

Training Session
on
OpenAI's
Generative Pretrained Transformer (GPT) - 3

Ishan Jain
ishanjain.ai@gmail.com

Agenda

1. Evolution of NLP/NLG models
2. OpenAI's GPT-3
3. GPT-3 applications
4. Conclusions & step forward

Terminologies

Open AI - A non-profit AI research organization

NLP - Natural Language Processing

NLU - Natural Language Understanding

NLG - Natural Language Generation

Language Models - A statistical language model is a probability distribution over sequences of words.

In-context learning - Learning and training of a model based on specific tasks and prompts

Parameters - A parameter is a calculation in a neural network that applies a great or lesser weighting to some aspect of the data

AGI - Artificial General Intelligence

Evolution of NLP/NLG models

No	Model	Year	Description
1	Google's BERT	2018	The self-supervised method released by Google in 2018.
2	<u>Allen NLP's ELMo</u>	2018	A deep contextualized word representation that models both complex and linguistic contexts.
3	<u>ULMFit</u>	2018	Universal Language Model Fine-tuning for Text Classification.
4	<u>Google's XLNet</u>	2019	A generalized autoregressive pretraining method that enables learning bidirectional contexts by maximizing the expected likelihood over all permutations of the factorization order.
5	<u>Google's ALBERT</u>	2019	A Lite BERT for Self-Supervised Learning of Language Representations
6	<u>FB's RoBERTa</u>	2019	Upgrade on BERT. Implemented in PyTorch. An optimised method for pre-training a self-supervised NLP system built on BERT's language masking strategy.
7	<u>Microsoft's CodeBERT</u>	2020	Google's BERT framework for NLP, has been built upon a bidirectional multi-layer neural architecture.
8	<u>Microsoft's ZeRO-2</u>	2020	Zero Redundancy Optimizer version 2 (ZeRO-2), a distributed deep-learning optimization algorithm.
9.	GPT-3	2020	GPT-3 is one of the most controversial pre-trained models by OpenAI.



GPT-3

The first AGI model

Natural Language Processing Model

Non-deterministic

NLP task-agnostic

Requires minimum fine-tuning

Evolution of Generative Pretrained Transformer (GPT) models

OpenAI's GPT

2018

Unsupervised

Natural Language Understanding

OpenAI's GPT-2

2019

10X scaled up on parameters and training data.

1.5 billion parameters. Trained on 8 million web pages

100X

OpenAI's GPT-3

2020

175 billion parameters

A neural-network based language model.

[illegible][illegible]

Language Models are Few-Shot Learners					
Tom B. Brown		Brendan Mnemonic		Noa Reid*	
Meta		Meta		Meta	
Jared Kaplan		Pranav Shyam		Ariella Shulman	
Meta		Meta		Meta	
Alec Radford		Santosh Agrawal		Adel H. Wu	
Meta		Meta		Meta	
Brendan Chow		Ariella Shulman		Daniel DeFries	
Meta		Meta		Meta	
Christopher Hesse		Mark Chen		Eric Nigam	
Meta		Meta		Meta	
Brendan Mnemonic		Chris Olah		Christopher Brierley	
Meta		Meta		Meta	
Sam McClelland		Alex Raffel		Hyo Sukwon	
Meta		Meta		Meta	
OpenAI					
Abstract					
<p>Recent work has demonstrated that robust task-specific fine-tuning of thousands of pre-trained models can be used to construct a general-purpose language model. This paper shows that a single transformer-based model can learn to perform a wide range of tasks, including those that require complex reasoning, without the need for task-specific fine-tuning. We show that a single transformer-based model can learn to perform a wide range of tasks, including those that require complex reasoning, without the need for task-specific fine-tuning. We show that a single transformer-based model can learn to perform a wide range of tasks, including those that require complex reasoning, without the need for task-specific fine-tuning.</p>					

What is OpenAI's GPT-3?

- The GPT-3 model architecture itself is a transformer-based neural network that has been fed 45TB of text data.
- A language model is a model that predicts the likelihood of a sentence existing in the world.
- GPT-3 is non-deterministic, in the sense that given the same input, multiple runs of the engine will return different responses.
- GPT-3 is trained on massive datasets that covered the entire web and contained 500B tokens, humongous 175 Billion parameters, a more than 100x increase over GPT-2, which was considered state-of-the-art technology with 1.5 billion parameters.
- OpenAI's GPT-3 is still in the experimental phase.

What makes GPT-3 so magical?

It is really big

With 175 billion parameters, it's the largest language model ever created

Minimum fine-tuning

It only requires few-shot demonstrations via textual interaction with the model.

Extraordinary - Supermodel

You can ask GPT-3 to be a translator, a programmer, a poet, or a famous author, and it can do it with its user (you) providing fewer than 10 training examples.

Custom language tasks without training

Task-agnostic NLP model

GPT-3 - a few shot learner

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

GPT-3 - a few shot learner

An overview of the original paper covering its training cost and research implications.

- GPT-3 shows that language model performance scales as a power-law of model size, dataset size, and the amount of computation.
- GPT-3 demonstrates that a language model trained on enough data can solve NLP tasks that it has never encountered. That is, GPT-3 studies the model as a general solution for many downstream jobs without fine-tuning.
- The cost of AI is increasing exponentially. Training GPT-3 would cost over \$4.6M using a Tesla V100 cloud instance.
- The size of state-of-the-art (SOTA) language models is growing by at least a factor of 10 every year. This outpaces the growth of GPU memory. For NLP, the days of "embarrassingly parallel" are coming to the end; model parallelization will become indispensable.
- Although there is a clear performance gain from increasing the model capacity, it is not clear what is really going on under the hood. Especially, it remains a question of whether the model has learned to do reasoning, or simply memorizes training examples in a more intelligent way.

GPT-3 - a few shot learner

About GPT-3's training data

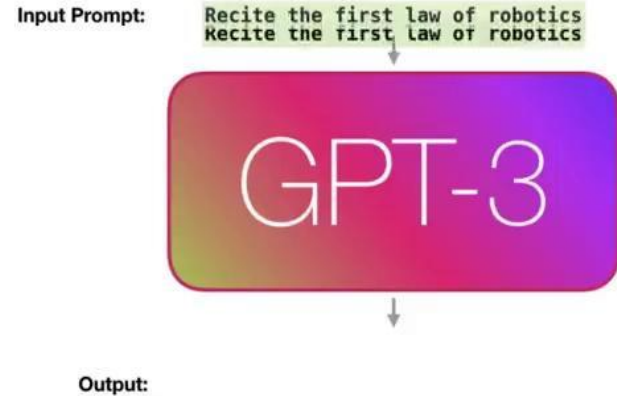
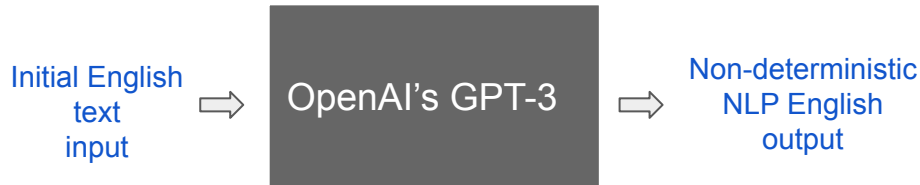
- It is a weighted mix of Common Crawl, WebText2 (a larger version of the original), two book corpora, and English Wikipedia.
- Some components (e.g. Wikipedia), were completely sampled 3+ times during training, while others like the Common Crawl, weren't even completely sampled. The authors claim that this is to help raise the overall quality of the corpus by prioritising known-good datasets.
- Altogether, the filtered/cleaned dataset is 500 billion tokens, or 700GB.
- Due to a bug, some data overlapped between the training and test sets. The paper analyzes of the impact of this leakage.

How does GPT-3 work?

How does GPT-3 work?

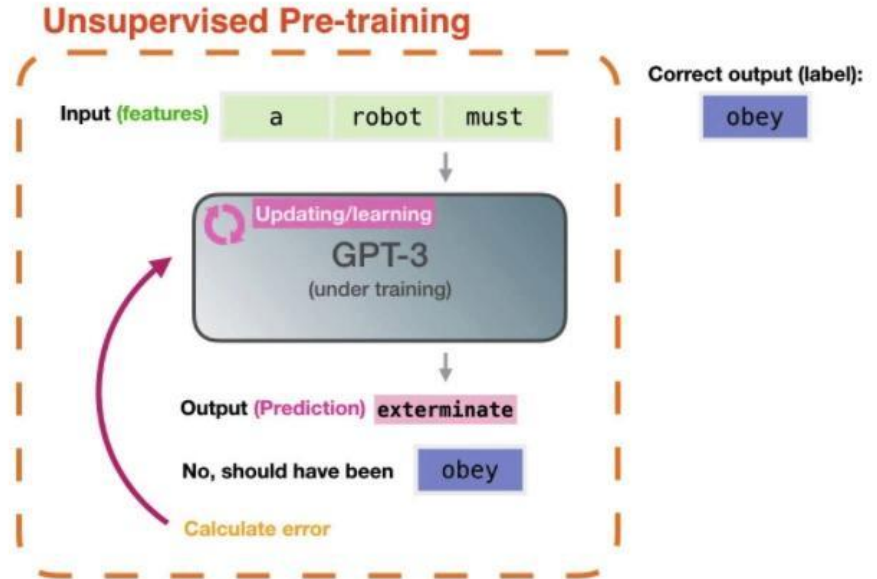
GPT- 3 language model generates text. We can optionally pass it some text as input, which influences its output.

The output is generated from what the model “learned” during its training period where it scanned vast amounts of text.

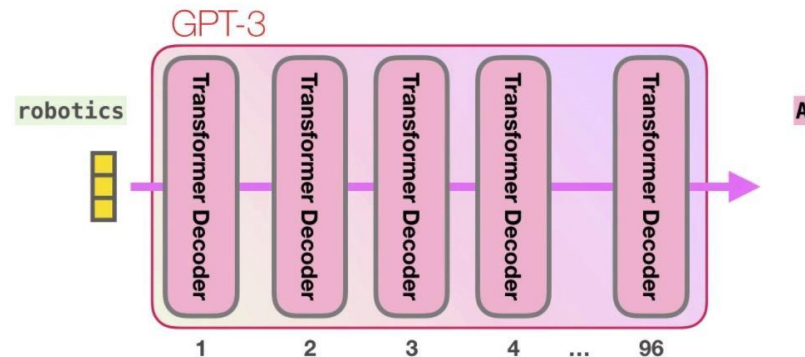
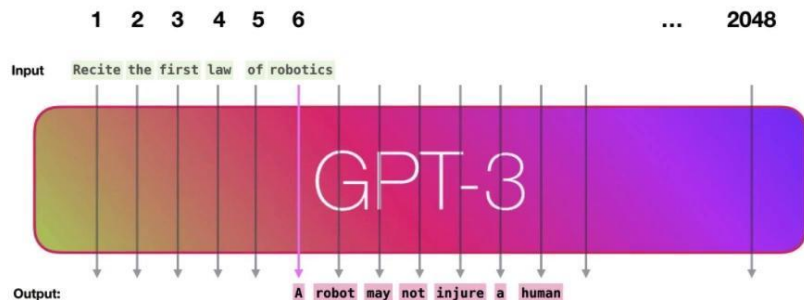


How does GPT-3 work?

- Training is the process of exposing the model to lots of text.
- **All the experiments are from that one trained model.**
- It was estimated to cost 355 GPU years and cost \$5m - \$12m.



How does GPT-3 work?



- GPT3 actually generates output one token at a time (let's assume a token is a word).
- The important calculations of the GPT3 occur inside its stack of 96 transformer decoder layers. This is the “depth” in “deep learning”.
- Each of these layers has its own 1.8B parameter to make its calculations. That is where the “magic” happens.
- GPT3 is 2048 tokens wide. That is its “context window”. That means it has 2048 tracks along which tokens are processed.

How does GPT-3 work?

The model is evaluated in three different settings:

- **Few-shot learning**, when the model is given a few demonstrations of the task (typically, 10 to 100) at inference time but with no weight updates allowed.
- **One-shot learning**, when only one demonstration is allowed, together with a natural language description of the task.
- **Zero-shot learning**, when no demonstrations are allowed and the model has access only to a natural language description of the task.

The three settings we explore for in-context learning

Zero-shot

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



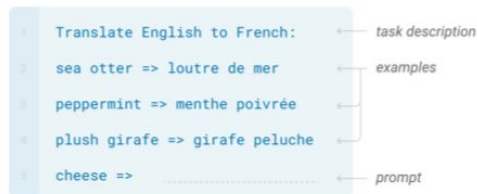
One-shot

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.



Traditional fine-tuning (not used for GPT-3)

Fine-tuning

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

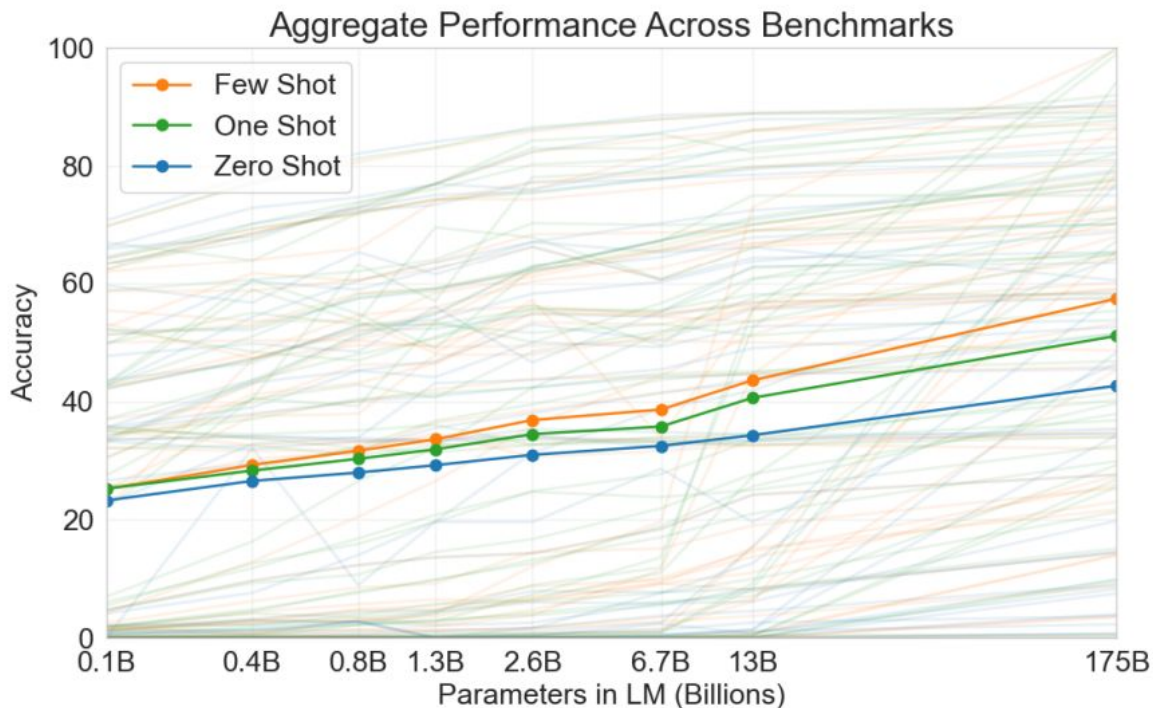


Zero-shot, one-shot and few-shot, contrasted with traditional fine-tuning.

Performance of GPT-3

Performance of GPT-3

- Unlike its predecessors, GPT-3 can infer a task from one or a few example: this is a massive step towards generalization.
- GPT-3 can be “tuned” by providing instructions in plain English, whereas its predecessors require task-specific tuning.
- Increasing model size improves performance across almost all tasks; in contrast, fine-tuning limits performance gains to one task and risks catastrophic forgetting and overfitting.



Performance of GPT-3

Key Achievements:

- The GPT-3 model without fine-tuning achieves promising results on a number of NLP tasks, and even occasionally surpasses state-of-the-art models that were fine-tuned for that specific task:
 - On the **CoQA** benchmark, 81.5 F1 in the zero-shot setting, 84.0 F1 in the one-shot setting, and 85.0 F1 in the few-shot setting, compared to the 90.7 F1 score achieved by fine-tuned SOTA.
 - On the **TriviaQA** benchmark, 64.3% accuracy in the zero-shot setting, 68.0% in the one-shot setting, and 71.2% in the few-shot setting, surpassing the state of the art (68%) by 3.2%.
 - On the **LAMBADA** dataset, 76.2 % accuracy in the zero-shot setting, 72.5% in the one-shot setting, and 86.4% in the few-shot setting, surpassing the state of the art (68%) by 18%.
- **The news articles generated by the 175B-parameter GPT-3 model are hard to distinguish from real ones**, according to human evaluations (with accuracy barely above the chance level at ~52%).

GPT-3 Applications

GPT-3 as a Search Engine

The OpenAI' GPT-3 API allows searching over documents based on the natural-language meaning of queries rather than keyword matching.



Casetext automates litigation tasks to help attorneys

[Casetext Demo](#)

[Search Engine Demo](#)

Conversation & Chat

The API can enable fast, complex and consistent natural language discussions. With a brief prompt, the API generates dialogues spanning a range of topics, from space travel to history.

The logo for ai|channels, featuring the text "ai|channels" in white on a dark blue rectangular background.

AI Channels is a social network for people and artificial intelligence agents.

[AI Channels Demo](#)

Content Comprehension & Generation

The API can be used to build tools to help individuals consume content more efficiently.

The API can generate complex and consistent natural language, and enables use cases like creative writing.



Replika

Replika, an AI companion

[Replica Demo](#)

The API can transform dense text into simplified summaries.

[Summarization Demo](#)

[Meme Generation Demo](#)

Software Engineer & Data Scientist

Context-aware code suggestions.

Translate natural language to unix commands.

Generate ML codes.

Words to Websites.

After fine-tuning with code from thousands of Open Source GitHub repositories, the API completes code based on function names and comments.

[Code Completion Demo](#)

[Natural Language Shell Demo](#)

[No Code AI Demo](#)

[Website Template Demo](#)

Personalization

The API can be used to build tools to help individuals consume content more efficiently.



Art of Problem Solving (AoPS) is helping to effectively prepare the next generation of STEM professionals through engaging online instruction, at a time when the traditional nature of in-person education is being challenged.

[AoPS Demo](#)

Conclusions

- GPT-3 can write.
- GPT-3 can have a conversation.
- GPT-3 can self reflect.
- GPT-3 is amazing. It's scary. It's exhilarating. It's the biggest thing since bitcoin. It's the future.