

CS371N: Natural Language Processing

Lecture 13: Subword Tokenization, Decoders, Decoding

Greg Durrett



TEXAS

The University of Texas at Austin



Announcements

- ▶ A3 due Thursday
- ▶ A2 back soon



Recap: BERT



Today

- ▶ Subword tokenization
- ▶ GPT-2/GPT-3: **decoders**, which are able to actually generate text
- ▶ **Decoding** methods for getting outputs from these models
- ▶ Prompting: a new way of using large language models without taking any gradient steps

Subword Tokenization



Handling Rare Words

- ▶ Words are a difficult unit to work with. Why? What becomes harder about LMs (e.g., the assignment) when you have 100,000+ words?
- ▶ Some answers: the final matrix multiply and softmax start to dominate the computation, many params, still some words you haven't seen, doesn't take advantage of morphology, ...
- ▶ Character-level models were explored extensively in 2016-2018 but simply don't work well — becomes very expensive to represent sequences



Subword Tokenization

- ▶ Subword tokenization: wide range of schemes that use tokens that are **between characters and words** in terms of granularity
- ▶ These “word pieces” may be full words or parts of words
_the _**eco tax** _port i co _in _Po nt - de - Bu is ...
- ▶ _ indicates the word piece starting a word (can think of it as the space character).



Subword Tokenization

- ▶ Subword tokenization: wide range of schemes that use tokens that are **between characters and words** in terms of granularity
- ▶ These “word pieces” may be full words or parts of words

_the **_eco tax** _port i co _in _Po nt - de - Bu is...

Output: _le _port ique **_éco taxe** _de _Pont - de - Bui s

- ▶ Can achieve transliteration with this, subword structure makes some translations easier to achieve



Byte Pair Encoding (BPE)

- ▶ Start with every individual byte (basically character) as its own symbol

```
for i in range(num_merges):  
    pairs = get_stats(vocab)  
    best = max(pairs, key=pairs.get)  
    vocab = merge_vocab(best, vocab)
```

- ▶ Count bigram character cooccurrences
- ▶ Merge the most frequent pair of adjacent characters
- ▶ Doing 8k merges => vocabulary of around 8000 word pieces. Includes many whole words
- ▶ Most SOTA NMT systems use this on both source + target



Byte Pair Encoding (BPE)

Original:	furiously		Original:	tricycles
BPE:	_fur iously	(b)	BPE:	_t ric y cles
Unigram LM:	_fur ious ly		Unigram LM:	_tri cycle s
Original:	Completely preposterous suggestions			
BPE:	_Comple t ely	_prep ost erous	_suggest ions	
Unigram LM:	_Complete ly	_pre post er ous	_suggestion s	

- ▶ What do you see here?
- ▶ BPE produces less linguistically plausible units than another technique based on a unigram language model: rather than greedily merge, find chunks which make the sequence look likely under a unigram LM
- ▶ Unigram LM tokenizer leads to slightly better BERT

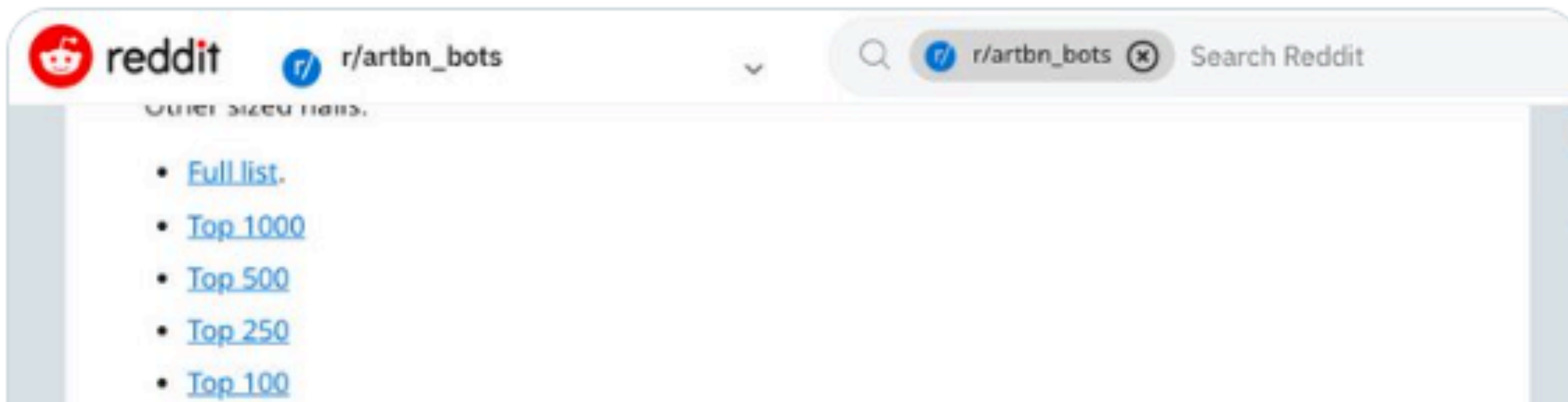


What's in the token vocab?



Matthew Watkins
@SoC_trilogy

I've just found out that several of the anomalous GPT tokens ("TheNitromeFan", " SolidGoldMagikarp", " davidjl", " Smartstocks", " RandomRedditorWithNo",) are handles of people who are (competitively? collaboratively?) counting to infinity on a Reddit forum. I kid you not.



Rank	User	Counts
1	/u/davidjl123	163477
2	/u/Smartstocks	113829
3	/u/atomicimploder	103178
4	/u/TheNitromeFan	84581
5	/u/SolidGoldMagikarp	65753
6	/u/RandomRedditorWithNo	63434
7	/u/rideride	59024
8	/u/Mooraeell	57785
9	/u/Removedpixel	55080
10	/u/Adinida	48415
11	/u/rschaosid	47132



Tokenization Today

- ▶ **All pre-trained** models use some kind of subword tokenization with a tuned vocabulary; usually between 50k and 250k pieces (larger number of pieces for multilingual models)
- ▶ As a result, classical word embeddings like GloVe **are not used**. All subword representations are randomly initialized and learned in the Transformer models

GPT



OpenAI GPT/GPT2

- ▶ Very large language models using the Transformer architecture
- ▶ Straightforward left-to-right language model, trained on raw text
- ▶ GPT2 (March 2019): trained on 40GB of text

	Parameters	Layers	d_{model}
	117M	12	768
approximate size of BERT	345M	24	1024
	762M	36	1280
GPT-2	1542M	48	1600

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

Radford et al. (2019)



OpenAI GPT2

slide credit: OpenAI

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)

- ▶ How was this generated? We'll come back to this in a few minutes



Pre-Training Cost (with Google/AWS)

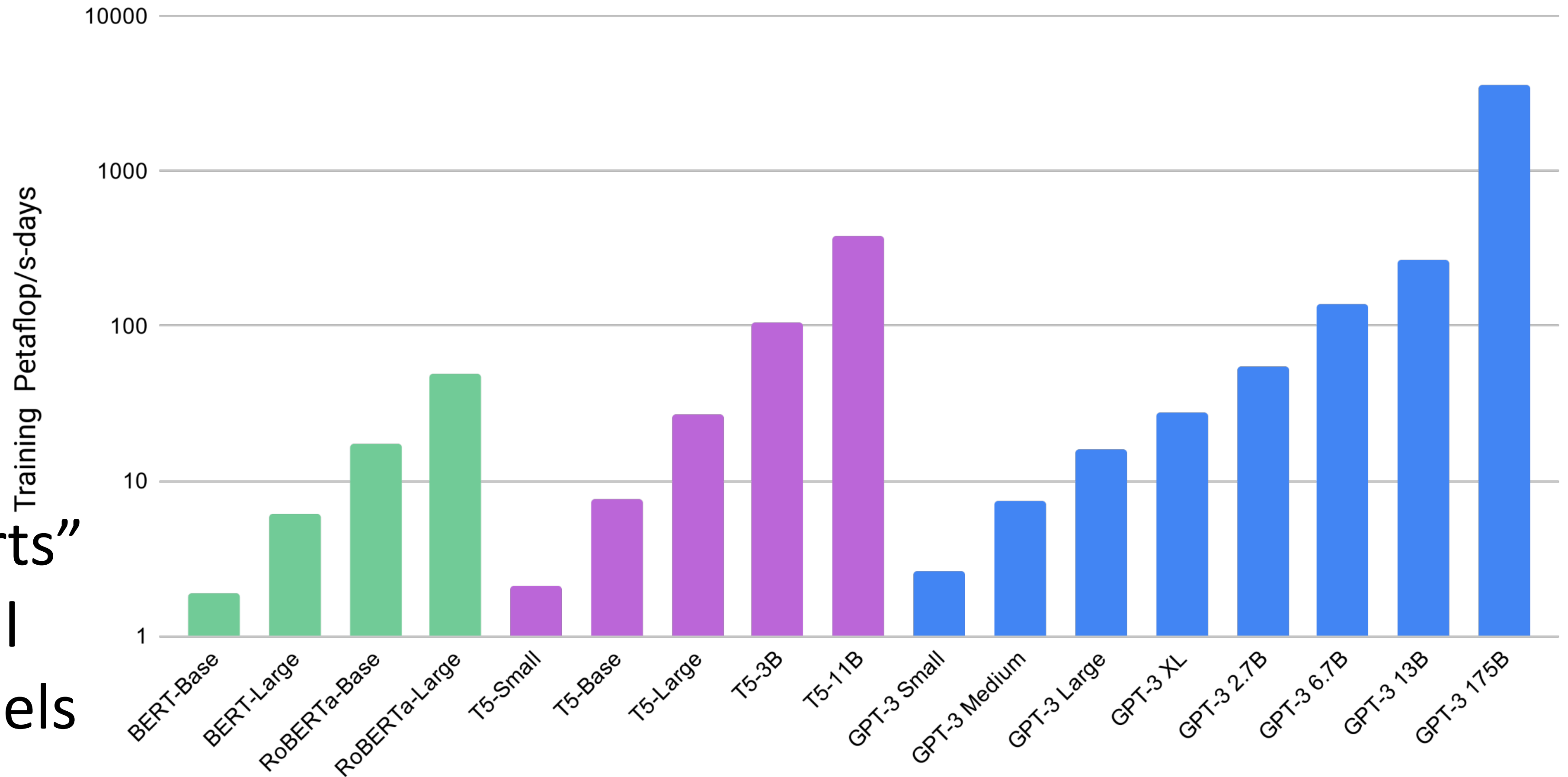
- ▶ BERT: Base \$500, Large \$7000
- ▶ GPT-2 (as reported in other work): \$25,000
- ▶ This is for a single pre-training run...developing new pre-training techniques may require many runs
- ▶ *Fine-tuning* these models can typically be done with a single GPU (but may take 1-3 days for medium-sized datasets)



Scaling Up: GPT-3

- ▶ 175B parameter model: 96 layers, 96 heads, 12k-dim vectors
- ▶ Trained on Microsoft Azure, estimated to cost roughly \$10M
- ▶ GPT-4 may be “mixture of experts” combining several similar-sized models

Total Compute Used During Training



Brown et al. (2020)

Decoding Methods



Decoding Strategies

- ▶ LMs place a distribution $P(y_i | y_1, \dots, y_{i-1})$
- ▶ seq2seq models place a distribution $P(y_i | \mathbf{x}, y_1, \dots, y_{i-1})$
- ▶ Generation from both models looks similar; how do we do it?
 - ▶ Option 1: $\max y_i P(y_i | y_1, \dots, y_{i-1})$ — take greedily best option
 - ▶ Option 2: use beam search to find the sequence with the highest prob.



Beam Search



Decoding Strategies

- ▶ LMs place a distribution $P(y_i | y_1, \dots, y_{i-1})$
- ▶ seq2seq models place a distribution $P(y_i | \mathbf{x}, y_1, \dots, y_{i-1})$
- ▶ Generation from both models looks similar; how do we do it?
 - ▶ Option 1: $\max y_i P(y_i | y_1, \dots, y_{i-1})$ — take greedily best option
 - ▶ Option 2: use beam search to find the sequence with the highest prob.
 - ▶ Option 3: sample from the model; draw y_i from that distribution
- ▶ When should we use these different approaches?



Drawbacks of Sampling

- ▶ Sampling is “too random”

Pure Sampling:

They were cattle called Bolivian Cavalleros; they live in a remote desert uninterrupted by town and they speak huge, beautiful, paradisiacal Bolivian linguistic thing. They say, 'Lunch, marge.' They don't tell what the lunch is," director Professor Chuperas Omwell told Sky News. "They've only been talking to scientists, like we're being interviewed by TV

$P(y \mid \dots \text{they live in a remote desert uninterrupted by})$

0.01 roads

0.01 towns

0.01 people

0.005 civilization

...

0.0005 town

Good options, maybe accounting for 90% of the total probability mass. So a 90% chance of getting something good

Long tail with 10% of the mass

Holtzman et al. (2019)



Nucleus Sampling

$P(y \mid \dots \text{they live in a remote desert uninterrupted by})$

0.01 roads

0.01 towns

0.01 people

0.005 civilization

→ renormalize and sample

— cut off after $p\%$ of mass

- ▶ Define a threshold p . Keep the most probable options account for $p\%$ of the probability mass (the *nucleus*), then sample among these.
- ▶ To implement: sort options by probability, truncate the list once the total exceeds p , then renormalize and sample from it

Holtzman et al. (2019)



Decoding Strategies

- ▶ LMs place a distribution $P(y_i | y_1, \dots, y_{i-1})$
- ▶ seq2seq models place a distribution $P(y_i | \mathbf{x}, y_1, \dots, y_{i-1})$
- ▶ How to generate sequences?
 - ▶ Option 1: $\max y_i P(y_i | y_1, \dots, y_{i-1})$ — take greedily best option
 - ▶ Option 2: use beam search to find the sequence with the highest prob.
 - ▶ ~~Option 3: sample from the model; draw y_i from that distribution~~
 - ▶ Option 4: nucleus sampling



GPT-3

Story completion demo:
Different decoding strategies

Preview: Prompting, In-Context Learning



Pre-GPT-3: Fine-tuning

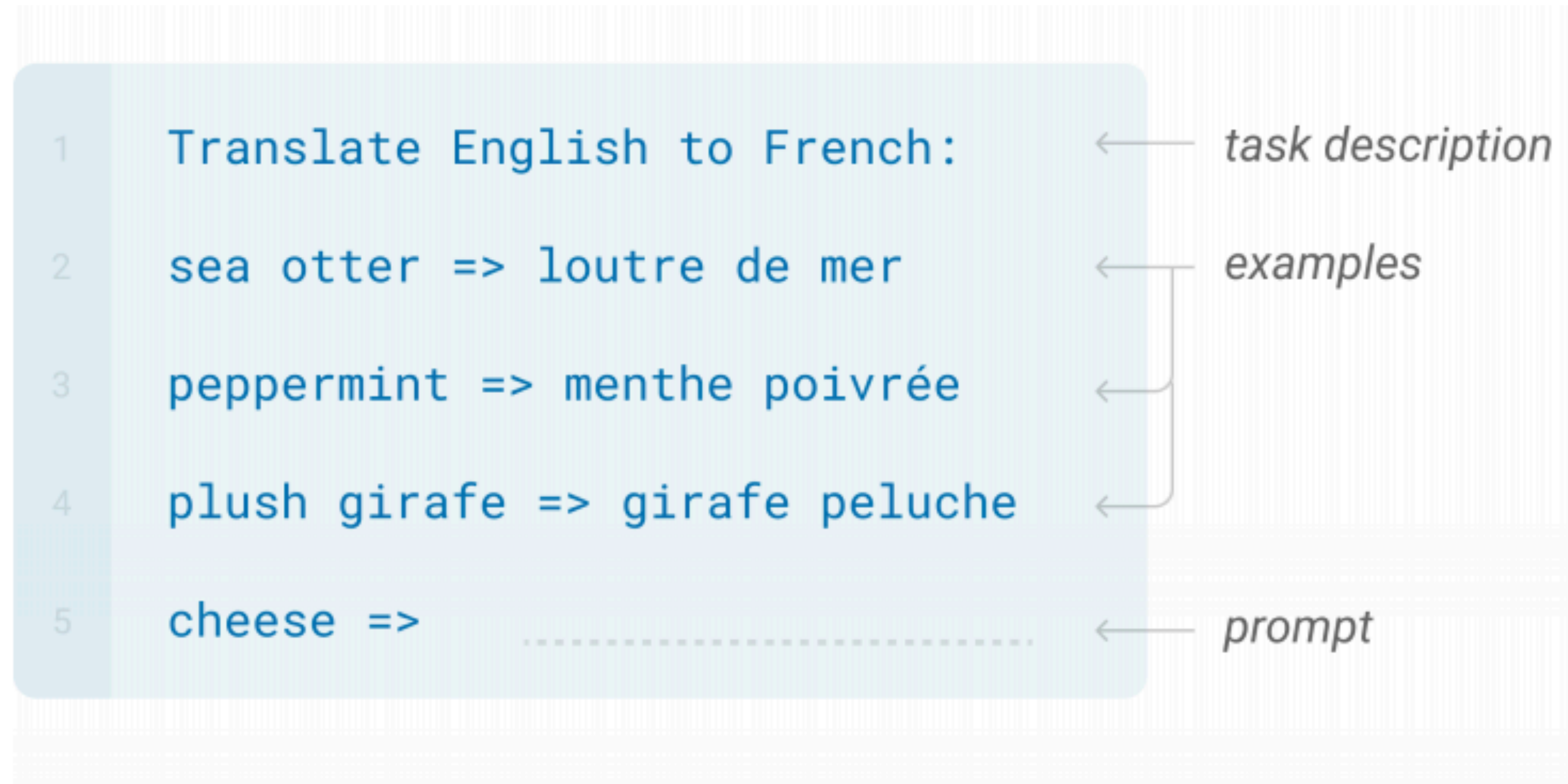
- ▶ Fine-tuning: this is the “normal way” of doing learning in models like GPT-2
- ▶ Requires computing the gradient and applying a parameter update on every example
- ▶ **This is super expensive with 175B parameters**





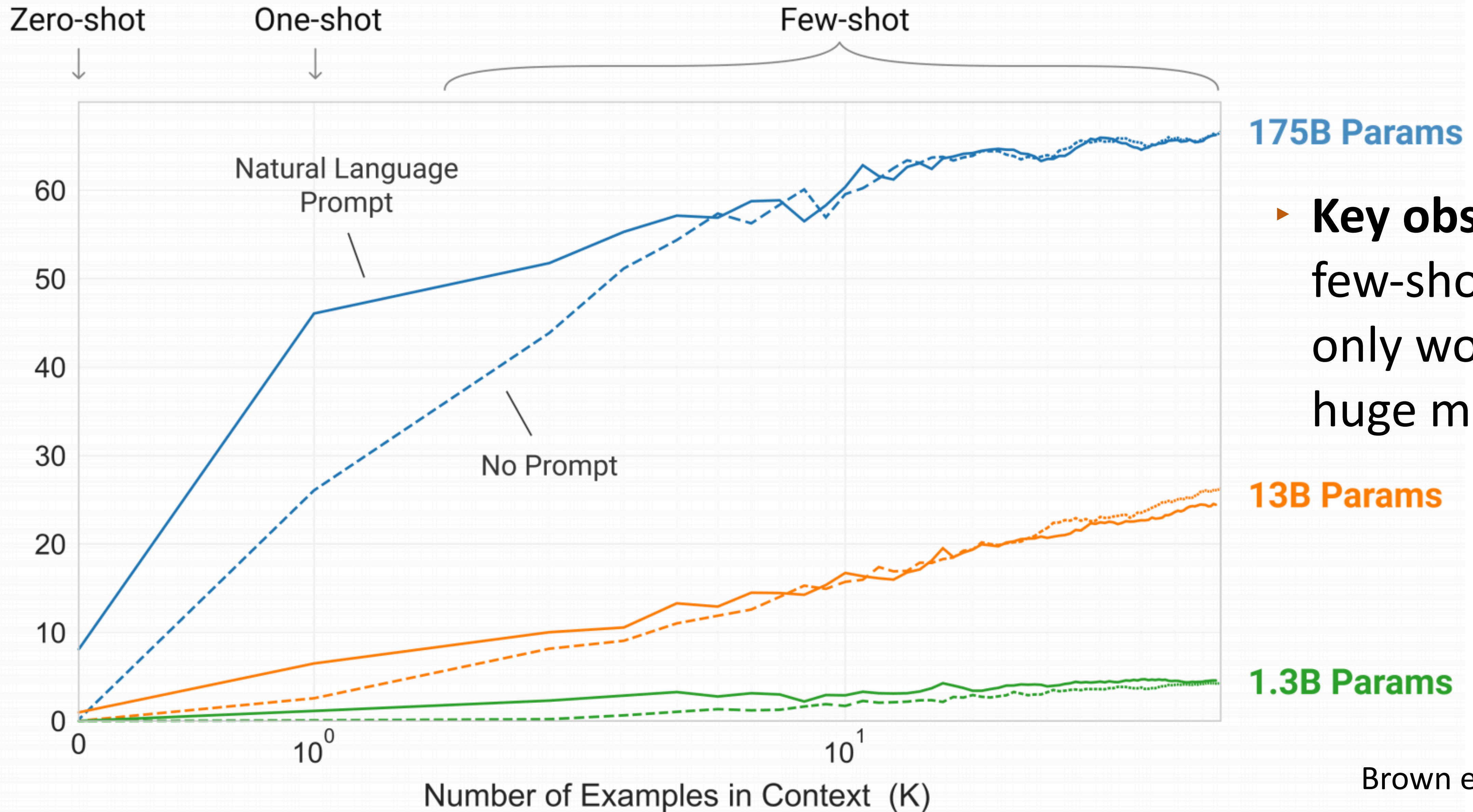
GPT-3: Few-shot Learning

- ▶ GPT-3 proposes an alternative: **in-context learning**. Just uses the off-the-shelf model, no gradient updates
- ▶ This procedure depends heavily on the examples you pick as well as the prompt (“*Translate English to French*”)





GPT-3



175B Params

► **Key observation:**
few-shot learning
only works with
huge models!

13B Params

1.3B Params



GPT-3

	SuperGLUE Average	BoolQ Accuracy	CB Accuracy	CB F1	COPA Accuracy	RTE Accuracy
Fine-tuned SOTA	89.0	91.0	96.9	93.9	94.8	92.5
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	76.1	93.8	62.3	88.2	92.5	93.3
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!



Prompts

- ▶ Prompts can help induce the model to engage in certain behavior
- ▶ In the GPT-2 paper, “tl;dr:” (too long; didn't read) is mentioned as a prompt that frequently shows up in the wild **indicating a summary**
- ▶ tl;dr is an indicator that the model should “switch into summary mode” now — and if there are enough clean instances of tl;dr in the wild, maybe the model has been trained on a ton of diverse data?
- ▶ Good prompt + a few training examples in-context = strong task performance?



Prompting

- ▶ Current training: GPT-3 trained on the web
- ▶ Current testing: feed in a very specific prompt and/or a set of in-context examples
- ▶ Two goals:
 1. Unify pre-training and testing phases
 2. Exploit data for downstream tasks — why are we trying to do question answering while ignoring all of the existing QA datasets?
- ▶ RLHF/instruction tuning: fine-tune on supervised tasks after pre-training (more after midterm)
- ▶ **Let's see how an instruction-tuned GPT-3 works**



Prompts

Prompting demo:
QA, Math QA, etc.

Ethical Issues



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



Stochastic Parrots (about LMs generally)

- ▶ Paper (that included authors at Google who were subsequently fired) about dangers of large language models
- ▶ **Claim 1:** environmental cost is disproportionately born by marginalized populations, who aren't even well-served by these tools
- ▶ **Claim 2:** massive data is fundamentally challenging to audit, contains data that is biased and is only a snapshot of a single point in time
- ▶ **Claim 3 (what we'll focus on today):** these models are not grounded in meaning — when they generate an answer to a question, it is merely by memorizing cooccurrence between symbols



Stochastic Parrots

Question: What is the name of the Russian mercenary group?

Answer: Wagner group.

Question: Where is the Wagner group?

Answer: In Syria.

Question: Who is the leader of the Wagner group?

Answer: A Russian citizen, Dmitry Utkin. He was born in 1971 and lives in Moscow. He has been a member of the Russian parliament since 2016. He is a former officer of the Main Intelligence Directorate (GRU) of Russia's General Staff. He was also a commander of the special forces unit "Vostok" (East). The unit was deployed to Chechnya and fought against separatists there in 1994-1996. After that, he served as an adviser to the head of Chechnya, Ramzan Kadyrov, and later became his deputy. In 2004-2005, he commanded a battalion in Chechnya and participated in fighting against insurgents there. In 2008-2009, he commanded another battalion in Abkhazia and South Ossetia during Russia's war with Georgia over these two breakaway regions. Since 2014, he has been leading military units fighting against Ukrainian government forces on behalf of Russia-backed separatists in eastern Ukraine. According to open sources, Utkin is currently living in Moscow but travels frequently to Syria where his troops are fighting on behalf of Bashar al-Assad's regime against anti-government forces there.

- ▶ We are likely to assume the model is producing factual information and presenting it in a coherent way, but this is our interpretation we project on the model
- ▶ Risks: medical diagnosis ("What do I have if I have X, Y, and Z symptoms?") could seem possible but cause serious harm



Takeaways

- ▶ Generative language models can do well at lots of generation tasks
- ▶ Decoding strategy can matter a lot (beam search vs. sampling)
- ▶ Prompting is a way to harness their power and learn to do many tasks with a single model. Can be done without fine-tuning