

DSME 6635: Artificial Intelligence for Business Research

Deep-Learning-based NLP: Pretraining

Renyu (Philip) Zhang

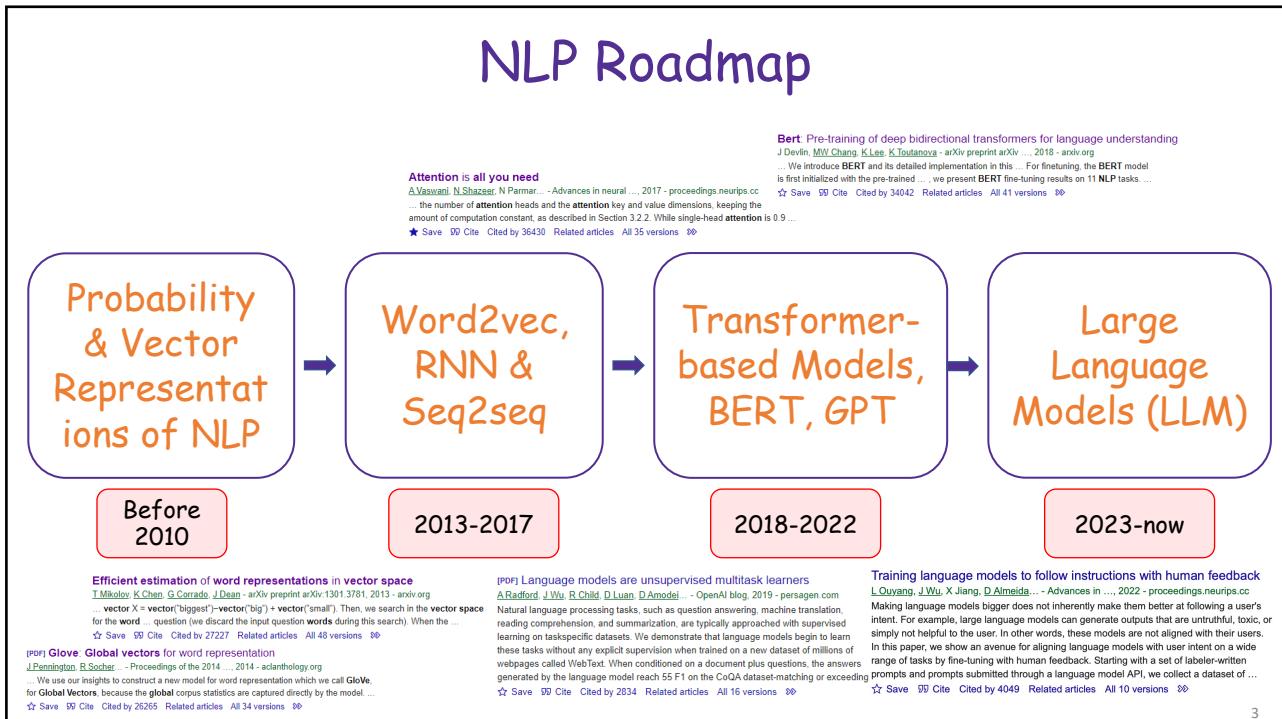
1

Agenda

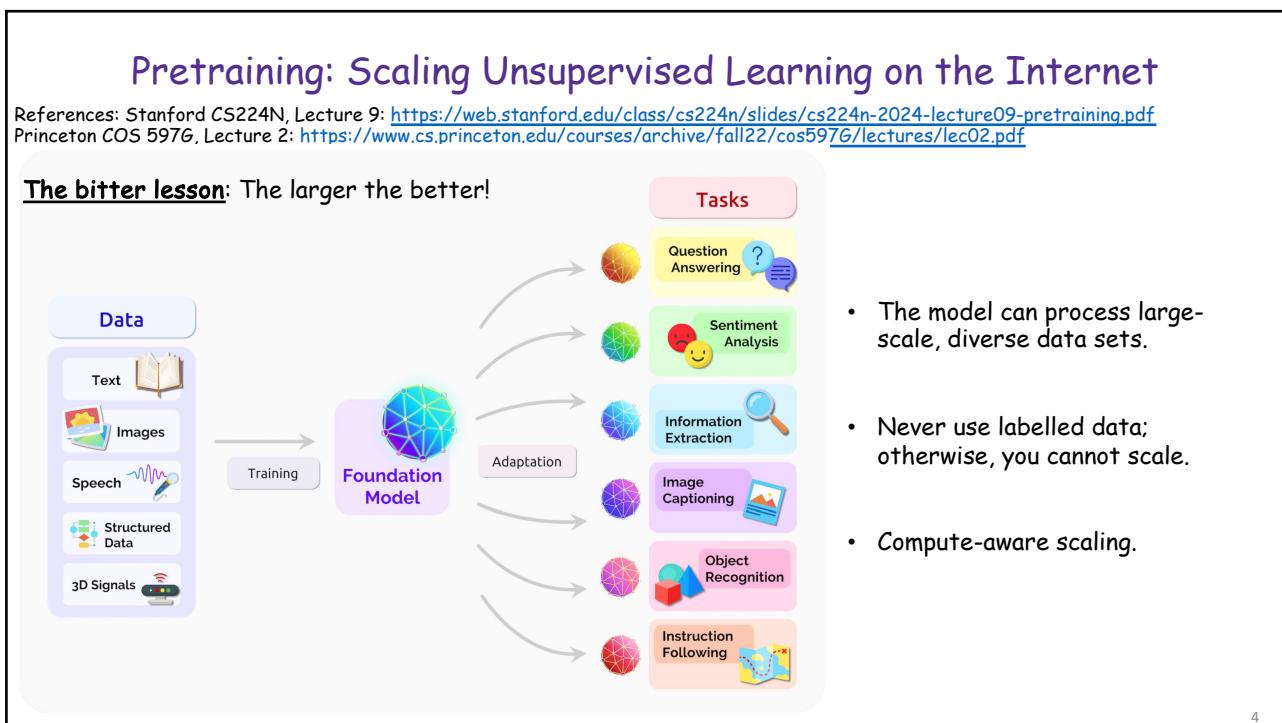
- BERT: Bidirectional Encoder Representations from Transformers
- GPT: Generative Pretrained Transformers

2

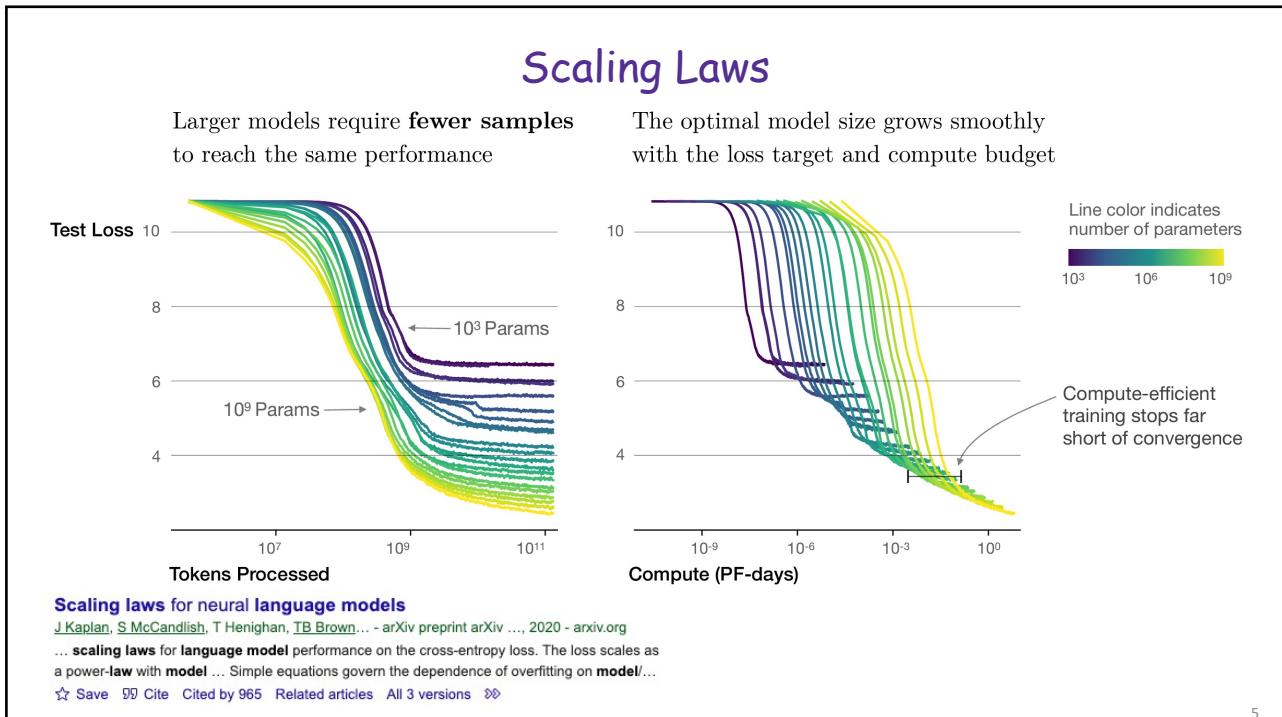
2



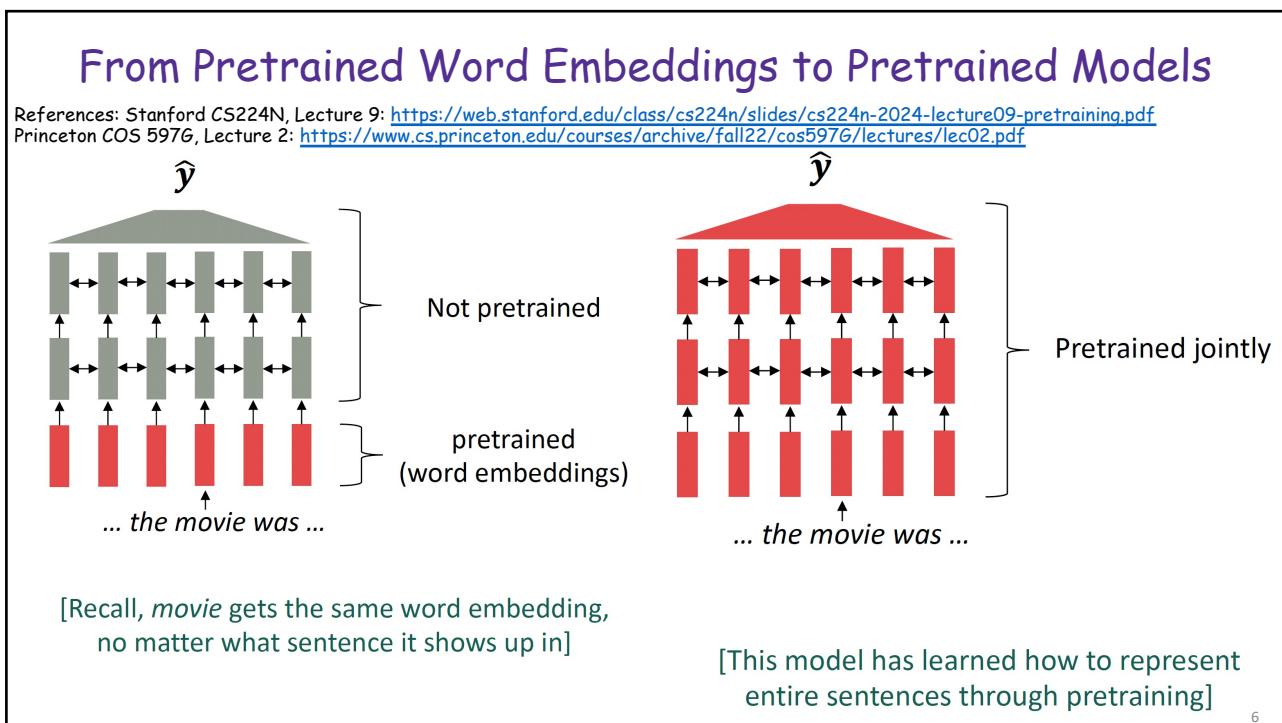
3



4



5



6

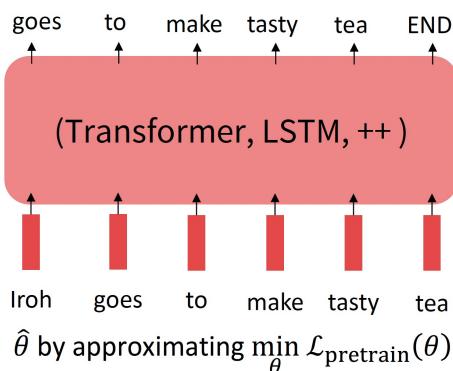
From Pretrained Word Embeddings to Pretrained Models

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**.

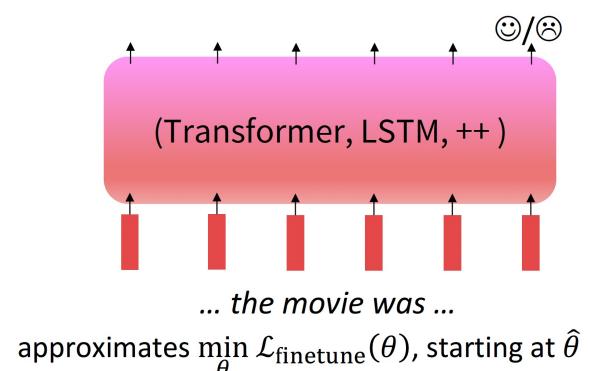
Step 1: Pretrain (on language modeling)

Lots of text; learn general things!



Step 2: Finetune (on your task)

Not many labels; adapt to the task!

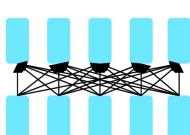


7

7

Three Pretraining Architectures

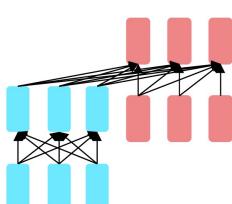
References: Stanford CS224N, Lecture 9: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture09-pretraining.pdf>
 Princeton COS 597G, Lecture 2: <https://www.cs.princeton.edu/courses/archive/fall22/cos597G/lectures/lec02.pdf>



Encoders

- Can condition on future.

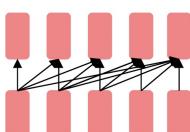
- Example: BERT.



Encoder-Decoders

- Combining encoder and decoder.

- Example: T5



Decoders

- Cannot condition on future.

- Example: GPT

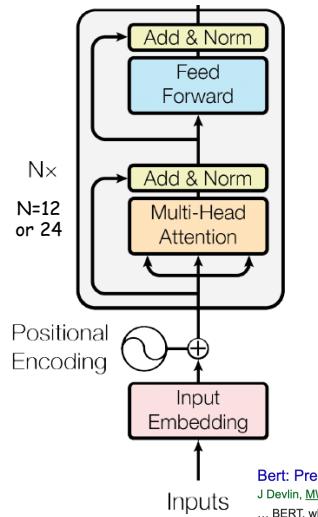
- All large language models are decoders.

8

8

BERT: Bidirectional Encoder Representations from Transformers

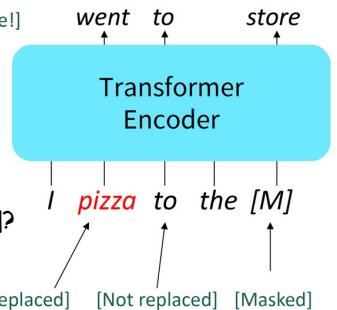
References: Stanford CS224N, Lecture 9: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture09-pretraining.pdf>
 Princeton COS 597G, Lecture 2: <https://www.cs.princeton.edu/courses/archive/fall22/cos597G/lectures/lec02.pdf>



- 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**. [Predict these!]
 - 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
 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 ...

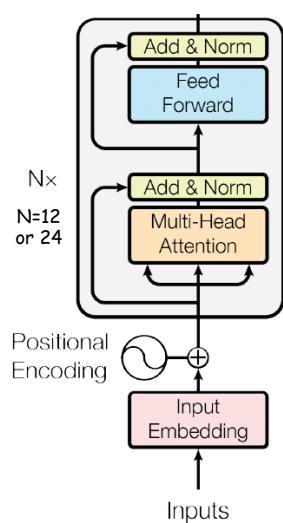
☆ Save 59 Cite Cited by 93230 Related articles All 46 versions ☺



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	$E_{[CLS]}$	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_A	E_A	E_A	E_A	E_B	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

Positional Encoding

Inputs

Input Embedding

Multi-Head Attention

Add & Norm

Feed Forward

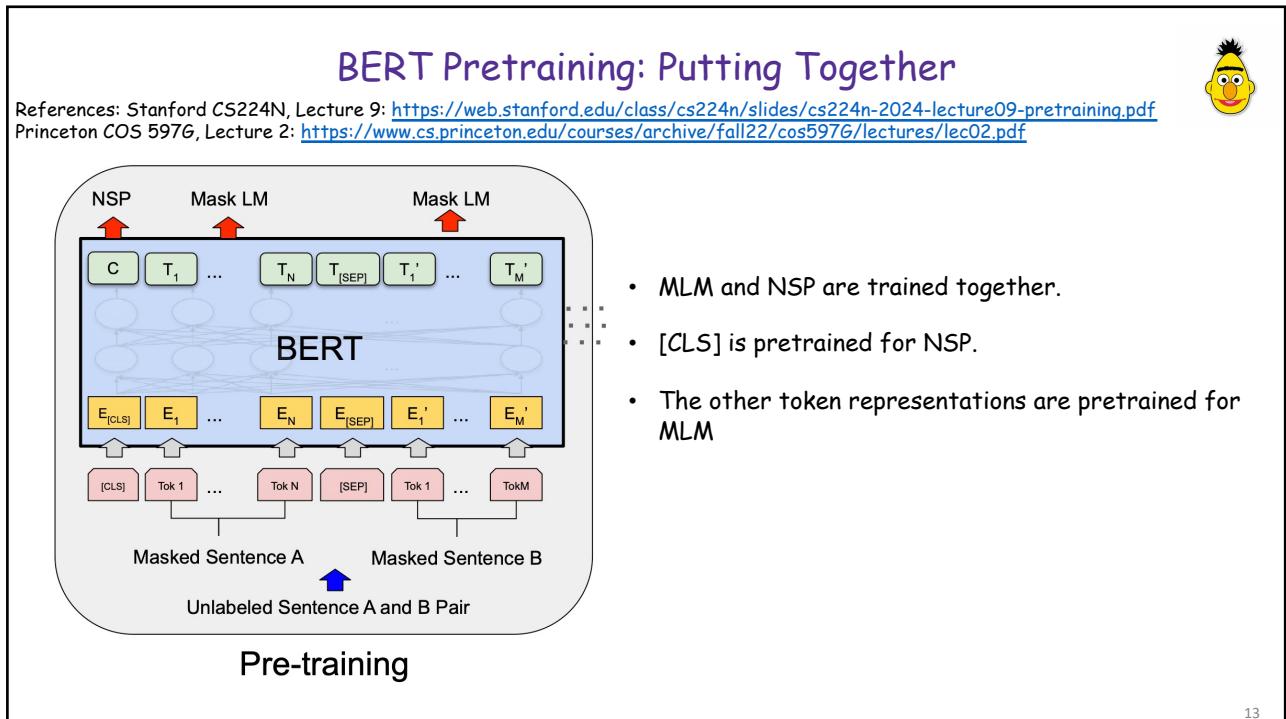
Add & Norm

Diagram of the BERT architecture showing the stack of N layers (12 or 24). Each layer takes inputs, adds position embeddings, and processes through a stack of Multi-Head Attention and Feed Forward blocks, followed by Add & Norm layers.

- 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

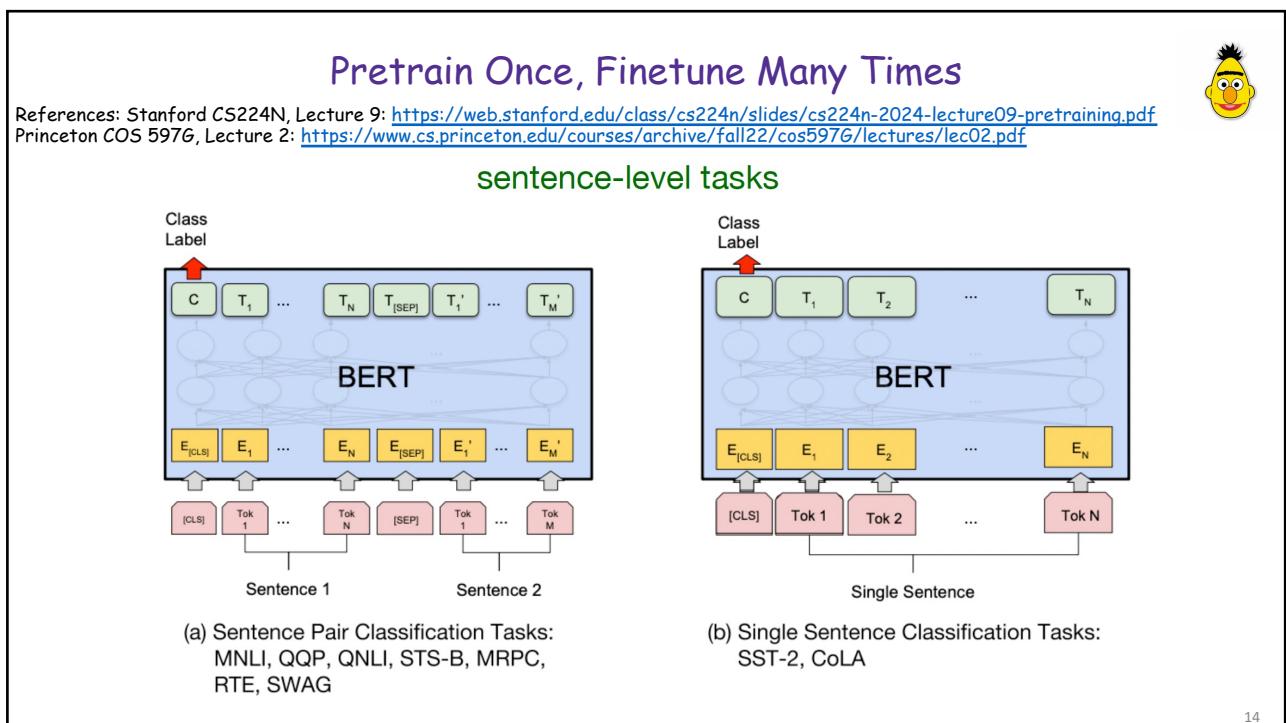
12

12



13

13



14

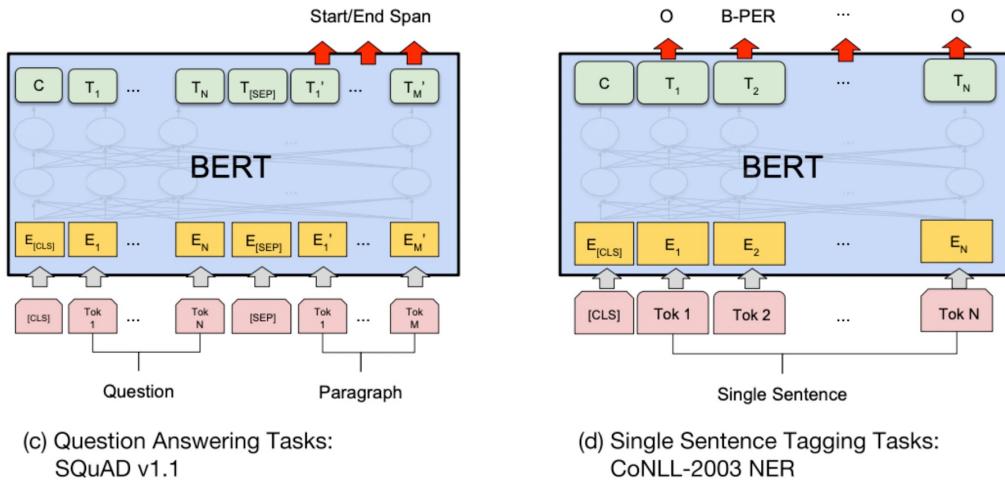
14

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>

token-level tasks



15

15

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>

Sentence-Level Task

- Sentence pair classification tasks:

MNU

Premise: A soccer game with multiple males playing.

Hypothesis: Some men are playing a sport

{entailment, contradiction, neutral}

QQP

Q1: Where can I learn to invest in stocks?

{duplicate, not duplicate}

• Single sentence classification tasks:

SST2

rich veins of funny stuff in this movie.

{positive, negative}

16

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>



Token-Level Task

- Extractive question answering e.g., SQuAD (Rajpurkar et al., 2016)

SQuAD

Question: The New York Giants and the New York Jets play at which stadium in NYC ?

Context: The city is represented in the National Football League by the New York Giants and the New York Jets , although both teams play their home games at MetLife Stadium in nearby East Rutherford , New Jersey , which hosted Super Bowl XLVIII in 2014 .

(Training example 29,883)

MetLife Stadium

- Named entity recognition (Tjong Kim Sang and De Meulder, 2003)

CoNLL 2003 NER

John Smith lives in New York

B-PER I-PER O O B-LOC I-LOC

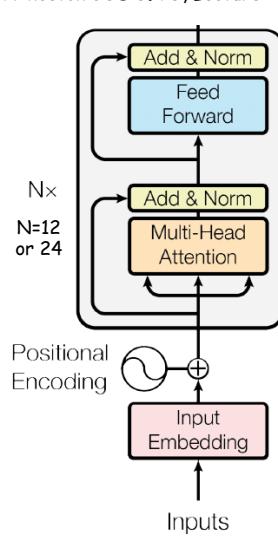
17

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

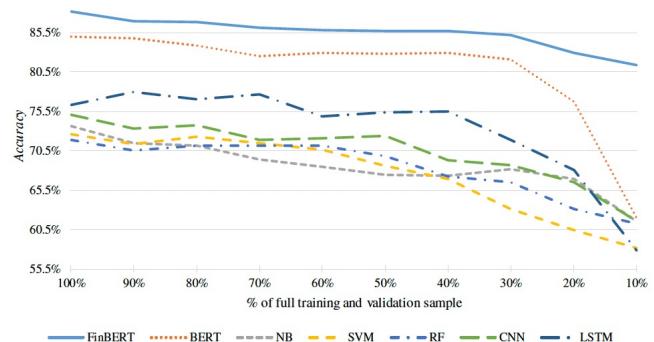
18

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
- The bitter lesson:** Once we have GPT-4 or Claud-3, what is the value of FinBERT?

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 99 Cite Cited by 144 Related articles Web of Science: 22 ☰

19

19

Agenda

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

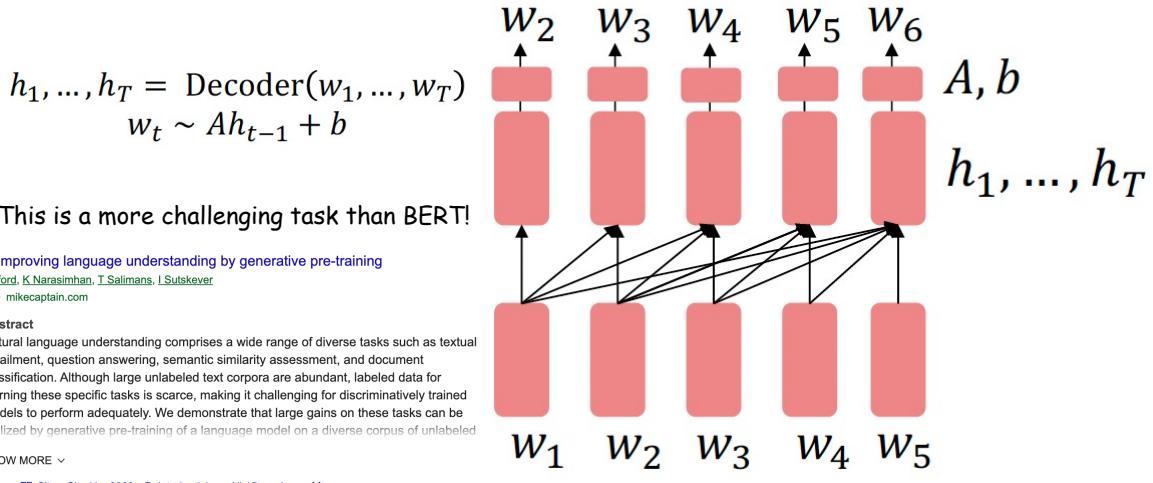
20

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.



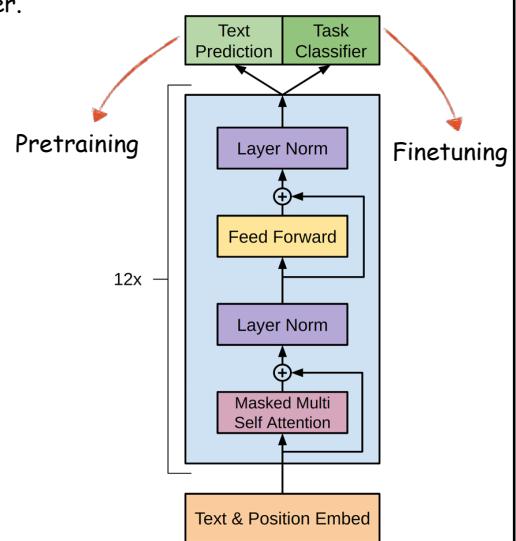
21

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.



[PDF] Improving language understanding by generative pre-training

A Radford, K Narasimhan, T Salimans, I Sutskever
2018 - mikedcaption.com

Abstract

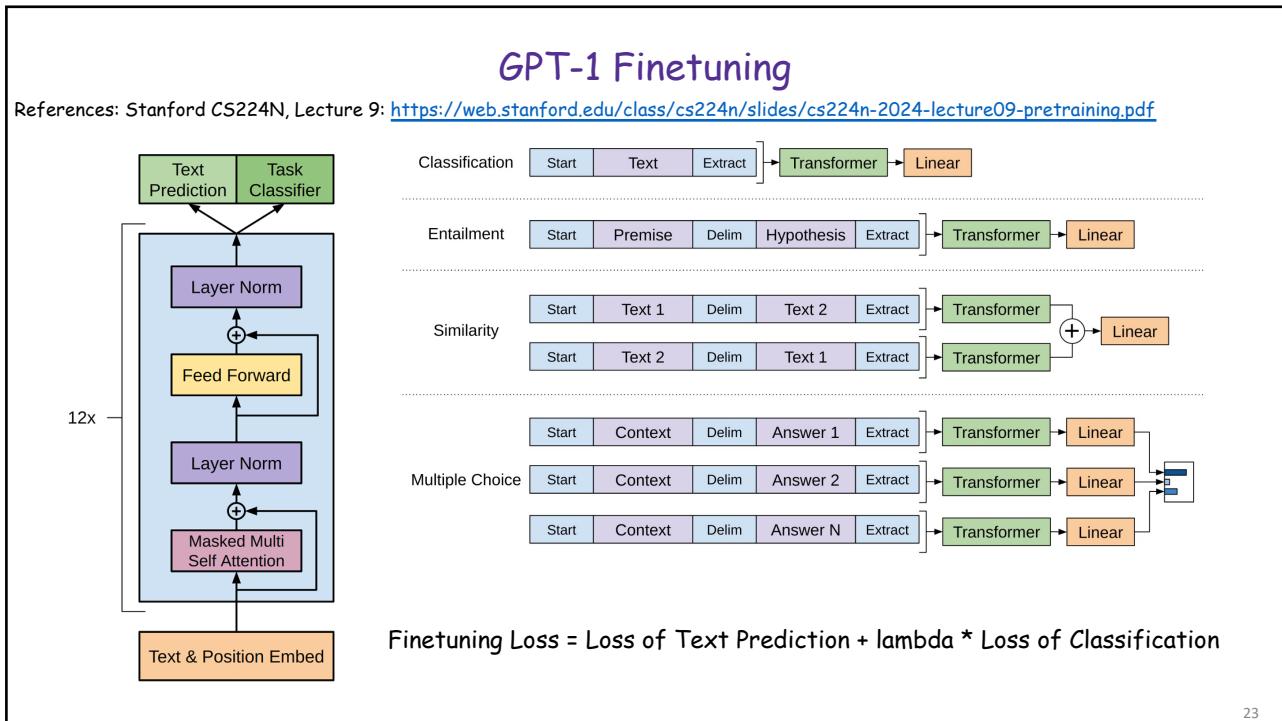
Natural language understanding comprises a wide range of diverse tasks such as textual entailment, question answering, semantic similarity assessment, and document classification. Although large unlabeled text corpora are abundant, labeled data for learning these specific tasks is scarce, making it challenging for discriminatively trained models to perform adequately. We demonstrate that large gains on these tasks can be realized by generative pre-training of a language model on a diverse corpus of unlabeled

SHOW MORE ▾

☆ Save 59 Cite Cited by 8363 Related articles All 15 versions ☰

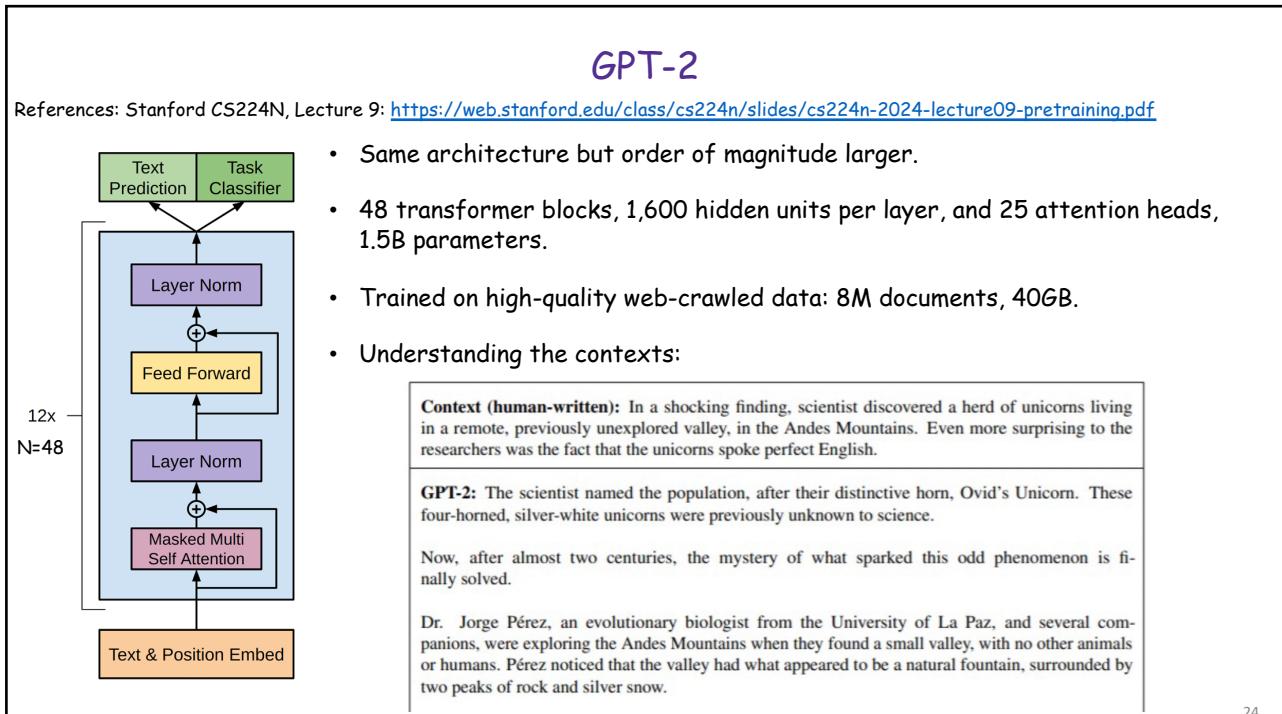
22

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>

