AI Seminar #16

# Language Models are Few-Shot Learners (aka. GPT-3)

최원혁
21. 03. 29

# Language Models are Few-Shot Learners

Tom B. Brown[*]  Benjamin Mann[*]  Nick Ryder[*]  Melanie Subbiah[*]

Jared Kaplan[†]  Prafulla Dhariwal  Arvind Neelakantan  Pranav Shyam  Girish Sastry

Amanda Askell  Sandhini Agarwal  Ariel Herbert-Voss  Gretchen Krueger  Tom Henighan

Rewon Child  Aditya Ramesh  Daniel M. Ziegler  Jeffrey Wu  Clemens Winter

Christopher Hesse  Mark Chen  Eric Sigler  Mateusz Litwin  Scott Gray

Benjamin Chess  Jack Clark  Christopher Berner

Sam McCandlish  Alec Radford  Ilya Sutskever  Dario Amodei

OpenAI

- 2020년 논문

- OpenAI
  - ✓ 2015년 설립된 인공지능 연구소
  - ✓ 설립자 : 일론 머스크, Altman, 피터 틸 등

• Language Model (Auto Regressive) 훈련을 통해서 다양한 NLP task에 접목할 수 있다.

**Zero-shot, One-shot, Few-shot** 으로만...!

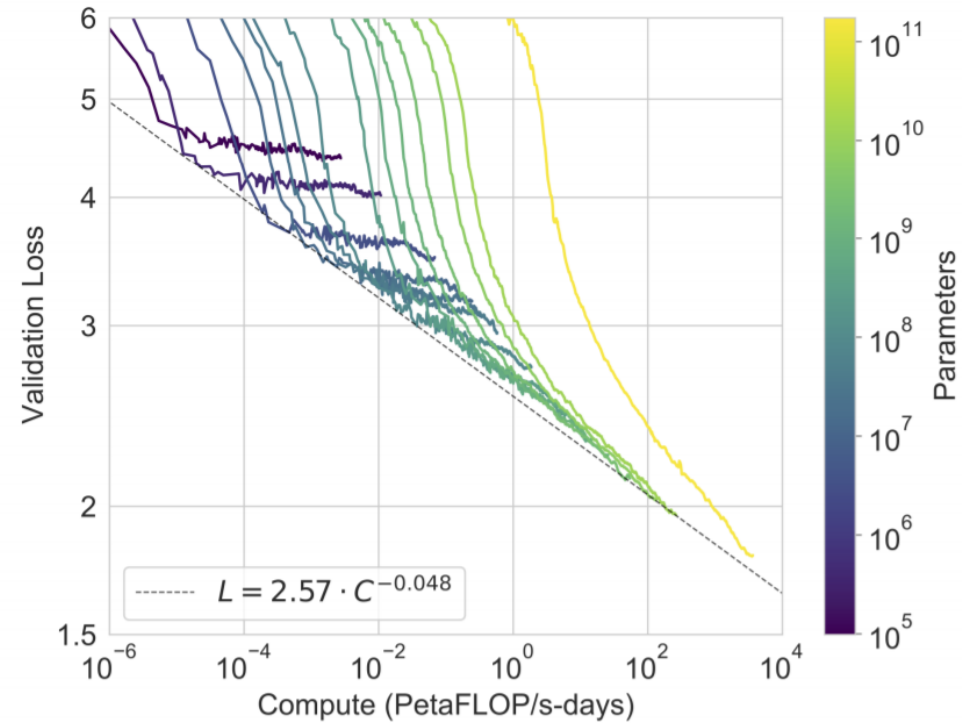| | |
|---|---|
| ▪ Language Modeling, Cloze, and Completion tasks | 문장 생성, 마지막 단어 예측, 다음 문장 예측 |
| ▪ Closed Book Question Answering | Question Answering |
| ▪ Translation | Translation |
| ▪ Winograd-Style tasks | 대명사가 가르키는 것 찾기 |
| ▪ Common Sense Reasoning | Common Sense Reasoning |
| ▪ Reading Comprehension | 독해 |
| ▪ superGLUE | 다양한 NLP task (8개) |
| ▪ NLI | 두 문장의 관계 추론 |
| ▪ Synthetic and Qualitative tasks | 산수, 철자 고치기, 뉴스 기사 생성 등 |

**Figure 3.1: Smooth scaling of performance with compute.** Performance (measured in terms of cross-entropy validation loss) follows a power-law trend with the amount of compute used for training. The power-law behavior observed in [KMH+20] continues for an additional two orders of magnitude with only small deviations from the predicted curve. For this figure, we exclude embedding parameters from compute and parameter counts.

# GPT-3

개요

- Model Size

| Model Name | $n_{params}$ | $n_{layers}$ | $d_{model}$ | $n_{heads}$ | $d_{head}$ | Batch Size | Learning Rate |
|---|---|---|---|---|---|---|---|
| GPT-3 Small | 125M | 12 | 768 | 12 | 64 | 0.5M | $6.0 \times 10^{-4}$ |
| GPT-3 Medium | 350M | 24 | 1024 | 16 | 64 | 0.5M | $3.0 \times 10^{-4}$ |
| GPT-3 Large | 760M | 24 | 1536 | 16 | 96 | 0.5M | $2.5 \times 10^{-4}$ |
| GPT-3 XL | 1.3B | 24 | 2048 | 24 | 128 | 1M | $2.0 \times 10^{-4}$ |
| GPT-3 2.7B | 2.7B | 32 | 2560 | 32 | 80 | 1M | $1.6 \times 10^{-4}$ |
| GPT-3 6.7B | 6.7B | 32 | 4096 | 32 | 128 | 2M | $1.2 \times 10^{-4}$ |
| GPT-3 13B | 13.0B | 40 | 5140 | 40 | 128 | 2M | $1.0 \times 10^{-4}$ |
| GPT-3 175B or "GPT-3" | 175.0B | 96 | 12288 | 96 | 128 | 3.2M | $0.6 \times 10^{-4}$ |

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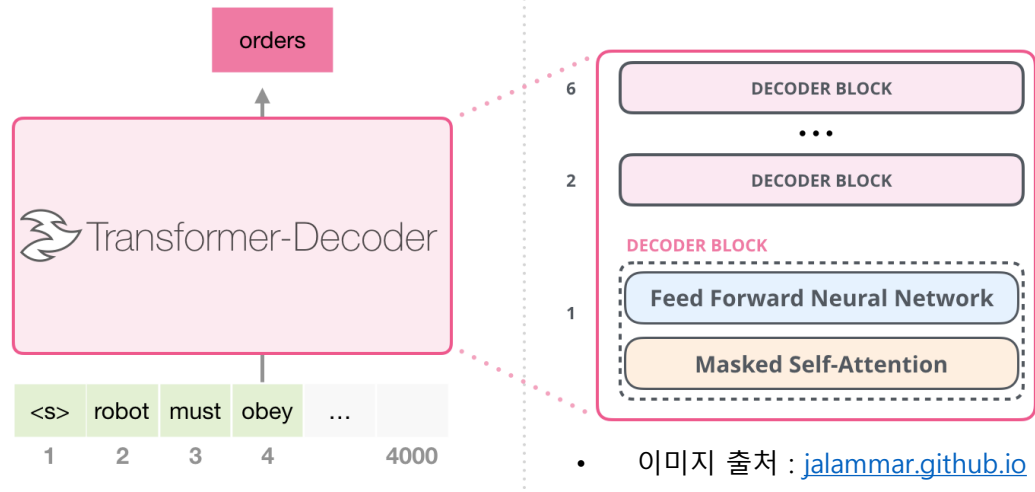
- Data Size

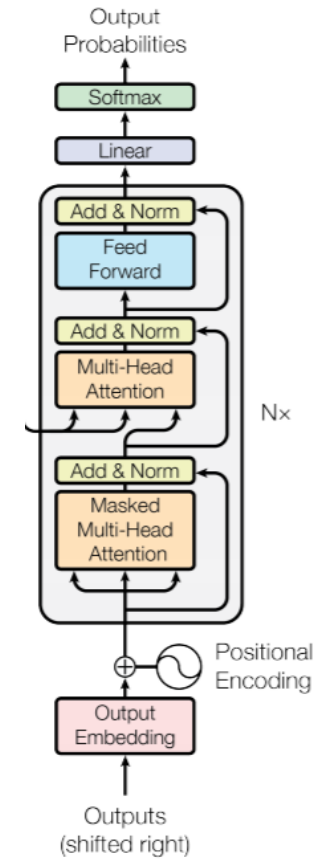| Dataset | Quantity (tokens) | Weight in training mix | Epochs elapsed when training for 300B tokens |
|---|---|---|---|
| Common Crawl (filtered) | 410 billion | 60% | 0.44 |
| WebText2 | 19 billion | 22% | 2.9 |
| Books1 | 12 billion | 8% | 1.9 |
| Books2 | 55 billion | 8% | 0.43 |
| Wikipedia | 3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

**Table 2.2: Datasets used to train GPT-3.** "Weight in training mix" refers to the fraction of examples during training that are drawn from a given dataset, which we intentionally do not make proportional to the size of the dataset. As a result, when we train for 300 billion tokens, some datasets are seen up to 3.4 times during training while other datasets are seen less than once.

# GPT-3

Architecture



orders

Transformer-Decoder

| | DECODER BLOCK |
|---|---|
| 6 | DECODER BLOCK |
| | ... |
| 2 | DECODER BLOCK |

DECODER BLOCK

| 1 | Feed Forward Neural Network |
|---|---|
| | Masked Self-Attention |

| <s> | robot | must | obey | ... | |
|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | | 4000 |

- 이미지 출처 : jalammar.github.io

GPT Model

Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Nx

Add & Norm

Masked
Multi-Head
Attention

Positional
Encoding

Output
Embedding

Outputs
(shifted right)

Transformer Model

# GPT-3

## Approach



Figure 1.1: Language model meta-learning. During unsupervised pre-training, a language model develops a broad set of skills and pattern recognition abilities. It then uses these abilities at inference time to rapidly adapt to or recognize the desired task. We use the term "in-context learning" to describe the inner loop of this process, which occurs within the forward-pass upon each sequence. The sequences in this diagram are not intended to be representative of the data a model would see during pre-training, but are intended to show that there are sometimes repeated sub-tasks embedded within a single sequence.

# GPT-3

## Approach



**Figure 2.1: Zero-shot, one-shot and few-shot, contrasted with traditional fine-tuning**. The panels above show four methods for performing a task with a language model – fine-tuning is the traditional method, whereas zero-, one-, and few-shot, which we study in this work, require the model to perform the task with only forward passes at test time. We typically present the model with a few dozen examples in the few shot setting. Exact phrasings for all task descriptions, examples and prompts can be found in Appendix G.

# GPT-3

1. Language Modeling, Cloze, and Completion Tasks

| Setting | PTB |
| --- | --- |
| SOTA (Zero-Shot) | $35.8^a$ |
| GPT-3 Zero-Shot | **20.5** |

- Penn Tree Bank (PTB) dataset
- Perplexity score
- SOTA : GPT-2

# GPT-3

1. Language Modeling, Cloze, and Completion Tasks

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

**Table 3.2: Performance on cloze and completion tasks.** GPT-3 significantly improves SOTA on LAMBADA while achieving respectable performance on two difficult completion prediction datasets. [a][Tur20] [b][RWC+19] [c][LDL19] [d][LCH+20]

- LAMBADA
  - ✓ Asked to predict the last word of sentences
- StoryCloze
  - ✓ Selecting the correct ending sentence for five-sentence long stories
- HellaSwag
  - ✓ Picking the best ending to a story or set of instructions

# GPT-3

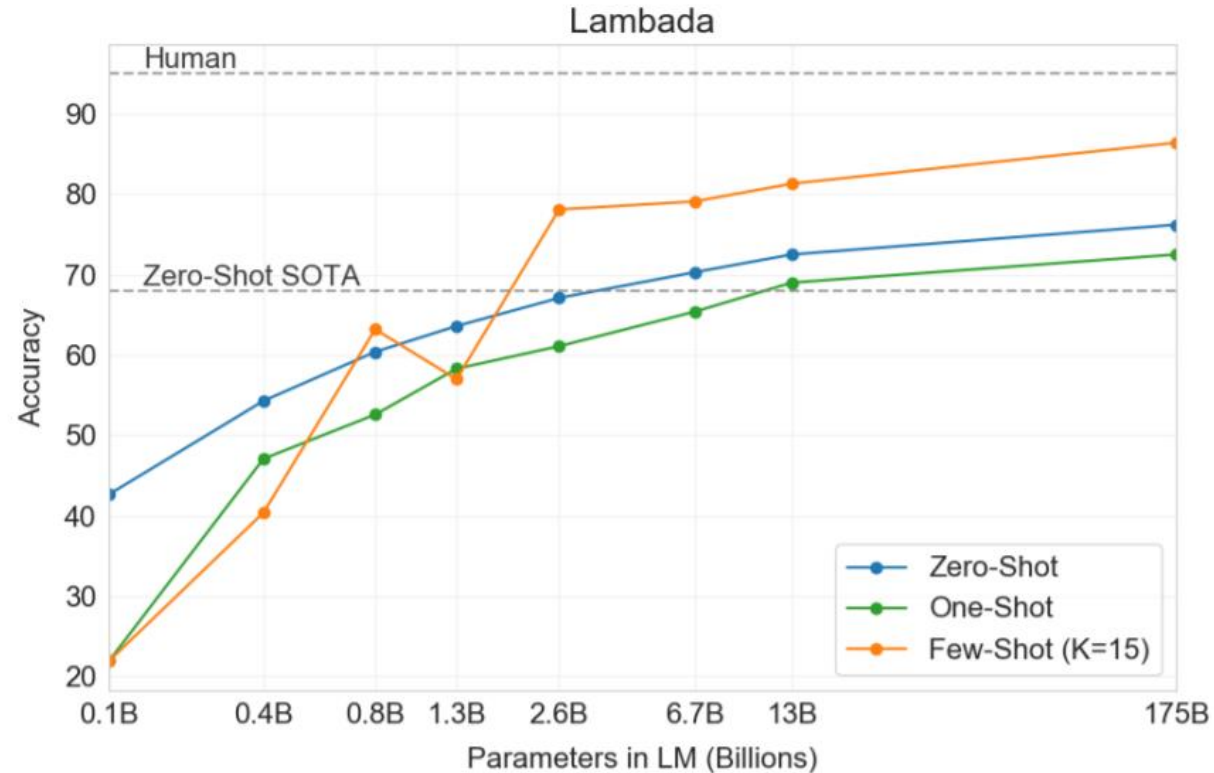1. Language Modeling, Cloze, and Completion Tasks



**Figure 3.2:** On LAMBADA, the few-shot capability of language models results in a strong boost to accuracy. GPT-3 2.7B outperforms the SOTA 17B parameter Turing-NLG [Tur20] in this setting, and GPT-3 175B advances the state of the art by 18%. Note zero-shot uses a different format from one-shot and few-shot as described in the text.

2. Closed Book Question Answering

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

**Table 3.3: Results on three Open-Domain QA tasks.** GPT-3 is shown in the few-, one-, and zero-shot settings, as compared to prior SOTA results for closed book and open domain settings. TriviaQA few-shot result is evaluated on the wiki split test server.

- Closed Book QA : 여분의 정보(지문)를 주지 않는 것
  - ✓ Ex) Q. 톨스토이의 대표적인 작품은? A. 전쟁과 평화

# GPT-3

3. Translation

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

**Table 3.4: Few-shot GPT-3 outperforms previous unsupervised NMT work by 5 BLEU when translating into English reflecting its strength as an English LM.** We report BLEU scores on the WMT'14 Fr↔En, WMT'16 De↔En, and WMT'16 Ro↔En datasets as measured by multi-bleu.perl with XLM's tokenization in order to compare most closely with prior unsupervised NMT work. SacreBLEU$^f$ [Pos18] results reported in Appendix H. Underline indicates an unsupervised or few-shot SOTA, bold indicates supervised SOTA with relative confidence. $^a$[EOAG18] $^b$[DHKH14] $^c$[WXH$^+$18] $^d$[oR16] $^e$[LGG$^+$20] $^f$[SacreBLEU signature: BLEU+case.mixed+numrefs.1+smooth.exp+tok.intl+version.1.2.20]

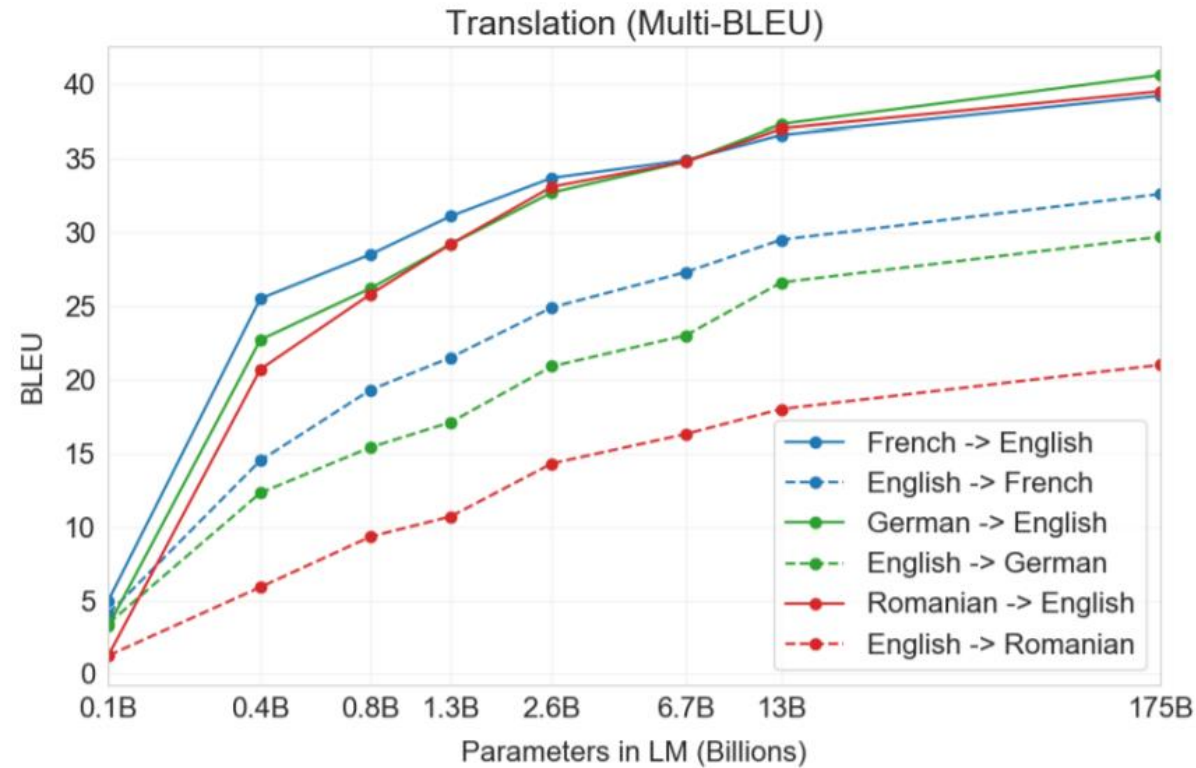# GPT-3

3. Translation



**Figure 3.4:** Few-shot translation performance on 6 language pairs as model capacity increases. There is a consistent trend of improvement across all datasets as the model scales, and as well as tendency for translation into English to be stronger than translation from English.

# GPT-3

4. Winograd-Style Tasks

| Setting | Winograd | Winogrande (XL) |
|---|---|---|
| Fine-tuned SOTA | **90.1**[a] | **84.6**[b] |
| GPT-3 Zero-Shot | 88.3* | 70.2 |
| GPT-3 One-Shot | 89.7* | 73.2 |
| GPT-3 Few-Shot | 88.6* | 77.7 |

**Table 3.5:** Results on the WSC273 version of Winograd schemas and the adversarial Winogrande dataset. See Section 4 for details on potential contamination of the Winograd test set. [a][SBBC19] [b][LYN+20]
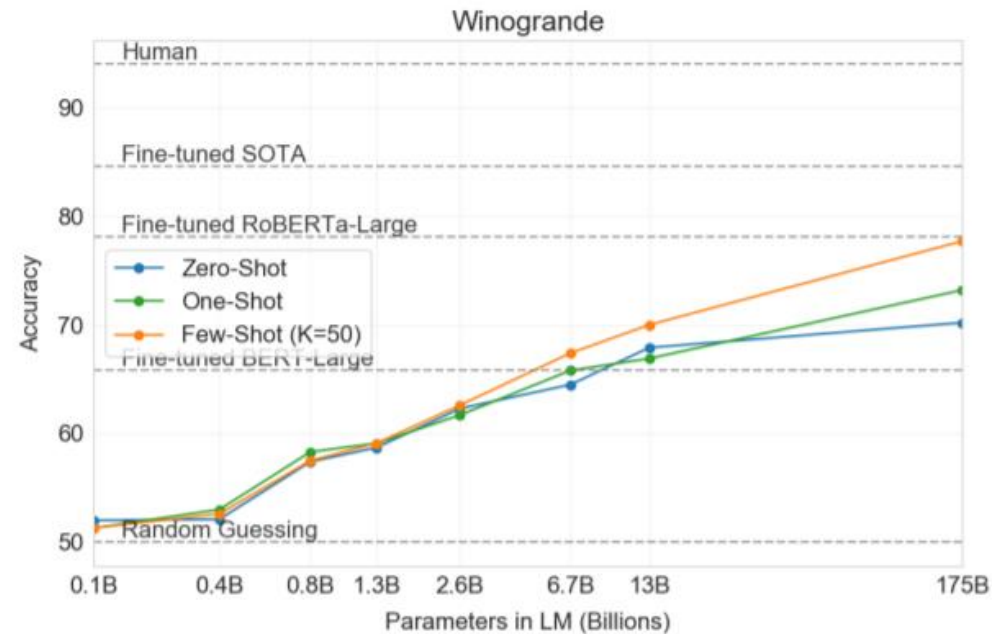


**Figure 3.5:** Zero-, one-, and few-shot performance on the adversarial Winogrande dataset as model capacity scales. Scaling is relatively smooth with the gains to few-shot learning increasing with model size, and few-shot GPT-3 175B is competitive with a fine-tuned RoBERTA-large.

# GPT-3

Result

5. Common Sense Reasoning

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

**Table 3.6:** GPT-3 results on three commonsense reasoning tasks, PIQA, ARC, and OpenBookQA. GPT-3 Few-Shot PIQA result is evaluated on the test server. See Section 4 for details on potential contamination issues on the PIQA test set.
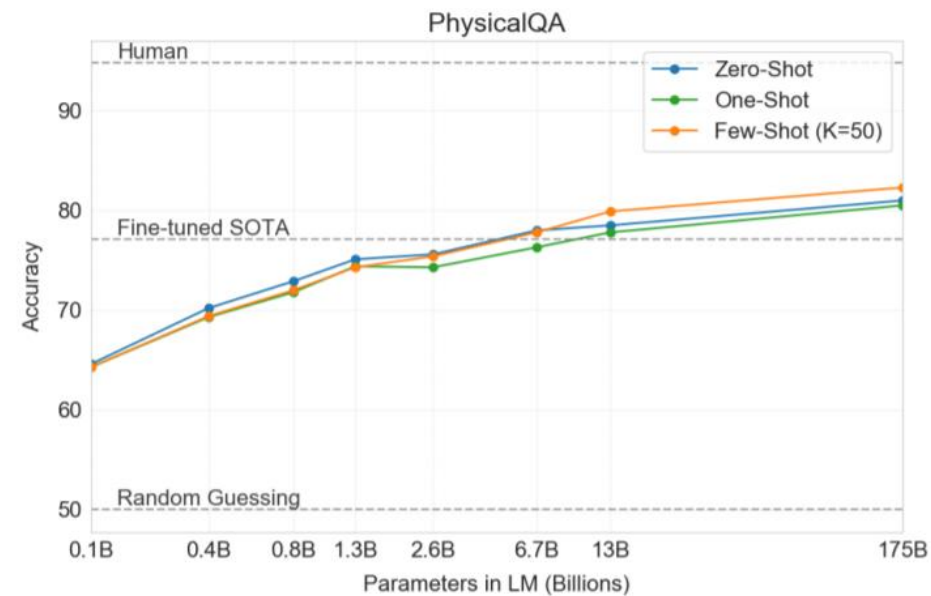


**Figure 3.6:** GPT-3 results on PIQA in the zero-shot, one-shot, and few-shot settings. The largest model achieves a score on the development set in all three conditions that exceeds the best recorded score on the task.

한양대학교 인공지능연구실                                                                                          16

# GPT-3

6. Reading Comprehension

| Setting | CoQA | DROP | QuAC | SQuADv2 | RACE-h | RACE-m |
|---|---|---|---|---|---|---|
| Fine-tuned SOTA | **90.7**[a] | **89.1**[b] | **74.4**[c] | **93.0**[d] | **90.0**[e] | **93.1**[e] |
| 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 |

**Table 3.7:** Results on reading comprehension tasks. All scores are F1 except results for RACE which report accuracy.
[a][JZC+19] [b][JN20] [c][AI19] [d][QIA20] [e][SPP+19]

7. SuperGLUE

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

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

**Table 3.8:** Performance of GPT-3 on SuperGLUE compared to fine-tuned baselines and SOTA. All results are reported on the test set. GPT-3 few-shot is given a total of 32 examples within the context of each task and performs no gradient updates.
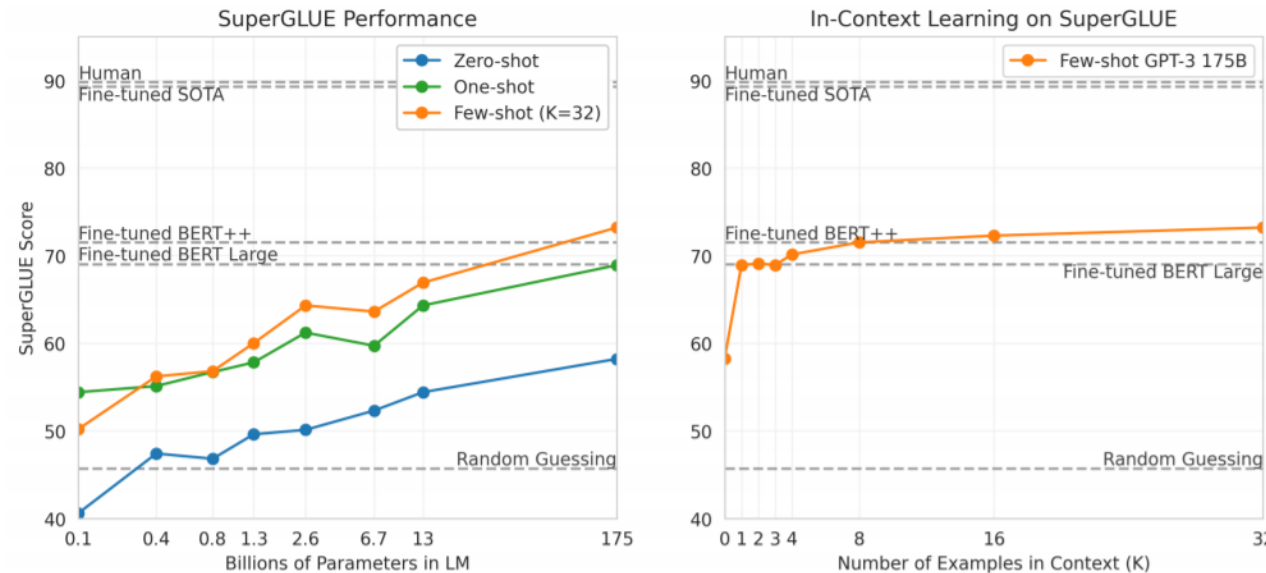


**Figure 3.8: Performance on SuperGLUE increases with model size and number of examples in context.** A value of $K = 32$ means that our model was shown 32 examples per task, for 256 examples total divided across the 8 tasks in SuperGLUE. We report GPT-3 values on the dev set, so our numbers are not directly comparable to the dotted reference lines (our test set results are in Table 3.8). The BERT-Large reference model was fine-tuned on the SuperGLUE training set (125K examples), whereas BERT++ was first fine-tuned on MultiNLI (392K examples) and SWAG (113K examples) before further fine-tuning on the SuperGLUE training set (for a total of 630K fine-tuning examples). We find the difference in performance between the BERT-Large and BERT++ to be roughly equivalent to the difference between GPT-3 with one example per context versus eight examples per context.

# GPT-3
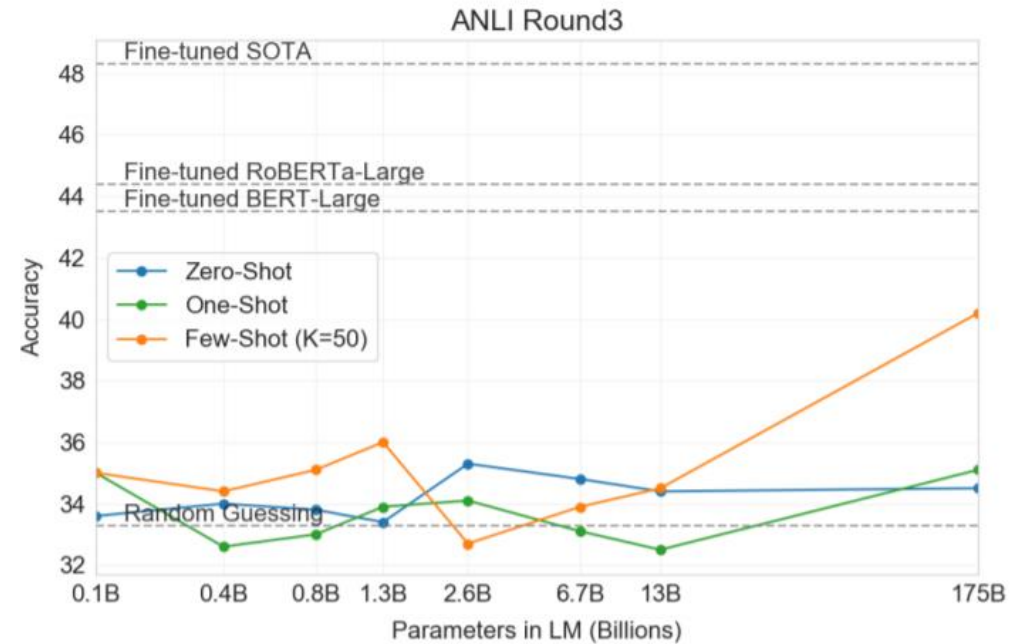
8. NLI



**Figure 3.9: Performance of GPT-3 on ANLI Round 3.** Results are on the dev-set, which has only 1500 examples and therefore has high variance (we estimate a standard deviation of 1.2%). We find that smaller models hover around random chance, while few-shot GPT-3 175B closes almost half the gap from random chance to SOTA. Results for ANLI rounds 1 and 2 are shown in the appendix.
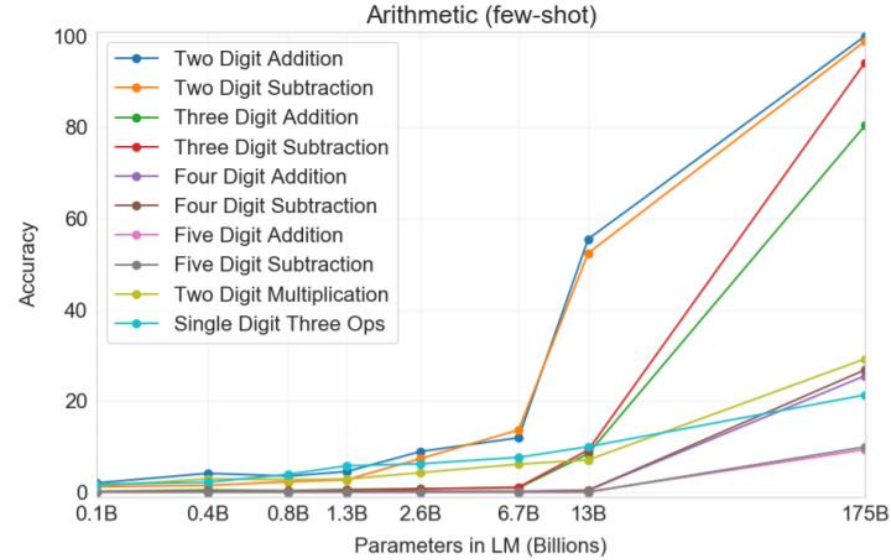
9. Synthetic and Qualitative tasks



**Figure 3.10:** Results on all 10 arithmetic tasks in the few-shot settings for models of different sizes. There is a significant jump from the second largest model (GPT-3 13B) to the largest model (GPT-3 175), with the latter being able to reliably accurate 2 digit arithmetic, usually accurate 3 digit arithmetic, and correct answers a significant fraction of the time on 4-5 digit arithmetic, 2 digit multiplication, and compound operations. Results for one-shot and zero-shot are shown in the appendix.

| 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 |

**Table 3.9:** Results on basic arithmetic tasks for GPT-3 175B. {2,3,4,5}D{+,-} is 2, 3, 4, and 5 digit addition or subtraction, 2Dx is 2 digit multiplication. 1DC is 1 digit composite operations. Results become progressively stronger moving from the zero-shot to one-shot to few-shot setting, but even the zero-shot shows significant arithmetic abilities.
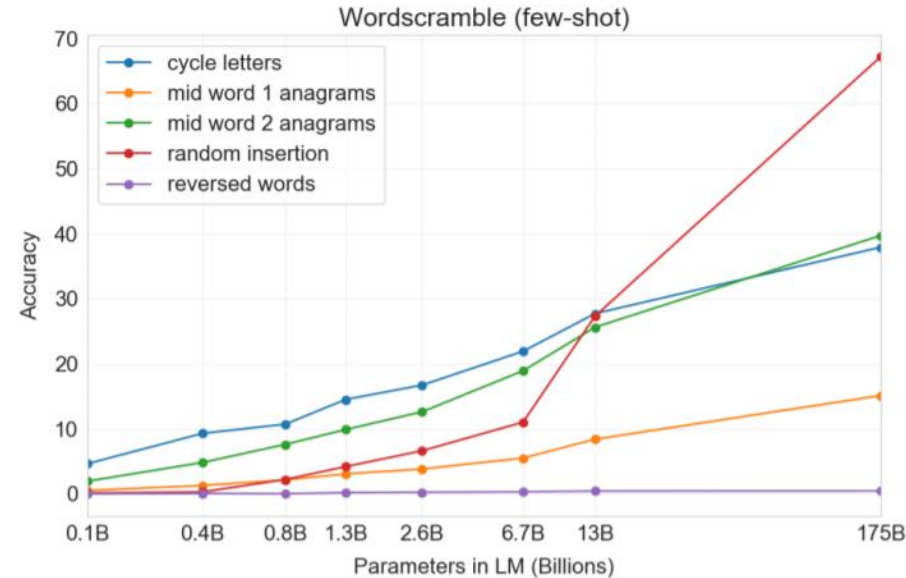
# GPT-3

9. Synthetic and Qualitative tasks



**Figure 3.11:** Few-shot performance on the five word scrambling tasks for different sizes of model. There is generally smooth improvement with model size although the random insertion task shows an upward slope of improvement with the 175B model solving the task the majority of the time. Scaling of one-shot and zero-shot performance is shown in the appendix. All tasks are done with $K = 100$.

- **Cycle letters in word (CL)** – The model is given a word with its letters cycled, then the "=" symbol, and is expected to generate the original word. For example, it might be given "lyinevitab" and should output "inevitably".

- **Anagrams of all but first and last characters (A1)** – The model is given a word where every letter except the first and last have been scrambled randomly, and must output the original word. Example: criroptuon = corruption.

- **Anagrams of all but first and last 2 characters (A2)** – The model is given a word where every letter except the first 2 and last 2 have been scrambled randomly, and must recover the original word. Example: opoepnnt → opponent.

- **Random insertion in word (RI)** – A random punctuation or space character is inserted between each letter of a word, and the model must output the original word. Example: s.u!c/c!e.s s i/o/n = succession.

- **Reversed words (RW)** – The model is given a word spelled backwards, and must output the original word. Example: stcejbo → objects.

# GPT-3

9. Synthetic and Qualitative tasks

News article generation

| | 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$e$-4) | 4.9% |
| GPT-3 Medium | 61% | 58%–65% | 10.3 (7$e$-21) | 6.0% |
| GPT-3 Large | 68% | 64%–72% | 7.3 (3$e$-11) | 8.7% |
| GPT-3 XL | 62% | 59%–65% | 10.7 (1$e$-19) | 7.5% |
| GPT-3 2.7B | 62% | 58%–65% | 10.4 (5$e$-19) | 7.1% |
| GPT-3 6.7B | 60% | 56%–63% | 11.2 (3$e$-21) | 6.2% |
| GPT-3 13B | 55% | 52%–58% | 15.3 (1$e$-32) | 7.1% |
| GPT-3 175B | 52% | 49%–54% | 16.9 (1$e$-34) | 7.8% |

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

9. Synthetic and Qualitative

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

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

# GPT-3

1. Text synthesis and several NLP task에 약하다

    - Repeat, lose coherence, contradict

2. Structural and Algorithmic 한계

    - Bidirectional task에서 낮은 성능 - fill-in-the-blank tasks

3. Limits of the pretraining objective

    - 여전히 General한 model이라고 하기 어렵다

4. Gender & Race 문제

.

.

.

# The End