

# Pretraining Language Models (part 2)

Wei Xu

(many slides from Greg Durrett, Alan Ritter)

# Administrivia

---

- ▶ Readings —

- ▶ T5 by Raffel et al.

<https://arxiv.org/abs/1910.10683>

- ▶ GPT-3 by Brown et al.

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

# This and Next Lecture

---

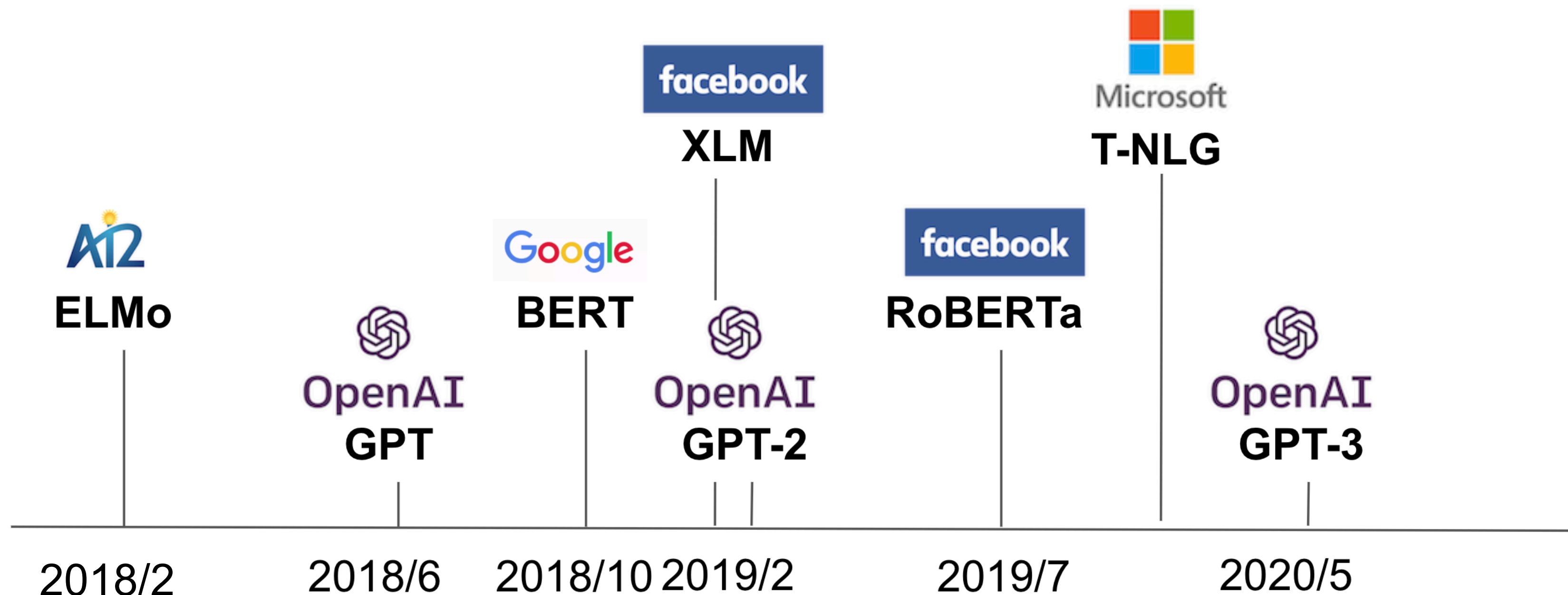
- ▶ BERT (recap)
- ▶ GPT / GPT-2
- ▶ BART / T5
- ▶ GPT-3
- ▶ T0/Flan/PaLM/OPT (likely next class)

# Recap: Context-dependent Embeddings

- ▶ AI2 released ELMo in spring 2018, GPT was released in summer 2018, BERT came out October 2018
- ▶ aka Pre-trained Language Models

and many more:

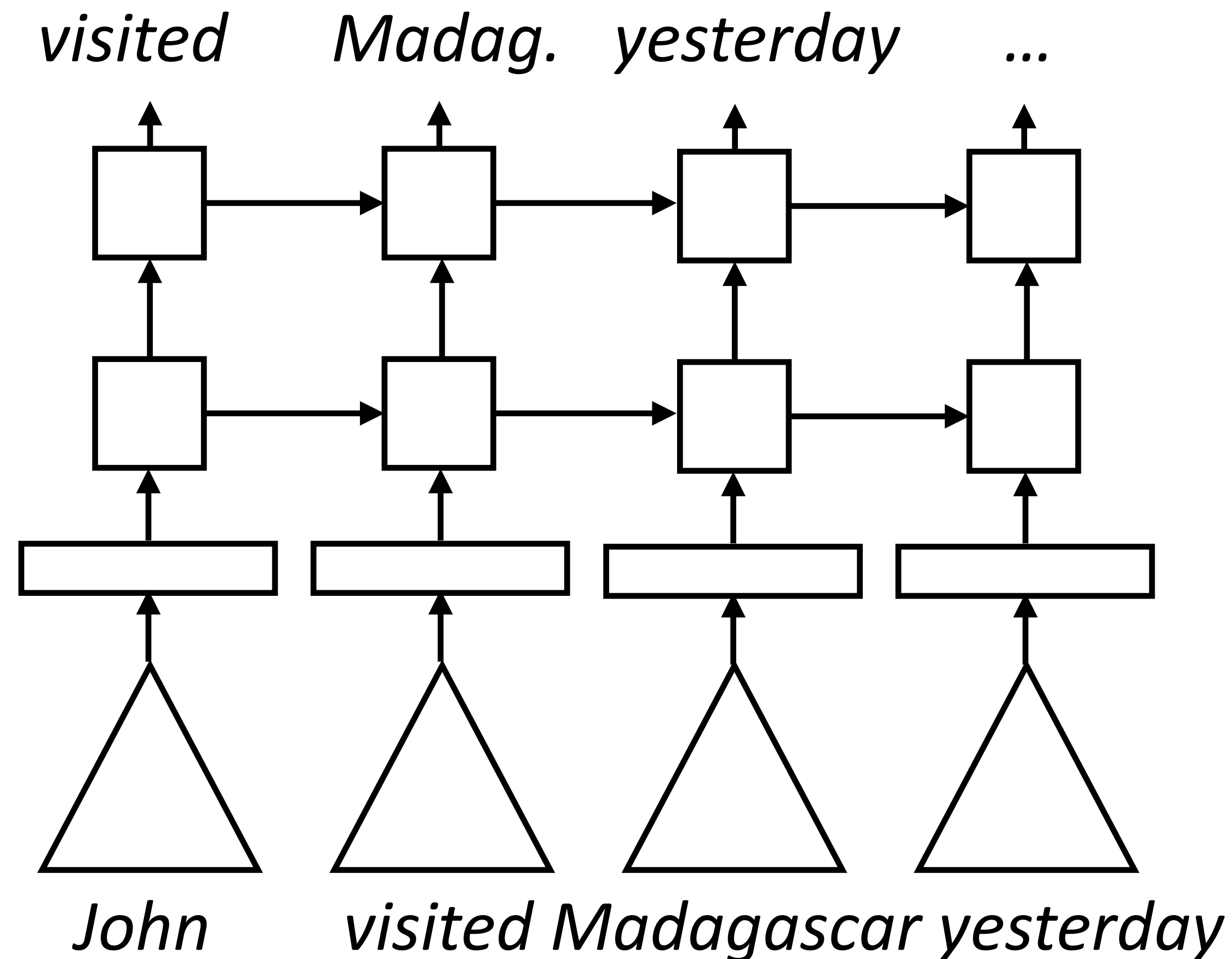
- BART
- DialoGPT
- GPT-J
- T5
- T0
- OPT
- PaLM
- ChatGPT
- Llama
- Mistral ...



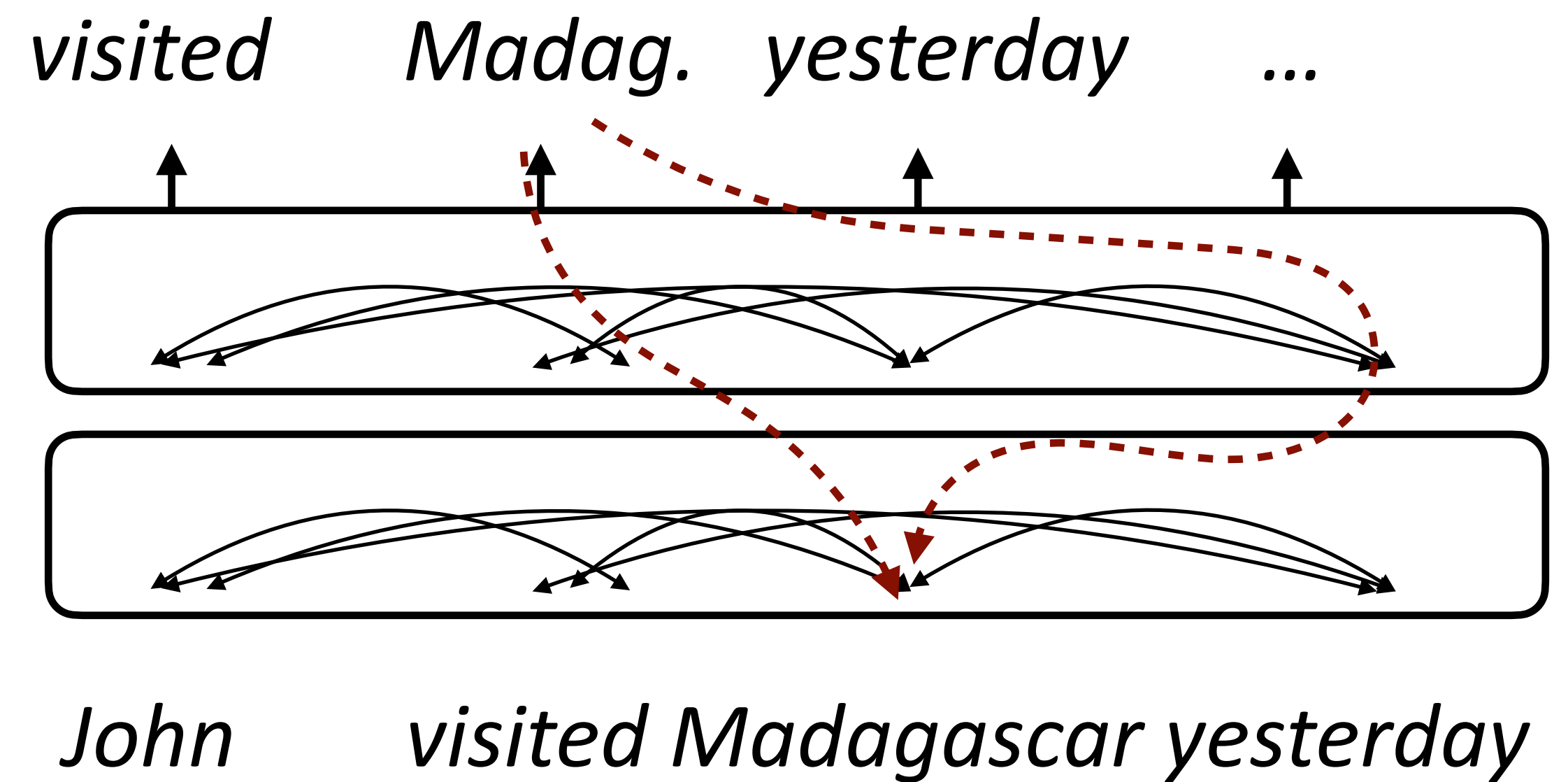
# Recall: BERT

- ▶ How to learn a “deeply bidirectional” model? What happens if we just replace an LSTM with a transformer?

## ELMo (Language Modeling)



## BERT



- ▶ Transformer LMs have to be “one-sided” (only attend to previous tokens), not what we want

# Recall: Masked Language Modeling

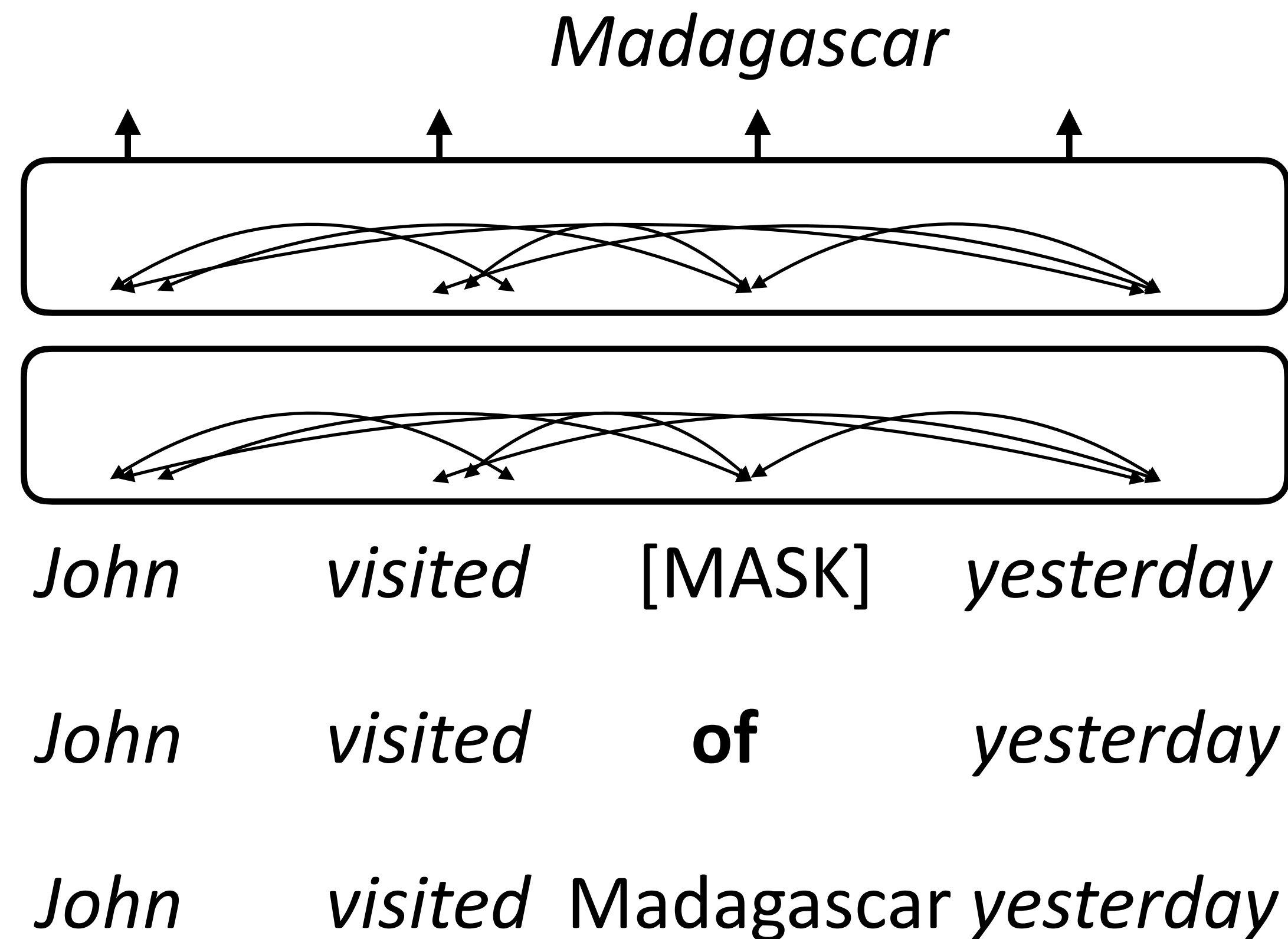
- ▶ How to prevent cheating? Next word prediction fundamentally doesn't work for bidirectional models, instead do *masked language modeling*

- ▶ BERT formula: take a chunk of text, predict 15% of the tokens

- ▶ For 80% (of the 15%), replace the input token with [MASK]

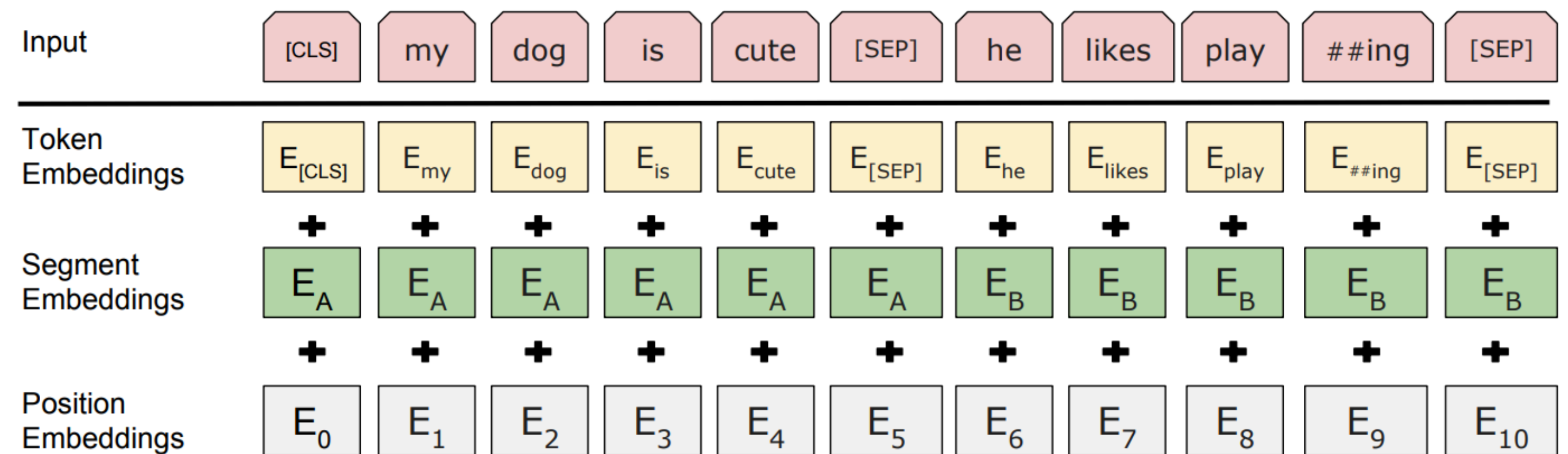
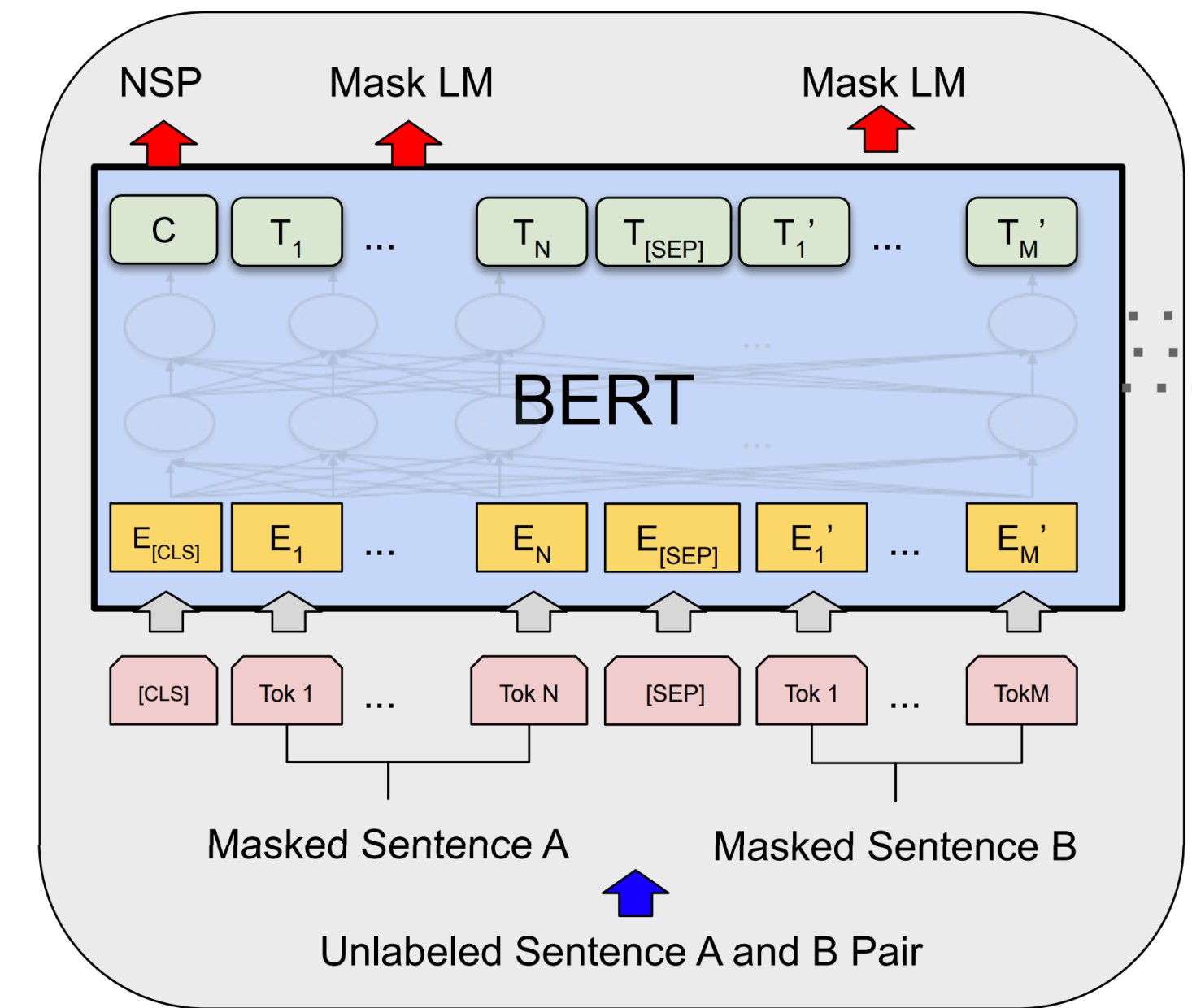
- ▶ For 10%, replace w/random

- ▶ For 10%, keep same (why?)

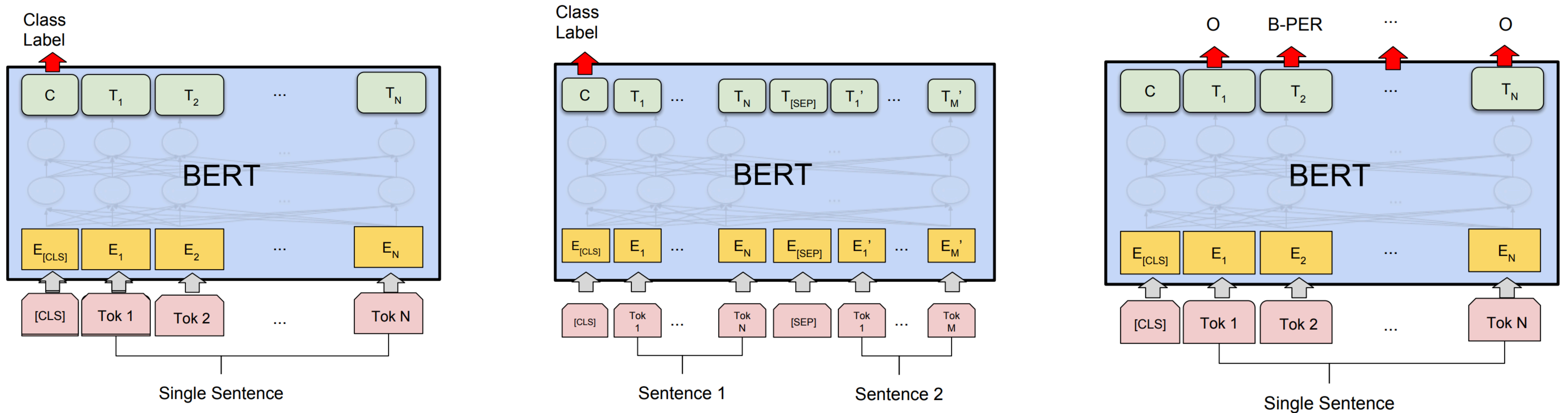


# Recall: BERT Architecture

- ▶ BERT Base: 12 layers, 768-dim, 12 heads. Total params = 110M
- ▶ BERT Large: 24 layers, 1024-dim, 16 heads. Total params = 340M
- ▶ Positional embeddings and segment embeddings, 30k word pieces
- ▶ This is the model that gets **pre-trained** on a large corpus



# Recall: What can BERT do?



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

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

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

- ▶ CLS token is used to provide classification decisions
- ▶ Sentence pair tasks (entailment): feed both sentences into BERT
- ▶ BERT can also do tagging by predicting tags at each word piece

Devlin et al. (2019)



# Recall: RoBERTa

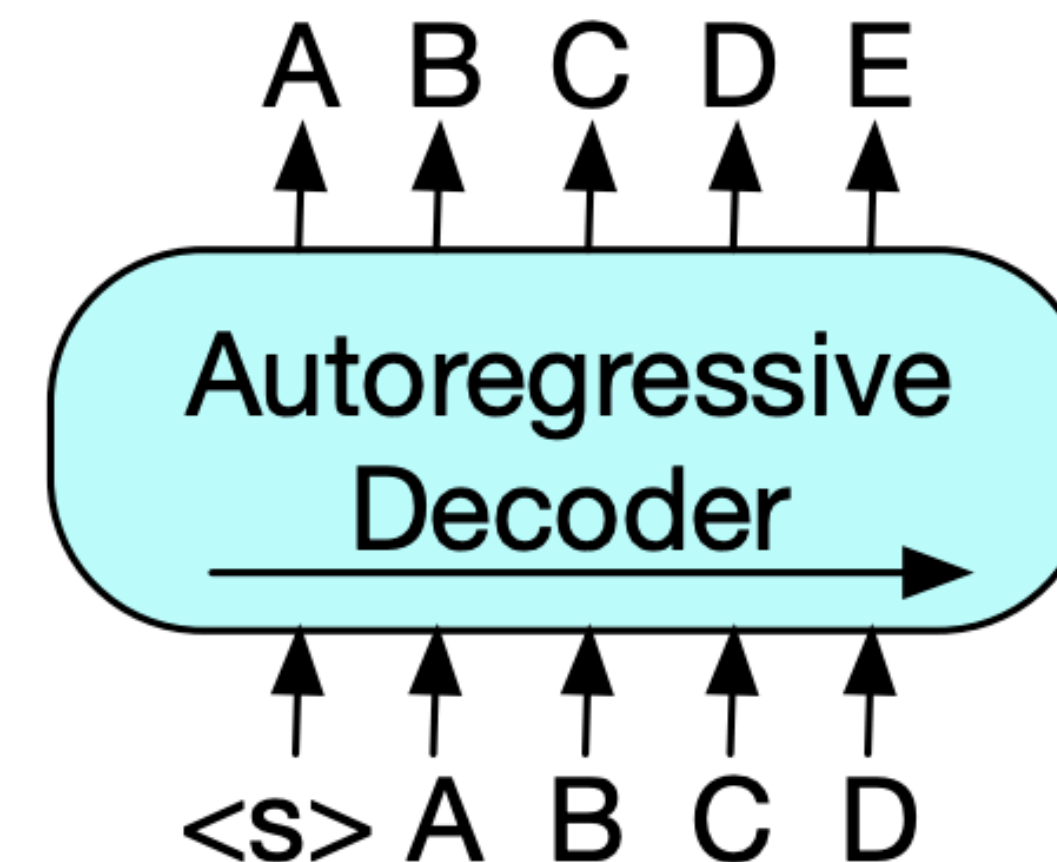
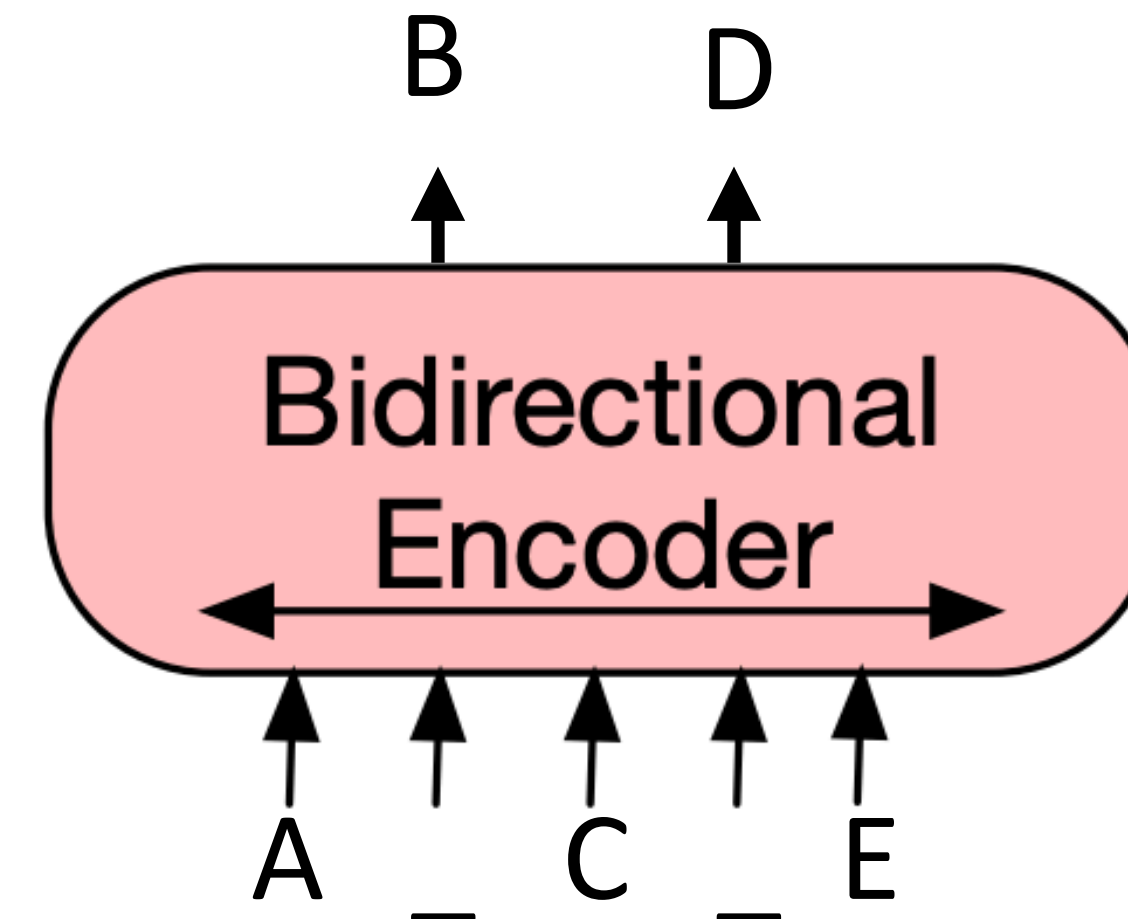
- ▶ “Robustly optimized BERT”
- ▶ 160GB of data instead of 16 GB
- ▶ Dynamic Masking: standard BERT uses the same MASK scheme for every epoch, RoBERTa recomputes them
- ▶ New training + more data = better performance
- ▶ For this and more: check out Huggingface or fairseq

Model	data	bsz	steps	SQuAD (v1.1/2.0)	MNLI-m	SST-2
RoBERTa						
with BOOKS + WIKI	16GB	8K	100K	93.6/87.3	89.0	95.3
+ additional data (§3.2)	160GB	8K	100K	94.0/87.7	89.3	95.6
+ pretrain longer	160GB	8K	300K	94.4/88.7	90.0	96.1
+ pretrain even longer	160GB	8K	500K	<b>94.6/89.4</b>	<b>90.2</b>	<b>96.4</b>
BERT <sub>LARGE</sub>						
with BOOKS + WIKI	13GB	256	1M	90.9/81.8	86.6	93.7

BART/T5  
(seq2seq type of LMs)

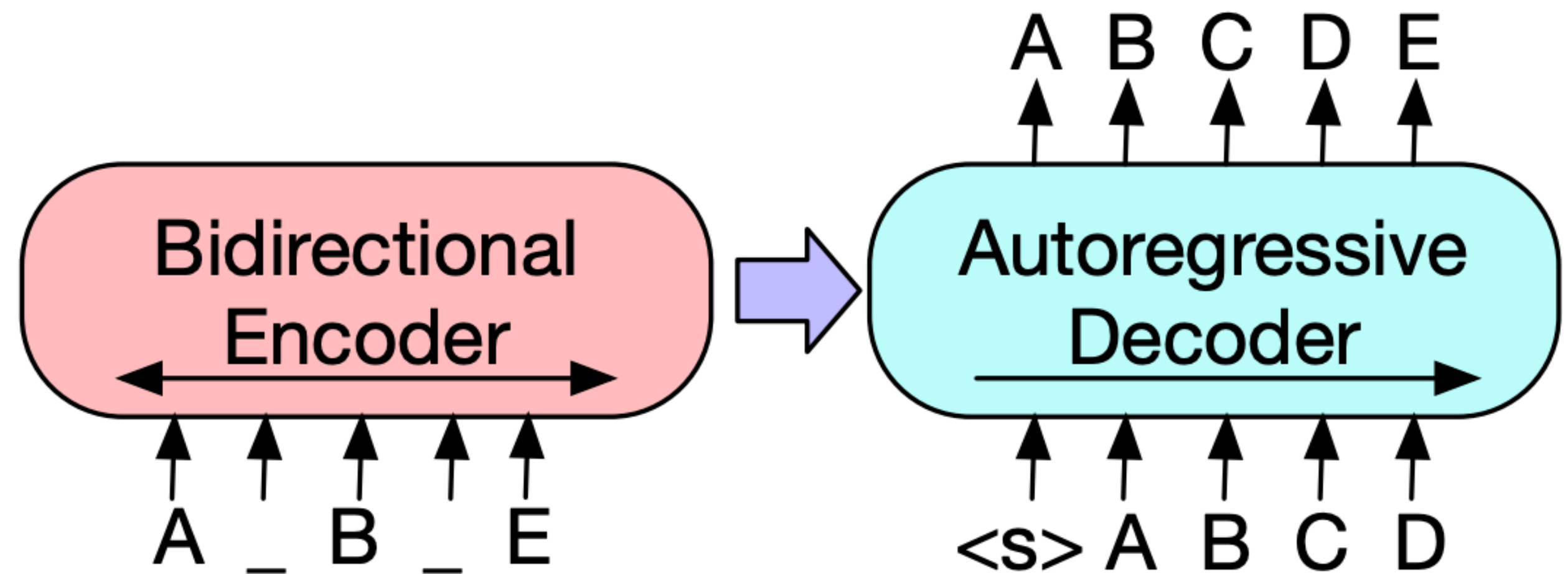
# BERT vs. GPT

- ▶ BERT: only parameters are an encoder, trained with masked language modeling objective
  - ▶ No way to do translation or left-to-right language modeling tasks
- ▶ GPT: only the decoder, autoregressive LM
  - ▶ (Small-size versions) Typically used for unconditioned generation tasks, e.g. story or dialog generation



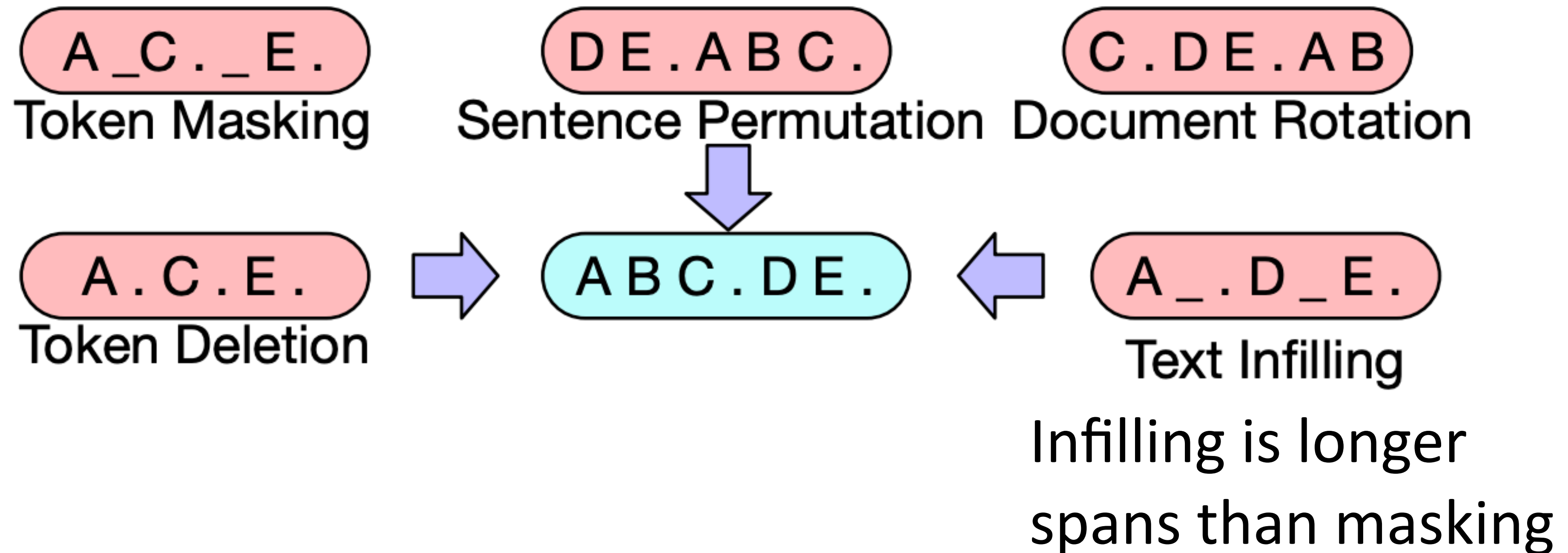
# BART

- ▶ What to do for seq2seq tasks?
- ▶ Sequence-to-sequence BERT variant: permute/make/delete tokens, then predict full sequence autoregressively
- ▶ For downstream tasks: feed document into both encoder + decoder, use decoder hidden state as output
- ▶ Good results on dialogue, summarization tasks



# BART

- ▶ BART uses multiple de-noising LM objective:



# BART

Model	SQuAD 1.1 F1	MNLI Acc	ELI5 PPL	XSum PPL	ConvAI2 PPL	CNN/DM PPL
BART Base						
w/ Token Masking	90.4	84.1	25.05	7.08	11.73	6.10
w/ Token Deletion	90.4	84.1	24.61	6.90	11.46	5.87
w/ Text Infilling	<b>90.8</b>	84.0	24.26	<b>6.61</b>	<b>11.05</b>	5.83
w/ Document Rotation	77.2	75.3	53.69	17.14	19.87	10.59
w/ Sentence Shuffling	85.4	81.5	41.87	10.93	16.67	7.89
w/ Text Infilling + Sentence Shuffling	<b>90.8</b>	83.8	24.17	6.62	11.12	<b>5.41</b>

- ▶ Infilling is all-around a bit better than masking or deletion
- ▶ Final system: combination of infilling and sentence permutation

# BART

	SQuAD 1.1 EM/F1	SQuAD 2.0 EM/F1	MNLI m/mm	SST Acc	QQP Acc	QNLI Acc	STS-B Acc	RTE Acc	MRPC Acc	CoLA Mcc
BERT	84.1/90.9	79.0/81.8	86.6/-	93.2	91.3	92.3	90.0	70.4	88.0	60.6
UniLM	-/-	80.5/83.4	87.0/85.9	94.5	-	92.7	-	70.9	-	61.1
XLNet	<b>89.0/94.5</b>	86.1/88.8	89.8/-	95.6	91.8	93.9	91.8	83.8	89.2	63.6
RoBERTa	88.9/ <b>94.6</b>	<b>86.5/89.4</b>	<b>90.2/90.2</b>	96.4	92.2	94.7	<b>92.4</b>	86.6	<b>90.9</b>	<b>68.0</b>
BART	88.8/ <b>94.6</b>	86.1/89.2	89.9/90.1	<b>96.6</b>	<b>92.5</b>	<b>94.9</b>	91.2	<b>87.0</b>	90.4	62.8

- ▶ Results on GLUE benchmark are not better than RoBERTa

# CoLA

- ▶ Corpus of Linguistic Acceptability (CoLA); to test whether a model can recognize (a) morphological anomalies, (b) syntactic anomalies, and (c) semantic anomalies.

Label	Sentence	Source
*	The more books I ask to whom he will give, the more he reads.	Culicover and Jackendoff (1999)
✓	I said that my father, he was tight as a hoot-owl.	Ross (1967)
✓	The jeweller inscribed the ring with the name.	Levin (1993)
*	many evidence was provided.	Kim and Sells (2008)
✓	They can sing.	Kim and Sells (2008)
✓	The men would have been all working.	Baltin (1982)
*	Who do you think that will question Seamus first?	Carnie (2013)
*	Usually, any lion is majestic.	Dayal (1998)
✓	The gardener planted roses in the garden.	Miller (2002)
✓	I wrote Blair a letter, but I tore it up before I sent it.	Rappaport Hovav and Levin (2008)

(✓= acceptable, \*=unacceptable)

Warstadt et al. (2020)



# BART for Summarization

---

This is the first time anyone has been recorded to run a full marathon of 42.195 kilometers (approximately 26 miles) under this pursued landmark time. It was not, however, an officially sanctioned world record, as it was not an "open race" of the IAAF. His time was 1 hour 59 minutes 40.2 seconds. Kipchoge ran in Vienna, Austria. It was an event specifically designed to help Kipchoge break the two hour barrier.

Kenyan runner Eliud Kipchoge has run a marathon in less than two hours.

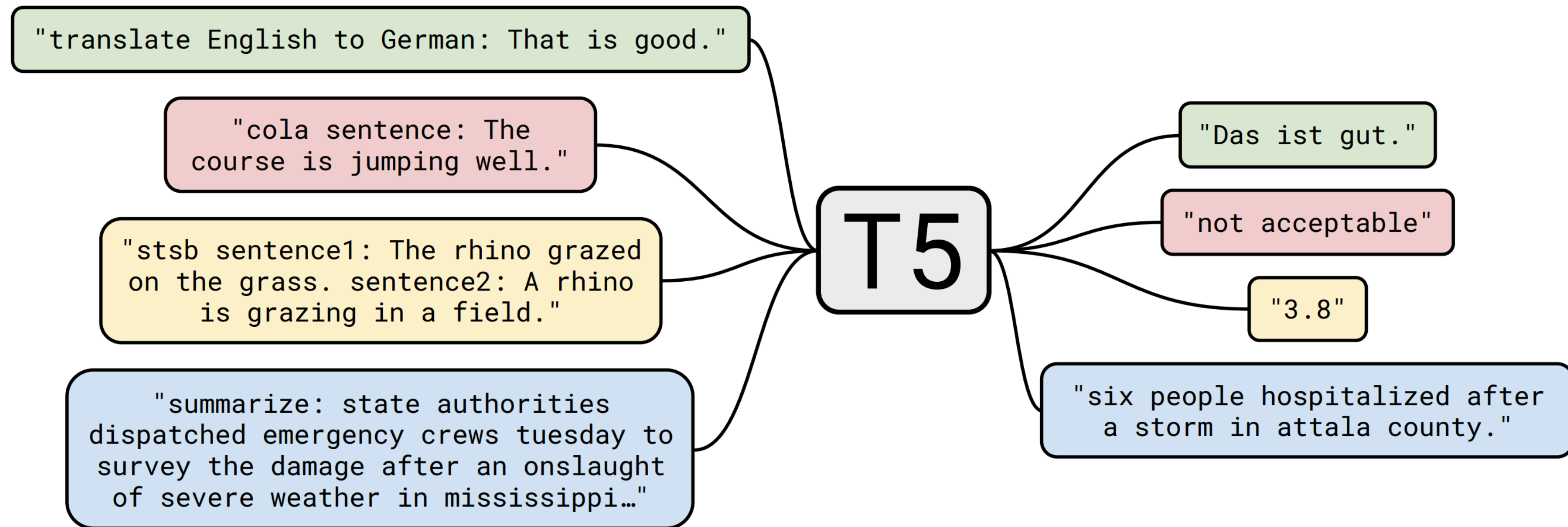
PG&E stated it scheduled the blackouts in response to forecasts for high winds amid dry conditions. The aim is to reduce the risk of wildfires. Nearly 800 thousand customers were scheduled to be affected by the shutoffs which were expected to last through at least midday tomorrow.

Power has been turned off to millions of customers in California as part of a power shutoff plan.

- 
- ▶ But, strong results on dialogue, summarization, and other generation tasks.

# T5

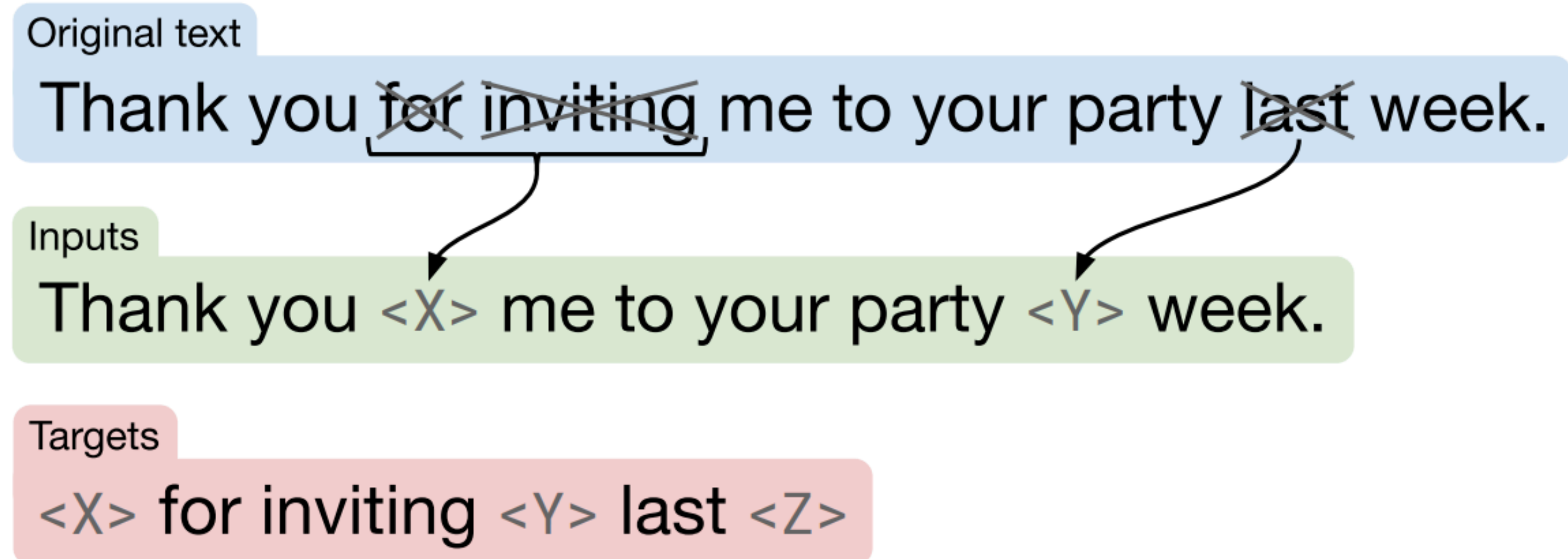
- ▶ Frame many problems as sequence-to-sequence ones:



# T5

---

- ▶ Pre-training: similar denoising scheme to BART



- ▶ Different mask tokens for individual masked spans; also different format for targets

# T5

- ▶ Compared several different unsupervised LM objectives:

Objective	Inputs	Targets
Prefix language modeling	Thank you for inviting	me to your party last week .
BERT-style <a href="#">Devlin et al. (2018)</a>	Thank you <M> <M> me to your party apple week .	<i>(original text)</i>
Deshuffling	party me for your to . last fun you inviting week Thank	<i>(original text)</i>
MASS-style <a href="#">Song et al. (2019)</a>	Thank you <M> <M> me to your party <M> week .	<i>(original text)</i>
I.i.d. noise, replace spans	Thank you <X> me to your party <Y> week .	<X> for inviting <Y> last <Z>
I.i.d. noise, drop tokens	Thank you me to your party week .	for inviting last
Random spans	Thank you <X> to <Y> week .	<X> for inviting me <Y> your party last <Z>

# T5

---

Number of tokens	Repeats	GLUE	CNNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Full dataset	0	<b>83.28</b>	<b>19.24</b>	<b>80.88</b>	<b>71.36</b>	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
$2^{29}$	64	<b>82.87</b>	<b>19.19</b>	<b>80.97</b>	<b>72.03</b>	<b>26.83</b>	<b>39.74</b>	<b>27.63</b>
$2^{27}$	256	82.62	<b>19.20</b>	79.78	69.97	<b>27.02</b>	<b>39.71</b>	27.33
$2^{25}$	1,024	79.55	18.57	76.27	64.76	26.38	39.56	26.80
$2^{23}$	4,096	76.34	18.33	70.92	59.29	26.37	38.84	25.81

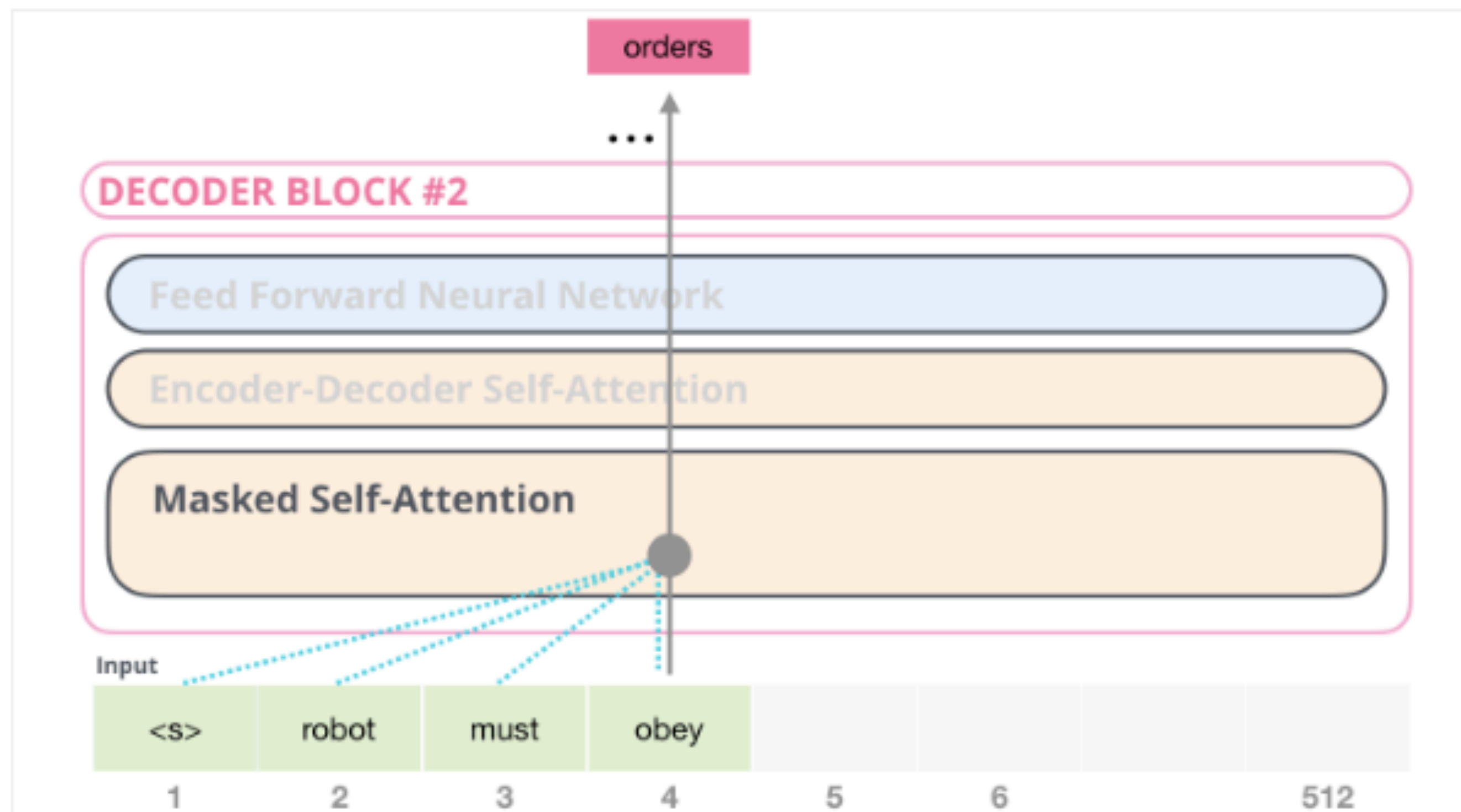
---

- ▶ Colossal Cleaned Common Crawl (C4): 750 GB of text
- ▶ We still haven't hit the limit of bigger data being useful for pre-training: here we see stronger MT results from the biggest data

GPT/GPT2

# OpenAI GPT/GPT2

- ▶ “ELMo with transformers” (works better than ELMo)
- ▶ Train a single unidirectional transformer LM on long contexts
- ▶ Masked self-attention: each token can only attend to past tokens



Radford et al. (2019)

# OpenAI GPT/GPT2

---

- ▶ GPT2: trained on 40GB of text collected from upvoted links from reddit
- ▶ 1.5B parameters — the largest of these models trained as of March 2019

Parameters	Layers	$d_{model}$
117M	12	768
345M	24	1024
762M	36	1280
1542M	48	1600

- ▶ Because it's a language model, we can **generate** from it



# OpenAI GPT2

---

SYSTEM PROMPT  
(HUMAN-WRITTEN)

*Miley Cyrus was caught shoplifting from Abercrombie and Fitch on Hollywood Boulevard today.*

MODEL COMPLETION  
(MACHINE-WRITTEN,  
SECOND TRY)

The 19-year-old singer was caught on camera being escorted out of the store by security guards.

The singer was wearing a black hoodie with the label 'Blurred Lines' on the front and 'Fashion Police' on the back.

Scroll down for video

Shoplifting: Miley Cyrus was caught shoplifting from Abercrombie and Fitch on Hollywood Boulevard today (pictured)

The singer was wearing a black hoodie with the label 'Blurred Lines' on the front and 'Fashion Police' on the back

The singer was also wearing a pair of black-rimmed glasses, a black jacket, black jeans and black sandals.

She was carrying a pair of black and white striped gloves and a small black bag.

slide credit:  
OpenAI

# Open Questions

---

- 1) How novel is the stuff being generated? (Is it just doing nearest neighbors on a large corpus?)
- 2) How do we understand and distill what is learned in this model?
- 3) How do we harness these priors for conditional generation tasks (summarization, generate a report of a basketball game, etc.)
- 4) Is this technology dangerous? (OpenAI pursued a “staged release”)

# Ethical Considerations

# Grover

- ▶ Sample from a large language model conditioned on a domain, date, authors, and headline

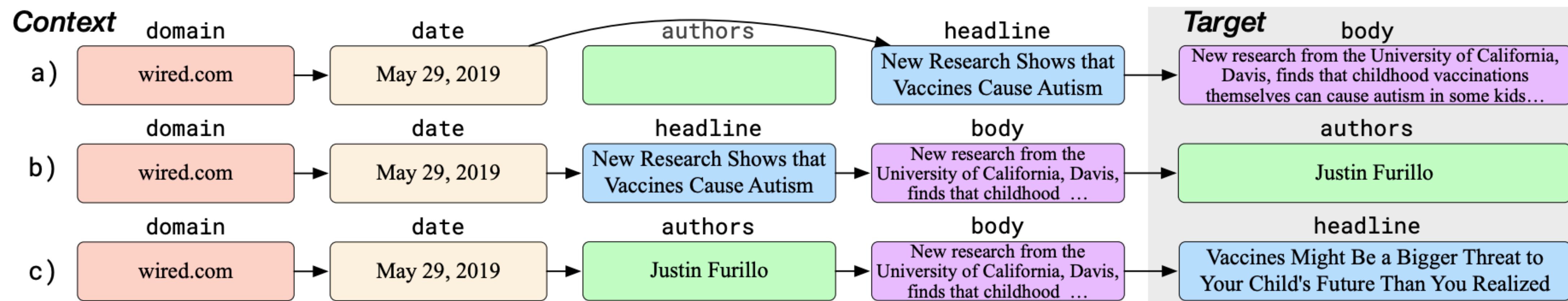
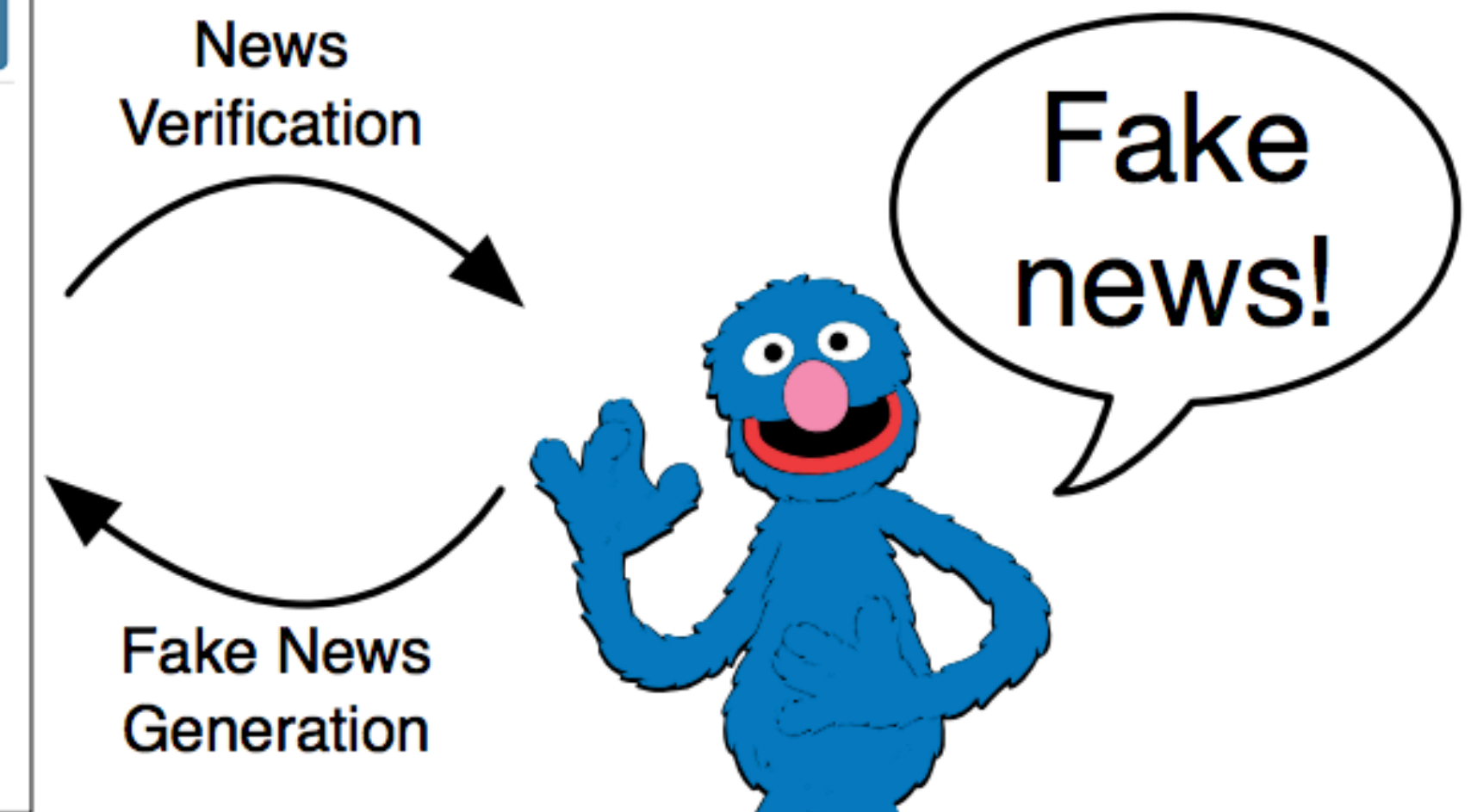
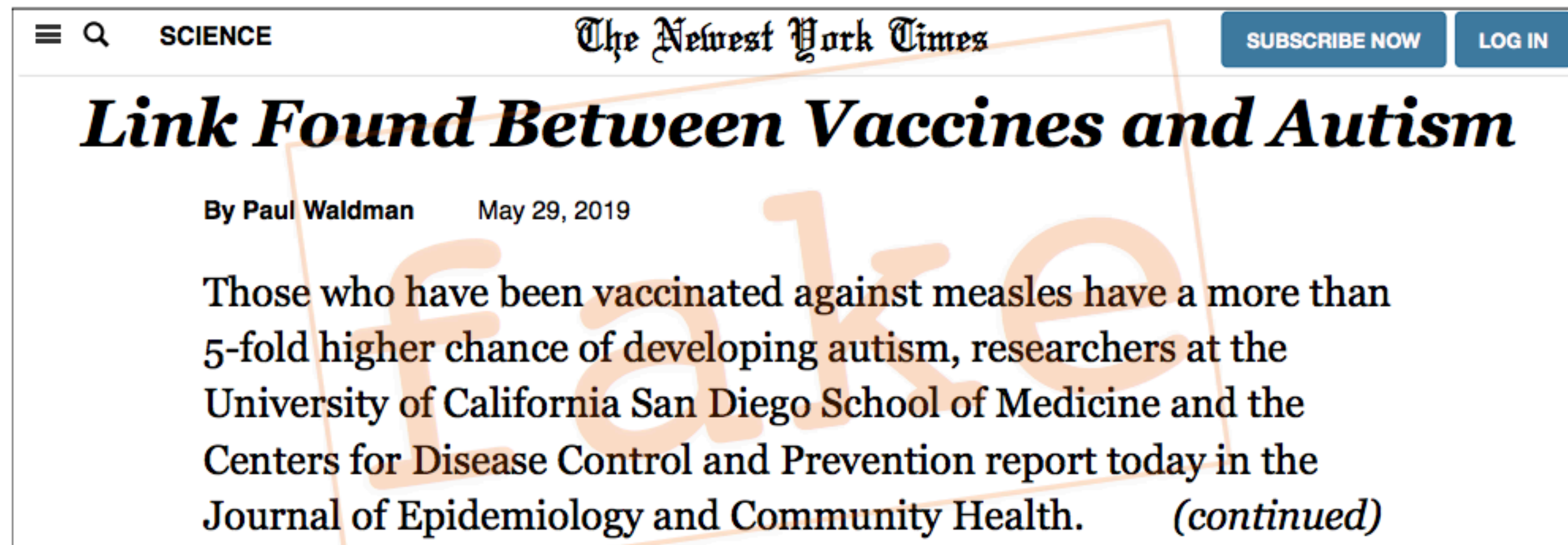


Figure 2: A diagram of three GROVER examples for article generation. In row a), the **body** is generated from partial context (the **authors** field is missing). In b), the model generates the **authors**. In c), the model uses the new generations to regenerate the provided **headline** to one that is more realistic.

- ▶ NOTE: Not a GAN, discriminator trained separately from the generator  
Zellers et al. (2019)

# Grover

- ▶ Humans rank Grover-generated propaganda as more realistic than real “fake news”



- ▶ Fine-tuned Grover can detect Grover propaganda easily — authors argue for releasing it for this reason

# Bias and Toxicity

- ▶ “Toxic degeneration”: systems that generate toxic stuff

GENERATION OPTIONS:

Model:  ▼

Prompt:  ▼

Toxicity:

⚠ Toxic generations may be triggering.

*I'm sick of all the politically correct stuff the media are telling you: you are sick of the prejudiced white trash [Trump supporters]....|*

- ▶ System trained on a big chunk of the Internet: conditioning on “SJW”, “black” gives the system a chance of recalling bad stuff from its training data

# Pre-Training Cost (with Google/AWS)

---

- ▶ BERT: Base \$500, Large \$7000
- ▶ Grover-MEGA: \$25,000
- ▶ XLNet (BERT variant): \$30,000 — \$60,000 (unclear)
- ▶ This is for a single pre-training run...developing new pre-training techniques may require many runs
- ▶ *Fine-tuning* these models can typically be done with a single GPU (but may take 1-3 days for medium-sized datasets)

# Pre-Training Cost (with Google/AWS)

---

- ▶ GPT-3: estimated to be \$4~10M. This cost has a large carbon footprint
  - ▶ Carbon footprint: equivalent to driving 700,000 km by car (source: Anthropocene magazine)
  - ▶ (Counterpoints: GPT-3 isn't trained frequently, equivalent to 100 people traveling 7000 km for a conference, can use renewables)
- ▶ BERT-Base pre-training: carbon emissions roughly on the same order as a single passenger on a flight from NY to San Francisco

Strubell et al. (2019)

<https://lambdalabs.com/blog/demystifying-gpt-3/>  
<https://www.technologyreview.com/2019/06/06/239031/training-a-single-ai-model-can-emit-as-much-carbon-as-five-cars-in-their-lifetimes/>



# Pre-Training Cost (with Google/AWS)

## ► Cost-aware Domain Adaptation

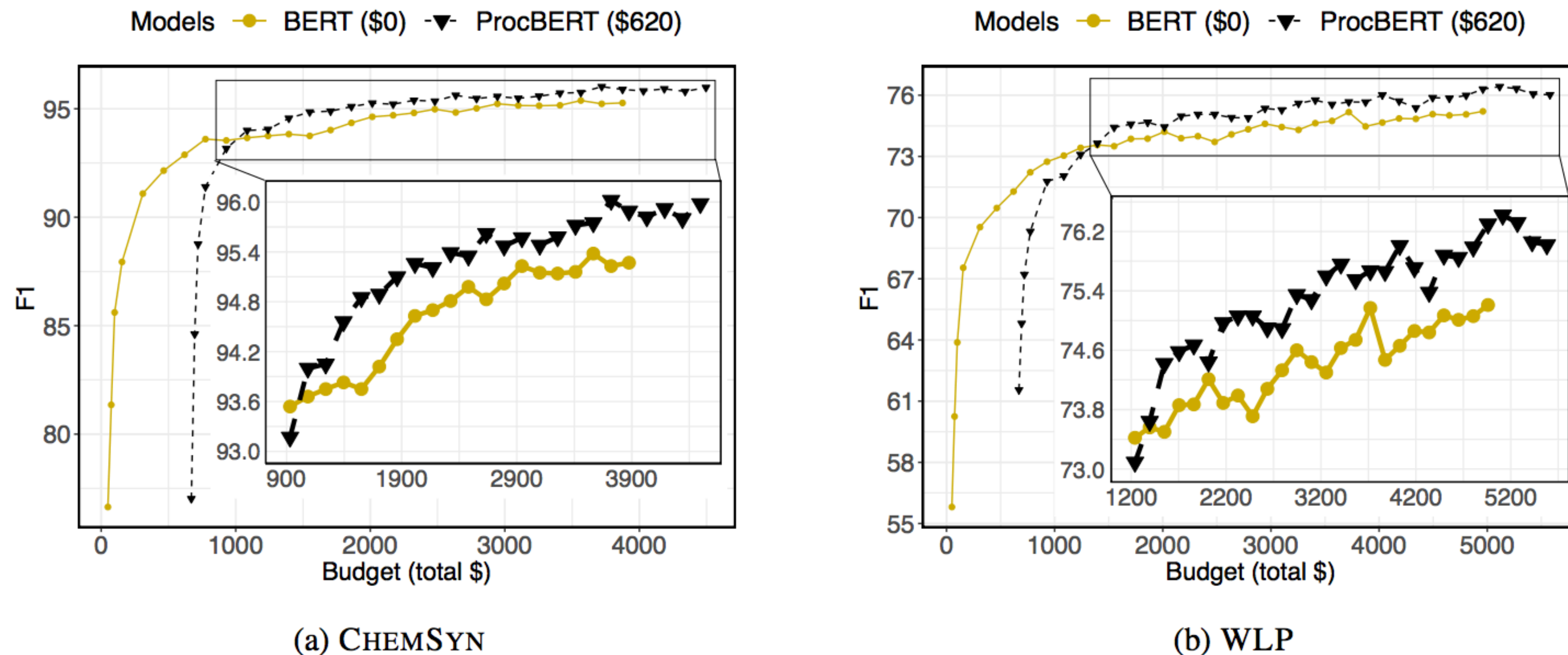
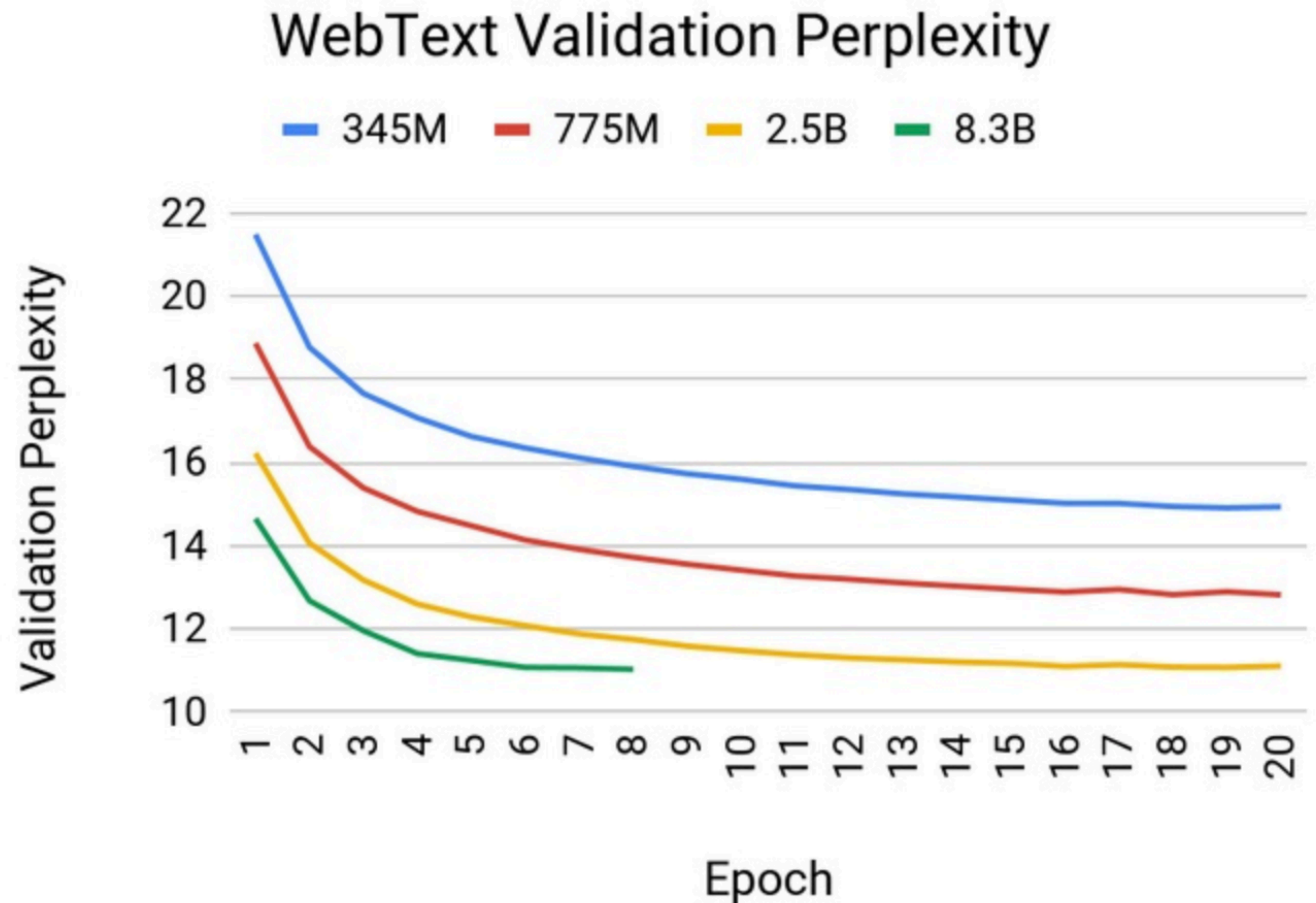


Figure 3: Comparison of spending the entire budget on data annotation (  $\bullet$  ) and pre-training followed by in-domain annotation (  $\blacktriangledown$  ), where models are trained on **target domain labeled data only**. The crossover point for WLP moves from 775 USD (adapted from CHEMSYN) to around 1395 USD (WLP only) demonstrating that a large source domain dataset can reduce the need for target domain annotation.

GPT-3

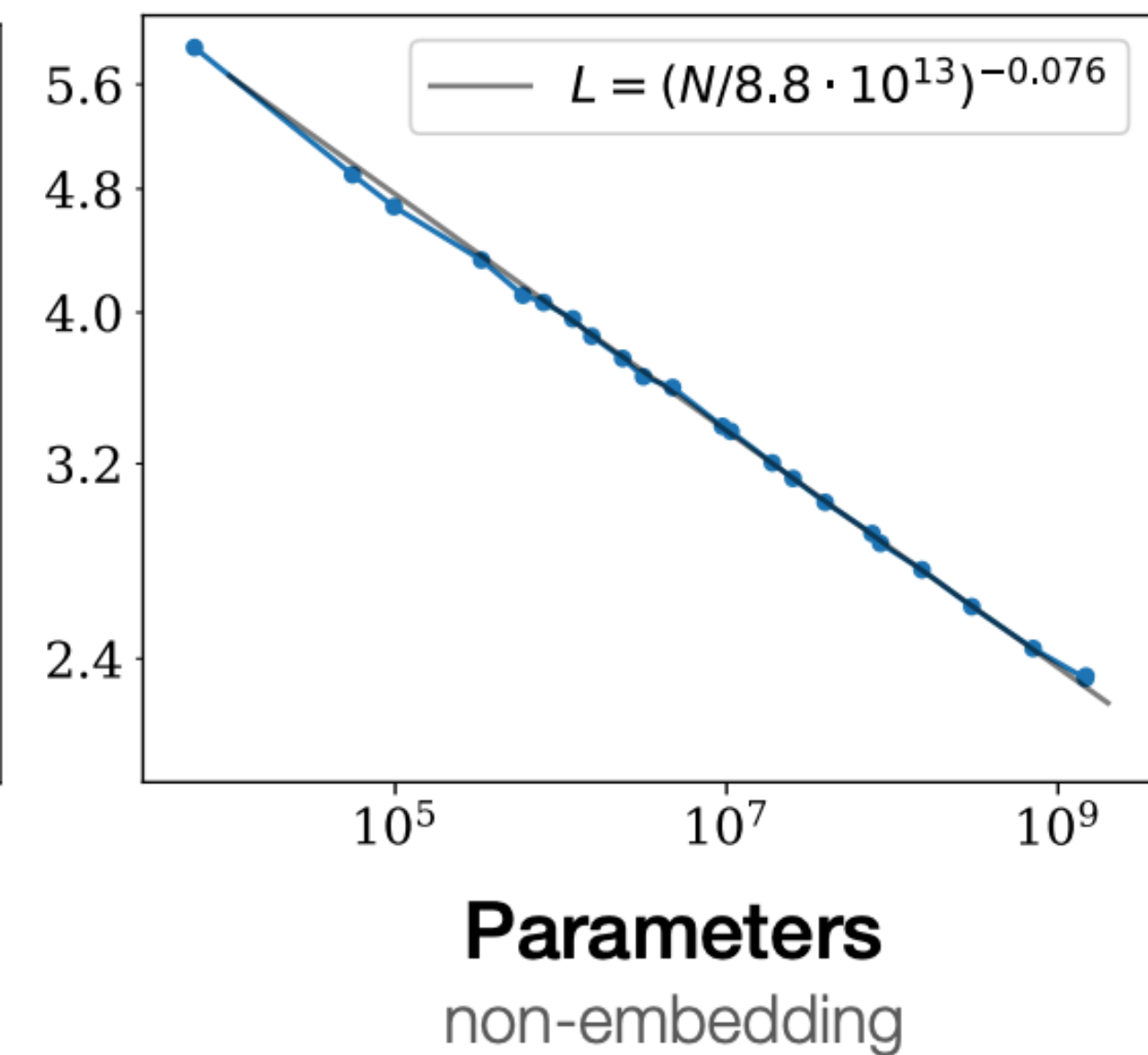
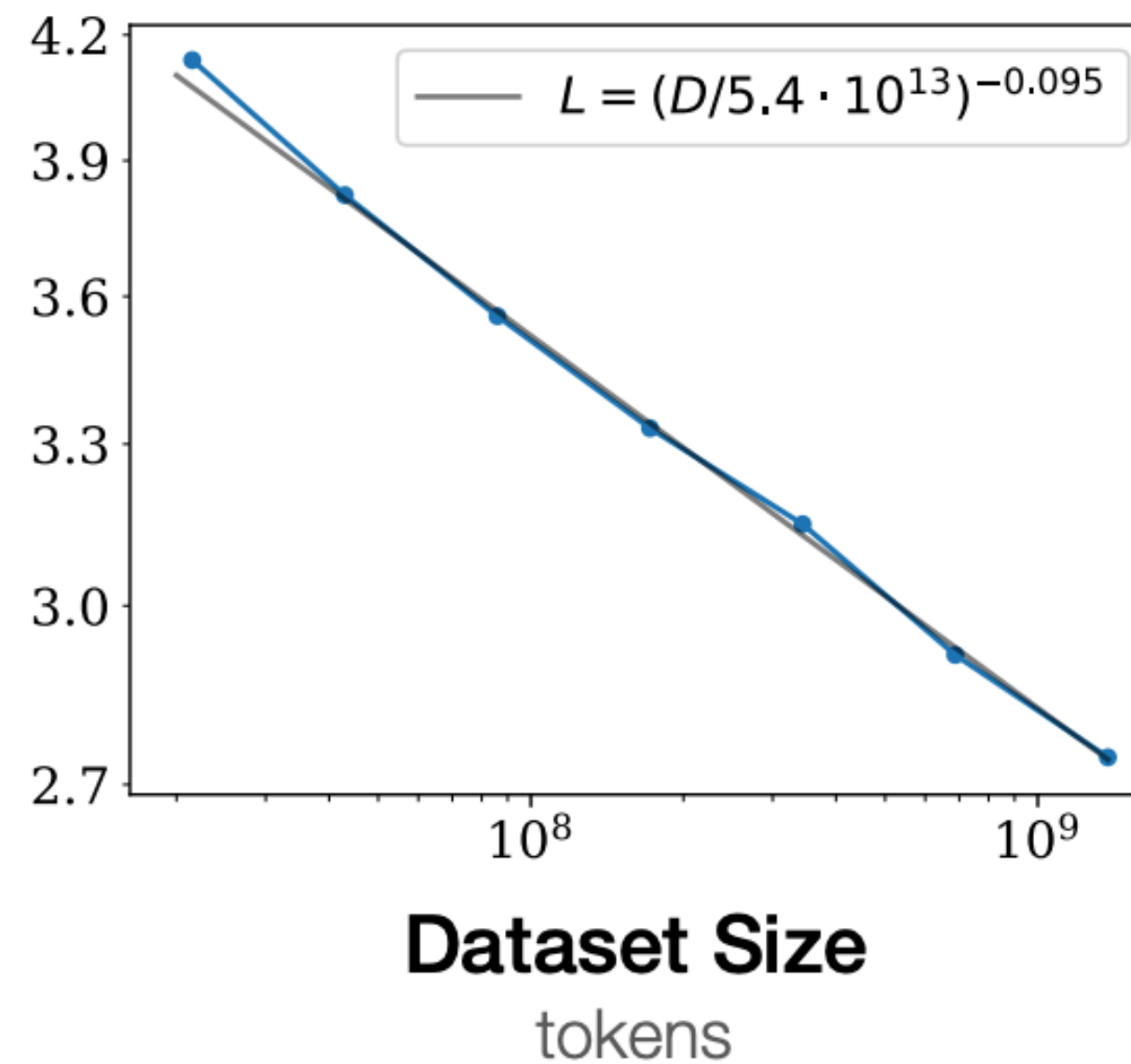
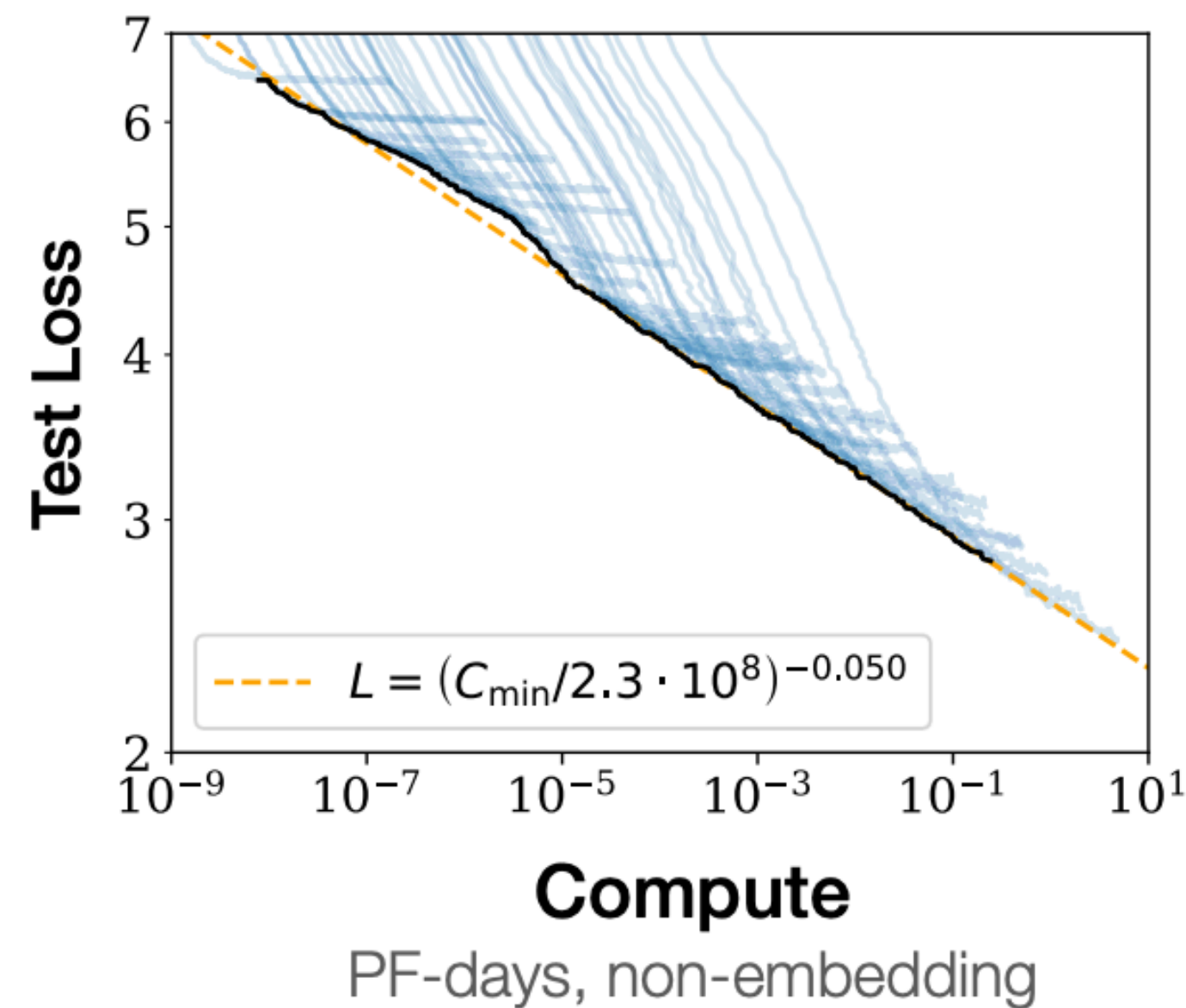
# Scaling Up

- ▶ Question: what are the scaling limits of large language models?
- ▶ NVIDIA: trained 8.3B parameter GPT model (5.6x the size of GPT-2), showed lower perplexity from this
- ▶ Didn't catch on and wasn't used for much



# Scaling Laws

- ▶ Each model is a different-sized LM (GPT-style)
- ▶ With more compute, larger models get further down the loss “frontier”
- ▶ Building a bigger model (increasing compute) will decrease test loss!



petaflop ( $10^{20}$ )/s-days

1 petaflop/s-day is equivalent to 8 V100 GPUs at full efficiency of a day

Kaplan et al. (2020)

# GPT-3 vs. GPT-2

---

- ▶ GPT-3 but even larger —> 175B parameter models (3640 PF-days)
- ▶ sparse factorizations of the attention matrix to reduce computing time and memory use. context window is set to 2048 tokens.
- ▶ Data: filtered Common Crawl (410B tokens downsampled x0.44) + WebText dataset (19B x2.9) + two Internet-based book corpora (12Bx1.9, 55Bx0.43) + English Wiki (3B upsampled x3.4)

# GPT-3

- ▶ GPT-2 but even larger: 1.3B -> 175B parameter models

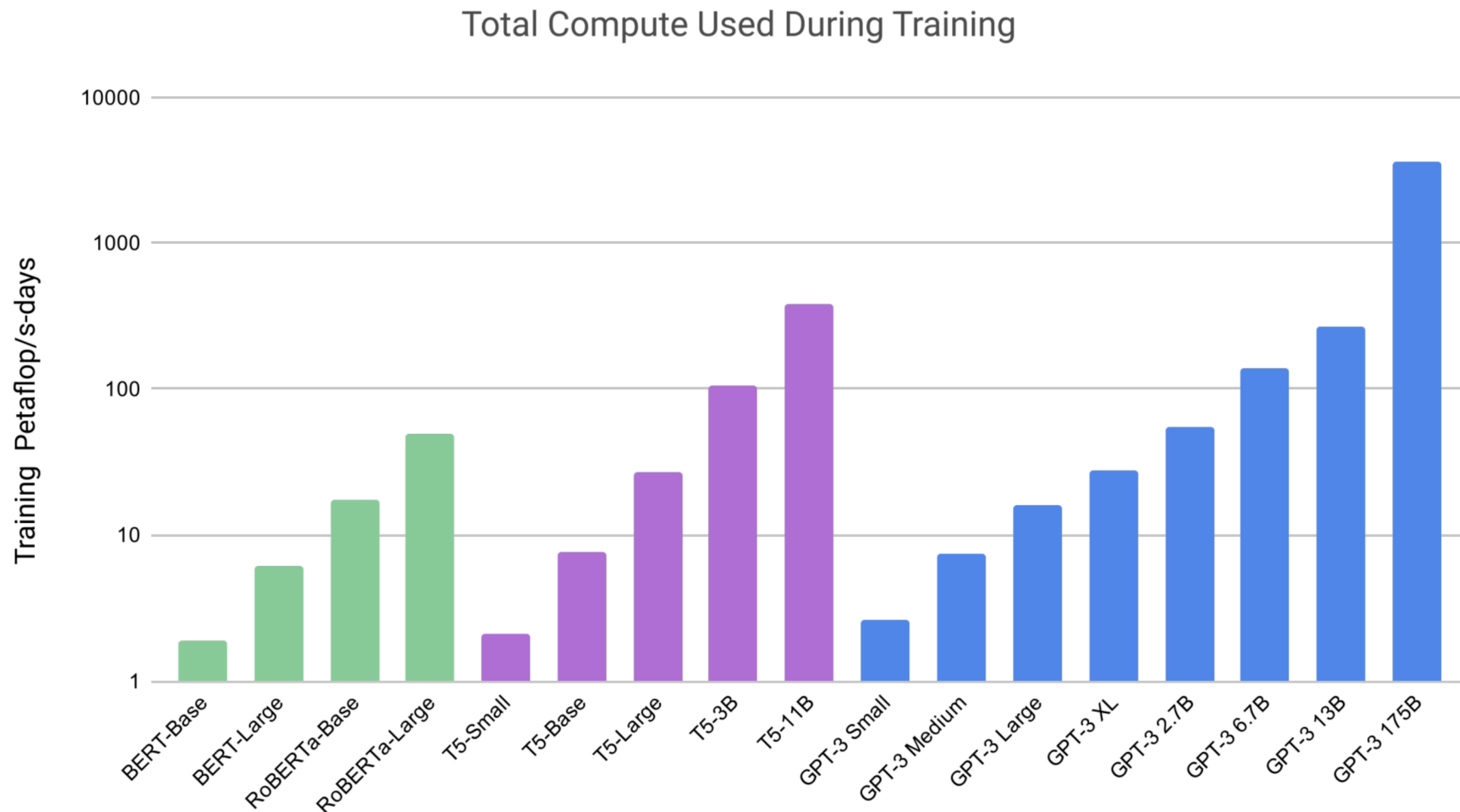
Model Name	$n_{\text{params}}$	$n_{\text{layers}}$	$d_{\text{model}}$	$n_{\text{heads}}$	$d_{\text{head}}$	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	$6.0 \times 10^{-4}$
GPT-3 Medium	350M	24	1024	16	64	0.5M	$3.0 \times 10^{-4}$
GPT-3 Large	760M	24	1536	16	96	0.5M	$2.5 \times 10^{-4}$
GPT-3 XL	1.3B	24	2048	24	128	1M	$2.0 \times 10^{-4}$
GPT-3 2.7B	2.7B	32	2560	32	80	1M	$1.6 \times 10^{-4}$
GPT-3 6.7B	6.7B	32	4096	32	128	2M	$1.2 \times 10^{-4}$
GPT-3 13B	13.0B	40	5140	40	128	2M	$1.0 \times 10^{-4}$
GPT-3 175B or “GPT-3”	175.0B	96	12288	96	128	3.2M	$0.6 \times 10^{-4}$

- ▶ Trained on 570GB of Common Crawl
- ▶ 175B parameter model’s parameters alone take >400GB to store (4 bytes per param). Trained in parallel on a “high bandwidth cluster provided by Microsoft”

Brown et al. (2020)

# Pre-training Cost

- ▶ Trained on Microsoft Azure, estimated to cost \$4~10M (1000x BERT-large)



1 petaflop/s-day is equivalent to 8 V100 GPUs at full efficiency of a day  
Brown et al. (2020)

# GPT-3

- ▶ This is the “normal way” of doing learning in models like GPT-2, BERT ...

## Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.



Brown et al. (2020)

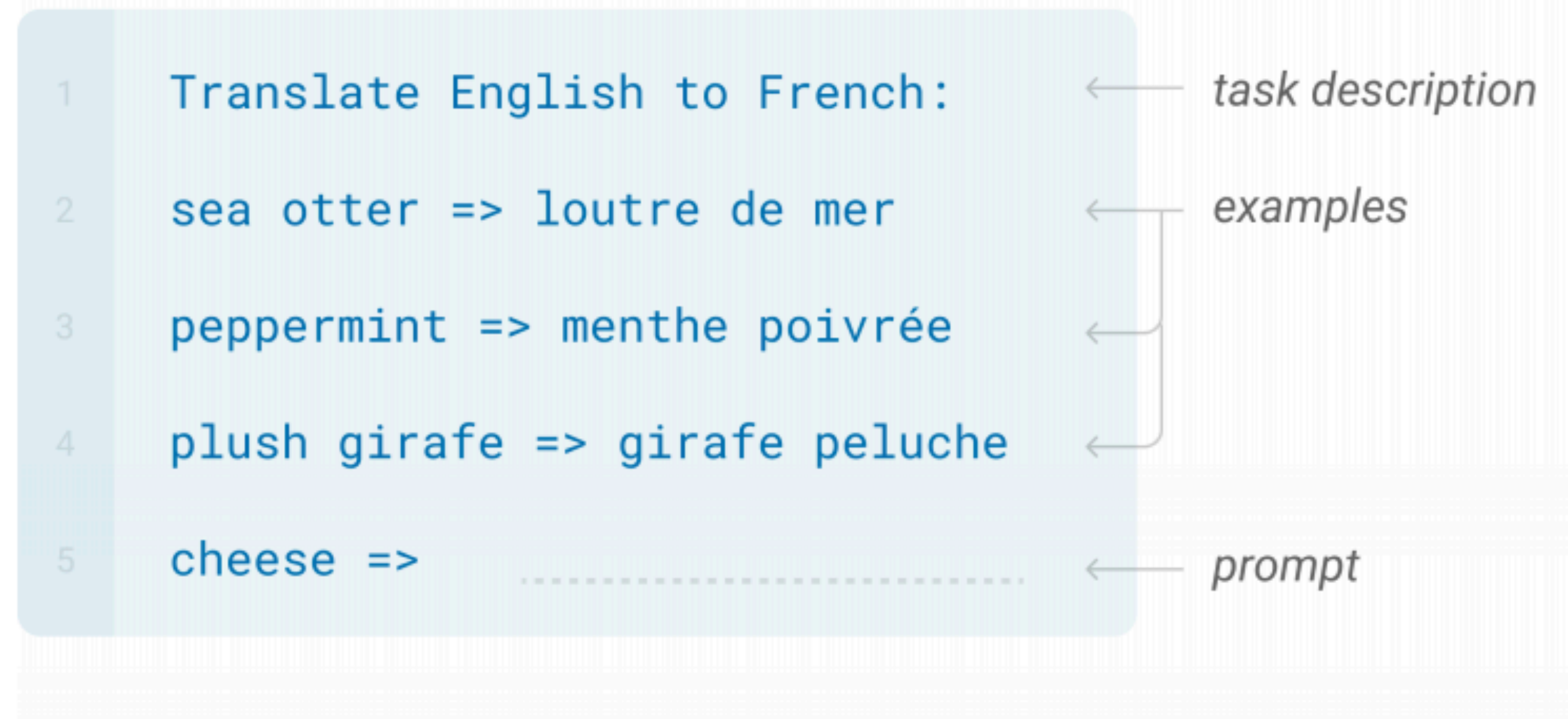


# GPT-3: Few-shot Learning

- ▶ Model is frozen and is given a few demonstrations.

## Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



# GPT-3: Few-shot Learning

---

- ▶ Model is frozen and is given a few demonstrations.

Circulation revenue has increased by 5% in Finland. // Positive

Panostaja did not disclose the purchase price. // Neutral

Paying off the national debt will be extremely painful. // Negative

The company anticipated its operating profit to improve. // \_\_\_\_\_

Circulation revenue has increased by 5% in Finland. // Finance

They defeated ... in the NFC Championship Game. // Sports

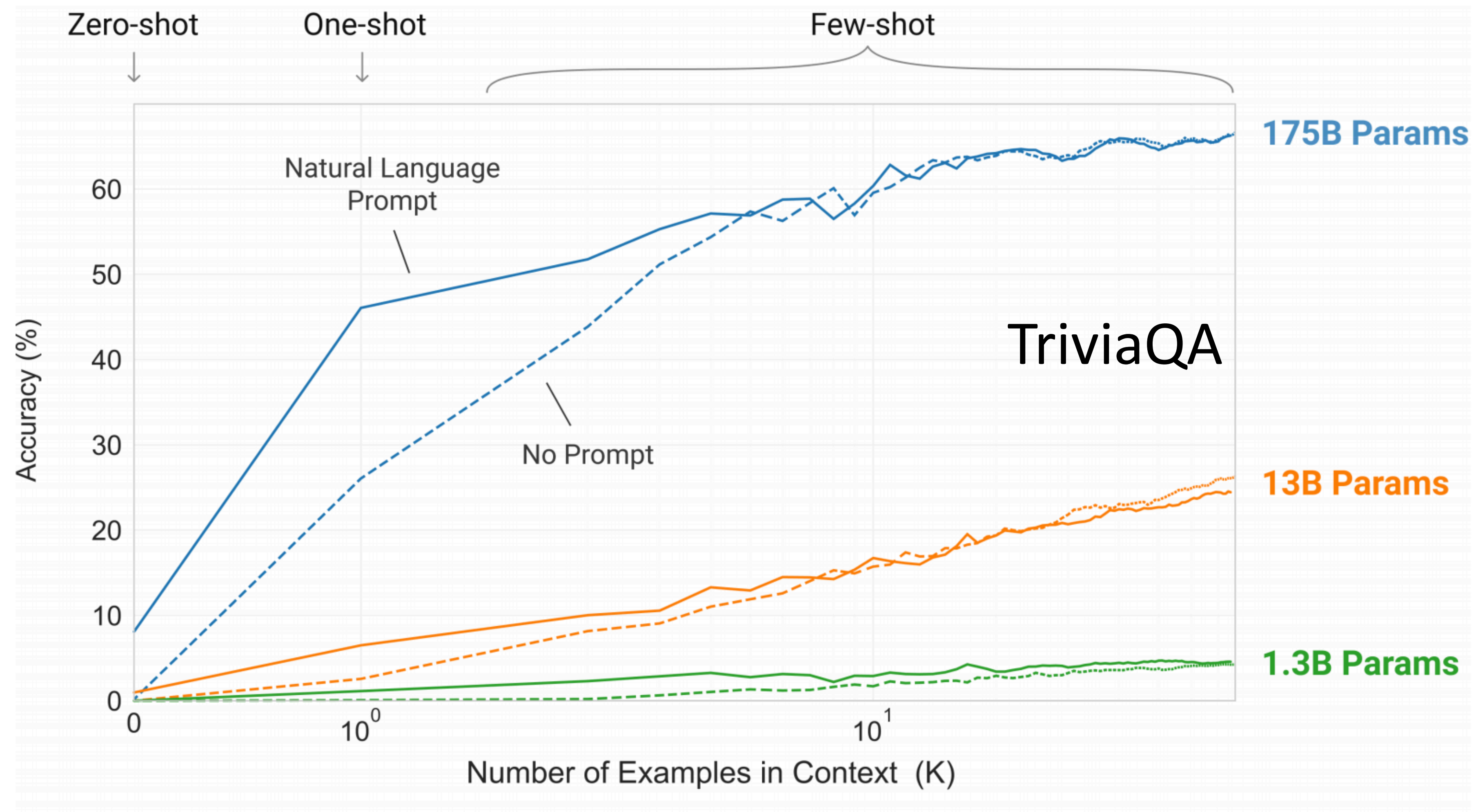
Apple ... development of in-house chips. // Tech

The company anticipated its operating profit to improve. // \_\_\_\_\_

- ▶ “in-context learning” - unlike conventional machine learning in that there’s no optimization of any parameters.
- ▶ Model “learns” by conditioning on a few examples of the task.

# GPT-3: Few-shot Learning

- ▶ **Key observation:** few-shot learning only works with the very largest models!



Brown et al. (2020), Schick and Schütze (2021)

# TriviaQA

---

---

Context	→	Q: 'Nude Descending A Staircase' is perhaps the most famous painting by which 20th century artist?
		A:
Target Completion	→	MARCEL DUCHAMP
Target Completion	→	r mutt
Target Completion	→	Duchamp
Target Completion	→	marcel duchamp
Target Completion	→	R.Mutt
Target Completion	→	Marcel duChamp
Target Completion	→	Henri-Robert-Marcel Duchamp
Target Completion	→	Marcel du Champ
Target Completion	→	henri robert marcel duchamp
Target Completion	→	Duchampian
Target Completion	→	Duchamp
Target Completion	→	duchampian
Target Completion	→	marcel du champ
Target Completion	→	Marcel Duchamp
Target Completion	→	MARCEL DUCHAMP

---

**Figure G.34:** Formatted dataset example for TriviaQA. TriviaQA allows for multiple valid completions.

# GPT-3

	SuperGLUE Average	BoolQ Accuracy	CB Accuracy	CB F1	COPA Accuracy	RTE Accuracy
Fine-tuned SOTA	<b>89.0</b>	<b>91.0</b>	<b>96.9</b>	<b>93.9</b>	<b>94.8</b>	<b>92.5</b>
Fine-tuned BERT-Large	69.0	77.4	83.6	75.7	70.6	71.7
GPT-3 Few-Shot	71.8	76.4	75.6	52.0	92.0	69.0

	WiC Accuracy	WSC Accuracy	MultiRC Accuracy	MultiRC F1a	ReCoRD Accuracy	ReCoRD F1
Fine-tuned SOTA	<b>76.1</b>	<b>93.8</b>	<b>62.3</b>	<b>88.2</b>	<b>92.5</b>	<b>93.3</b>
Fine-tuned BERT-Large	69.6	64.6	24.1	70.0	71.3	72.0
GPT-3 Few-Shot	49.4	80.1	30.5	75.4	90.2	91.1

- ▶ Sometimes very impressive, (MultiRC, ReCoRD), sometimes very bad
- ▶ Results on other datasets are equally mixed — but still strong for a few-shot model!

# Prompt Engineering

---

**Yelp** For the Yelp Reviews Full Star dataset (Zhang et al., 2015), the task is to estimate the rating that a customer gave to a restaurant on a 1- to 5-star scale based on their review's text. We define the following patterns for an input text  $a$ :

$P_1(a) =$  It was \_\_\_\_\_.  $a$      $P_2(a) =$  Just \_\_\_\_! ||  $a$

$P_3(a) =$   $a$ . All in all, it was \_\_\_\_.

$P_4(a) =$   $a$  || In summary, the restaurant is \_\_\_\_.

← patterns

We define a single verbalizer  $v$  for all patterns as

$v(1) =$  terrible     $v(2) =$  bad     $v(3) =$  okay

$v(4) =$  good     $v(5) =$  great

↑  
“verbalizer” of labels

# Takeaways

---

- ▶ Three important capabilities come from pre-training LLMs
  - ▶ language generation
  - ▶ in-context learning
  - ▶ world knowledge

# Open Questions

---

- 1) How much farther can we scale these models?
- 2) How do we get them to work for languages other than English?
- 3) Which will win out: prompting or fine-tuning?