

# BERT and GPT-1

Graeme Bates, Mingyang Cen  
10/27/25

# Review of Transformers

- Embeddings
  - Embedding Tokens into High Dimensional Euclidean Space
- Attention
  - Context plays a major role in Natural Language. Attention allows embedding to take into account the context of a given token.

# Tokenization

BPE (byte pair encoding).

Start with taking each letter as a token, then merge together the tokens that appear the most together into one. Repeat until the desired vocabulary size is reached.

Ex: then that fox ate the rabbit. → (th)en (th)at fox ate (th)e rabbit.

→ (th)en (th)(at) fox (at)e (th)e rabbit. → (th)en( th)(at) fox (at)e( th)e rabbit.

Tokens: ( th) (th) (at) [space] a b e f h i n o r t x

# GPT-1(Radford et al) Motivation

- Labeled Data is hard to come by, but there is an abundance of unlabeled text.
- Natural Language has inherent structures that are common to all tasks.
- Pre-training on unlabeled text allows for a “regularization scheme” for more fine tuned tasks downstream.
- Task Specific Fine-Tuning of a Pre-Trained model requires less labeled data then traditional NLP schemes.
- Text Embeddings with attention allow for long-range linguistic structure.

# GPT-1 in action

- Broadly speaking, GPT-1 is pre-training with transformers to train next token prediction  $\min_{\theta} \mathbb{E}_{x \sim D} [-\log P_{\theta}(x_{n+1} | x_1, x_2, \dots, x_n)]$  here theta represents the learned parameters of the NN
- Minimizing the loss between the actual next token  $y_{n+1}$  and predicted token  $P_{\theta}(x_{n+1} | x_1, \dots, x_n)$

# GPT-1 Pre-training Data Sets

- Books Corpus Dataset containing 7,000 unpublished books
  - Long structure of books helps with long-range information training.
- Banord Benchmark containing shuffled sentences
  - This lowered token level perplexity

# GPT-1 Pre-Training Architecture

- 12-Layer Decoder only Transformer with masked self attention heads
  - left to right nature of GPT-1, each token only vectorized with attention to previous tokens
- 768 Dimensional states (ie token embedding)
- Position-wise feedforward networks: 3072 dimensional inner states.
- Extensive use of Layer Normalization with weight initialization of  $N(0,0.2)$

# GPT-1 Pre-training Details

- Adam Optimizer
  - Momentum and Adaptive Learning Rate
- Max Learning Rate of  $2.5e-4$
- Cosine Scheduler for warm-up and cool down of Learning Rate.
  - Smoothly increase/decrease Learning rate
- Objective:  $\min_{\theta} -\log P_{\theta}(x_{n+1}|x_{1:n})$



# GPT-1 Supervised Fine Tuning

Pre-trained model ready for next token prediction ie. text generation

But this is of limited use, so the goal was to fine tune this pre-trained model to perform various more specialized tasks.

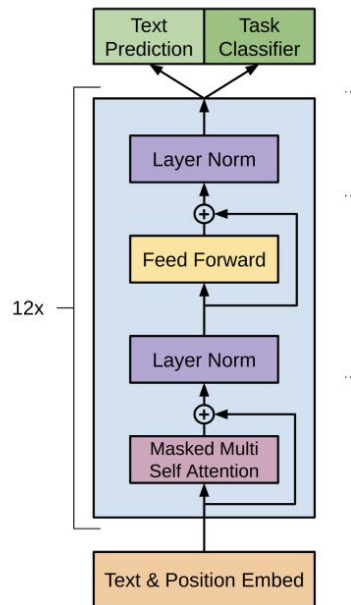
Final Linear Layer is added for fine-tuning and smaller learning rate is used to prevent canceling out pre-training.

## Types of Tasks for Fine Tuning

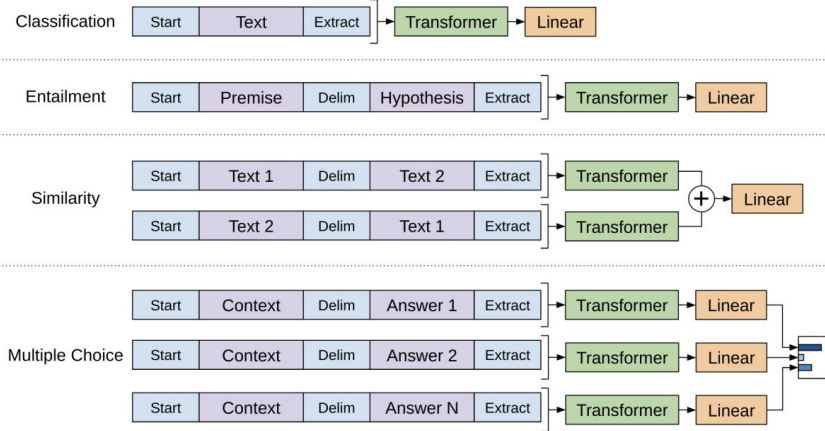
- Natural Language Inference
  - Recognizing the textual entailment
  - ie. judging relationship between two sentences; contradictory, supportive no relation etc...
- Question-Answering
  - Selecting best ending to multi-sentence stories
  - Completing middle/high school level exams.
- Classification
- Semantic Similarity

# GPT-1 Structure

The transformer structure is as follows:



These are the transformations applied to turn a labeled dataset into a sequence of natural language tokens for fine-tuning. Special blue tokens indicate delimiters.



# GPT-1 Was it Worth It?

- GPT-1 set state of the art benchmarks on 9 of the 12 studied tasks.

# GPT-1 Was it Worth It?

- Natural Language Inference

Method	MNLI-m	MNLI-mm	SNLI	SciTail	QNLI	RTE
ESIM + ELMo [44] (5x)	-	-	<u>89.3</u>	-	-	-
CAFE [58] (5x)	80.2	79.0	<u>89.3</u>	-	-	-
Stochastic Answer Network [35] (3x)	<u>80.6</u>	<u>80.1</u>	-	-	-	-
CAFE [58]	78.7	77.9	88.5	<u>83.3</u>		
GenSen [64]	71.4	71.3	-	-	<u>82.3</u>	59.2
Multi-task BiLSTM + Attn [64]	72.2	72.1	-	-	82.1	<b>61.7</b>
Finetuned Transformer LM (ours)	<b>82.1</b>	<b>81.4</b>	<b>89.9</b>	<b>88.3</b>	<b>88.1</b>	56.0

GPT-1 Outperforms on 4 of 5 baselines  
RTE dataset happens to be the smallest

# GPT-1 Was it Worth It?

- Question-Answering

Method	Story Cloze	RACE-m	RACE-h	RACE
val-LS-skip [55]	76.5	-	-	-
Hidden Coherence Model [7]	<u>77.6</u>	-	-	-
Dynamic Fusion Net [67] (9x)	-	55.6	49.4	51.2
BiAttention MRU [59] (9x)	-	<u>60.2</u>	<u>50.3</u>	<u>53.3</u>
Finetuned Transformer LM (ours)	<b>86.5</b>	<b>62.9</b>	<b>57.4</b>	<b>59.0</b>

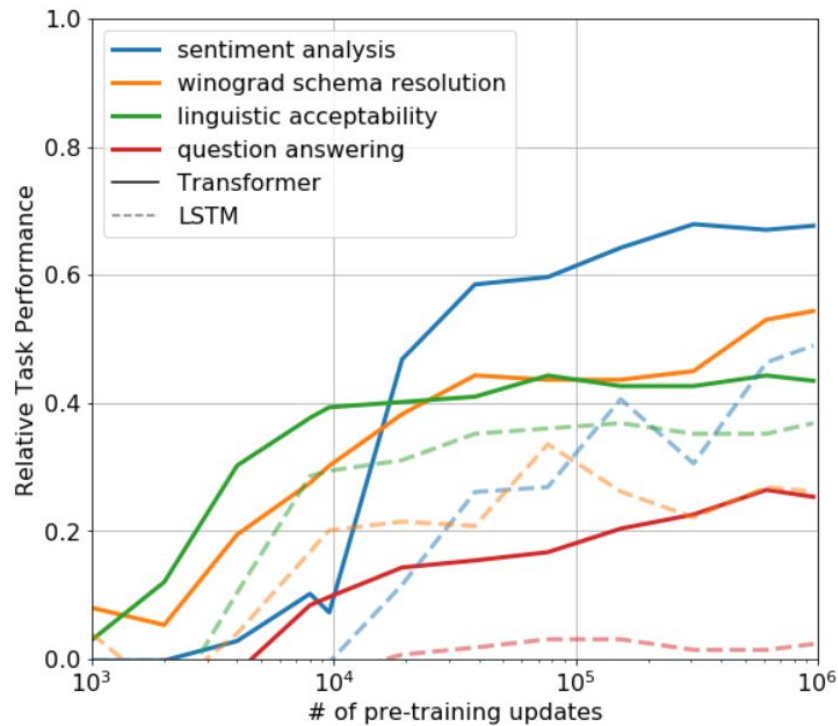
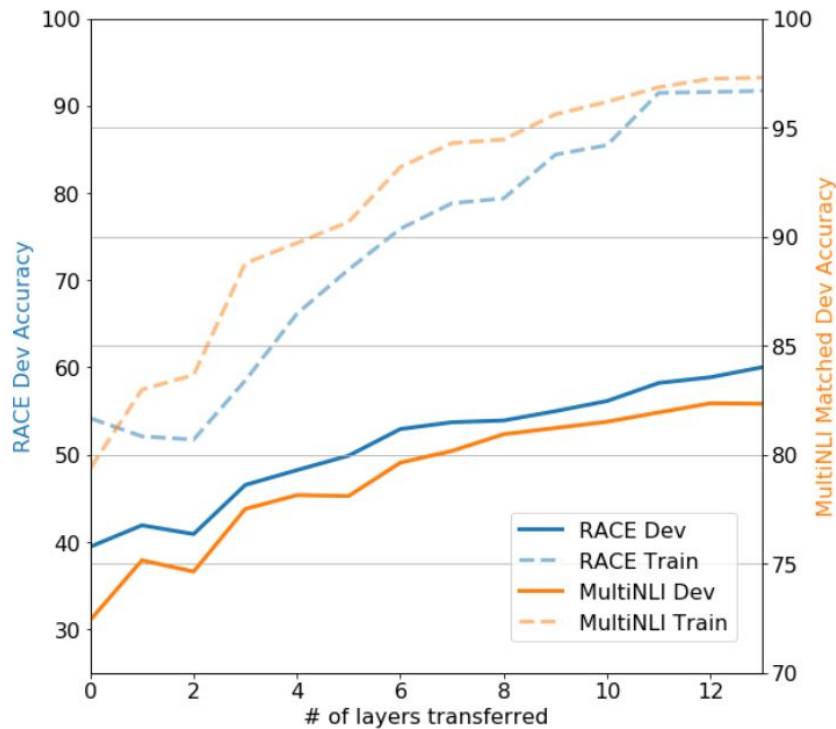
GPT-1 Outperforms on all datasets

# GPT-1 Was it Worth It?

- Semantic Similarity and Classification

Method	Classification		Semantic Similarity			GLUE
	CoLA (mc)	SST2 (acc)	MRPC (F1)	STSB (pc)	QQP (F1)	
Sparse byte mLSTM [16]	-	<b>93.2</b>	-	-	-	-
TF-KLD [23]	-	-	<b>86.0</b>	-	-	-
ECNU (mixed ensemble) [60]	-	-	-	<u>81.0</u>	-	-
Single-task BiLSTM + ELMo + Attn [64]	<u>35.0</u>	90.2	80.2	55.5	<u>66.1</u>	64.8
Multi-task BiLSTM + ELMo + Attn [64]	18.9	91.6	83.5	72.8	<u>63.3</u>	<u>68.9</u>
Finetuned Transformer LM (ours)	<b>45.4</b>	91.3	82.3	<b>82.0</b>	<b>70.3</b>	<b>72.8</b>


# GPT-1 Analysis



# BERT Motivation (Devlin et al.)

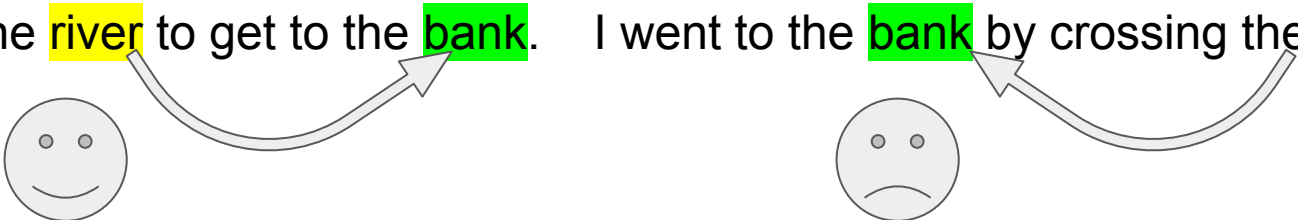
GPT-1 was trained on **constrained** self-attention, where every token can only have its attention computed with ones to its left. BERT fixes this and allows every token to attend to every other.

I crossed the **street** to get to the **bank**. I went to the **bank** by crossing the **street**.



The diagram illustrates the attention mechanism for the two sentences. In the first sentence, a curved arrow points from the word 'street' (highlighted in yellow) to the word 'bank' (highlighted in green). In the second sentence, a curved arrow points from the word 'bank' (highlighted in green) to the word 'street' (highlighted in yellow).

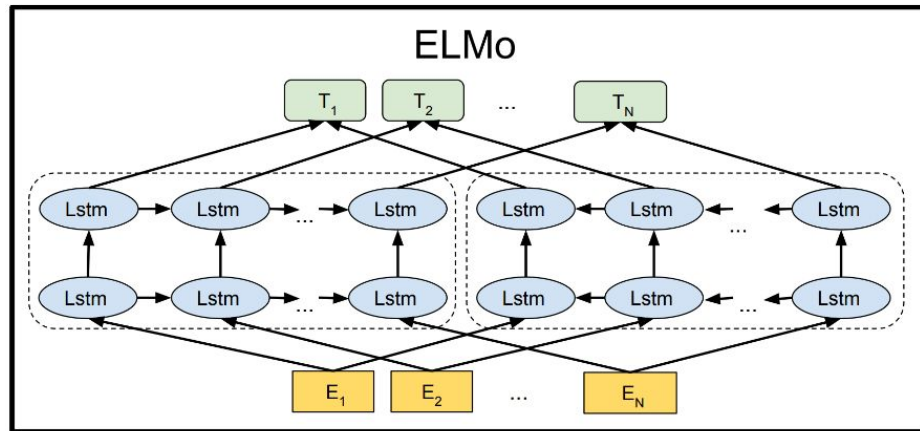
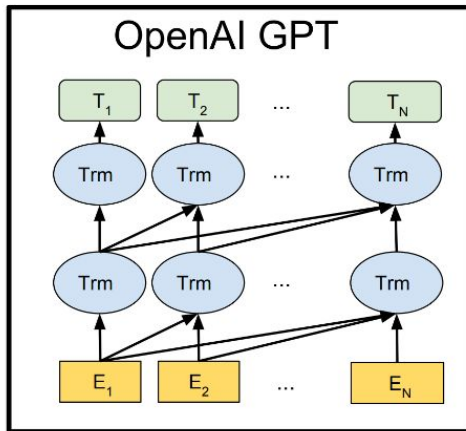
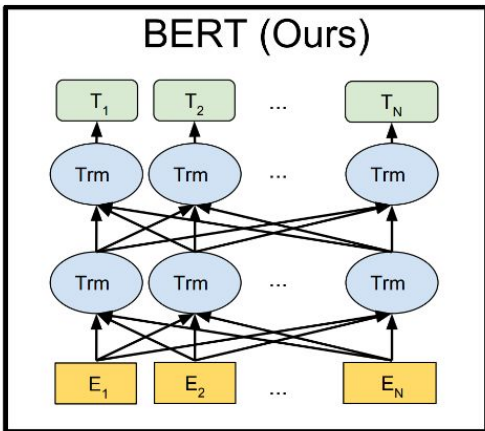
I crossed the **river** to get to the **bank**. I went to the **bank** by crossing the **river**.



The diagram illustrates the attention mechanism for the two sentences. In the first sentence, a curved arrow points from the word 'river' (highlighted in yellow) to the word 'bank' (highlighted in green). In the second sentence, a curved arrow points from the word 'bank' (highlighted in green) to the word 'river' (highlighted in yellow). A sad face icon is positioned below the first sentence, and another sad face icon is positioned below the second sentence.



# BERT Structure



ELMo/BiLSTM was a previous “feature-based” model: essentially a left-to-right RNN-style (long short-term memory) model combined with a right-to-left model. As you can see, BERT is bidirectional in every layer, making it superior.

# BERT Details

## **B**idirectional **E**ncoder **R**epresentations from **T**ransformers

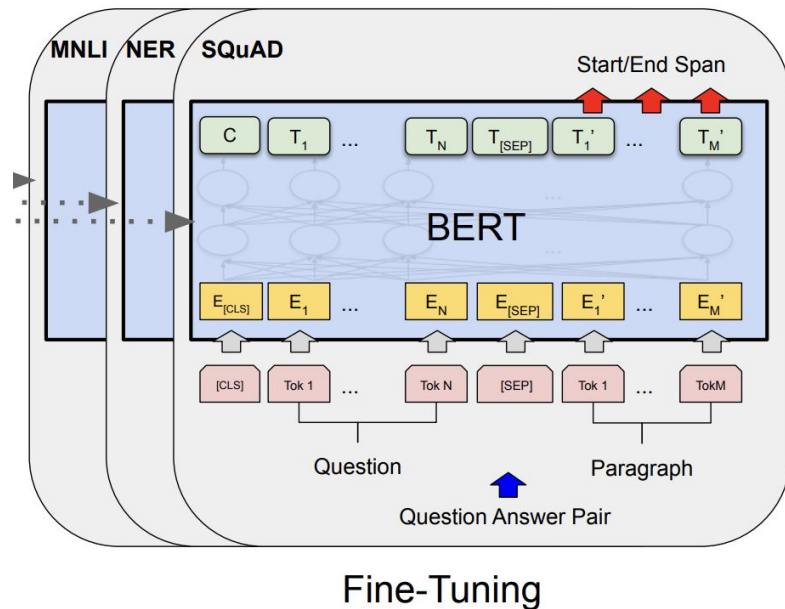
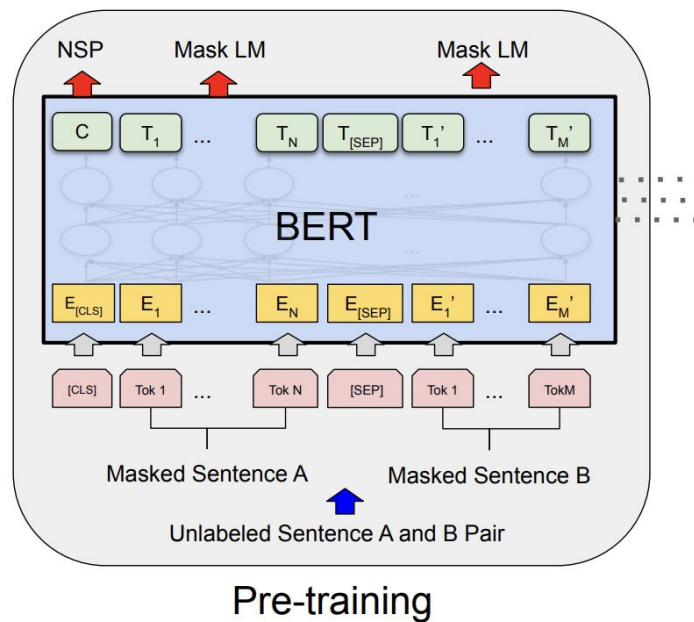
Effectively an encoder-only transformer.

Instead of next-token prediction, the model is pre-trained on two tasks:

- Masked LM/Cloze: Give the model a sequence of tokens, but choose some randomly to hide or replace with random. The model is challenged to find the correct words.
- A sequence of two sentences (one after the other) is chosen from the training data. The second is replaced with a random unrelated sentence with a 50% probability.

Same fine-tuning process as GPT.

# BERT Details



# BERT Pre-training: Cloze task

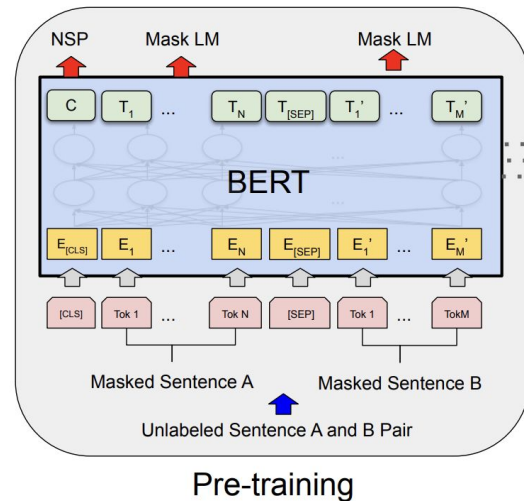
My dog is hairy.

→ My dog is **[MASK]**. (80%)

→ My dog is **hairy**. (10%)

→ My dog is **apple**. (10%)

The model's task is to output the probability distribution of the actual, non-masked tokens as the  $T_i$ .



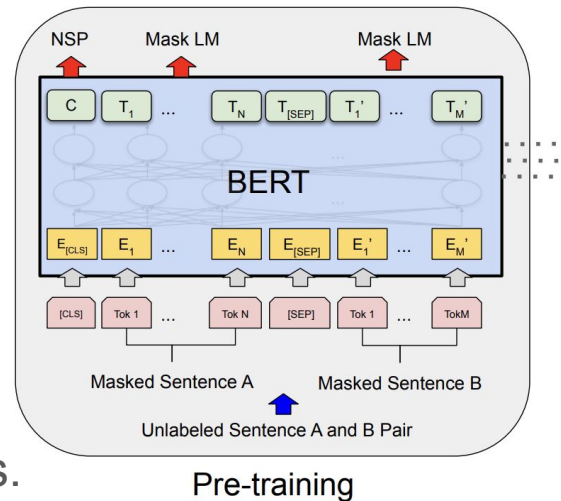
# BERT Pre-training: Next Sentence Prediction

The man went to the store. [SEP] He bought a gallon of milk.

→ **IsNext**

The man went to the store. [SEP] Penguins are flightless birds.

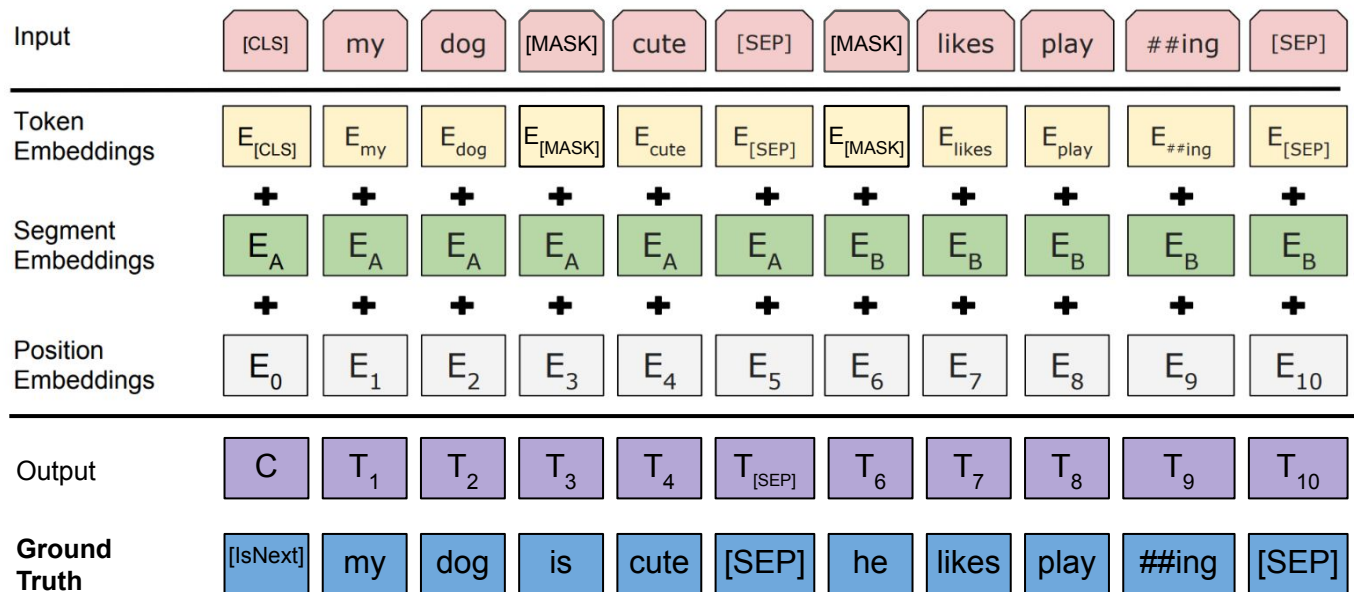
→ **NotNext**



# BERT Embedding and Output Format

The two tasks, Cloze and NSP, are combined.

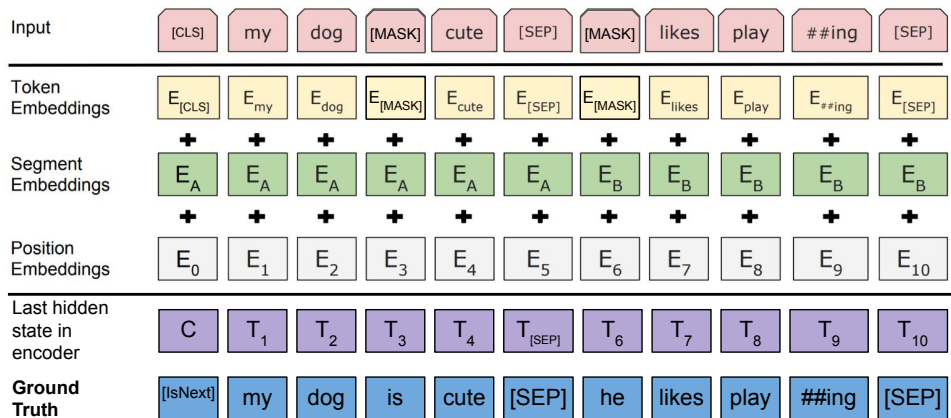
BERT uses a combination of three embeddings: a token embedding, a segment embedding, and a position embedding.



# BERT Embedding and Output Format

This embedding allows fine-tuning on several natural language processing tasks:

Is this text grammatical? Is this movie review positive? Do these two sentences mean the same thing? Which part of this article answers this question? Does it even answer the question? Does this statement imply the other?



# BERT Ablation Study Results

Every part of BERT is important!

The number of attention heads, embedding dimension, parameters, etc. cause smooth improvement.

Tasks	Dev Set				
	MNLI-m (Acc)	QNLI (Acc)	MRPC (Acc)	SST-2 (Acc)	SQuAD (F1)
BERT <sub>BASE</sub>	84.4	88.4	86.7	92.7	88.5
No NSP	83.9	84.9	86.5	92.6	87.9
LTR & No NSP	82.1	84.3	77.5	92.1	77.8
+ BiLSTM	82.1	84.1	75.7	91.6	84.9

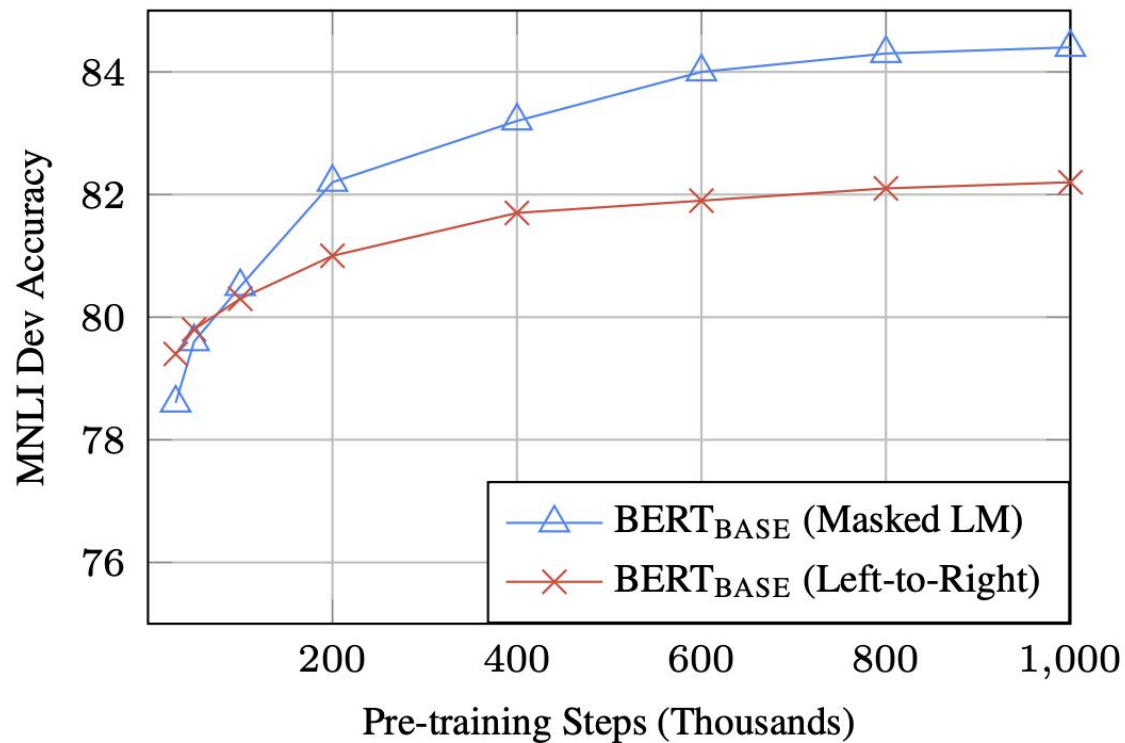
Cloze task only

GPT-1-style Cloze

ELMo +  
GPT-1-style Cloze



# BERT Ablation Study Results



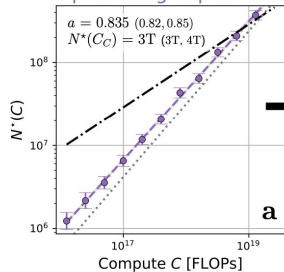
# Scaling Laws

Kaplan et al.

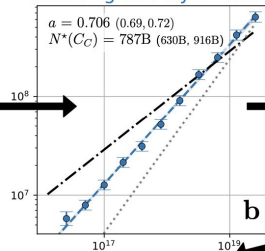
$N^{0.74} \sim D$  where  $N$  = params,  $D$  = dataset size.

$10^3$ – $10^9$  non-embedding parameters in tests

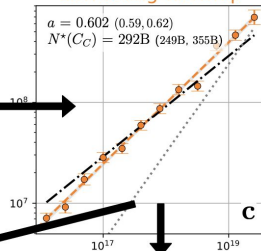
Reproducing Kaplan et al.



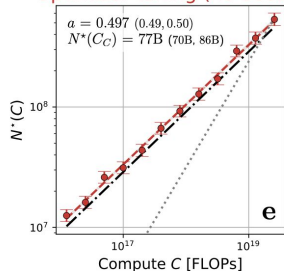
Counting last layer FLOPs



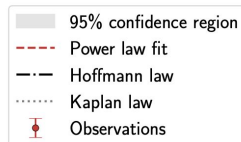
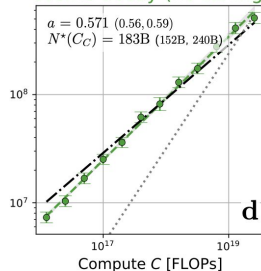
Correcting warmup



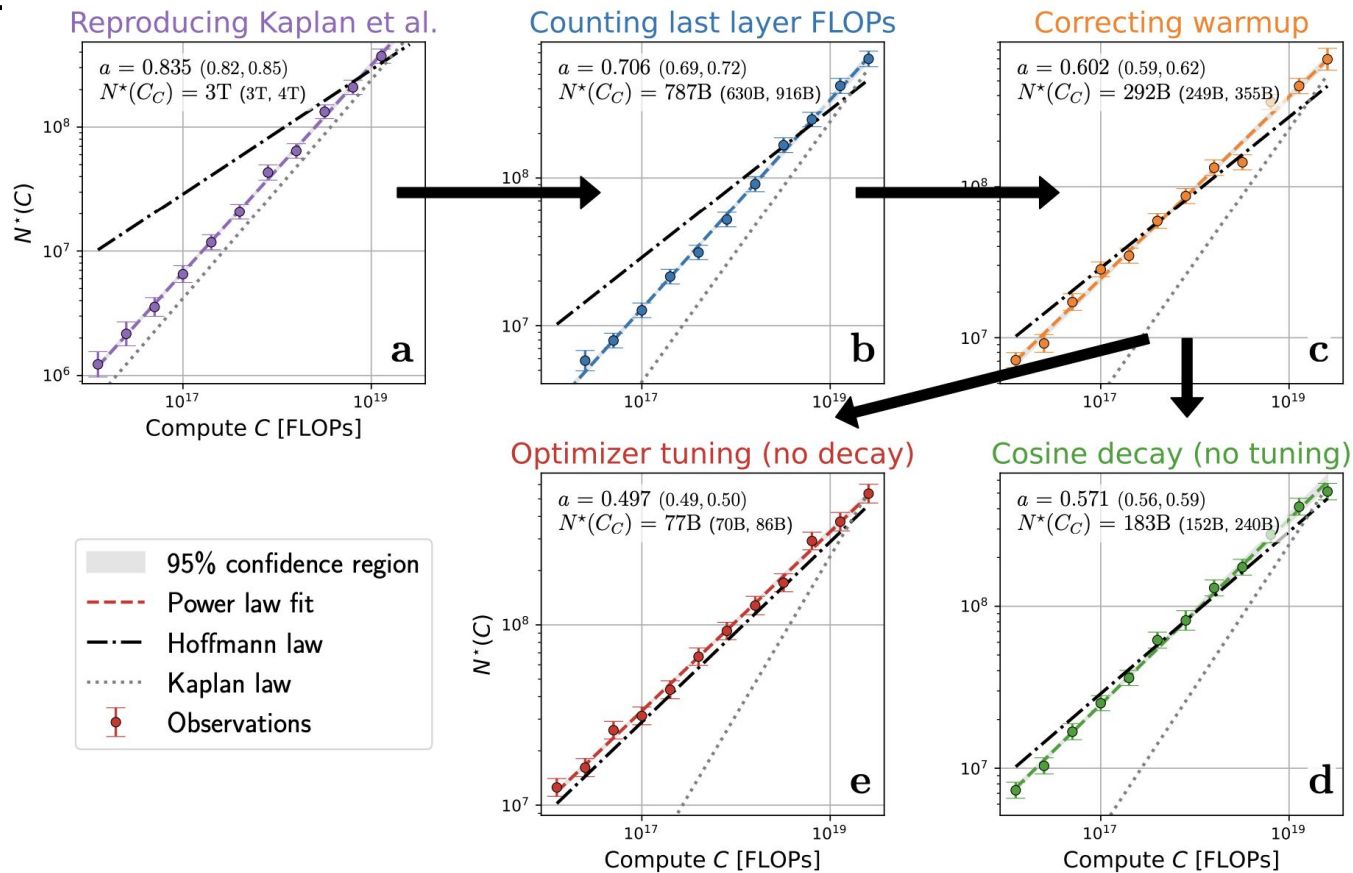
Optimizer tuning (no decay)



Cosine decay (no tuning)



# Scaling Laws



# Data Sourcing

WebText (GPT-2): Scraped content of websites receiving outward links from trusted Reddit posts.

Common Crawl: Scraped websites

Various books, trusted websites, Wikipedia

# Data Processing

- Deduplication using various hashing algorithms
- Heuristics to remove bad data
- Training small classifier models to tell apart good data and trash.
- Weighting different sources: better sources like books and Wikipedia can get 2-3 passes while less filtered sources like Common Crawl get less than 1.
- Separating different languages, removing formatting

# Zero-shot and Few-shot Tasks

Very large models experience a sort of emergent behavior: if the task to be performed can be expressed in language, then the text (task description, examples of acceptable output, input) can be fed directly into a language model to produce good output. This gets even better after fine-tuning.

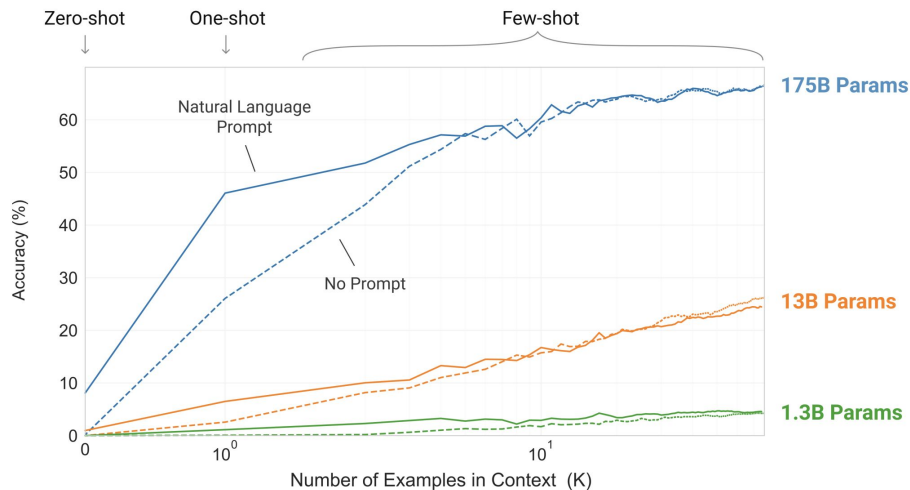
Language Models are Unsupervised Multitask Learners

	LAMBADA (PPL)	LAMBADA (ACC)	CBT-CN (ACC)	CBT-NE (ACC)	WikiText2 (PPL)	PTB (PPL)	enwik8 (BPB)	text8 (BPC)	WikiText103 (PPL)	1BW (PPL)
SOTA	99.8	59.23	85.7	82.3	39.14	46.54	0.99	1.08	18.3	<b>21.8</b>
117M	<b>35.13</b>	45.99	<b>87.65</b>	<b>83.4</b>	<b>29.41</b>	65.85	1.16	1.17	37.50	75.20
345M	<b>15.60</b>	55.48	<b>92.35</b>	<b>87.1</b>	<b>22.76</b>	47.33	1.01	<b>1.06</b>	26.37	55.72
762M	<b>10.87</b>	<b>60.12</b>	<b>93.45</b>	<b>88.0</b>	<b>19.93</b>	<b>40.31</b>	<b>0.97</b>	<b>1.02</b>	22.05	44.575
1542M	<b>8.63</b>	<b>63.24</b>	<b>93.30</b>	<b>89.05</b>	<b>18.34</b>	<b>35.76</b>	<b>0.93</b>	<b>0.98</b>	<b>17.48</b>	42.16

Various benchmarks for GPT-2 before fine-tuning. Bold is better than state of the art at the time.

# Zero-shot and Few-shot Tasks

(Task description, examples of task, input) → language model → Output



GPT-3: accuracy plotted against number of examples given.

# References

GPT-1,2,3:

[https://cdn.openai.com/research-covers/language-unsupervised/language\\_understanding\\_paper.pdf](https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf)

[https://cdn.openai.com/better-language-models/language\\_models\\_are\\_unsupervised\\_multitask\\_learners.pdf](https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf)

<https://arxiv.org/pdf/2005.14165>

BERT:

<https://arxiv.org/pdf/1810.04805>

WordPiece tokenization:

<https://huggingface.co/learn/llm-course/chapter6/6>



# References

Kaplan scaling laws:

<https://arxiv.org/pdf/2001.08361>

Chinchilla scaling laws:

<https://arxiv.org/pdf/2203.15556>

Reconciliation:

<https://arxiv.org/pdf/2406.12907>

<https://arxiv.org/abs/2406.19146>