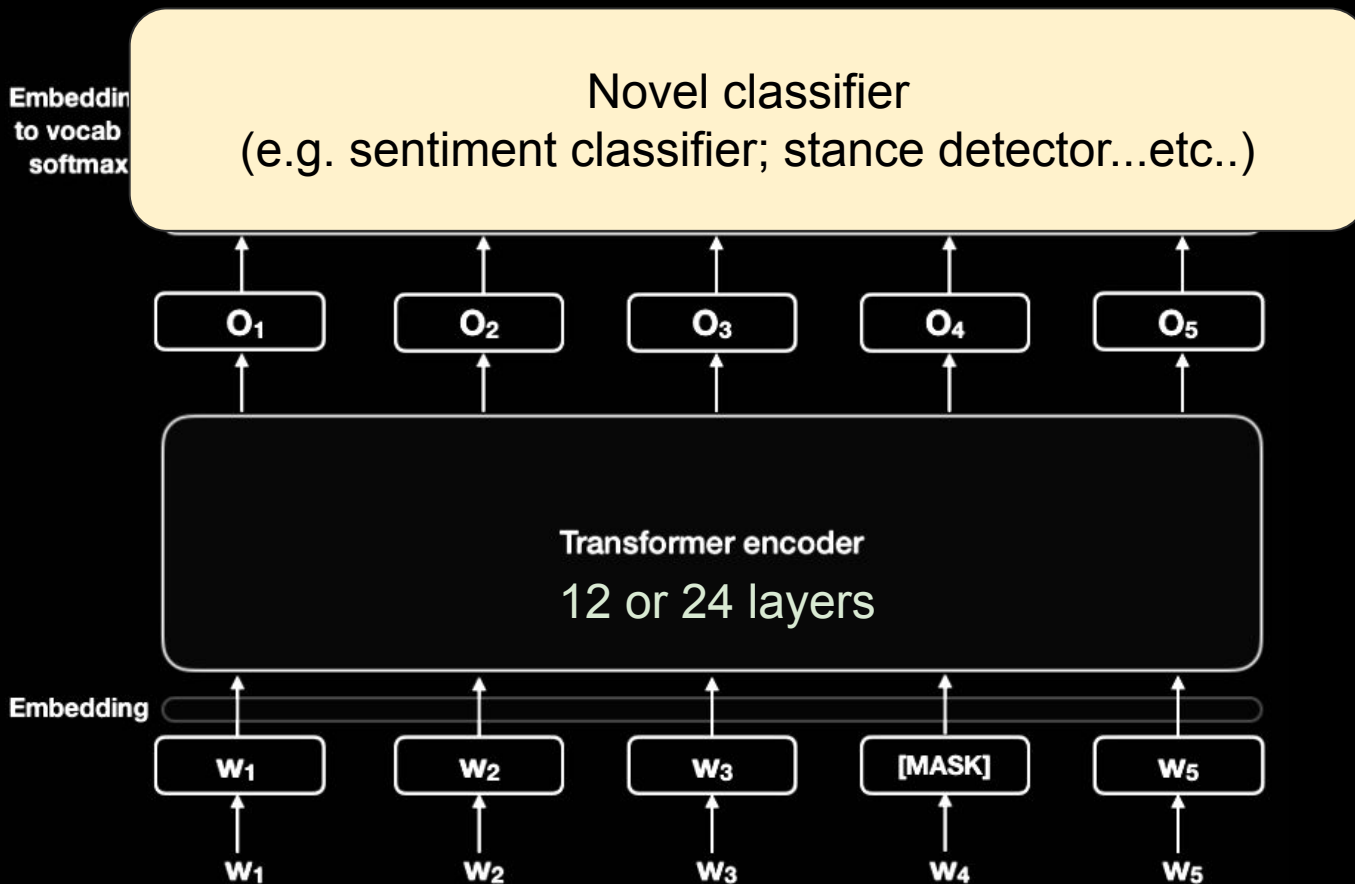
BERT: Pre-training; Fine-tuning



BERT: Pre-training; Fine-tuning

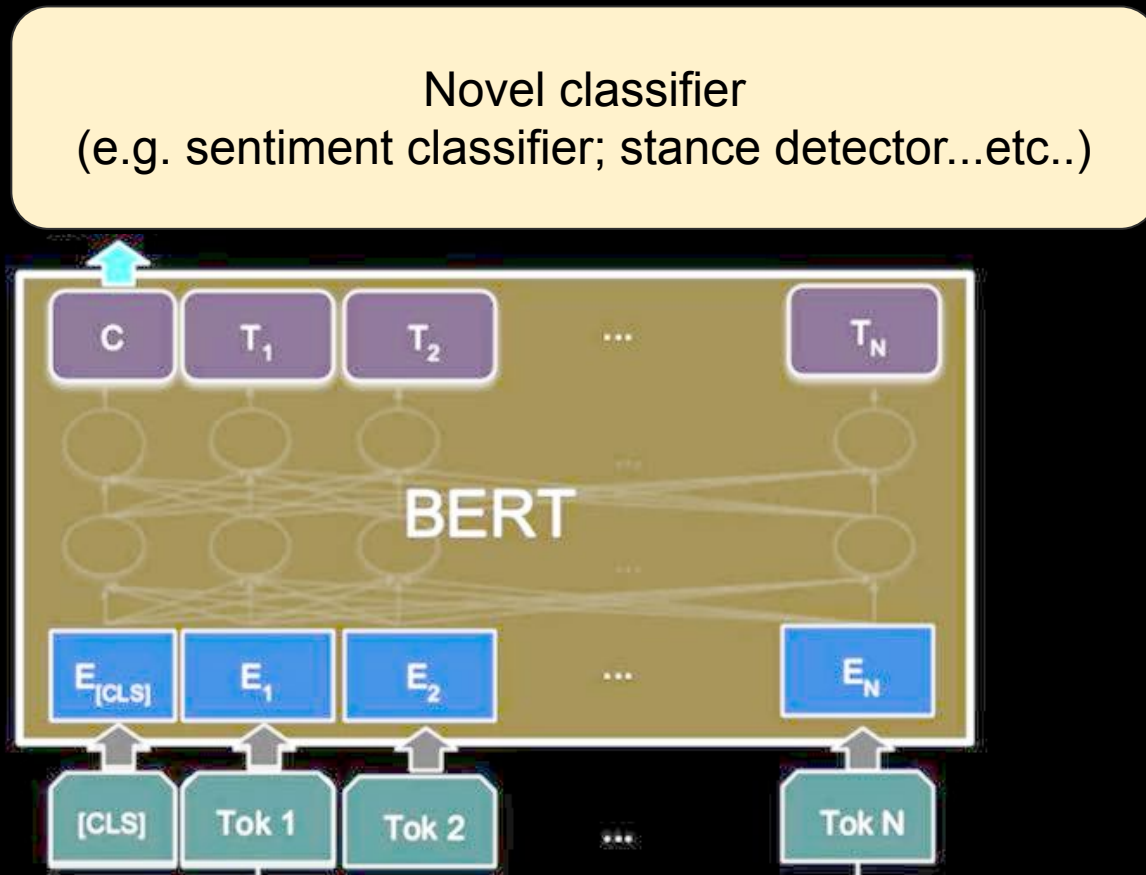


BERT: Pre-training; Fine-tuning



BERT: Pre-training; Fine-tuning

[CLS] vector at start is supposed to capture meaning of whole sequence.



BERT: Pre-training; Fine-tuning

[CLS] vector at start is supposed to capture meaning of whole sequence.

Average of top layer (or second to top) also often used.

