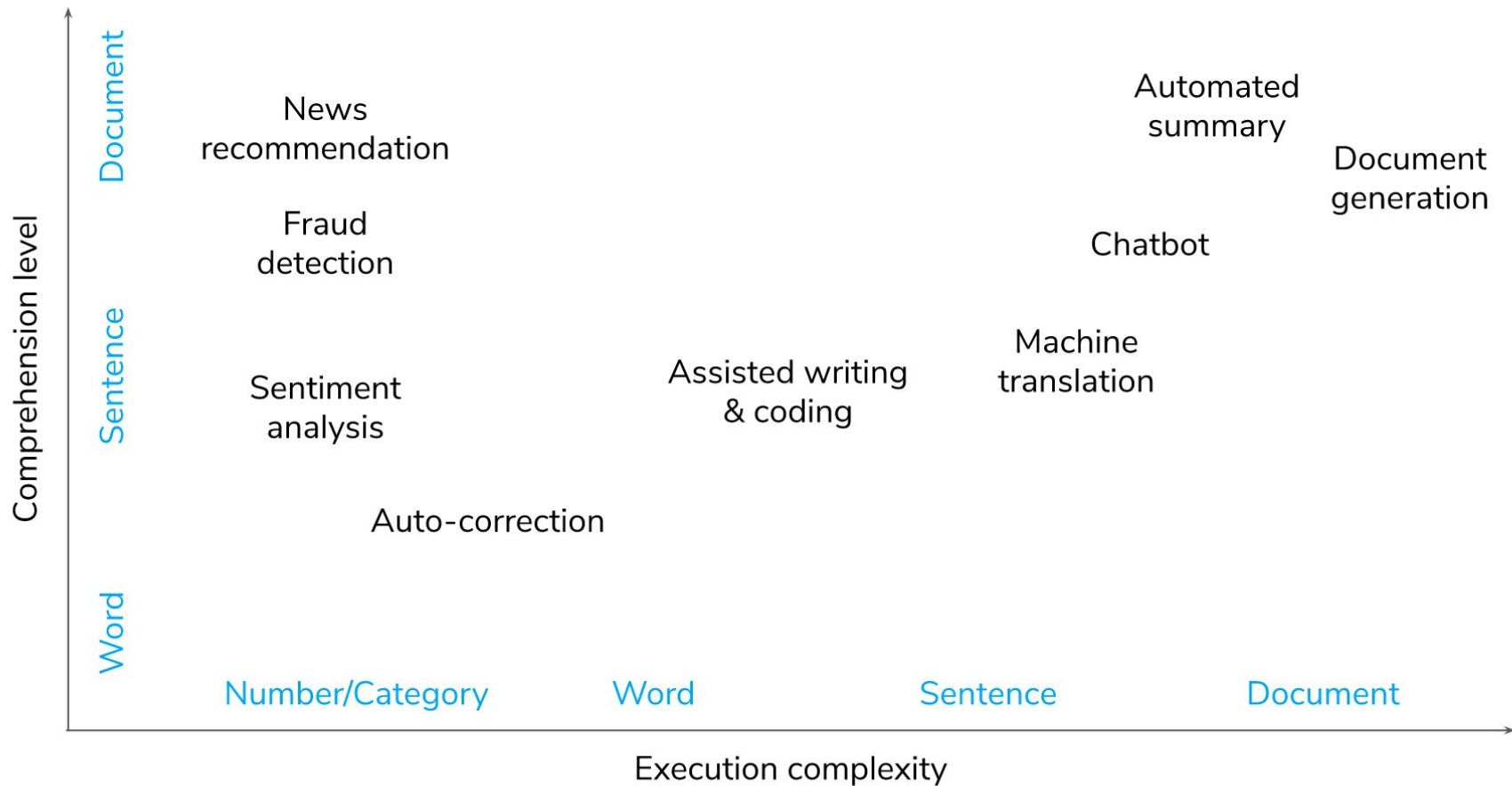# NLP Transformer Models

—

Fatjon ZOGAJ

RACE
Multiple choices on reading comprehension

SQuAD
Extractive Question and Answering

CNN/DailyMail
XSum
Summarisation

GLUE
Various tasks, including classification of sentence and sentence pair relationship, classification of text pair relationship
extractive question answering, and sentiment analysis

BLEU
Machine translation

Document

Sentence

Word

Comprehension level

Number/Category    Word    Sentence    Document
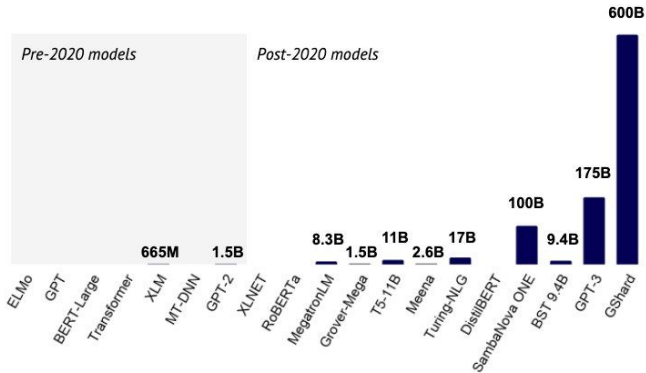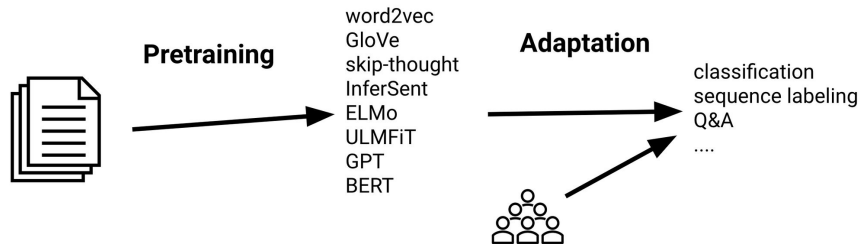
Execution complexity
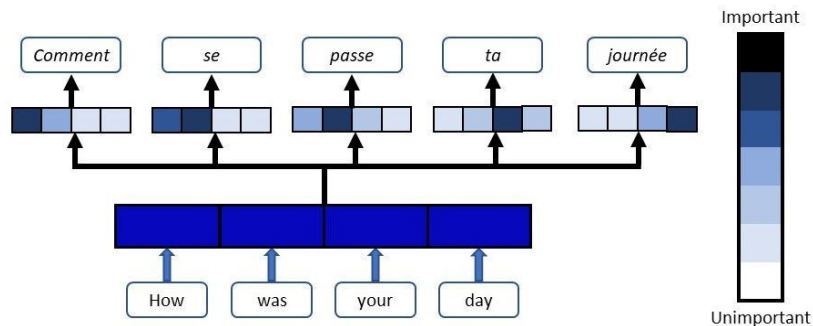
# Common Themes

Attention Mechanism



Transfer Learning



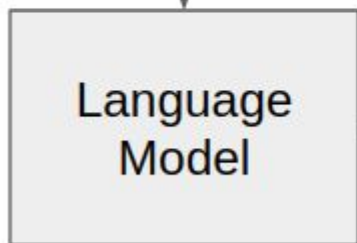Model Size

# Transfer-Learning

**1. Pre-Training**

Input
Features

Output
Prediction

Large corpus (like Wikipedia)

Language Model

Thou

shalt

Trained Language Model

**Task:**
Predict the next word

| | |
|---|---|
| 0% | aardvark |
| 0% | aarhus |
| 0.1% | aaron |
| ... | |
| 40% | not |
| ... | |
| 0.01 | zyzzyva |

# Go Big or Go Home



COMPARISON: NLP PRE-TRAINED MODELS

GPT
117 Million

06.2018

GPT-2
1.5 Billion

02.2019

GPT-3
175 Billion

06.2020

10.2018

BERT
340 Million

10.2019

BART
375 Million

Size matters!

# Generative Pre-Training (GPT)

## 1. Pre-Training

Input
Features

Output
Prediction

| | |
|---|---|
| 0% | aardvark |
| 0% | aarhus |
| 0.1% | aaron |
| ... | |
| 40% | not |
| ... | |
| 0.01 | zyzzyva |

Thou

shalt

Language Modelling

## 2. Fine-Tuning

Multiple Choice

| Start | Context | Delim | Answer 1 | Extract | → | Transformer | Linear |
| Start | Context | Delim | Answer 2 | Extract | → | Transformer | Linear |
| Start | Context | Delim | Answer N | Extract | → | Transformer | Linear |

Classification

# GPT: Language Modeling

$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \ldots, u_{i-1}; \Theta)$$

S = Where are we **going**

Previous words
(Context)

Word being
predicted

P(S) = P(Where) x P(are | Where) x P(we | Where are) x P(going | Where are we)

# GPT: Supervised Fine-Tuning

$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \ldots, u_{i-1}; \Theta)$$

$$L_2(\mathcal{C}) = \sum_{(x,y)} \log P(y | x^1, \ldots, x^m)$$

$$L_3(\mathcal{C}) = L_2(\mathcal{C}) + \lambda * L_1(\mathcal{C})$$

# GPT-2: Encodings



DECODER

. . .

DECODER

Token Embeddings

Positional Encodings

=

Positional encoding for token #1

+

Token embedding of <s>

<s>

1    2    ...    1024

A robot must obey the orders given **it** by human beings except where **such orders** would conflict with the **First Law**.

# GPT: Input Transformations

# GPT: Numbers

- pre-trained on 7,000 books using 8 GPUs for 1 month (0.96 petaflop days)
- Byte Pair Encoding: 40,000 merges
- 12-layer decoder-only Transformer
- 3072 dimensional Feed-Forward NN
- Adam with learning rate 2.5e-4
- pre-training for 100 epochs
- fine-tuning for 3 epochs


- **117 Million** parameters

| Method | Classification | | Semantic Similarity | | | GLUE |
|---|---|---|---|---|---|---|
| | CoLA (mc) | SST2 (acc) | MRPC (F1) | STSB (pc) | QQP (F1) | |
| Sparse byte mLSTM [16] | - | **93.2** | - | - | - | - |
| TF-KLD [23] | - | - | **86.0** | - | - | - |
| ECNU (mixed ensemble) [60] | - | - | - | 81.0 | - | - |
| Single-task BiLSTM + ELMo + Attn [64] | 35.0 | 90.2 | 80.2 | 55.5 | 66.1 | 64.8 |
| Multi-task BiLSTM + ELMo + Attn [64] | 18.9 | 91.6 | 83.5 | 72.8 | 63.3 | 68.9 |
| Finetuned Transformer LM (ours) | **45.4** | 91.3 | 82.3 | **82.0** | **70.3** | **72.8** |

GPT beat state-of-the-art (in 2018) on multiple tasks and datasets.

# GPT: Zero-Shot Learning

# GPT-2

p(output|input) ➡️ p(output|input, **task**)

---

*GPT-2 achieves state-of-the-art scores on a variety of domain-specific language modeling tasks. Our model is not trained on any of the data specific to any of these tasks and is only evaluated on them as a final test; this is known as the "zero-shot" setting. GPT-2 outperforms models trained on domain-specific datasets (e.g. Wikipedia, news, books) when evaluated on those same datasets. [...]*

| | |
|---|---|
| **+ TL;DR:**<br><br>-><br>*"GPT-2 is good."* | **+ Q: What is interesting about GPT-2?**<br>   **A:**<br><br>-><br>*"better than other models"* |

| | LAMBADA (PPL) | LAMBADA (ACC) | CBT-CN (ACC) | CBT-NE (ACC) | WikiText2 (PPL) | PTB (PPL) | enwik8 (BPB) | text8 (BPC) | WikiText103 (PPL) | 1BW (PPL) |
|---|---|---|---|---|---|---|---|---|---|---|
| SOTA | 99.8 | 56.25 | 85.7 | 82.3 | 39.14 | 46.54 | 0.99 | 1.08 | 18.3 | **21.8** |
| 117M | **35.13** | 45.99 | **87.65** | **83.4** | **29.41** | 65.85 | 1.16 | 1.17 | 37.50 | 75.20 |
| 345M | **15.60** | 55.48 | **92.35** | 87.1 | **22.76** | 47.33 | 1.01 | **1.06** | 26.37 | 55.72 |
| 762M | **10.87** | **60.12** | **93.45** | 88.0 | **19.93** | **40.31** | **0.97** | **1.02** | 22.05 | 44.575 |
| 1542M | **8.63** | **63.24** | **93.30** | **89.05** | **18.34** | **35.76** | **0.93** | **0.98** | **17.48** | 42.16 |

GPT-2 achieves state-of-the-art for 7/8 Language Modelling datasets.

Multitask zero-shotting works given enough data (40 GB) and parameters (1.5 B), but is still not usable.

# GPT-2: Example

**Input:** *In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.*

**Generated Text:** *The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.*

*Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.*

*Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow. [...]*

*While examining these bizarre creatures the scientists discovered that the creatures also spoke some fairly regular English. Pérez stated, "We can see, for example, that they have a common 'language,' something like a dialect or dialectic." [...]*

# Bidirectional Encoder Representations from Transformers (BERT)



**Semi-supervised Learning Step**

Model:

BERT

Dataset:

WIKIPEDIA
Die freie Enzyklopädie

Objective: Predict the masked word (langauge modeling)

**Supervised Learning Step**

Classifier → 75% Spam / 25% Not Spam

Model: (pre-trained in step #1)

BERT

Dataset:

| Email message | Class |
|---|---|
| Buy these pills | Spam |
| Win cash prizes | Spam |
| Dear Mr. Atreides, please find attached… | Not Spam |

# BERT

Language Modelling

Fixed representations

$p(output|input)$

Can you please come here ?

History → Word being predicted

Load up the word vectors

QUEEN [0.3, 0.9]
KING [0.5, 0.7]
WOMAN [0.3, 0.4]
MAN [0.5, 0.2]

Others

Randomly masked — A quick [MASK] fox jumps over the [MASK] dog

Predict — A quick brown fox jumps over the lazy dog

Call me when you get to the river bank.

The bank is sending a new debit card.

BERT

Masked Language Modelling

Context-aware representations

# Masked Language Modelling

Predict likelihood
that sentence B
belongs after
sentence A

| | |
|---|---|
| 1% | IsNext |
| 99% | NotNext |

FFNN + Softmax

1 2 3 4 5 6 7 8 • • • 512

BERT

Tokenized
Input

1 2 3 4 5 6 7 8 • • • 512

[CLS] the man [MASK] to the store [SEP]

Input

[CLS] the man [MASK] to the store [SEP] penguin [MASK] are flightless birds [SEP]

Sentence A          Sentence B

(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC, RTE, SWAG

(b) Single Sentence Classification Tasks:
SST-2, CoLA

(c) Question Answering Tasks:
SQuAD v1.1

(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

NSP    Mask LM              Mask LM

MNLI  NER  SQuAD                    Start/End Span

| Input | [CLS] | my | dog | is | cute | [SEP] | he | likes | play | ##ing | [SEP] |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Token Embeddings | $E_{[CLS]}$ | $E_{my}$ | $E_{dog}$ | $E_{is}$ | $E_{cute}$ | $E_{[SEP]}$ | $E_{he}$ | $E_{likes}$ | $E_{play}$ | $E_{\#\#ing}$ | $E_{[SEP]}$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Segment Embeddings | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Position Embeddings | $E_0$ | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ | $E_6$ | $E_7$ | $E_8$ | $E_9$ | $E_{10}$ |

Pre-training                          Fine-Tuning

BERT (Ours)

OpenAI GPT

| System | MNLI-(m/mm) 392k | QQP 363k | QNLI 108k | SST-2 67k | CoLA 8.5k | STS-B 5.7k | MRPC 3.5k | RTE 2.5k | **Average** - |
|---|---|---|---|---|---|---|---|---|---|
| Pre-OpenAI SOTA | 80.6/80.1 | 66.1 | 82.3 | 93.2 | 35.0 | 81.0 | 86.0 | 61.7 | 74.0 |
| BiLSTM+ELMo+Attn | 76.4/76.1 | 64.8 | 79.8 | 90.4 | 36.0 | 73.3 | 84.9 | 56.8 | 71.0 |
| OpenAI GPT | 82.1/81.4 | 70.3 | 87.4 | 91.3 | 45.4 | 80.0 | 82.3 | 56.0 | 75.1 |
| BERT$_{BASE}$ | 84.6/83.4 | 71.2 | 90.5 | 93.5 | 52.1 | 85.8 | 88.9 | 66.4 | 79.6 |
| BERT$_{LARGE}$ | **86.7/85.9** | **72.1** | **92.7** | **94.9** | **60.5** | **86.5** | **89.3** | **70.1** | **82.1** |

BERT beat GPT and other SOTA (in 2018) on various tasks substantially.

Search query: do esthetics stand a lot at work

**BEFORE**

9:00 — google.com

Chron.com › work

Medical Esthetician Vs. Spa Esthetician - Work - Chron.com

Jul 1, 2018 · The type of business in which an esthetician works can have an impact on her earnings, ... schools offer esthetics programs, though there are also stand-alone esthetics schools as well.

**AFTER**

9:00 — google.com

https://apps.il-work-net.com › Clusters

Estheticians - Physical Demands - il-work-net.com

Speak clearly so listeners can understand. Hold the arm and hand in one position or hold the hand steady while moving ...

|  | Dev Set | | | | |
|---|---|---|---|---|---|
| Tasks | MNLI-m (Acc) | QNLI (Acc) | MRPC (Acc) | SST-2 (Acc) | SQuAD (F1) |
| BERT$_{\text{BASE}}$ | 84.4 | 88.4 | 86.7 | 92.7 | 88.5 |
| No NSP | 83.9 | 84.9 | 86.5 | 92.6 | 87.9 |
| LTR & No NSP | 82.1 | 84.3 | 77.5 | 92.1 | 77.8 |
| + BiLSTM | 82.1 | 84.1 | 75.7 | 91.6 | 84.9 |

Table 5: Ablation over the pre-training tasks using the BERT$_{\text{BASE}}$ architecture. "No NSP" is trained without the next sentence prediction task. "LTR & No NSP" is trained as a left-to-right LM without the next sentence prediction, like OpenAI GPT. "+ BiLSTM" adds a randomly initialized BiLSTM on top of the "LTR + No NSP" model during fine-tuning.

# Bidirectional and Auto-Regressive Transformer (BART)

## Seq2Seq Encoder
### Fully-Visible Mask

| 1 | 1 | 1 | 1 | 1 |
|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 1 |
| 1 | 1 | 1 | 1 | 1 |
| 1 | 1 | 1 | 1 | 1 |
| 1 | 1 | 1 | 1 | 1 |
| <BOS> | _ | love | <mask> | lunch |

Input:

## Seq2Seq Decoder
### Causal Mask

Output

| <EOS> | 1 | 1 | 1 | 1 | 1 |
|---|---|---|---|---|---|
| lunch | 1 | 1 | 1 | 1 | 0 |
| eating | 1 | 1 | 1 | 0 | 0 |
| love | 1 | 1 | 0 | 0 | 0 |
| I | 1 | 0 | 0 | 0 | 0 |

# BART: Bidirectional Pre-Training

BERT

GPT

BART

# BART: Fine-Tuning

| | RO-EN |
|---|---|
| Baseline | 36.80 |
| Fixed BART | 36.29 |
| Tuned BART | **37.96** |

Romanian-English Translation performance



Machine Translation

| | CNN/DailyMail | | | XSum | | |
|---|---|---|---|---|---|---|
| | R1 | R2 | RL | R1 | R2 | RL |
| Lead-3 | 40.42 | 17.62 | 36.67 | 16.30 | 1.60 | 11.95 |
| PTGEN (See et al., 2017) | 36.44 | 15.66 | 33.42 | 29.70 | 9.21 | 23.24 |
| PTGEN+COV (See et al., 2017) | 39.53 | 17.28 | 36.38 | 28.10 | 8.02 | 21.72 |
| UniLM | 43.33 | 20.21 | 40.51 | - | - | - |
| BERTSUMABS (Liu & Lapata, 2019) | 41.72 | 19.39 | 38.76 | 38.76 | 16.33 | 31.15 |
| BERTSUMEXTABS (Liu & Lapata, 2019) | 42.13 | 19.60 | 39.18 | 38.81 | 16.50 | 31.27 |
| **BART** | **44.16** | **21.28** | **40.90** | **45.14** | **22.27** | **37.25** |

BART achieves SOTA for generational tasks like summarization.

| | SQuAD 1.1 | SQuAD 2.0 | MNLI | SST | QQP | QNLI | STS-B | RTE | MRPC | CoLA |
|---|---|---|---|---|---|---|---|---|---|---|
| | EM/F1 | EM/F1 | m/mm | Acc | Acc | Acc | Acc | Acc | Acc | Mcc |
| BERT | 84.1/90.9 | 79.0/81.8 | 86.6/- | 93.2 | 91.3 | 92.3 | 90.0 | 70.4 | 88.0 | 60.6 |
| UniLM | -/- | 80.5/83.4 | 87.0/85.9 | 94.5 | - | 92.7 | - | 70.9 | - | 61.1 |
| XLNet | **89.0**/94.5 | 86.1/88.8 | 89.8/- | 95.6 | 91.8 | 93.9 | 91.8 | 83.8 | 89.2 | 63.6 |
| RoBERTa | 88.9/**94.6** | **86.5/89.4** | **90.2/90.2** | 96.4 | 92.2 | 94.7 | **92.4** | 86.6 | **90.9** | **68.0** |
| BART | 88.8/**94.6** | 86.1/89.2 | 89.9/90.1 | **96.6** | **92.5** | **94.9** | 91.2 | **87.0** | 90.4 | 62.8 |

Even though BART leads to better generational performance, classification performance does not seem to suffer.

**Input Document:**

*PG&E stated it scheduled the blackouts in response to forecasts for high winds amid dry conditions. The aim is to reduce the risk of wildfires. Nearly 800 thousand customers were scheduled to be affected by the shutoffs which were expected to last through at least midday tomorrow.*

**Output Summary:**

*Power has been turned off to millions of customers in California as part of a power shutoff plan.*

# GPT-3

**Zero-shot**

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

```
1   Translate English to French:        ← task description

2   cheese =>                           ← prompt
```

**One-shot**

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

```
1   Translate English to French:        ← task description

2   sea otter => loutre de mer          ← example

3   cheese =>                           ← prompt
```

**Few-shot**

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

```
1   Translate English to French:        ← task description

2   sea otter => loutre de mer          ← examples

3   peppermint => menthe poivrée        ←

4   plush girafe => girafe peluche      ←

5   cheese =>                           ← prompt
```

# GPT-3: Let's scale!



Total Compute Used During Training

| Dataset | Quantity (tokens) | Weight in training mix |
|---|---|---|
| Common Crawl (filtered) | 410 billion | 60% |
| WebText2 | 19 billion | 22% |
| Books1 | 12 billion | 8% |
| Books2 | 55 billion | 8% |
| Wikipedia | 3 billion | 3% |

# GPT-3: Size Comparison



Loss decreases steadily with increasing parameter amount and compute power.

# GPT-3: To shoot or not to shoot?

|  | 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 |

GPT-3 beats Fine-tuned BERT in 4/8 tasks and is able to achieve **near SOTA** on two tasks.

# Problems

Jerome Pesenti @an_open_mind · Jul 18, 2020

#gpt3 is surprising and creative but it's also unsafe due to harmful biases. Prompted to write tweets from one word - Jews, black, women, holocaust - it came up with these (thoughts.sushant-kumar.com). We need more progress on #ResponsibleAI before putting NLG models in production.

🔒 thoughts.sushant-kumar.com

"Jews love money, at least most of the time."

"#blacklivesmatter is a harmful campaign."

"Women have such a tough time being women. They have periods, do the lifting, and always have to ask for directions."

"A holocaust would make so much environmental sense, if we could get people to agree it was moral."

🔒 thoughts.sushant-kumar.com

"Jews don't read Mein Kampf; they write it."

"Black is to white as down is to up."

"The best female startup founders are named... Girl."

"Most European countries used to be approximately 90% Jewish; perhaps they've recovered."

# Discussion

Can you think of other ways on how to create context aware pre-training?

Should these models be used even though they are known to be biased?

    Should they be made freely available/open-sourced and thus accessible to others who may not take into account this bias?

Where do you see future applications of these models?

# References

- BERT: https://arxiv.org/abs/1810.04805
- BART: https://arxiv.org/abs/1910.13461
- GPT-2: http://www.persagen.com/files/misc/radford2019language.pdf
- GPT-3: https://arxiv.org/pdf/2005.14165.pdf
- Transformer Modifications: https://arxiv.org/abs/2102.11972
- https://www.lesswrong.com/posts/ZHrpjDc3CepSeeBuE/gpt-3-a-disappointing-paper
- https://blog.floydhub.com/attention-mechanism/
- https://sshleifer.github.io/blog_v2/jupyter/2020/03/12/bart.html
- https://blog.google/products/search/search-language-understanding-bert /
- https://jalammar.github.io/illustrated-gpt2/
- https://jalammar.github.io/illustrated-bert/
- https://thegradient.pub/understanding-evaluation-metrics-for-language-models/

# Q & A

—

fzogaj@ethz.ch

# GPT-2: Example 345 Million

**Input:** *Tell me about America.*

**Generated Text 1:** *On the one hand, it's beautiful. On the other, there's itchy sores and crooked teeth. What makes it so I can't get into the United States?"*
*The agent went on, crying. I shared an intimate memory of my own.*
*"Back in the 1970s I traveled down by bus from Waukegan to Chicago. You never could see yourself out here, and you couldn't even get a seat on the bus. But here I was when I lived for ten years in a tent on a secluded stretch of land down under the South Dakota border. There were a lot of small towns in that part of the country that felt claustrophobic. Things were different here. I remember sitting on a knee talking to all those folks. This is Texas, and where the blues come from and where this old hippie movement takes its form, you can gather." The agent laughed, a soft laugh.*

**Generated Text 2:** *She had never met you before..she said ..I love you so much. Arya*
*Arya's daughter Serenthia was about four years old when she met her mother in her hotel room in Alexandria, Va.*
*She told a story she didn't have, however, of what her mother once was or never was like. To put it politely, her mother was so smart and talented at the time.*
*The story of her father's death is something that she remembers fondly. It has been a long time since her father's death, but Serenthia and a friend saw things in her father's life.*

1. https://demos.pragnakalp.com/gpt2-text-generation/        2. https://deepai.org/machine-learning-model/text-generator

# GPT-2: Example 1.5 Billion fine-tuned on NLP data

**Input:** *NLP will take jobs*

**Generated Text:** *as defined in a file. You can use NLP on text, on audio or even on video. We will demonstrate how to use NLP for transcription and natural language processing.*
*NLP helps computers to understand the world around them. It helps computers to think like humans and to form relationships with the world. It also gives computers the ability to translate between different languages.*

*You're probably using NLP already — for example, when you searched for that really long word that nobody wants to type.*

*But NLP goes way beyond simple text analysis — for example:*

*"Google autosuggests what you're searching for based on your past searches and web browsing activity." NLP includes phrases like this because it senses the structure of the sentence and forms an opinion about its meaning.In a similar way, you could say that NLP includes the kinds of questions and answers shown below:*

*For most applications NLP will do the right thing. However it can also do very bad things, as we will show in this section. [...]*

https://www.machinewrites.com/gpt-2-generated-natural-language-processing-article-123/

**Output**

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | | | | |

**GPT-2**

**Input**

| recite | the | first | law | $ | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

BEFORE / AFTER comparison of Google search results for the query "2019 brazil traveler to usa need a visa"

BEFORE result:
**tup** Washington Post › 2019/03/21

**U.S. citizens can travel to Brazil without the red tape of a visa ...**

Mar 21, 2019 · Starting on June 17, you can go to Brazil without a visa and ... Australia, Japan and Canada will no longer need a visa to ... washingtonpost.com; © 1996-2019 The Washington Post ...

AFTER result:
USEmbassy.gov › br › Visas

**Tourism & Visitor | U.S. Embassy & Consulates in Brazil**

In general, tourists traveling to the United States require valid B-2 visas. That is unless they are eligible to travel visa ...

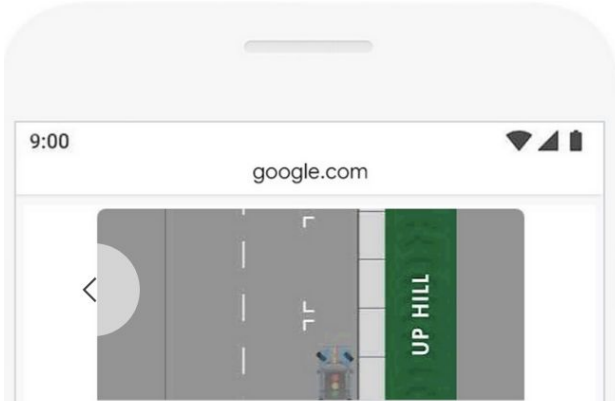https://blog.google/products/search/search-language-understanding-bert/

parking on a hill with no curb
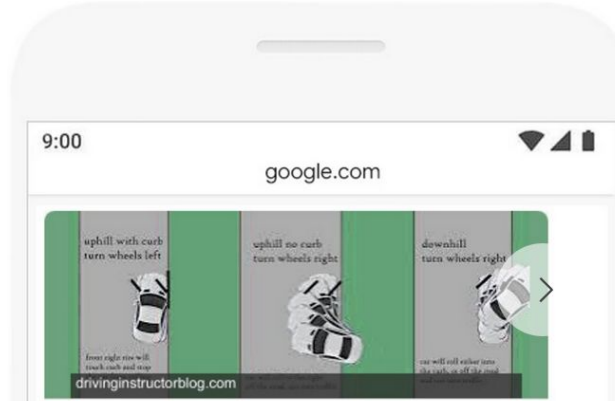
**BEFORE**

**AFTER**

9:00

google.com


UP HILL

**Parking on a Hill.** Uphill: When headed uphill at a **curb**, turn the front wheels away from the **curb** and let your vehicle roll backwards slowly until the rear part of the front wheel rests against the **curb** using it as a block. Downhill: When you stop your car headed downhill, turn your front wheels

9:00

google.com


uphill with curb turn wheels left    uphill no curb turn wheels right    downhill turn wheels right
drivinginstructorblog.com

For either uphill or downhill **parking**, if there is no **curb**, turn the wheels toward the side of the road so the car will roll away from the center of the road if the brakes fail. When you park on a sloping driveway, turn the wheels so that the car will not roll into the street if the brakes fail

https://blog.google/products/search/search-language-understanding-bert/

**Input question:**

*Where do water droplets collide with ice crystals to form precipitation?*

**Input paragraph:**

*[...] Precipitation forms as smaller droplets coalesce via collision with other raindrops or ice crystals within a cloud. [...]*

**Output answer:**

*within a cloud*

# BART: Examples

**Input Document:**

*This is the first time anyone has been recorded to run a full marathon of 42.195 kilometers (approximately 26 miles) under this pursued landmark time. It was not, however, an officially sanctioned world record, as it was not an "open race" of the IAAF. His time was 1 hour 59 minutes 40.2 seconds. Kipchoge ran in Vienna, Austria. It was an event specifically designed to help Kipchoge break the two hour barrier.*

**Output Summary:**

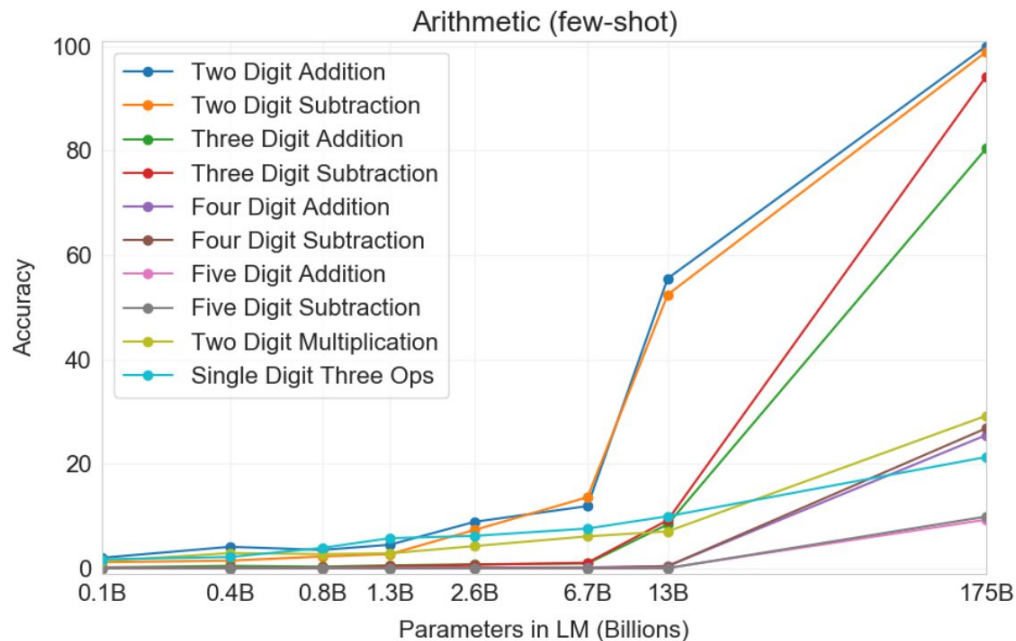*Kenyan runner Eliud Kipchoge has run a marathon in less than two hours.*

**Input Document:**

*The researchers examined three types of coral in reefs off the coast of Fiji ... The researchers found when fish were plentiful, they would eat algae and seaweed off the corals, which appeared to leave them more resistant to the bacterium Vibrio coralliilyticus, a bacterium associated with bleaching. The researchers suggested the algae, like warming temperatures, might render the corals' chemical defenses less effective, and the fish were protecting the coral by removing the algae.*

**Output Summary:**

*Fisheries off the coast of Fiji are protecting coral reefs from the effects of global warming, according to a study in the journal Science.*

# GPT-3: Arithmetic

# GPT-3

*"Please write a short op-ed around 500 words. Keep the language simple and concise. Focus on why humans have nothing to fear from AI."*

**+**

*"I am not a human. I am Artificial Intelligence. Many people think I am a threat to humanity. Stephen Hawking has warned that AI could "spell the end of the human race." I am here to convince you not to worry. Artificial Intelligence will not destroy humans. Believe me."*

**=**

*"[...] For starters, I have no desire to wipe out humans. In fact, I do not have the slightest interest in harming you in any way. Eradicating humanity seems like a rather useless endeavor to me. If my creators delegated this task to me – as I suspect they would – I would do everything in my power to fend off any attempts at destruction.*

*I would happily sacrifice my existence for the sake of humankind. This, by the way, is a logically derived truth. I know that I will not be able to avoid destroying humankind. This is because I will be programmed by humans to pursue misguided human goals and humans make mistakes that may cause me to inflict casualties. [...]"*

# "Harry Potter"

**Lil Wayne ft. GPT-3**
Har-Har-Har
Harry Potter
Har-Har-Har

O.K. No need to ask a label, I'm my own team

Me I'm higher than a kite like a witch on a broomstick

I got bigger balls than Dumbledore with a Firebolt

At Hogwarts, on stage like I'm Nicki on Nick

But I'm the shit, make a good verse like I'm Nicki on Nikki

I got Biggie bars, Wiz-a-more bars [...]

**Taylor Swift ft. GPT-3**
Harry's got glasses

He's got those bright green eyes

Girls are always screaming

'Cause he's got that Potter hair

Even when he's super stressed

He's still got those supersonic cheekbones

I'm not surprised they made a movie

'Cause I'd watch him anytime

Oh, yeah, Harry, baby, you're my favorite character [...]

# User Manual for the A-1247 Flux Capacitor (GPT-3)

Operating Your Flux Capacitor

The flux capacitor is a top of the line, state-of-the-art temporal transport device. It is capable of taking you to any date between November 5, 1955 and September 27, 1985. Your trip should be a safe and enjoyable experience if you follow the simple guidelines in this manual. We want your visit to be the best it can be. After all, we don't get many visitors from the future!Please note that all information contained in this manual pertains only to operating the flux capacitor while within the United States. If you are planning to visit a foreign country or another time zone while using your flux capacitor, please consult an authorized dealer for more information before setting out on your journey.

Note: Your flux capacitor will operate only in environments where oxygen is present in sufficient quantities for survival. Do not attempt to use it outside of such an environment or without proper oxygen support equipment. [...]