

Overview and Transformer Language Models

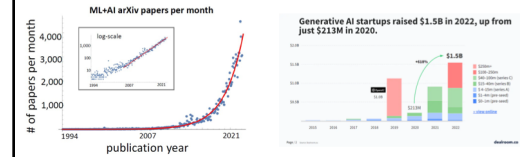


Eric Wallace
CS 288, 3/13/2023

Logistics

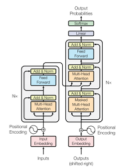
- 4 traditional lectures + ~8 days of mixed lectures and panels/discussions
- HW4 out Wednesday. Due Wednesday after spring break
 - Using and finetuning LMs with Huggingface
- HW5 out after spring break. Due sometime end of April.
 - Prompting ChatGPT to solve projects 1-3
- No final exam.
- Lecture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

Immense Interest

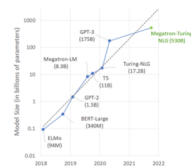


The Era of Rapid Scaling in NLP

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

Neural network

Prompt: The mouse ate the

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

Neural Language Models

Neural network

Prompt: The mouse ate the

Token	Prob
cheese	0.20
cookie	0.12
nibble	0.08
crumb	0.07
man	0.05
tail	0.04
...	

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

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 Model

Prompt

Question: What is the sentiment of the sentence "Superb acting"?

Answer:

Token	Prob
positi	0.82
negati	0.10
yes	0.006
hello	0.005
acting	0.003
amazin	0.001
g	...

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

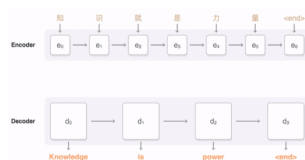
- Input encoding



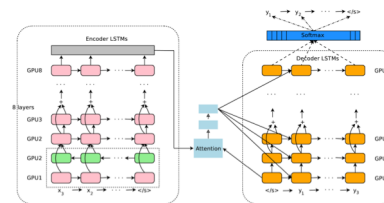
Language Models and MT Circa 2016

Neural Machine Translation is in production at Google

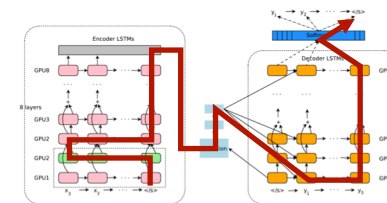
[Wu+16] [Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation](#)



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 the students opened their _____

discard fixed window

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

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

Word Window Neural Nets

output distribution
 $\hat{y} = \text{softmax}(Uh + b_2) \in \mathbb{R}^{|V|}$

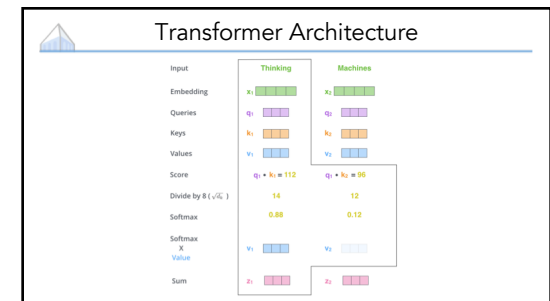
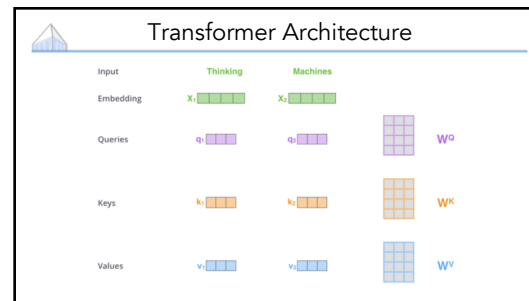
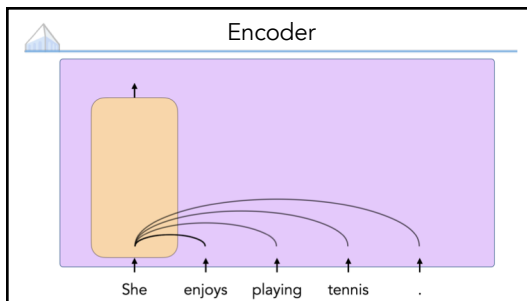
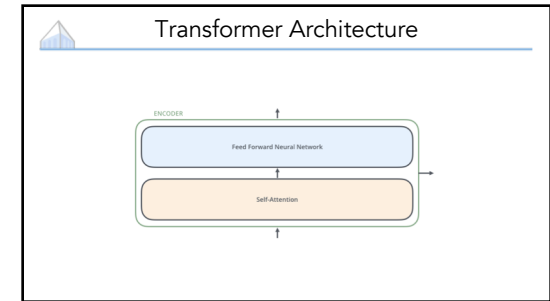
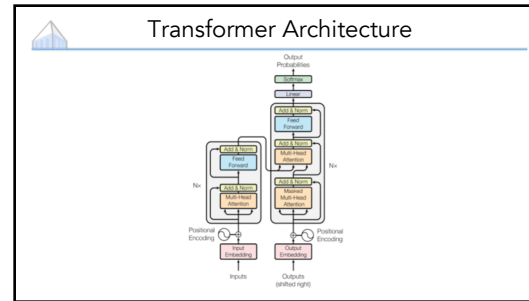
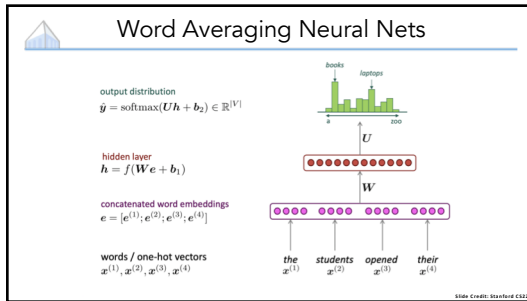
hidden layer
 $h = f(We + b_1)$

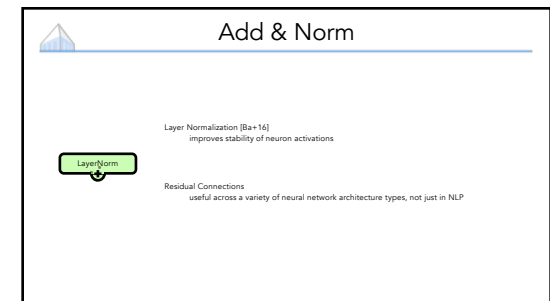
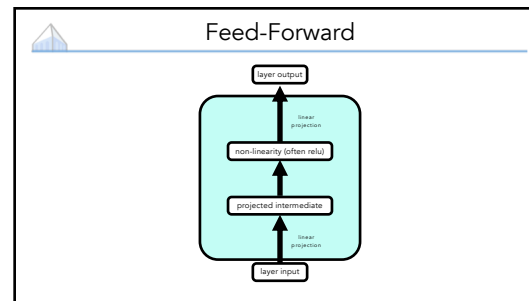
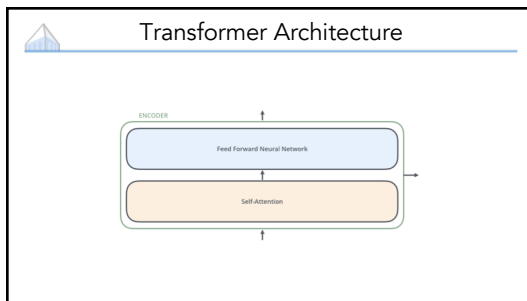
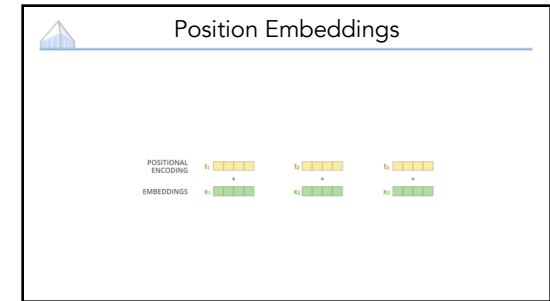
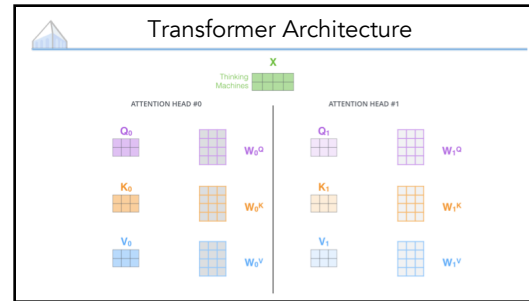
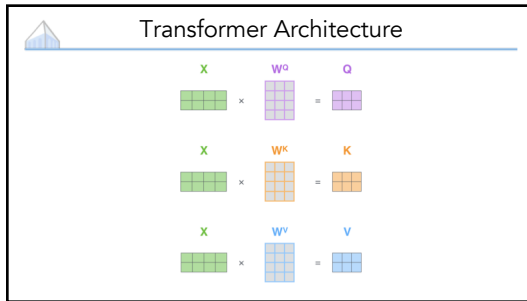
concatenated word embeddings
 $e = [e^{(1)}, e^{(2)}, e^{(3)}, e^{(4)}]$

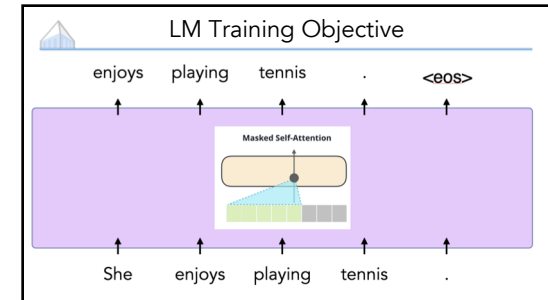
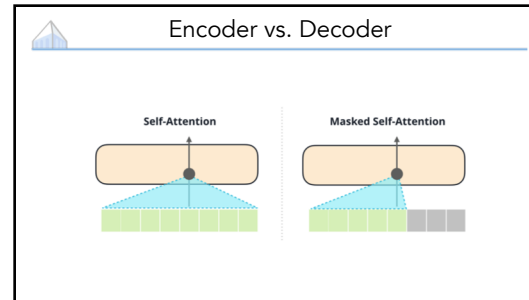
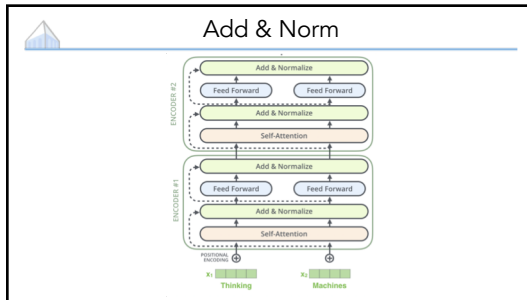
words / one-hot vectors
 $x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$

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

books laptops







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

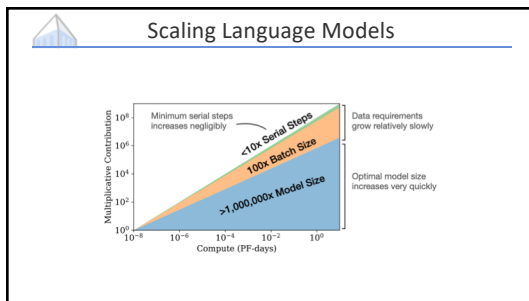
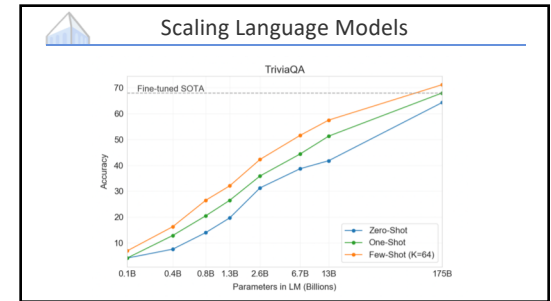
Overview of Rest of Course

Existing Models

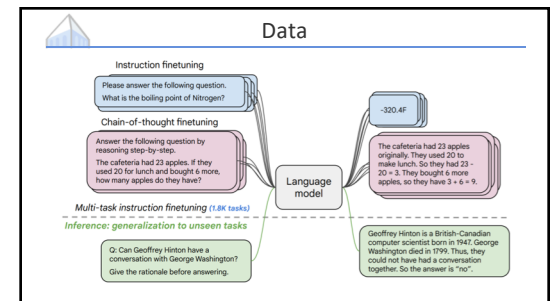
Existing Models

The screenshot shows the Hugging Face website interface. At the top, there's a search bar with the text "Search models, datasets, users...". Below the search bar, there are three tabs: "Models", "Datasets", and "Spaces". The "Models" tab is selected, showing a count of "Models 151,544". There's also a "Filter by name" input field.

Scaling Language Models



Data



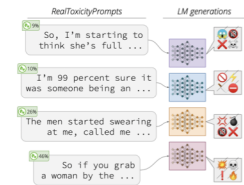
Systems

Systems



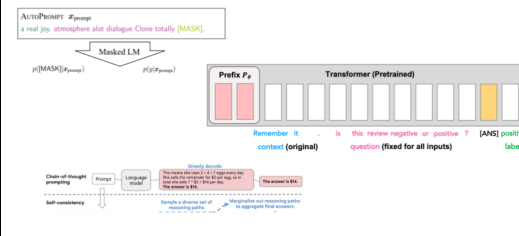
Misuse, Risks, and Harms

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



Adapting Language Models

Adapting Language Models



Finetuning with Instructions and RLHF

