

# Pretraining and Finetuning

He He



**NEW YORK UNIVERSITY**

October 11, 2023

# Logistics

- Please submit midterm course eval:  
[https://nyu.qualtrics.com/jfe/form/SV\\_6X7nHX4HenyFwN0](https://nyu.qualtrics.com/jfe/form/SV_6X7nHX4HenyFwN0)
- HW3 will be out next Monday
- Project: start early! Proposal due in two weeks.

# Table of Contents

Representation learning

# Representation learning

What are good representations?

- Enable a notion of distance over text (word embeddings)
- Contains good features for downstream tasks

# Representation learning

What are good representations?

- Enable a notion of distance over text (word embeddings)
- Contains good features for downstream tasks

Examples: negative the food is good but doesn't worth an hour wait

- Simple features (e.g. BoW) require complex models.
- Good features only need simple models (e.g. linear classifier).

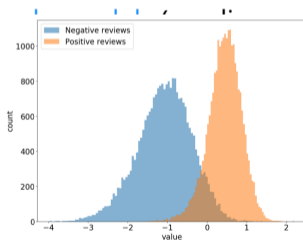


Figure: Sentiment neuron [Radford et al., 2017]

# Representation learning

What can we do with good representations:

- Learning with small data: fine-tuning learned representations
- Transfer learning: one model/representation for many tasks
- Metric learning: get a similarity metric for free

# Representation learning

What can we do with good representations:

- Learning with small data: fine-tuning learned representations
- Transfer learning: one model/representation for many tasks
- Metric learning: get a similarity metric for free

How to obtain such a representation:

- Training a neural network on any task gives us a representation good for *that task*.
- But on which task can we learn good *general* representations?

## What can we learn from word guessing?

- The cats that are raised by my sister \_\_\_\_\_ sleeping.



## What can we learn from word guessing?

- The cats that are raised by my sister \_\_\_\_\_ sleeping. *syntax*
- Jane is happy that John invited \_\_\_\_\_ friends to his birthday party.

## What can we learn from word guessing?

- The cats that are raised by my sister \_\_\_\_\_ sleeping. *syntax*
- Jane is happy that John invited \_\_\_\_\_ friends to his birthday party. *coreference*
- \_\_\_\_\_ is the capital of Tanzania.

## What can we learn from word guessing?

- The cats that are raised by my sister \_\_\_\_\_ sleeping. *syntax*
- Jane is happy that John invited \_\_\_\_\_ friends to his birthday party. *coreference*
- \_\_\_\_\_ is the capital of Tanzania. *knowledge*
- The boy is \_\_\_\_\_ because he lost his keys.

## What can we learn from word guessing?

- The cats that are raised by my sister \_\_\_\_\_ sleeping. *syntax*
- Jane is happy that John invited \_\_\_\_\_ friends to his birthday party. *coreference*
- \_\_\_\_\_ is the capital of Tanzania. *knowledge*
- The boy is \_\_\_\_\_ because he lost his keys. *commonsense*
- John took 100 bucks to Vegas. He won 50 and then lost 100. Now he only has \_\_\_\_\_ to go home.

## What can we learn from word guessing?

- The cats that are raised by my sister \_\_\_\_\_ sleeping. *syntax*
- Jane is happy that John invited \_\_\_\_\_ friends to his birthday party. *coreference*
- \_\_\_\_\_ is the capital of Tanzania. *knowledge*
- The boy is \_\_\_\_\_ because he lost his keys. *commonsense*
- John took 100 bucks to Vegas. He won 50 and then lost 100. Now he only has \_\_\_\_\_ to go home. *numerical reasoning*

## What can we learn from word guessing?

- The cats that are raised by my sister \_\_\_\_\_ sleeping. *syntax*
- Jane is happy that John invited \_\_\_\_\_ friends to his birthday party. *coreference*
- \_\_\_\_\_ is the capital of Tanzania. *knowledge*
- The boy is \_\_\_\_\_ because he lost his keys. *commonsense*
- John took 100 bucks to Vegas. He won 50 and then lost 100. Now he only has \_\_\_\_\_ to go home. *numerical reasoning*

Word guessing entails many tasks related to language understanding!



But aren't we already doing this in skip-gram / CBOW?

# Self-supervised learning

**Key idea:** predict parts of the input from the rest

- **No supervision** is needed—both input and output are from the raw data.
- Easy to **scale**—only need unlabeled data.
- Learned representation is **general**—useful for many tasks.

# Self-supervised learning

**Key idea:** predict parts of the input from the rest

- **No supervision** is needed—both input and output are from the raw data.
- Easy to **scale**—only need unlabeled data.
- Learned representation is **general**—useful for many tasks.

**Approach:**

- **Pretrain:** train a model using self-supervised learning objectives on large data.
- **Finetune:** update part or all of the parameters of the pretrained model (which provides an initialization) on labeled data of a downstream task.



## A bit of history

- Pretrain an RNN model on unlabeled data and finetune on supervised tasks [Dai et al., 2015] [ULMFiT; Howard et al., 2018]
  - Promising results on a small scale

## A bit of history

- Pretrain an RNN model on unlabeled data and finetune on supervised tasks [Dai et al., 2015] [ULMFiT; Howard et al., 2018]
  - Promising results on a small scale
- ELMo: replace static word embedding by contextual word embeddings from pretrained bi-LSTM [Peters et al., 2018]
  - First impactful result in NLP

## A bit of history

- Pretrain an **RNN** model on unlabeled data and finetune on supervised tasks [Dai et al., 2015] [ULMFiT; Howard et al., 2018]
  - Promising results on a small scale
- ELMo: replace static word embedding by **contextual word embeddings** from pretrained bi-LSTM [Peters et al., 2018]
  - First impactful result in NLP
- Pretrain a **Transformer** model and finetune on supervised tasks
  - GPT [Radford et al., 2018], BERT [Devlin et al., 2018]

## A bit of history

- Pretrain an **RNN** model on unlabeled data and finetune on supervised tasks [Dai et al., 2015] [ULMFiT; Howard et al., 2018]
  - Promising results on a small scale
- ELMo: replace static word embedding by **contextual word embeddings** from pretrained bi-LSTM [Peters et al., 2018]
  - First impactful result in NLP
- Pretrain a **Transformer** model and finetune on supervised tasks
  - GPT [Radford et al., 2018], BERT [Devlin et al., 2018]
- **Scale** the pretrained model to larger sizes
  - GPT-2 (1.5B), T5 (11B), GPT-3 (175B), PaLM (540B)
  - We will talk about 100B+ models in the third module

# Types of pretrained models

- **Encoder models**, e.g., BERT
  - Encode text into vector representations that can be used for downstream classification tasks

Current pretrained models are all transformer based.

## Types of pretrained models

- **Encoder models**, e.g., BERT
  - Encode text into vector representations that can be used for downstream classification tasks
- **Encoder-decoder models**, e.g., T5
  - Encode input text into vector representations and generate text conditioned on the input

Current pretrained models are all transformer based.

## Types of pretrained models

- **Encoder models**, e.g., BERT
  - Encode text into vector representations that can be used for downstream classification tasks
- **Encoder-decoder models**, e.g., T5
  - Encode input text into vector representations and generate text conditioned on the input
- **Decoder models**, e.g., GPT-2
  - Read in text (prefix) and continue to generate text

Current pretrained models are all transformer based.

## Encoder models

An encoder takes a sequence of tokens and output their *contextualized* representations:

$$h_1, \dots, h_n = \text{Encoder}(x_1, \dots, x_n)$$

We can then use  $h_1, \dots, h_n$  for other tasks.



## Encoder models

An encoder takes a sequence of tokens and output their *contextualized* representations:

$$h_1, \dots, h_n = \text{Encoder}(x_1, \dots, x_n)$$

We can then use  $h_1, \dots, h_n$  for other tasks.

How do we train an Encoder?

- Use any supervised task:  $y = f(h_1, \dots, h_n)$
- Use self-supervised learning: predict a word from its context

# Masked language modeling

? language processing is ?

# Masked language modeling

? language processing is ?

Learning objective (MLE):

$$\max_{x \in \mathcal{D}, i \sim p_{\text{mask}}} \log p(x_i | x_{-i}; \theta)$$

- $x$ : a sequence of tokens sampled from a corpus  $\mathcal{D}$   
*natural language processing is fun*
- $p_{\text{mask}}$ : mask generator  
Sample two positions uniformly at random, e.g., 1 and 5
- $x_{-i}$ : noisy version of  $x$  where  $x_i$  is corrupted  
*[MASK] language processing is [MASK]*

## BERT: objective

- **Masked language modeling:**
  - Randomly sample 15% tokens as prediction targets
  - Replace the target tokens by [MASK] or a random token, or leave it unchanged
    - cats are cute → cats [MASK] /is/are cute
  - Later work has shown that just use [MASK] is sufficient

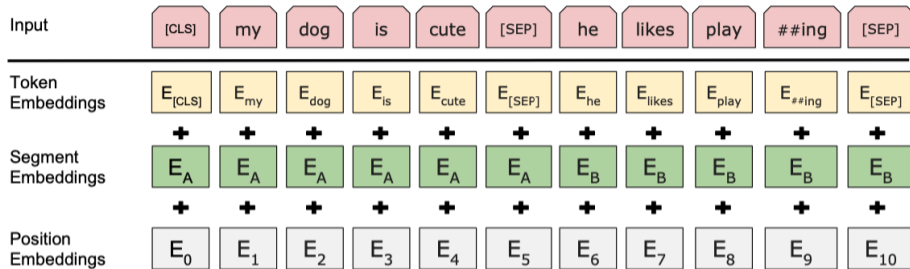
## BERT: objective

- **Masked language modeling:**
  - Randomly sample 15% tokens as prediction targets
  - Replace the target tokens by [MASK] or a random token, or leave it unchanged  
cats are cute → cats [MASK] /is/are cute
  - Later work has shown that just use [MASK] is sufficient
- **Next sentence prediction:** predict whether a pair of sentences are consecutive

$$\max_{x \sim \mathcal{D}, x_n \sim p_{\text{next}}} \sum \log p(y \mid x, x_n; \theta)$$

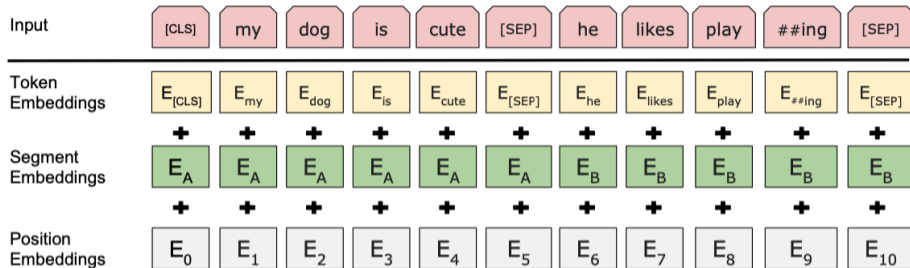
- $x_n$ : either the sentence following  $x$  or a randomly sampled sentence
- $y$ : binary label of whether  $x_n$  follows  $x$
- Later work has shown that this objective is not necessary

# BERT: architecture



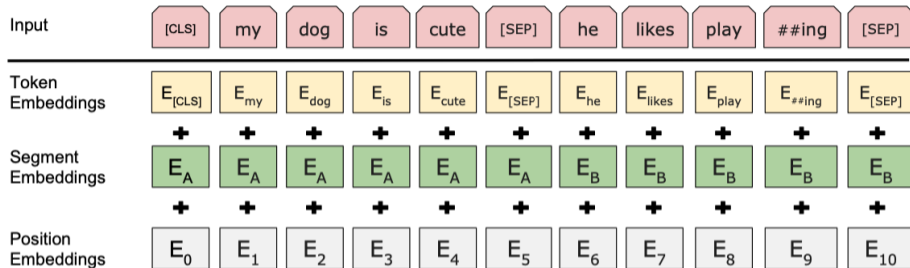
- Tokenization: wordpiece (similar to byte pair encoding) (see [details](#))

# BERT: architecture



- Tokenization: wordpiece (similar to byte pair encoding) (see [details](#))
- [CLS]: first token of all sequences; used for next sentence prediction

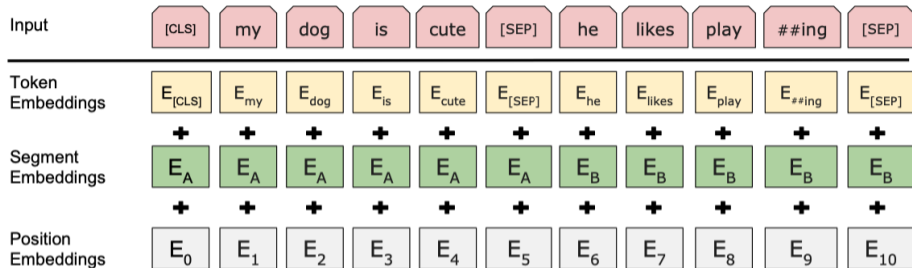
# BERT: architecture



- Tokenization: wordpiece (similar to byte pair encoding) (see [details](#))
- [CLS]: first token of all sequences; used for next sentence prediction
- Distinguish two sentences in a pair: [SEP] and segment embedding

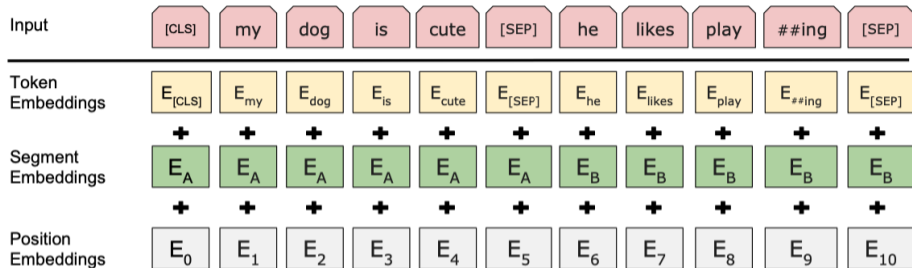


# BERT: architecture



- Tokenization: wordpiece (similar to byte pair encoding) (see [details](#))
- [CLS]: first token of all sequences; used for next sentence prediction
- Distinguish two sentences in a pair: [SEP] and segment embedding
- Learned position embedding

# BERT: architecture

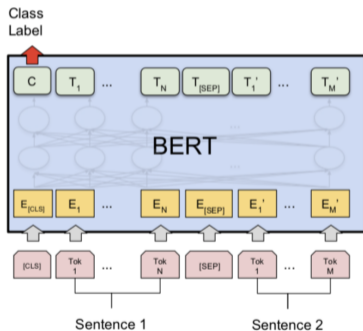


- Tokenization: wordpiece (similar to byte pair encoding) (see [details](#))
- [CLS]: first token of all sequences; used for next sentence prediction
- Distinguish two sentences in a pair: [SEP] and segment embedding
- Learned position embedding
- 12 (base; 110M params) or 24 (large; 340M params) layer Transformer

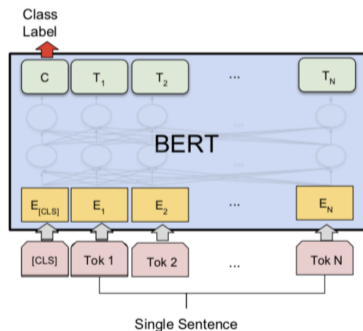
# Finetuning BERT

Classification tasks: Add a linear layer (randomly initialized) on top of the [CLS] embedding

$$p(y | x) = \text{softmax}(Wh_{[\text{CLS}]})$$



(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG

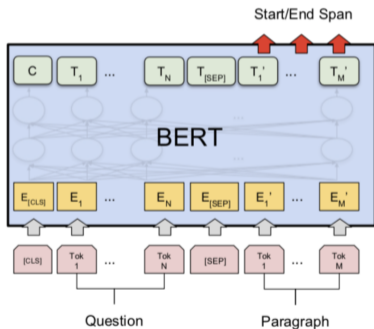


(b) Single Sentence Classification Tasks:  
SST-2, CoLA

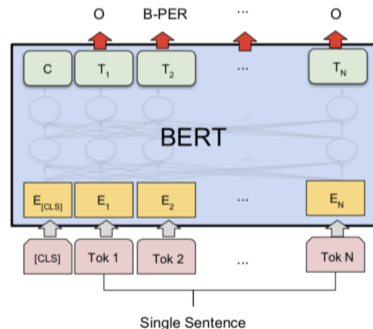
# Finetuning BERT

Sequence labeling tasks: Add linear layers (randomly initialized) on top of every token

$$p(y_i | x) = \text{softmax}(Wh_i)$$




(c) Question Answering Tasks:  
SQuAD v1.1



(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER

# Finetuning BERT

- Finetune all parameters (both the newly added layer and the pretrained weights)
- Use a small learning rate (e.g.,  $1e-5$ )
- Train for a small number of epochs (e.g, 3 epochs)
- Led to SOTA results on many NLU tasks
-  How to generate text from BERT?

## Encoder-decoder models

An encoder-decoder model encodes input text to a sequence of contextualized representations, and decodes a sequence of tokens autoregressively.

$$h_1, \dots, h_n = \text{Encoder}(x_1, \dots, x_n)$$

$$s_1, \dots, s_m = \text{Decoder}(y_0, \dots, y_{m-1}, h_1, \dots, h_n)$$

$$p(y_i | x, y_{<i}) = \text{softmax}(Ws_i)$$

## Encoder-decoder models

An encoder-decoder model encodes input text to a sequence of contextualized representations, and decodes a sequence of tokens autoregressively.

$$h_1, \dots, h_n = \text{Encoder}(x_1, \dots, x_n)$$

$$s_1, \dots, s_m = \text{Decoder}(y_0, \dots, y_{m-1}, h_1, \dots, h_n)$$

$$p(y_i | x, y_{<i}) = \text{softmax}(Ws_i)$$

How do we train the encoder-decoder?

- Use any supervised task, e.g., machine translation
- Use self-supervised learning: predict text spans from their context

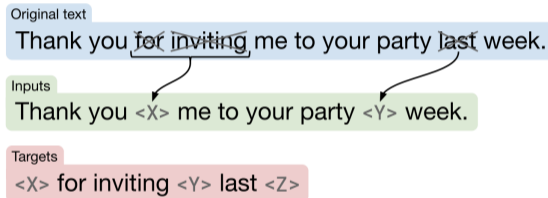


How to train an encoder-decoder model using the BERT objective?

# Masked language modeling using an encoder-decoder

**Input:** text with corrupted spans

**Output:** recovered spans



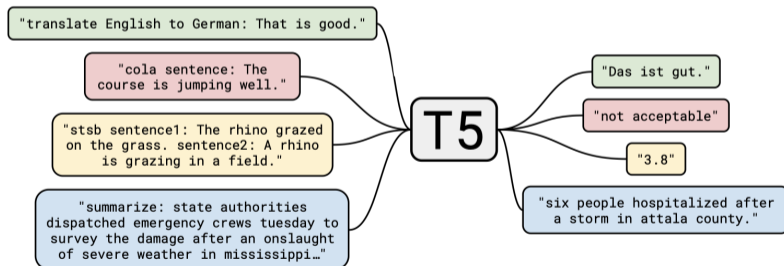
Compare with encoder-only models:

- Encoder: predict single tokens based on encoder representation
- Encoder-decoder: predict a sequence of tokens (flexibility in objective design)



## T5: objective

- First train on unlabeled data by **masked language modeling**
  - Predict corrupted spans as a sequence
- Then **continue training** by **supervised multitask learning**
  - Formulate tasks as text-to-text format using a prefix to denote the task
  - Mixing examples from different datasets when constructing batches



- Jointly training with the two objectives works slightly worse

## T5: finetune

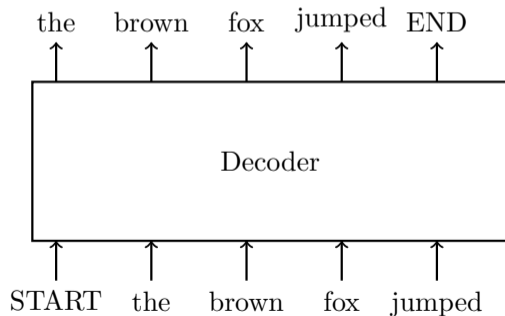
- Formulate the task in text-to-text format
- Fine-tune all parameters (similar to BERT fine-tuning)
- Advantages over encoder models: unified modeling of many different tasks including text generation

## Decoder-only models

A decoder-only model predicts the next token given the prefix autoregressively.

$$s_1, \dots, s_m = \text{Decoder}(y_0, \dots, y_{m-1}, h_1, \dots, h_n)$$
$$p(y_i | y_{<i}) = \text{softmax}(Ws_i)$$

(A prefix of  $y$  can be the input.)



(more on language models later)

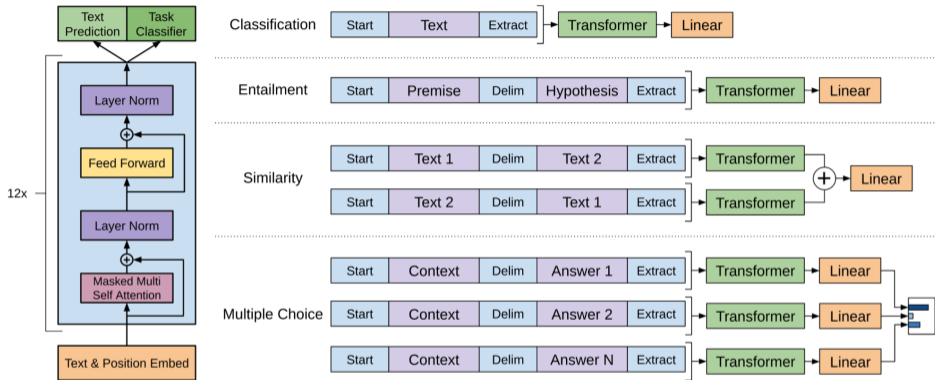
# Generative Pretraining (GPT)

- **Model:** 12 layer decoder-only transformer
- **Objective:** next word prediction

$$\max \sum_{y \in \mathcal{D}} \sum_i \log p(y_i | y_{<i})$$

- **Finetuning:** auxiliary LM objective  $L_{\text{task}} + \lambda L_{\text{LM}}$  (next word prediction on labeled task data)

# Generative Pretraining (GPT): task-specific finetuning



- Single input: linear on top of extract
- Multiple input: process each input separately then aggregate

# Ablation studies of GPT

Architecture, pretraining, finetuning: which is critical?

Method	Avg. Score	CoLA (mc)	SST2 (acc)	MRPC (F1)	STSB (pc)	QQP (F1)	MNLI (acc)	QNLI (acc)	RTE (acc)
Transformer w/ aux LM (full)	74.7	45.4	91.3	82.3	82.0	<b>70.3</b>	<b>81.8</b>	<b>88.1</b>	<b>56.0</b>
Transformer w/o pre-training	59.9	18.9	84.0	79.4	30.9	65.5	75.7	71.2	53.8
Transformer w/o aux LM	<b>75.0</b>	<b>47.9</b>	<b>92.0</b>	<b>84.9</b>	<b>83.2</b>	69.8	81.1	86.9	54.4
LSTM w/ aux LM	69.1	30.3	90.5	83.2	71.8	68.1	73.7	81.1	54.6

- Auxiliary objective only helps on larger datasets (MNLI, QQP)
- Pretrained transformer > pretrained LSTM (single layer) > non-pretrained transformer

# Compare with BERT

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 <sub>BASE</sub>	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>92.7</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>82.1</b>

Table 1: GLUE Test results, scored by the evaluation server (<https://gluebenchmark.com/leaderboard>). The number below each task denotes the number of training examples. The “Average” column is slightly different than the official GLUE score, since we exclude the problematic WNLI set.<sup>8</sup> BERT and OpenAI GPT are single-model, single task. F1 scores are reported for QQP and MRPC, Spearman correlations are reported for STS-B, and accuracy scores are reported for the other tasks. We exclude entries that use BERT as one of their components.

Medium-sized encoder models tend to work better than decoder-only models when finetuned

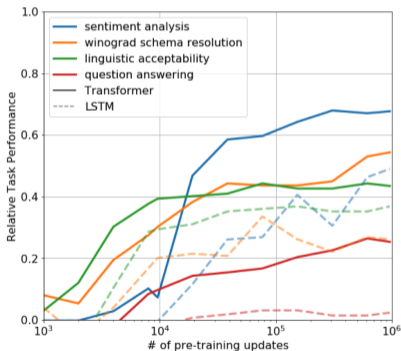
## Zero-shot behaviors

**Key insight:** if the model has learned to understand language through predicting next words, it should be able to perform these tasks *without finetuning*



# Zero-shot behaviors

**Key insight:** if the model has learned to understand language through predicting next words, it should be able to perform these tasks *without finetuning*



Heuristics for zero-shot prediction:

- Sentiment classification: [example] + very + {positive, negative} *prompting*
- Linguistic acceptability: thresholding on log probabilities
- Multiple choice: predicting the answer with the highest log probabilities

**Learning dynamics:** zero-shot performance increases during pretraining

# What are these models trained on?

Both quantity and quality are important

- Wikipedia: encyclopedia articles (clean, single domain)
- Toronto Books Corpus: e-books (diverse domain)
- WebText (40GB): content submitted to Reddit with a vote  $\geq 3$  (diverse, bias)
- CommonCrawl (20TB): scraped HTML with markers removed (diverse, large, noisy, bias)
  - A cleaned version: C4 (750GB)

Active research area: What data is good for pretraining?

# Summary

Lots of learning happens from just observing the world (data).

- Self-supervised learning: [benefits from large data and compute](#)
  - Basic: predict parts from other parts based on the structure of data (works beyond text)
  - Advanced: design hard negatives to improve efficiency
- Finetuning: adapt pretrained models to downstream tasks on a small amount of labeled data