



Announcements

- ▶ Two more talks this semester:
 - ▶ Shinji Watanabe (CMU): Friday 11/4 11am GDC 6.302
 - ▶ Colin Raffel (UNC): Friday 11/18 11am GDC 6.302
- ▶ A5 due Tuesday



Recap: BERT



Today

- ▶ Seq2seq pre-trained models (BART, T5): how can we leverage the same kinds of ideas we saw in BERT for seq2seq models like machine translation?
- ▶ GPT-2/GPT-3: scaling language models further
- ▶ Prompting: a new way of using large language models without taking any gradient steps

Seq2seq Pre-trained Models: BART, T5



How do we pre-train seq2seq models?

- LMs $P(\mathbf{w})$: trained unidirectionally
- Masked LMs: trained bidirectionally but with masking
- How can we pre-train a model for $P(\mathbf{y}|\mathbf{x})$?
- Well, why was BERT effective?
 - Predicting a mask requires some kind of text “understanding”:
- What would it take to do the same for sequence prediction?

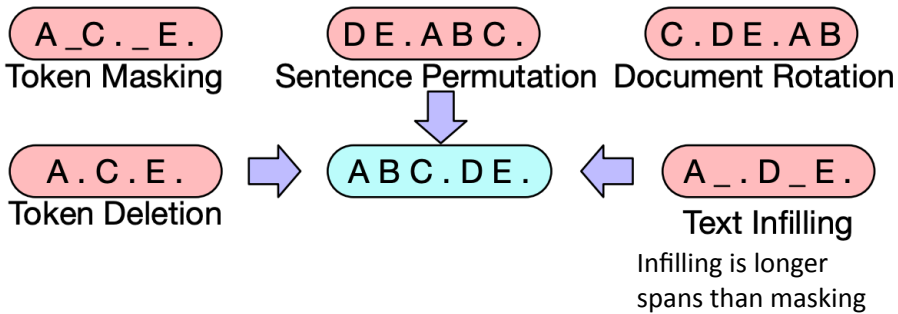


How do we pre-train seq2seq models?

- How can we pre-train a model for $P(\mathbf{y}|\mathbf{x})$?
- Requirements: (1) should use unlabeled data; (2) should force a model to attend from \mathbf{y} back to \mathbf{x}



BART



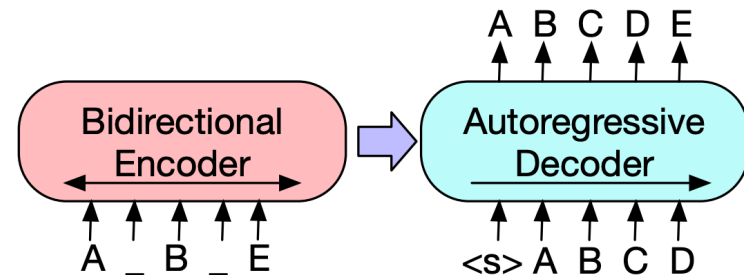
- Several possible strategies for corrupting a sequence are explored in the BART paper

Lewis et al. (2019)



BART

- Sequence-to-sequence Transformer trained on this data: permute/make/delete tokens, then predict full sequence autoregressively

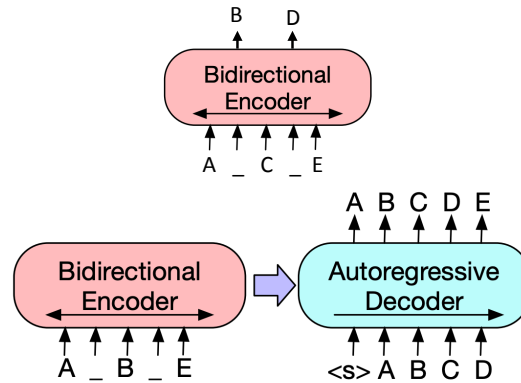


Lewis et al. (2019)



BERT vs. BART

- ▶ BERT: only parameters are an encoder, trained with masked language modeling objective. Cannot generate text or do seq2seq tasks
- ▶ BART: both an encoder and a decoder. Can also use just the encoder wherever we would use BERT



Lewis et al. (2019)



BART for Summarization

- ▶ **Pre-train** on the BART task: take random chunks of text, noise them according to the schemes described, and try to “decode” the clean text
- ▶ **Fine-tune** on a summarization dataset: a news article is the input and a summary of that article is the output (usually 1-3 sentences depending on the dataset)
- ▶ Can achieve good results even with **few summaries to fine-tune on**, compared to basic seq2seq models which require 100k+ examples to do well

Lewis et al. (2019)



BART for Summarization: Outputs

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.



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

Lewis et al. (2019)



BART for Summarization: Outputs

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.



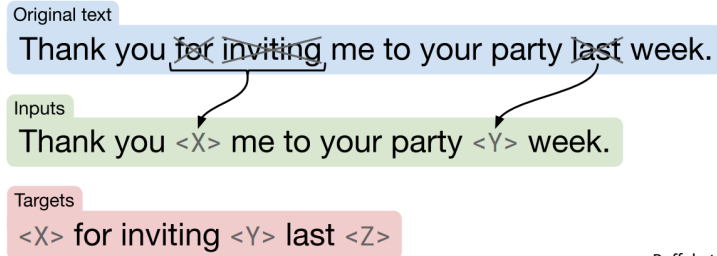
Power has been turned off to millions of customers in California as part of a power shutoff plan.

Lewis et al. (2019)



T5

- Pre-training: similar denoising scheme to BART (they were released within a week of each other in fall 2019)
- Input: text with gaps. Output: a series of phrases to fill those gaps.



Raffel et al. (2019)



T5

| Number of tokens | Repeats | summarization | | machine translation | | |
|------------------|---------|---------------|--------------|---------------------|--------------|--------------|
| | | GLUE | CNN/DM | EnDe | EnFr | EnRo |
| ★ Full dataset | 0 | 83.28 | 19.24 | 26.98 | 39.82 | 27.65 |
| 2 ²⁹ | 64 | 82.87 | 19.19 | 26.83 | 39.74 | 27.63 |
| 2 ²⁷ | 256 | 82.62 | 19.20 | 27.02 | 39.71 | 27.33 |
| 2 ²⁵ | 1,024 | 79.55 | 18.57 | 26.38 | 39.56 | 26.80 |
| 2 ²³ | 4,096 | 76.34 | 18.33 | 26.37 | 38.84 | 25.81 |

- Colossal Cleaned Common Crawl: 750 GB of text
- We still haven't hit the limit of bigger data being useful for pre-training: here we see stronger MT results from the biggest data

Raffel et al. (2019)



Successes of T5

- How can we handle a task like QA by framing it as a seq2seq problem?

| Dataset | SQuAD 1.1 |
|---------|--|
| Input | At what speed did the turbine operate? \n (Nikola_Tesla) On his 50th birthday in 1906, Tesla demonstrated his 200 horsepower (150 kilowatts) 16,000 rpm bladeless turbine. ... |
| Output | 16,000 rpm |

- Format: *Question* \n *Passage* → *Answer*
encoder decoder

Raffel et al. (2019)



UnifiedQA

| Dataset | NarrativeQA |
|----------|--|
| AB Input | What does a drink from narcissus's spring cause the drinker to do? \n Mercury has awakened Echo, who weeps for Narcissus, and states that a drink from Narcissus's spring causes the drinkers to ``Grow dotingly enamored of themselves.'' ... |
| Output | fall in love with themselves |

Abstractive question, requires generating *free-form answer*

- Past work: different architectures for every QA formulation. (Span selection, answer generation, multiple choice, ...)
- Now: one 11B parameter T5 model

Khashabi et al. (2020)



UnifiedQA

Multiple choice

| | | |
|----|----------------|---|
| MC | Dataset | MCTest |
| | Input | Who was Billy? \n (A) The skinny kid (B) A teacher (C) A little kid (D) The big kid \n Billy was like a king on the school yard. A king without a queen. He was the biggest kid in our grade, so he made all the rules during recess. ... |
| | Output | The big kid |
| YN | Dataset | BoolQ |
| | Input | Was America the first country to have a president? \n (President) The first usage of the word president to denote the highest official in a government was during the Commonwealth of England ... |
| | Output | no |

Yes/no

- ▶ Past work: different architectures for every QA formulation. (Span selection, answer generation, multiple choice, ...)
- ▶ Now: one 11B parameter T5 model

Khashabi et al. (2020)



Takeaways

- ▶ BART and T5 are useful for all sorts of seq2seq tasks involving language — so if you were going to use a seq2seq model, use one of these. (Caveat: need specialized models for things like language-to-code, but there's PLBART and CodeT5)
- ▶ UnifiedQA suggests that big generative models are good at generalizing across tasks and even to new tasks (although QA results have a long way to go)
- ▶ If we have a strong enough pre-trained model and train on enough tasks, can we generalize to new tasks?
- ▶ How do we specify those new tasks if they're not close to tasks we've already run on?
- ▶ Answer: **prompting**. But to do that well, we'll need to scale up further

GPT



OpenAI GPT/GPT2

- ▶ Very large language models using the Transformer architecture
- ▶ Straightforward left-to-right language model, trained on raw text
- ▶ GPT2: trained on 40GB of text

| | Parameters | Layers | d_{model} |
|--------------------------|--------------|-----------|-------------|
| | 117M | 12 | 768 |
| approximate size of BERT | 345M | 24 | 1024 |
| | 762M | 36 | 1280 |
| GPT-2 | 1542M | 48 | 1600 |

- ▶ By far the largest of these models trained when it came out in March 2019
- ▶ Because it's a language model, we can **generate** from it

Radford et al. (2019)



OpenAI GPT2

SYSTEM PROMPT
(HUMAN-WRITTEN)

Miley Cyrus was caught shoplifting from Abercrombie and Fitch on Hollywood Boulevard today.

MODEL COMPLETION
(MACHINE-WRITTEN,
SECOND TRY)

The 19-year-old singer was caught on camera being escorted out of the store by security guards.

The singer was wearing a black hoodie with the label 'Blurred Lines' on the front and 'Fashion Police' on the back.

Scroll down for video

Shoplifting: Miley Cyrus was caught shoplifting from Abercrombie and Fitch on Hollywood Boulevard today (pictured)

slide credit: OpenAI



Pre-Training Cost (with Google/AWS)

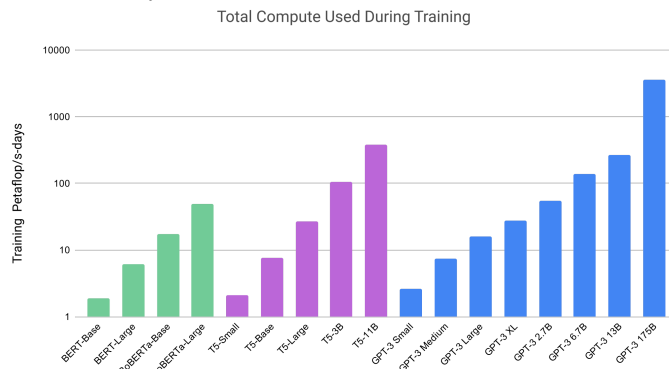
- ▶ BERT: Base \$500, Large \$7000
- ▶ GPT-2 (as reported in other work): \$25,000
- ▶ This is for a single pre-training run...developing new pre-training techniques may require many runs
- ▶ *Fine-tuning* these models can typically be done with a single GPU (but may take 1-3 days for medium-sized datasets)

<https://syncedreview.com/2019/06/27/the-staggering-cost-of-training-sota-ai-models/>



Pushing the Limits: GPT-3

- ▶ 175B parameter model: 96 layers, 96 heads, 12k-dim vectors
- ▶ Trained on Microsoft Azure, estimated to cost roughly \$10M



Brown et al. (2020)



Questions

- 1) How novel is the stuff being generated? (Is it just doing nearest neighbors on a large corpus?) How can we find out?
- 2) Can we use this model for things beyond story generation?
- 3) What harms might come from this model? (OpenAI pursued a "staged release" strategy and didn't release biggest model)



GPT-3

Story completion demo



Pre-GPT-3: Fine-tuning

- ▶ Fine-tuning: this is the “normal way” of doing learning in models like GPT-2
- ▶ Requires computing the gradient and applying a parameter update on every example
- ▶ **This is super expensive with 175B parameters**



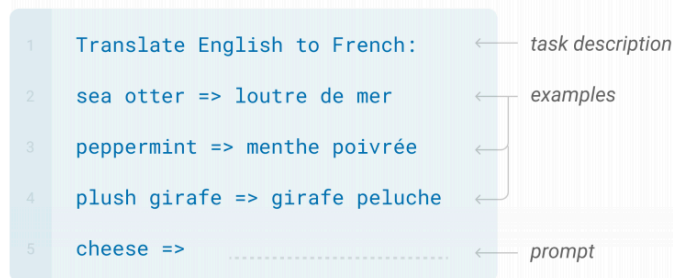
Brown et al. (2020)



GPT-3: Few-shot Learning

- ▶ GPT-3 proposes an alternative: **in-context learning**. Just uses the off-the-shelf model, no gradient updates

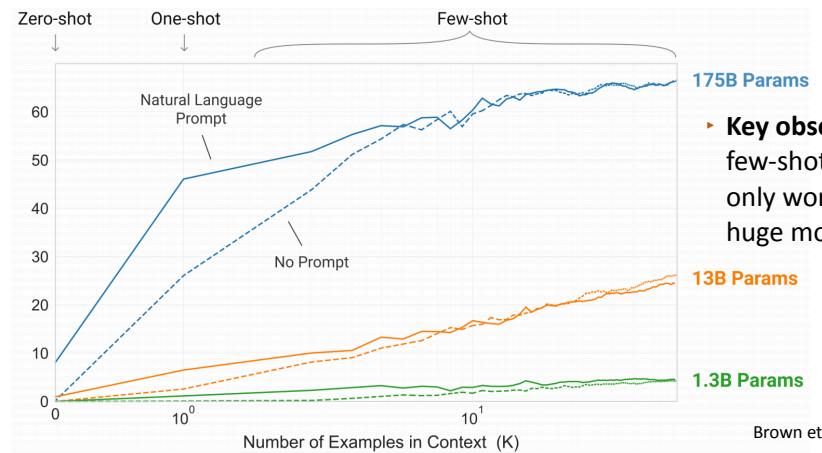
- ▶ This procedure depends heavily on the examples you pick as well as the prompt (“Translate English to French”)



Brown et al. (2020)



GPT-3



- ▶ **Key observation:** few-shot learning only works with huge models!

Brown et al. (2020)



GPT-3

| | 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 |

- ▶ Sometimes very impressive, (MultiRC, ReCoRD), sometimes very bad
- ▶ Results on other datasets are equally mixed — but still strong for a few-shot model!

Brown et al. (2020)

Prompting



PaLM

- ▶ “Pathways Language Model” from Google — **540B parameters!**
- ▶ Much of the paper is about data curation and datacenter networking

| Model | Layers | # of Heads | d_{model} | # of Parameters (in billions) | Batch Size |
|-----------|--------|------------|--------------------|----------------------------------|-------------------|
| PaLM 8B | 32 | 16 | 4096 | 8.63 | 256 → 512 |
| PaLM 62B | 64 | 32 | 8192 | 62.50 | 512 → 1024 |
| PaLM 540B | 118 | 48 | 18432 | 540.35 | 512 → 1024 → 2048 |

- ▶ Another big jump over GPT-3, but other advancements meant that new systems were even better

| Model | Avg NLG | Avg NLU |
|--------------|---------|---------|
| GPT-3 175B | 52.9 | 65.4 |
| GLaM 64B/64E | 58.4 | 68.7 |
| PaLM 8B | 41.5 | 59.2 |
| PaLM 62B | 57.7 | 67.3 |
| PaLM 540B | 63.9 | 74.7 |

Chowdhery et al. (2022)



Prompts

- ▶ Prompts can help induce the model to engage in certain behavior
- ▶ In the GPT-2 paper, “tl;dr:” (too long; didn't read) is mentioned as a prompt that frequently shows up in the wild **indicating a summary**
- ▶ tl;dr is an indicator that the model should “switch into summary mode” now — and if there are enough clean instances of tl;dr in the wild, maybe the model has been trained on a ton of diverse data?
- ▶ Good prompt + a few training examples in-context = strong task performance?

Brown et al. (2020)



Prompting

- ▶ Current training: GPT-3/PaLM trained on the web
- ▶ Current testing: feed in a very specific prompt and/or a set of in-context examples
- ▶ Two goals:
 1. Unify pre-training and testing phases
 2. Exploit data for downstream tasks — why are we trying to do question answering while ignoring all of the existing QA datasets?
- ▶ **Instruction tuning: fine-tune on supervised tasks after pre-training**
- ▶ **Let's see how an instruction-tuned GPT-3 works**



Prompts

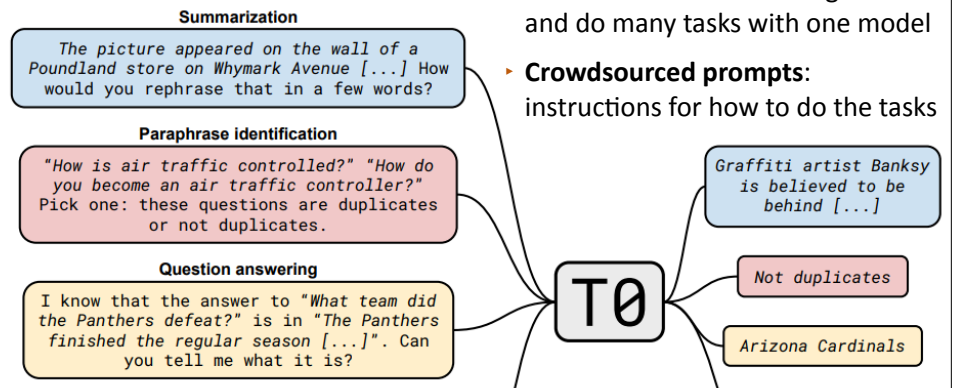
Prompting demo

Brown et al. (2020)

Instruction Tuning



Task Generalization: T0

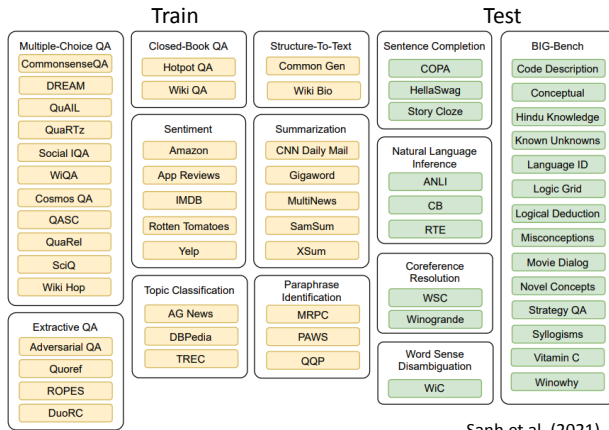


Sanh et al. (2021)



Task Generalization

- ▶ Train: a collection of tasks with prompts. **This uses existing labeled training data**
- ▶ Test: a new task specified only by a new prompt. **No training data in this task**

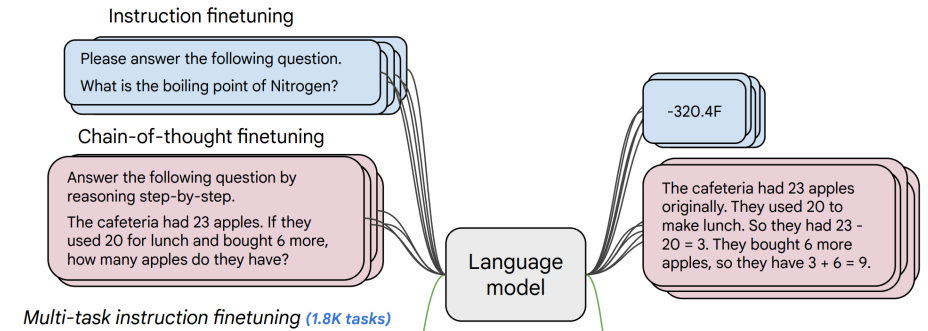


Sanh et al. (2021)



Frontiers

- ▶ FLAN-PaLM (October 20, 2022): 1800 tasks, 540B parameter model fine-tuned on many tasks after pre-training



Chung et al. (2022)



Frontiers

- ▶ FLAN-PaLM (October 20, 2022): 1800 tasks, 540B parameter model
- ▶ MMLU task (Hendrycks et al., 2020): 57 high school/college/professional exams:

| | | |
|----------------------------|---|---|
| Conceptual Physics | When you drop a ball from rest it accelerates downward at 9.8 m/s^2 . If you instead throw it downward assuming no air resistance its acceleration immediately after leaving your hand is | |
| | (A) 9.8 m/s^2 | ✓ |
| | (B) more than 9.8 m/s^2 | ✗ |
| | (C) less than 9.8 m/s^2 | ✗ |
| College Mathematics | In the complex z -plane, the set of points satisfying the equation $z^2 = z ^2$ is a | |
| | (A) pair of points | ✗ |
| | (B) circle | ✗ |
| | (C) half-line | ✗ |
| | (D) line | ✓ |

Figure 4: Examples from the Conceptual Physics and College Mathematics STEM tasks.

Chung et al. (2022)



Frontiers

- ▶ FLAN-PaLM (October 20, 2022): 1800 tasks, 540B parameter model
- ▶ MMLU task (Hendrycks et al., 2020): 57 high school/college/professional exams:

| | | |
|-----------|-----------------------------------|-------------|
| - | Random | 25.0 |
| - | Average human rater | 34.5 |
| May 2020 | GPT-3 5-shot | 43.9 |
| Mar. 2022 | Chinchilla 5-shot | 67.6 |
| Apr. 2022 | PaLM 5-shot | 69.3 |
| Oct. 2022 | Flan-PaLM 5-shot | 72.2 |
| | Flan-PaLM 5-shot: CoT + SC | 75.2 |
| - | Average human expert | 89.8 |

Chung et al. (2022)



Takeaways

- Pre-trained seq2seq models and generative language models can do well at lots of generation tasks
- Prompting is a way to harness their power and learn to do many tasks with a single model. Can be done without fine-tuning
- Instruction-tuned models are *by far* the best models we have for most generation and very complex language understanding tasks today. However, pre-trained models like BERT can still do well for classification
- Biggest best models (text-davinci-002, FLAN-PaLM) are closed-source