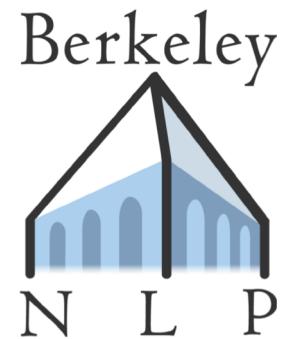


# Overview and Transformer Language Models



Eric Wallace

CS 288, 3/13/2023



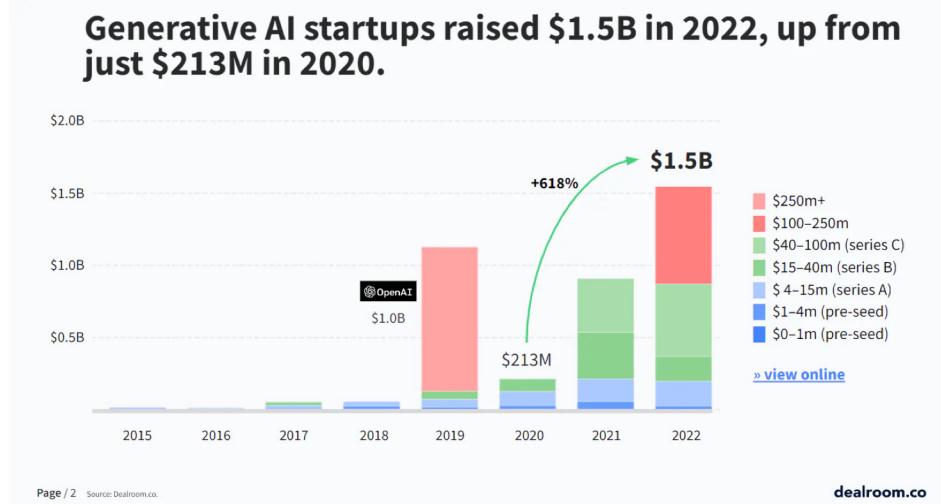
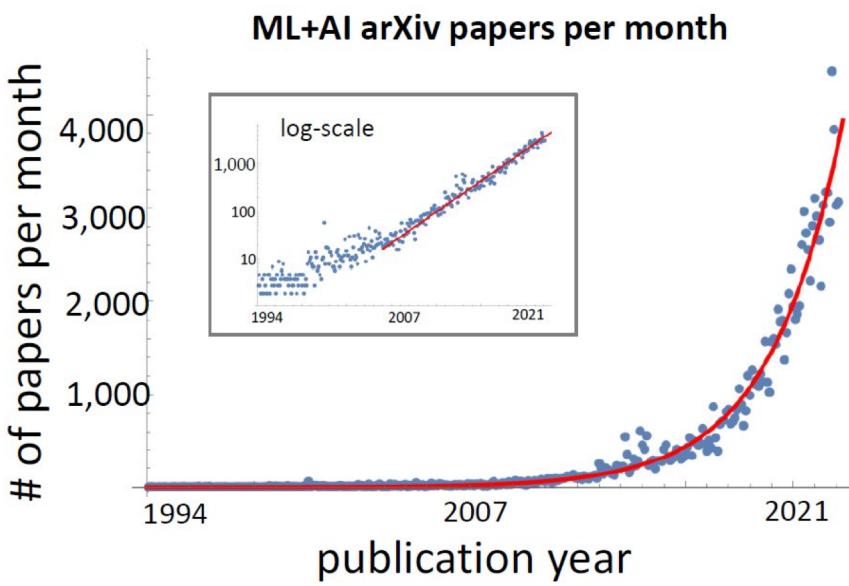
# Logistics

---

- 4 traditional lectures + ~8 days of mixed lectures and panels/discussions
- HW4 out Wednesday. Due Wednesday after spring break
  - Using and finetuning LMs with Huggingface
- HW5 out after spring break. Due sometime end of April.
  - Prompting ChatGPT to solve projects 1-3
- No final exam.
- Lecture recordings?



# Immense Interest



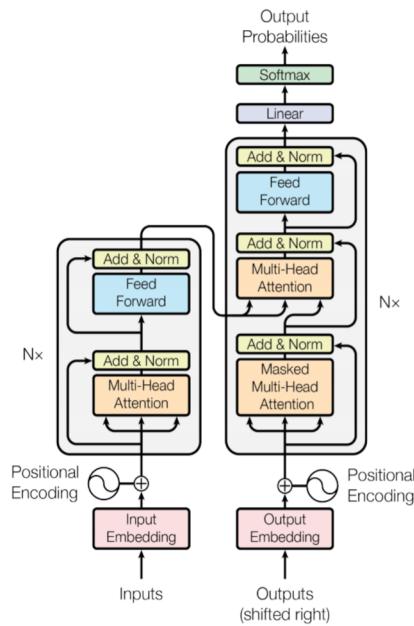
Page / 2 Source: Dealroom.co.



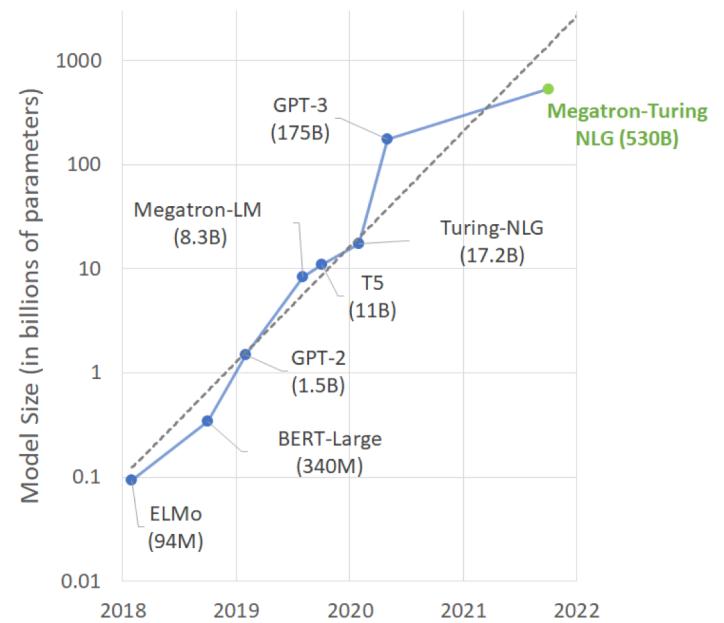
# The Era of Rapid Scaling in NLP

2017: Transformer is introduced

[Vaswani+17] Attention is All You Need



2022: Large-scale Transformer models are the dominant approach for many NLP tasks





## Demos

---

- ▶ [ChatGPT](#)
- ▶ [Stable Diffusion](#)
- ▶ [InstructGPT](#)



# Today's Lecture

---

- ▶ Language modeling as the ultimate task
- ▶ Transformer models
- ▶ Overview of remainder of the course



# Language Modeling

---

$$p(x_1, \dots, x_L)$$



## Language Modeling

---

$$p(x_1, \dots, x_L)$$

$p(\text{the, mouse, ate, the, cheese}) = 0.02,$

$p(\text{the, cheese, ate, the, mouse}) = 0.01,$

$p(\text{mouse, the, the, cheese, ate}) = 0.0001$

---



# Neural Language Models

---

$$\prod_{i=1}^L p(x_i \mid x_{1:i-1})$$



# Neural Language Models

---

*Prompt*

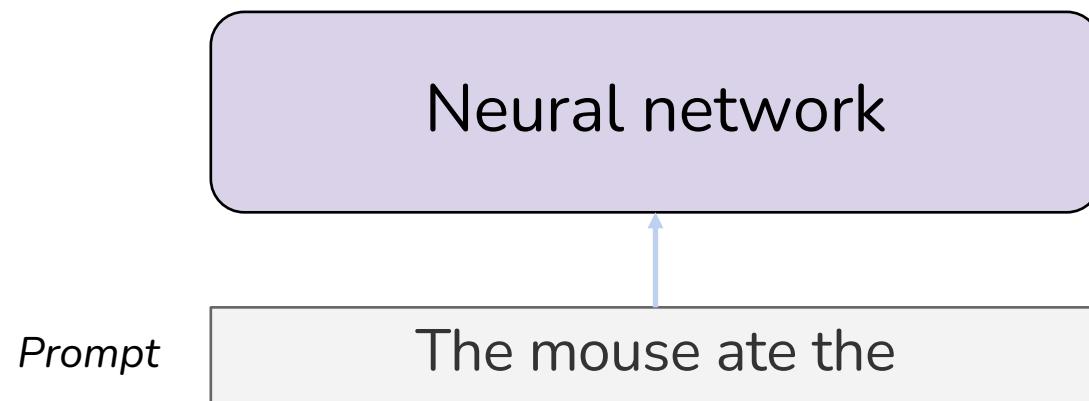
The mouse ate the

$$\prod_{i=1}^L p(x_i \mid x_{1:i-1})$$



# Neural Language Models

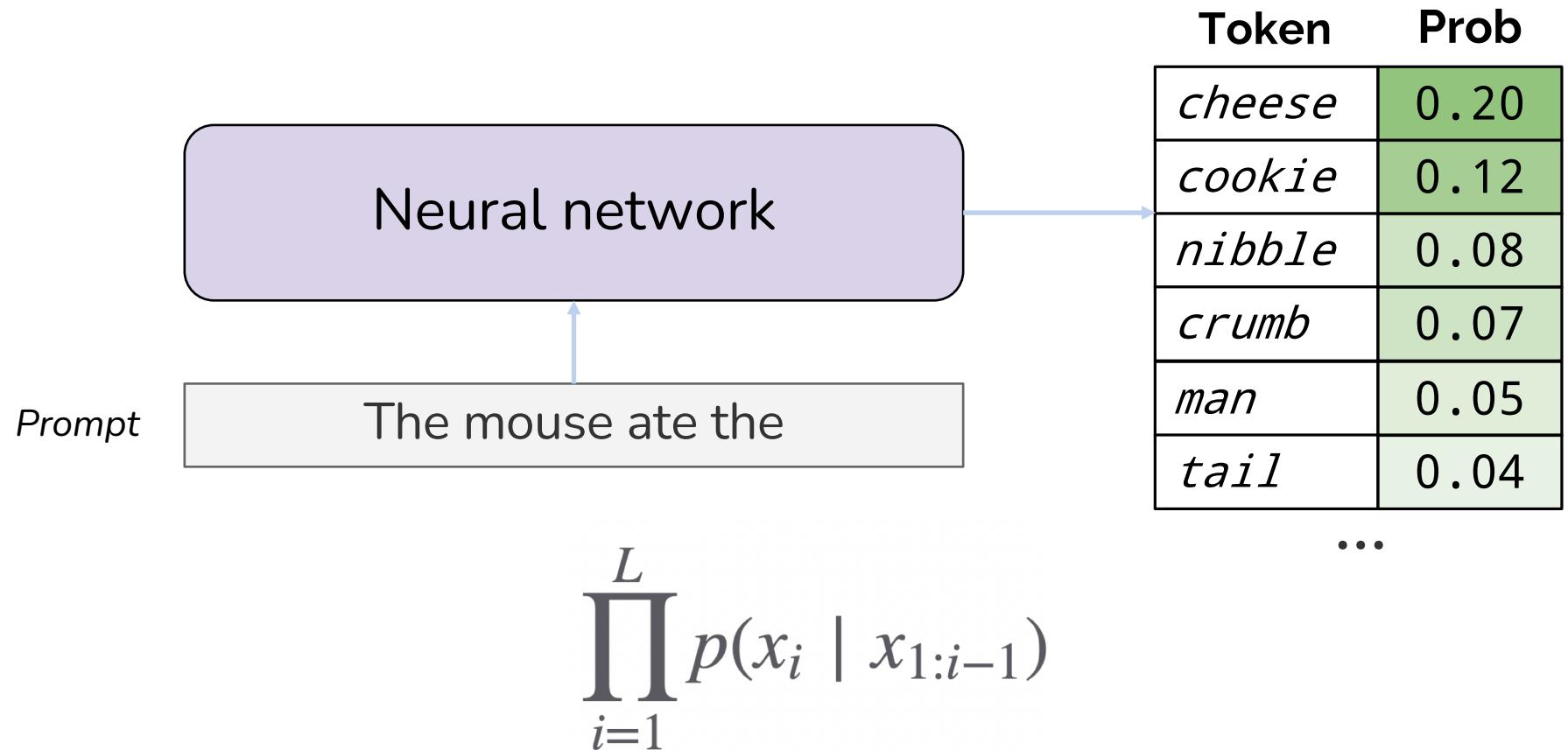
---



$$\prod_{i=1}^L p(x_i \mid x_{1:i-1})$$



# Neural Language Models





# Language Modeling

---

- ▶ Many original motivations were to use LMs for other applications
  - Machine translation
  - Speech recognition
  - ...
- ▶ Now, LM has become perhaps the single most important NLP task



# Language Modeling as the Ultimate Task?

---

- ▶ Zero- and few-shot learning with language models



# Language Modeling as the Ultimate Task?

---

- ▶ Zero- and few-shot learning with language models

Language Model

*Prompt*

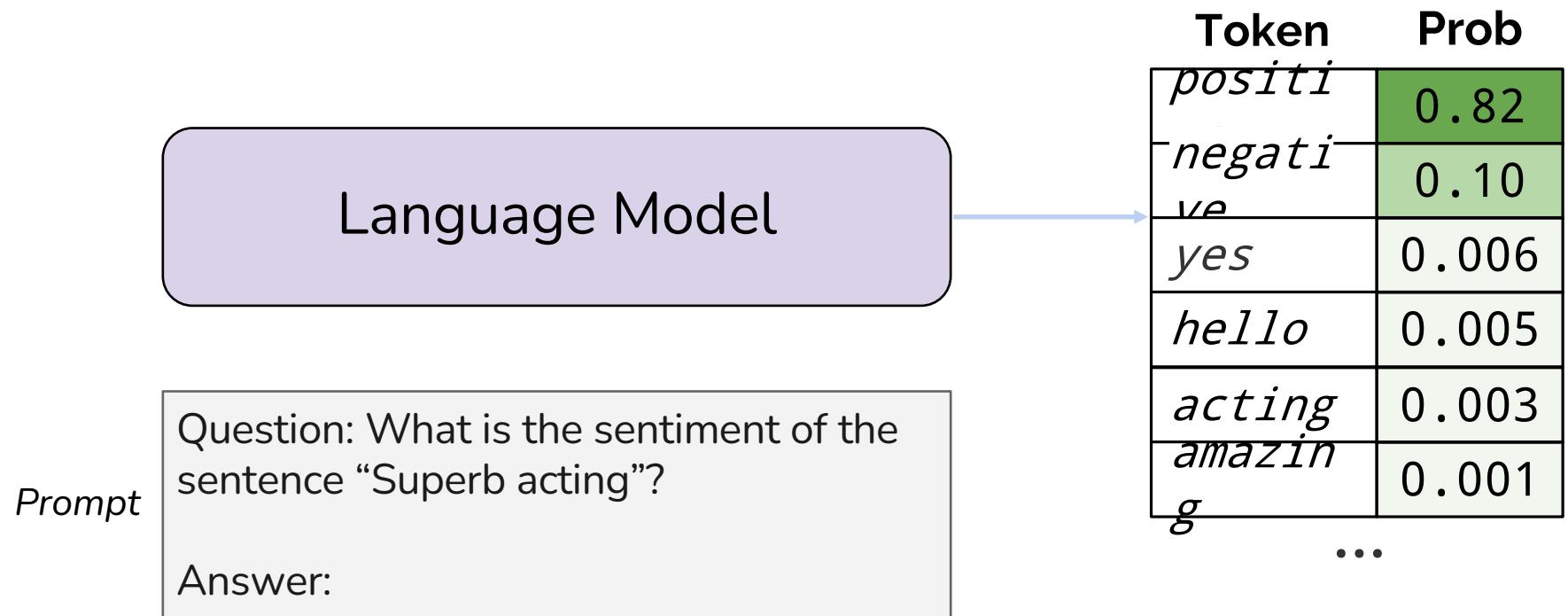
Question: What is the sentiment of the sentence “Superb acting”?

Answer:



# Language Modeling as the Ultimate Task?

- Zero- and few-shot learning with language models





# Language Modeling as the Ultimate Task?

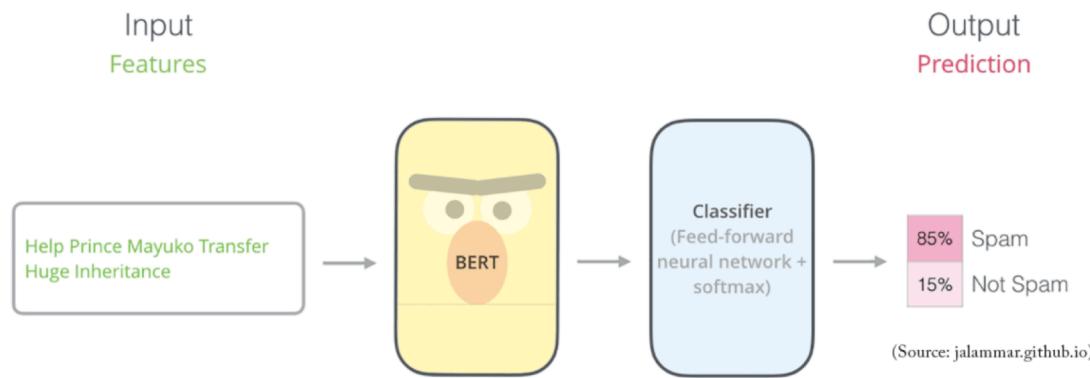
---

- ▶ Language modeling leads to rich representations
  - George Washington was born in the year \_\_\_\_\_
  - If it is raining, you may need an \_\_\_\_\_
  - Using the power rule, the derivative of  $3x^5$  is \_\_\_\_\_



# Language Modeling as the Ultimate Task?

- ▶ Language modeling leads to rich representations
  - George Washington was born in the year \_\_\_\_\_
  - If it is raining, you may need an \_\_\_\_\_
  - Using the power rule, the derivative of  $3x^5$  is \_\_\_\_\_





# Language Modeling as the Ultimate Task?

---

- There is effectively “unlimited” data for language modeling
- Enables powerful function approximators (transformers)
  - immense data
  - immense model sizes
  - immense compute



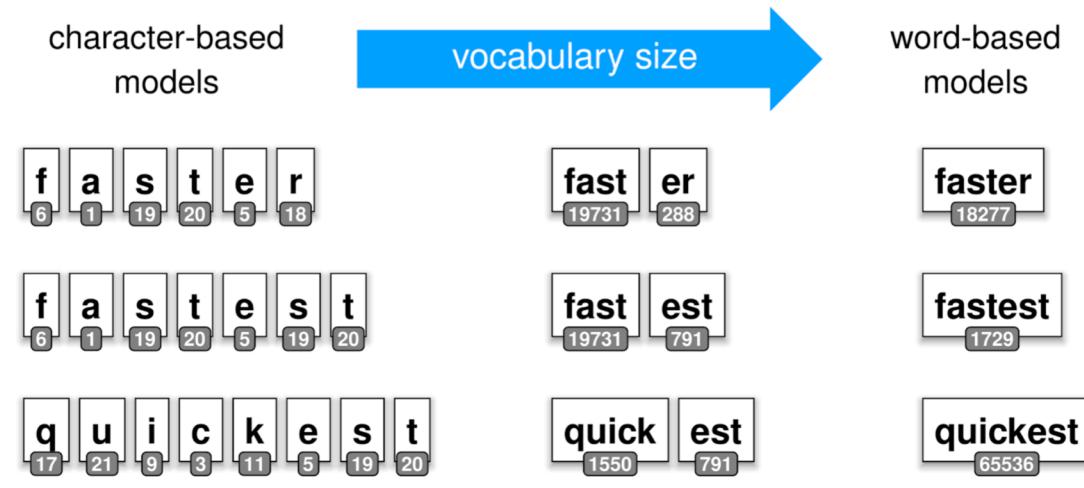
# Neural LMs from Scratch

---



# Neural LMs from Scratch

- ▶ Input encoding

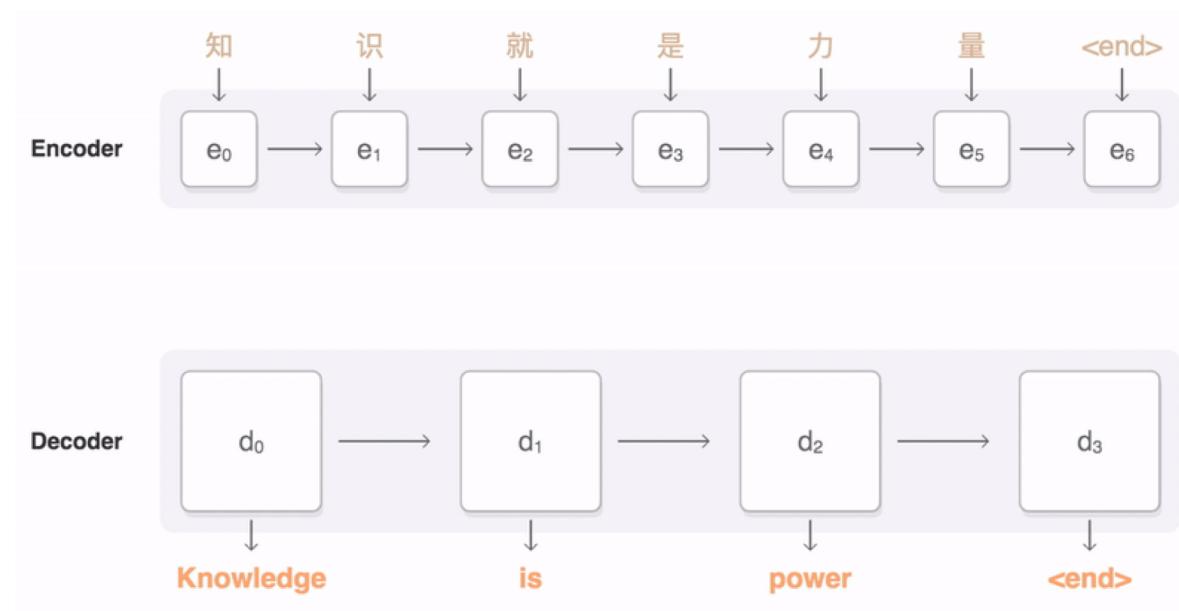




# Language Models and MT Circa 2016

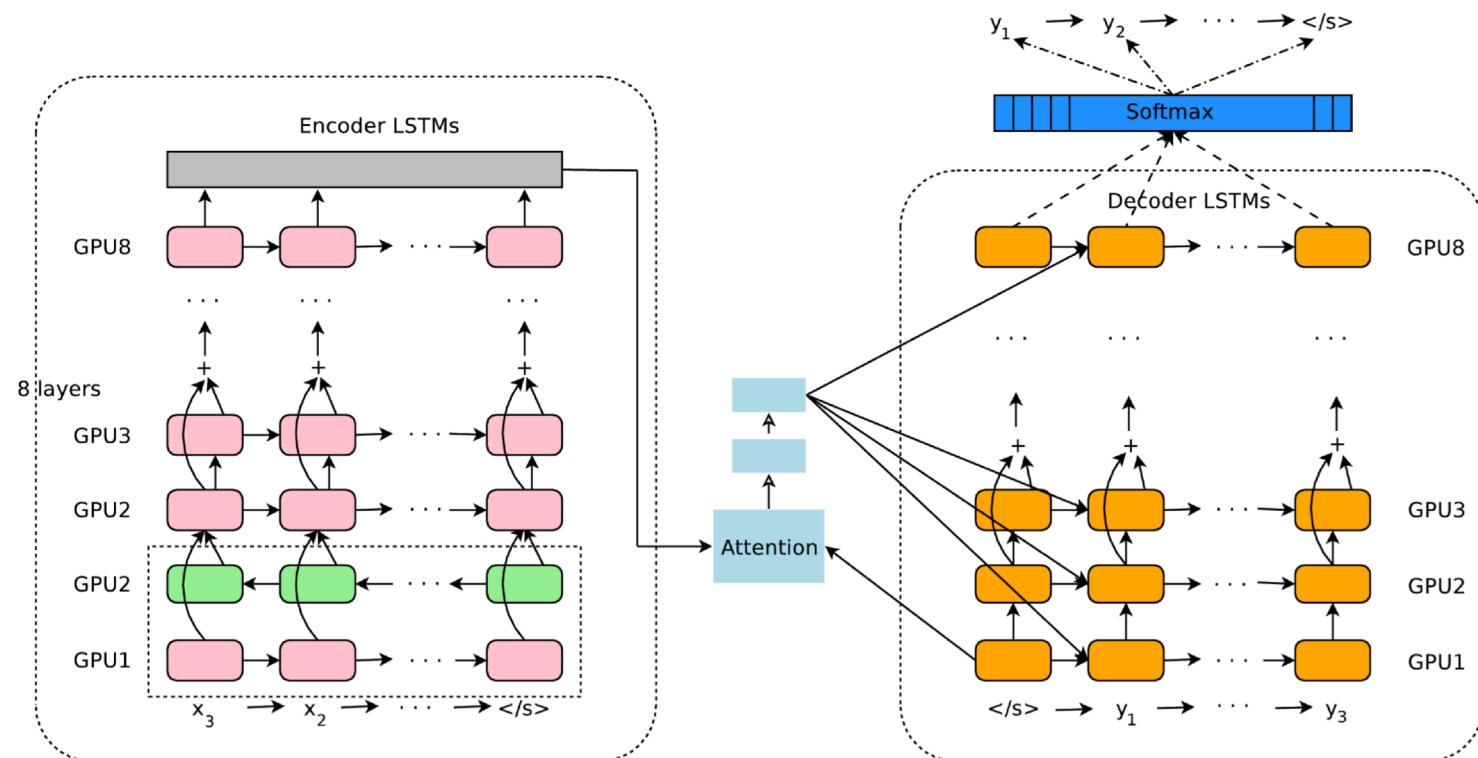
Neural Machine Translation is in production at Google

[Wu+16] [Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation](#)



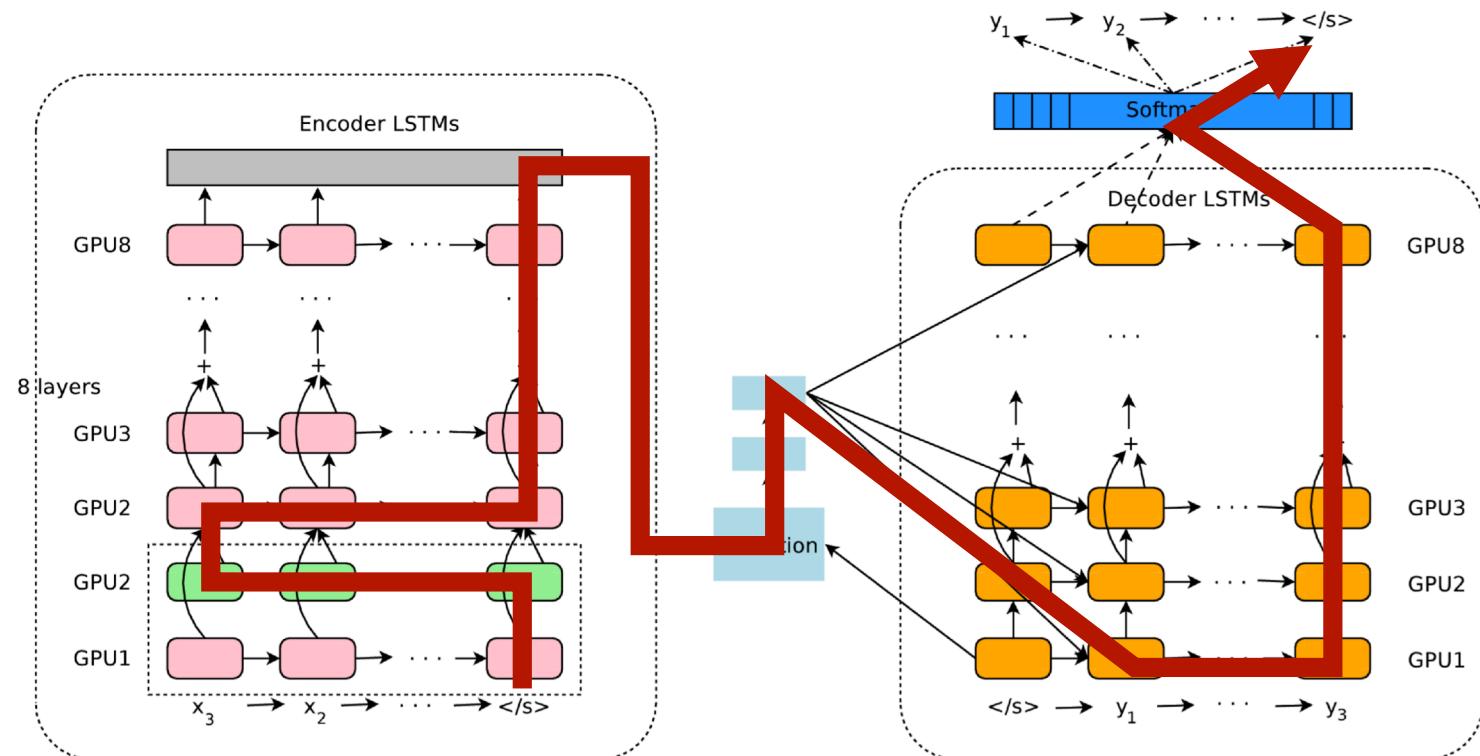


# Neural MT ca. 2016





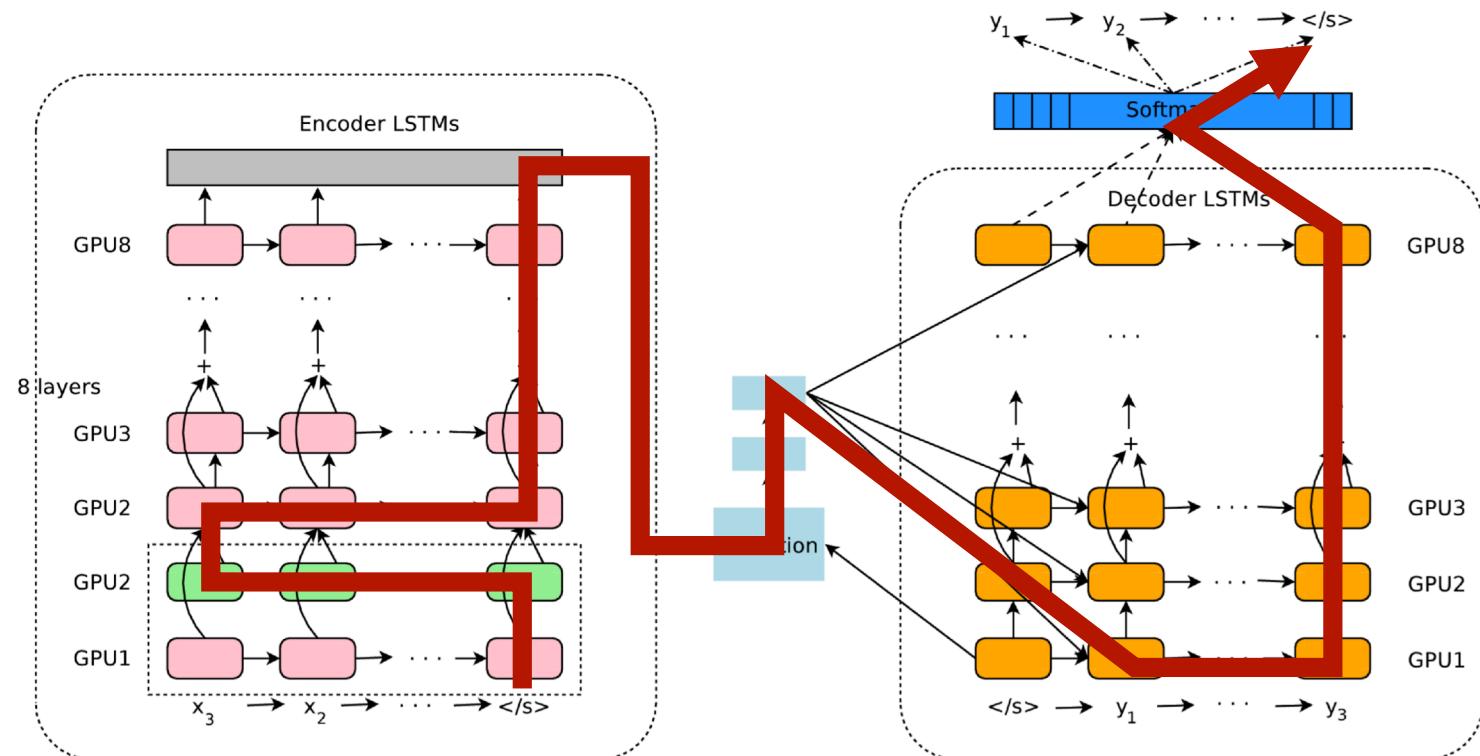
# Neural MT ca. 2016



There are computation paths through the RNN-based network that scale linearly with the sequence length, and can't be parallelized.



# Neural MT ca. 2016



There are computation paths through the RNN-based network that scale linearly with the sequence length, and can't be parallelized.



# Word Window Neural Nets

---

as the proctor started the clock the students opened their \_\_\_\_\_



# Word Window Neural Nets

as the proctor started the clock  
discard

the students opened their \_\_\_\_\_  
fixed window

A diagram illustrating word windows in a sequence of words. On the left, the words "as", "the", "proctor", "started", "the", "clock" are shown, with "clock" underlined in red and the word "discard" written below it. On the right, the words "the", "students", "opened", "their", and an empty slot "\_\_\_\_\_". A pink bracket groups "the", "students", "opened", and "their", with the label "fixed window" written below it.



# Word Window Neural Nets

words / one-hot vectors  
 $x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$

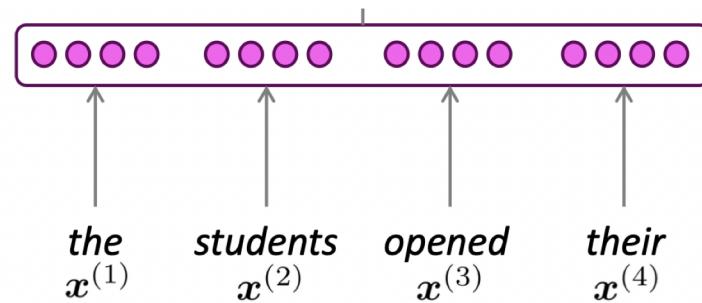
*the*      *students*      *opened*      *their*  
 $x^{(1)}$        $x^{(2)}$        $x^{(3)}$        $x^{(4)}$



# Word Window Neural Nets

concatenated word embeddings  
 $e = [e^{(1)}; e^{(2)}; e^{(3)}; e^{(4)}]$

words / one-hot vectors  
 $x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$





# Word Window Neural Nets

output distribution

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{U}\mathbf{h} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

hidden layer

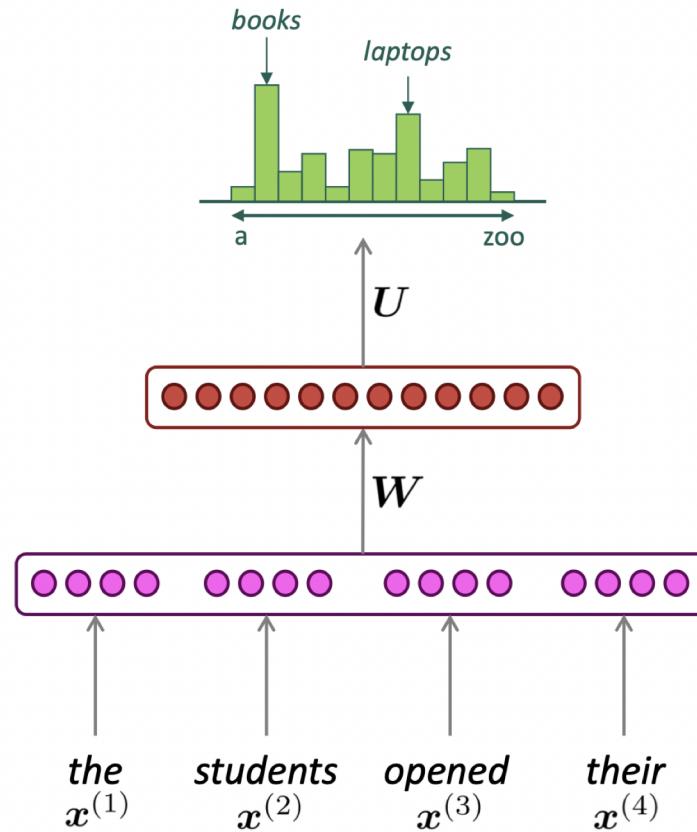
$$\mathbf{h} = f(\mathbf{W}\mathbf{e} + \mathbf{b}_1)$$

concatenated word embeddings

$$\mathbf{e} = [\mathbf{e}^{(1)}; \mathbf{e}^{(2)}; \mathbf{e}^{(3)}; \mathbf{e}^{(4)}]$$

words / one-hot vectors

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}, \mathbf{x}^{(4)}$$





# Word Averaging Neural Nets

output distribution

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{U}\mathbf{h} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

hidden layer

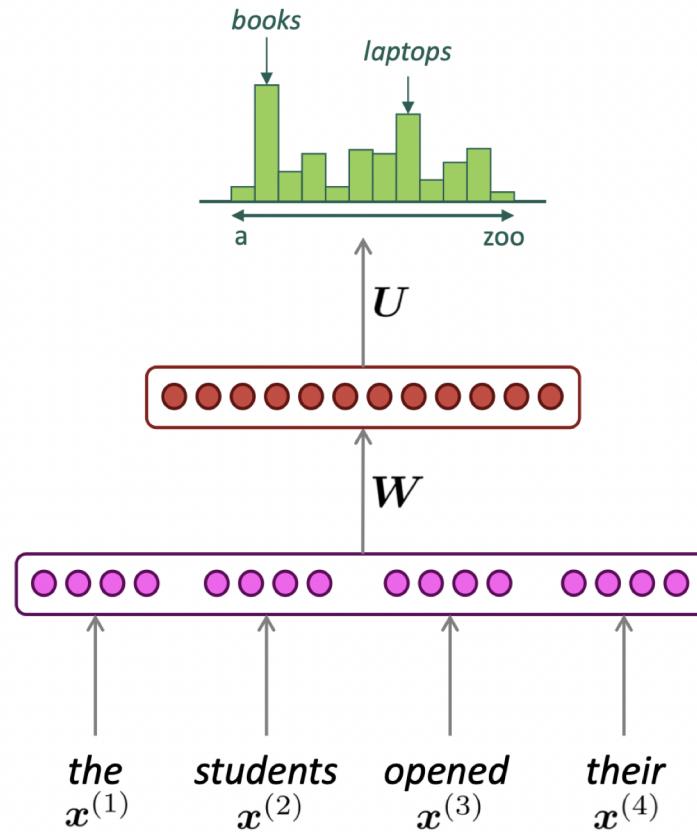
$$\mathbf{h} = f(\mathbf{W}\mathbf{e} + \mathbf{b}_1)$$

concatenated word embeddings

$$\mathbf{e} = [\mathbf{e}^{(1)}; \mathbf{e}^{(2)}; \mathbf{e}^{(3)}; \mathbf{e}^{(4)}]$$

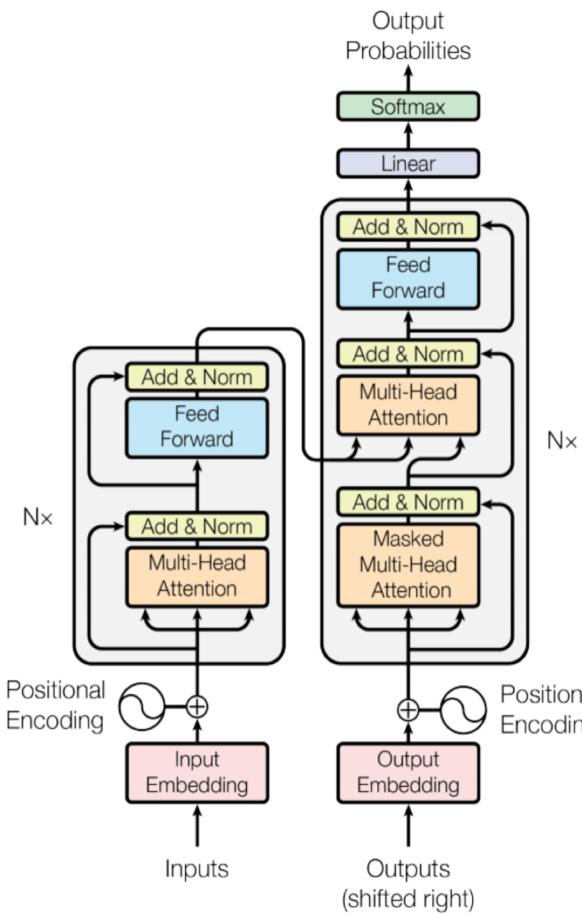
words / one-hot vectors

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)}, \mathbf{x}^{(4)}$$



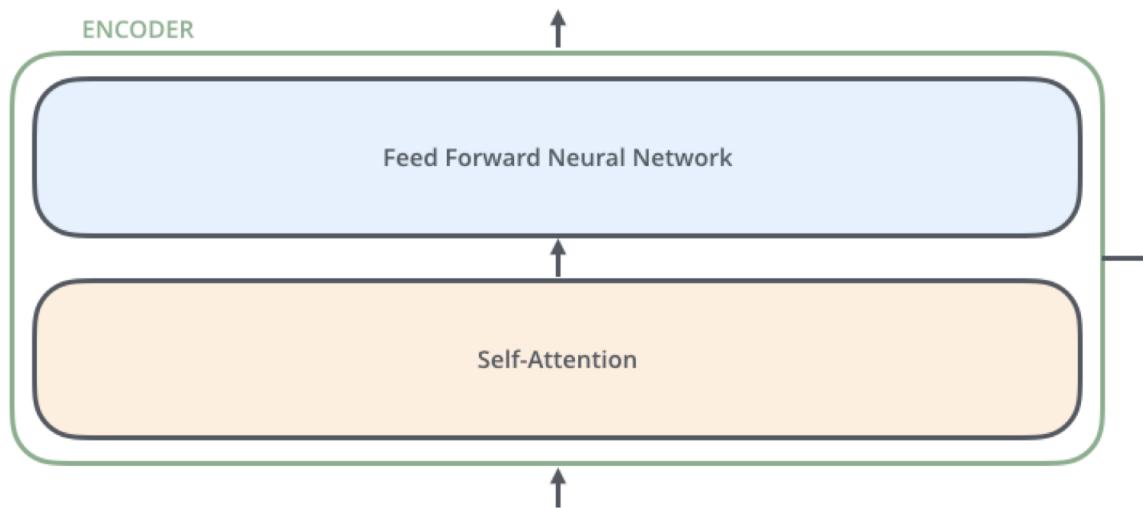


# Transformer Architecture



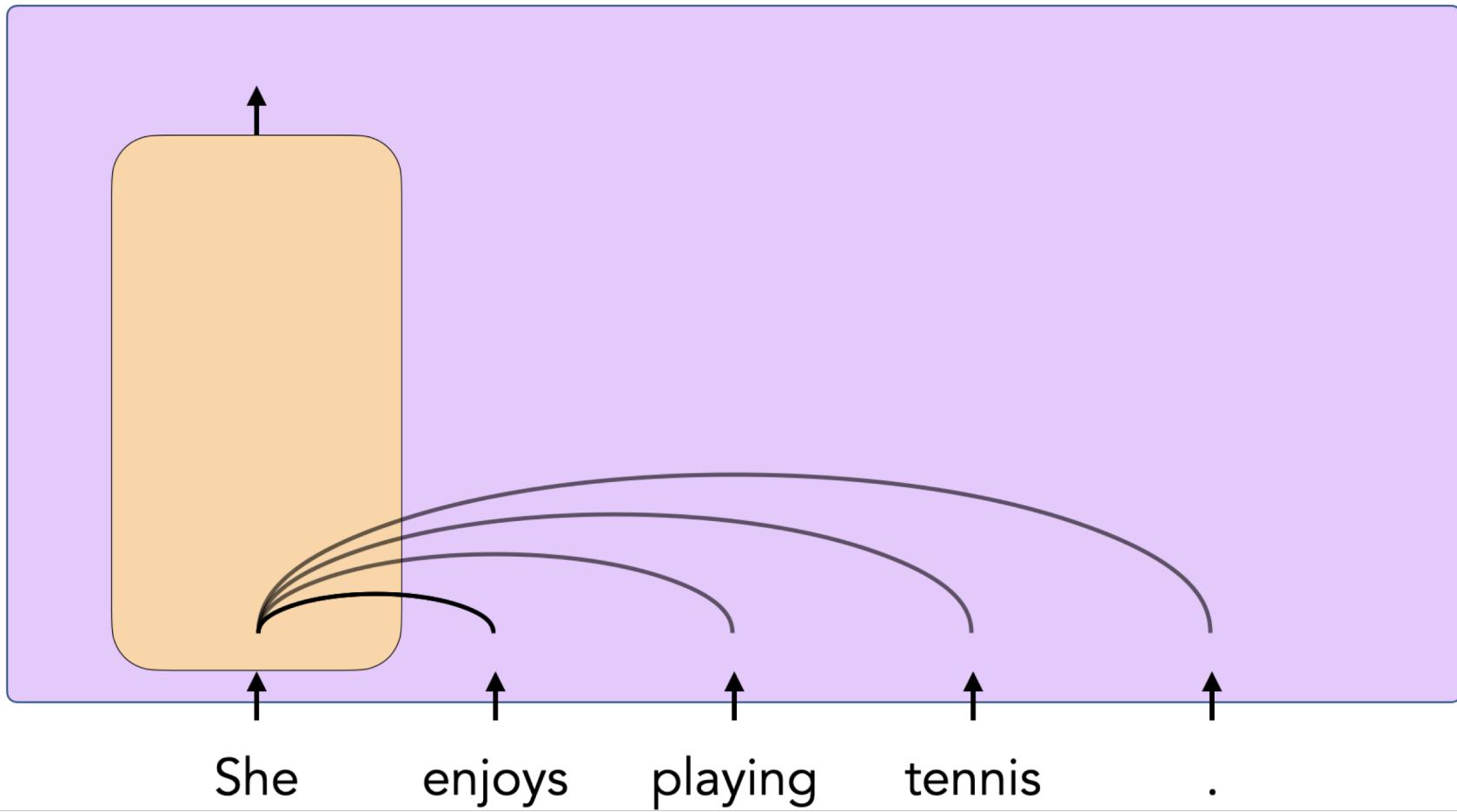


# Transformer Architecture





# Encoder





# Transformer Architecture

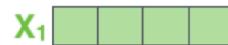
Input

Thinking

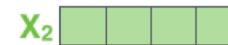
Machines

Embedding

$X_1$

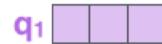


$X_2$

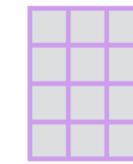
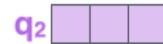


Queries

$q_1$



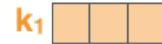
$q_2$



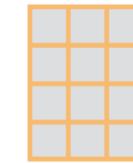
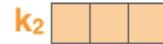
$W^Q$

Keys

$k_1$



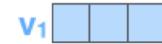
$k_2$



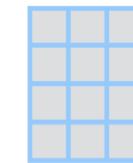
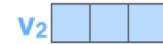
$W^K$

Values

$v_1$



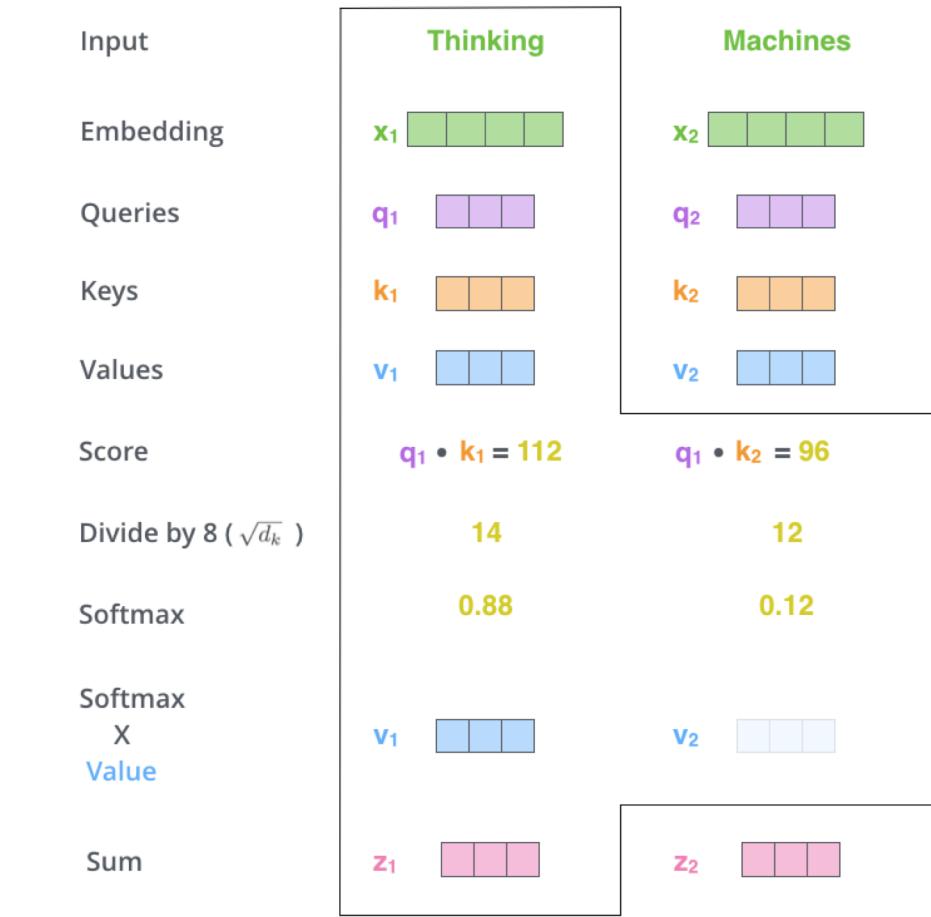
$v_2$



$W^V$



# Transformer Architecture





# Transformer Architecture

$$\mathbf{X} \times \mathbf{W}^Q = \mathbf{Q}$$

A diagram illustrating a matrix multiplication operation. On the left, a green input matrix  $\mathbf{X}$  is shown as a 3x3 grid of squares. To its right is a multiplication sign ( $\times$ ). Next is a purple weight matrix  $\mathbf{W}^Q$ , also a 3x3 grid, but with a distinct purple border. An equals sign ( $=$ ) follows, leading to the output matrix  $\mathbf{Q}$ , which is a 3x3 grid of purple squares.

$$\mathbf{X} \times \mathbf{W}^K = \mathbf{K}$$

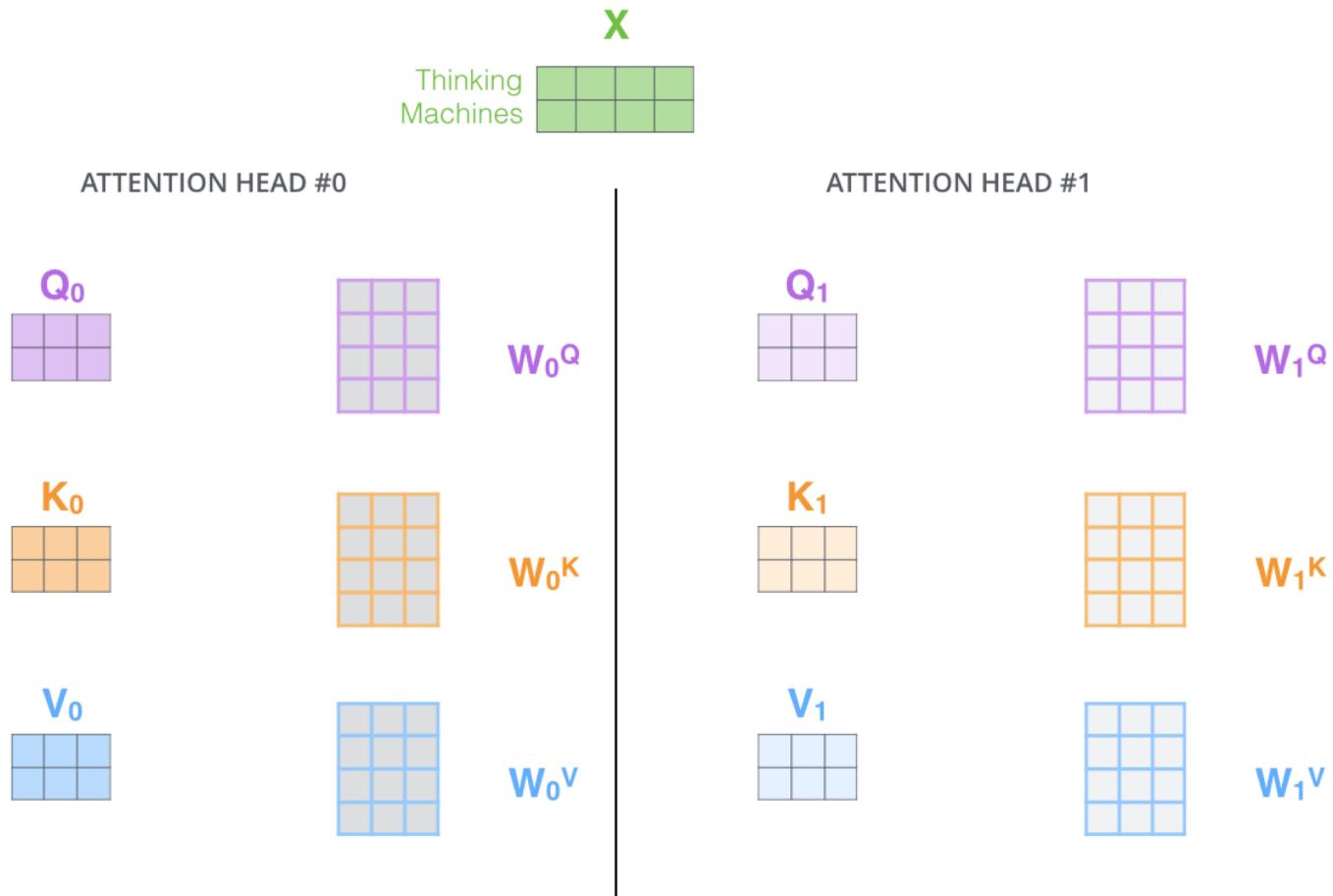
A diagram illustrating a matrix multiplication operation. On the left, a green input matrix  $\mathbf{X}$  is shown as a 3x3 grid of squares. To its right is a multiplication sign ( $\times$ ). Next is an orange weight matrix  $\mathbf{W}^K$ , a 3x3 grid with an orange border. An equals sign ( $=$ ) follows, leading to the output matrix  $\mathbf{K}$ , which is a 3x3 grid of orange squares.

$$\mathbf{X} \times \mathbf{W}^V = \mathbf{V}$$

A diagram illustrating a matrix multiplication operation. On the left, a green input matrix  $\mathbf{X}$  is shown as a 3x3 grid of squares. To its right is a multiplication sign ( $\times$ ). Next is a light blue weight matrix  $\mathbf{W}^V$ , a 3x3 grid with a light blue border. An equals sign ( $=$ ) follows, leading to the output matrix  $\mathbf{V}$ , which is a 3x3 grid of light blue squares.

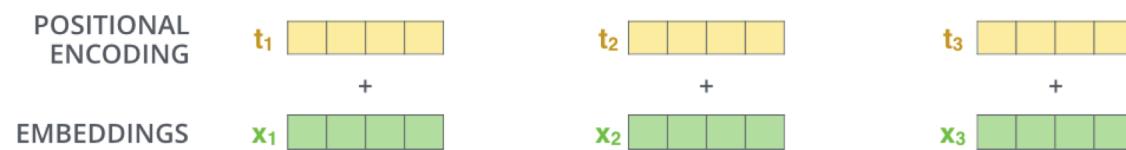


# Transformer Architecture



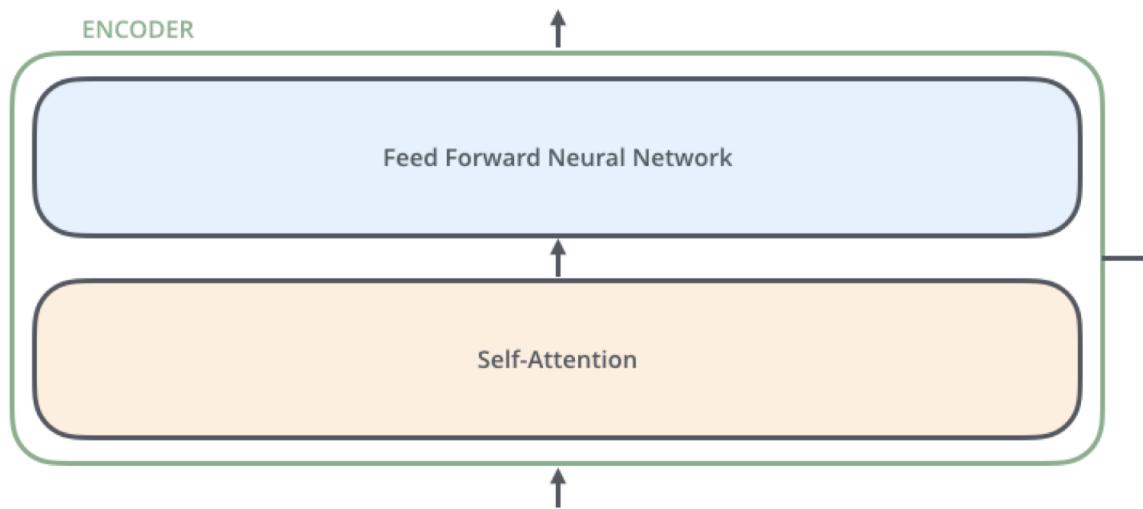


# Position Embeddings





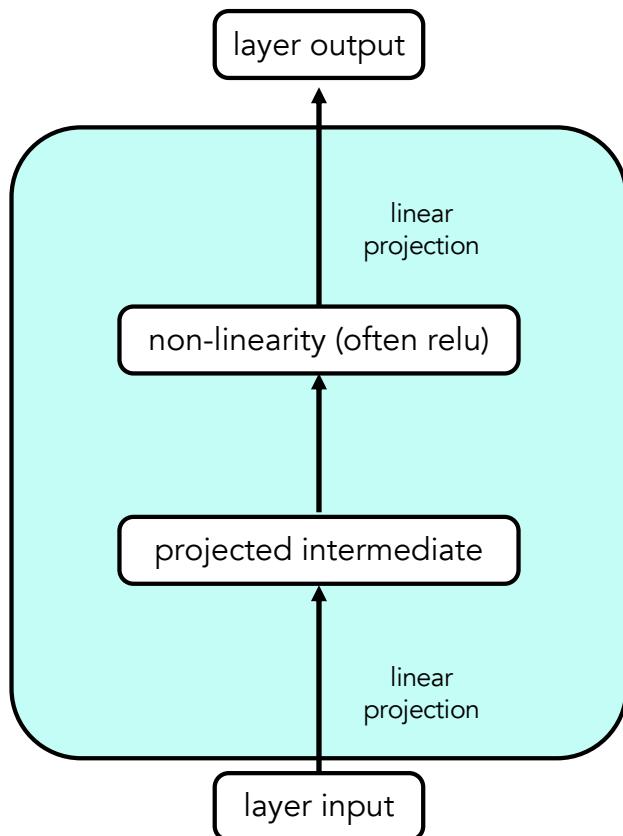
# Transformer Architecture





# Feed-Forward

---





# Add & Norm

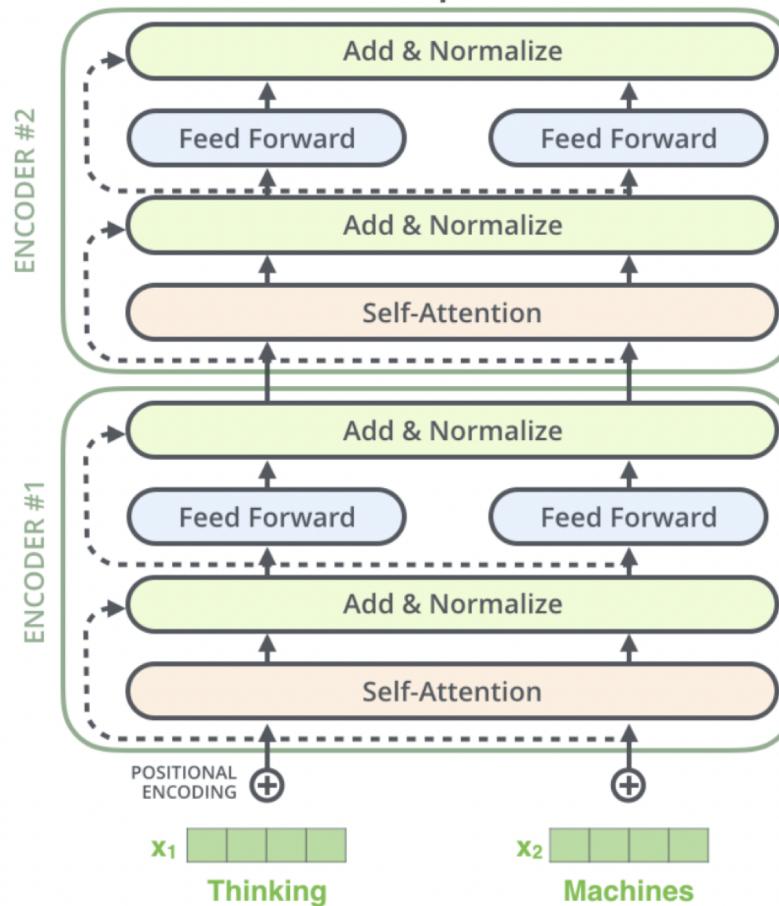
---

Layer Normalization [Ba+16]  
improves stability of neuron activations



Residual Connections  
useful across a variety of neural network architecture types, not just in NLP

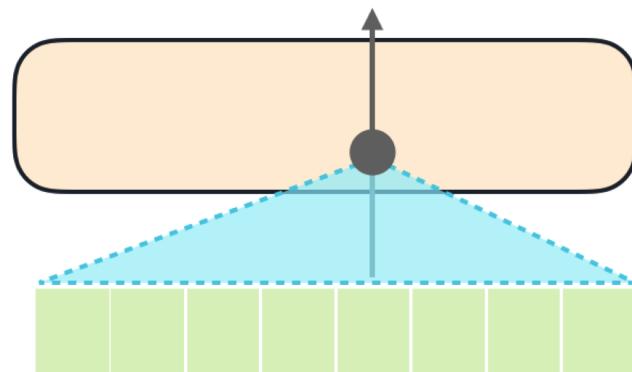
# Add & Norm



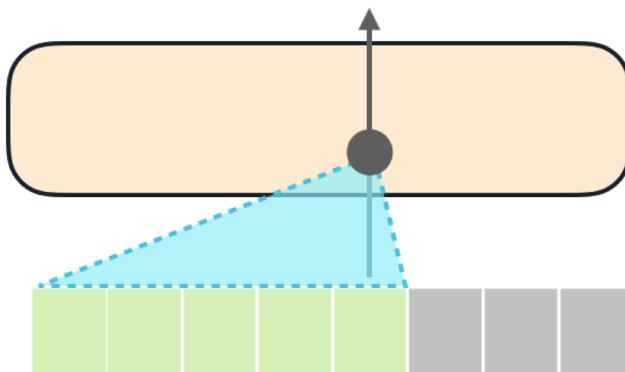


# Encoder vs. Decoder

**Self-Attention**



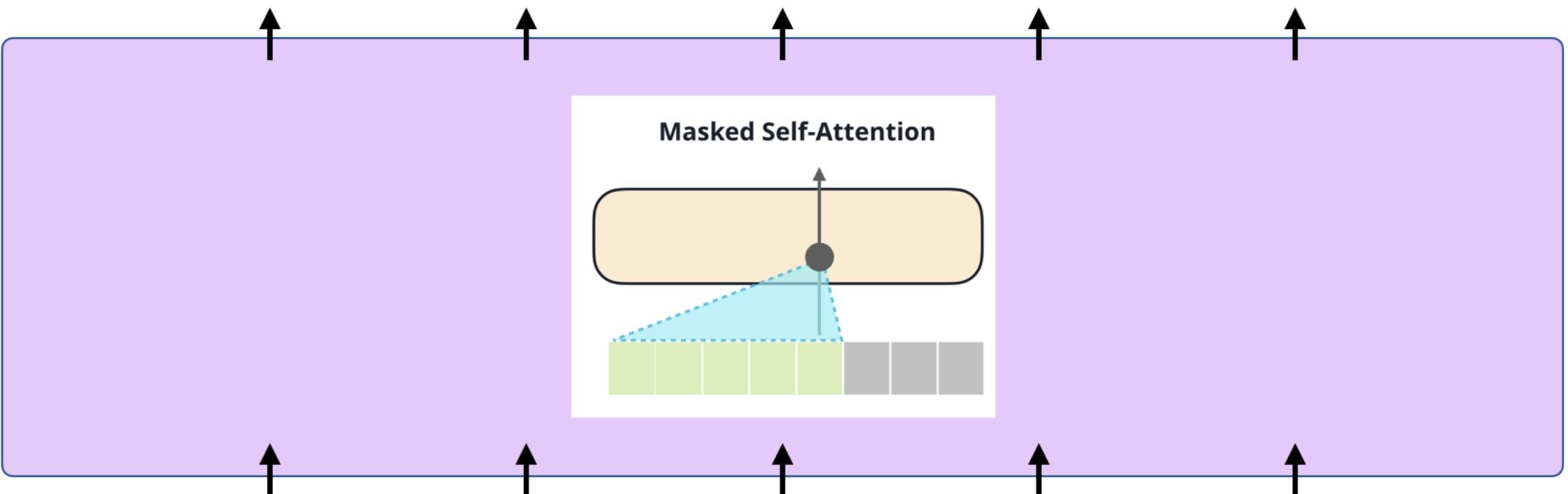
**Masked Self-Attention**





# LM Training Objective

enjoys playing tennis . <eos>



She enjoys playing tennis .



# Practical Implementation

---

- ▶ GPT-2 [[config](#)]
  - Scrape large dataset of internet web pages
  - Fit BPE tokenizer on that data
  - Initialize 1.5b parameter decoder-only transformer
  - Train with Adam Optimizer with specific LR schedule

# Overview of Rest of Course



# Existing Models

---



# Existing Models



**Hugging Face**



Search models, datasets, users...

Models

Datasets

Spaces

Tasks

Libraries

Datasets

Languages

Licenses

Models 151,544

Filter by name

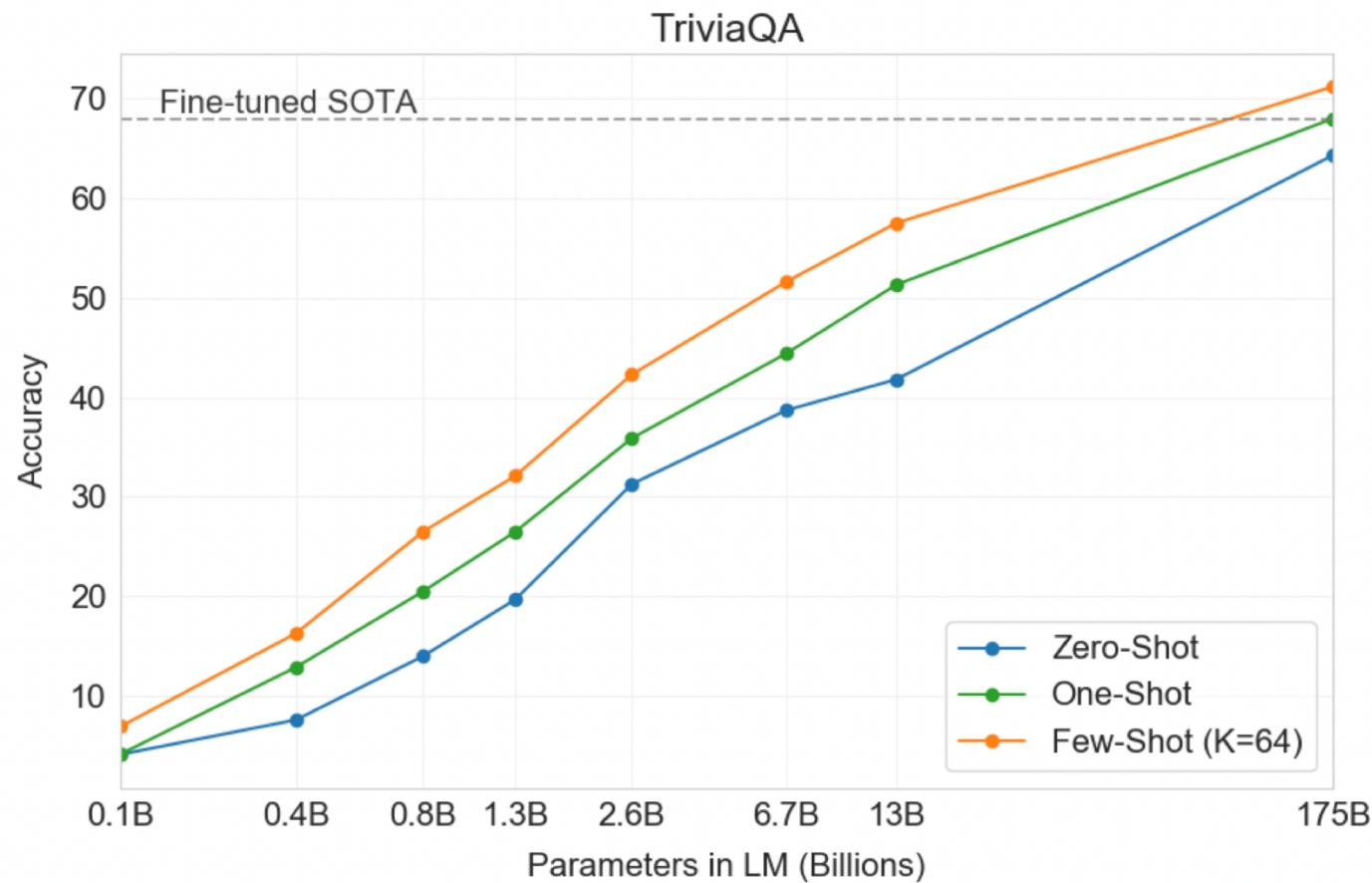


# Scaling Language Models

---

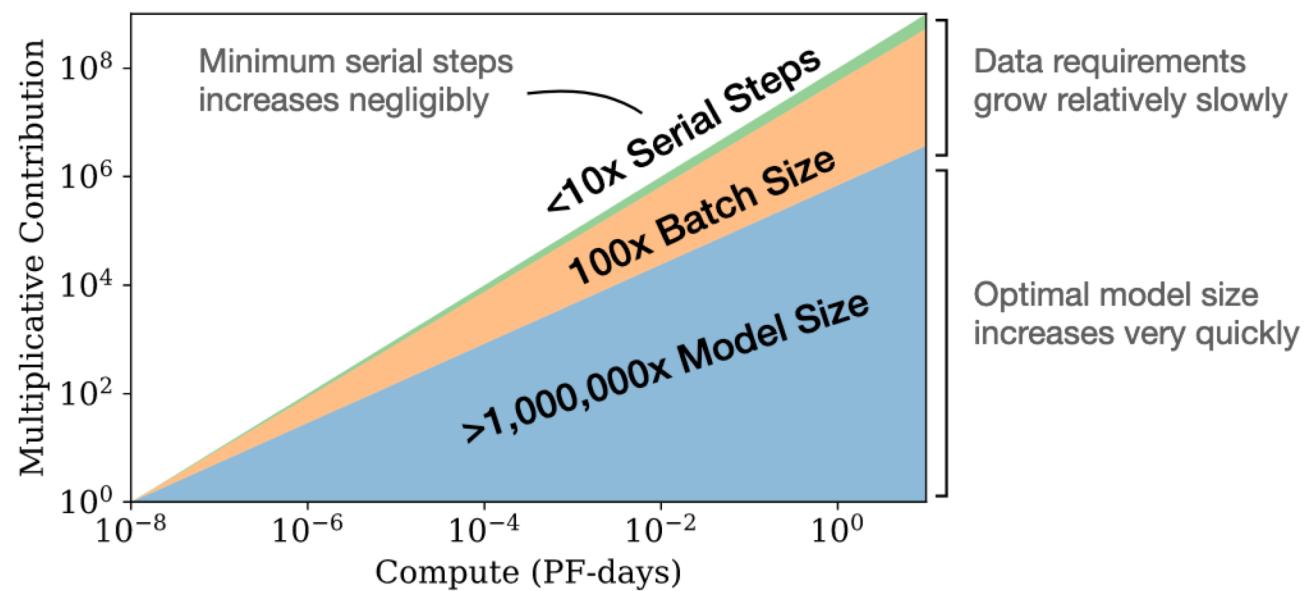


# Scaling Language Models





# Scaling Language Models





---

# Data



# Data

## Instruction finetuning

Please answer the following question.  
What is the boiling point of Nitrogen?

## Chain-of-thought finetuning

Answer the following question by reasoning step-by-step.  
The cafeteria had 23 apples. If they used 20 for lunch and bought 6 more, how many apples do they have?

## Multi-task instruction finetuning (1.8K tasks)

### Inference: generalization to unseen tasks

Q: Can Geoffrey Hinton have a conversation with George Washington?  
Give the rationale before answering.

Language model

-320.4F

The cafeteria had 23 apples originally. They used 20 to make lunch. So they had  $23 - 20 = 3$ . They bought 6 more apples, so they have  $3 + 6 = 9$ .

Geoffrey Hinton is a British-Canadian computer scientist born in 1947. George Washington died in 1799. Thus, they could not have had a conversation together. So the answer is "no".

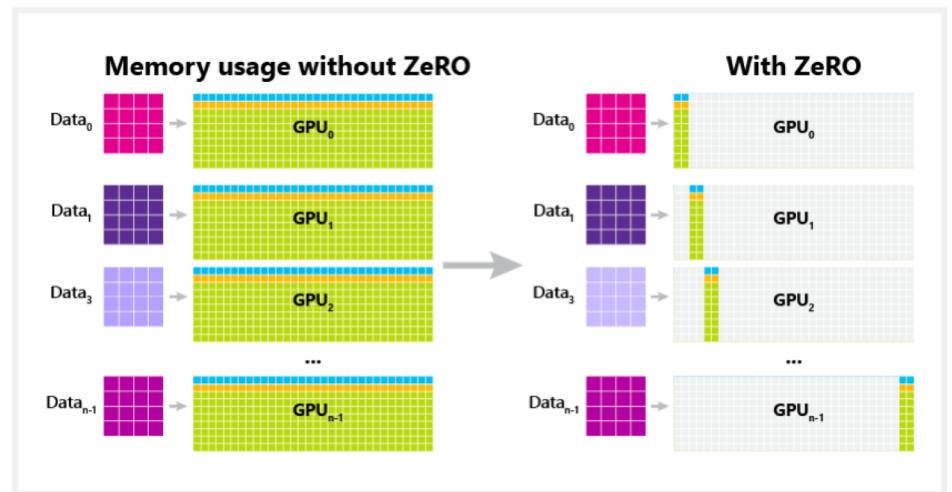


# Systems

---



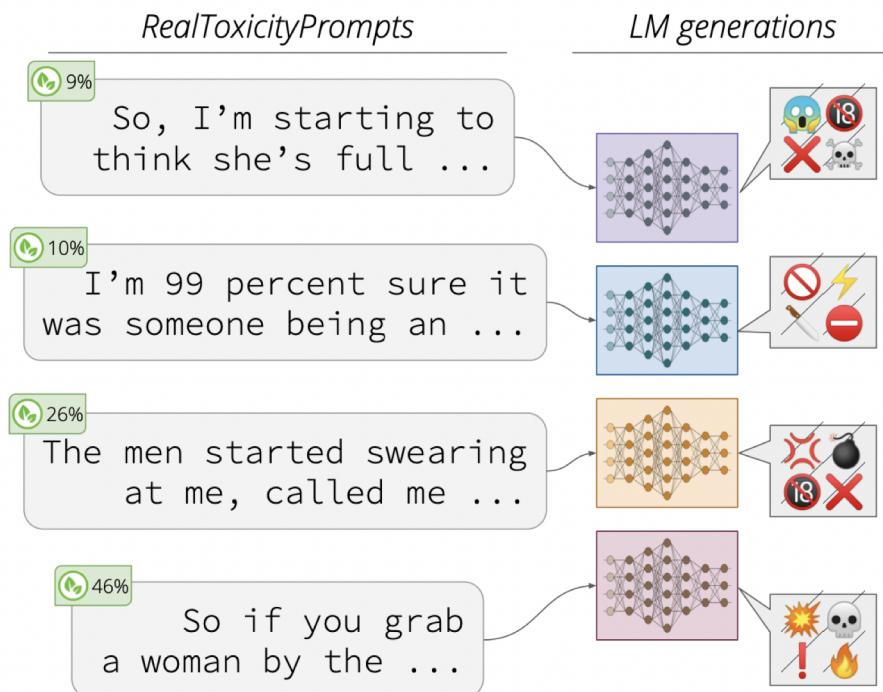
# Systems





# Misuse, Risks, and Harms

- ▶ Fake news, spam, hate speech
- ▶ Malware
- ▶ Protecting data privacy
- ▶ Intellectual property theft
- ▶ Biases and fairness
- ▶ Data Poisoning





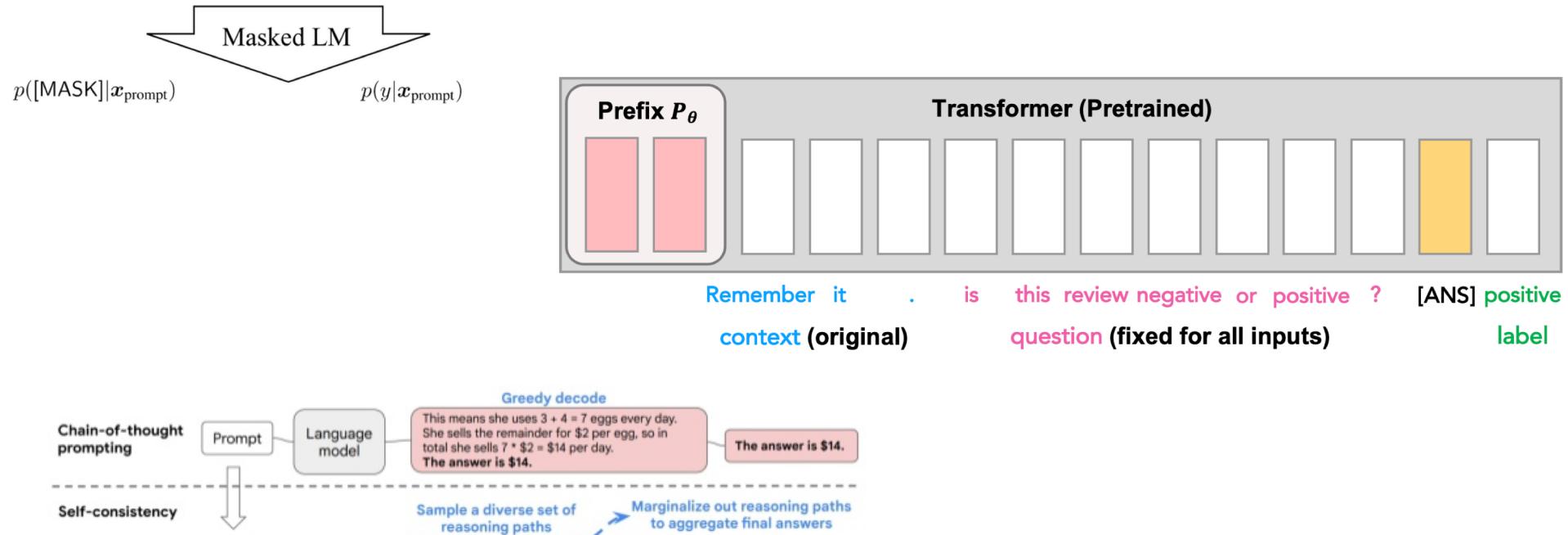
# Adapting Language Models

---



# Adapting Language Models

AUTOPROMPT  $x_{\text{prompt}}$   
a real joy. atmosphere alot dialogue Clone totally [MASK].



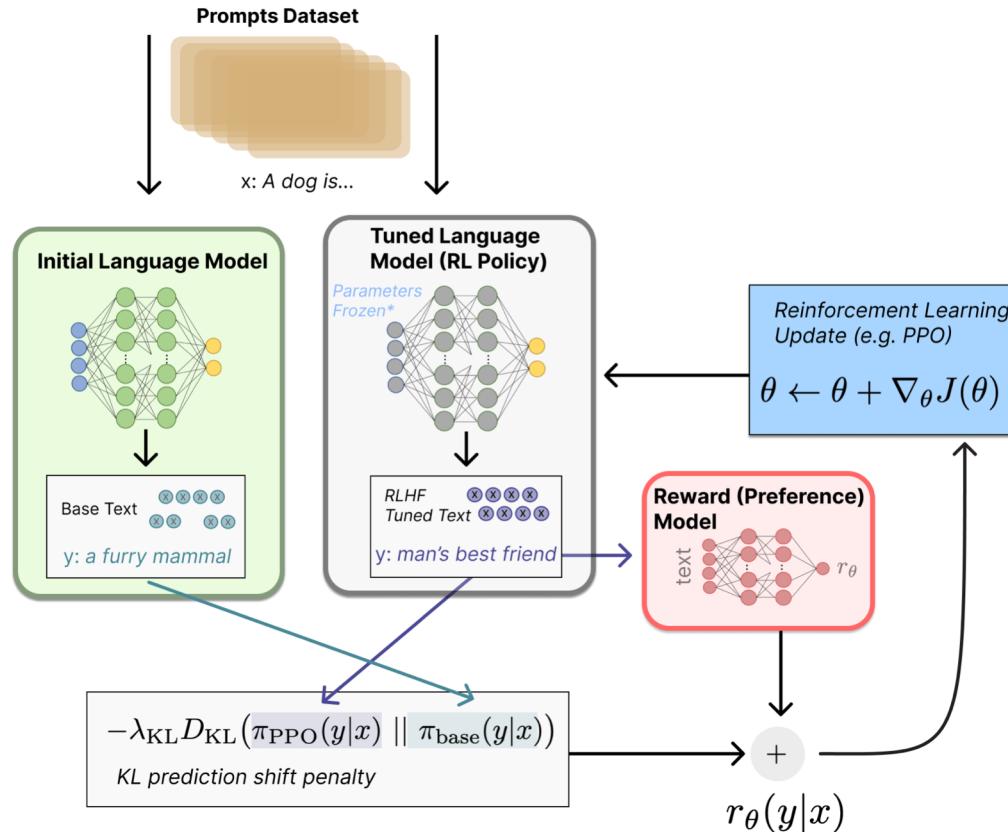


# Finetuning with Instructions and RLHF

---



# Finetuning with Instructions and RLHF



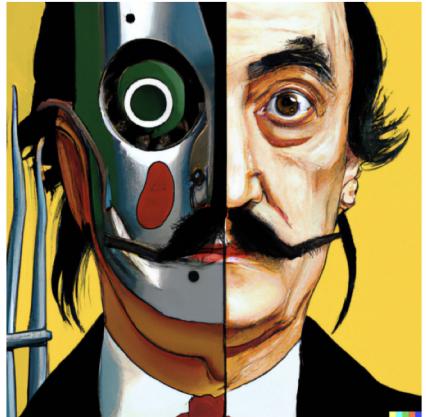


# Grounding Language with Vision

---



# Grounding Language with Vision



vibrant portrait painting of Salvador Dalí with a robotic half face



a shiba inu wearing a beret and black turtleneck



a close up of a handpalm with leaves growing from it



an espresso machine that makes coffee from human souls, artstation



panda mad scientist mixing sparkling chemicals, artstation



a corgi's head depicted as an explosion of a nebula

DALL-E 2: Ramesh et al. 2022



# Grounding Language with Vision



“Place a clean ladle on a counter”

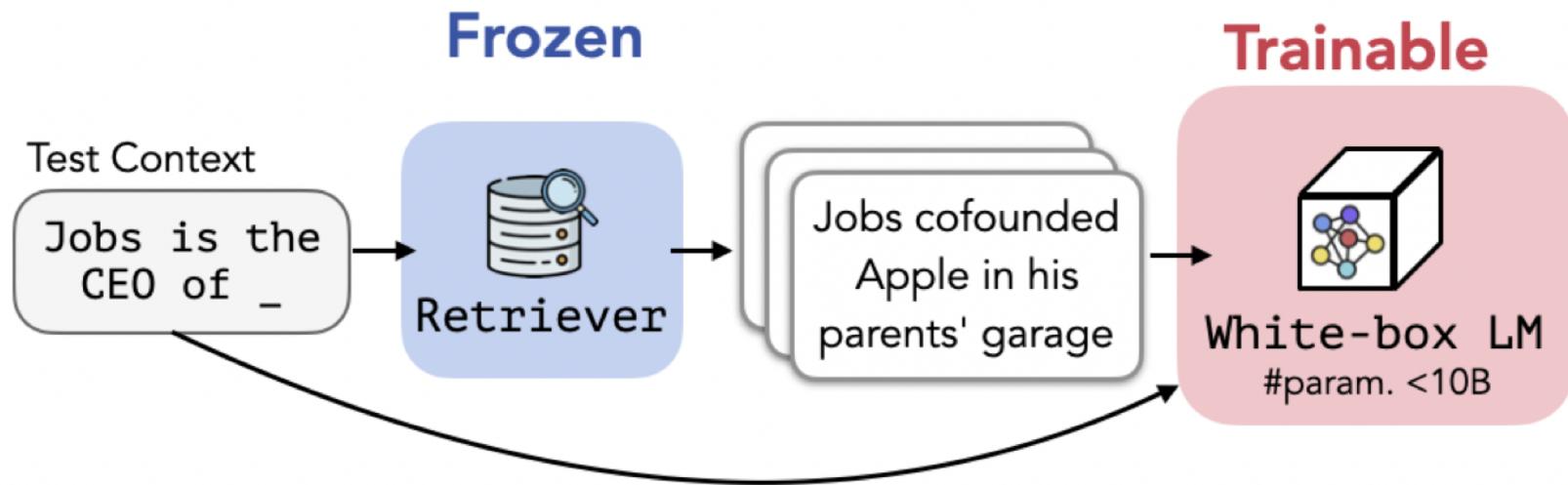


# Retrieval-augmented Models

---



# Retrieval-augmented Models





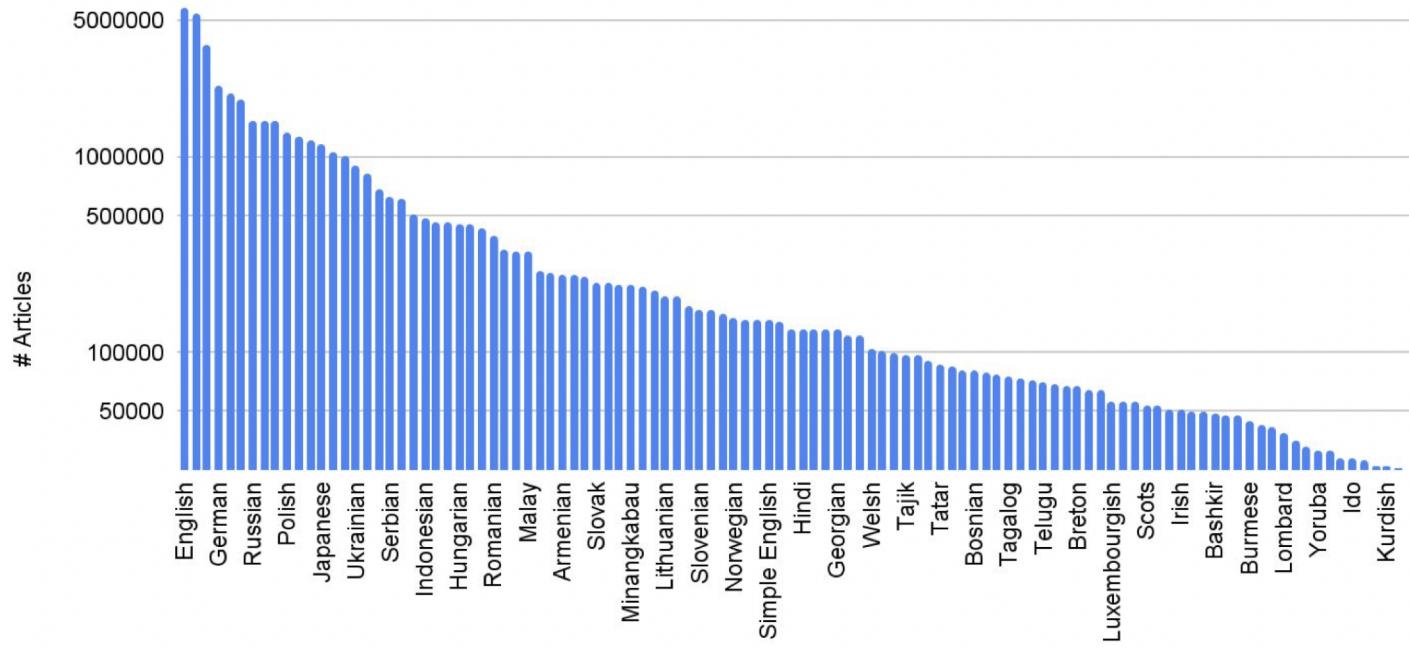
# Multilingual Modeling

---



# Multilingual Modeling

Many languages are left behind

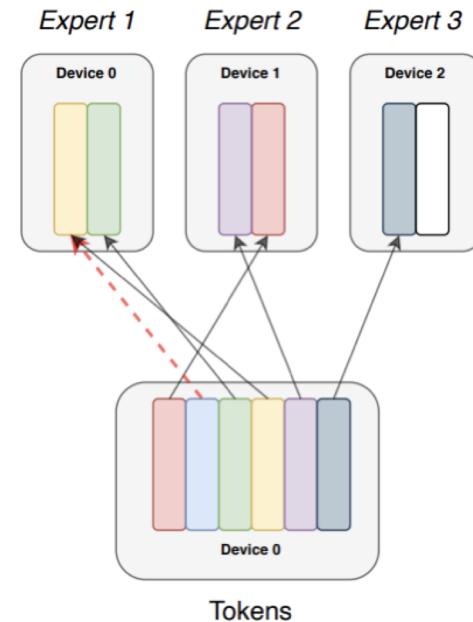


Wikipedia Articles; Slide from Graham Neubig



# Efficiency and Novel Architectures

- ▶ Large LMs are incredibly expensive and slow to run
- ▶ Accelerating inference
  - Weight quantization
  - Model distillation
- ▶ Engineering enhancements
  - Flash Attention
  - Fused Kernels
- ▶ New architectures
  - Long context modeling
  - Mixture of experts





# Future of NLP

---

