
Week 8 Reading Summaries: GPT-3 and Chinchilla

Amogh Patankar

University of California, San Diego
3869 Miramar Street Box #3037, La Jolla, CA- 92092
apatankar@ucsd.edu

Abstract

With these readings, we dive deep into the most famous generative model in recent years, GPT-3, as well as a paper from DeepMind looking at the training strategies for compute-optimal LLMs. The OpenAI paper specifically trains GPT-3, a 175B parameter model, and tests its performance in a few-shot environment; the testing is done in many scenarios. Chinchilla, a compute-optimal model from DeepMind, is tested on a range of downstream evaluations, and outperforms Gopher, a 280B parameter model, reaching a SOTA level on a famous benchmark, Massive Multitask Language Understanding (MMLU).

1 GPT-3

1.1 Introduction

In recent memory, ML has continued to trend in the direction of NLP systems, specifically, models with large quantity of parameters, being fine-tuned for challenging tasks like Q+A, comprehension, etc. However, for these large models, the quantity of data often limits performance, and pre-training and fine-tuning struggles with a narrowed training distribution. The solution, is meta-learning- models learning a wider range of abilities, and applying specific ones at inference. This paper uses a 175B parameter model (GPT-3), and evaluates on many novel tasks, achieving SOTA performance on most evaluations.

1.2 Approach

For this model, they identify four techniques to improve their power- fine-tuning, zero-shot, one-shot, and few-shot learning. They test OS, 1S, and FS learning as opposed to traditional finetuning.

1.2.1 Model + Architecture

The paper utilizes the same architecture as GPT-2, but uses dense and sparse attention layers in an alternating pattern; they also test eight different models, all with different sizes and hyperparameters.

1.2.2 Training Dataset

They train the model on a variety of datasets- primarily though, the Common Crawl dataset, which is roughly 1 trillion words. However, they filter it, using high quality corpora, perform deduplication, and add higher quality corpora to CommonCrawl.

1.2.3 Training Process

The authors use gradient noise scale to guide their batch size choices, and use a mixture of model parallelism, inter- and intra- network, while using V100 GPUs.

31 1.2.4 Evaluation

32 They perform a variety of evaluations for GPT-3 for few-shot learning, they evaluate by drawing
33 random samples and use to condition the model. For multiple choice tasks, they use K random
34 samples, with context with correct completion, then calculate LM likelihood. For binary classification
35 tasks, they give more meaningful options, and use multiple choice framework, while with free-form
36 responses, they use beam search, and score using F1 scoring. Then, for model size and learning
37 settings, they use a developmental set for the model size(s).

38 1.3 Results

39 With respect to GPT-3, the authors note that language modeling performance is restricted by a
40 power-law, especially with regards to training compute. In the various sections, they evaluate various
41 evaluation frameworks and tests.

42 1.3.1 Language Modeling, Cloze, and Completion Tasks

43 For language modeling, they calculate 0S perplexity on the PTB dataset, setting a new SOTA, and
44 because PTB is a traditional dataset, it doesn't have clear example separation. LAMBADA is intended
45 to test long-range dependency modeling, and is the main usage of FS modeling- and hits the SOTA
46 again, by over 18%. Testing HellaSwag and StoryCloze is meant for few-shot ending prediction, and
47 GPT is close but doesn't exceed SOTA.

48 1.3.2 Closed Book Q + A

49 The authors measure GPT-3 against factual knowledge Q + A; usage of 0S, 1S, and FS modeling is
50 even stricter than previous work, and fine-tuning isn't allowed. GPT-3 achieves matches or exceeds
51 SOTA for TriviaQA, and WebQS, while coming close for NQs. This suggests a limitation of GPT-3's
52 knowledge capacity.

53 1.3.3 Translation

54 For GPT-3, the authors note that the training is 93% in English, and 7% other, meaning that GPT-3
55 learns from a mixture of languages. This results in GPT-3 performing at a level close to SOTA,
56 but not quite equivalent, for obvious reasons, such as another model being specifically trained on
57 50%/50% English and French, for ENG to FRE translation, for example.

58 1.3.4 Winograd-Style Tasks

59 Winograd tasks, referring to word determination, were used as evaluators for GPT-3, and as such,
60 GPT-3 trails fine-tuned SOTA in other commonsense reasoning challenges.

61 1.3.5 Common Sense Reasoning

62 For common sense reasoning, they evaluate on PIQA, where GPT-3 achieves SOTA accuracy for
63 0S, 1S, and FS; But on ARC (easy, challenge), and OpenBookQA, it falls short of fine-tuned SOTA.
64 Overall, they show mixed results for commonsense reasoning by GPT-3 vs SOTA.

65 1.3.6 Reading Comprehension

66 For reading comprehension tasks, GPT-3 is best on freeform datasets, and worst on structured dialog
67 acts with respect to interaction based data. This trend follows for discrete reasoning datasets like
68 DROP too.

69 1.3.7 SuperGLUE

70 For the superGLUE benchmark, GPT samples a new set of examples in context for each problem;
71 GPT varies in performance fro the varous benchmarks in SuperBLUE.

72 1.3.8 Natural Language Inference

73 For NLI tasks, i.e. understanding the relationship between sentences, GPT is much better in settings
74 that aren't few-shot; GPT-3 performs better than smaller models, but still has room to improve for
75 NLI tasks.

76 1.3.9 Synthetic + Qualitative Tasks

77 The authors also test GPT-3 on tasks to perform computations, such as arithmetic, word-scrambling
78 and manipulation, SAT analogies, etc. Across these tasks, GPT-3 is able to perform excessively well,
79 beating SOTA and achieving human level average scores in some cases like SAT analogies.

80 1.4 Measuring and Preventing Memorization of Benchmarks

81 The authors also ensure that GPT-3 doesn't simply memorize information that exists in train sets,
82 and regurgitate it during test time/inference time, and this is backed by previous literature. GPT-3
83 doesn't overfit by a significant amount, meaning the memorization is limited; they also use a clean
84 version of the dataset(s) that removes all leaked examples that match (under 13 grams). Performance
85 in the following: reading comprehension, German translation, reversed words, PIQA, Winograd,
86 and language modeling, only marginally decreases after cleaning. The authors do work for data
87 contamination and clean datasets in order to have a fair, performant model.

88 1.5 Limitations

89 GPT-3 has a few limitations. Namely, it struggles with text synthesis and NLP tasks; it seems to
90 repeat itself semantically, and has difficulty with in-context learning. Similarly, it has constraints with
91 respect to algorithm and architecture. Since it's an autoregressive model, it doesn't use bidirectionality,
92 and it may explain mediocre performance in tasks. Other limitations of GPT-3 include sampling
93 efficiency, ambiguity of memorization vs learning, and compute cost and practicality.

94 1.6 Broader Impacts

95 The authors discuss the misuse of language models, fairness, bias, and representation, and energy
96 usage. They note that language models can be misused as text generation from GPT is very human-
97 like, and that threat actors pose high danger. They analyze gender, race, and other biases that GPT-3
98 characterizes biases by, and discuss energy-intensive compute, and note efficiency.

99 1.7 Conclusion

100 The authors present a 175B parameter generative language model that shows SOTA performance on
101 many benchmarks and evaluations, while nothing the flexibility of the model.

102 2 Chinchilla Scaling Law

103 2.1 Introduction

104 The authors introduce the problem- that LLMs have large parameter counts (in excess of 500B), and
105 as such, the compute and energy rises. While literature exists showing power law relationship for
106 parameters and performance, model size and training tokens should be scaled in proportions. The
107 authors look at FLOPs budgets, and try to model pre-training loss to work backwards to find compute
108 numbers. That is, in this paper, the authors come up with a compute-optimal model, Chinchilla, that
109 outperforms Gopher with reduced size and cost.

110 2.2 Related Work

111 The authors discuss large language models, scaling behavior for those LLMs, hyperparameter
112 judgment, and improved model architectures. They note that the advent of LLMs has introduced a
113 need for high quality data, and they analyze scaling properties for LLMs. Moreover, they note that
114 tuning hyperparameter combinations and heuristics to determine those result in good performance.

115 2.3 Estimating The Optimal Parameter/Training Tokens Allocation

116 The authors look at three different approaches in order to solve the question of model size and training
117 token tradeoff.

118 With approach 1, they fix the model sizes, as well as the # of training steps required to train the
119 model, with varying training sequences (# of training tokens). In doing so, they find the model size
120 that achieves the lowest loss, after tuning FLOP count, learning rate, etc.

121 With approach 2, the authors vary the model size for varying FLOP counts, and compute loss. For
122 each FLOP count, they train diverse model sizes; they utilize a power law between FLOPS, model
123 size, and # of training tokens to find the right values for α and β .

124 With approach 3, the authors model loss as a function of parameter count and seen tokens, then
125 minimize the Huber loss, then plotting as contours, and finding optimal α and β values.

126 2.3.1 Optimal Model Scaling

127 With those approaches, they find that the predictions all seem similar- that as compute increases,
128 model size and training data must also increase proportionally. The amount of training data projected
129 is more than currently used for LLMs, meaning that datasets are the big changes to see engineering
130 improvements.

131 2.4 Chinchilla

132 The authors discuss Chinchilla, their newer model after Gopher, already a compute-optimal model;
133 Chinchilla is more performant, and has a smaller memory and inference footprint.

134 2.4.1 Model + Training details

135 Chinchilla is trained on similar data as compared to Gopher; the distribution, optimizer, tokenizer,
136 and quantization of the weights is different, just as the model architecture is.

137 2.4.2 Results

- 138 • Language Modeling: Chinchilla outperforms Gopher on all evaluations of the Pile dataset,
139 and is more performant than existing models. However, when comparing Chinchilla to other
140 models, train/test leakage affects results, similar to the GPT-3 model from the previous
141 paper.
- 142 • MMLU: Regarding the MMLU task, Chinchilla performs Gopher by roughly 8%, and greater
143 accuracy than other existing models in specific tasks.
- 144 • Reading Comprehension: Chinchilla outperforms other models on LAMBADA.
- 145 • BIG-bench: Chinchilla outperforms Gopher and other models by roughly 10%.
- 146 • Common sense: Chinchilla is able to reach performant metrics for 0, 5, and 10-shot on
147 TruthfulQA dataset/evaluation.
- 148 • Closed-book Q + A: Chinchilla also achieves closed-book SOTA performances on accura-
149 cies.
- 150 • Gender bias + toxicity: Regarding gender and toxicity, the authors conduct a thorough
151 evaluation for Chinchilla. They note that for gender, Chinchilla is better at gender dis-
152 crimination/categorization than its predecessor, Gopher. Similar performance is noticed for
153 toxicity recognition.

154 2.5 Discussion and Conclusion

155 The authors suggest that based on the results of their performant solution, Chinchilla, that LLM trends
156 will move towards optimally setting model size, training duration, and # of training tokens. They
157 note that predicting the scaling of LLMs has limitations, such as intermediate scaling performance,
158 and that high-quality data is required for scaling laws to apply. All in all, a large variety of factors is
159 required for scaling laws to apply correctly for LLMs, as seen by their analysis of Chinchilla, Gopher,
160 and all the other models they evaluated.