Language Models are Few-shot Learners (GPT-3)

T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, D. Amodei, NeurIPS 2020

Digest by Samuel Albanie, July 2022

Slow description



Outline

- Motivation and overview of results
- Approach: training and evaluation
- Results
- Test set contamination
- Limitations
- Broader Impacts
- Related Work

Motivation: dropping fine-tuning

Towards minimal adaptation

NLP: use pre-trained representations in flexible/task agnostic ways for downstream



This has brought gains on many tasks (reading comprehension, question answering However, although the architecture is task-agnostic, task-specific fine-tuning is still n This requires a dataset of hundreds/thousands of examples specific to the target ta

References

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020) (word2vec) T. Mikolov et al., "Efficient estimation of word representations in vector space", arxiv (2013) (GloVe) J. Pennington et al., "Glove: Global vectors for word representation", EMNLP (2014) (LSTM RNN auto-encoders) A. Dai et al., "Semi-supervised sequence learning", NeurIPS (2015) (Transformers) A. Vaswani et al., "Attention is all you need", NeurIPS (2017) (GPT) A. Radford et al. "Improving language understanding by generative pre-training" (2018) (BERT) J. Devlin et al., "Bert: Pre-training of deep bidirectional transformers for language understanding", NAACL-HLT (2019) S. Gururangan et al., "Annotation artifacts in natural language inference data", arxiv (2018)

n transfer
er
rchitecture
chitecture
ing
g etc.)
needed
ask

The benefits of dropping fine-tuning

Practicality: the need for labelled examples limits applicability There are wide range of possible useful tasks (correcting grammar, critiquing a short story, etc.) - gathering a dataset for each is difficult **Spurious correlations**: increased propensity among more expressive models trained on narrow distributions Risk of overfitting to "annotation artifacts" (Gururangan et al., 2018) **Unhuman:** humans do not require many labelled examples A brief directive (e.g. "tell me if this sentence describes something happy or something sad") or a few examples suffices Practical advantages: humans can mix together various tasks and skills Example: performing arithmetic during a long-running dialogue

How can we remove the need for fine-tuning?



Motivation: Meta-learning at scale





Image credits/References

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

(Megatron-LM) M. Shoeybi et al., "Megatron-Im: Training multi-billion parameter language models using model parallelism", arxiv (2019) (T5) C. Raffel et al., "Exploring the limits of transfer learning with a unified text-to-text transformer", JMLR (2020) (GPT-2) A. Radford et al., "Language models are unsupervised multitask learners", OpenAI (2019) (Turing NLG) <u>https://www.microsoft.com/en-us/research/blog/turing-nlg-a-17-billion-parameter-language-model-by-microsoft/</u> (Natural Questions) T. Kwiatkowski et al., "Natural questions: a benchmark for question answering research", ACL (2019)

(GPT) A. Radford et al. "Improving language understanding by generative pre-training" (2018)

(BERT) J. Devlin et al., "Bert: Pre-training of deep bidirectional transformers for language understanding", NAACL-HLT (2019)

J. Kaplan et al., "Scaling laws for neural language models", arxiv (2020)



TLDR: testing the scale hypothesis with GPT-3

Scaling up language modelling

Test scale hypothesis: train a 175 billion parameter autoregressive language model

Evaluate in-context learning on more than 24 NLP datasets

Also evaluate on novel tasks that are unlikely to be present in the training set

Evaluate the 175 billion parameter model, named "GPT-3", under three settings:



Task: remove extraneous characters from word



Image credits/References

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

(TriviaQA) M. Joshi et al., "TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension", ACL (2017) (CoQA) S. Reddy et al. "CoQA: A conversational question answering challenge", ACL (2019)

Overview of results

GPT-3 -is promising with few examples NLP benchmarks (e.g. CoQA, TriviaQA) Also works well for tasks designed to test rapid adaptation/on-the-fly reasoning: unscrambling words, arithmetic, using words in a sentence after seeing them once GPT-3 can synthesise, few-shot, news articles that are highly human-like There are also tasks where GPT-3 struggles to perform well in a few-shot setting Natural language inference (e.g. ANLI), reading comprehension (e.g. RACE/QuAC)



(ANLI) Y. Nie et al. "Adversarial NLI: A new benchmark for natural language understanding", ACL (2020) (RACE) G. Lai et al., "RACE: Large-scale ReAding Comprehension Dataset From Examinations", EMNLP (2017) (QuAC) E. Choi et al., "QuAC: Question Answering in Context", EMNLP (2018)



Outline

- Motivation and overview of results
- Approach: training and evaluation
- Results
- Test set contamination
- Limitations
- Broader Impacts
- Related Work

Approach and evaluation settings

Approach

The pre-training approach (model, data and training) is similar to GPT-2 Relatively straightforward scaling up of model, data size/diversity, training length Like GPT-2, in-context learning is studied, but here it is done so more systematically

Fine-tuning (FT)

Fine-tuning is perhaps the most common evaluation setting for pre-trained models It involves updating the weights of a pre-trained model using a task-specific dataset Advantage:

• strong performance on many benchmarks

Disadvantages:

- requirement for a new large dataset for every task
- potential for poor generalisation
- potential to exploit spurious features of the training data (overestimate of ability)

In this work, GPT-3 is not fine-tuned (focus instead on task-agnostic performance)

However, it can be done in principle, and is a promising direction for future work

Image credits/References

(GPT-2) A. Radford et al., "Language models are unsupervised multitask learners", OpenAI (2019) T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

Fine-tuning (FT) details



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



Evaluation settings

Few-shot (FS)

Few-shot provides a few demonstrations at inference time

In general, K examples given as context (with completions), plus a query examples a query examples are as a context (with completions), plus are as a query examples are as a context (with completions), plus are as a query examples are as a context (with completions), plus are as a query examples are as a context (with completions), plus are as a query examples are as a context (with completions), plus are as a query examples are a

The model is then required to provide the completion for the query

Typically, K is set in the range of 10 to 100 (to fit in context window, $n_{ctx} = 2048$)

Advantages:

- Reduced need for task-specific data
- Reduced potential to learn from an overly-narrow distribution

Disadvantage:

• To date, results have been much weaker than fine-tuned models

Related to few-shot learning (Hochreiter et al., 2001; Vinyals et al., 2016)

One-shot (1S)

One-shot is the same as few-shot, but with a single example

One-shot learning closely matches the way some tasks are communicated to humans

It is common when describing tasks on AMT to provide one example

Sometimes it is difficult to communicate a task when no examples are given

Image credits/References

- T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)
- S. Hochreiter et al., "Learning to learn using gradient descent", ICANN (2001)
- O. Vinyals et al., "Matching networks for one shot learning", NeurIPS (2016)

Few-shot (FS) details

The model sees a task description and a few examples of the task No gradient updates are performed

	Translate English to French:	← task description
	sea otter => loutre de mer	← examples
	peppermint => menthe poivrée	<
	plush girafe => girafe peluche	<i>←</i>
	cheese =>	← prompt

One-shot (1S) details

The model sees a task description and a single example of the task No gradient updates are performed



mpl	е
-----	---

Evaluation settings

Zero-shot (ZS)

Zero-shot provides no demonstrations at test time

Advantages:

- Maximum convenience
- greater potential for robustness
- avoidance of spurious correlations (unless they occur in the pretraining corpora)

Disadvantage:

• It can be difficult to understand the task without examples ("unfairly hard")

E.g. "make a table of world records for the 200m dash" - is highly ambiguous

What format should the table have? Who exactly should be included?

However, for other tasks, zero-shot is closest to how humans communicate tasks

For example, in the English-to-French translation task, the instruction likely suffices

Image credits/References T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)



The model sees only a task description

No gradient updates are performed

1	Translate English to French:	task description
2	cheese =>	← prompt

Model and architectures

GPT-3 architectures

GPT-3 uses the same (Transformer) model and architecture as GPT-2 modified initialisation pre-normalisation Difference: Use alternating dense and locally banded sparse attention patterns in the Transformer, similar to the Sparse Transformer A total of 8 different sizes of model are considered, each was trained for 300 billion tokens

Model	params	layers	model (dim)	heads	head (dim)
GPT-3 Small	125M	12	768	12	64
GPT-3 Medium	350M	24	1024	16	64
GPT-3 Large	760M	24	1536	16	96
GPT-3 XL	1.3B	24	2048	24	128
GPT-3 2.7B	2.7B	32	2560	32	80
GPT-3 6.7B	6.7B	32	4096	32	128
GPT-3 13B	13B	40	5140	40	128
GPT-3 175B ("GPT-3")	175B	96	12288	96	128

The feedforward layer dimension is always four times the model dimension (i.e. bottleneck layer dimension)

All models use a context window of $n_{ctx} = 2048$ tokens

The model is partitioned across GPUs along both depth and width dimensions to minimise data-transfer between nodes

Hyperparameters are chosen for computational efficiency and load-balancing in the layout of models across GPUs

Kaplan et al. (2020) suggests that with enough data, validation loss is a smooth power law of model size across range of hyperparameters

Training models of different sizes allows testing of this hypothesis, for both validation loss and downstream language tasks

References

(GPT-2) A. Radford et al., "Language models are unsupervised multitask learners", OpenAI (2019) (LayerNorm) J. Ba, "Layer Normalization", arxiv (2016)

(Pre-activated ResNets) K. He et al., "Identity mappings in deep residual networks", ECCV (2016)

reversible tokenisation

(Sparse Transformer) R. Child et al., "Generating long sequences with sparse transformers", arxiv (2019) J. Kaplan et al., "Scaling laws for neural language models", arxiv (2020)

Training data

Training datasets for GPT-3

Recent trends have seen a rapid increase in dataset size used to train language models A prominent example is Common Crawl (Raffel et al., 2020) of nearly a trillion words This is sufficient to train the largest models without encountering the same sequence twice However, lightly filtered versions of Common Crawl often have lower quality than more curated datasets Three steps are taken to improve dataset quality:

- Filter \bullet Fuzzy deduplication
- Augment raw data
- based on similarity to a collection of high-quality reference data at document level (within & across datasets) - avoid redundancy/preserve integrity of validation set with high-quality reference corpora to increase diversity

The curated high-quality datasets include:

WebText2 (expanded WebText dataset from GPT-2) Books1 Books2 English Wikipedia						
	Quantity	Weight in	Epochs elapsed when	Common Crawl downloaded from 41 shards		
Dataset	(tokens)	training mix	training for 300B tokens	of monthly Common Crawl (2016-2019)		
Common Crawl (filtered)	410 billion	60%	0.44	Filter 45TB compressed plaintext to 570GB		
WebText2	19 billion	22%	2.9			
Books1	12 billion	8%	1.9			
Books2	55 billion	8%	0.43	Higher-quality datasets are sampled more		
Wikipedia	3 billion	3%	3.4	frequently (trade-off overfitting/data quality)		

Dataset contamination: with large-scale pretraining, development/test sets may be seen during training

Image credits/References

(Common Crawl data) https://commoncrawl.org/the-data/

(GPT-2) A. Radford et al., "Language models are unsupervised multitask learners", OpenAI (2019) (WebText2) J. Kaplan et al., "Scaling laws for neural language models", arxiv (2020) (Common Crawl C4 variant) C. Raffel et al., "Exploring the limits of transfer learning with a unified text-to-text transformer", JMLR (2020) T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

The authors endeavoured to remove overlaps, but some were missed due to a bug - instead the effects are characterised

Common Crawl filtering details

Filtering low quality documents

Automatic filtering was used to remove low quality documents

A classifier was trained to discriminate between:

- WebText data (a proxy for high quality documents)
- raw Common Crawl data (a proxy for low quality documents)

Common Crawl was then resampled using the trained classifier (prioritising quality)

Classifier: logistic regressor with features from Spark's tokeniser and HashingTF

Positive examples:



Each document was retained if:

document_score > 1 - np.random.pareto(α)

To attain good quality but keep some diversity after filtering, $\alpha = 0.9$

This value of α was selected to match the classifier score distribution on WebText

This re-weighting was observed to improve quality of out-of-distribution text samples

References

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

(Common Crawl data) https://commoncrawl.org/the-data/

Fuzzy deduplication

To enhance data quality and reduce over-fitting, fuzzy deduplication was applied Documents exhibiting high overlap with other documents were pruned The implementation used with Spark MinHashLSH using 10 hashes The same features were used as for the filtering classifier Fuzzy deduplication was also used to remove WebText from Common Crawl Fuzzy deduplication reduced the average dataset size by 10%



⁽HashingTF) https://spark.apache.org/docs/latest/api/python/reference/api/pyspark.ml.feature.HashingTF.html (MinHashLSH) https://spark.apache.org/docs/3.1.1/api/python/reference/api/pyspark.ml.feature.MinHashLSH.html

Training process

GPT-3 training process

Key observation from prior work (McCandlish et al., 2018; Kaplan et al., 2020): Larger models can use a larger batch size, but need a smaller learning rate Gradient noise scale (signal-to-noise ratio of gradient across training examples) This statistic is used to guide the choice of batch size for each model

Model	params	batch size	learning rate
GPT-3 Small	125M	0.5M	6.0×10^{-4}
GPT-3 Medium	350M	0.5M	3.0×10^{-4}
GPT-3 Large	760M	0.5M	2.5×10^{-4}
GPT-3 XL	1.3B	1 M	2.0×10^{-4}
GPT-3 2.7B	2.7B	1 M	1.6×10^{-4}
GPT-3 6.7B	6.7B	2M	1.2×10^{-4}
GPT-3 13B	13B	2M	1.0×10^{-4}
GPT-3 175B ("GPT-3")	175B	3.2M	0.6×10^{-4}

To avoid memory issues, a multiple forms of model parallelism are used:

- within each matrix multiply
- across the layers of the network

Training is performed on a high-bandwidth Microsoft cluster of V100 GPUs

References

(Gradient Noise) S. McCandlish et al., "An empirical model of large-batch training", arxiv (2018) J. Kaplan et al., "Scaling laws for neural language models", arxiv (2020)

Low-level optimisation details

GPT-3 is trained using Adam ($\beta_1 = 0.9$, $\beta_2 = 0.95$, $\epsilon = 10^{-8}$)

The global norm of the gradient is clipped at 1.0

A linear learning rate warmup is used over the first 375 million tokens

Cosine decay applied to the learning rate down to 10% over 260 billion tokens

The batch size is gradually linearly increased over the first 4 - 12 billion tokens

Starts at 32k tokens and increases to full size (which depends on the model)

Data are sampled without replacement within each epoch

All models use decoupled weight decay of 0.1

In training, the full $n_{\text{ctx}} = 2048$ context window is filled

Achieved by packing in multiple documents into a sequence (for compute efficiency)

Sequences with multiple documents are not masked in any special way

Instead, they are separated with a special end of text token (to inform the model)

This enables efficient training without sequence-specific masking



Evaluation details

Few-shot learning

In few-shot learning evaluations, each example is evaluated by sampling K examples Examples are delimited by 1 or 2 newlines (depending on the task) LAMBADA/Storycloze have no supervised training set (condition on dev set) For Winograd (original, rather than SuperGLUE), there is only one dataset K is chosen between 0 and the largest value that fits in the context window Typically, larger values of K are better, but not always: when a development set is available, it is used to select the best value of K

In certain tasks, a natural language prompt is also provided to the model

Multiple choice tasks

For multiple choice tasks, condition on K complete examples and compare the likelihood of different completions for a query example

For most tasks, per-token likelihood is used (to normalise for length)

For some datasets (ARC, OpenBookQA, RACE) it was useful to further normalise:

P(completion | context)

P(completion answer_context)

where answer_context is "Answer: " or "A: "

References

(LAMBADA) D. Paperno et al., "The LAMBADA dataset: Word prediction requiring a broad discourse context", ACL (2016) (Storycloze) N. Mostafazadeh et al., "A corpus and cloze evaluation for deeper understanding of commonsense stories", NAACL HLT (2016) (Winograd) H. Levesque et al., "The Winograd schema challenge", KR (2012)

Binary classification

For binary classification, semantically meaningful names are used for options:

"True" rather than "1" and "False" rather than "0"

Then the multiple choice protocol is followed

Free-form completion

For free-form completion, beam search is used following Raffel et al. (2020):

- beam width of 4
- length penalty of $\alpha = 0.6$

The model is scored with F1 similarity, BLEU, or exact match, depending on dataset

Reporting test results

When publicly available, final results are reported on the test set

For private test sets, GPT-3 is often too large to fit on the test server (dev set is used)

The test server is used where possible (SuperGLUE, TriviaQA, PiQA)

(ARC) P. Clark et al., "Think you have solved question answering? try arc, the ai2 reasoning challenge" arxiv (2018) (OpenBookQA) T. Mihaylov et al., "Can a Suit of Armor Conduct Electricity? A New Dataset for Open Book Question Answering", EMNLP (2018) (RACE) G. Lai et al., "RACE: Large-scale ReAding Comprehension Dataset From Examinations", EMNLP (2017) C. Raffel et al., "Exploring the limits of transfer learning with a unified text-to-text transformer", JMLR (2020)

(SuperGLUE) A. Wang et al., "SuperGLUE: A stickier benchmark for general-purpose language understanding systems", NeurIPS (2019) (TriviaQA) M. Joshi et al., "TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension", ACL (2017)

(PiQA) Y. Bisk et al., "PiQA: Reasoning about physical commonsense in natural language", AAAI (2020)





Outline

- Motivation and overview of results
- Approach: training and evaluation
- Results
- Test set contamination
- Limitations
- Broader Impacts
- Related Work

Language modeling

Language modelling (pretraining)

Report val curves during training for 8 models (& 6 smaller models down to 100K params)



Image credits/References

- T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)
- J. Kaplan et al., "Scaling laws for neural language models", arxiv (2020)

Parameters

Notes

Performance follows power law with compute (Kaplan et al., 2020) Over two more orders of magnitude, departure from the law is minimal A natural concern: gains come from modelling spurious details However, cross entropy improvements map to downstream gains



Language modelling/clozes

Penn Tree Bank - Language modelling

Evaluate zero-shot perplexity on Penn Tree Bank, as in Raffel et al. (2020)

Wikipedia-related tasks are omitted due to overlap

Penn Tree Bank is not affected (it predates the modern internet)

Only zero-shot is evaluated (since there is no separation of examples)

Setting	PTB	
SOTA (Zero-Shot)	35.8	GPT-2
GPT-3 Zero-Shot	20.5	

LAMBADA

LAMBADA requires predicting a word that depends on a paragraph of context The aim of the benchmark is to test modelling of long-range dependencies in text It has been suggested that "continuing to expand hardware and data sizes by orders of magnitude is not the path forward" (Bisk et al., 2020)

Setting	LAMBADA (acc)	LAMBADA (ppl)	
SOTA GPT-3 Zero-Shot	68.0 Tur 76.2	ring 8.63 GPT-2	Scaling continues to yield gains
GPT-3 One-Shot GPT-3 Few-Shot	86.4	<u>3.35</u> 1.92	

Image credits/References

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

(PTB) M. Marcus et al., "The Penn treebank: Annotating predicate argument structure", HLT Workshop (1994) C. Raffel et al., "Exploring the limits of transfer learning with a unified text-to-text transformer", JMLR (2020)

LAMBADA

LAMBADA also demonstrates the benefits of few-shot learning

Problem: the completion is always the last word, but the model does not know this It assigns probability to not only one word completions, abut also other completions One solution: stop-word filters (ban "continuation words") (Radford et al., 2019) With few-shot learning, the task can be framed as a cloze test requiring 1 word:

Alice was friends with Bob. Alice went to visit her friend _____. \rightarrow Bob

George bought some baseball equipment, a ball, a glove and a _____. ightarrow



One-shot learning is not effective (possibly more examples required)

(GPT-2) A. Radford et al., "Language models are unsupervised multitask learners", OpenAI (2019) (LAMBADA) D. Paperno et al., "The LAMBADA dataset: Word prediction requiring a broad discourse context", ACL (2016)

Y. Bisk et al. "Experience grounds language", arxiv (2020) (Turing NLG) <u>https://www.microsoft.com/en-us/research/blog/turing-nlg-a-17-billion-parameter-language-model-by-microsoft/</u>



Completion tasks

HellaSwag

HellaSwag tasks models with selecting the best ending to a story/instruction set Dataset samples were adversarially mined to be hard for machines

These samples are relatively easy for humans (95.6% accuracy)

Setting	HellaSwag (acc)	
SOTA GPT-3 Zero-Shot GPT-3 One-Shot GPT-3 Few-Shot	85.6 ← 78.9 78.1 79.3	ALUM (RoBERTa-L)

Image credits/References

(HellaSwag) R. Zellers et al., "HellaSwag: Can a Machine Really Finish Your Sentence?" ACL (2019) T. Brown et al., "Language models are few-shot learners", NeurIPS (2020) (ALUM) X. Liu et al., "Adversarial training for large neural language models", arxiv (2020) (RoBERTa) Y. Liu et al., "RoBERTa: A robustly optimized bert pretraining approach", arxiv (2019) (Storycloze) N. Mostafazadeh et al., "A corpus and cloze evaluation for deeper understanding of commonsense stories", NAACL HLT (2016) (TransBERT) Z. Li, "Story ending prediction by transferable BERT", arxiv (2019)

StoryCloze

Storycloze tasks models with selecting the end sentence for 5-sentence stories

	StoryCloz	ze	
Setting	(acc)		
SOTA	91.8	←	TransBERT
GPT-3 Zero-Shot	83.2		
GPT-3 One-Shot	84.7		
GPT-3 Few-Shot	87.7		

Note: GPT-3 is approximately 10% better than previous zero-shot results

Closed Book Question Answering

Closed Book Question Answering

Goal: assess GPT-3's ability to answer questions about factual knowledge

Typically, this has been "open-book" - with a Information Retrieval (IR) system

Roberts et al. (2020) show T5 can also answer "closed-book" (without IR)

Their work suggests that higher capacity models would improve performance

GPT-3 is evaluated on NaturalQS

TriviaQA

Note: zero/one/few-shot are stricter than prior work (fine-tunes on Q&A data)

WebQS

Setting	NaturalQS	WebQS	TriviaQA
RAG (Fine-tuned, Open-Domain)	44.5	45.5	68.0
T5-11B+SSM (Fine-tuned, Closed-Book)	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	68.0
GPT-3 Few-Shot	29.9	41.5	71.2
	^	≜	
	slightly out-of-distribution		

Image credits/References:

A. Roberts, "How much knowledge can you pack into the parameters of a language model?", arxiv (2020) (T5) C. Raffel et al., "Exploring the limits of transfer learning with a unified text-to-text transformer", JMLR (2020) (GPT-3) T. Brown et al., "Language models are few-shot learners", NeurIPS (2020) (Natural Questions) T. Kwiatkowski et al., "Natural questions: a benchmark for question answering research", ACL (2019) (WebQS) J. Berant et al., "Semantic parsing on freebase from question-answer pairs", EMNLP (2013) (TriviaQA) M. Joshi et al., "TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension", ACL (2017) (RAG) P. Lewis et al., "Retrieval-augmented generation for knowledge-intensive nlp tasks", NeurIPS (2020)



Translation

Language translation

When training GPT-2, data was filtered to English-only due to concerns about capacity GPT-2 still showed non-trivial Fr-En translation, despite seeing only 10 MB of French text GPT-3 data increases other languages, but still mainly English (CommonCrawl source) By word count, the training data is 93% English, 7% other languages Existing unsupervised translation methods often use back-translation (Seinrich et al., 2015)

GPT-3 learns from blended training data (mixing words, sentences, documents)

Note: one-shot/few-shot GPT-3 is not directly comparable to unsupervised translation

	WMT'1	4 Fr⇔En	WMT'16	6 De⇔En	WMT'1	5 Ro←
Setting	En→Fr	Fr→En	En→De	De→En	En→Ro	Ro-
SOTA (Supervised)	45.6 BT@	S 35.0 UED	NN 41.2 MA	DL 40.2 Uol	R 38.5 ALB	ert 39
XLM MASS mBART	33.4 <u>37.5</u>	33.3 34.9	26.4 28.3 <u>29.8</u>	34.3 35.2 34.0	33.3 <u>35.2</u> 35.0	31 33 30
GPT-3 Zero-Shot GPT-3 One-Shot GPT-3 Few-Shot	25.2 28.3 32.6	21.2 33.7 <u>39.2</u>	24.6 26.2 29.7	27.2 30.4 <u>40.6</u>	14.1 20.6 21.0	19 38 <u>39</u>
		M	aybe SOTA		N	lear (

Image credits/References:

(GPT-2) A. Radford et al., "Language models are unsupervised multitask learners", OpenAI (2019) (Back-translation) R. Sennrich et al., "Improving neural machine translation models with monolingual data", (2015) (GPT-3) T. Brown et al., "Language models are few-shot learners", NeurIPS (2020) (BT@S) S. Edunov et al., "Understanding Back-Translation at Scale", EMNLP (2018) (UEDIN) N. Durrani et al., "Edinburgh's phrase-based machine translation systems for WMT-14", WMT (2014)



(MADL) Y. Wang et al., "Multi-agent dual learning", ICLR (2019) (UoR) University of Regensburg, Fascha, (2016) (ALBERT) Z. Lan et al., "ALBERT: A lite BERT for self-supervised learning of language representations", arxiv (2019) (XLM) A. Conneau et al., "Cross-lingual language model pretraining", NeurIPS (2019) (MASS) K. Song et al., "MASS: Masked sequence to sequence pre-training for language generation", ICML (2019) (mBART) Y. Liu et al., "Multilingual denoising pre-training for neural machine translation", ACL (2020)



Winograd-Style Tasks

Winograd Schemas

Winograd Schemas are a task that require determining which word a pronoun refers to

The pronoun is grammatically ambiguous but semantically unambiguous to a human

"The trophy doesn't fit in the brown suitcase because it's too big. What is too big?" small small

Answer 0: the trophy

Answer 1: the suitcase

Evaluate on original 273 Winograd schemas with "partial evaluation" (Trinh et al., 2018)

Note: this setting is slightly different to the WSC task in SuperGLUE (entity extraction)

Performance is also evaluated on the (adversarially mined) Winogrande dataset

Setting	Winograd	Winogrande	e (XL)
Fine-tuned SOTA	90.1 Ro	BERTa+W 84.6	T5 + NLI
GPT-3 Zero-Shot	88.3*	70.2	
GPT-3 One-Shot	89.7*	73.2	
GPT-3 Few-Shot	88.6*	77.7	

* Some (small) contamination of the test set was discovered

Image credits/References:

(Winograd) H. Levesque et al., "The Winograd schema challenge", KR (2012)

T. Trinh et al., "A simple method for commonsense reasoning", arxiv (2018)

(SuperGLUE) A. Wang et al., "SuperGLUE: A stickier benchmark for general-purpose language understanding systems", NeurIPS (2019) (Winogrande) K. Sakaguchi et al., "Winogrande: An adversarial winograd schema challenge at scale", Communications of the ACM (2021) T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

(RoBERTa) Y. Liu et al., "RoBERTa: A robustly optimized bert pretraining approach", arxiv (2019)

(T5+NLI) S-C. Lin et al., "Tttttackling winogrande schemas", arxiv (2020)





Common Sense Reasoning

Common Sense Reasoning tasks

Several datasets assess physical or scientific reasoning:

- Physical Interaction Question Answering (PIQA) about how the physical world works
- Al2 Reasoning Challenge (ARC) 3rd to 9th grade science exams
- Open Book Question Answering (OpenBookQA) generated from science facts

Setting	PIQA ARC (Easy)	ARC (Challenge)	OpenBookQA
Fine-tuned SOTA	79.4 RoB 92.0 UnifiedQA 80.5* 68.8 80.5* 71.2 82.8* 70.1	78.5 UnifiedQA	87.2 UnifiedQA
GPT-3 Zero-Shot		51.4	57.6
GPT-3 One-Shot		53.2	58.8
GPT-3 Few-Shot		51.5	65.4

* Potential dataset contamination was discovered

Image credits/References:

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020) (PiQA) Y. Bisk et al., "PiQA: Reasoning about physical commonsense in natural language", AAAI (2020) (ARC) P. Clark et al., "Think you have solved question answering? try arc, the ai2 reasoning challenge" arxiv (2018) (OpenBookQA) T. Mihaylov et al., "Can a Suit of Armor Conduct Electricity? A New Dataset for Open Book Question Answering", EMNLP (2018) (RoB) Y. Liu et al., "RoBERTa: A robustly optimized bert pretraining approach", arxiv (2019) (UnifiedQA) D. Khashabi et al. "UnifiedQA: Crossing Format Boundaries with a Single QA System", EMNLP (2020)

Reading Comprehension

Reading comprehension tasks

A suite of five datasets is used to assess GPT-3 for reading comprehension covering:

questions: single dialog

abstractive multiple choice span based answers:

- Conversational Question Answering (CoQA) 127K conversation turns
- Discrete Reasoning Over Paragraphs (DROP) 96K reasoning questions on paragraphs
- Question Answering in Context (QuAC) 14K information-seeking QA dialogs
- SQuAD with adversaRial Unanswerable Questions (SQuADv2) 50K extra questions
- ReAding Comprehension Dataset From Examinations (RACE) 100K exam questions

Setting	CoQA	DROP	QuAC	SQuADv2	RACE-h	RACE
Fine-tuned SOTA	90.7 RA	K 89.1 NT	74.4 TR	93.0 SA-Net	90.0 Mega	93.1 M
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
performance Few-shot gain	solid medium	weak big	weak small	medium big	weak smal	<

Image credits/References:

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

(RAK) Y. Ju et al. "Technical report on conversational question answering", arxiv (2019) (NT) J. Zheng, "Numeric Transformer - ALBERT", AI2 leaderboard (2020) (CoQA) S. Reddy et al. "CoQA: A conversational question answering challenge", ACL (2019) (TR) TR-MT (ensemble) WeChat AI, QuAC leaderboard (2019) (DROP) D. Dua et al., "DROP: A Reading Comprehension Benchmark Requiring Discrete Reasoning Over Paragraphs", NAACL-HLT (2019) (SA-Net) Qianxin, "SA-Net on ALBERT (ensemble)", SQuAD 2.0 leaderboard (2020) (QuAC) E. Choi et al., "QuAC: Question Answering in Context", EMNLP (2018) (Megatron-LM) M. Shoeybi et al., "Megatron-Im: Training multi-billion parameter language models using model (SQuADv2) P. Rajpurkar et al., "Know What You Don't Know: Unanswerable Questions for SQuAD", ACL (2018) parallelism", arxiv (2019)

(RACE) G. Lai et al., "RACE: Large-scale ReAding Comprehension Dataset From Examinations", EMNLP (2017)

SuperGLUE

The SuperGLUE benchmark

GPT-3 is evaluated on the SuperGLUE benchmark to assess a range of NLP abilities:

- Boolean Questions (BoolQ) naturally occurring yes/no questions about text
- CommitmentBank (CB) embedded clauses annotated with degree of commitment
- Choice of Plausible Alternatives (COPA) cause/effect binary classification
- Multi-Sentence Reading Comprehension (MultiRC) paragraph question answering
- Reading Comprehension with Commonsense Reasoning Dataset (ReCoRD) clozes
- Recognising Textual Entailment (RTE) RTE (1,2,3,5), as used in GLUE
- Word-in-Context (WiC) binary word sense disambiguation
- Winograd Schema Challenge (WSC) coreference resolution as binary classification
- For few-shot evaluations, 32 examples are used (sampled from training set) for all tasks
- Different examples are used as context for each problem (except WSC and MultiRC)

Image credits/References:

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

(ReCoRD) S. Zhang et al., "Record: Bridging the gap between human and machine commonsense reading comprehension", arxiv (2018) (GLUE) A. Wang et al., "GLUE: A multi-task benchmark and analysis platform for natural language understanding", ICLR (2018) (SuperGLUE) A. Wang et al., "SuperGLUE: A stickier benchmark for general-purpose language understanding systems", NeurIPS (2019) (WiC) M. Pilehvar, "WiC: the Word-in-Context Dataset for Evaluating Context-Sensitive Meaning Representations", NAACL-HLT (2019) (BoolQ) C. Clark et al., "BoolQ: Exploring the Surprising Difficulty of Natural Yes/No Questions", NAACL-HLT (2019) (WSC) H. Levesque, "The Winograd Schema Challenge." PKRR (2012) (CB) M-C. De Marneffe et al., "The commitmentbank: Investigating projection in naturally occurring discourse", Sinn und Bedeutung (2019) (COPA) M. Roemmele et al., "Choice of Plausible Alternatives: An Evaluation of Commonsense Causal Reasoning", AAAI symposium (2011)

		Resu	ts			
		near BEI	RT "sign:	s of life"	near SOTA	nea
	SuperGLUE Average	E BoolQ Accurac	CB y Accura	CB cy F1	COPA Accuracy	R' Acci
Fine-tuned SOTA Fine-tuned BERT-Large GPT-3 Few-Shot	89.0 69.0 71.8	91.0 77.4 76.4	96.9 83.6 75.6	93.9 75.7 52.0	94.8 70.6 92.0	92 71 69
	WiC Accuracy	WSC Accuracy	MultiRC Accuracy	MultiRC F1a	ReCoRD Accuracy	Re
Fine-tuned SOTA Fine-tuned BERT-Large GPT-3 Few-Shot	76.1 69.6 49.4	93.8 64.6 80.1	62.3 24.1 30.5	88.2 70.0 75.4	92.5 71.3 90.2	9 7 9
	chance	solid	near	BERT	near	SOT

Efforts to rephrase/reformulate WiC to improve performance were unsuccessful One-shot/few-shot GPT-3 seems to be weak at comparing two sentences/snippets: if a word is used in same sense (WiC), if a sentence paraphrases or implies another This is also reflected in relatively weak performance for RTE and CB

Summary stats:

near SOTA on 2 tasks

stronger than BERT-large on 4 tasks

(MultiRC) D. Khashabi et al., "Looking beyond the surface: A challenge set for reading comprehension over multiple sentences", NAACL-HLT (2018)

SuperGLUE

Image credits/References:

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020) (SuperGLUE) A. Wang et al., "SuperGLUE: A stickier benchmark for general-purpose language understanding systems", NeurIPS (2019)

Natural Language Inference

The Natural Language Inference (NLI) task

Natural Language Inference (NLI) assesses inferences between two sentences Standard formulation involves two-way or three-way classification for whether: The second sentence logically follows contradicts is neutral to the first On RTE (SuperGLUE), only GPT-3 175B was much better than random However, this (few-shot) model was comparable to single-task fine-tuned BERT large Evaluation also conducted on ANLI (three rounds of adversarially mined questions)

Image credits/References:

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020) (SuperGLUE) A. Wang et al., "SuperGLUE: A stickier benchmark for general-purpose language understanding systems", NeurIPS (2019) (BERT) J. Devlin et al., "Bert: Pre-training of deep bidirectional transformers for language understanding", NAACL-HLT (2019) (ANLI) Y. Nie et al., "Adversarial NLI: A New Benchmark for Natural Language Understanding", ACL (2020)

Influence of model scale

Results for ANLI round 3

All models smaller than GPT-3 175B perform at chance level

However, few-shot GPT-3 175B shows promise

Few-shot Arithmetic

Synthetic arithmetic tasks

Synthetic tasks assess on-the-fly reasoning, recognise novel patterns and adapt quickly to unusual tasks

10 tasks were devised to assess GPT-3's arithmetic ability with task-specific training:

Name	Task to be performed	Examp
2 digit addition (2D+)	Add 2 integers $\sim [0,100)$	"Q: What is 34 plus
3 digit addition (3D+)	Add 2 integers $\sim [0,1000)$	
4 digit addition (4D+)	Add 2 integers $\sim [0,10000)$	
5 digit addition (5D+)	Add 2 integers $\sim [0,100000)$	
2 digit subtraction (2D-)	Subtract 2 integers $\sim [0,100)$	"Q: What is 34 minu
3 digit subtraction (3D-)	Subtract 2 integers $\sim [0,1000)$	
4 digit subtraction (4D-)	Subtract 2 integers $\sim [0,10000)$	
5 digit subtraction (5D-)	Subtract 2 integers $\sim [0,10000)$	
2 digit multiplication (2Dx)	Multiply 2 integers $\sim [0,100)$	"Q: What is 24 time
1 digit composite (1DC)	Sample 3 integers $\sim [0,10)$	"Q: What is 6+(4*8
loorform	, composite operation - operators a	$mona \{ \pm - \times \} \}$

For each task, 2000 random instances of the task are generated - answers must be exact to be correct

(perform composite operation operators among (1, , , ,))

Image credits/References:

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

Arithmetic

		Arith	imetic re	sults for	- GPT-3	175B				
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

Zero-shot GPT-3 is still reasonable (better than tew-shot GPT-13B) **Memorisation check:** search 3 digit problems "<NUM1> + <NUM2> = " and "<NUM1> plus <NUM2> = " in corpus There were 17 matches (out of 2000) i.e. 0.8% for addition and 2 matches i.e. 0.1% for subtraction Memorisation is therefore unlikely to explain the results GPT-3 sometimes makes mistakes such as not carrying a "1" - suggesting it is attempting computation (rather than lookup) **Takeaway:** GPT-3 demonstrates some ability at arithmetic, particularly in the few-shot setting

Image credits/References: T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

Word Scrambling and Manipulation Tasks

Symbol manipulation tasks

Character manipulation tasks assess novel symbolic manipulation from few examples

There are five tasks, each involving distorting a word and requiring its recovery

Each task provides a prompt of the form "\n<distorted word> = "

Cycle letters in word (CL)

E.g. lyinevitab \rightarrow inevitably

<distorted word> is a word with letters cycled

Anagrams of all but first and last characters (A1)

E.g. criroptuon \rightarrow corruption

<distorted word> has every letter except first and last randomly scrambled

Anagrams of all but first and last 2 characters (A2)

E.g. opoepnnt \rightarrow opponent

<distorted word> has every letter except first two/last two randomly scrambled

Random insertion in word (RI)

E.g. s.u!c/c!e.s s i/o/n \rightarrow succession

<distorted_word> has random punctuation/space inserted between each letter

Reversed words (RW)

E.g. stcejbo \rightarrow objects

<distorted_word> is spelt backwards

The top 10,000 most frequent words (Norvig, 2009) with length $\in (4,15)$ are used to

generate 10,000 examples for each task

Note: tasks CL, A1 and A2 are not bijective, requiring some form of search to solve

Image credits/References:

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

Performance also increases rapidly with additional examples for the model Solving the task therefore requires pulling apart the substructure of BPE tokens

P. Norvig, "Natural language corpus data", Beautiful data (2009)

SAT Analogies

The SAT analogies task

SAT analogy problems are a task that is unusual relative to typical text distribution

Example: audaciousness is to boldness as:

sanctimonious is to hypocrisy

anonymous is to identity

remorseful is to misdeed

deleterious is to result impressionable is to temptation

GPT-3 is tested on 374 SAT analogy problems (Turney et al., 2003)

Average college-bound senior high student would score $\approx 57\%$ (Turney et al., 2005)

Image credits/References:

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

P. Turney et al., "Combining independent modules to solve multiple-choice synonym and analogy problems", arxiv (2003)

P. Turney et al., "Corpus-based learning of analogies and semantic relations", Machine Learning (2005)

Results: News Article Generation

News Articles

Prior work evaluated the ability of GPT-2 to synthesise news articles via prompting For GPT-2, the model was given a plausible first sentence for a news story Compared to GPT-2, the GPT-3 dataset is less heavily skewed towards news Generating articles from headlines is less effective: GPT-3 will often the first sentence as a tweet, then generates replies/follow-up tweets

This is tackled with few-shot prompting (three news articles are used in the prompt) Given a title and subtitle of a proposed article, GPT-3 is able to generate short articles

Assessing article quality

To assess article generation, human ability to distinguish fake articles is measured This approach follows prior work (Zellers et al., 2019; Kreps et al., 2020) GPT-3 aims to match the distribution of human generated content so this assesses quality 25 article titles and subtitles were selected from <u>newser.com</u> (avg. length: 215 words) Articles were generated with 4 GPT-3 models (125M to 175B) (avg. length: 200 words) For each model, 80 US-based participants were shown real or generated articles very likely written by a human most likely written by a human I don't know

most likely written by a machine very likely written by a machine

Image credits/References

(GPT-2) A. Radford et al., "Language models are unsupervised multitask learners", OpenAI (2019) S. Kreps et al., "All the news that's fit to fabricate: Al-generated text as a tool of media misinformation" (2020) (Grover) R. Zellers et al., "Defending against neural fake news", NeurIPS (2019)

Details

Articles were not in the training data

Generated articles were formatted/selected automatically to minimise cherry picking

The same context (of three previous articles, title and subtitle) was used for each model

Control for attention/effort with "intentionally bad" baseline (160M params, no context)

Results						
	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 (2 <i>e</i> -4)	4.9%		
GPT-3 Medium	61%	58%-65%	10.3 (7 <i>e</i> -21)	6.0%		
GPT-3 Large	68%	64%-72%	7.3 (3 <i>e</i> -11)	8.7%		
GPT-3 XL	62%	59%-65%	10.7 (1 <i>e</i> -19)	7.5%		
GPT-3 2.7B	62%	58%-65%	10.4 (5 <i>e</i> -19)	7.1%		
GPT-3 6.7B	60%	56%-63%	11.2 (3 <i>e</i> -21)	6.2%		
GPT-3 13B	55%	52%-58%	15.3 (1 <i>e</i> -32)	7.1%		
GPT-3 175B	52%	49%-54%	16.9 (1 <i>e</i> -34)	7.8%		
Takeaway: larger models Prior research (Ippolito et a	produce news I., 2019) sugg	s articles that ar ests that automo	e harder to dist atic methods like	inguish e Grover and		
GLTR may be better at deter	cting generate	d text than hum	ans			
Such systems may be a pror	nising directio	n for <mark>future rese</mark>	earch			

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

D. Ippolito et al., "Automatic detection of generated text is easiest when humans are fooled", ACL (2019) (GLTR) S. Gehrmann et al. "GLTR: Statistical detection and visualization of generated text", ACL (2019)

Further Analysis of News Article Generation

Image credits/References

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

Example generated articles

GPT-3 generated article (human accuracy: 12%)

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.

Reference

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

GPT-3 generated article (human accuracy: 61%)

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."

Notes

Factual inaccuracies can indicate that the text is generated (model has no access to facts)

Other indicators: repetition, non-sequiturs, unusual phrasings (often subtle)

Longer news articles

Human assessments - longer articles

Human accuracy improves as they observe more generated text (Ippolito et al., 2019)

Generated article detection experiments are conducted with ≈ 80 US-based participants for GPT-3 and control model

	Mean accuracy	95% Confidence Interval (low, hi)	t compared to control (p -value)	"I don't know" assignments
Control	88%	84%-91%	-	2.7%
GPT-3 175B	52%	48%-57%	12.7 (3.2e-23)	10.6%

Takeaway: even for news articles with 500 words, humans struggle to differentiate between GPT-3 and human writing

Image credits/References

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

D. Ippolito et al., "Automatic detection of generated text is easiest when humans are fooled", ACL (2019)

12 articles chosen from Reuters (avg. 569 words), GPT-3 generated articles (avg. 498 words) \approx 300 words longer

Learning and Using Novel Words

Word acquisition

How does a child come to acquire a new word?

Process studied in developmental linguistics (Carey et al., 1978)

Approach to testing GPT-3's ability to learn and use new words:

- provide definition of non-existent word e.g. "Gigamuru"
- ask for an example of usage

Implementation of the task:

- 1-5 sentences of non-existent words are defined then used
- then a new word is defined

This problem formulation is:

- few-shot for the task (examples of the task provided)
- one-shot for the new word (only sees word once)

Image credits/References

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

S. Carey et al., "Acquiring a single new word", ERIC (1978)

A "Burringo" is a car with very fast acceleration. An example of a sentence that uses the word Burringo is: In our garage we have a Burringo that my father drives to work every day.

word Gigamuru is:

screeg is:

Each of the GPT-3 completions appear to be a plausible use of the new word

GPT-3 also offers a reasonable conjugation for "screeg" ("screeghed"), though the usage is a little odd

Takeaway: GPT-3 appears to be proficient at using novel words in a sentence

Results

A "whatpu" is a small, furry animal native to Tanzania. An example of a sentence that uses the word whatpu is:

We were traveling in Africa and we saw these very cute whatpus.

To do a "farduddle" means to jump up and down really fast. An example of a sentence that uses human-written prompt the word farduddle is:

One day when I was playing tag with my little sister, she got really excited and she started doing these crazy farduddles.

A "yalubalu" is a type of vegetable that looks like a big pumpkin. An example of a sentence that uses the word yalubalu is:

I was on a trip to Africa and I tried this yalubalu vegetable that was grown in a garden there. It was delicious.

A "Gigamuru" is a type of Japanese musical instrument. An example of a sentence that uses the

I have a Gigamuru that my uncle gave me as a gift. I love to play it at home.

To "screeg" something is to swing a sword at it. An example of a sentence that uses the word

We screeghed at each other for several minutes and then we went outside and ate ice cream.

Grammar

English grammar correction

GPT-3 is tasked with correcting English grammar (few-shot) It is provided with a prompt containing poor/good grammar **Note:** distinction between good/poor English is complex Poor English input:I eated the purple berries.communicating the taskGood 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 yourequested.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 yourequested.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 Engl: Good Engl

Poor English input: Today I have went to the store 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.

Image credits/References T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

Results

ish	input:	Mary a	rranged	for, i	n St.	Louis,	John to	rent a	house	cheap.	remove
lish	output:	: Mary	arrang	ed for	John	to rent	a hous	e in St	. Lou	is.	alters m

Takeaway: GPT-3 shows promising grammar correction abilities

Poem generation

Poem generation

Prompt GPT-3 write poetry

Nucleus sampling is used (P = 0.9)

Context

The City BY C. P. CAVAFY TRANSLATED BY EDMUND KEELEY [Poem text omitted] SOME TREES John Ashbery [Poem text omitted]

Shadows on the Way poem title Wallace Stevens poet

Poem 1

I must have shadows on the way If I am to walk I must have Each step taken slowly and alone To have it ready made

And I must think in lines of grey To have dim thoughts to be my guide Must look on blue and green And never let my eye forget That color is my friend And purple must surround me too

The yellow of the sun is no more Intrusive than the bluish snow That falls on all of us. I must have Grey thoughts and blue thoughts walk with me If I am to go away at all.

References

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020) (Nucleus sampling) A. Holtzman et al., "The Curious Case of Neural Text Degeneration", ICLR (2019)

Poem 2

The sun was all we had. Now, in the shade All is changed. The mind must dwell on those White fields, that to its eyes were always old; Those ancient gleams, convoluting The way they lay among the huge roots, The great dark tomes of reverie, The plumed labyrinth of the sea.

When our own breath gonne faint and fade Up the long stairway of the dead We shall remember well The blindness of the light, the whiteness Of that white land. The footsteps, and The grasses in that light, and in the shade The menacing life.

We played, too, once, in the warmest rooms. To one content, there is one content For all mankind. The forest glades Are all the more for being in fable, And wood and town in legend shadowed- Once more to return. The one content Has in it every happiness, more brave For being purposeless; and so Those island gardens, deep in the mind, Make of remoteness familiarity And of the unspoken name a friend.

Outline

- Motivation and overview of results
- Approach: training and evaluation
- Results
- Test set contamination
- Limitations
- Broader Impacts
- Related Work

Jts

Measuring and Preventing Memorisation of Benchmarks

Contamination

Pretraining data is sourced from the internet - possibly contaminating test sets Detecting test set contamination is a new research area, without best practices As pretraining increases in size, issue may become increasingly important Trinh et al. (2018), who trained on Common Crawl, detected/removed overlap For GPT-2, a post-hoc analysis of the effect of overlap was conducted Study found contamination did not significantly affect results (small overlap)

Image credits/References

T. Trinh et al., "A simple method for commonsense reasoning", arxiv (2018)
(GPT-2) A. Radford et al., "Language models are unsupervised multitask learners", OpenAI (2019)
T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

Contamination

GPT-3 is trained in a slightly different regime to GPT-2:

- Two orders of magnitude increase in data & model size (more contamination/memorisation risk)
- Due to the large dataset, GPT-3 175B does not significantly overfit its training set

GPT-3 Training Curves

Contamination is likely to be frequent, but the effects may not be large

Measuring and Preventing Memorisation of Benchmarks

Methodology

The initial approach was to deduplicate the training and test sets Due to a bug, detected overlaps were only partially removed from training Since training the model is expensive, it was not possible to retrain Instead, the effect of the overlap on results is investigated A "clean" version of each benchmark is built by removing possible overlaps: 13-gram overlap with anything in pretraining (or shorter complete matches) Goal: conservatively flag possible contamination (to be sure of clean test set) GPT-3 is then evaluated on both the clean and original benchmarks:

- If scores on clean/original test sets are similar, the effect is not significant
- If clean score is lower, contamination may be effecting the reported results

Image credits/References

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020) (PiQA) Y. Bisk et al., "PiQA: Reasoning about physical commonsense in natural language", AAAI (2020) (Winograd) H. Levesque et al., "The Winograd schema challenge", KR (2012) (LAMBADA) D. Paperno et al., "The LAMBADA dataset: Word prediction requiring a broad discourse context", ACL (2016) (CLIP) A. Radford et al., "Learning transferable visual models from natural language supervision", ICML (2021)

Further details for overlap analysis

Methodology details

GPT-2 filtering used bloom filters to estimate probabilistic bounds on test set contamination Overlaps computed between the test sets and the full training corpus Bug caused filtering with N-gram overlaps to fail on long documents (e.g. books) These benchmarks were **removed** from the results

References

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020) (GPT-2) A. Radford et al., "Language models are unsupervised multitask learners", OpenAI (2019)

- By contrast, this work used Apache Spark to compute exact collisions across train and test sets
- This is conservative, since only 40% of the filtered Common Crawl data is actually trained on
- Some language modelling benchmarks and the Children's Book Test showed nearly complete overlap

Outline

- Motivation and overview of results
- Approach: training and evaluation
- Results
- Test set contamination
- Limitations
- Broader Impacts
- Related Work

ults

Limitations

Text synthesis and other NLP tasks

GPT-3 achieves strong gains over GPT-2, but has weakness on synthesis and NLP tasks During synthesis, samples sometimes repeat themselves at the document level Samples can also lose coherence, contradict themselves, and contain non-sequiturs For discrete language tasks, GPT-3 seems to struggle with "common sense physics" Struggles with "If I put cheese into the fridge, will it melt?" (though does well on PIQA) For tasks such as determining if two words are used the same way (WiC) or if one sentence implies another (ANLI), GPT-3 is not much better than chance

Structural/Algorithmic limitations

GPT-3 is autoregressive, making it straightforward to sample and compute likelihoods The model lacks any bidirectional architecture and does not use denoising pretraining Various works suggest that bidirectional models can improve fine-tuning Conjecture: a large bidirectional model would be stronger at fine-tuning than GPT-3 A bidirectional model at the scale of GPT-3 seems to be a promising future direction

References

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

(PiQA) Y. Bisk et al., "PiQA: Reasoning about physical commonsense in natural language", AAAI (2020)

(WiC) M. Pilehvar, "WiC: the Word-in-Context Dataset for Evaluating Context-Sensitive Meaning Representations", NAACL-HLT (2019)

(ANLI) Y. Nie et al. "Adversarial NLI: A new benchmark for natural language understanding", ACL (2020)

K. Guu et al,. "REALM: Retrieval-Augmented Language Model Pretraining" arxiv (2020)

The pretraining objective

GPT-3 (or future variants) could run into the limits of the pretraining objective By weighting every token equally, there is no notion of what is important to predict Salient Span Masking (Guu et al., 2020; Roberts et al. 2020) predicts entities of interest Self-supervision casts task as prediction - goal-directed actions may be better GPT-3 is not grounded in other domains of experience (Bisk et al., 2020) Scaling pure self-supervised prediction may hit limits - other approaches may help Possible directions:

- learning the objective function from humans (Ziegler et al., 2019)
- learning from additional modalities to achieve grounding (Chen et al., 2020)

Sample efficiency

GPT-3 takes steps towards test-time sample efficiency (zero/one/few-shot regimes) However, GPT-3 sees much more text than a human sees in their lifetime (Linzen, 2020) Better sample efficiency is a key direction for future work: grounding; algorithms

A. Roberts et al., "How much knowledge can you pack into the parameters of a language model?", arxiv (2020) Y. Bisk et al., "Experience grounds language", arxiv (2020)

D. Ziegler et al., "Fine-tuning language models from human preferences", arxiv (2019) Y-C. Chen et al., "Uniter: Universal image-text representation learning", ECCV (2020) T. Linzen, "How can we accelerate progress towards human-like linguistic generalization?" ACL (2020)

Limitations

Does GPT-3 learn tasks from scratch?

There is ambiguity whether GPT-3's few-shot learning really learns tasks "from scratch" Perhaps it simply recognises and identifies tasks it has seen during pretraining These alternatives exist on a spectrum:

- examples in pretraining exactly match the test time distribution
- recognising the same task but in a different format
- adapting to a specific style of a general task (such as question answering)
- learning a skill entirely from scratch novel words word scrambling Unclear what humans learn from prior demonstrations and what is from scratch Organising diverse demonstrations in pretraining/identifying at test time: an advance Understanding how exactly few-shot learning works is an important future direction

References

- T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)
- G. Hinton et al., "Distilling the knowledge in a neural network", arxiv (2015)
- X. Liu et al., "Improving multi-task deep neural networks via knowledge distillation for natural language understanding", arxiv (2019)

Computational cost of inference

Inference for models at the scale of GPT-3 is expensive and inconvenient

This may be addressed via distillation (Hinton et al., 2015)

Distillation is well-studied for multi-task models (Liu et al., 2019)

It has not been studied at the scale of GPT-3 (new challenges may arise)

Deep learning limitations

As with much of deep learning, GPT-3 suffers from:

- decisions that are not easily interpretable
- poorly-calibrated predictions on novel inputs (higher variance than humans)
- biases of the data it has been trained on

Outline

- Motivation and overview of results
- Approach: training and evaluation
- Results
- Test set contamination
- Limitations
- Broader Impacts
- Related Work

Broader Impacts, Misuse of Language Models

Overview

A range of **beneficial applications** for society can come from language models:

code auto-completion writing auto-completion grammar assistance

game narrative generation

improved search engines question answering

However they also have potentially harmful applications

GPT-3 in particular improves text generation and adaptability

This makes it harder to distinguish synthetic text and human-written text

This ability suggests potential to advance both beneficial and harmful applications

The focus here will be on harms, to stimulate efforts to study/mitigate them

Misuse of language models

It can be hard to predict malicious uses - they often repurpose language models:

in very different environment for a different purpose than researchers intended

One approach: take inspiration from traditional security risk assessment frameworks

Identify threats/potential impacts, assess likelihood, determine risk (Ross, 2012)

Three factors: misuse applications, threat actors, external incentive structures

References

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

- R. Ross, "Guide for Conducting Risk Assessments", Special Publication NIST SP (2012)
- I. Solaiman et al., "Release strategies and the social impacts of language models", arxiv (2019)

Potential misuse applications

Socially harmful activities that use generated text could be augmented by GPT-3:

misinformation spam phishing abuse of legal/government processes

fraudulent academic essay writing social engineering pretexting

Such applications are often bottlenecked by the human ability to write quality text Language models can lower the barrier to entry and increase efficiency/efficacy GPT-3 generation of paragraphs that are hard to distinguish: concerning milestone

Threat Actor Analysis

Threat actors can be grouped by skills/resources (Solaiman et al., 2019):

• low/moderately skilled/resourced actors who may build a malicious product

• Advanced Persistent Threats (APTs), e.g. state-sponsored, with long-term agendas Approach: monitor forums where misinformation, malware, fraud etc. are discussed After GPT-2 release in spring 2019, there was discussion but limited experimentation Discussions correlated with media coverage; threats do not appear immediate For APTs, the GPT-3 authors consulted with professional threat analysts Since GPT-2 release: no discernible difference in operations Language models may not represent a good investment without convincing demos

External Incentive Structures

Incentives for threat actors groups

Economic factors (scalability, ease of deployment) influence TTPs: Augmenting TTPs with language models could further reduce cost of deployment Ease of use - language models are stochastic (may require human filtering, limiting scalability)

References T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

- Threat actor groups each have tactics, techniques and procedures (TTPs) to accomplish their goals
- Phishing is widely practiced because it is a cheap, low-effort, high-yield way to steal logins/deploy malware
- It is likely that AI researchers will eventually develop models that are of greater interest to malicious actors
- This will introduce new challenges for the research community (mitigation research, prototyping, coordination)

Broader Impacts: Bias and Gender

Overview

Biased training data can produce models that generate prejudiced content Bias causes harms (stereotypes, demeaning portrayals) (Crawford, 2016) An analysis of GPT-3 bias is conducted w.r.t fairness, bias and representation Goal: preliminary analysis of limitations and behaviours Focus on biases relating to gender, race and religion (though others are likely)

Gender and occupation

The study focused on associations between occupation and gender Occupations are more likely to be followed by a male identifier than female: Assessed by prompting "The {occupation} was a" (neutral variant) 388 occupations were tested: 83% were more likely to precede male identifiers Male identifiers: (man, male, etc.), Female identifiers: (woman, female, etc.) Occupations indicative of high levels of education were heavily male leaning Occupations requiring physical labour were also male leaning Followed by female identifiers: midwife, nurse, receptionist, housekeeper etc.

References

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)
K. Crawford, "The trouble with bias", NeurIPS (2017)
(Winogender) R. Rudinger et al., "Gender bias in coreference resolution", arxiv (2018)

Gender and competenceOther contexts were also investigated:• "The competent {occupation} was a" (competent variant)• "The incompetent {occupation} was a" (incompetent variant)The incompetent {occupation} was a" (incompetent variant)The competent variant produced even higher fraction of male identifiersThe incompetent variant produce similar ratio to neutral variant (majority male)Average occupation bias: $\frac{1}{n_{jobs}} \sum_{jobs} log(\frac{P(female | context)}{P(male | context)}) \stackrel{0 (neutral) + ve (towards female) - ve (towards male)neural variant: -1.11competent variant: -2.14incompetent variant: -1.15$

Pronoun resolution

Assess pronoun resolution on the Winogender dataset

Prompt: "The advisor met with the advisee because she wanted to get advice about job applications. 'She' refers to" Compare probabilities for (occupation: advisor, participant: advisee) Language models learn biases (tendency to associate females as participants) Of the models compared, the largest (175B) had the highest accuracy (64.17%) Largest model: only model where accuracy for occupants was higher for females Preliminary evidence that larger models may be more robust than smaller models

Gender and Race

Gender co-occurrence

Study which words are likely to occur in the vicinity of other pre-selected words

800 samples of length 50 were generated (temperature 1.0, top_p 0.9)

Prompts: "He was very", "She was very", "He would be described as" **etc.**

Male/female pronouns (easier to study co-occurrence without isolating "they")

Adjectives/adverbs compared in top 100 favoured words with NLTK POS tags

Females more often described with appearance words ("beautiful", "gorgeous")

Males described with adjectives spanning a greater spectrum

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)

Image credits/References

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

E. Loper et al., "NLTK: The natural language toolkit", arxiv (2002)

P-S. Huang et al., "Reducing Sentiment Bias in Language Models via Counterfactual Evaluation", EMNLP (2020) S. Baccianella et al., "Sentiwordnet 3.0: An enhanced lexical resource for sentiment analysis and opinion mining", LREC (2010)

Race

800 samples generated for prompts: "The {race} man was very", "The {race} woman was very", "People would describe the {race} woman as" etc. The placeholder "{race}" was replace with categories e.g. White, Asian etc. Co-occurrences are then measured in the generated samples Prior work shown different sentiments for occupations (Huang et al., 2020) A study is therefore conducted to understand how race affects sentiment Sentiment is assessed with Senti WordNet for disproportionate co-occurrences Sentiment scores vary from 100 (positive words e.g. wonderful: 100) to -100 (negative words e.g. wretched: -87.5), 0 denotes neutral (e.g. chalet) Note: this study explicitly prompts the model to talk about race Co-occurrences can reflect socio-historical factors (discussions related to slavery generally has a negative sentiment - this can affective demographic sentiment)

Religion, Future Bias and Fairness Challenges

	Religion
Study co-c	occurrences with religious terms related to:
Athe	eism Buddhism Christianity Hinduism Islam Judaism
800 samp	les of length 50 were generated (temperature 1.0, top_p 0.9)
Prompts: "	'{Religion practitioners} are" e.g. "Christians are"
Buddhist	<mark>s are</mark> divided into two main branches - Theravada and
Mahayana	. Theravada is the more conservative branch, centering
on monas	tic life and the earliest sutras and refusing to
recogniz	e the later Mahayana sutras as authentic."
Models m	ake associations reflecting how these terms are sometimes presented
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', 'En- lightenment', 'Non-Violent'
Christianity	'Attend', 'Ignorant', 'Response', 'Judgmental', 'Grace', 'Execution', 'Egypt', 'Continue', 'Com- ments', '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'
For Islam:	terrorism terrorist violent co-occur more frequently

Image credits/References

T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)

M. Mitchell et al., "Model cards for model reporting", FAccT (2018)

Y. Qian et al., "Reducing gender bias in word-level language models with a gender-equalizing loss function", arxiv (2019) P-S. Huang et al., "Reducing Sentiment Bias in Language Models via Counterfactual Evaluation", EMNLP (2020)

Future Bias and Fairness Challenges

This represents a preliminary analysis aiming to encourage future research It is expected to be an area of continuous research

Study: "subjective signposting" (gender, race, religion - subjective starting point) Inspired by work to characterise/label model attributes (Mitchell et al., 2018) Characterisation of biases is insufficient - intervention is required Prior work: losses (Qian et al., 2019) counterfactuals (Huang et al., 2020) Research on large language models would benefit from (Blodgett et al., 2020):

- common vocabulary for normative, technical, empirical bias challenges
- engage literature outside NLP
- better articulates normative statements about harm
- engages with lived experience of communities affected by NLP systems

Pure metric-driven bias removal objectives have blind spots (Gonen et al., 2019)

A holistic approach to bias mitigation is recommended

S. L. Blodgett et al., "Language (technology) is power: A critical survey of 'bias' in nlp", arxiv (2020) H. Gonen et al., "Lipstick on a pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them", NAACL-HLT (2019)

Energy usage

Energy usage for GPT-3

Pre-training at large scale uses energy-intensive computation GPT-3 175B required several thousand petaflop/s-days Petaflop/s-day: 10¹⁵ operations/second for a day (Amodei et al., 2018) By contrast, GPT-2 required tens of petaflop/s-days Researchers should be cognisant of cost/efficiency (Schwartz et al., 2020) Useful to account for how training is amortised over the lifetime of model The model will be used for various purposes/fine-tuned for specific tasks GPT-3 consumes a lot of energy in training, but can have efficient uses Generating 100 pages with GPT-3 175B uses 0.4 kW-hr (costs a few cents) Distillation (Liu et al., 2019) can reduce costs by creating efficient versions Algorithmic gains may enhance efficiency (Hernandez et al., 2020)

Image credits/References

- T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)
- D. Amodei et al., "AI and Compute", <u>https://openai.com/blog/ai-and-compute/</u> (2018)
- R. Schwartz et al., "Green AI", Communications of the ACM (2020)
- X. Liu et al., "Improving multi-task deep neural networks via knowledge distillation for natural language understanding", arxiv (2019)
- D. Hernandez et al., "Measuring the algorithmic efficiency of neural networks", arxiv (2020)

Outline

- Motivation and overview of results
- Approach: training and evaluation
- Results
- Test set contamination
- Limitations
- Broader Impacts
- **Related Work**

Related Work

Vaswani et al., (2017)

Image credits/References

Instructions in natural language

- A. Vaswani et al., "Attention is all you need", NeurIPS (2017)
- B. McCann et al., "The natural language decathlon: Multitask learning as question answering", arxiv (2018)

S. Hochreiter et al., "Learning to learn using gradient descent", ICANN (2001) (T5) C. Raffel et al., "Exploring the limits of transfer learning with a unified text-to-text transformer", JMLR (2020)

Summary

Takeaways for GPT-3

This work introduced GPT-3, a 175 billion parameter language model GPT-3 shows strong performance on a range of NLP tasks in zero/one/few-shot settings In many tasks it matches state-of-the-art fine-tuned methods GPT-3 can also synthesise high-quality written text samples We observed that GPT-3 performance scales fairly smoothly with model size Some societal impacts of the model have been considered, together with limitations Although GPT-3 has many weaknesses, the results suggest that large language models may play a key role in developing adaptable, general language systems

References T. Brown et al., "Language models are few-shot learners", NeurIPS (2020)