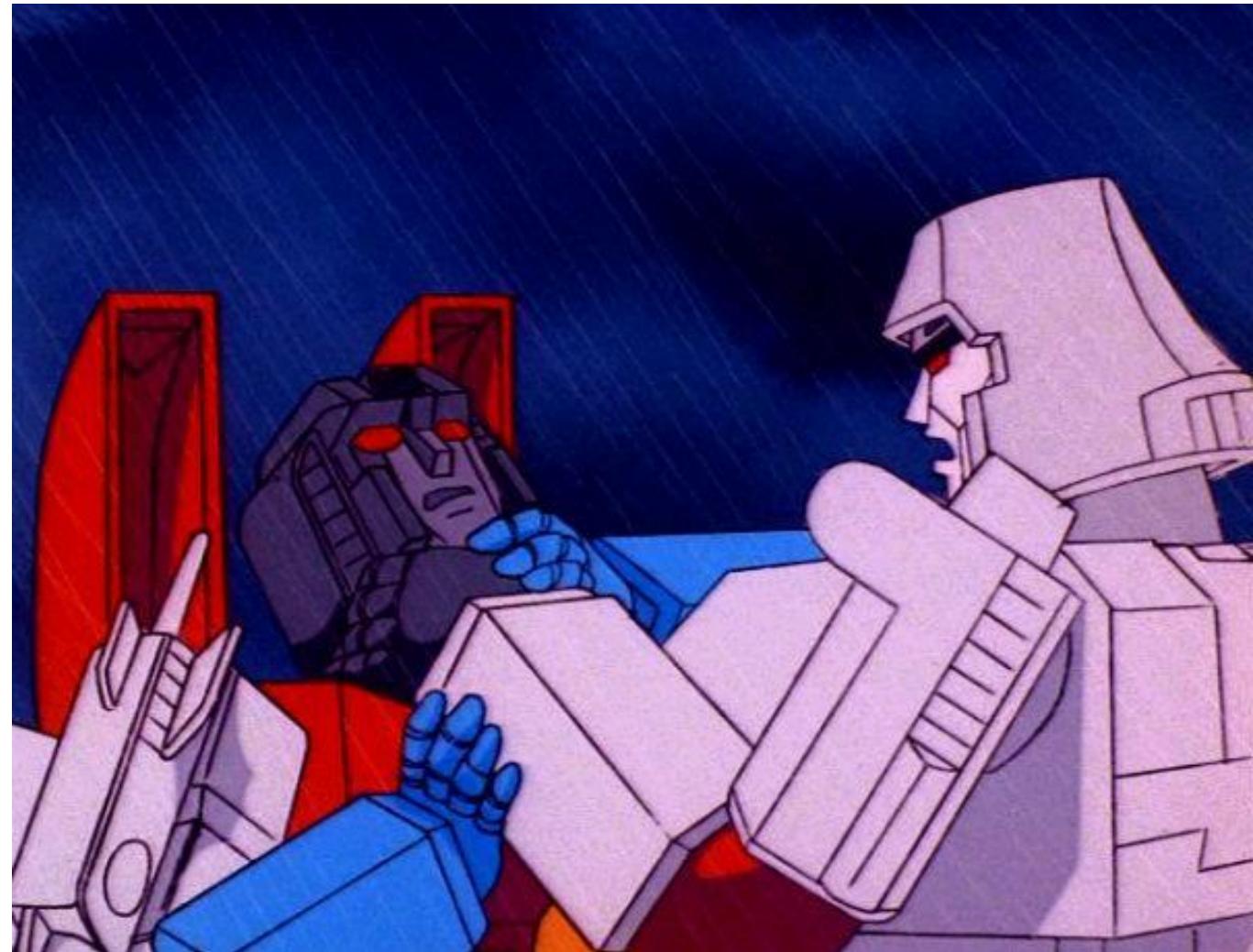


Transformers

Lecture 6





Language Modelling Task

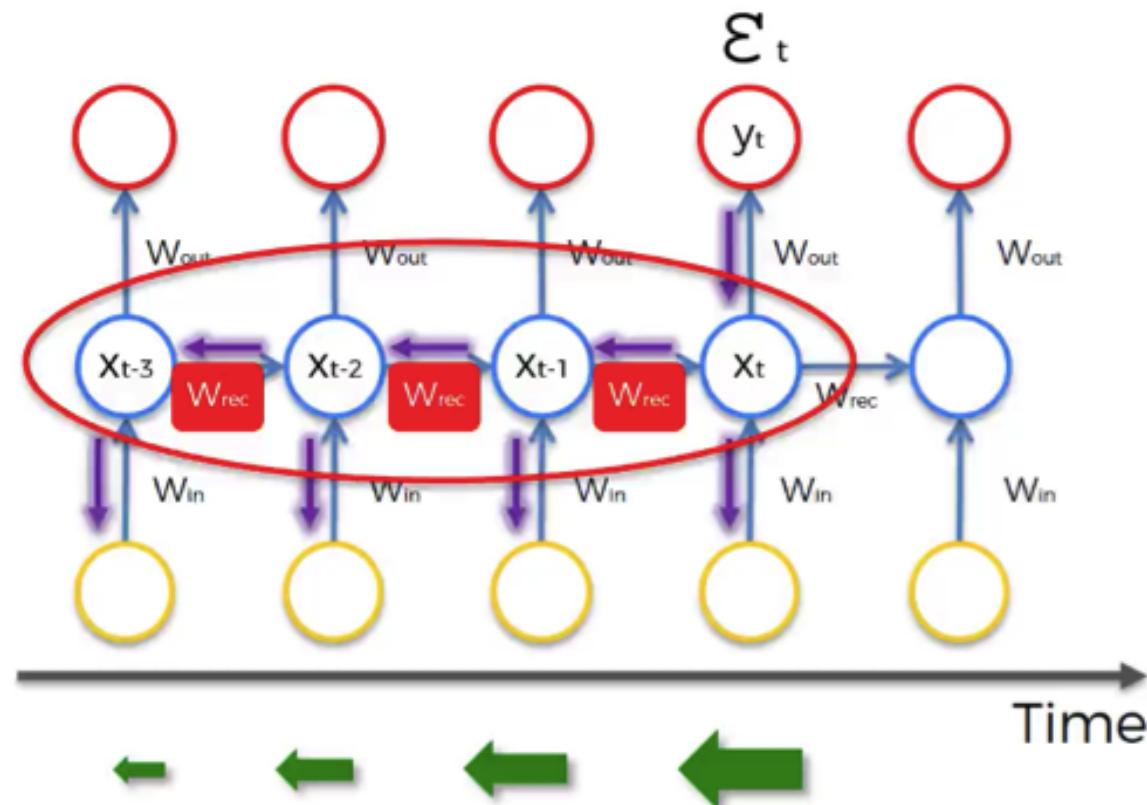
- Final goal: predict/estimate the probability of a sequence

Probability(*Some sentence over here.*)

- Actual task:
 - Predict the next word
 - MLM
- In a perfect world:
 - The RNN hidden states should be able capture all contextual information
 - Right?



The Vanishing Gradient Problem



$$\frac{\partial \mathcal{E}}{\partial \theta} = \sum_{1 \leq t \leq T} \frac{\partial \mathcal{E}_t}{\partial \theta} \quad (3)$$

$$\frac{\partial \mathcal{E}_t}{\partial \theta} = \sum_{1 \leq k \leq t} \left(\frac{\partial \mathcal{E}_t}{\partial \mathbf{x}_t} \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} \frac{\partial^+ \mathbf{x}_k}{\partial \theta} \right) \quad (4)$$

$$\frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} = \prod_{t \geq i > k} \frac{\partial \mathbf{x}_i}{\partial \mathbf{x}_{i-1}} = \prod_{t \geq i > k} \mathbf{W}_{rec}^T diag(\sigma'(\mathbf{x}_{i-1})) \quad (5)$$

$W_{rec} \sim \text{small} \rightarrow$ Vanishing
 $W_{rec} \sim \text{large} \rightarrow$ Exploding

Formula Source: Razvan Pascanu et al. (2013)

Again, we need more! What of larger semantic units?

- How can we know when larger units are similar in meaning?
 - *CTV News*: Poilievre-led attempt to bring down Trudeau minority over carbon tax fails.
 - *CBC News*: Liberals survive non-confidence vote on carbon tax with Bloc, NDP backing.
 - *The Beaverton*: Co-worker that everyone hates surprised he can't get colleagues to do what he wants.



NATIONAL - 2 WEEKS AGO

Co-worker that everyone hates surprised he can't get colleagues to do what he wants

OTTAWA – Local man Pierre Poilievre, an employee at an Ottawa small business named the House of Commons, was surprised that none of the colleagues who despise him were willing to support hi...



SHARE

RNN & next word prediction: Not good compositional representation

- Next word prediction:

$$P(t_i | t_1, t_2, \dots, t_{i-1})$$

- The hidden state i is encoding information of everything from the beginning (index 0) to the very end (index i).
- We want some bigger semantic units
 - Poilievre-led attempt to **bring down Trudeau minority over carbon tax** fails.
- Some hacks may work, but not really

Again, we need more! What of larger semantic units?

- How can we know when larger units are similar in meaning?
 - *CTV News*: Poilievre-led attempt to bring down Trudeau minority over carbon tax fails.
 - *CBC News*: Liberals survive non-confidence vote on carbon tax with Bloc, NDP backing.
 - *The Beaverton*: Co-worker that everyone hates surprised he can't get colleagues to do what he wants.

People interpret the meaning of larger text units – entities, descriptive terms, facts, arguments, stories – by **semantic composition** of smaller elements.

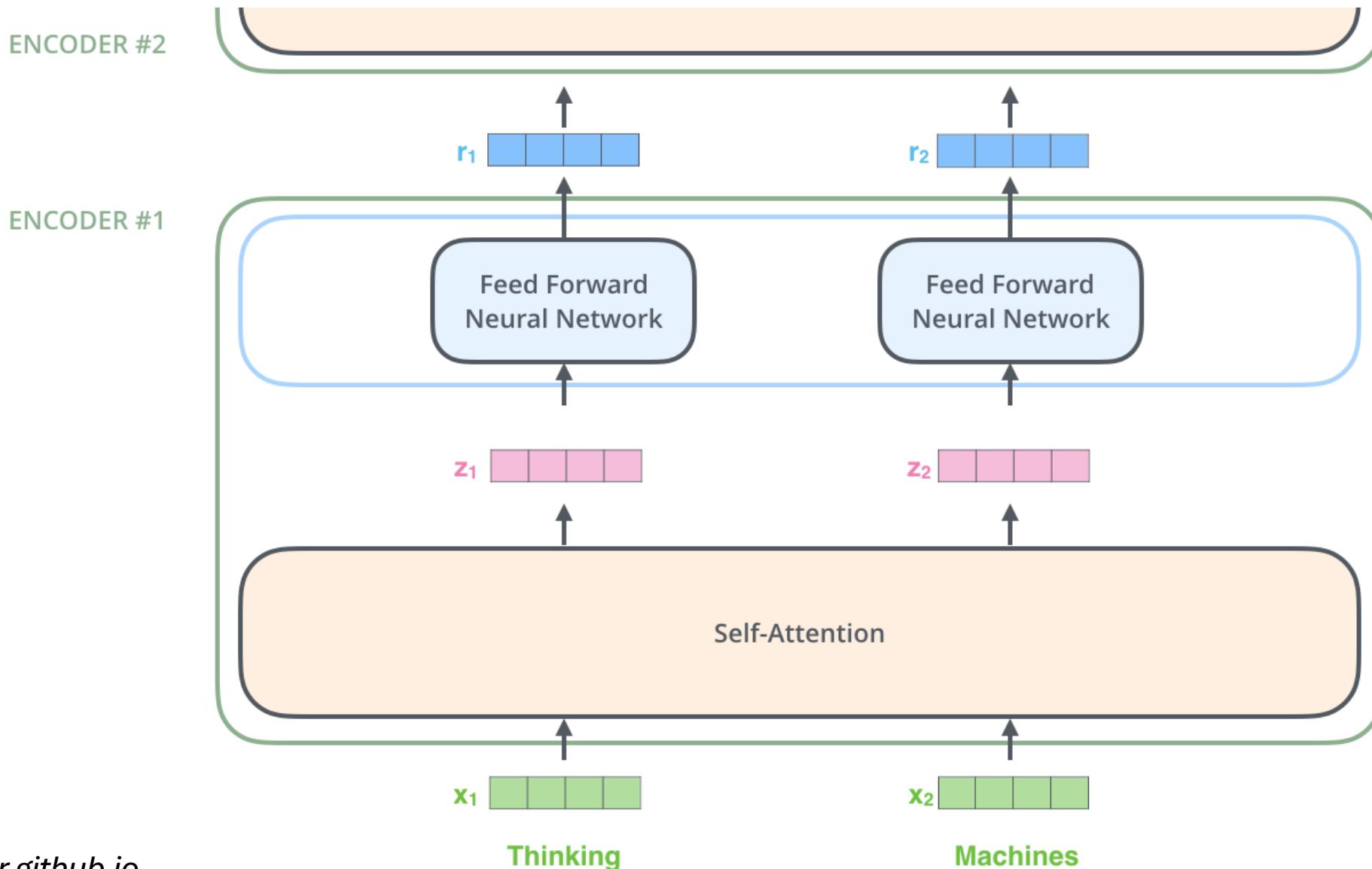
Beyond Word2Vec and RNN-LMs

- Mitigation #1:
 - Long short-term memory (LSTM).
 - Won't cover in this class, but covered DL4NLP.
- Mitigation #2:
 - Attention Mechanism -> Transformer

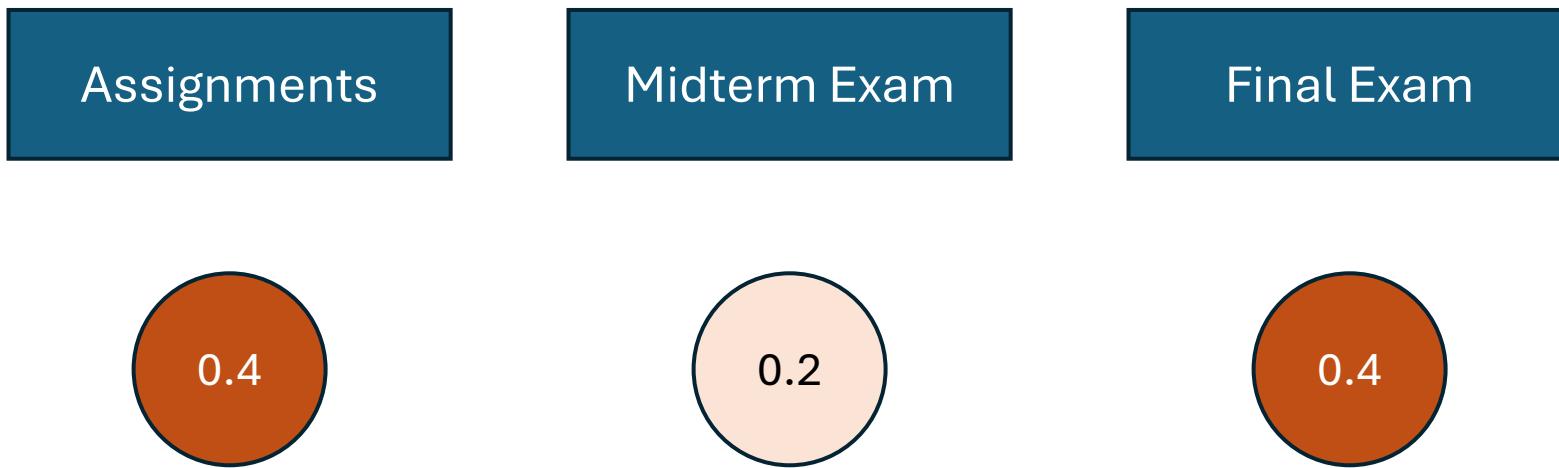
LSTM

- (Won't appear on your exam)
- Basic idea:
 - Reset the hidden state once a while.
 - When to know how to reset? Train a separate model (actually, a part of the model) to do it.

Transformers

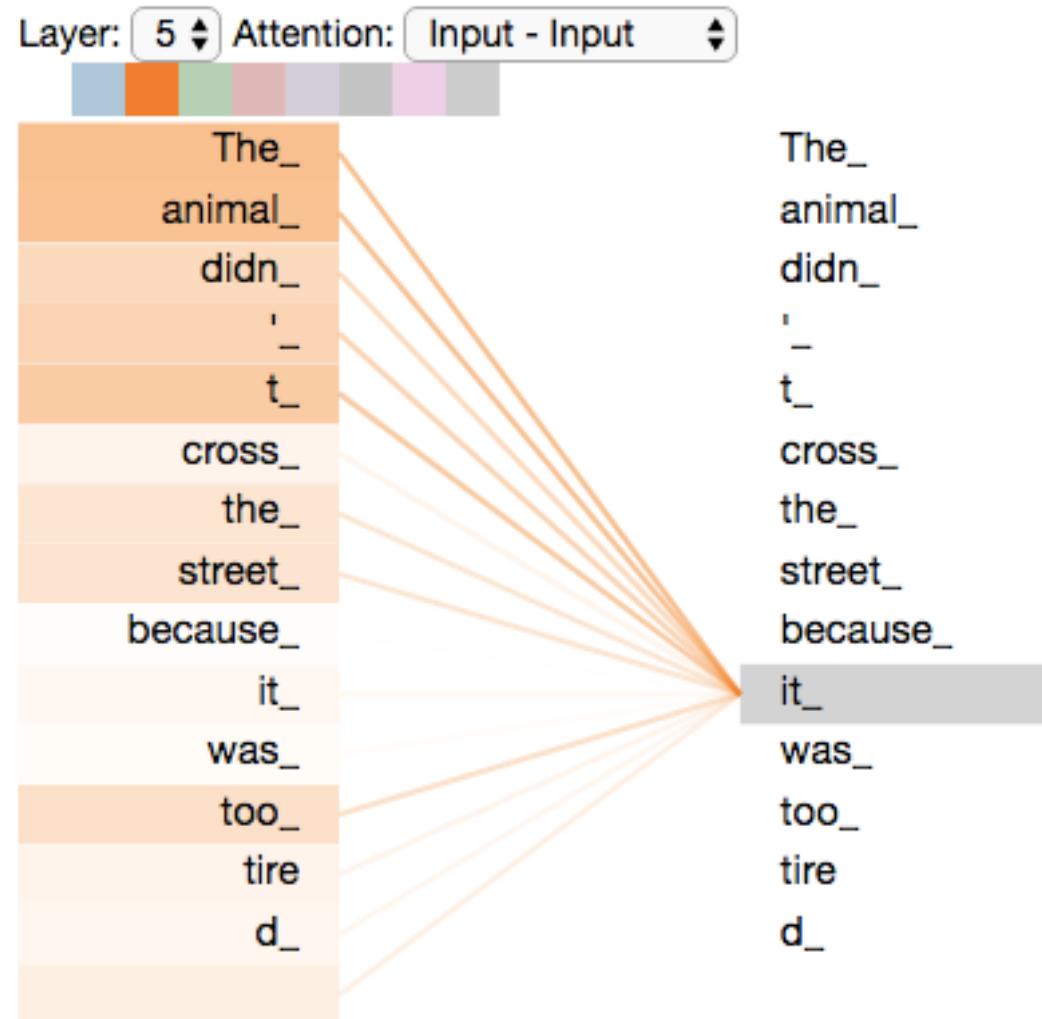


Attention Mechanism

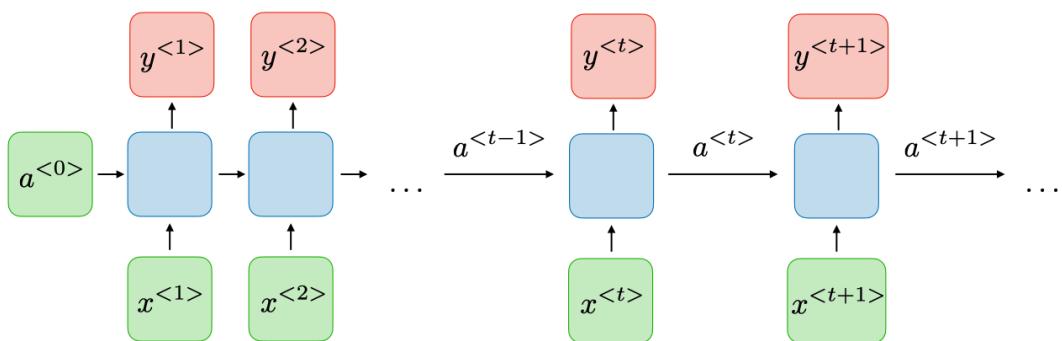


$$\text{Total} = \boxed{0.75} * A + \boxed{0.1} * Q + \boxed{0.15} * E$$

Attention Mechanism

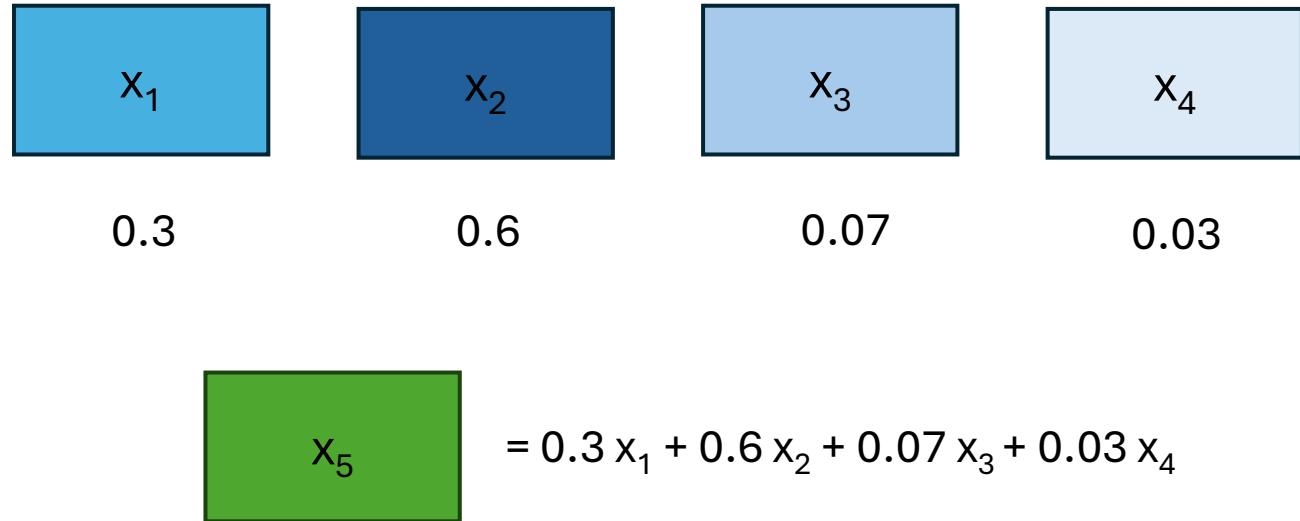


RNN vs. Attention



Recurrent neural network.
Almost “recursive.”

$$h_i = f(h_{i-1} + x_i)$$



$$x_i = \sum_j a_j x_j$$

????

$$h_i = f(x_i)$$

Recall GloVe

Encoding meaning in vectors

- How can we capture ratios of co-occurrence probability components in a word vector space?
- Solution:

- Log-bilinear model:

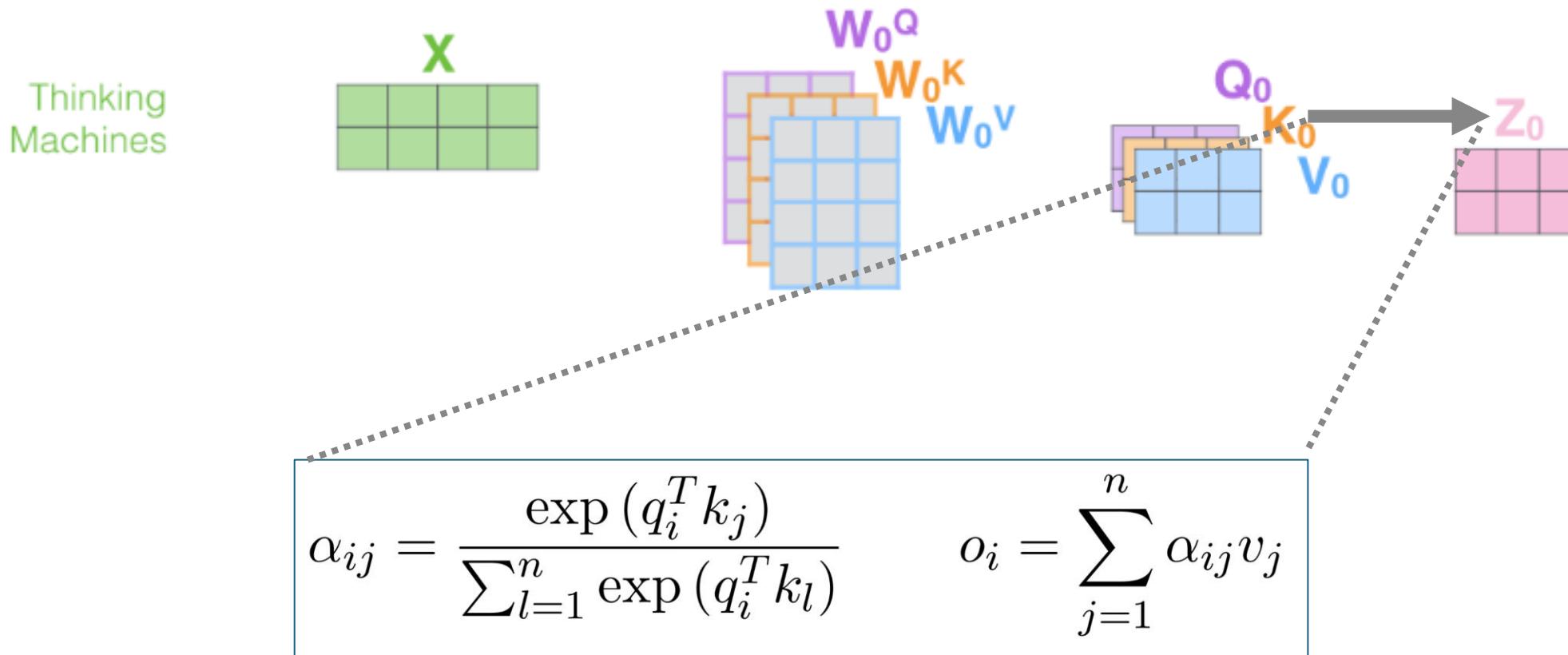
$$w_i \cdot w_j = \log P(i|j)$$

- with vector differences:

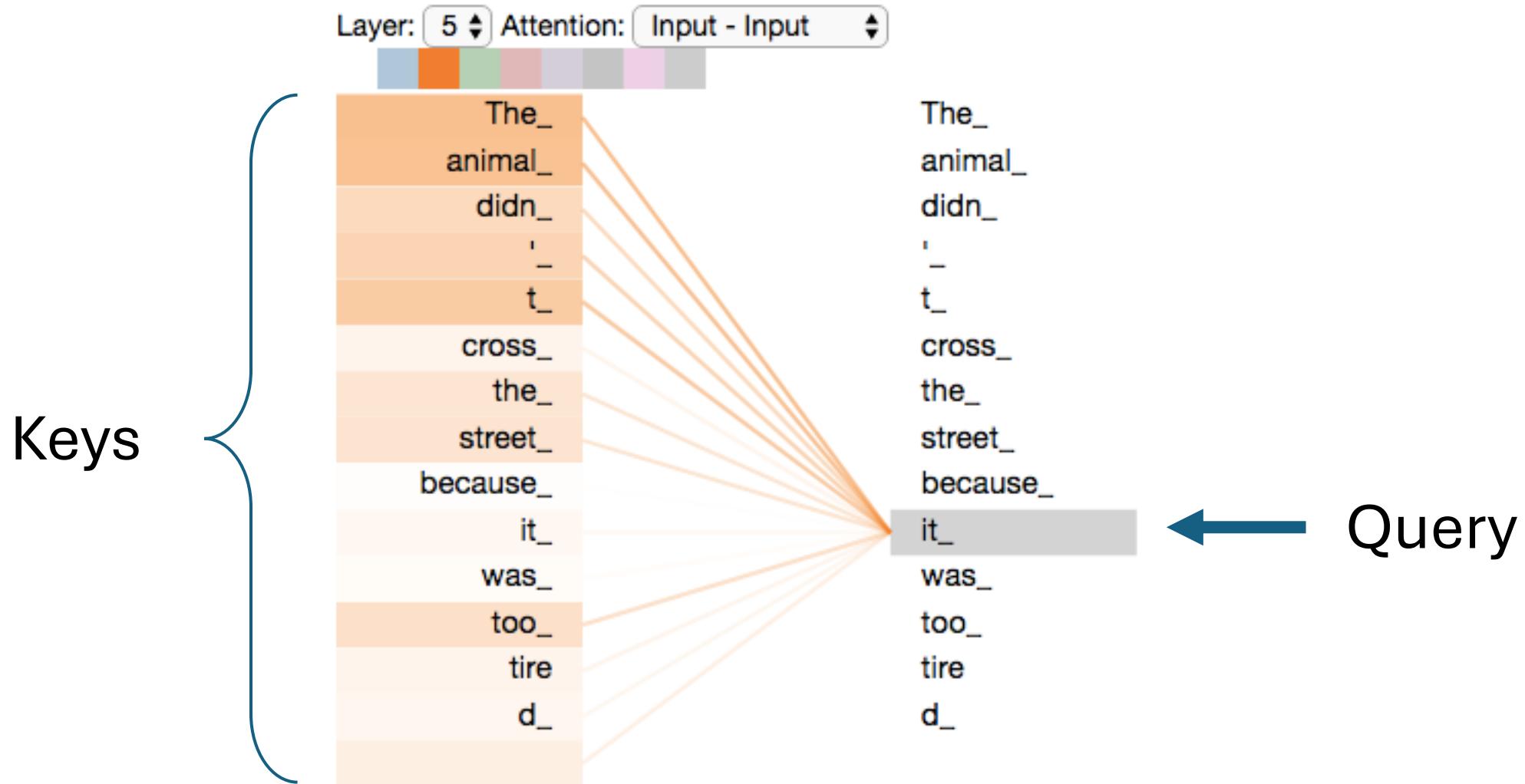
$$w_x \cdot (w_a - w_b) = \log \frac{P(x|a)}{P(x|b)}$$

Pennington et al. (2014)

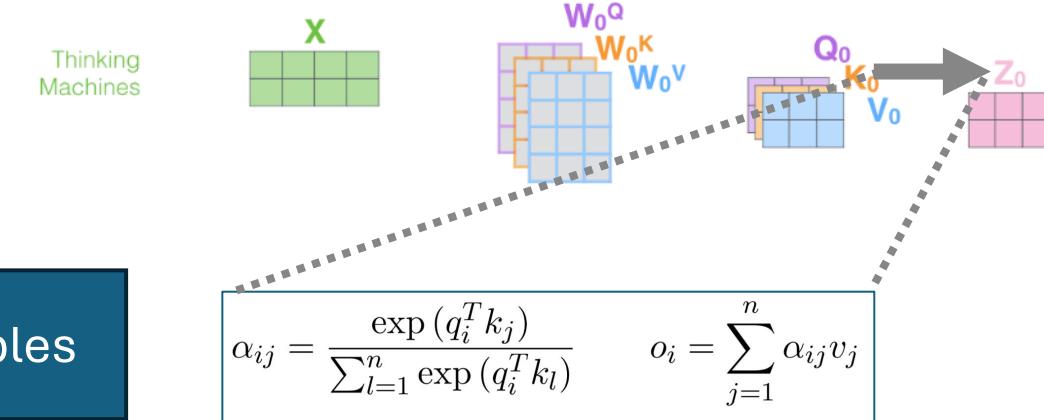
Self Attention



Attention Mechanism



Self Attention



I like eating apples

What is $o(\text{like})$? i.e., the attention score from like to the sentence.

query

Q_2
“like”

key

| | | | |
|-----------|--------------|----------------|----------------|
| K1 “I” | K2 “like” | K3 “eating” | K4 “apples” |
|-----------|--------------|----------------|----------------|

attention
weights*

| | | | |
|------|------|------|-----|
| 14.2 | 18.1 | 10.3 | 7.9 |
|------|------|------|-----|

attention
scores* (o)

| | | | |
|-----|-----|------|------|
| 0.3 | 0.6 | 0.07 | 0.03 |
|-----|-----|------|------|

*: made-up numbers, not real.

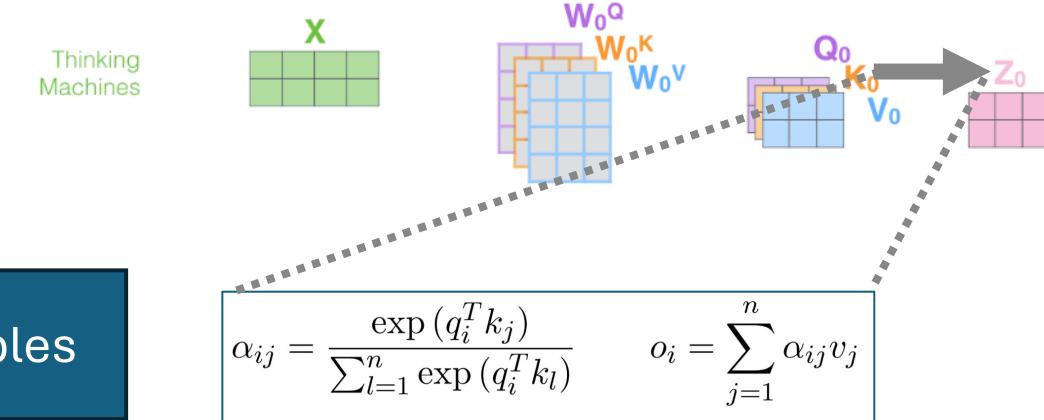
$$X \times W^Q = Q$$

$$X \times W^K = K$$

$$Q \times K^T$$

$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}} \right)$$

Self Attention



I like eating apples

What is $z(\text{like})$?

attention scores* (o) 0.3 0.6 0.07 0.03

values

| | | | |
|-----------|--------------|----------------|----------------|
| V1 “I” | V2 “like” | V3 “eating” | V4 “apples” |
|-----------|--------------|----------------|----------------|

Weighted sum

Z2
“like”

$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}} \right)$$

$X \times W^V = V$

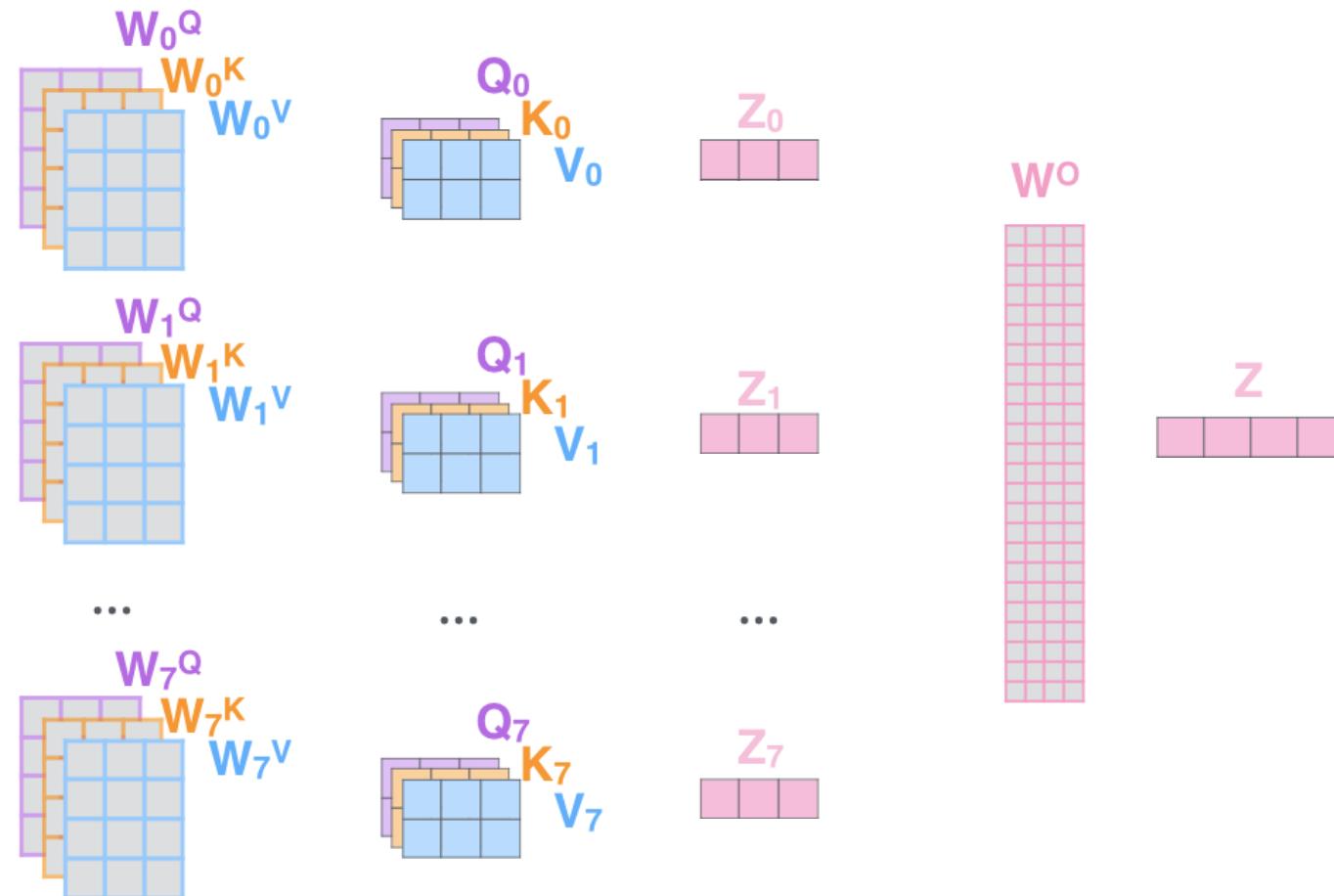
*: made-up numbers, not real.

Multi-Head Self Attention

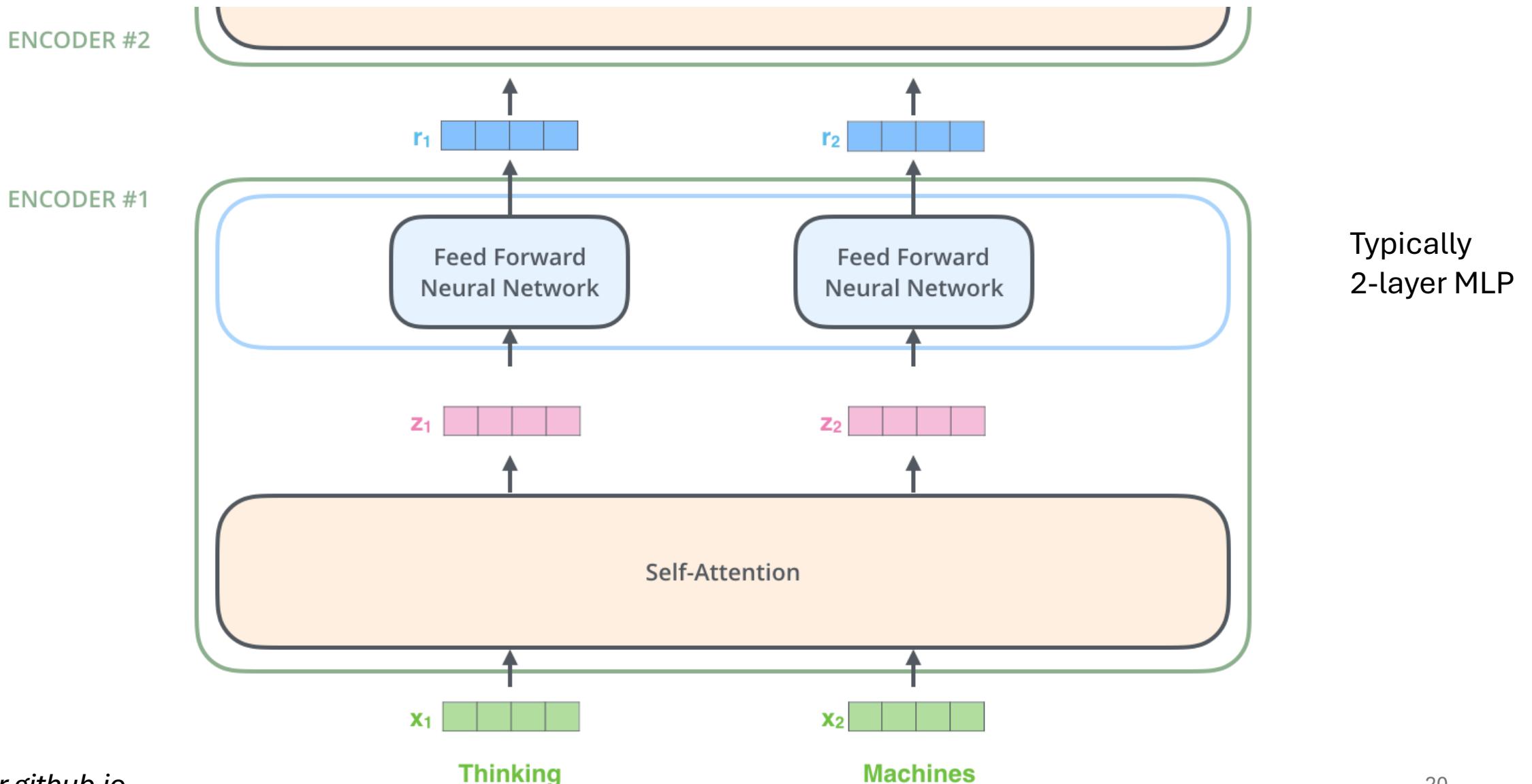
- 1) This is our input sentence* each word*
- 2) We embed each word*
- 3) Split into 8 heads. We multiply X or R with weight matrices
- 4) Calculate attention using the resulting $Q/K/V$ matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer



* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



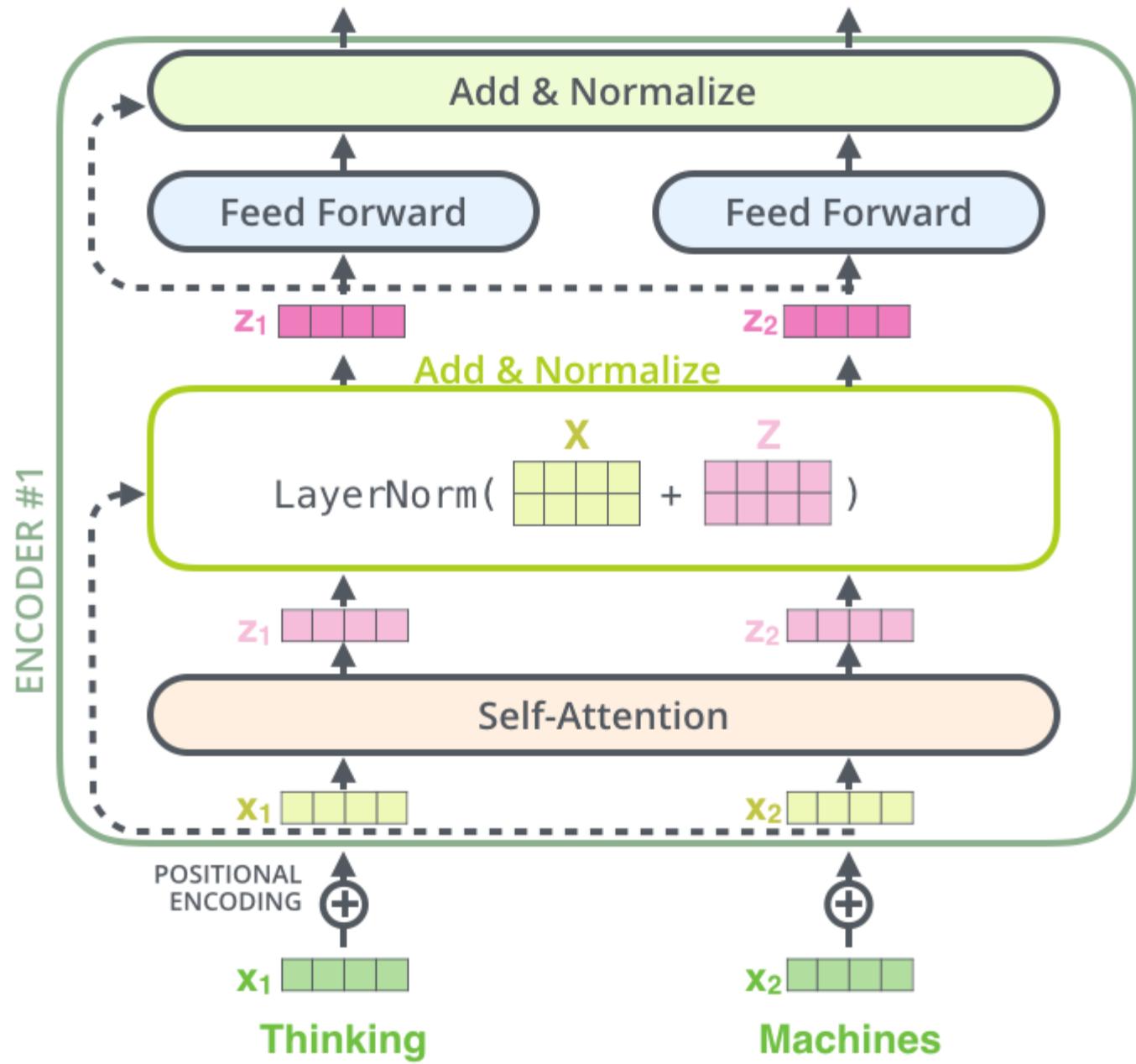
Transformers



Residuals

residual_mid

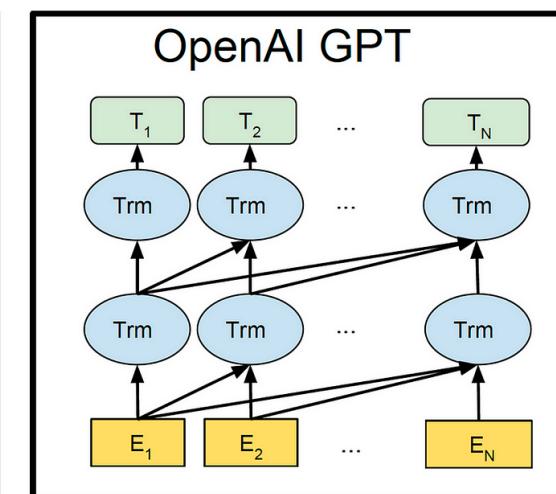
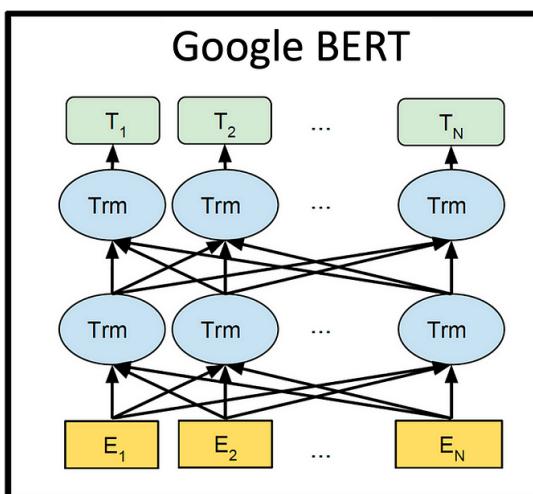
mlp_out + residual_mid

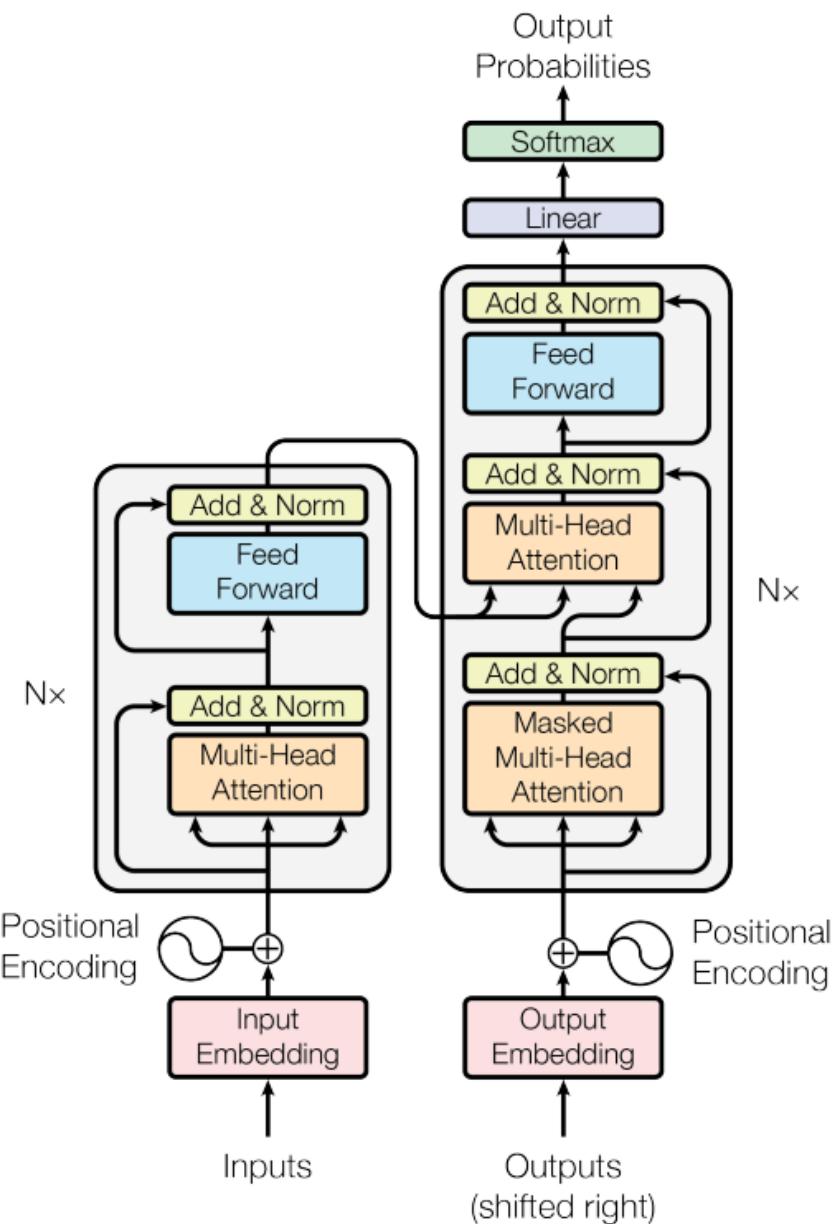


“Ogres are like onions

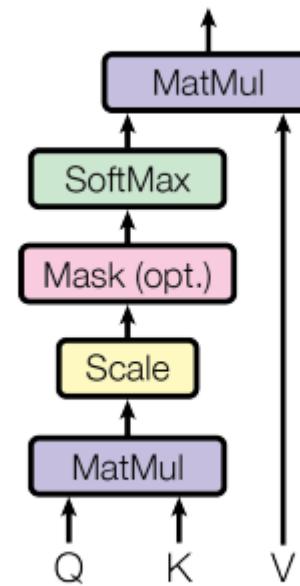
they have layers”

- Shrek

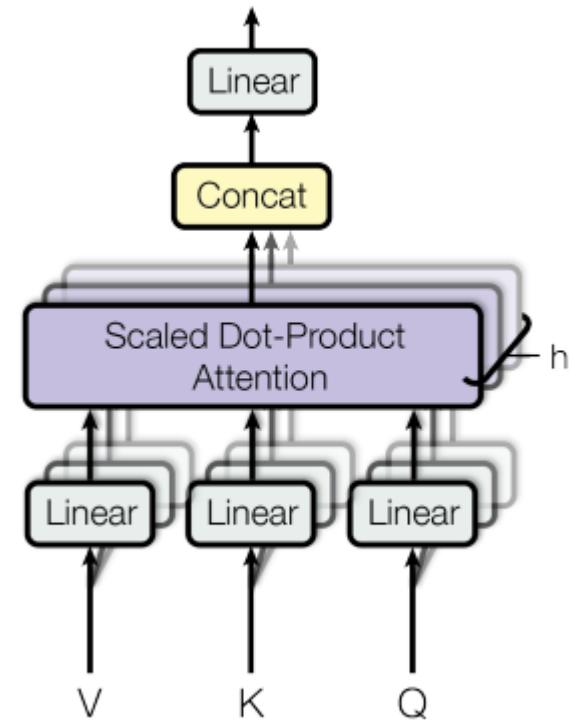




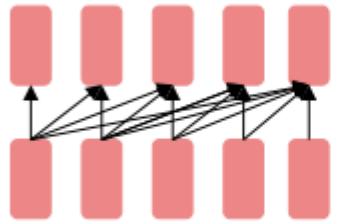
Scaled Dot-Product Attention



Multi-Head Attention

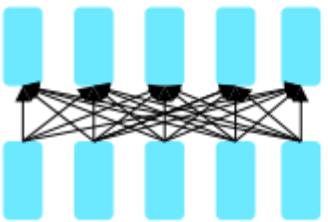


Three types of architectures



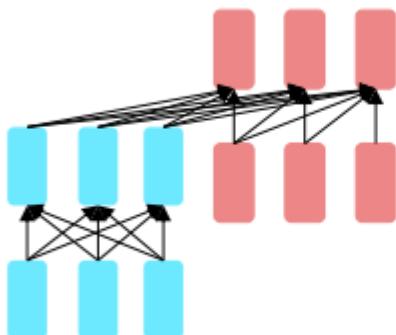
Decoders

- Next word prediction.
- Easy to train. Abundant amount of data.
- Nice to generate from; can't condition on future words.



Encoders

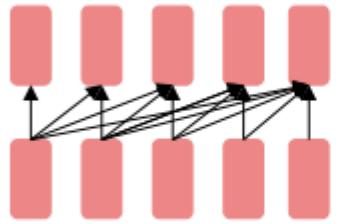
- Gets bidirectional context – can condition on future!
- Good word embeddings.
- MLM, BERT.



**Encoder-
Decoders**

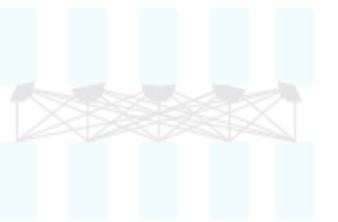
- Good parts of decoders and encoders?
- What's the best way to pretrain them?

Three types of architectures



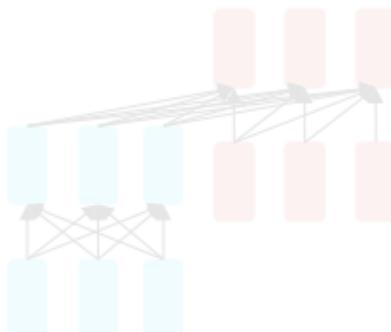
Decoders

- Next word prediction.
- Easy to train. Abundant amount of data.
- Nice to generate from; can't condition on future words.



Encoders

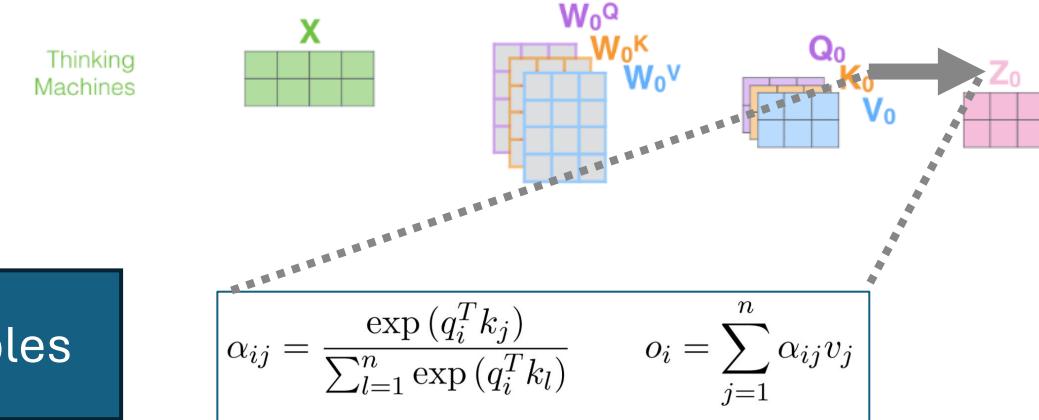
- Gets bidirectional context – can condition on future!
- Good word embeddings.
- MLM, BERT.



Encoder-Decoders

- Good parts of decoders and encoders?
- What's the best way to pretrain them?

Self Attention



I like eating apples

What is $o(\text{like})$? i.e., the attention score from like to the sentence.

query

Q4
“apples”

$$X \times W^Q = Q$$

key

| | | | |
|-----------|--------------|----------------|----------------|
| K1 “I” | K2 “like” | K3 “eating” | K4 “apples” |
|-----------|--------------|----------------|----------------|

$$X \times W^K = K$$

attention
weights*

| | | | |
|------|------|------|-----|
| 14.2 | 18.1 | 10.3 | 7.9 |
|------|------|------|-----|

$$Q \times K^T$$

attention
scores* (o)

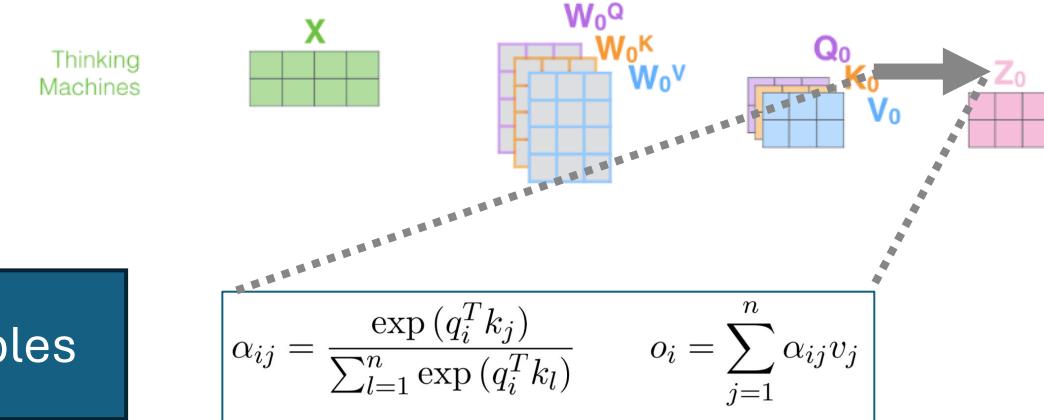
| | | | |
|-----|-----|------|------|
| 0.3 | 0.6 | 0.07 | 0.03 |
|-----|-----|------|------|

$$\text{softmax}\left(\frac{\text{Q} \times \text{K}^T}{\sqrt{d_k}} \right)$$

*: made-up numbers, not real.

Self Attention

I like eating apples



What is $z(\text{like})$?

attention scores* (o) 0.3 0.6 0.07 0.03

values

| | | | |
|-----------|--------------|----------------|----------------|
| V1 “I” | V2 “like” | V3 “eating” | V4 “apples” |
|-----------|--------------|----------------|----------------|

Weighted sum

Z4
“like”

$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}} \right)$$

$X \times W^V = V$

*: made-up numbers, not real.

GPT-2

Language Models are Unsupervised Multitask Learners

| | LAMBADA (PPL) | LAMBADA (ACC) | CBT-CN (ACC) | CBT-NE (ACC) | WikiText2 (PPL) | PTB (PPL) | enwik8 (BPB) | text8 (BPC) | WikiText103 (PPL) | 1BW (PPL) |
|-------|------------------|------------------|-----------------|-----------------|--------------------|--------------|-----------------|----------------|----------------------|--------------|
| SOTA | 99.8 | 59.23 | 85.7 | 82.3 | 39.14 | 46.54 | 0.99 | 1.08 | 18.3 | 21.8 |
| 117M | 35.13 | 45.99 | 87.65 | 83.4 | 29.41 | 65.85 | 1.16 | 1.17 | 37.50 | 75.20 |
| 345M | 15.60 | 55.48 | 92.35 | 87.1 | 22.76 | 47.33 | 1.01 | 1.06 | 26.37 | 55.72 |
| 762M | 10.87 | 60.12 | 93.45 | 88.0 | 19.93 | 40.31 | 0.97 | 1.02 | 22.05 | 44.575 |
| 1542M | 8.63 | 63.24 | 93.30 | 89.05 | 18.34 | 35.76 | 0.93 | 0.98 | 17.48 | 42.16 |

Table 3. Zero-shot results on many datasets. No training or fine-tuning was performed for any of these results. PTB and WikiText-2 results are from (Gong et al., 2018). CBT results are from (Bajgar et al., 2016). LAMBADA accuracy result is from (Hoang et al., 2018) and LAMBADA perplexity result is from (Grave et al., 2016). Other results are from (Dai et al., 2019).

Perplexity

$$\text{PPL}(W) = P(w_1 w_2 \dots w_N)^{-\frac{1}{N}}$$

Inverse probability of a corpus, according to a LM.

Normalized by the number of words.



Perplexity... Boring!

OpenAI's GPT-2 has been promoted as “an AI” that exemplifies an emergent understanding of language after mere unsupervised training on about 40GB of webpage text. It sounds really convincing in interviews:

- Q: *Which technologies are worth watching in 2020?*

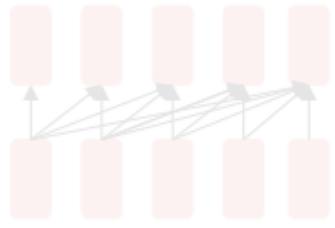
A: I would say it is hard to narrow down the list. The world is full of disruptive technologies with real and potentially huge global impacts. The most important is artificial intelligence, which is becoming exponentially more powerful. There is also the development of self-driving cars. There is a lot that we can do with artificial intelligence to improve the world....

- Q: *Are you worried that ai [sic] technology can be misused?*

A: Yes, of course. But this is a global problem and we want to tackle it with global solutions.....

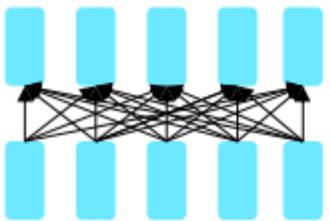
--- “AI can do that”, *The World in 2020 – The Economist*

Three types of architectures



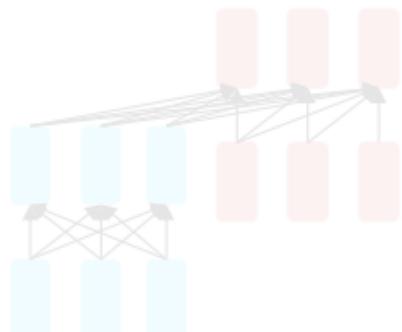
Decoders

- Next word prediction.
- Easy to train. Abundant amount of data.
- Nice to generate from; can't condition on future words.



Encoders

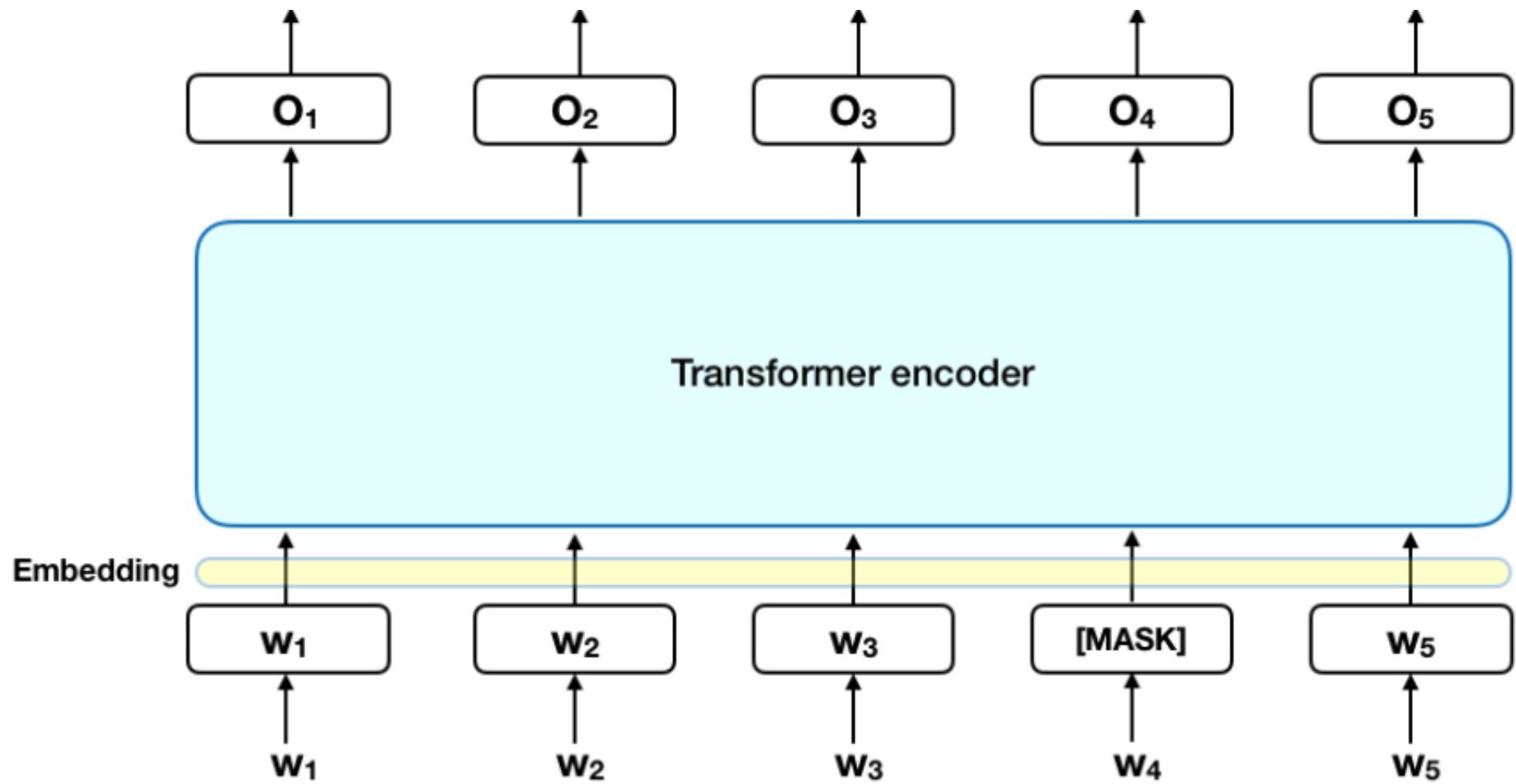
- Gets bidirectional context – can condition on future!
- Good word embeddings.
- MLM, BERT.



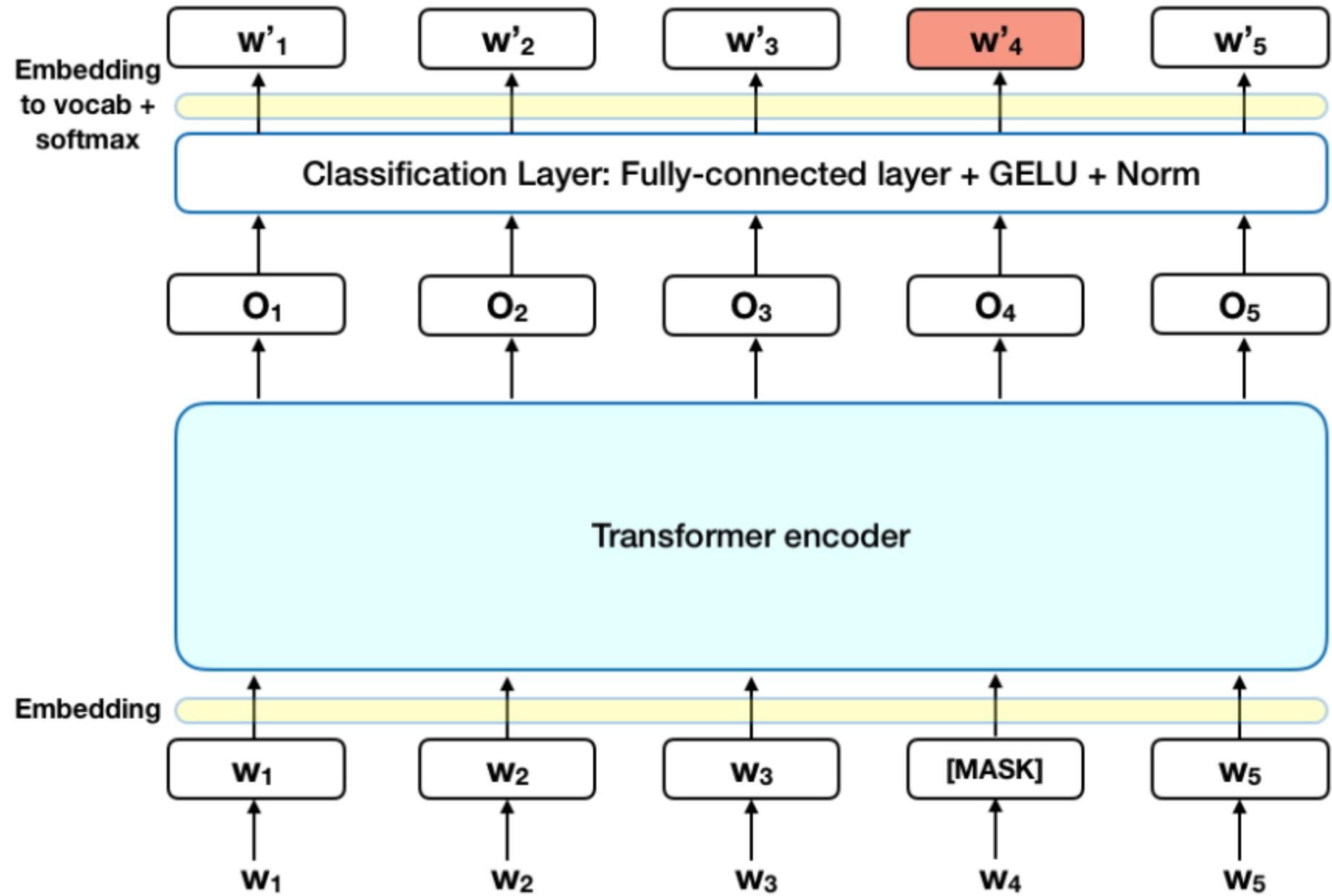
Encoder-Decoders

- Good parts of decoders and encoders?
- What's the best way to pretrain them?

BERT



BERT



Position Encoding

| Input | [CLS] | my | dog | is | cute | [SEP] | he | likes | play | #ing | [SEP] |
|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|
| Token Embeddings | E _[CLS] | E _{my} | E _{dog} | E _{is} | E _{cute} | E _[SEP] | E _{he} | E _{likes} | E _{play} | E _{#ing} | E _[SEP] |
| Segment Embeddings | + E _A | + E _B |
| Position Embeddings | E ₀ | E ₁ | E ₂ | E ₃ | E ₄ | E ₅ | E ₆ | E ₇ | E ₈ | E ₉ | E ₁₀ |

$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \sin(\omega_2 \cdot t) \\ \cos(\omega_2 \cdot t) \\ \vdots \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$

where

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$

$$\omega_k = \frac{1}{10000^{2k/d}}$$

Position Encoding

| Input | [CLS] | my | dog | is | cute | [SEP] | he | likes | play | #ing | [SEP] |
|---------------------|-------------|-----------------|------------------|-----------------|-------------------|--------------------|-----------------|--------------------|-------------------|--------------------|--------------------|
| Token Embeddings | $E_{[CLS]}$ | E_{my} | E_{dog} | E_{is} | E_{cute} | $E_{[\text{SEP}]}$ | E_{he} | E_{likes} | E_{play} | $E_{\#\text{ing}}$ | $E_{[\text{SEP}]}$ |
| Segment Embeddings | $+ E_A$ | $+ E_A$ | $+ E_A$ | $+ E_A$ | $+ E_A$ | $+ E_A$ | $+ E_B$ | $+ E_B$ | $+ E_B$ | $+ E_B$ | $+ E_B$ |
| Position Embeddings | E_0 | E_1 | E_2 | E_3 | E_4 | E_5 | E_6 | E_7 | E_8 | E_9 | E_{10} |

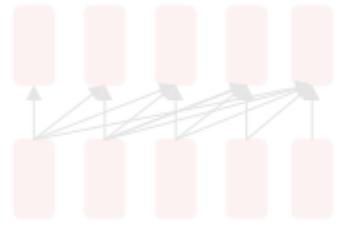
- Encodings of any two distinct positions are distinct
- Each position maps to only one encoding
- Test sentences may be longer than training
- Distance between two positions should be constant across sentences (of varying lengths).

$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \sin(\omega_2 \cdot t) \\ \cos(\omega_2 \cdot t) \\ \vdots \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1} \quad \vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$

where

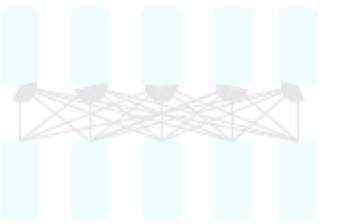
$$\omega_k = \frac{1}{10000^{2k/d}}$$

Three types of architectures



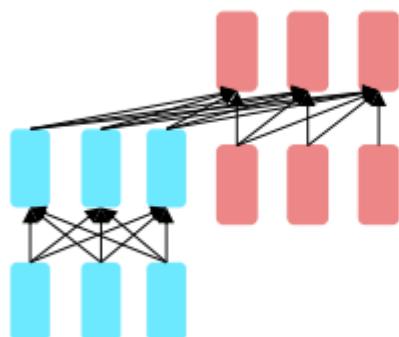
Decoders

- Next word prediction.
- Easy to train. Abundant amount of data.
- Nice to generate from; can't condition on future words.



Encoders

- Gets bidirectional context – can condition on future!
- Good word embeddings.
- MLM, BERT.



Encoder-
Decoders

- Good parts of decoders and encoders?
- What's the best way to pretrain them?

Pretraining encoder-decoders: What pretraining objective to use?



- What Raffel et al. (2018) found to work best was span corruption. Their model: T5.
- Replace different-length spans from the input with unique placeholders; decode out the spans that were removed!

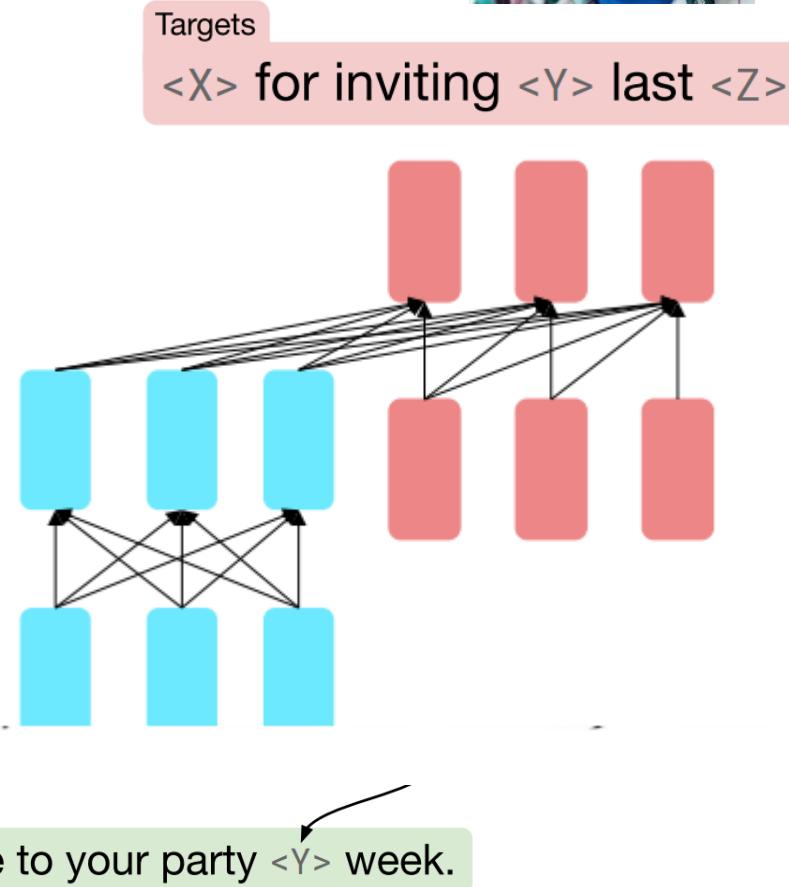
Original text

Thank you ~~for~~ ~~inviting~~ me to your party ~~last~~ week.

- This is implemented in text preprocessing: it's still an objective that looks like language modeling at the decoder side.

Inputs

Thank you <X> me to your party <Y> week.

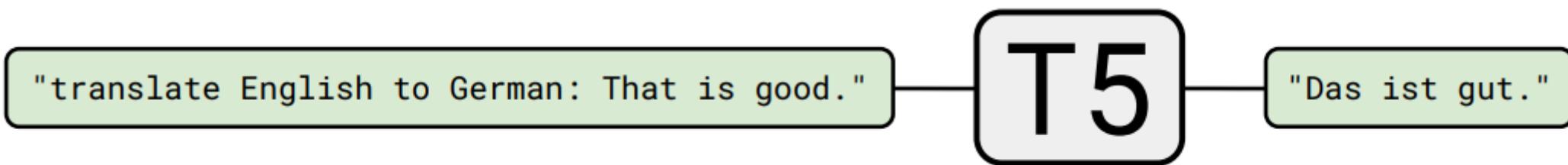


Pretraining encoder-decoders: What pretraining objective to use?

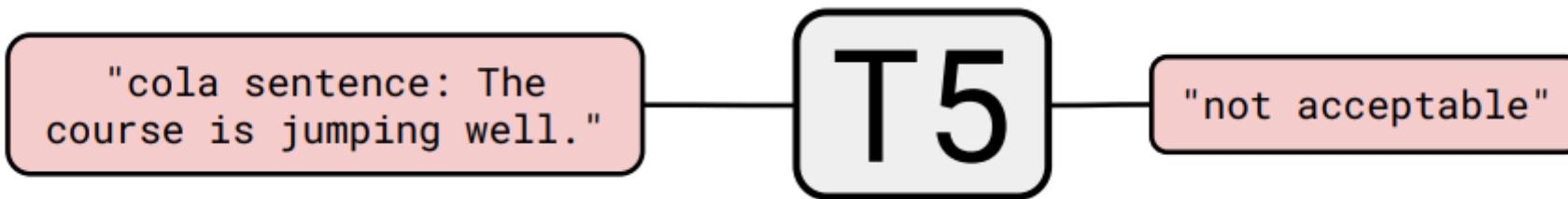
- Raffel et al., (2018) found encoder-decoders to work better than decoders for their tasks, and span corruption (denoising) to work better than language modeling.

| Architecture | Objective | Params | Cost | GLUE | CNNDM | SQuAD | SGLUE | EnDe | EnFr | EnRo |
|-------------------|-----------|--------|-------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| ★ Encoder-decoder | Denoising | $2P$ | M | 83.28 | 19.24 | 80.88 | 71.36 | 26.98 | 39.82 | 27.65 |
| | Denoising | P | M | 82.81 | 18.78 | 80.63 | 70.73 | 26.72 | 39.03 | 27.46 |
| | Denoising | P | $M/2$ | 80.88 | 18.97 | 77.59 | 68.42 | 26.38 | 38.40 | 26.95 |
| | Denoising | P | M | 74.70 | 17.93 | 61.14 | 55.02 | 25.09 | 35.28 | 25.86 |
| | Denoising | P | M | 81.82 | 18.61 | 78.94 | 68.11 | 26.43 | 37.98 | 27.39 |
| Encoder-decoder | LM | $2P$ | M | 79.56 | 18.59 | 76.02 | 64.29 | 26.27 | 39.17 | 26.86 |
| | LM | P | M | 79.60 | 18.13 | 76.35 | 63.50 | 26.62 | 39.17 | 27.05 |
| | LM | P | $M/2$ | 78.67 | 18.26 | 75.32 | 64.06 | 26.13 | 38.42 | 26.89 |
| | LM | P | M | 73.78 | 17.54 | 53.81 | 56.51 | 25.23 | 34.31 | 25.38 |
| | LM | P | M | 79.68 | 17.84 | 76.87 | 64.86 | 26.28 | 37.51 | 26.76 |

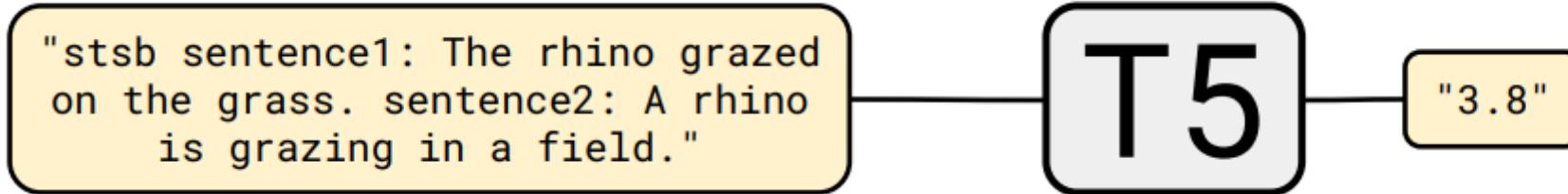
One surprising finding



One surprising finding



One surprising finding



One surprising finding

"translate English to German: That is good."

"cola sentence: The course is jumping well."

"stsbs sentence1: The rhino grazed on the grass. sentence2: A rhino is grazing in a field."

"summarize: state authorities dispatched emergency crews tuesday to survey the damage after an onslaught of severe weather in mississippi..."

A fascinating property of T5:

- It can be finetuned to answer a wide range of questions, retrieving knowledge from its parameters.
- With natural language!

T5

"Das ist gut."

"not acceptable"

"3.8"

"six people hospitalized after a storm in attala county."

One surprising finding

Text-to-Text Transfer Transformer

"translate English to German: That is good."

"cola sentence: The course is jumping well."

"stsbs sentence1: The rhino grazed on the grass. sentence2: A rhino is grazing in a field."

"summarize: state authorities dispatched emergency crews tuesday to survey the damage after an onslaught of severe weather in mississippi..."

T5

"Das ist gut."

"not acceptable"

"3.8"

"six people hospitalized after a storm in attala county."

Input

Embedding

Queries

Keys

Values

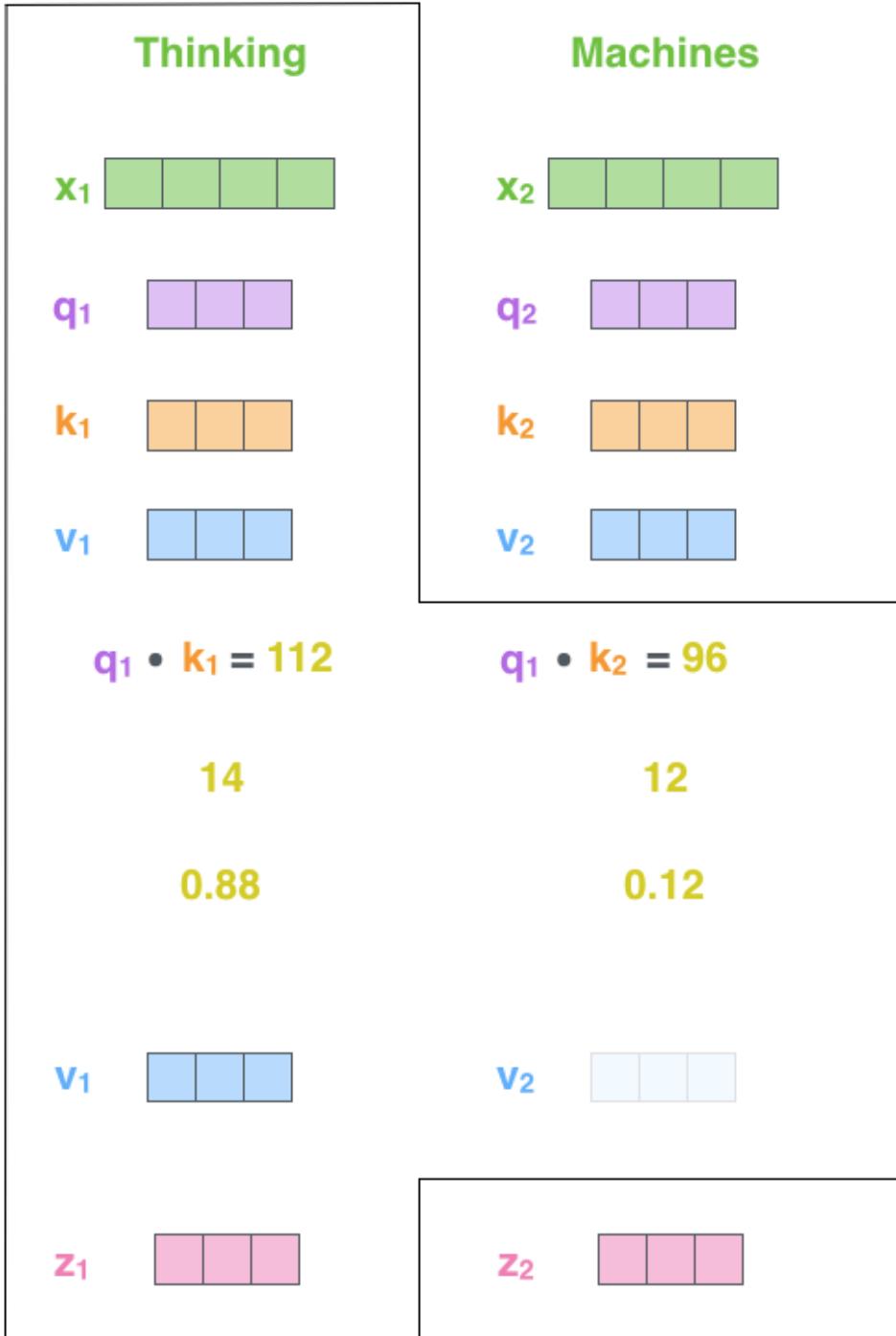
Score

Divide by 8 ($\sqrt{d_k}$)

Softmax

Softmax
X
Value

Sum



Review

- Embeddings:
 - Token + Sentence + Position
- Multi-Head Attention
- Feed forward module (MLP)
- Layers