

DSME 6635: Artificial Intelligence for Business Research

Deep-Learning-based NLP: Pretraining

Renyu (Philip) Zhang

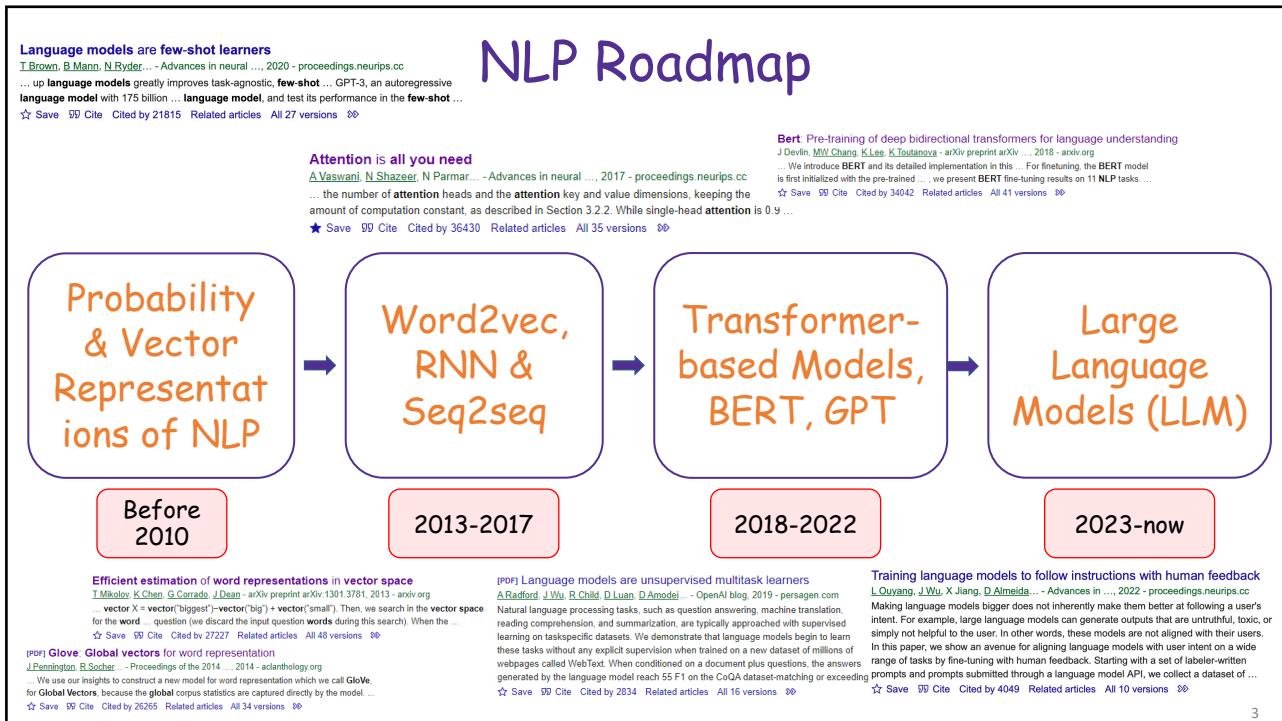
1

Agenda

- BERT: Bidirectional Encoder Representations from Transformers
- GPT: Generative Pretrained Transformers

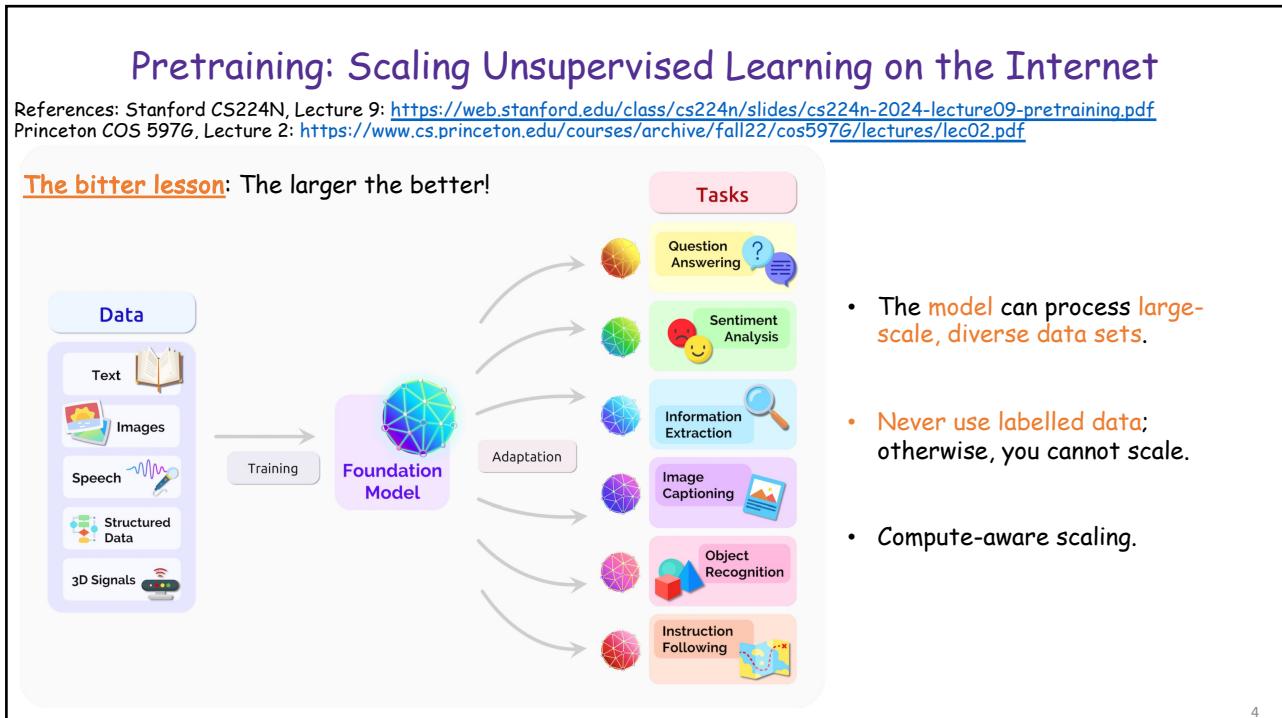
2

2



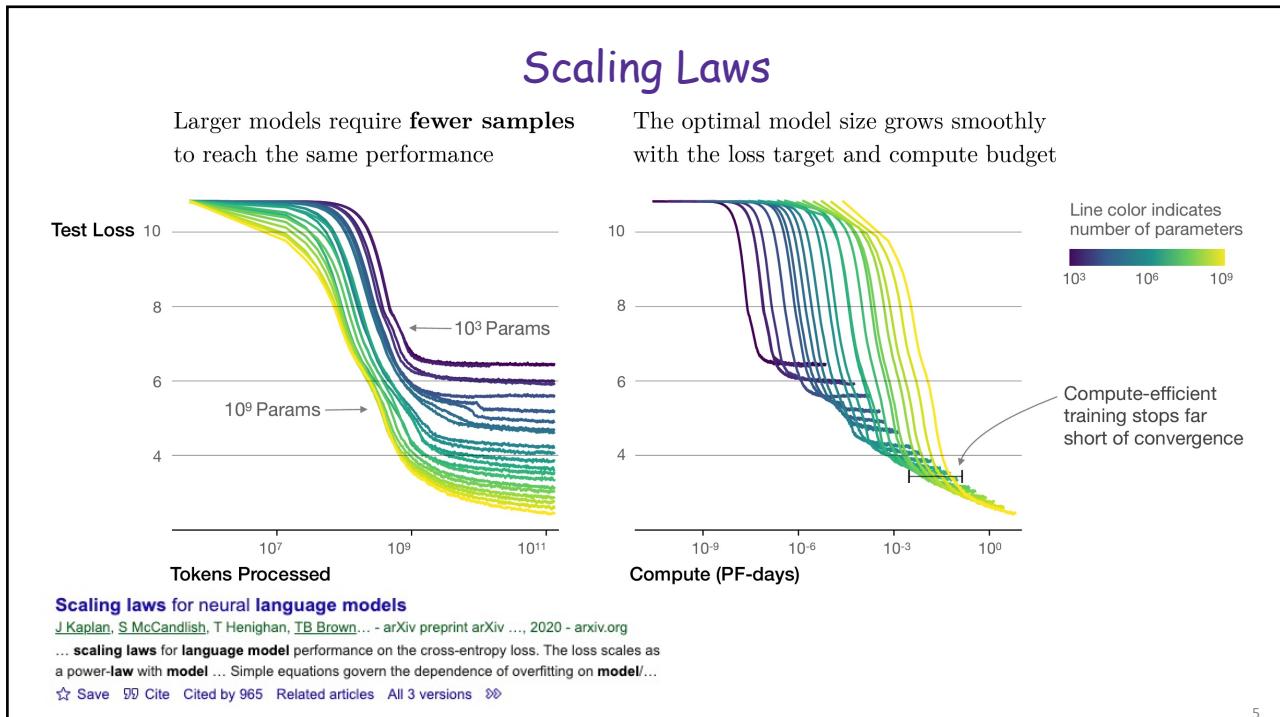
3

3

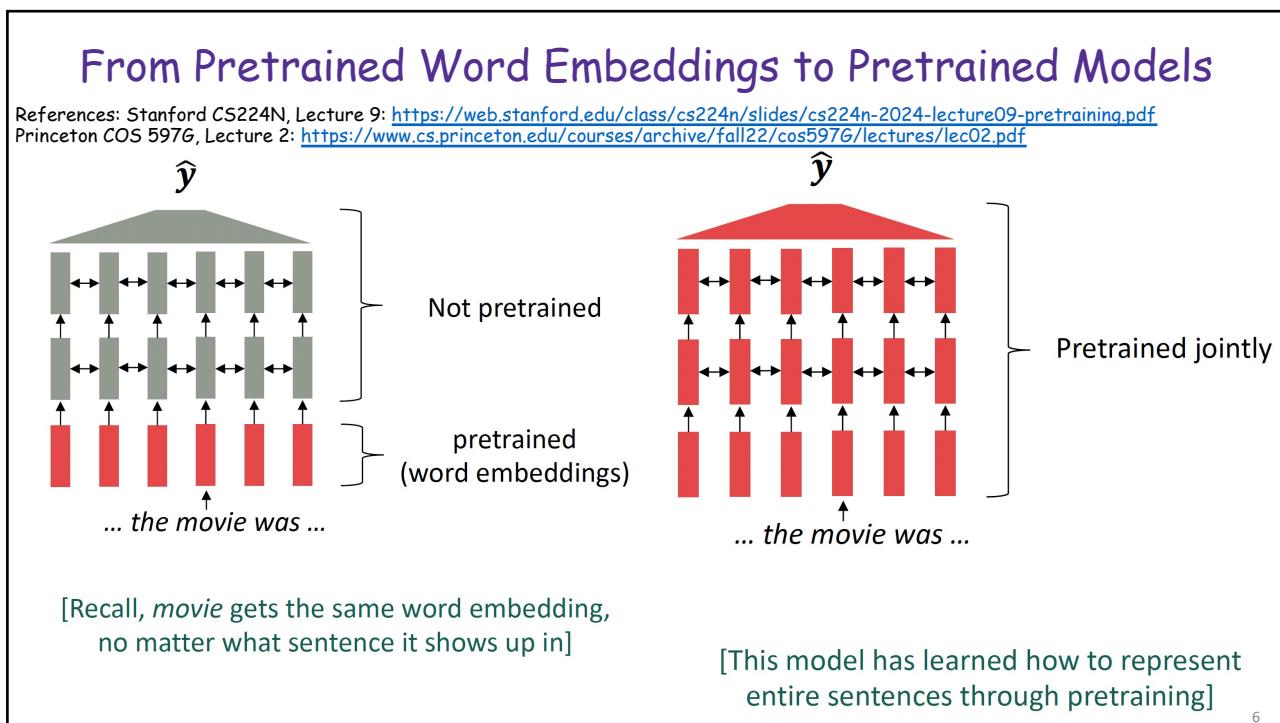


4

4



5



6

6

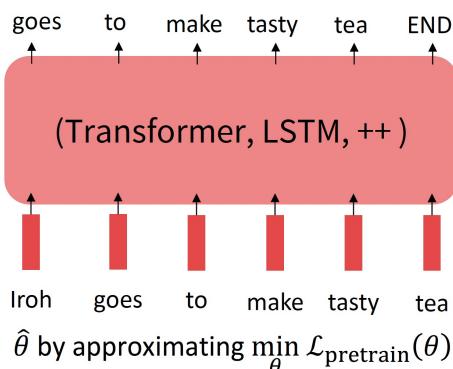
From Pretrained Word Embeddings to Pretrained Models

References: Stanford CS224N, Lecture 9: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture09-pretraining.pdf>
 Princeton COS 597G, Lecture 2: <https://www.cs.princeton.edu/courses/archive/fall22/cos597G/lectures/lec02.pdf>

- Pretraining can improve downstream NLP applications by serving as **parameter initialization**.

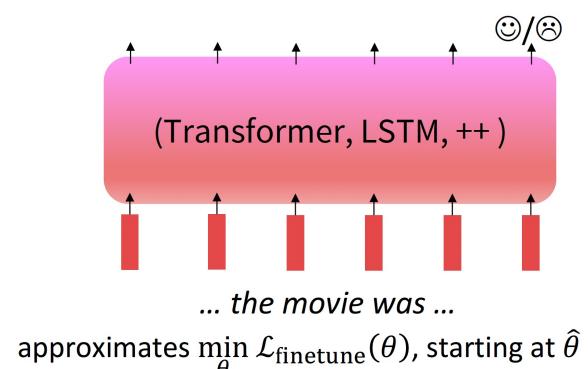
Step 1: Pretrain (on language modeling)

Lots of text; learn general things!



Step 2: Finetune (on your task)

Not many labels; adapt to the task!

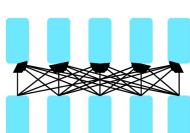


7

7

Three Pretraining Architectures

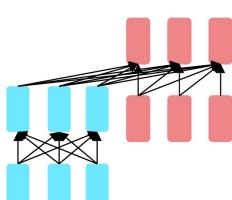
References: Stanford CS224N, Lecture 9: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture09-pretraining.pdf>
 Princeton COS 597G, Lecture 2: <https://www.cs.princeton.edu/courses/archive/fall22/cos597G/lectures/lec02.pdf>



Encoders

- Can condition on future.

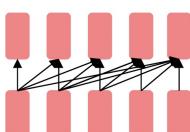
- Example: BERT.



Encoder-Decoders

- Combining encoder and decoder.

- Example: T5



Decoders

- Cannot condition on future.

- Example: GPT

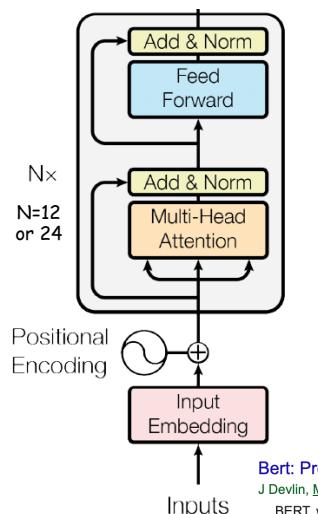
- All (very) large language models are decoders.

8

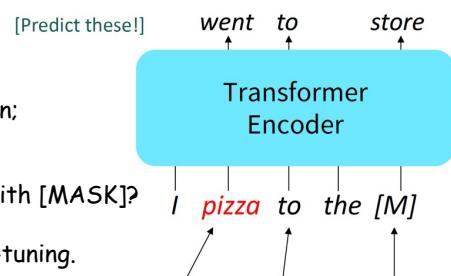
8

BERT: Bidirectional Encoder Representations from Transformers

References: Stanford CS224N, Lecture 9: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture09-pretraining.pdf>
 Princeton COS 597G, Lecture 2: <https://www.cs.princeton.edu/courses/archive/fall22/cos597G/lectures/lec02.pdf>



- Key idea: Learn representations based on **bidirectional context**.
 - We went to the river **bank**. vs. I need to go to the **bank** to make a deposit.
- Pretraining objectives: **masked language modeling + next sentence prediction**
- 15% of tokens are randomly masked.
- The masked tokens in the inputs:
 - 80% replaced with [MASK];
 - 10% replaced with a random token;
 - 10% no change.
- Why not all masked tokens replaced with [MASK]?
- [MASK] tokens are never seen in fine-tuning.



Bert: Pre-training of deep **bidirectional** transformers for language understanding [Replaced] [Not replaced] [Masked]
 J Devlin, MW Chang, K Lee, K Toutanova - arXiv preprint arXiv ..., 2018 - arxiv.org
 ... BERT, which stands for **Bidirectional Encoder Representations** from Transformers. Unlike ...
 2018, BERT is designed to pretrain deep **bidirectional representations** from unlabeled text by ...

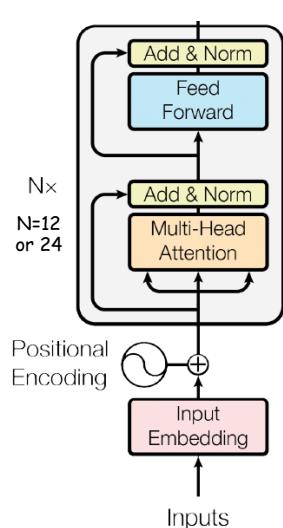
☆ Save 99 Cite Cited by 93230 Related articles All 46 versions ☰

9

9

Next Sentence Prediction (NSP)

References: Stanford CS224N, Lecture 9: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture09-pretraining.pdf>
 Princeton COS 597G, Lecture 2: <https://www.cs.princeton.edu/courses/archive/fall22/cos597G/lectures/lec02.pdf>



- Understanding the **relationships between two sentences** are also important.
- Reduce the gap between pretraining and finetuning.

[CLS]: a special token always at the beginning

[SEP]: a special token used to separate two segments

Input = [CLS] the man went to [MASK] store [SEP]

he bought a gallon [MASK] milk [SEP]

Label = IsNext

They sample two contiguous segments for 50% of the time and another random segment from the corpus for 50% of the time

Input = [CLS] the man [MASK] to the store [SEP]

penguin [MASK] are flight ##less birds [SEP]

Label = NotNext

10

10

Subwords and Input Embeddings

References: Stanford CS224N, Lecture 9: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture09-pretraining.pdf>
Princeton COS 597G, Lecture 2: <https://www.cs.princeton.edu/courses/archive/fall22/cos597G/lectures/lec02.pdf>



- To make sure training and testing vocabularies are consistent, uncommon words are split into components.

word	vocab mapping	embedding
Common words	hat → hat	█
Variations	learn → learn	█
misspellings	taaaaasty → taa## aaa## sty	█
novel items	laern → la## ern##	█
	Transformerify → Transformer## ify	█

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

Which of the two segments?

11

11

BERT Pretraining: Putting Together

References: Stanford CS224N, Lecture 9: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture09-pretraining.pdf>
Princeton COS 597G, Lecture 2: <https://www.cs.princeton.edu/courses/archive/fall22/cos597G/lectures/lec02.pdf>



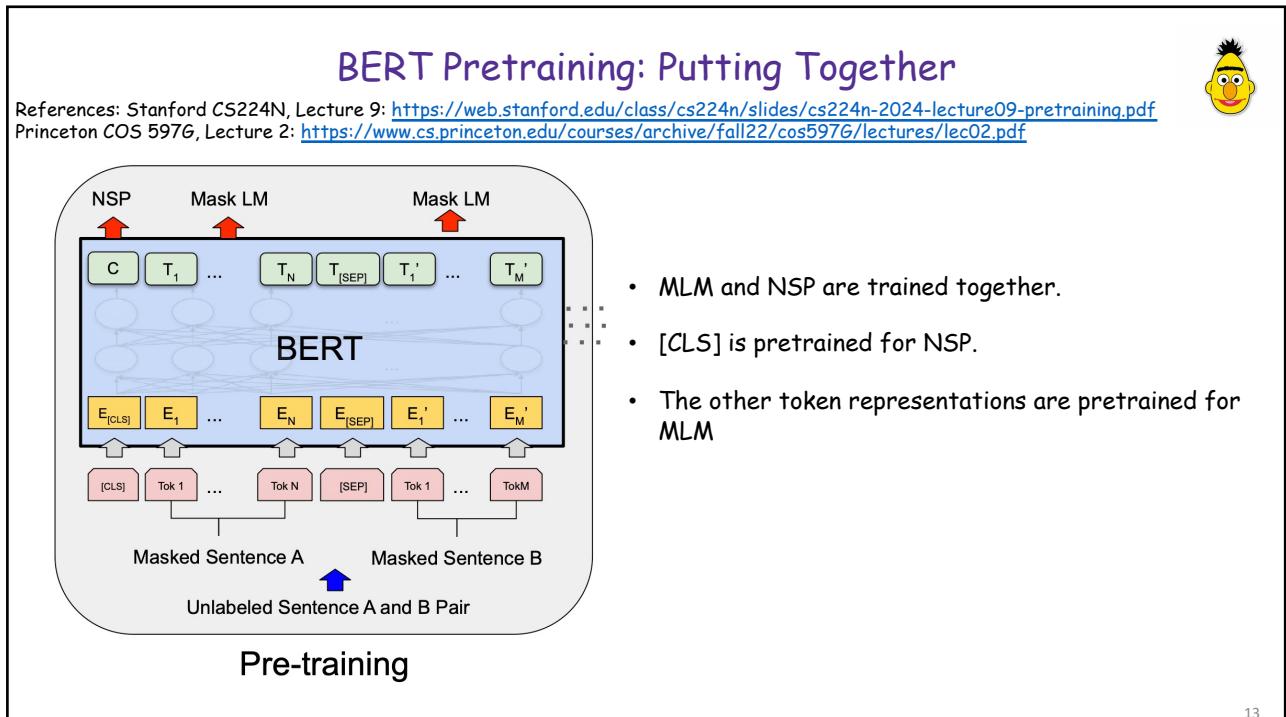
Nx
N=12 or 24

Input Embedding → Positional Encoding → N x (Multi-Head Attention → Add & Norm → Feed Forward → Add & Norm)

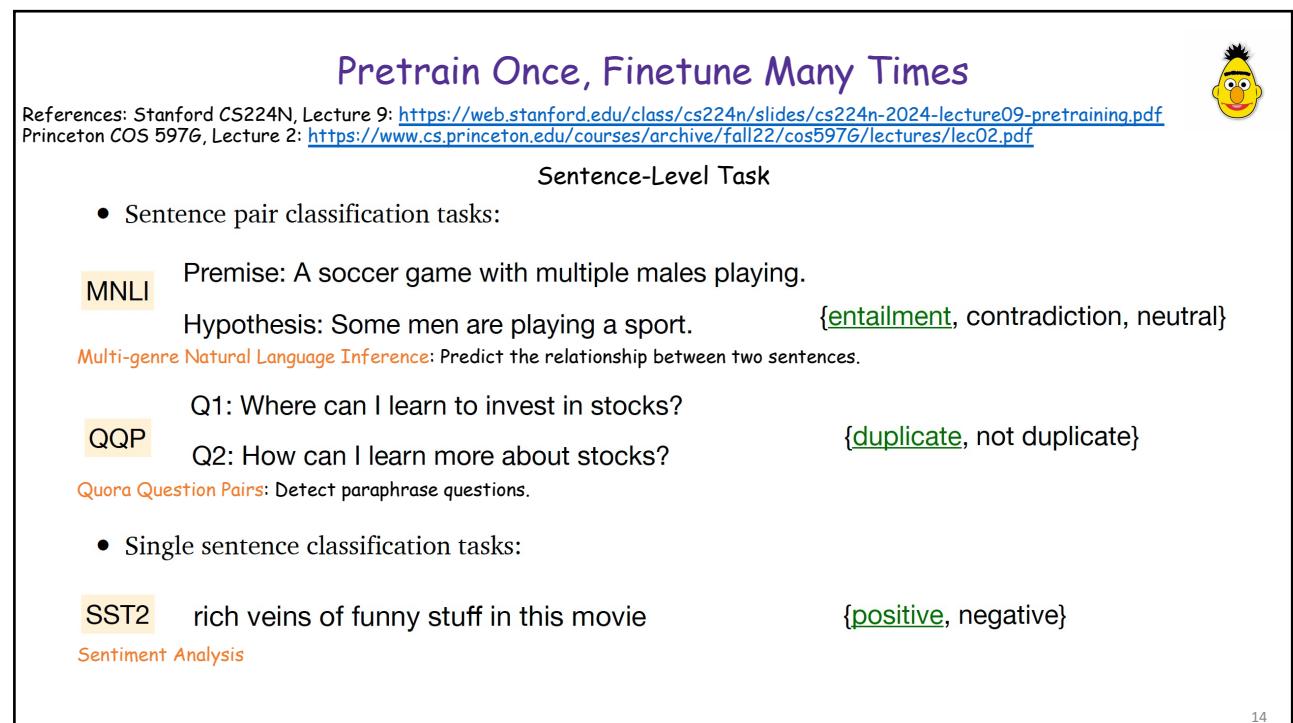
- BERT-base: 12 layers, 768-dim hidden size, 12 attention heads, 110M parameters
- BERT-large: 24 layers, 1024-dim hidden size, 16 attention heads, 340M parameters
- Trained on: Wikipedia (2.5B) + BookCorpus (0.8B)
- Max sequence size: 512 word pieces (roughly 256 + 256 non-contiguous sequences)
- Trained for 1M steps, batch size = 128K
- Pretrained with 64 TPUs for 4 days

12

12

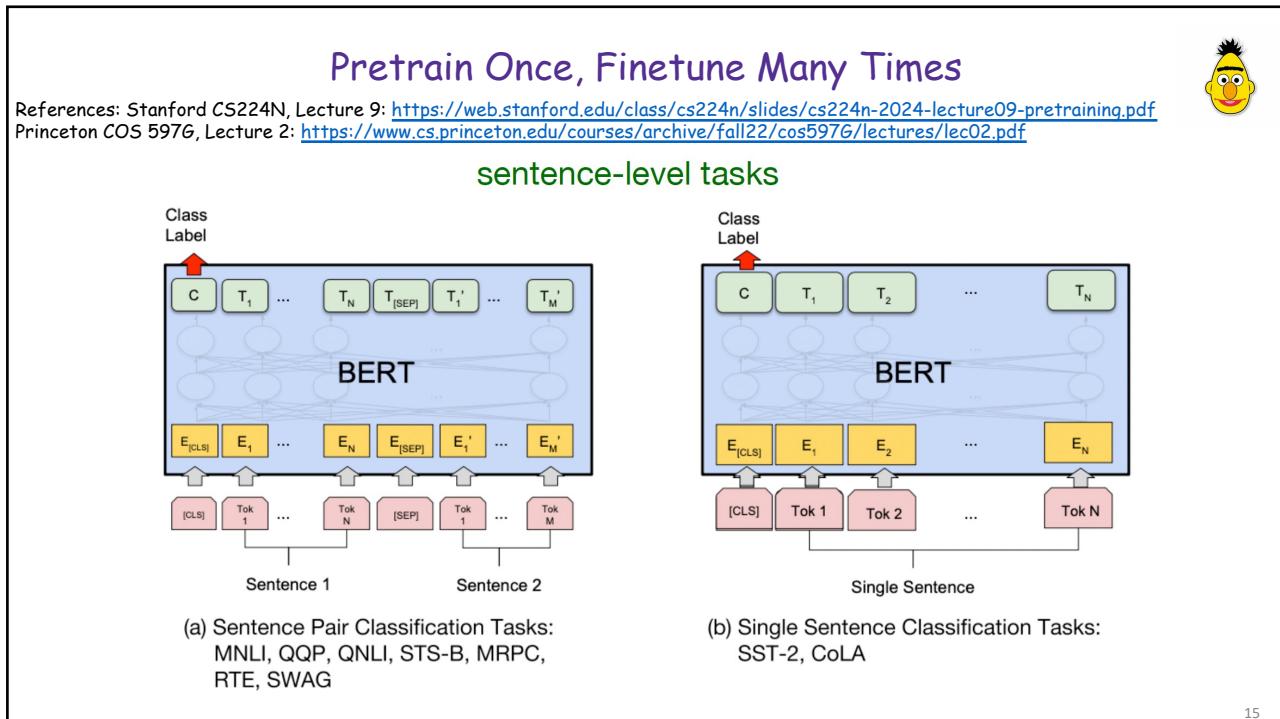


13

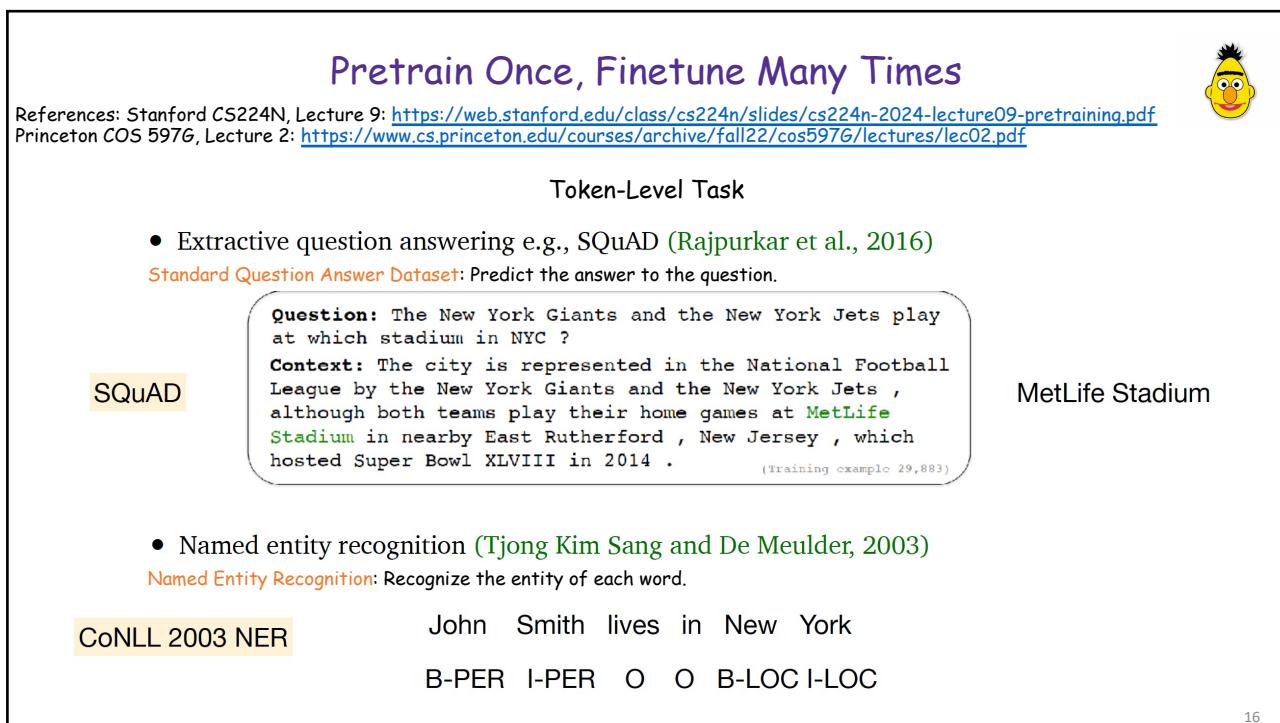


14

14



15

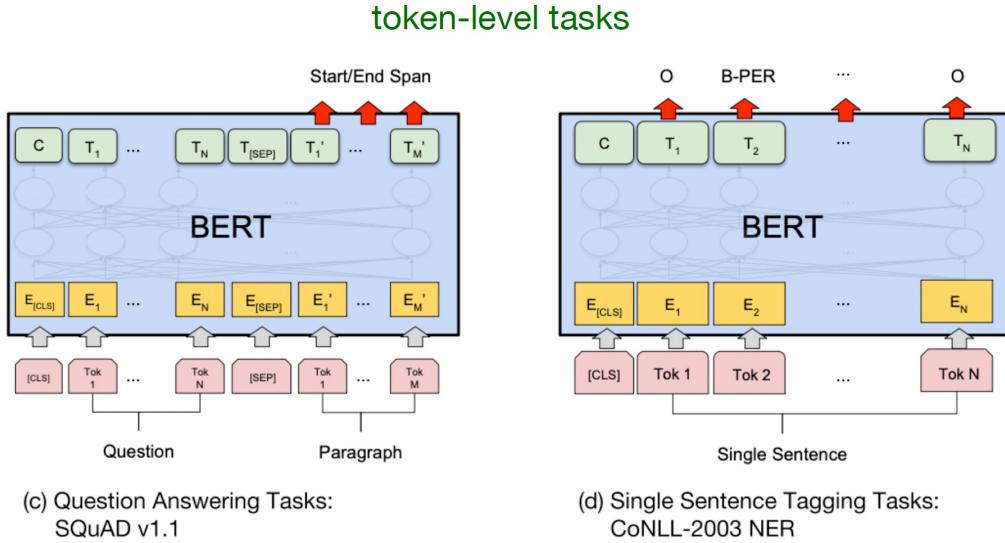


16

Pretrain Once, Finetune Many Times



References: Stanford CS224N, Lecture 9: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture09-pretraining.pdf>
Princeton COS 597G, Lecture 2: <https://www.cs.princeton.edu/courses/archive/fall22/cos597G/lectures/lec02.pdf>



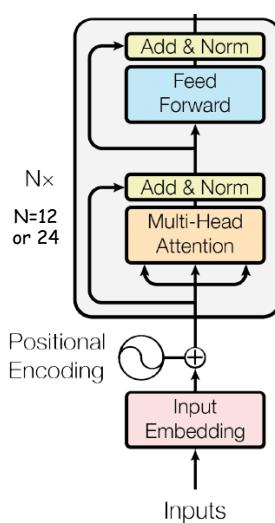
17

17

BERT was the State-of-The-Art



References: Stanford CS224N, Lecture 9: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture09-pretraining.pdf>
Princeton COS 597G, Lecture 2: <https://www.cs.princeton.edu/courses/archive/fall22/cos597G/lectures/lec02.pdf>



- BERT-base: 12 layers, 768-dim hidden size, 12 attention heads, 110M parameters
 - BERT-large: 24 layers, 1024-dim hidden size, 16 attention heads, 340M parameters

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average -
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT _{BASE}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1

- **Key issue with encoders:** Not a language model, i.e., does not naturally lead to autoregressive generation methods.

18

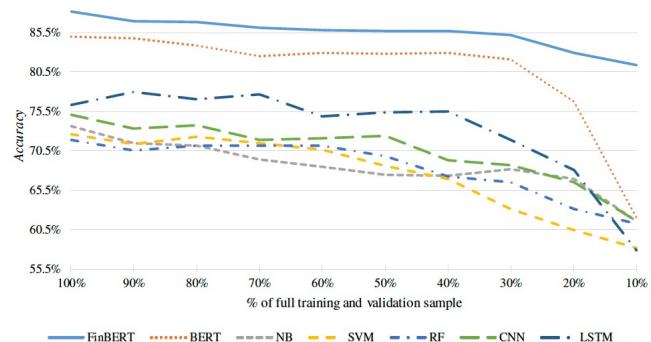
18

Revisit FinBERT



- Pretrain BERT-base using financial datasets (4.9B tokens in total) with 4 P100 GPUs (100G memory):
 - Corporate annual and quarterly filings from SEC's EDGAR website (1994-2019).
 - Financial analyst reports from Thomson Intestext database (2003-2012).
 - Earnings conference call transcripts from the Sekking Alpha website (2004-2019).
- Finetuning and evaluation:**
 - Sentiment analysis 10,000 sentences
 - 36% positive
 - 46% neutral
 - 18% negative
- Can FinBERT beat GPT-4 or Claude-3 in tasks related to financial texts?
 • How can we make **fair comparisons?**

Figure 1 Sentiment classification accuracy across sample sizes



FinBERT: A large language model for extracting information from financial text

AH Huang, H Wang, Y Yang - Contemporary Accounting ..., 2023 - Wiley Online Library

... model that adapts to the **finance** domain. We show that **FinBERT** incorporates **finance** knowledge and can better summarize contextual **information** in **financial texts**. Using a sample of ...

☆ Save ⌂ Cite Cited by 144 Related articles Web of Science: 22 ☰

19

19

Agenda

- BERT: Bidirectional Encoder Representations from Transformers
- GPT: Generative Pretrained Transformers

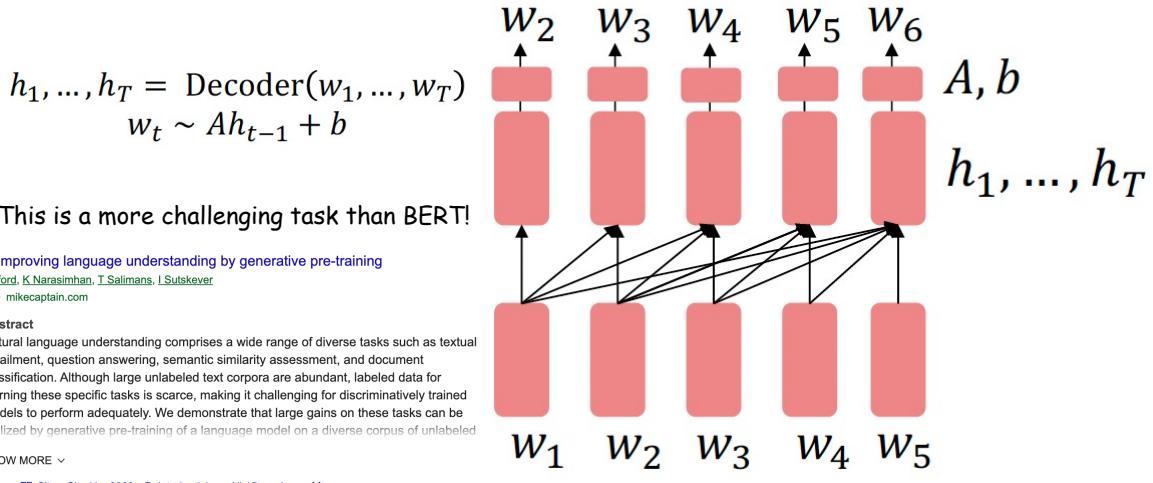
20

20

Pretraining Decoders

References: Stanford CS224N, Lecture 9: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture09-pretraining.pdf>

- Key idea: Pretrain decoders as language models $\Pr(W_n | W_1, W_2, \dots, W_{n-1})$ via autoregression.



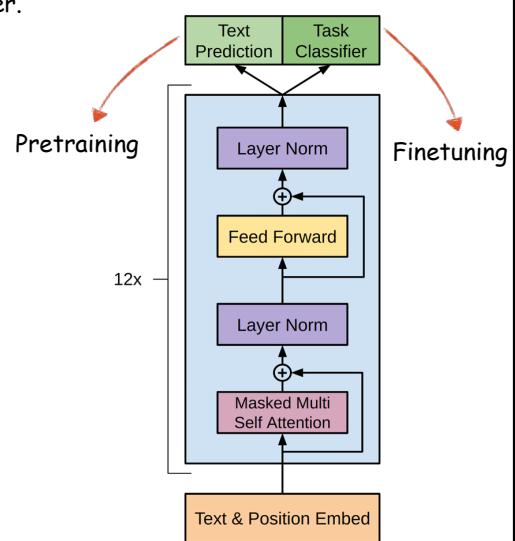
21

21

GPT-1

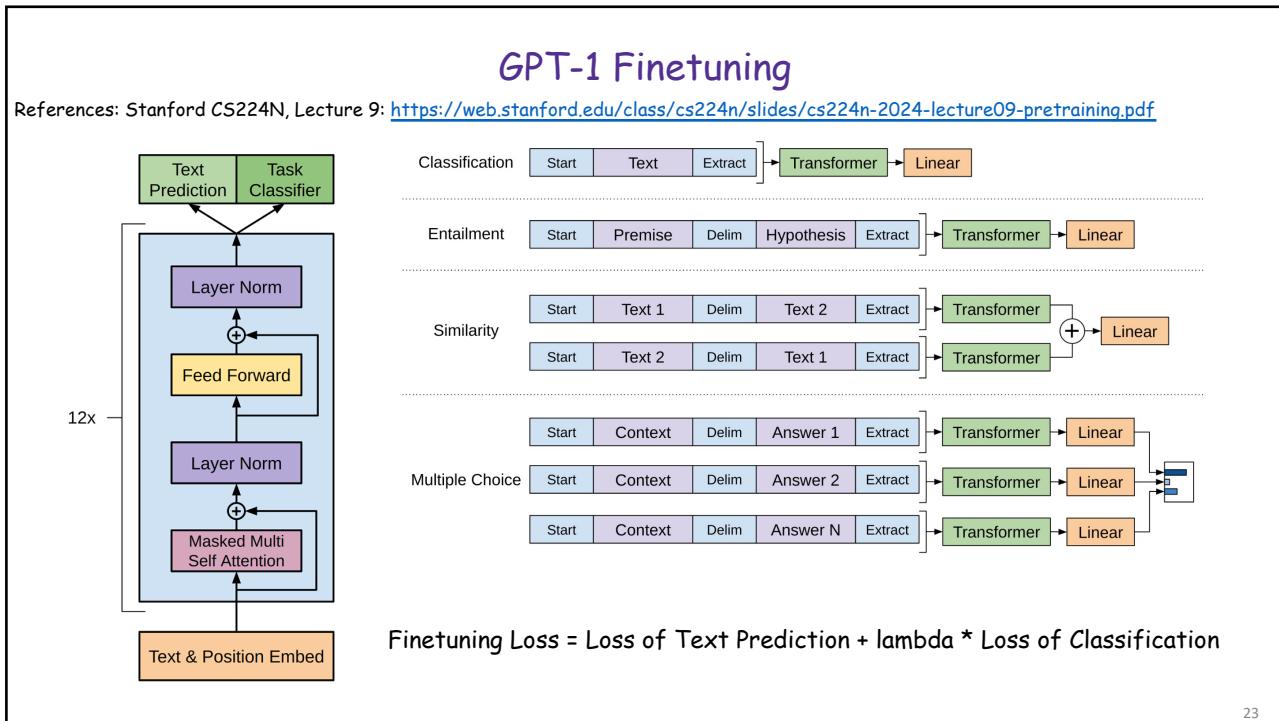
References: Stanford CS224N, Lecture 9: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture09-pretraining.pdf>

- Architecture: Only masked self-attention, but deeper and larger.
- 12 layers of transformer decoders, 117M parameters.
- 768-dim hidden states, 3072-dim MLP hidden layers.
- Byte-pair encoding with 40,000 merges.
- Trained on BooksCorpus of over 7,000 unique books.



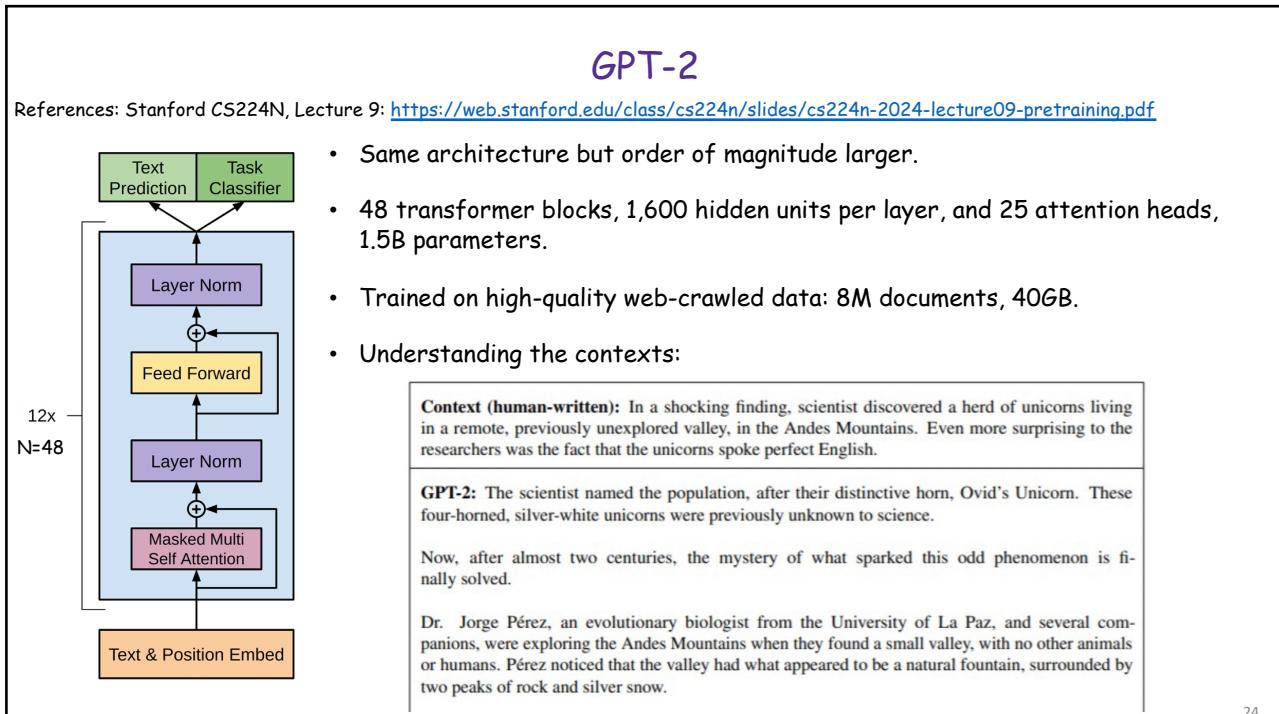
22

22



23

23



24

24

GPT-3

References: Stanford CS224N, Lecture 9: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture09-pretraining.pdf>
<https://www.cs.princeton.edu/courses/archive/fall22/cos597G/lectures/lec04.pdf>

