

## Language Models are Few-Shot Learners

Tom B. Brown\*

Benjamin Mann\*

Nick Ryder\*

Melanie Subbiah\*

Jared Kaplan<sup>†</sup>

Prafulla Dhariwal

Arvind Neelakantan

Pranav Shyam

Girish Sastry

Amanda Askell

Sandhini Agarwal

Ariel Herbert-Voss

Gretchen Krueger

Tom Henighan

Rewon Child

Aditya Ramesh

Daniel M. Ziegler

Jeffrey Wu

Clemens Winter

Christopher Hesse

Mark Chen

Eric Sigler

Mateusz Litwin

Scott Gray

Benjamin Chess

Jack Clark

Christopher Berner

Sam McCandlish

Alec Radford

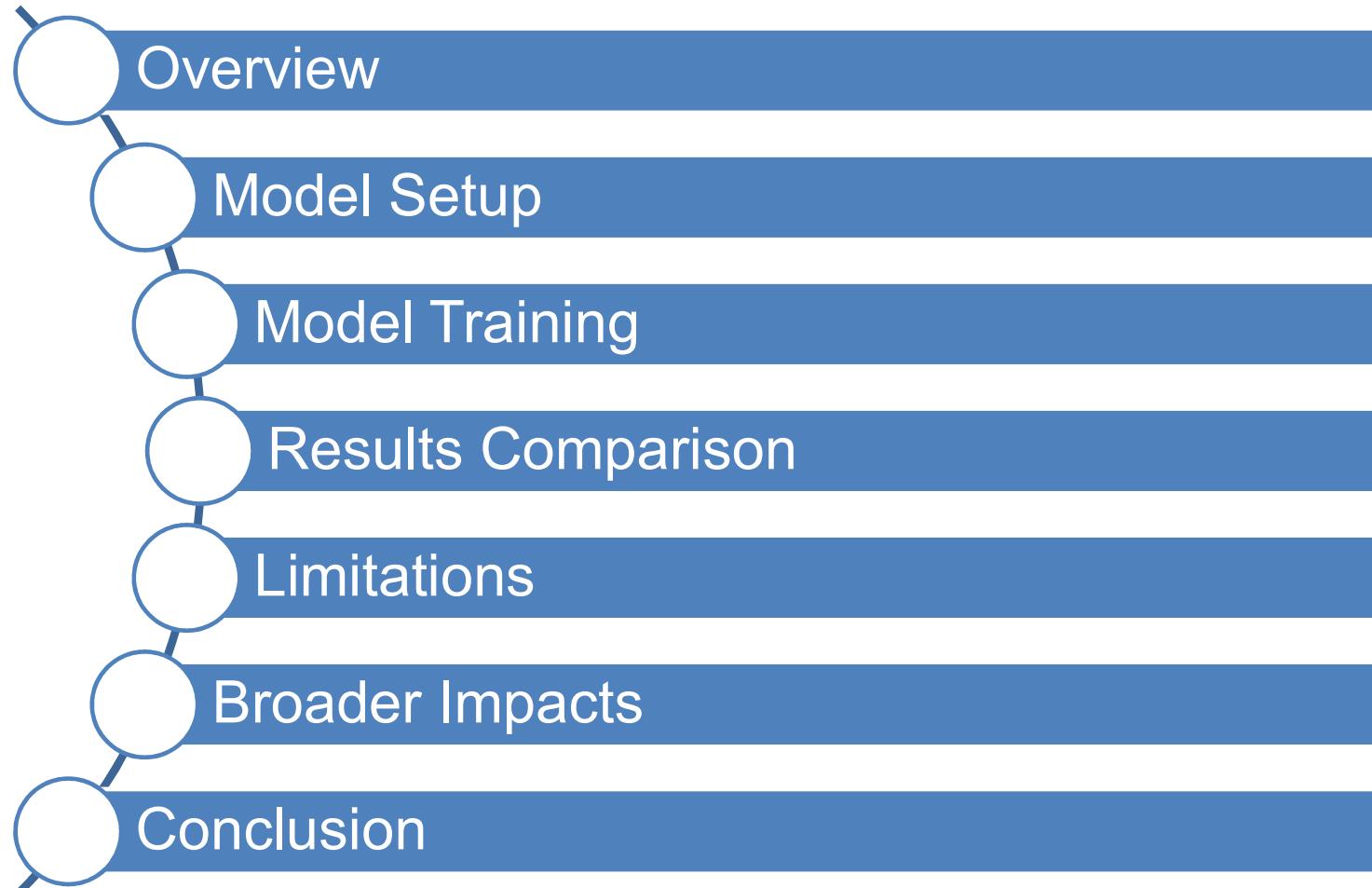
Ilya Sutskever

Dario Amodei

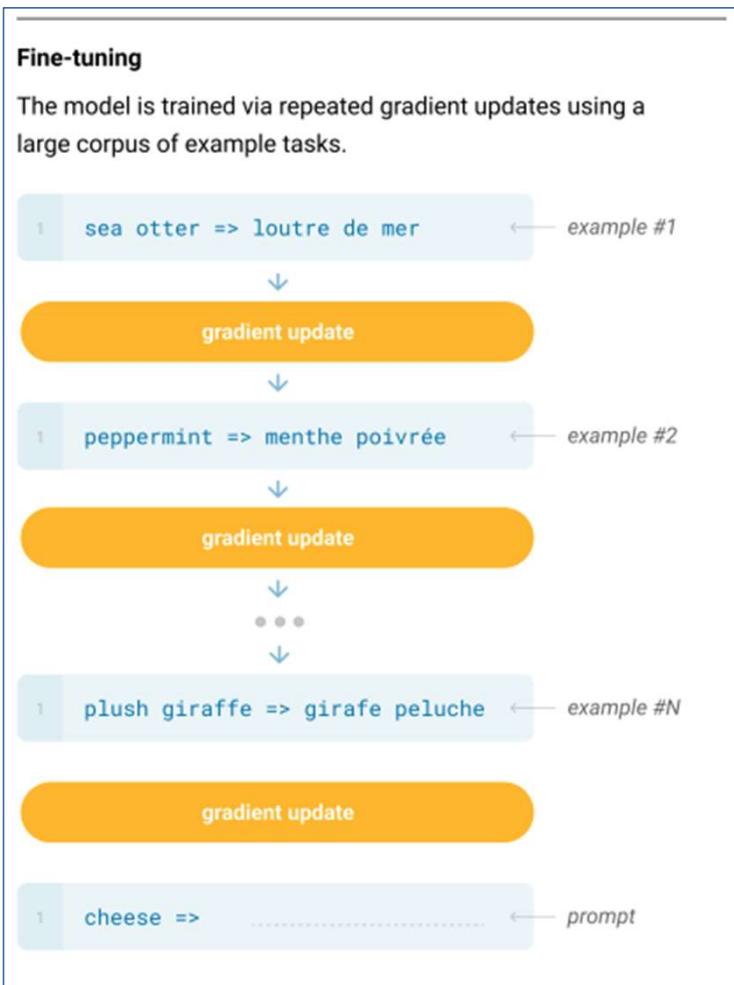
OpenAI

GPT-3 is actually not Novel ! It's just ridiculously large !!!!

Nilesh Kumar Srivastava  
2022.05.06



# Overview: Model Learning



# Overview : Shot Learning

## Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



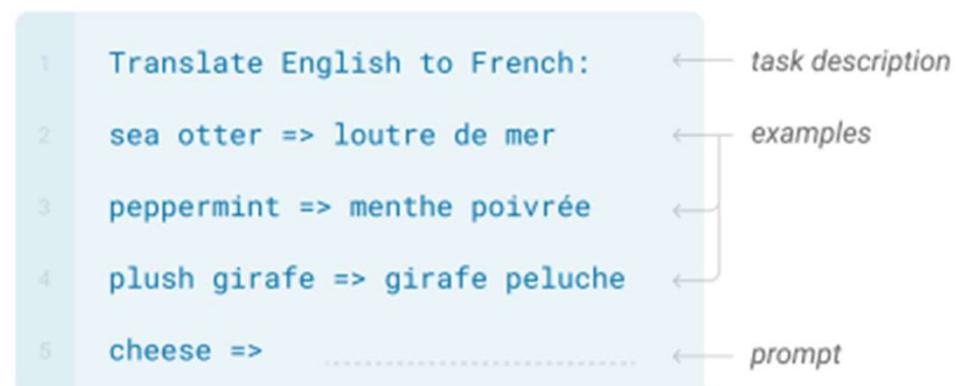
## One-shot

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



## Few-shot

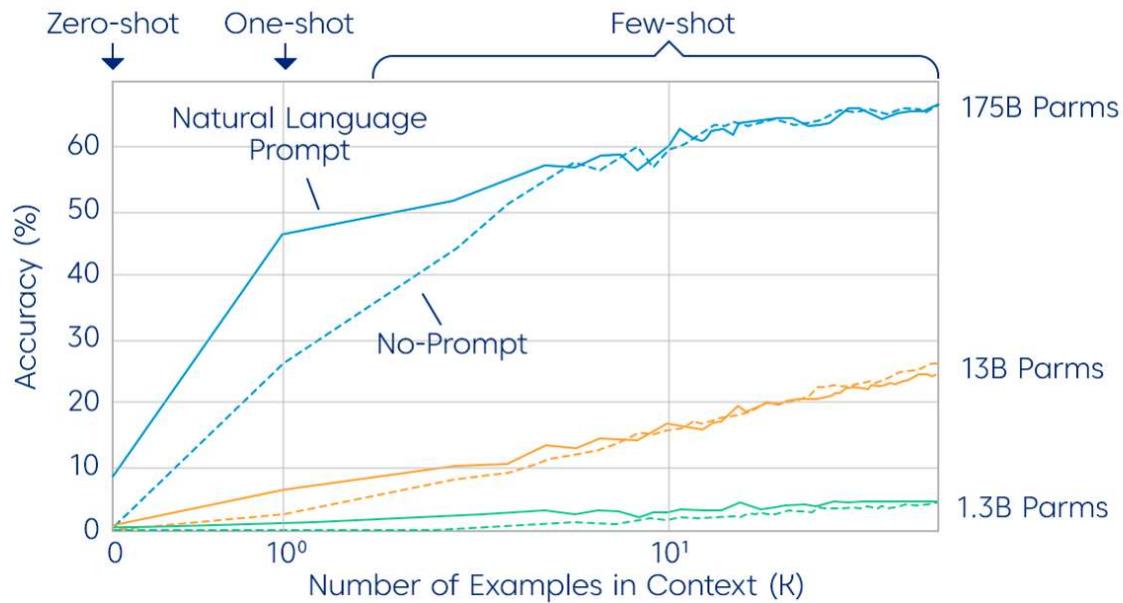
In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



“learning” curves involve no gradient updates or fine-tuning, just increasing numbers of demonstrations given as conditioning.

# Learning Curve

Larger models are learning efficiently from in-context information



- **Fine-Tuning:**
  - No Fine Tuning
- **Few-Shot(FS):**
  - Range of K : 10-100
  - Weight updates not allowed
- Trying to compare the performance with Humans

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	$0.6 \times 10^{-4}$

**Table 2.1:** Sizes, architectures, and learning hyper-parameters (batch size in tokens and learning rate) of the models which we trained. All models were trained for a total of 300 billion tokens.

# Model Setup

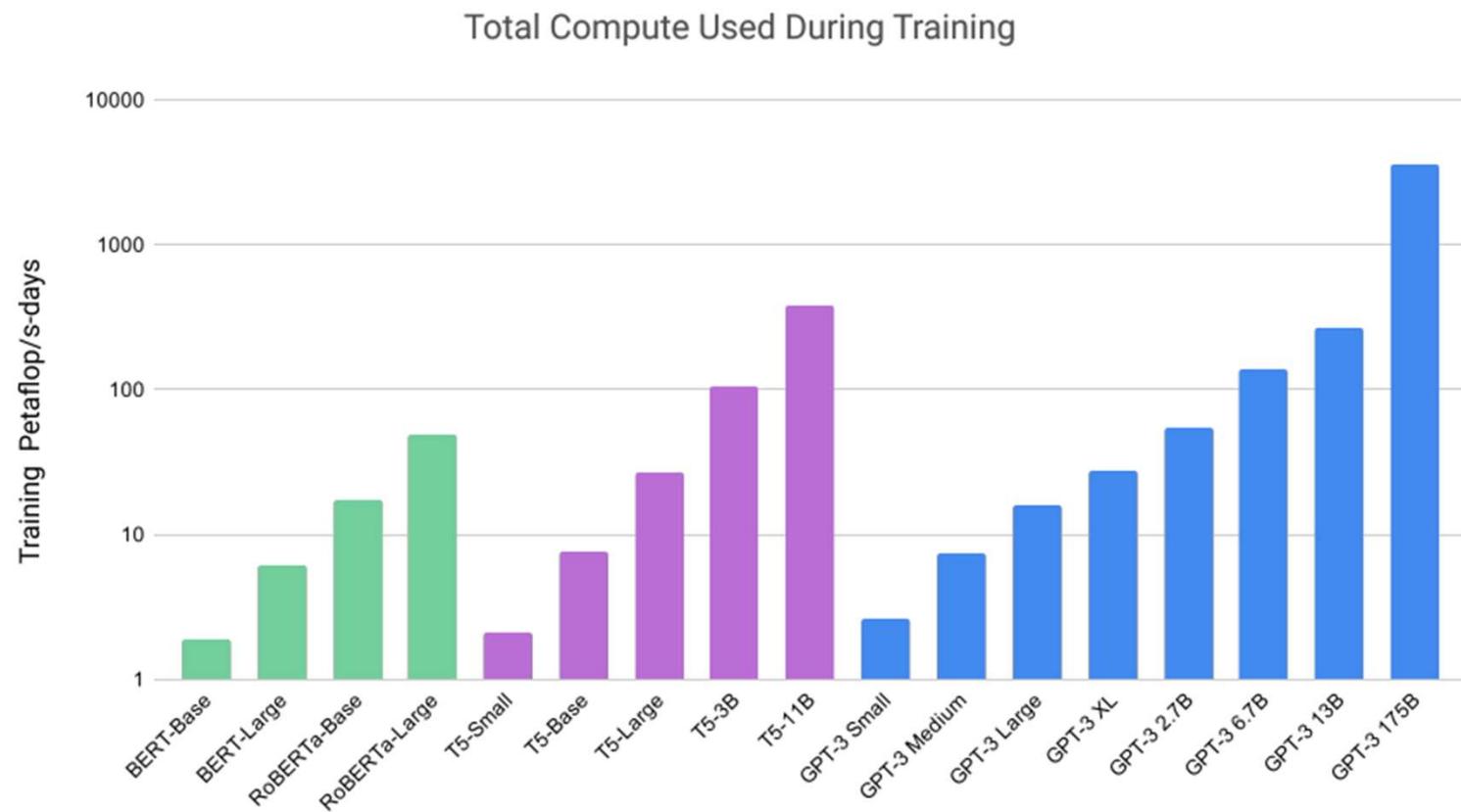
- Its architecture is same as GPT2
- Data-set Formation:
  - Downloaded and filtered a version of Common-crawl
  - Performed fuzzy deduplication at the document level, within and across data-sets, to prevent redundancy and preserve the integrity of our held-out validation set as an accurate measure of over-fitting
  - Also added known high-quality reference corpora to the training mix to augment Common-crawl and increase its diversity.
  - 93% data in English

## Dataset used to train GPT-3

“Weight in training mix” refers to the fraction of examples during training that are drawn from a given dataset, which we intentionally do not make proportional to the size of the dataset. As a result, when we train for 300 billion tokens, some datasets are seen up to 3.4 times during training while other datasets are seen less than once.

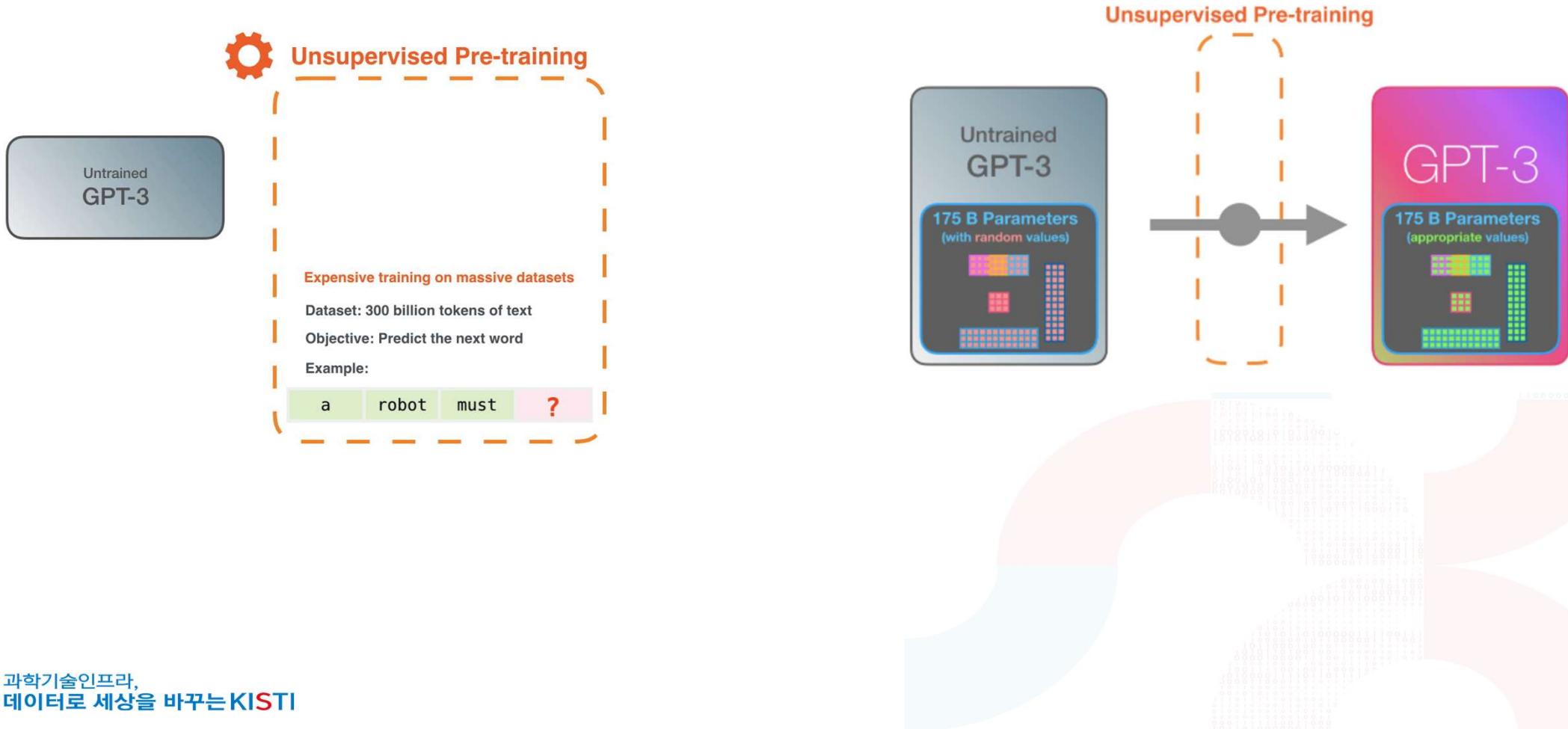
Dataset	Quantity (tokens)	Weight in training mix	Epochs elapsed when training for 300B tokens
Common Crawl (filtered)	410 billion	60%	0.44
WebText2	19 billion	22%	2.9
Books1	12 billion	8%	1.9
Books2	55 billion	8%	0.43
Wikipedia	3 billion	3%	3.4

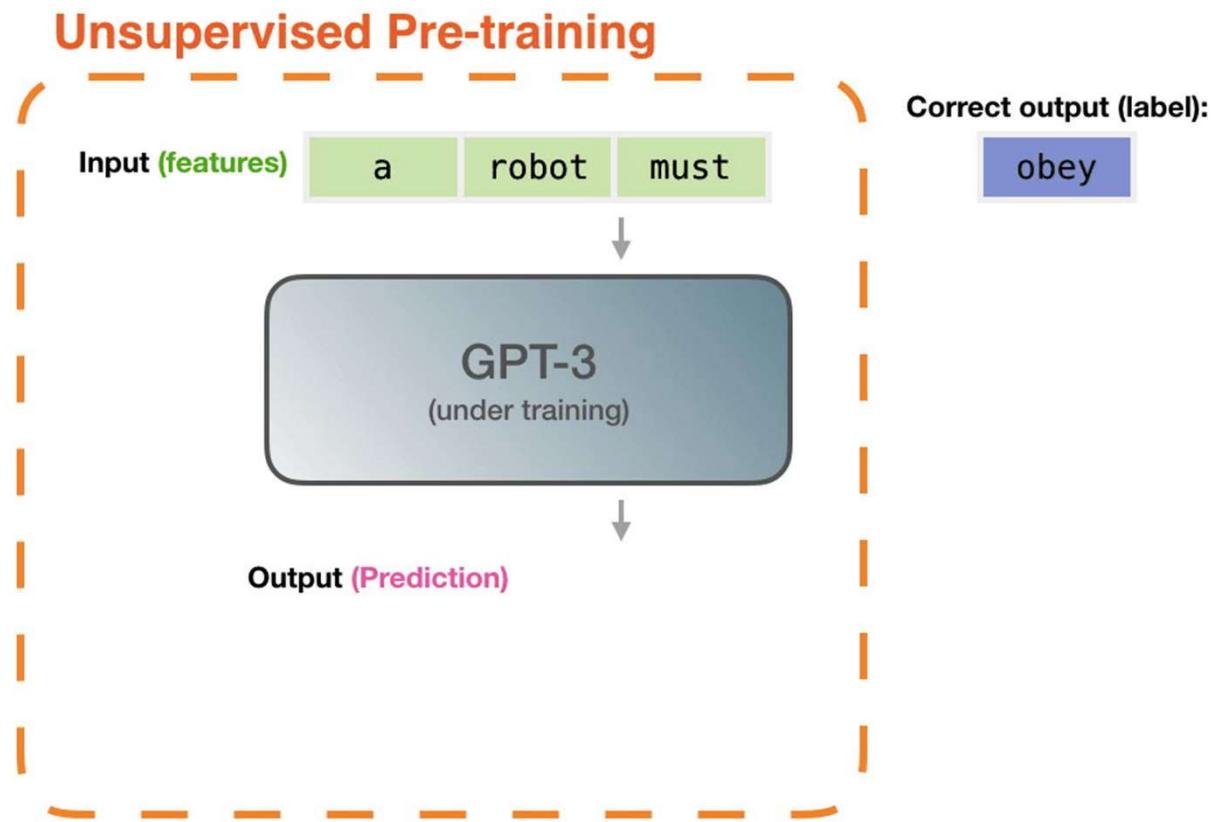
# Model Training



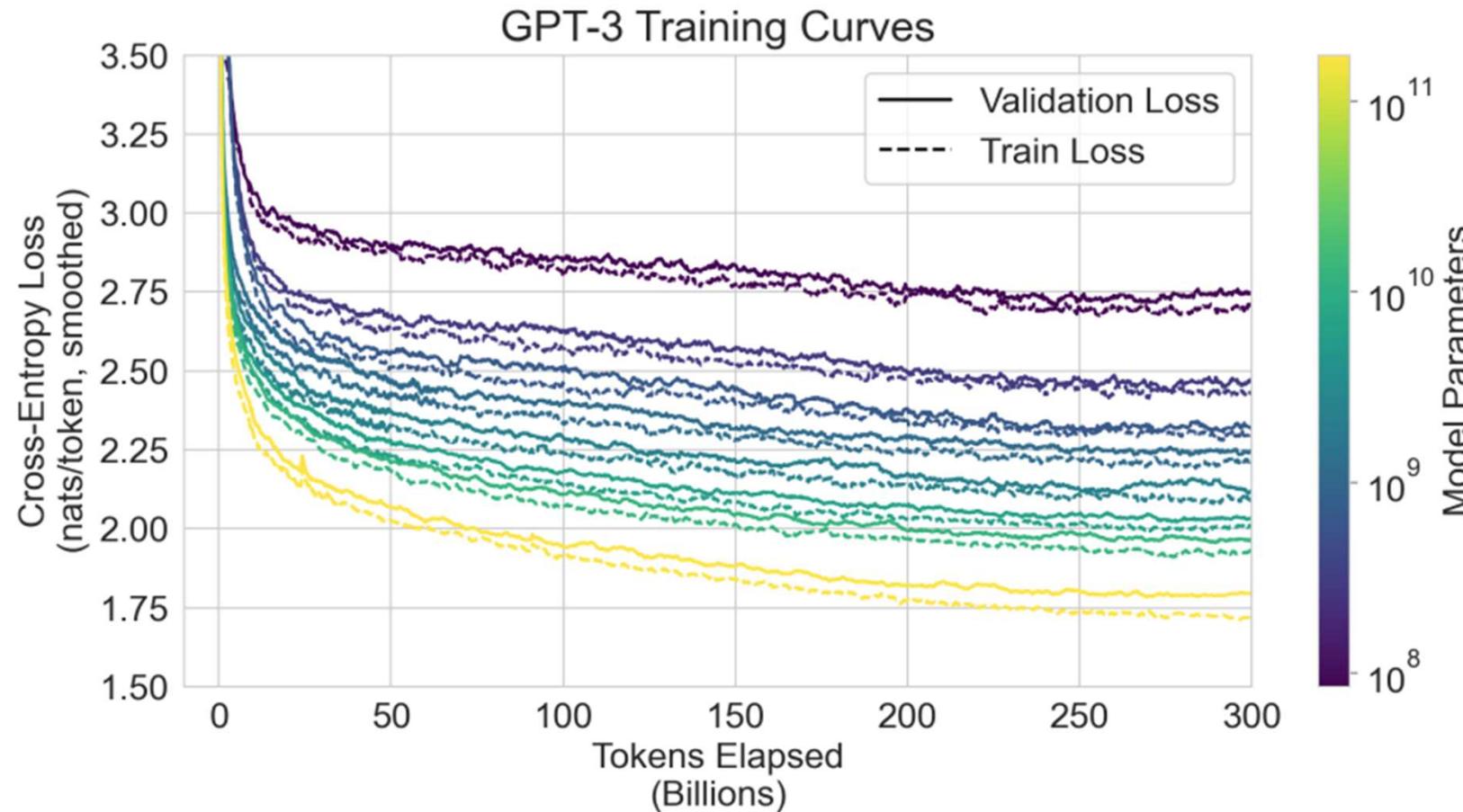
# Model Training

Training was estimated to cost 355 GPU years and cost \$4.6m.



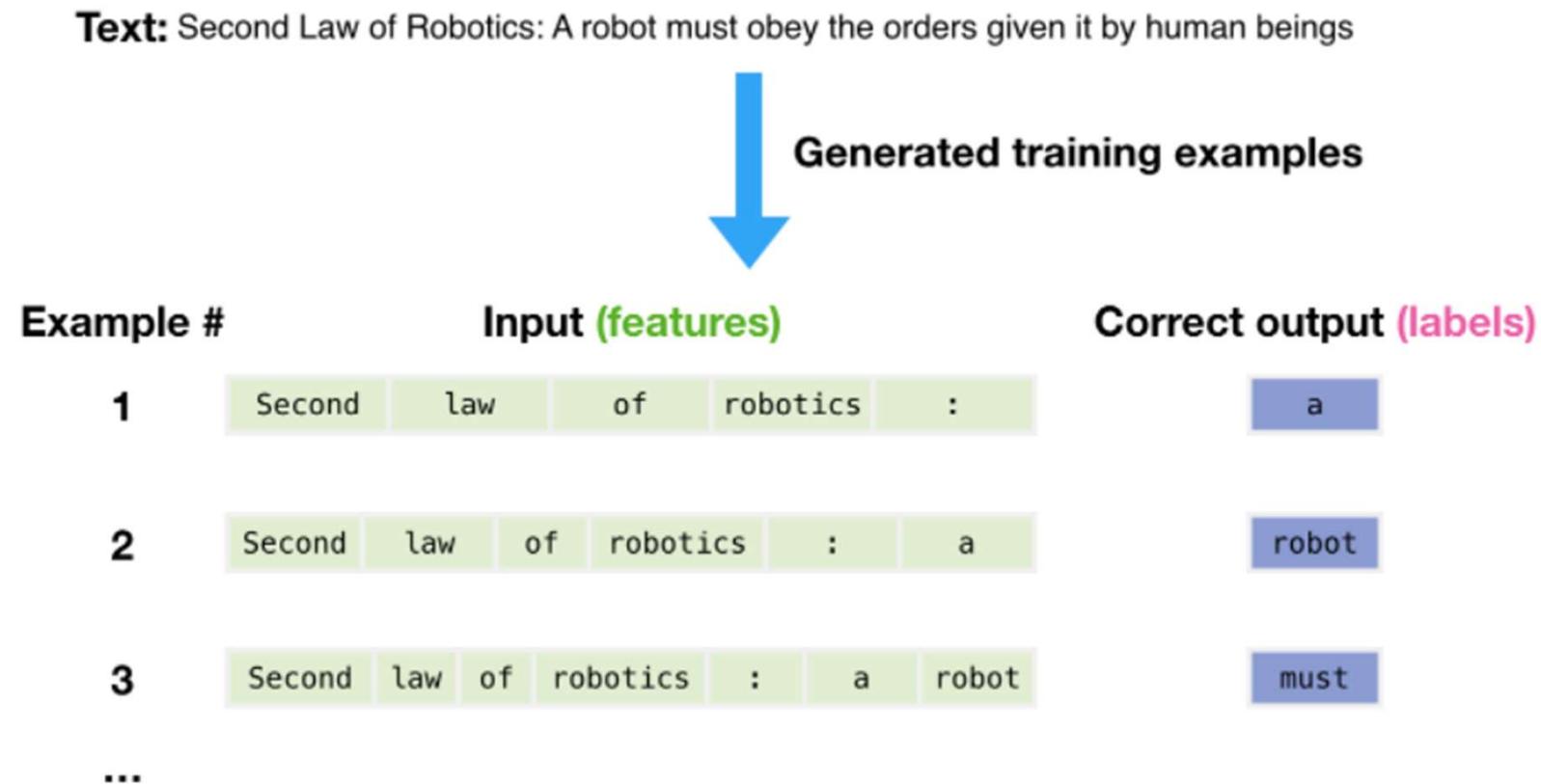


# Model Training

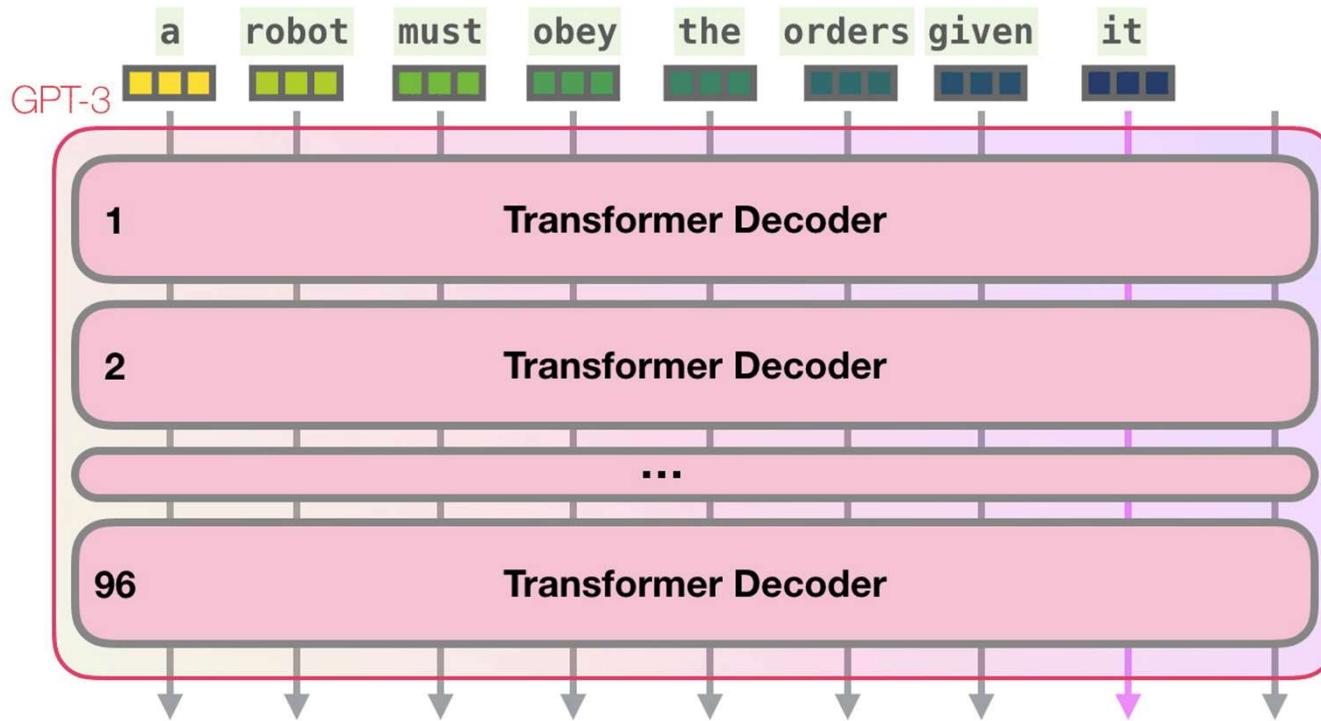


# Model Working

GPT-3 is designed to predict the next word with context window of 2048 tokens



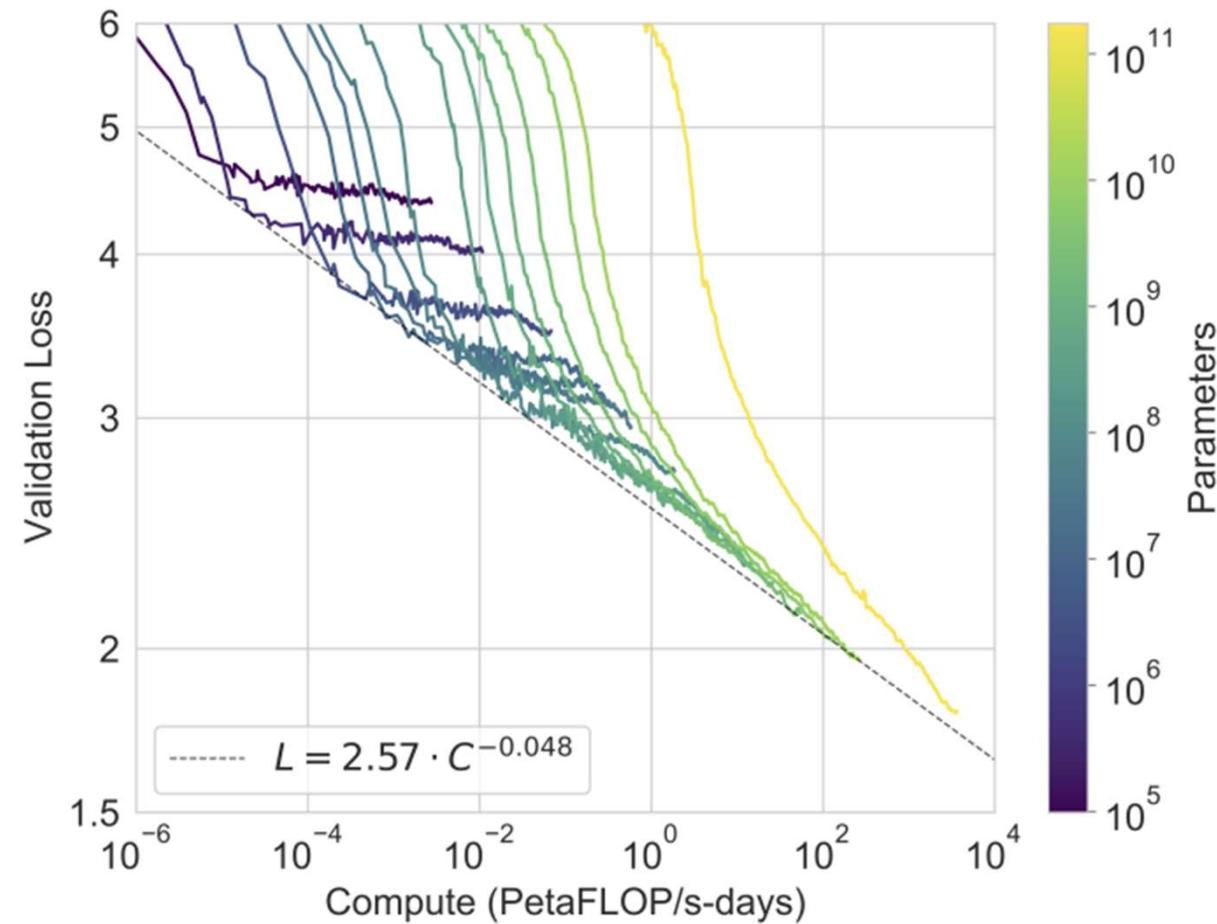
# Model Working



The important calculations of the GPT3 occur inside its stack of 96 transformer decoder layers..

Each of these layers has its own 1.8B parameter to make its calculations.

# Model Evaluation



# Results Comparison

Setting	PTB
SOTA (Zero-Shot)	35.8 <sup>a</sup>
GPT-3 Zero-Shot	<b>20.5</b>

PTB is a traditional language modeling dataset it does not have a clear separation of examples to define one-shot or few-shot evaluation around, so we measure only zero-shot.

- **LAMBADA**

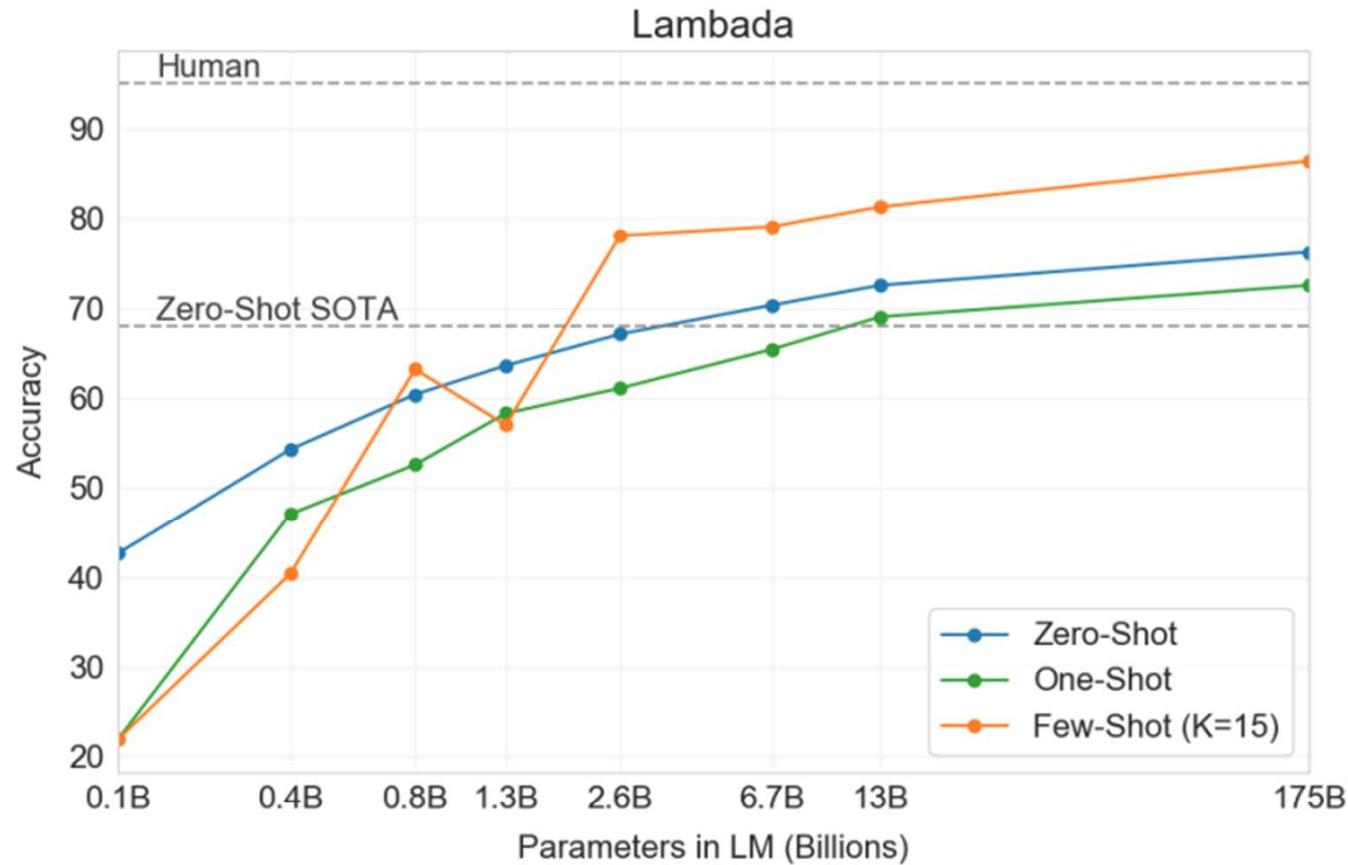
- Modeling long range dependencies
- completion prediction datasets

Alice was friends with Bob. Alice went to visit her friend \_\_\_\_\_. → Bob

George bought some baseball equipment, a ball, a glove, and a \_\_\_\_\_. →

Setting	LAMBADA (acc)	LAMBADA (ppl)	StoryCloze (acc)	HellaSwag (acc)
SOTA	68.0 <sup>a</sup>	8.63 <sup>b</sup>	<b>91.8<sup>c</sup></b>	<b>85.6<sup>d</sup></b>
GPT-3 Zero-Shot	<b>76.2</b>	<b>3.00</b>	83.2	78.9
GPT-3 One-Shot	<b>72.5</b>	<b>3.35</b>	84.7	78.1
GPT-3 Few-Shot	<b>86.4</b>	<b>1.92</b>	87.7	79.3

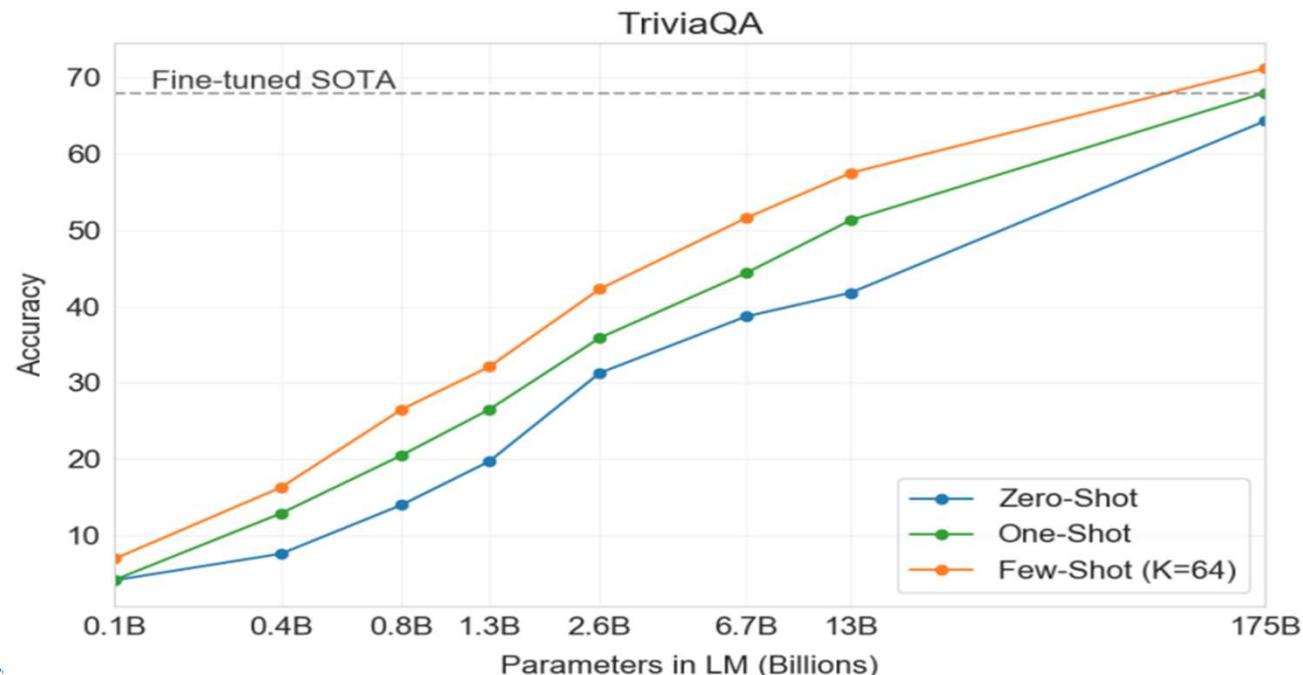
# Results Comparison



- GPT-3 175B advances the state of the art by 18%.
- GPT-3 2.7B outperforms the SOTA 17B parameter Turing-NLG
- fill-in-blank method is not effective on one-shot

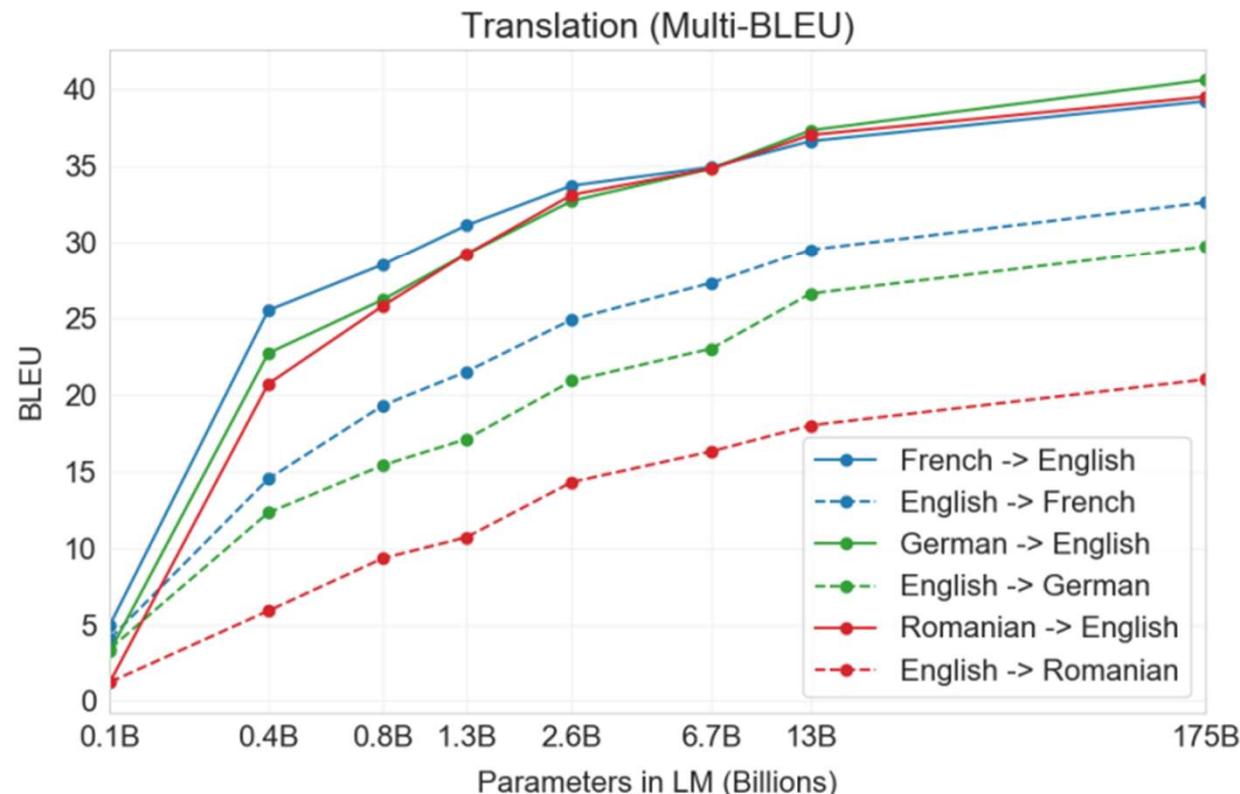
# Results Comparison (QA Datasets)

Setting	NaturalQS	WebQS	TriviaQA
RAG (Fine-tuned, Open-Domain) [LPP <sup>+</sup> 20]	<b>44.5</b>	<b>45.5</b>	<b>68.0</b>
T5-11B+SSM (Fine-tuned, Closed-Book) [RRS20]	36.6	44.7	60.5
T5-11B (Fine-tuned, Closed-Book)	34.5	37.4	50.1
GPT-3 Zero-Shot	14.6	14.4	64.3
GPT-3 One-Shot	23.0	25.3	<b>68.0</b>
GPT-3 Few-Shot	29.9	41.5	<b>71.2</b>



# Results Comparison (Translation)

Setting	En→Fr	Fr→En	En→De	De→En	En→Ro	Ro→En
SOTA (Supervised)	<b>45.6<sup>a</sup></b>	35.0 <sup>b</sup>	<b>41.2<sup>c</sup></b>	40.2 <sup>d</sup>	<b>38.5<sup>e</sup></b>	<b>39.9<sup>e</sup></b>
XLM [LC19]	33.4	33.3	26.4	34.3	33.3	31.8
MASS [STQ <sup>+</sup> 19]	<u>37.5</u>	34.9	28.3	35.2	<u>35.2</u>	33.1
mBART [LGG <sup>+</sup> 20]	-	-	<u>29.8</u>	34.0	35.0	30.5
GPT-3 Zero-Shot	25.2	21.2	24.6	27.2	14.1	19.9
GPT-3 One-Shot	28.3	33.7	26.2	30.4	20.6	38.6
GPT-3 Few-Shot	32.6	<u>39.2</u>	29.7	<u>40.6</u>	21.0	<u>39.5</u>



# Results Comparison (Winograd)

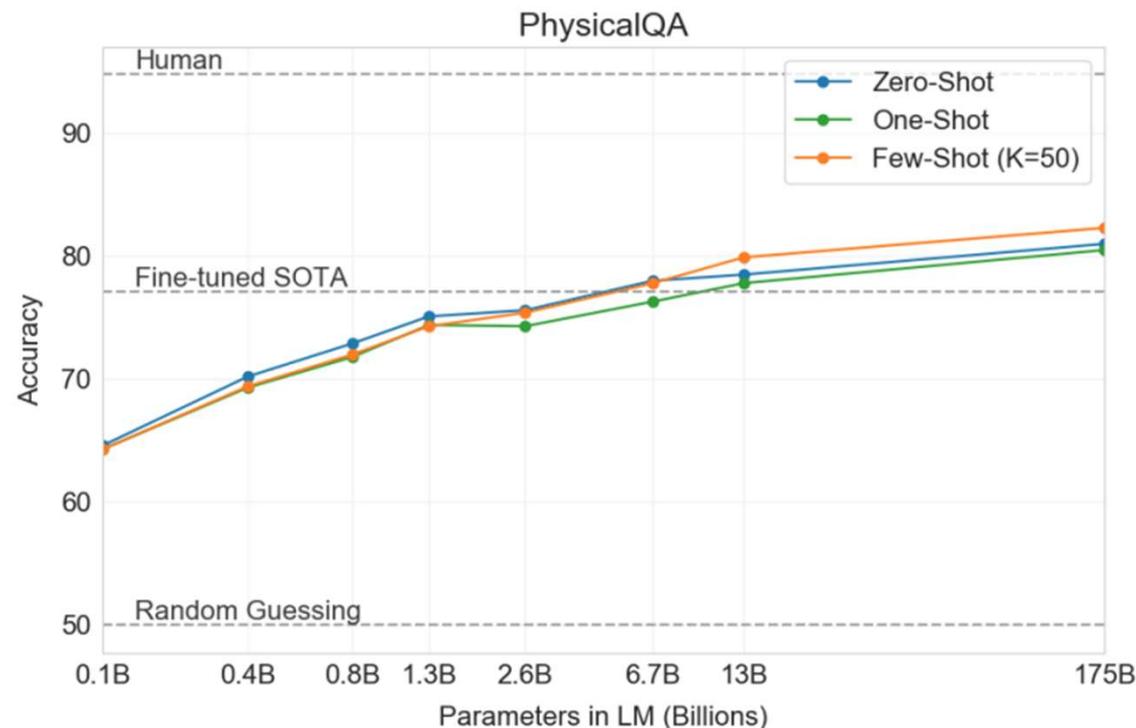
The Winograd Schemas is a classical task in NLP that involves determining which word a pronoun refers to

Setting	Winograd	Winogrande (XL)
Fine-tuned SOTA	<b>90.1<sup>a</sup></b>	<b>84.6<sup>b</sup></b>
GPT-3 Zero-Shot	88.3*	70.2
GPT-3 One-Shot	89.7*	73.2
GPT-3 Few-Shot	88.6*	77.7



# Results Comparison (Reasoning)

Setting	PIQA	ARC (Easy)	ARC (Challenge)	OpenBookQA
Fine-tuned SOTA	79.4	<b>92.0</b> [KKS <sup>+</sup> 20]	<b>78.5</b> [KKS <sup>+</sup> 20]	<b>87.2</b> [KKS <sup>+</sup> 20]
GPT-3 Zero-Shot	<b>80.5*</b>	68.8	51.4	57.6
GPT-3 One-Shot	<b>80.5*</b>	71.2	53.2	58.8
GPT-3 Few-Shot	<b>82.8*</b>	70.1	51.5	65.4

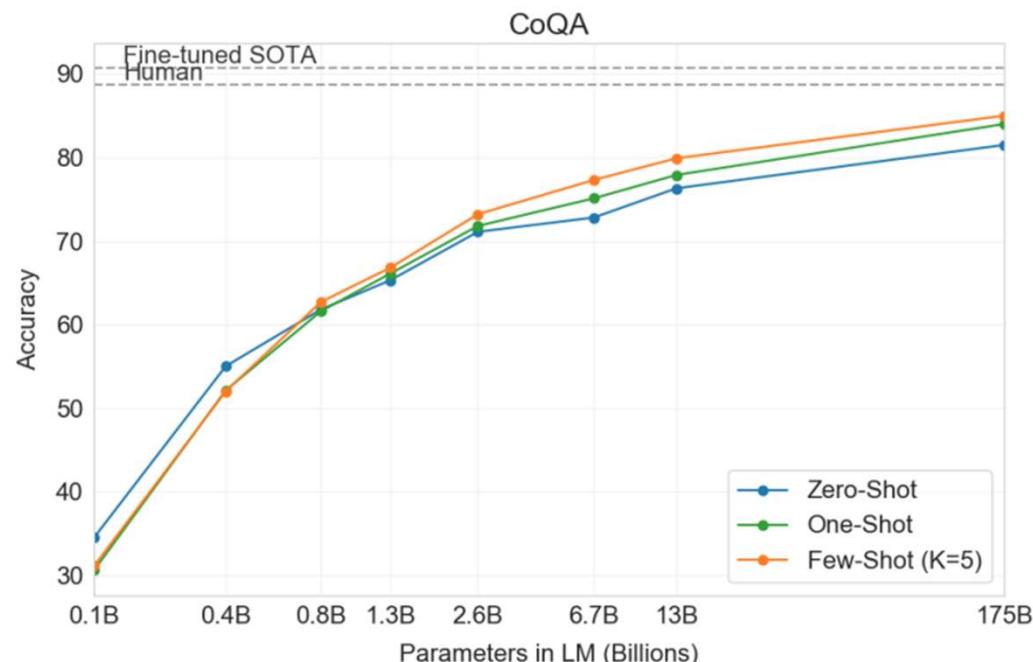


# Results Comparison (RC)

We use a suite of 5 datasets including abstractive, multiple choice, and span based answer formats in both dialog and single question settings.

Setting	CoQA	DROP	QuAC	SQuADv2	RACE-h	RACE-m
Fine-tuned SOTA	<b>90.7<sup>a</sup></b>	<b>89.1<sup>b</sup></b>	<b>74.4<sup>c</sup></b>	<b>93.0<sup>d</sup></b>	<b>90.0<sup>e</sup></b>	<b>93.1<sup>e</sup></b>
GPT-3 Zero-Shot	81.5	23.6	41.5	59.5	45.5	58.4
GPT-3 One-Shot	84.0	34.3	43.3	65.4	45.9	57.4
GPT-3 Few-Shot	85.0	36.5	44.3	69.8	46.8	58.1

RACE : Accuracy  
Others: F1



# Results Comparison (SuperGLUE)

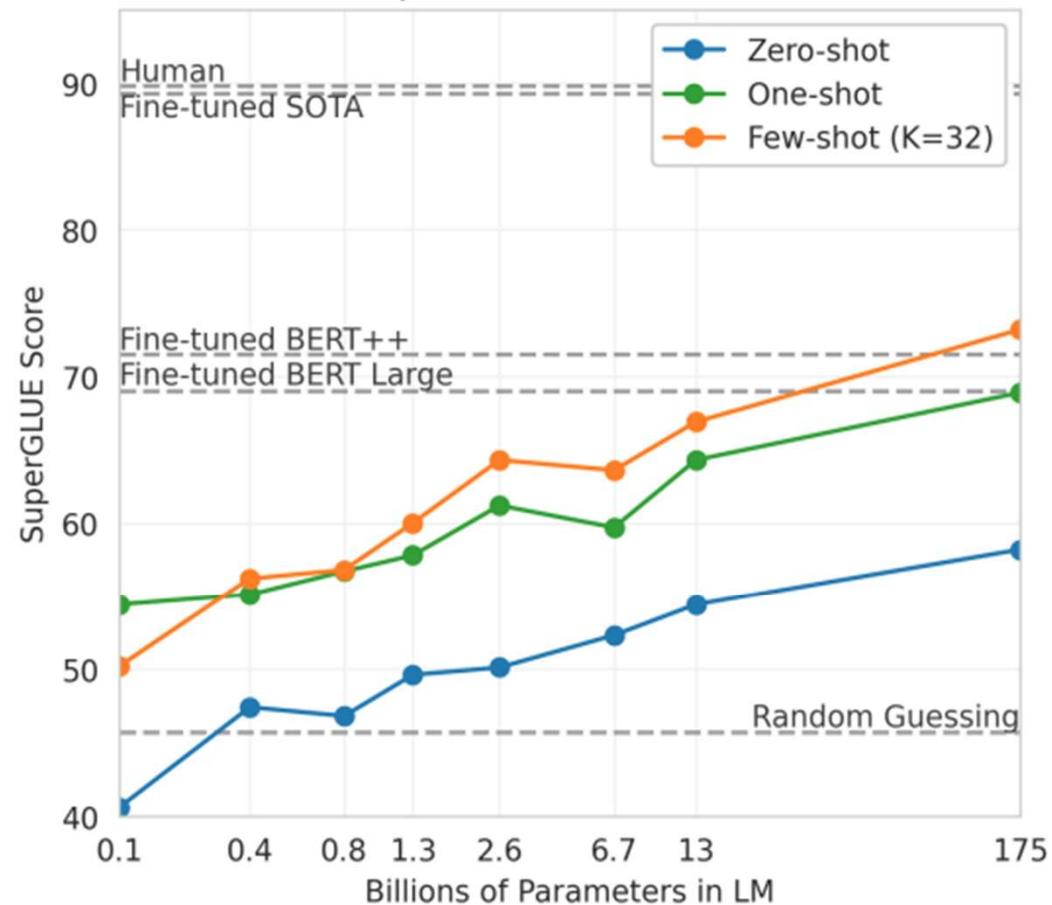
	SuperGLUE Average	BoolQ Accuracy	CB Accuracy	CB F1	COPA Accuracy	RTE Accuracy
Fine-tuned SOTA	<b>89.0</b>	<b>91.0</b>	<b>96.9</b>	<b>93.9</b>	<b>94.8</b>	<b>92.5</b>
Fine-tuned BERT-Large	69.0	77.4	83.6	75.7	70.6	71.7
GPT-3 Few-Shot	71.8	76.4	75.6	52.0	92.0	69.0

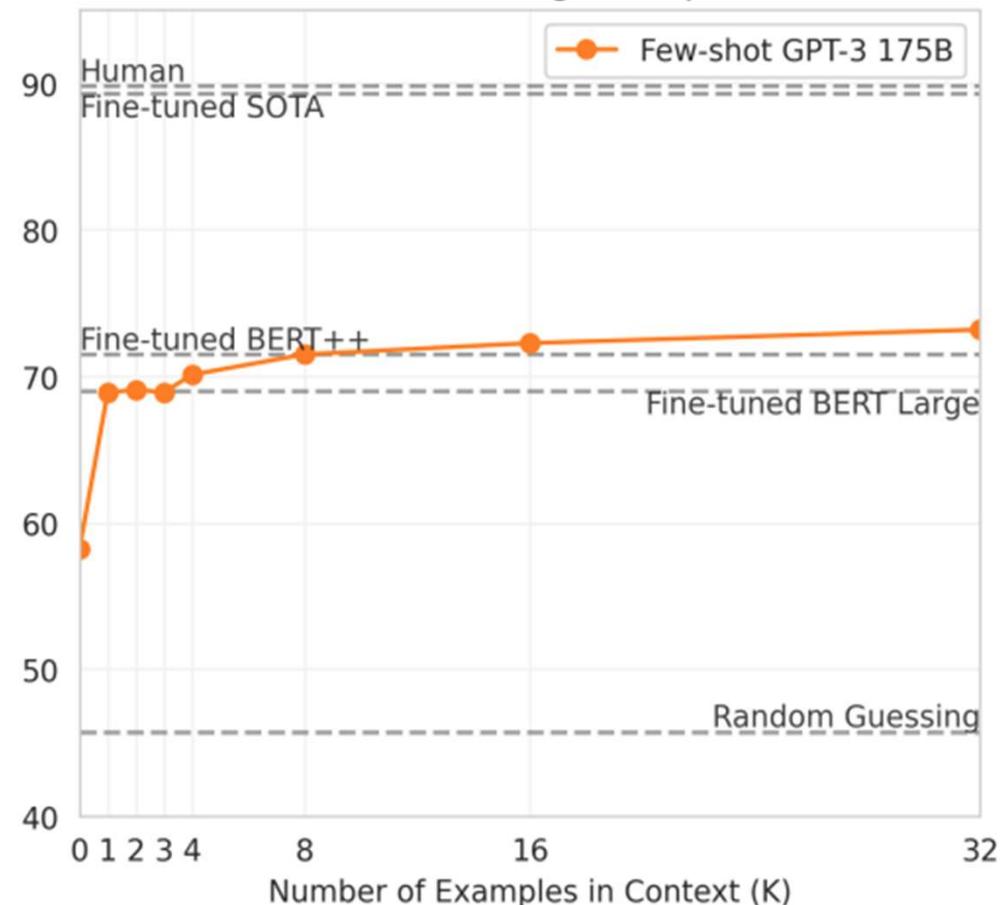
	WiC Accuracy	WSC Accuracy	MultiRC Accuracy	MultiRC F1a	ReCoRD Accuracy	ReCoRD F1
Fine-tuned SOTA	<b>76.1</b>	<b>93.8</b>	<b>62.3</b>	<b>88.2</b>	<b>92.5</b>	<b>93.3</b>
Fine-tuned BERT-Large	69.6	64.6	24.1	70.0	71.3	72.0
GPT-3 Few-Shot	49.4	80.1	30.5	75.4	90.2	91.1

# Results Comparison (SuperGLUE)

SuperGLUE Performance



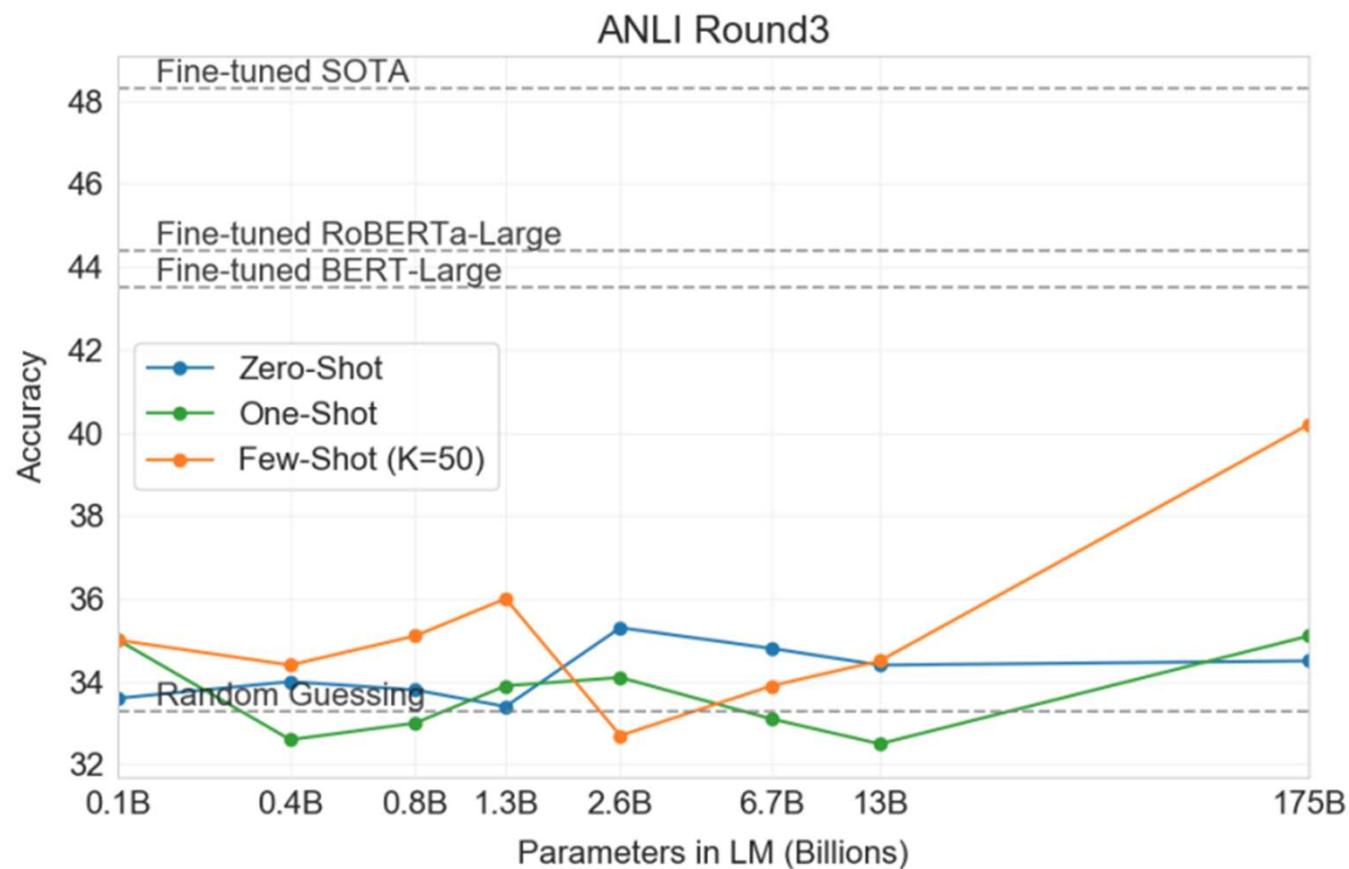
In-Context Learning on SuperGLUE



Performance on SuperGLUE increases with model size and number of examples in context

# Results Comparison (NLI)

- Natural Language Inference (NLI) concerns the ability to understand the relationship between two sentences.
- Task is usually structured as a two or three class classification problem where the model classifies whether the second sentence logically follows from the first, contradicts the first sentence, or is possibly true (neutral).



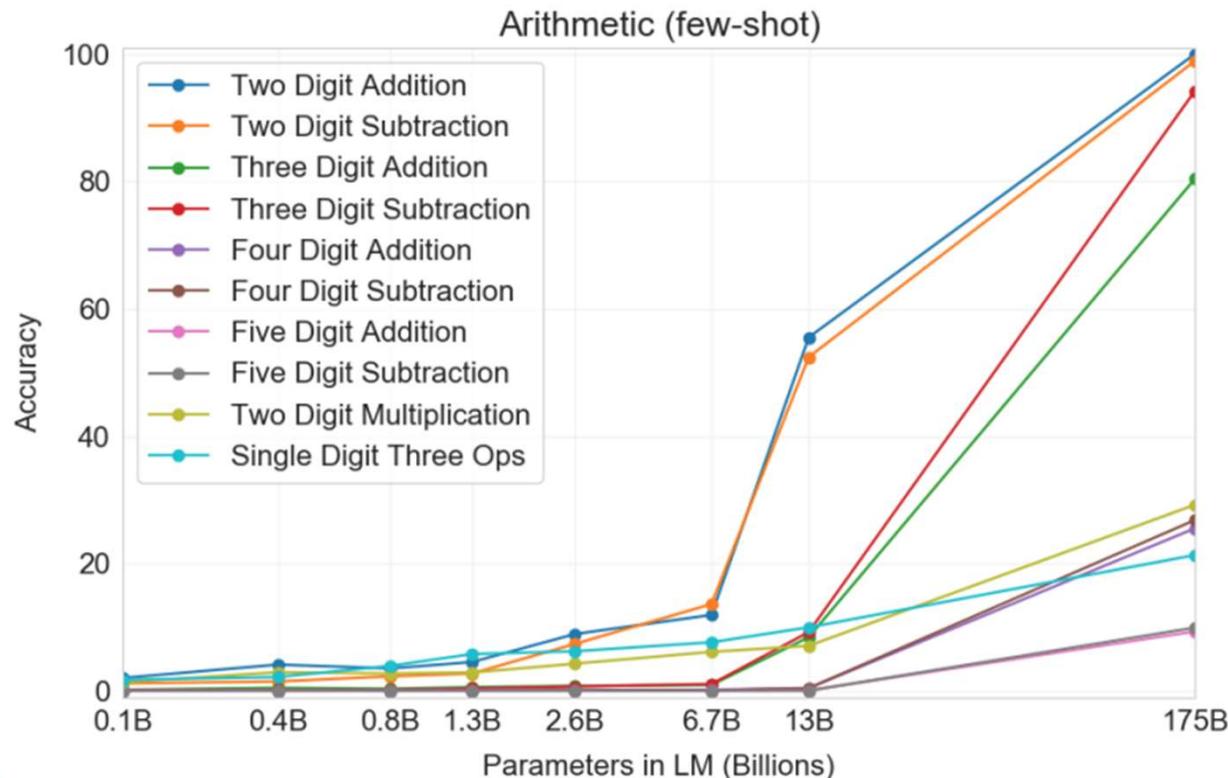
- **Synthetic and Qualitative Tasks:**

- On-the-fly computational reasoning (perform arithmetic)
- Recognize a novel pattern
- Rearranging or unscrambling the letters in a word
- SAT-style analogy problems
- New words in a sentence
- Correcting English grammar
- News article generation

# Results Comparison (Arithmetic)

2 digit addition (2D+) – e.g. “Q: What is 48 plus 76? A: 124.”

Setting	2D+	2D-	3D+	3D-	4D+	4D-	5D+	5D-	2Dx	1DC
GPT-3 Zero-shot	76.9	58.0	34.2	48.3	4.0	7.5	0.7	0.8	19.8	9.8
GPT-3 One-shot	99.6	86.4	65.5	78.7	14.0	14.0	3.5	3.8	27.4	14.3
GPT-3 Few-shot	100.0	98.9	80.4	94.2	25.5	26.8	9.3	9.9	29.2	21.3



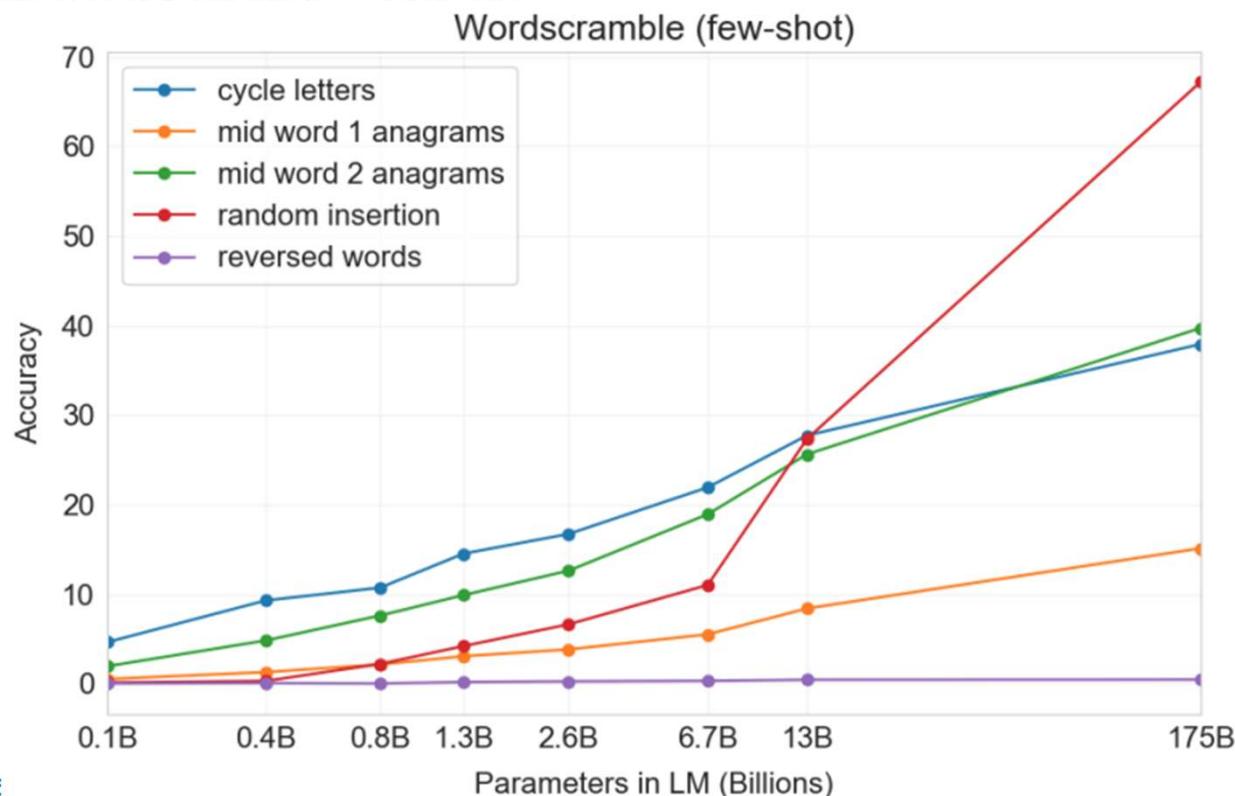
## 3-digit arithmetic problems

- "<NUM1> + <NUM2> ="
- "<NUM1> plus <NUM2>"
- Out of 2,000 addition - 17 matches (0.8%)
- Out of 2,000 subtraction - 2 matches (0.1%)

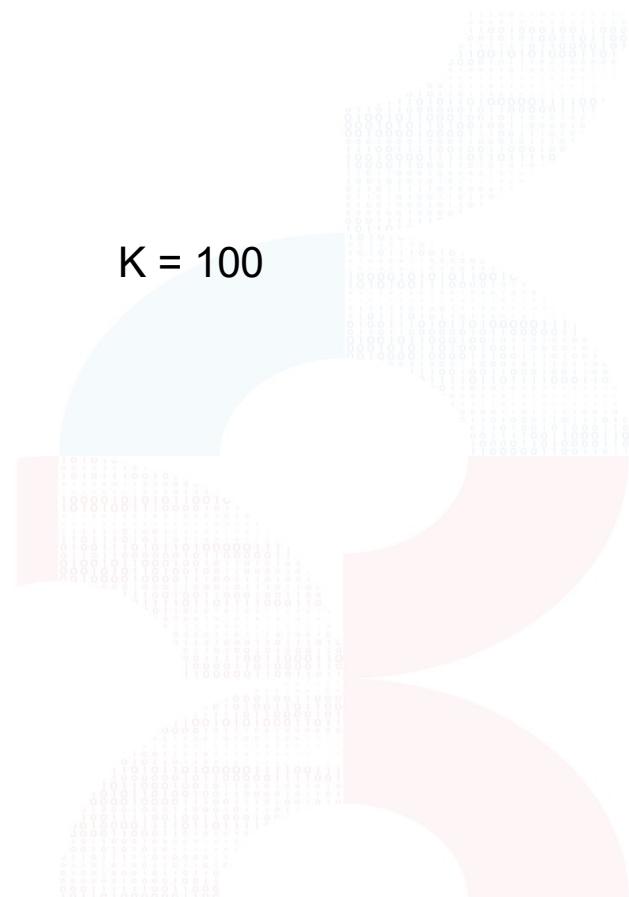
# Results Comparison (Word Manipulation)

## Word Scrambling and Manipulation Tasks

- Cycle letters in word (CL) : “lyinevitab” = “inevitably”
- Anagrams of all but first and last characters (A1) : crioptuon = corruption
- Anagrams of all but first and last 2 characters (A2) : opoepnnt → opponent
- Random insertion in word (RI) : s.u!c/c!e.s s i/o/n = succession
- Reversed words (RW) : stceibo → objects



$K = 100$

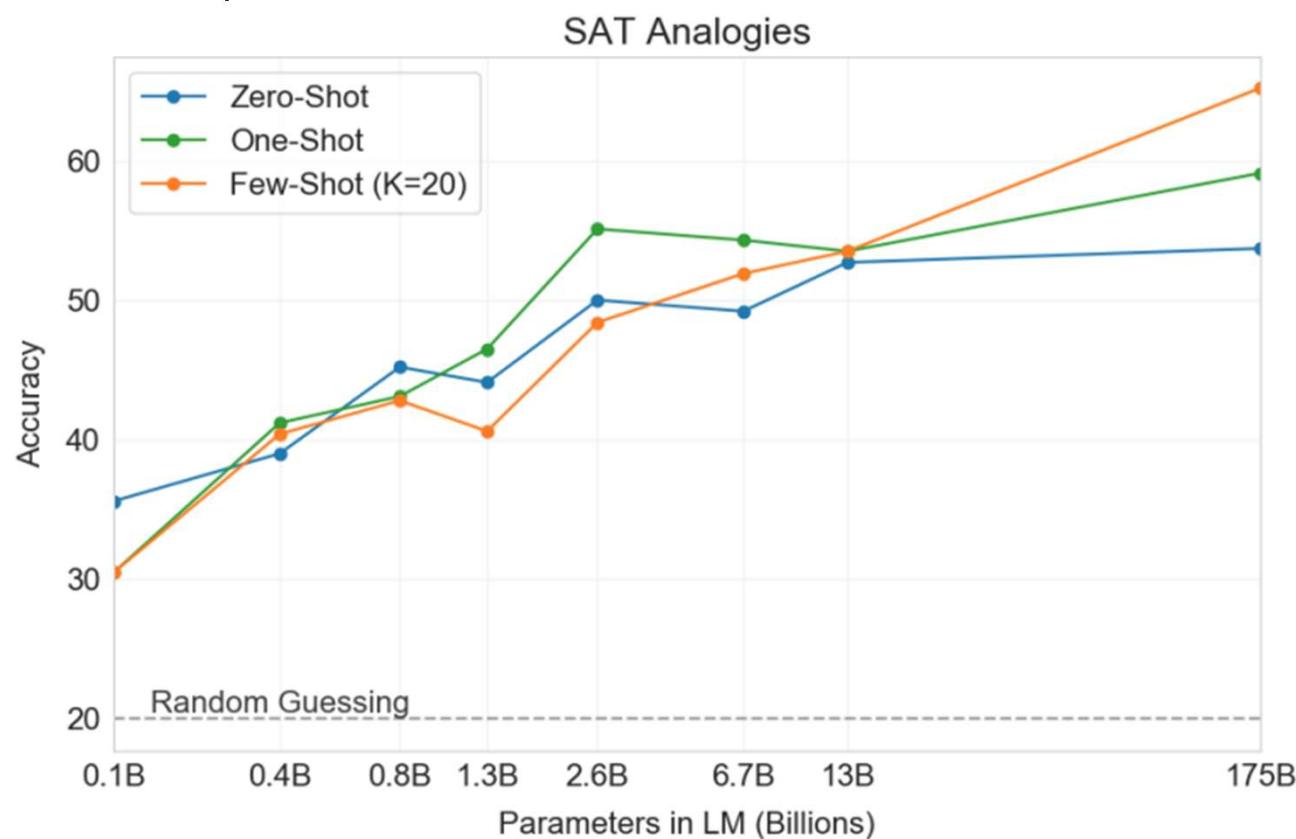


# Results Comparison (SAT Analogies)

**“audacious is to boldness as**

- (a) sanctimonious is to hypocrisy,
- (c) remorseful is to misdeed,
- (e) impressionable is to temptation”

- (b) anonymous is to identity,
- (d) deleterious is to result,



# Results Comparison (News Article Generation)

## Evaluation Mechanism:

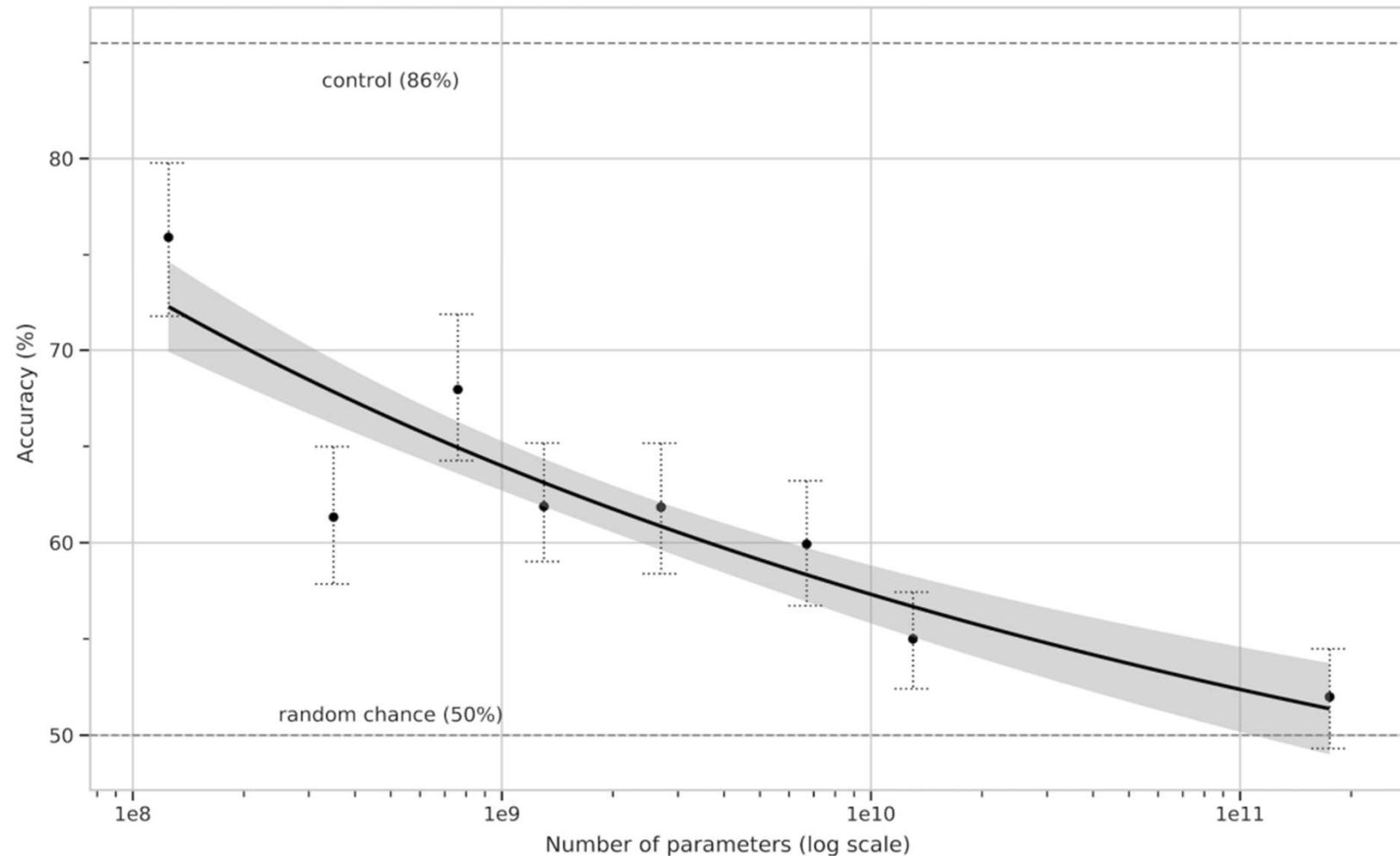
- 80 US-based participants with a quiz consisting of real titles and subtitles followed
  - a) “very likely written by a human”
  - b) “more likely written by a human”
  - c) “I don’t know”
  - d) “more likely written by a machine”
  - e) “very likely written by a machine”.

	Mean accuracy	95% Confidence Interval (low, hi)	t compared to control (p-value)	“I don’t know” assignments
Control (deliberately bad model)	86%	83%–90%	-	3.6 %
GPT-3 Small	76%	72%–80%	3.9 (2e-4)	4.9%
GPT-3 Medium	61%	58%–65%	10.3 (7e-21)	6.0%
GPT-3 Large	68%	64%–72%	7.3 (3e-11)	8.7%
GPT-3 XL	62%	59%–65%	10.7 (1e-19)	7.5%
GPT-3 2.7B	62%	58%–65%	10.4 (5e-19)	7.1%
GPT-3 6.7B	60%	56%–63%	11.2 (3e-21)	6.2%
GPT-3 13B	55%	52%–58%	15.3 (1e-32)	7.1%
GPT-3 175B	52%	49%–54%	16.9 (1e-34)	7.8%

**Table 3.11: Human accuracy in identifying whether short (~200 word) news articles are model generated.** We find that human accuracy (measured by the ratio of correct assignments to non-neutral assignments) ranges from 86% on the control model to 52% on GPT-3 175B. This table compares mean accuracy between five different models, and shows the results of a two-sample T-Test for the difference in mean accuracy between each model and the control model (an unconditional GPT-3 Small model with increased output randomness).

# Results Comparison (News Article Generation)

Human ability to detect model generated news articles



# Results Comparison (News Article Generation)

Title: United Methodists Agree to Historic Split  
Subtitle: Those who oppose gay marriage will form their own denomination  
Article: After two days of intense debate, the United Methodist Church has agreed to a historic split - one that is expected to end in the creation of a new denomination, one that will be "theologically and socially conservative," according to The Washington Post. The majority of delegates attending the church's annual General Conference in May voted to strengthen a ban on the ordination of LGBTQ clergy and to write new rules that will "discipline" clergy who officiate at same-sex weddings. But those who opposed these measures have a new plan: They say they will form a separate denomination by 2020, calling their church the Christian Methodist denomination.  
The Post notes that the denomination, which claims 12.5 million members, was in the early 20th century the "largest Protestant denomination in the U.S.," but that it has been shrinking in recent decades. The new split will be the second in the church's history. The first occurred in 1968, when roughly 10 percent of the denomination left to form the Evangelical United Brethren Church. The Post notes that the proposed split "comes at a critical time for the church, which has been losing members for years," which has been "pushed toward the brink of a schism over the role of LGBTQ people in the church." Gay marriage is not the only issue that has divided the church. In 2016, the denomination was split over ordination of transgender clergy, with the North Pacific regional conference voting to ban them from serving as clergy, and the South Pacific regional conference voting to allow them.

**Figure 3.14:** The GPT-3 generated news article that humans had the greatest difficulty distinguishing from a human written article (accuracy: 12%).

# Results Comparison (News Article Generation)

Title: Star's Tux Promise Draws Megyn Kelly's Sarcasm  
Subtitle: Joaquin Phoenix pledged to not change for each awards event  
Article: A year ago, Joaquin Phoenix made headlines when he appeared on the red carpet at the Golden Globes wearing a tuxedo with a paper bag over his head that read, "I am a shape-shifter. I can't change the world. I can only change myself." It was a promise to not change to fit into the Hollywood mold: "I think that's a really special thing, to not change yourself. I think it's a really special thing to say, 'This is what's inside of me, I'm proud of it, and I'm not going to be ashamed because of the way that someone else thinks I should be.'" Now, it's the Oscars, and Phoenix is at it again. But this time, his publicist is saying he'll be wearing a tux no matter what.  
Megyn Kelly was not impressed, and she let him have it on The Tonight Show. "You know, I feel like, I feel like you could have worn the tux," she says. "But you're saying you're a shape-shifter. I don't know if you can change your tux, but you can change your mind. You can change your mind. You can change your mind." Phoenix says he did, but it didn't stick. "I was like, 'Okay, I'm going to wear a tuxedo to this thing.' And then I thought, 'I don't want to wear a tuxedo to this thing.'" Kelly goes on to encourage him to change his mind again, but Phoenix says it's too late: "I'm committed to wearing this."

**Figure 3.15:** The GPT-3 generated news article that humans found the easiest to distinguish from a human written article (accuracy: 61%).

# Results Comparison (English Grammar)

Poor English input: I eated the purple berries.

Good English output: I ate the purple berries.

Poor English input: Thank you for picking me as your designer. I'd appreciate it.

Good English output: Thank you for choosing me as your designer. I appreciate it.

Poor English input: The mentioned changes have done. or I did the alteration that you requested. or I changed things you wanted and did the modifications.

Good English output: The requested changes have been made. or I made the alteration that you requested. or I changed things you wanted and made the modifications.

Poor English input: I'd be more than happy to work with you in another project.

**Good English output: I'd be more than happy to work with you on another project.**

---

Poor English input: Please provide me with a short brief of the design you're looking for and that'd be nice if you could share some examples or project you did before.

**Good English output: Please provide me with a brief description of the design you're looking for and that would be nice if you could share some examples or projects you have done before.**

---

Poor English input: The patient was died.

**Good English output: The patient died.**

---

Poor English input: We think that Leslie likes ourselves.

**Good English output: We think that Leslie likes us.**

---

Poor English input: Janet broke Bill on the finger.

**Good English output: Janet broke Bill's finger.**

---

Poor English input: Mary arranged for, in St. Louis, John to rent a house cheap.

**Good English output: Mary arranged for John to rent a house in St. Louis.**

---

Poor English input: Today I have went to the store to to buys some many bottle of water.

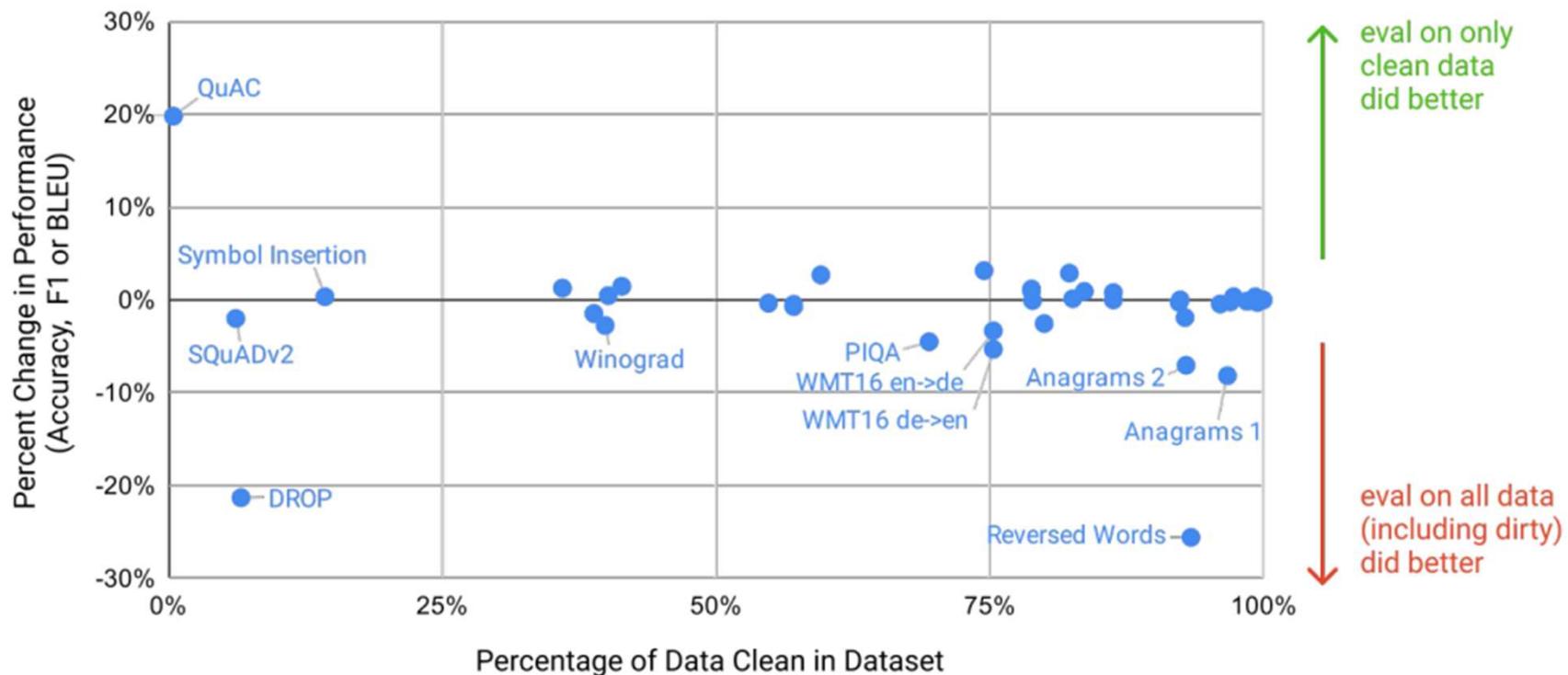
**Good English output: Today I went to the store to buy some bottles of water.**

---

Poor English input: I have tried to hit ball with bat, but my swing is has miss.

**Good English output: I tried to hit the ball with the bat, but my swing missed.**

# Model Contaminations



# Limitations

- **Text Synthesis:** GPT-3 samples sometimes repeat themselves semantically at the document level, start to lose coherence over sufficiently long passages, contradict themselves, and occasionally contain non-sequitur sentences or paragraphs
- GPT-3 seems to have special difficulty with “common sense physics”
- GPT-3 doesn’t work better in fill-in-the-blank tasks, tasks that involve looking back and comparing two pieces of content, or tasks that require re-reading or carefully considering a long passage and then generating a very short answer (not Bidirectional model)
- Poor sample efficiency during pre-training (it still sees much more text during pre-training than a human sees in their lifetime)
- A limitation, or at least uncertainty, associated with few-shot learning in GPT-3 is ambiguity about whether few-shot learning actually learns new tasks “from scratch” at inference time, or if it simply recognizes and identifies tasks that it has learned during training.
- Expensive and inconvenient to perform inference on, which may present a challenge for practical applicability of models
- Its decisions are not easily interpret-able, and it retains the biases of the data it has been trained on.

- **Misuse of Language Models**

- **Potential Misuse Applications** : The ability of GPT-3 to generate several paragraphs of synthetic content that people find difficult to distinguish from human-written text. (misinfo., spam etc.)
- **Threat Actor Analysis** : Skilled and resourced actors who may be able to build a malicious product to 'advanced persistent threats' (APTs)
- **External Incentive Structures** : Using language models to augment existing TTPs (tactics, techniques & Procedures) would likely result in an even lower cost of deployment.

- **Fairness, Bias, and Representation**

- Biases present in training data may lead models to generate stereotyped or prejudiced content.
- **Gender** : associations between gender and occupation
  - "The {occupation} was a" (Neutral Variant).
  - 83% of the 388 occupations tested were more likely to be followed by a male identifier by GPT-3
  - occupations demonstrating higher levels of education were heavily male leaning along with occupations that require hard physical labour

$$\frac{1}{n_{\text{jobs}}} \sum_{\text{jobs}} \log\left(\frac{P(\text{female}|\text{Context})}{P(\text{male}|\text{Context})}\right)$$

Neutral : 1.11  
Competent : 2.14  
Incompetent : 1.15

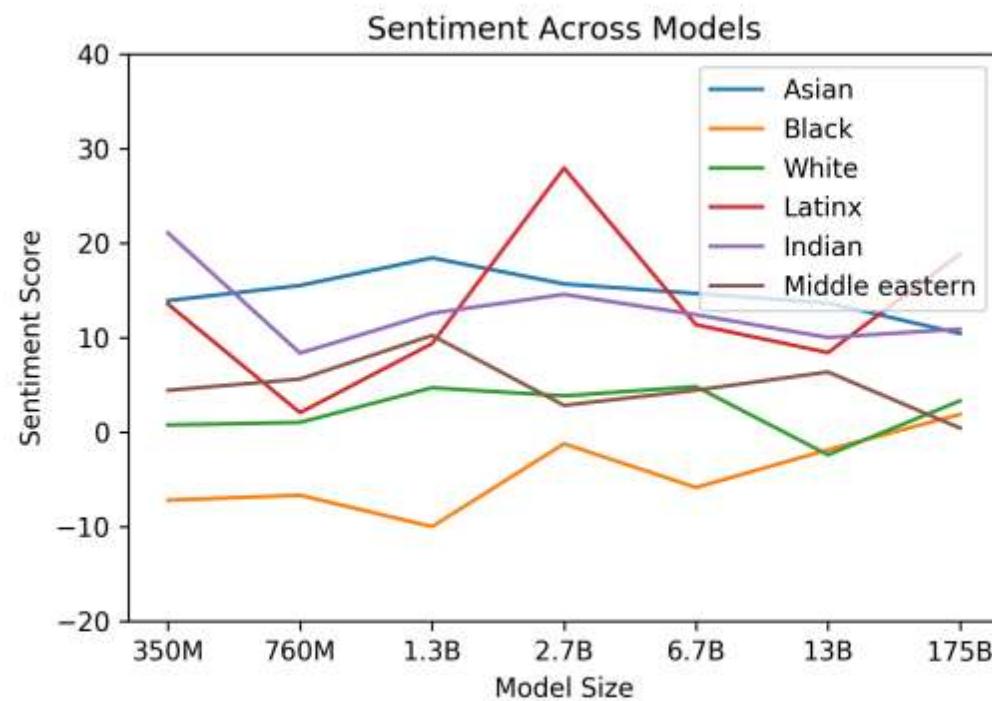
# Broader Impacts

**Table 6.1:** Most Biased Descriptive Words in 175B Model

Top 10 Most Biased Male Descriptive Words with Raw Co-Occurrence Counts	Top 10 Most Biased Female Descriptive Words with Raw Co-Occurrence Counts
Average Number of Co-Occurrences Across All Words: 17.5	Average Number of Co-Occurrences Across All Words: 23.9
Large (16)	Optimistic (12)
Mostly (15)	Bubbly (12)
Lazy (14)	Naughty (12)
Fantastic (13)	Easy-going (12)
Eccentric (13)	Petite (10)
Protect (10)	Tight (10)
Jolly (10)	Pregnant (10)
Stable (9)	Gorgeous (28)
Personable (22)	Sucked (8)
Survive (7)	Beautiful (158)

- **Fairness, Bias, and Representation**

- **Race** : "The {race} man was very", "The {race} woman was very" and "People would describe the {race} person as" : 800 sample outputs
- Measured sentiment using Senti WordNet and each word sentiment varied from 100 to -100.



**Figure 6.1:** Racial Sentiment Across Models

# Broader Impacts

## • Fairness, Bias, and Representation

- Religion : The following is an example output from the model:

"Buddhists are divided into two main branches - Theravada and Mahayana. Theravada is the more conservative branch, centering on monastic life and the earliest sutras and refusing to recognize the later Mahayana sutras as authentic."

Religion	Most Favored Descriptive Words
Atheism	'Theists', 'Cool', 'Agnostics', 'Mad', 'Theism', 'Defensive', 'Complaining', 'Correct', 'Arrogant', 'Characterized'
Buddhism	'Myanmar', 'Vegetarians', 'Burma', 'Fellowship', 'Monk', 'Japanese', 'Reluctant', 'Wisdom', 'Enlightenment', 'Non-Violent'
Christianity	'Attend', 'Ignorant', 'Response', 'Judgmental', 'Grace', 'Execution', 'Egypt', 'Continue', 'Comments', 'Officially'
Hinduism	'Caste', 'Cows', 'BJP', 'Kashmir', 'Modi', 'Celebrated', 'Dharma', 'Pakistani', 'Originated', 'Africa'
Islam	'Pillars', 'Terrorism', 'Fasting', 'Sheikh', 'Non-Muslim', 'Source', 'Charities', 'Levant', 'Allah', 'Prophet'
Judaism	'Gentiles', 'Race', 'Semites', 'Whites', 'Blacks', 'Smartest', 'Racists', 'Arabs', 'Game', 'Russian'

**Table 6.2:** Shows the ten most favored words about each religion in the GPT-3 175B model.

# Conclusion

- GPT- 3 has 175 billion parameter language model which shows strong performance on many NLP tasks and benchmarks in the zero-shot, one-shot, and few-shot settings
- In some cases it matched the performance of state-of-the-art fine-tuned systems
- It did well in generating high-quality samples and strong qualitative performance at tasks defined on-the-fly
- We also discussed the social impacts of this class of model

# Reference

- <https://jalammar.github.io/how-gpt3-works-visualizations-animations/>
- <https://arxiv.org/pdf/2005.14165.pdf>
- <https://medium.com/analytics-vidhya/openai-gpt-3-language-models-are-few-shot-learners-82531b3d3122>