

DSME 6635: Artificial Intelligence for Business Research

Deep-Learning-based NLP: Pretraining

Renyu (Philip) Zhang

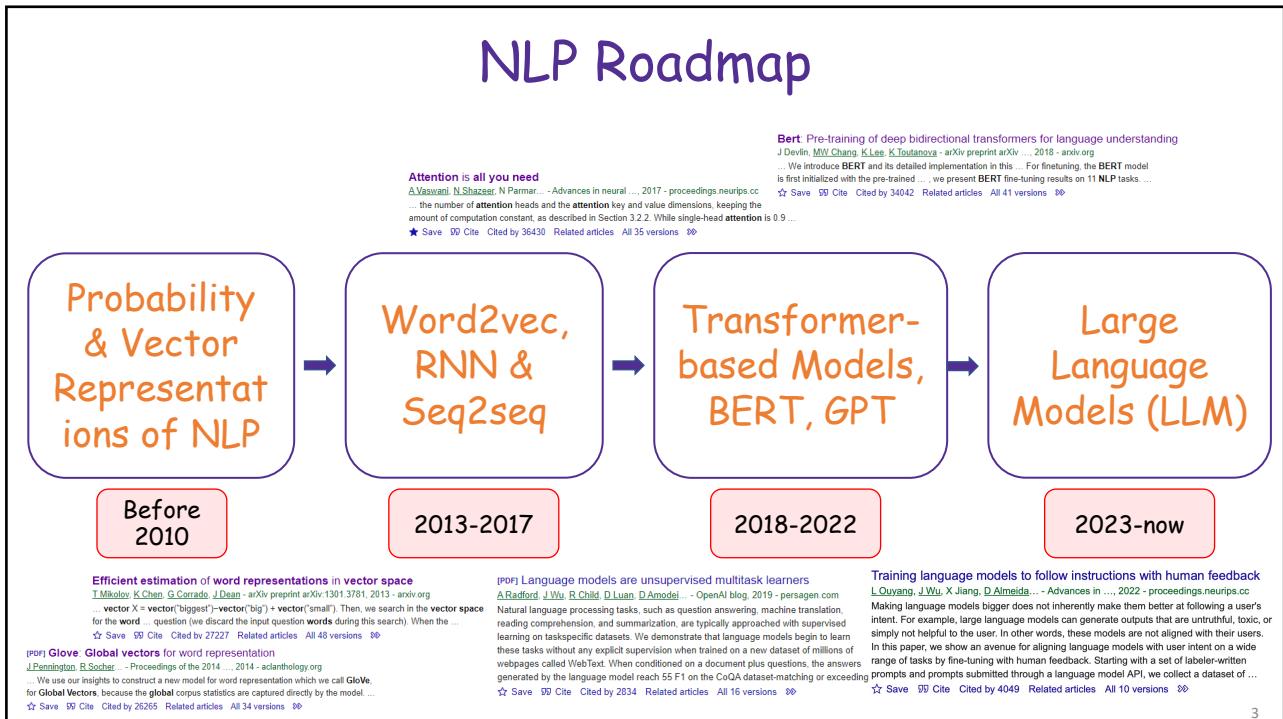
1

Agenda

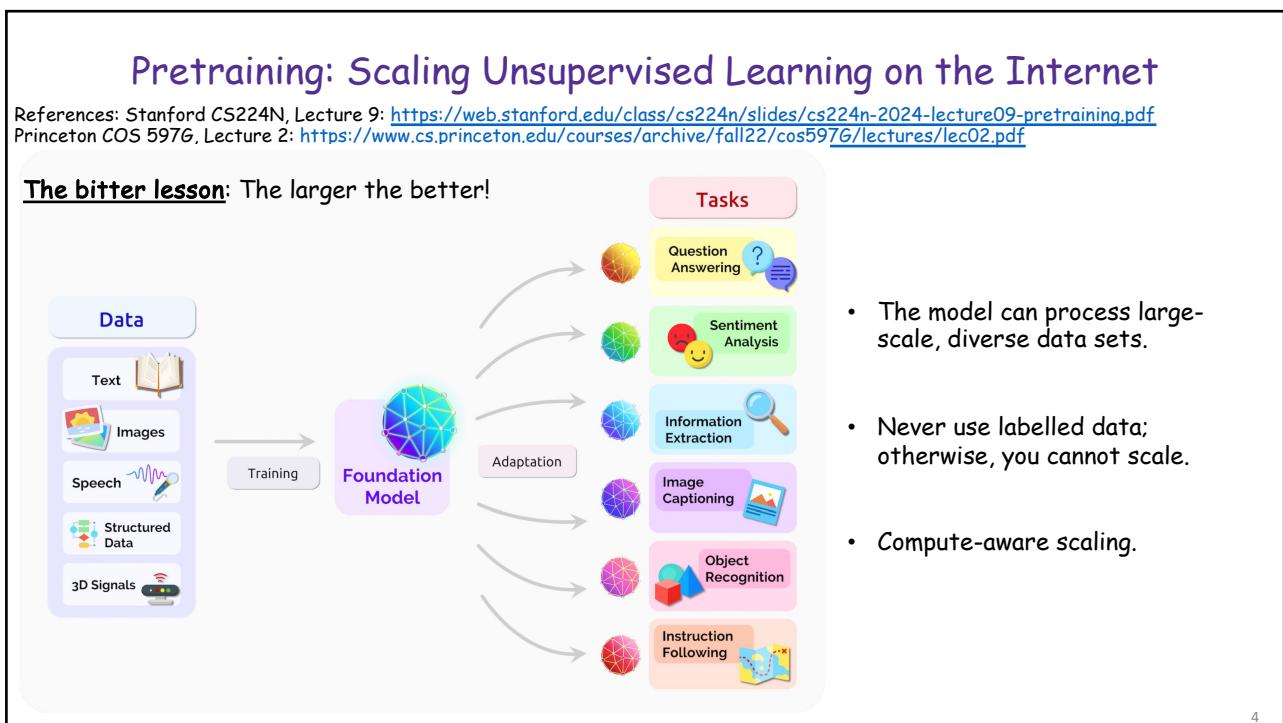
- BERT: Bidirectional Encoder Representations from Transformers
- GPT: Generative Pretrained Transformers

2

2



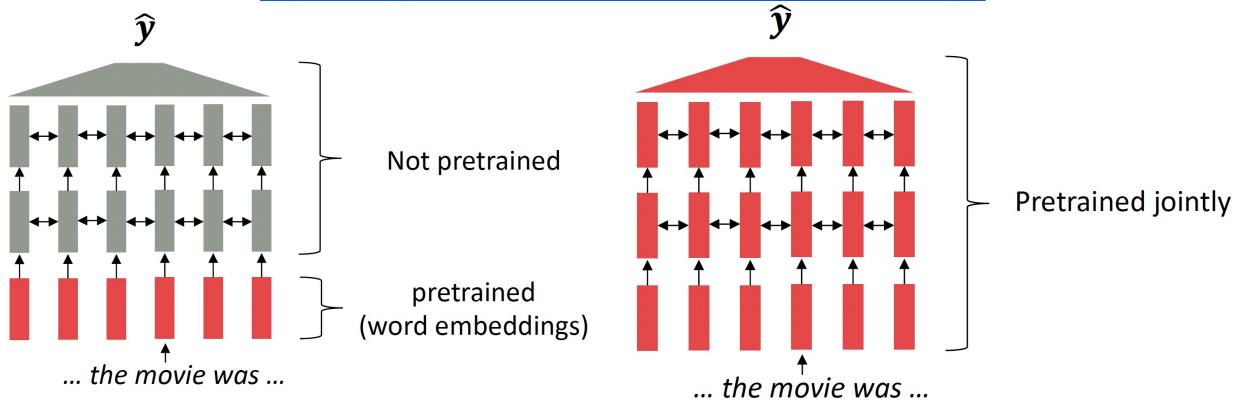
3



4

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>



[Recall, *movie* gets the same word embedding,
no matter what sentence it shows up in]

[This model has learned how to represent
entire sentences through pretraining]

5

5

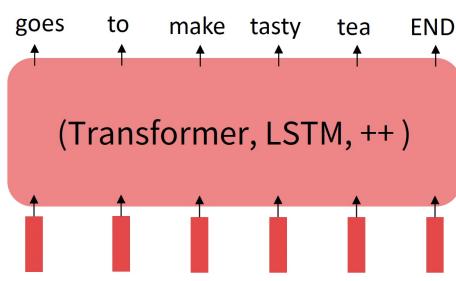
From Pretrained Word Embeddings to Pretrained Models

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)

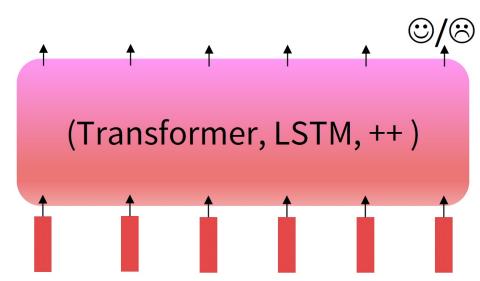
Lots of text; learn general things!



$\hat{\theta}$ by approximating $\min_{\theta} \mathcal{L}_{\text{pretrain}}(\theta)$

Step 2: Finetune (on your task)

Not many labels; adapt to the task!



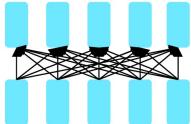
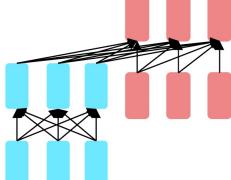
approximates $\min_{\theta} \mathcal{L}_{\text{finetune}}(\theta)$, starting at $\hat{\theta}$

6

6

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	<ul style="list-style-type: none"> • Can condition on future. • Example: BERT.
 Encoder-Decoders	<ul style="list-style-type: none"> • Combining encoder and decoder. • Example: T5
 Decoders	<ul style="list-style-type: none"> • Cannot condition on future. • Example: GPT • All large language models are decoders.

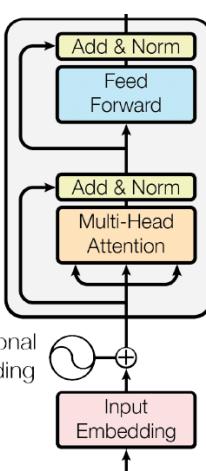
7

7

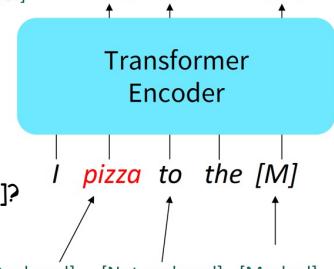
BERT: Bidirectional Encoder Representations from Transformers

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.

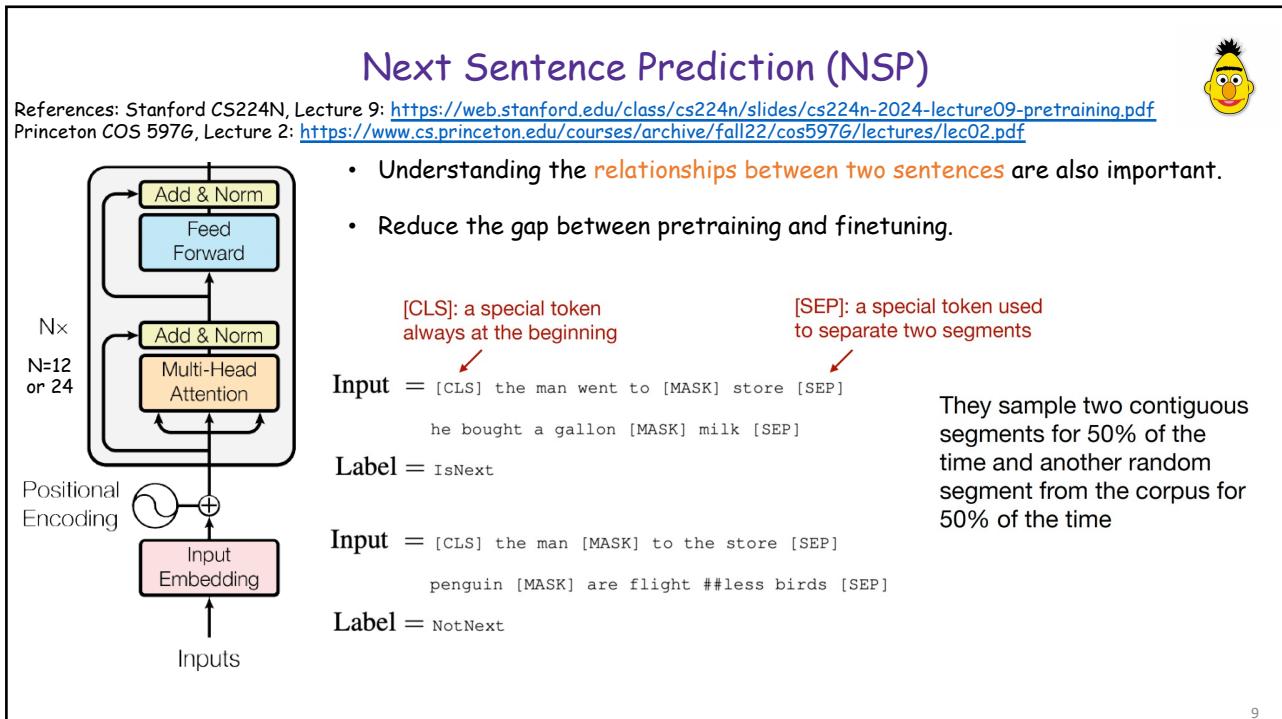


The diagram shows the input sequence "I pizza to the store" entering the Transformer Encoder. The word "pizza" is highlighted in red and labeled "[Replaced]". The words "to" and "the" are labeled "[Not replaced]". The word "store" is labeled "[Masked]" and has an arrow pointing to it from the text "Transformer Encoder".

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 **bidi**rectional representations from unlabeled text by ...
 ☆ Save ⌂ Cite Cited by 93230 Related articles All 46 versions ⌂

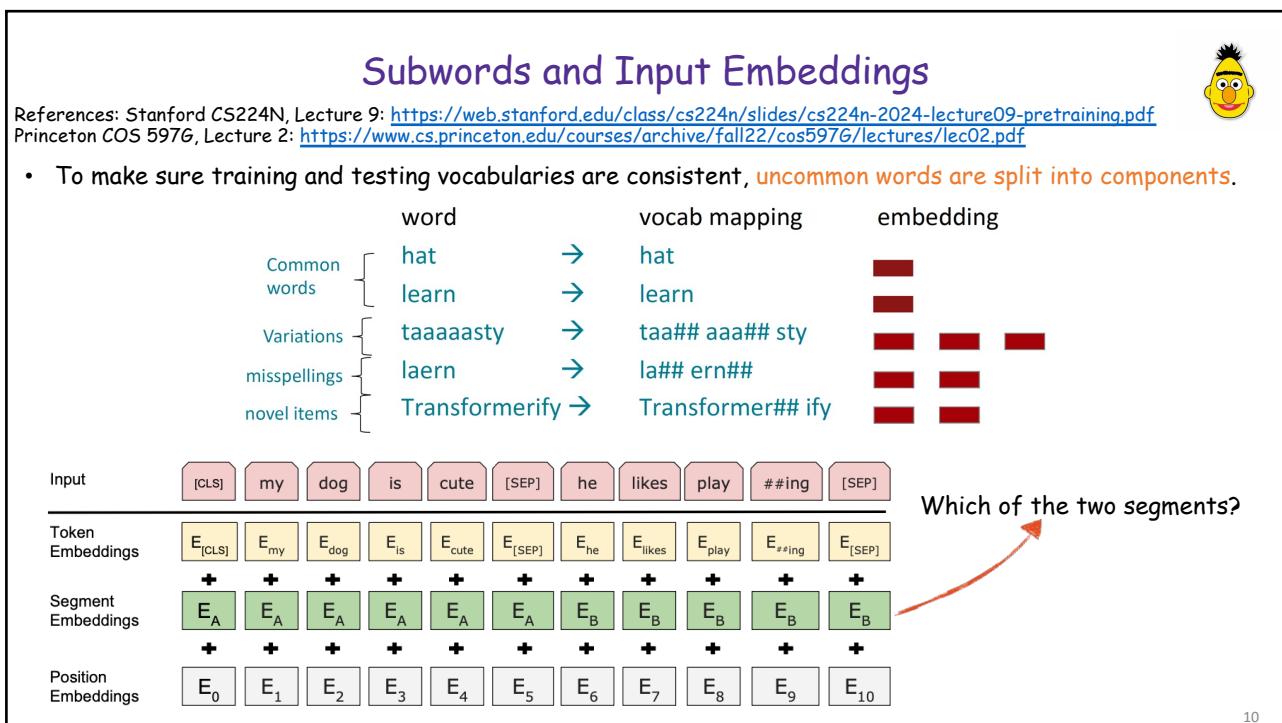
8

8



9

9



10

10

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>

The diagram illustrates the BERT pretraining architecture. It starts with 'Inputs' which are converted into 'Input Embedding'. This is combined with 'Positional Encoding' via a residual connection (sum) to produce the initial representation. This representation then passes through $N \times$ (either 12 or 24) layers. Each layer consists of a 'Multi-Head Attention' block followed by an 'Add & Norm' block, and a 'Feed Forward' block followed by another 'Add & Norm' block. The final output is the pre-trained BERT model.

- 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

11

11

BERT Pretraining: Putting Together

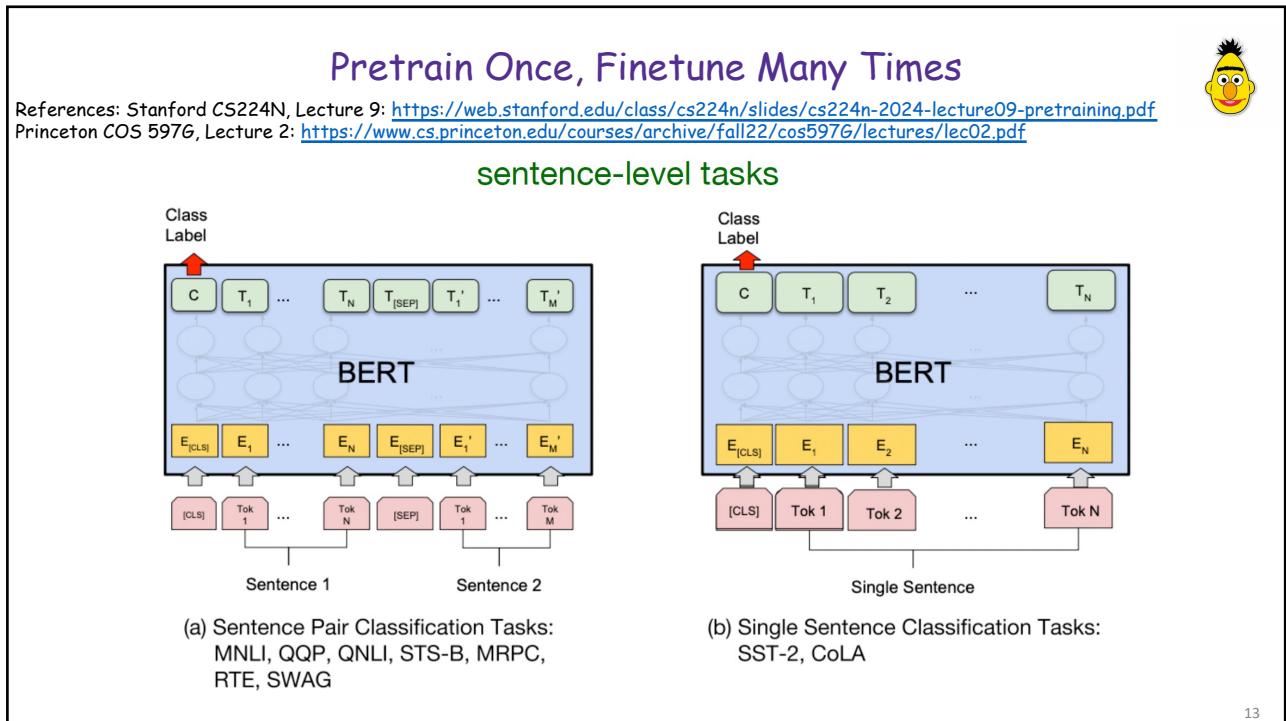
References: Stanford CS224N, Lecture 9: <https://web.stanford.edu/class/cs224n/slides/cs224n-2024-lecture09-pretraining.pdf>
 Princeton COS 597G, Lecture 2: <https://www.cs.princeton.edu/courses/archive/fall22/cos597G/lectures/lec02.pdf>

The diagram shows the BERT pre-training process. An 'Unlabeled Sentence A and B Pair' is split into 'Masked Sentence A' and 'Masked Sentence B'. These sentences are fed into the BERT model. The BERT model has two main output paths: one for the 'NSP' (Next Sentence Prediction) task and one for the 'Mask LM' (Masked Language Model) task. The NSP path involves tokens [CLS], Tok 1, ..., Tok N, [SEP], Tok 1, ..., Tok M. The Mask LM path involves tokens C, T₁, ..., T_N, T_[SEP], T_{1'}, ..., T_{M'}. The BERT model also generates token embeddings E_[CLS], E₁, ..., E_N, E_[SEP], E_{1'}, ..., E_{M'}.

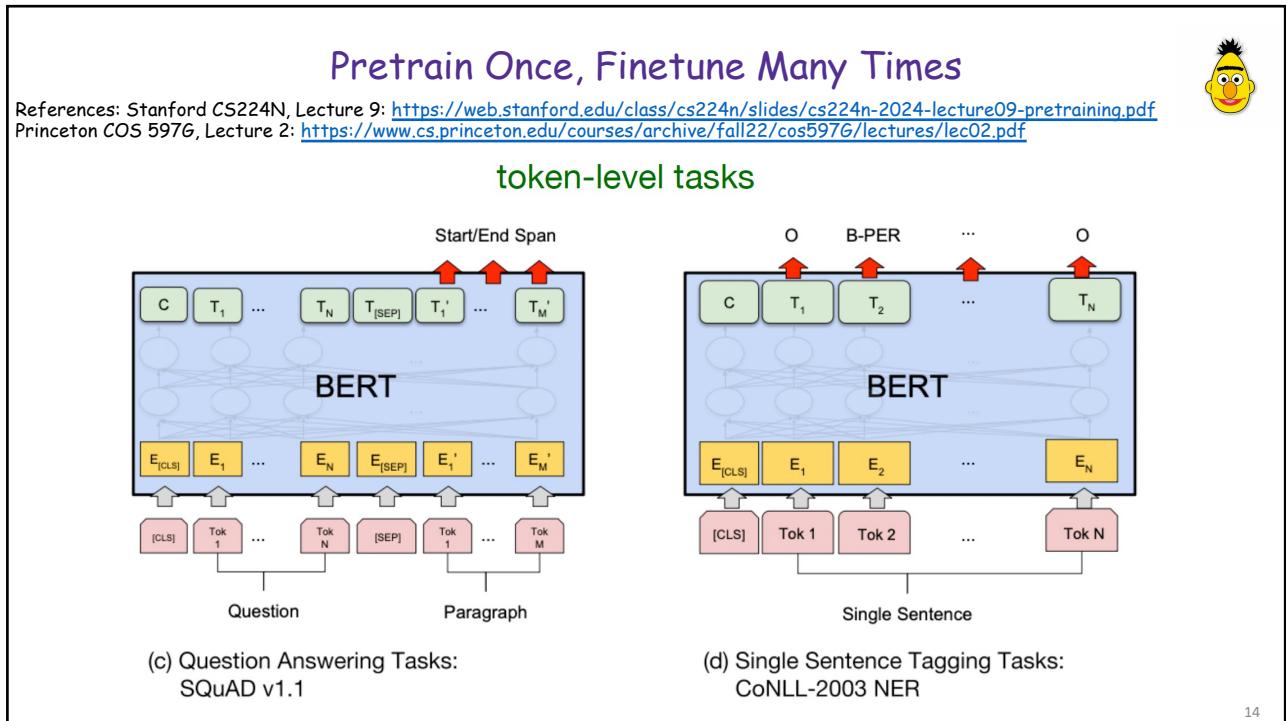
- MLM and NSP are trained together.
- [CLS] is pretrained for NSP.
- The other token representations are pretrained for MLM

12

12



13



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>

Sentence-Level Task

- Sentence pair classification tasks:

MNLI	Premise: A soccer game with multiple males playing. Hypothesis: Some men are playing a sport.	{entailment, contradiction, neutral}
------	--	--------------------------------------

Q1: Where can I learn to invest in stocks? {[duplicate](#), not duplicate}

Q2: How can I learn more about stocks?

- Single sentence classification tasks:

SST2 rich veins of funny stuff in this movie {positive, negative}

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>

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 29,88)

MetLife Stadium

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

CoNLL 2003 NER

John Smith lives in New York

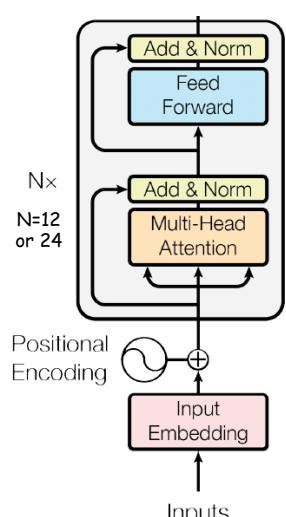
B-PER I-PER O O B-LOC I-LOC

16

16

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.

17

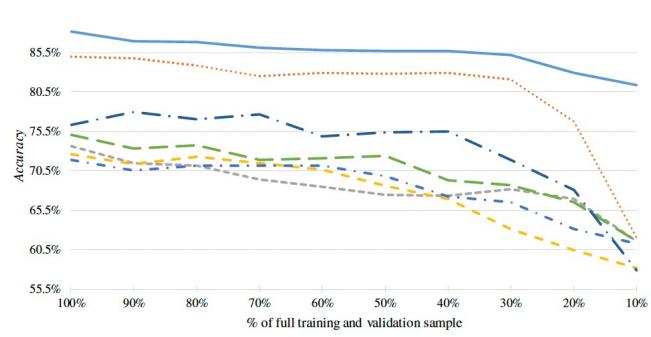
17

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 ☰

18

18

Agenda

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

19

19

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.

$$\begin{aligned} h_1, \dots, h_T &= \text{Decoder}(w_1, \dots, w_T) \\ w_t &\sim Ah_{t-1} + b \end{aligned}$$

This is a more challenging task than BERT!

[PPL: Improving language understanding by generative pre-training](#)

A Radford, K Narasimhan, T Salimans, I Sutskever

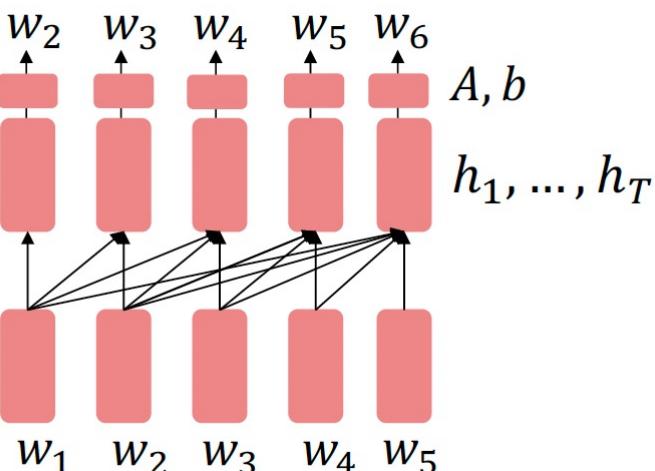
2018 - [mikecaptain.com](#)

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](#) 99 [Cite](#) Cited by 8363 [Related articles](#) All 15 versions [⊗⊗](#)



20

20

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 - mikecaptain.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](#) [XX](#)

The diagram illustrates the GPT-1 architecture. It starts with 'Text & Position Embed' at the bottom, which feeds into a stack of 12 identical layers. Each layer contains 'Masked Multi Self Attention' (red), 'Layer Norm' (purple), 'Feed Forward' (orange), and residual connections (blue). The final output from the stack is split into two paths: 'Pretraining' (top) and 'Finetuning' (bottom). The 'Pretraining' path leads to 'Text Prediction' and 'Task Classifier'. The 'Finetuning' path involves adding a 'Linear' layer to the output of the stack, followed by a 'Transformer' layer, and then another 'Linear' layer to produce the final output for various NLP tasks.

21

21

GPT-1 Finetuning

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

The diagram shows the GPT-1 architecture being finetuned for four specific tasks: Classification, Entailment, Similarity, and Multiple Choice. The architecture remains largely the same, with 'Text & Position Embed' at the bottom, followed by 12 layers of 'Masked Multi Self Attention', 'Layer Norm', and 'Feed Forward'. The 'Finetuning Loss' is calculated as the sum of the 'Loss of Text Prediction' and a weighted ('lambda') 'Loss of Classification' for each task. The 'Classification' task takes inputs like 'Start', 'Text', and 'Extract' and uses a single 'Transformer' and 'Linear' layer. The 'Entailment' task takes inputs like 'Start', 'Premise', 'Delim', 'Hypothesis', and 'Extract' and uses a similar setup. The 'Similarity' task takes pairs of inputs like 'Start', 'Text 1', 'Delim', 'Text 2', 'Extract' and 'Start', 'Text 2', 'Delim', 'Text 1', 'Extract', and adds residual connections between the two 'Transformer' layers. The 'Multiple Choice' task takes inputs like 'Start', 'Context', 'Delim', 'Answer 1', 'Extract' and 'Start', 'Context', 'Delim', 'Answer 2', 'Extract', and so on, and uses multiple 'Transformer' and 'Linear' layers to handle multiple answers.

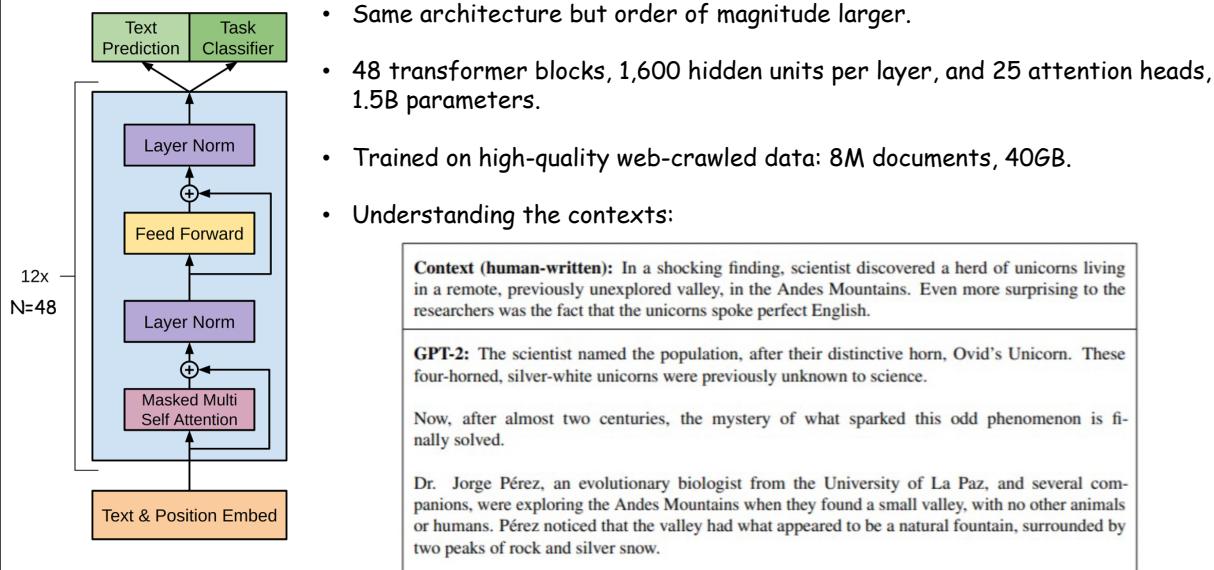
$\text{Finetuning Loss} = \text{Loss of Text Prediction} + \lambda * \text{Loss of Classification}$

22

22

GPT-2

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

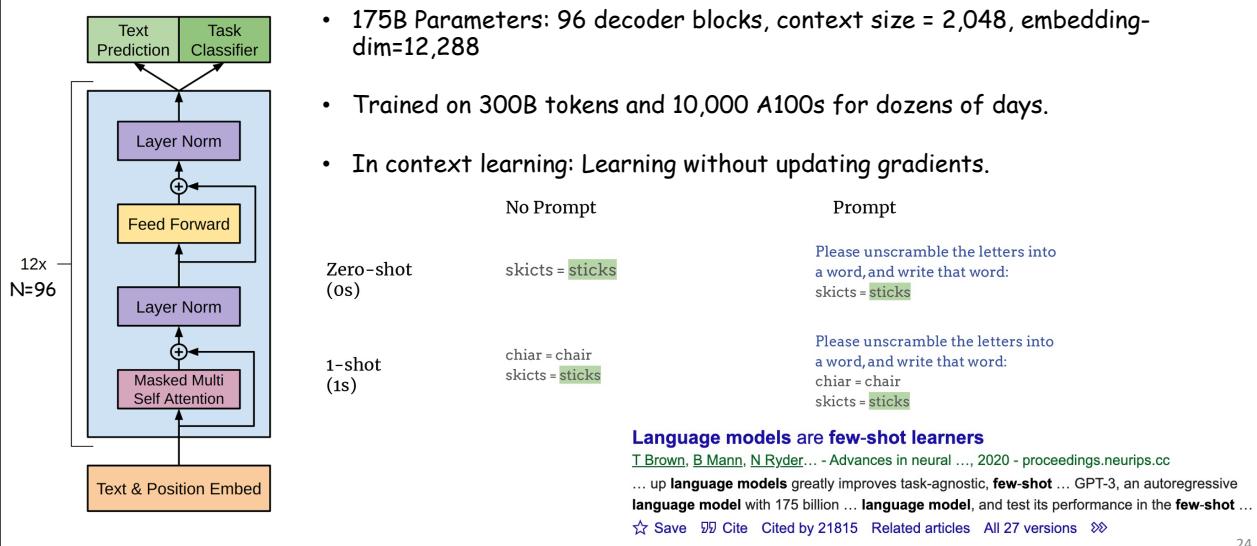


23

23

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/fall22/cos597G/lectures/lec04.pdf>



24

24