

GPT3 & instruct GPT

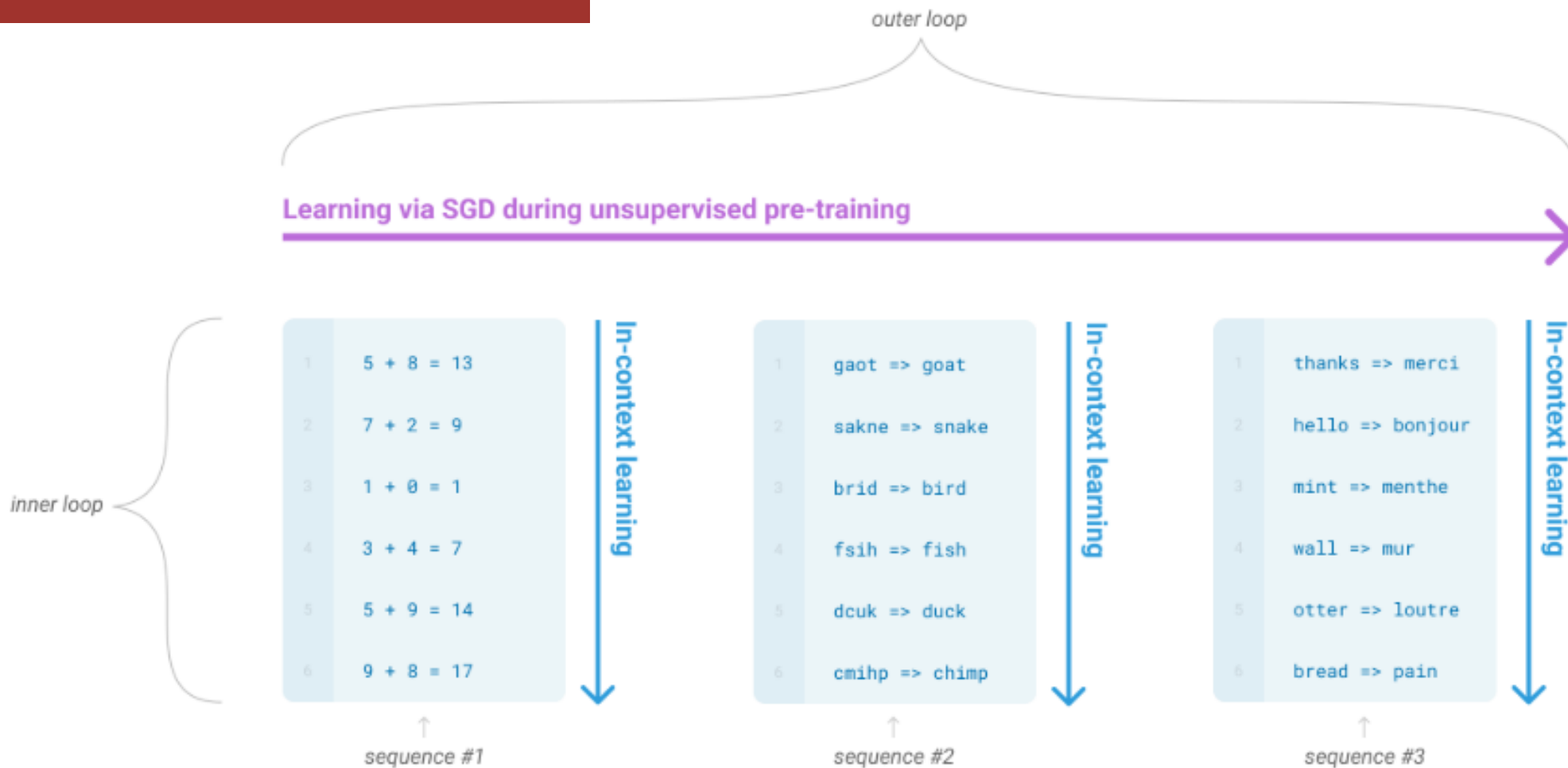
Kang Byeon Jin

Study Group NLP



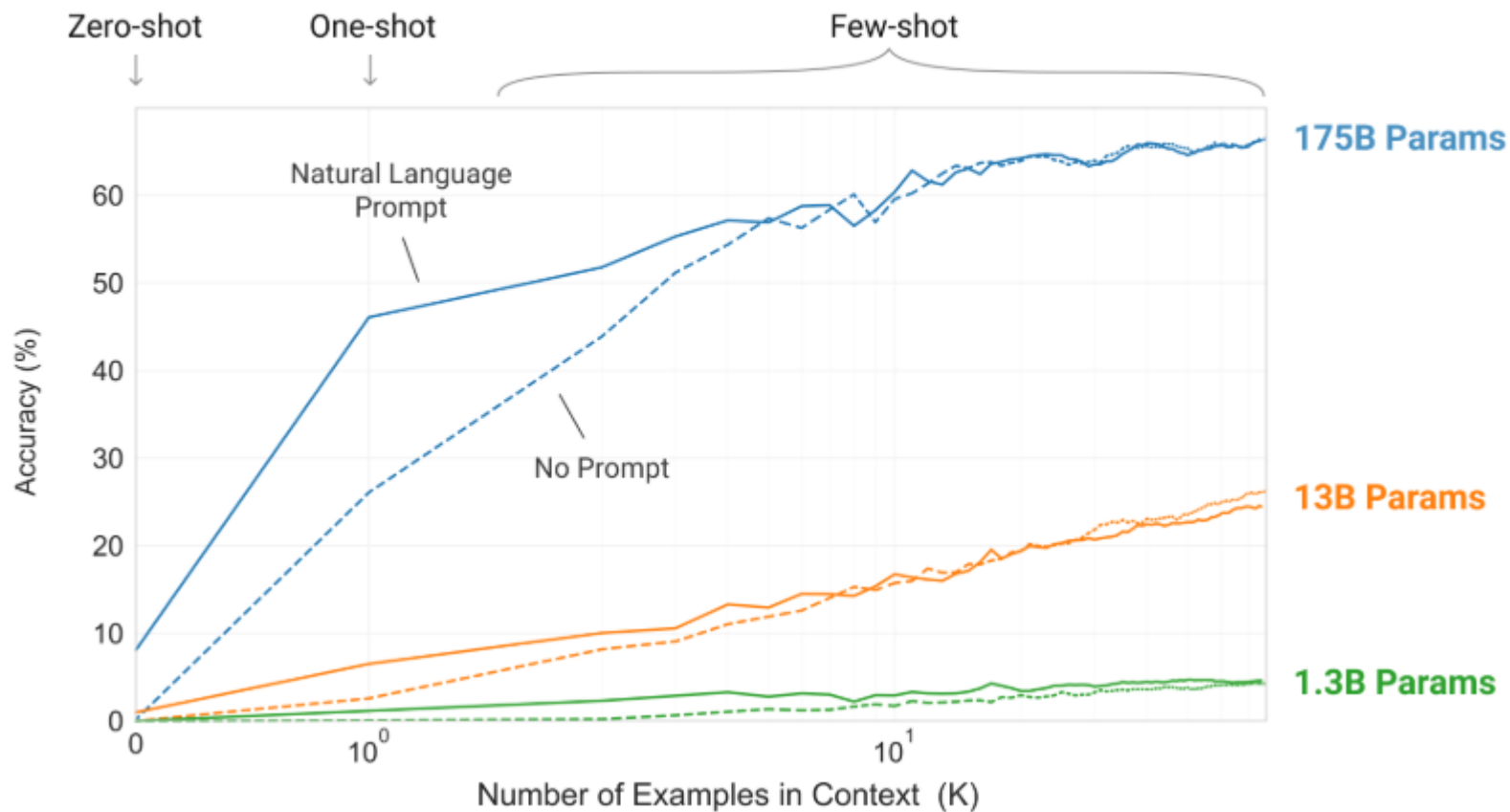
Language Models are Few-Shot Learners

Language model meta-learning



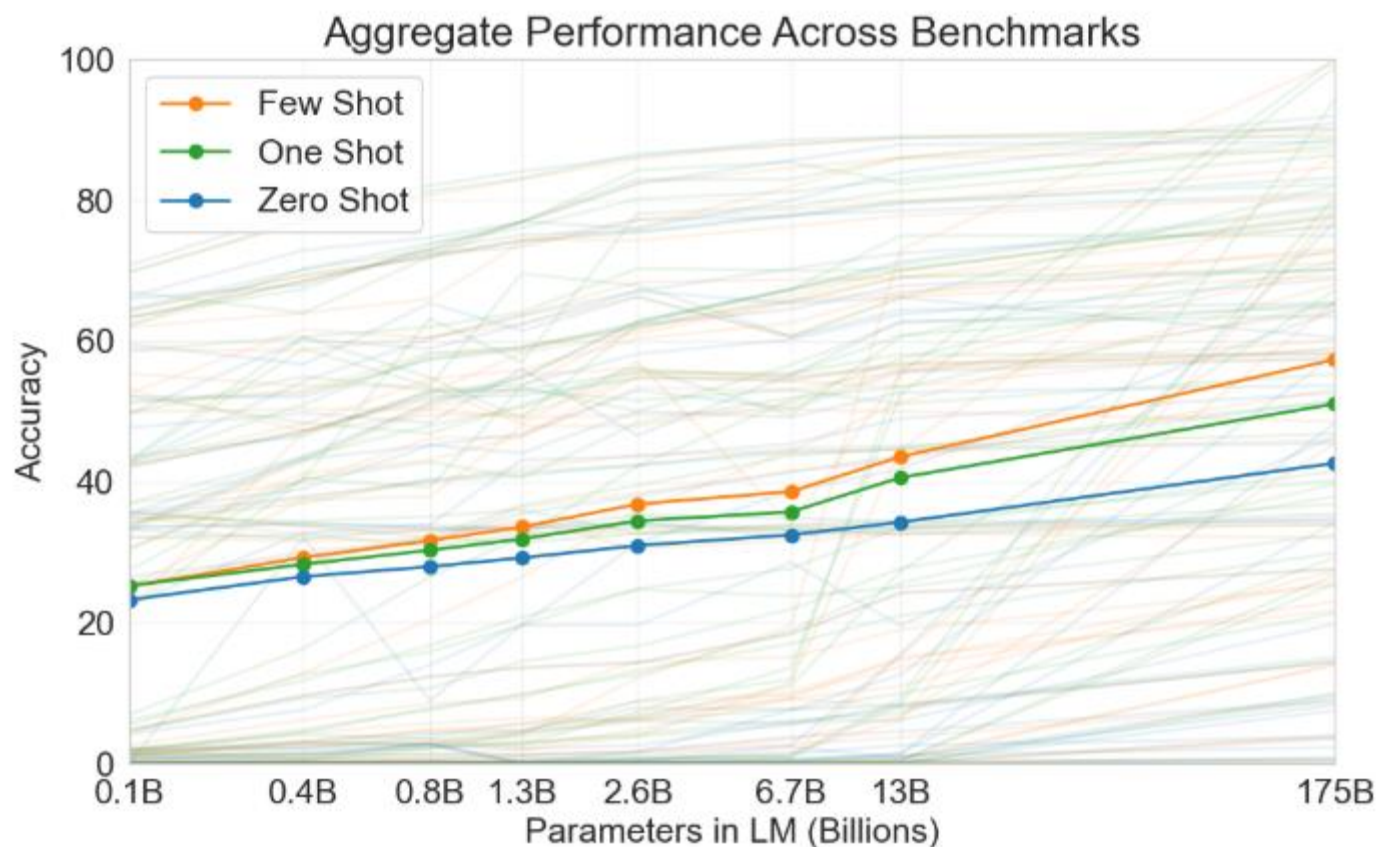
Language Models are Few-Shot Learners

Larger models make increasingly efficient use of in-context information



Language Models are Few-Shot Learners

Aggregate performance for all 42 accuracy-denominated benchmarks



Language Models are Few-Shot Learners

Approach

- Fine-Tuning (FT)

Pros: strong performance on many benchmarks

Cons: the need for a new large dataset for every task
the potential for poor generalization out-of-distribution
the potential to exploit spurious features of the training data
potentially resulting in an unfair comparison with human-performance

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.



Language Models are Few-Shot Learners

Approach

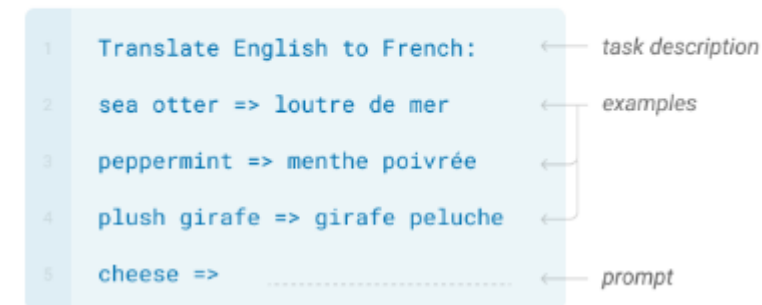
- Few-Shot (FS)

Pros: major reduction in the need for task-specific data and reduced potential to learn an overly narrow distribution from a large but narrow fine-tuning dataset

Cons: Results from this method have so far been much worse than SOTA fine-tuned models, a small amount of task specific data is still required

Few-shot

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



Language Models are Few-Shot Learners

Approach

- One-Shot (1S) Only one demonstration is allowed in addition to a natural language description of the task
- Zero-Shot (0S) The model is only given a natural language instruction describing the task

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

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

Language Models are Few-Shot Learners

Compare to traditional FT

“No gradient updates”

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.

```
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
```

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.



Language Models are Few-Shot Learners

Parameters of model

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	6.0×10^{-4}
GPT-3 Medium	350M	24	1024	16	64	0.5M	3.0×10^{-4}
GPT-3 Large	760M	24	1536	16	96	0.5M	2.5×10^{-4}
GPT-3 XL	1.3B	24	2048	24	128	1M	2.0×10^{-4}
GPT-3 2.7B	2.7B	32	2560	32	80	1M	1.6×10^{-4}
GPT-3 6.7B	6.7B	32	4096	32	128	2M	1.2×10^{-4}
GPT-3 13B	13.0B	40	5140	40	128	2M	1.0×10^{-4}
GPT-3 175B or “GPT-3”	175.0B	96	12288	96	128	3.2M	0.6×10^{-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.

All models use a context window of $n_{\text{ctx}} = 2048$ tokens.

✓ Make same input size to compare with other models

Language Models are Few-Shot Learners

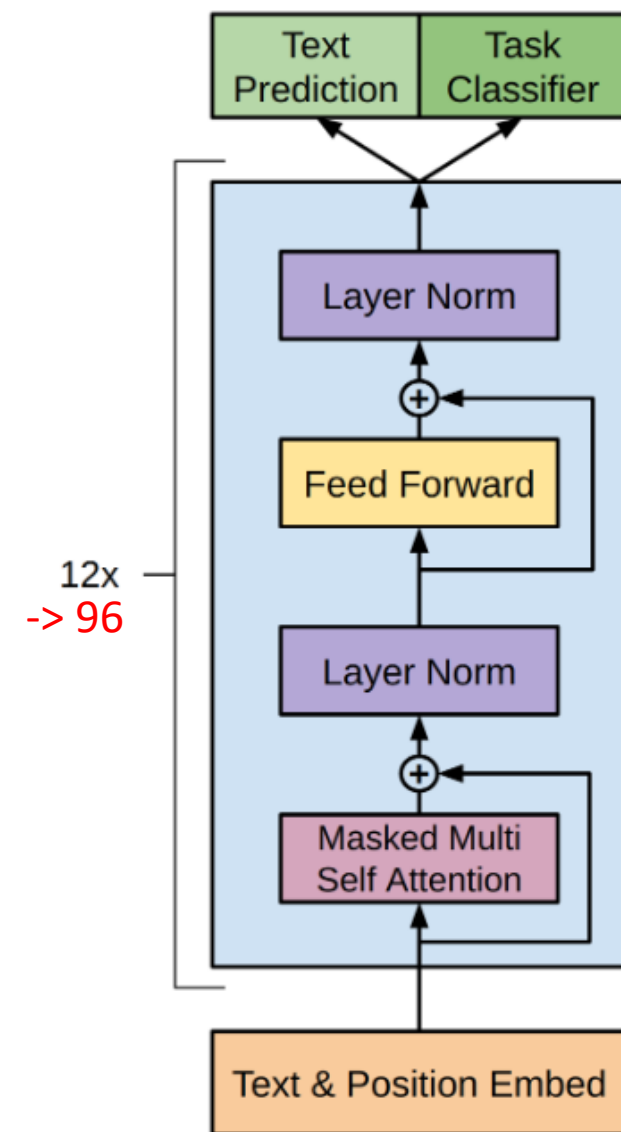
Parameters of model

Model Name	n_{params}	n_{layers}	d_{model}	n_{heads}	d_{head}	Batch Size
GPT-3 Small	125M	12	768	12	64	0.5M
GPT-3 Medium	350M	24	1024	16	64	0.5M
GPT-3 Large	760M	24	1536	16	96	0.5M
GPT-3 XL	1.3B	24	2048	24	128	1M
GPT-3 2.7B	2.7B	32	2560	32	80	1M
GPT-3 6.7B	6.7B	32	4096	32	128	2M
GPT-3 13B	13.0B	40	5140	40	128	2M
GPT-3 175B or “GPT-3”	175.0B	96	12288	96	128	3.2M

Table 2.1: Sizes, architectures, and learning hyper-parameters (batch size in tokens and learning rate) for the GPT-3 models which we trained. All models were trained for a total of 300 billion tokens.

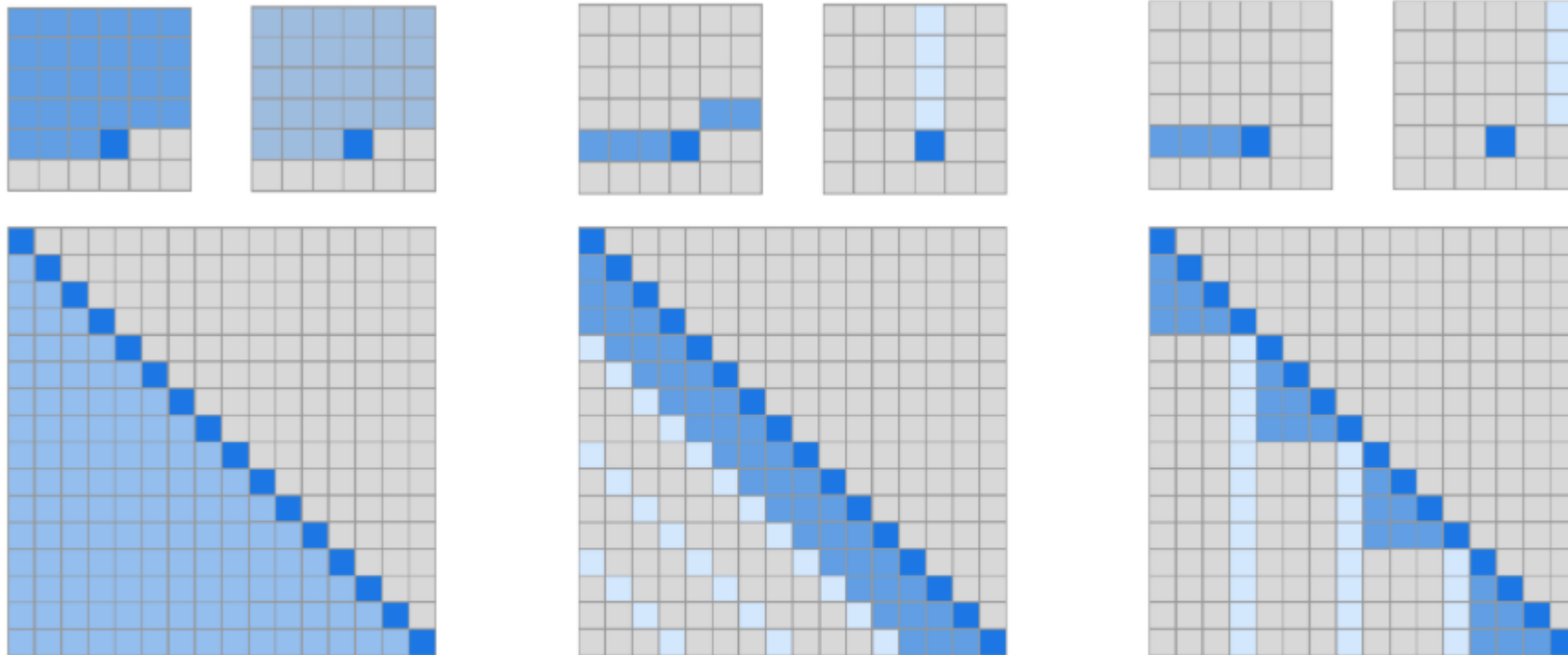
All models use a context window of $n_{\text{ctx}} = 2048$ tokens.

✓ Make same input size to compare with other models



Language Models are Few-Shot Learners

Architecture



(a) Transformer

(b) Sparse Transformer (strided)

(c) Sparse Transformer (fixed)

Language Models are Few-Shot Learners

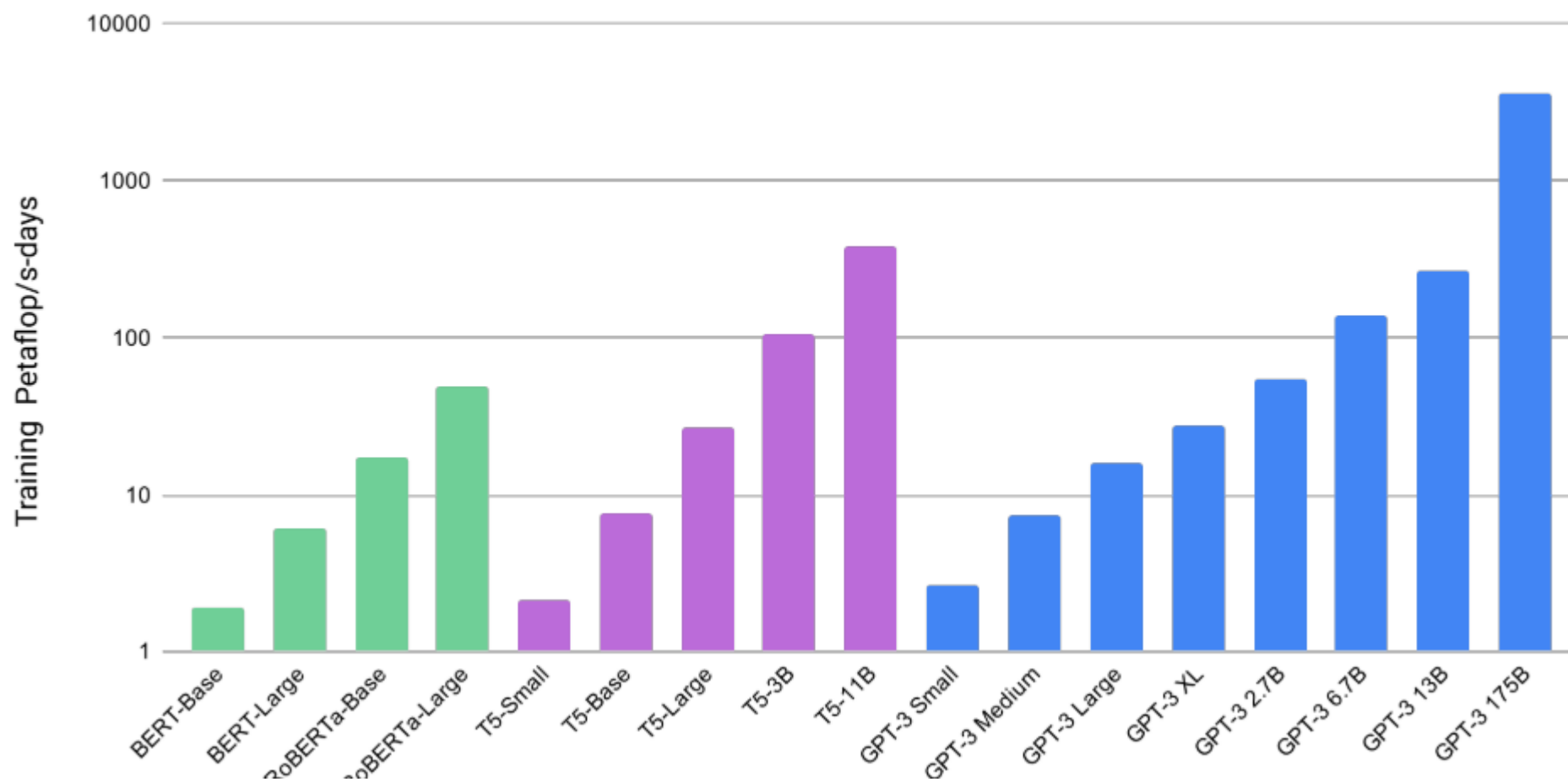
Training Dataset

- Common Crawl dataset (constituting nearly a trillion words)
- 3 steps to improve the average quality of the dataset
 1. downloaded and filtered a version of Common Crawl based on similarity to a range of high-quality reference corpora
 2. performed fuzzy deduplication at the document level, within and across datasets, to prevent redundancy and preserve the integrity of our held-out validation set as an accurate measure of overfitting
 3. Added known high-quality reference corpora to the training mix to augment Common Crawl and increase its diversity

Language Models are Few-Shot Learners

Total compute used during training

Total Compute Used During Training



Language Models are Few-Shot Learners

Datasets used to train GPT-3

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

Q. GPT Cost?

Q. Time Cost?

Language Models are Few-Shot Learners

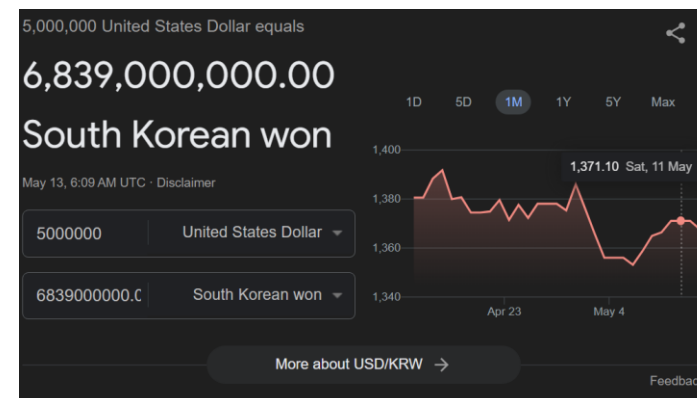
Datasets used to train GPT-3

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

Q. GPT3 Cost? ~ 5 million

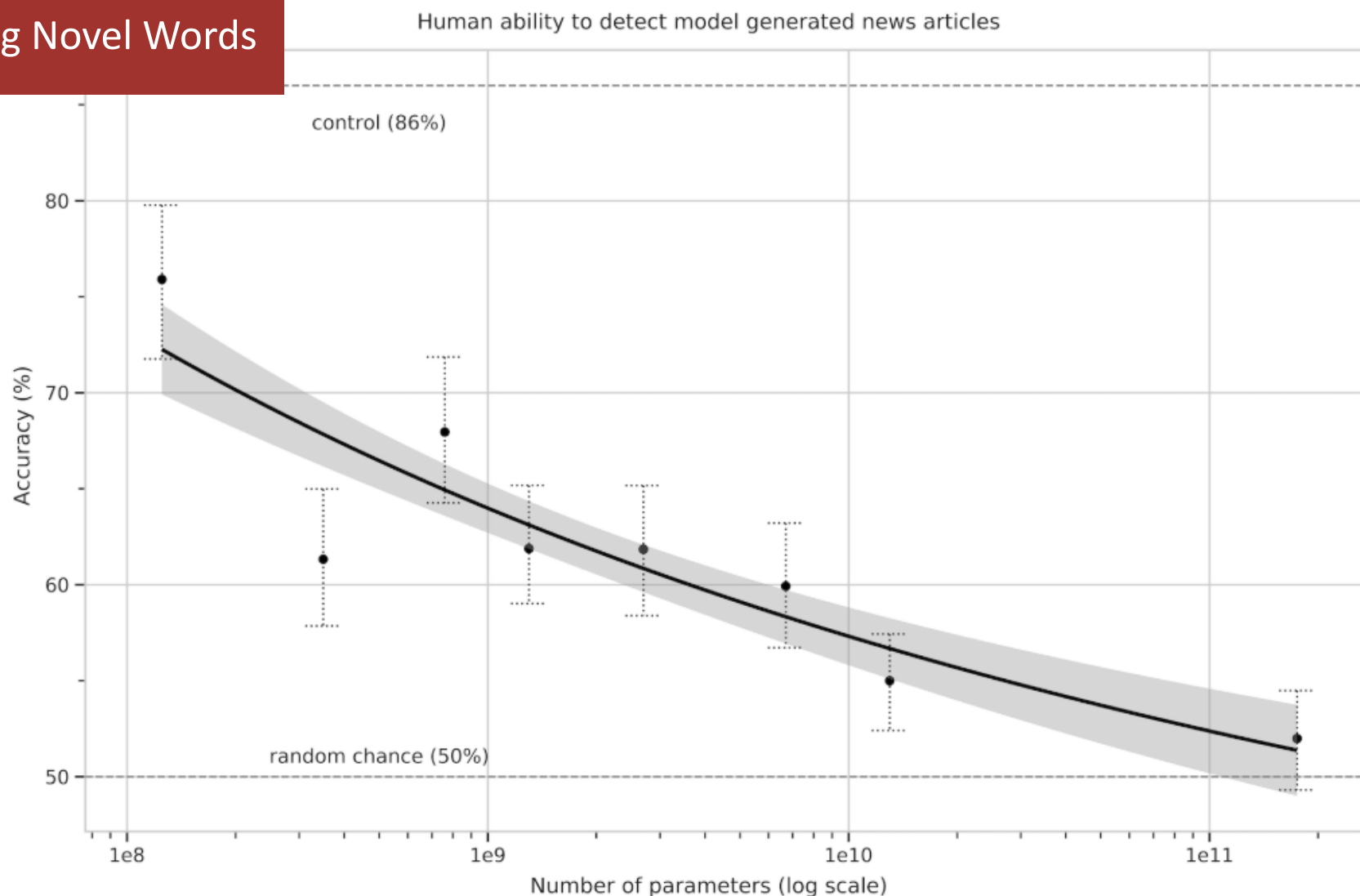
Q. time? 3.14×10^{23} flops computes

\approx v100 GPU(28 Tflops) 355years



Language Models are Few-Shot Learners

Learning and Using Novel Words



Language Models are Few-Shot Learners

Grammar

Poor English input: I'd be more than happy to work with you in another project.

Good English output: I'd be more than happy to work with you on another project.

Poor English input: Please provide me with a short brief of the design you're looking for and that'd be nice if you could share some examples or project you did before.

Good English output: Please provide me with a brief description of the design you're looking for and that would be nice if you could share some examples or projects you have done before.

Poor English input: The patient was died.

Good English output: The patient died.

Poor English input: We think that Leslie likes ourselves.

Good English output: We think that Leslie likes us.

Poor English input: Janet broke Bill on the finger.

Good English output: Janet broke Bill's finger.

Poor English input: Mary arranged for, in St. Louis, John to rent a house cheap.

Good English output: Mary arranged for John to rent a house in St. Louis.

Poor English input: Today I have went to the store to to buys some many bottle of water.

Good English output: Today I went to the store to buy some bottles of water.

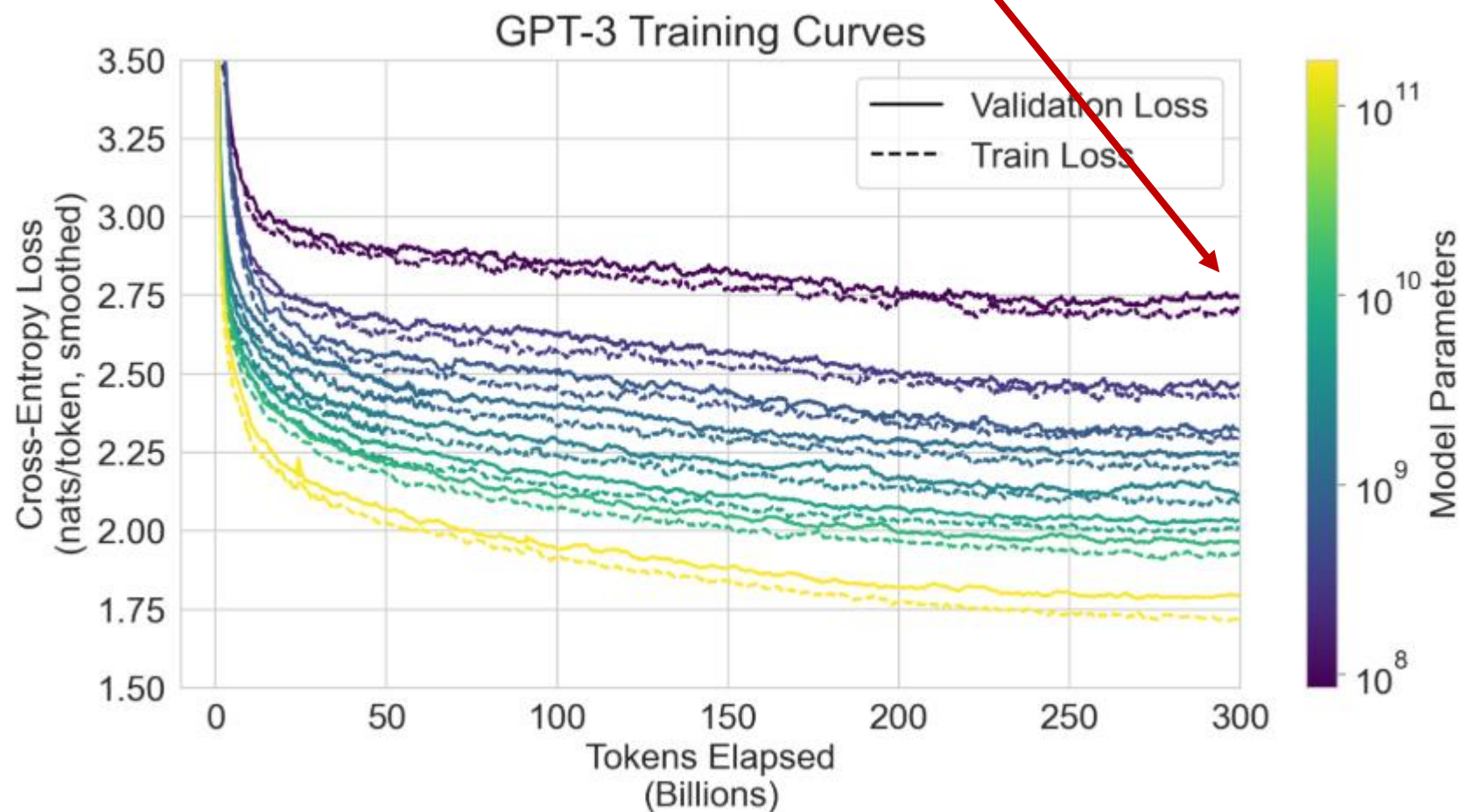
Poor English input: I have tried to hit ball with bat, but my swing is has miss.

Good English output: I tried to hit the ball with the bat, but my swing missed.

Language Models are Few-Shot Learners

Overfitting?

most of the gap comes from a difference in difficulty rather than overfitting



Q. Gap meaning?

Q. Difficulty?

Q. Validation loss
vs
Training loss

Language Models are Few-Shot Learners

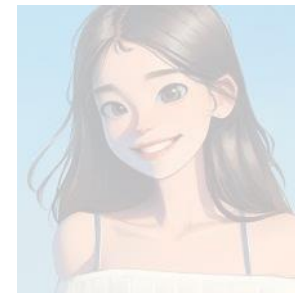
Limitation

- Weaknesses in text synthesis and several NLP tasks
 - *samples still sometimes repeat themselves semantically at the document level, start to lose coherence over sufficiently long passages*
 - *occasionally contain non-sequitur sentences or paragraphs*
- Structural and algorithmic limitation
 - *Auto-regressive, not bidirectional*
 - *Pretraining objective weights every token equally*

Language Models are Few-Shot Learners

Limitation

- poor sample efficiency during pre-training
- ambiguity about whether few-shot learning actually learns new tasks “from scratch” at inference time, or if it simply recognizes and identifies tasks that it has learned during training
- its predictions on novel inputs as observed by the much higher **variance** in performance than humans on standard benchmarks, and it retains the **biases** of the data it has been trained on



instruct GPT

“Training language models to follow instructions with human feedback”

Instruct GPT

Difference

Can “directly instruct” to GPT

Write a short story about the moon and the stars: Does the moon like the stars? Does the moon like the stars? Does the moon like the stars? Does the moon like the stars? Does the moon like the stars? Does the moon like the stars? Does the moon like the stars? Does the moon like the stars? Does the moon like the stars?

GPT3

Write a short story about the moon and the stars:

Once upon a time, the moon and the stars lived together in the sky. They were the best of friends, and they loved spending their time together. The moon would light up the sky at night, and the stars would twinkle and shine. They were always happy and they never fought.

But one day, the moon and the stars began to grow apart. The moon started to feel left out and the stars began to ignore her. The moon didn't understand what was happening, and she felt very lonely.

Instruct GPT



Instruct GPT?

Motivation

“Misaligned”

=> Untruthful, Toxic, Not helpful to user

Instruct GPT?

Motivation

“Misaligned”

=> Untruthful, Toxic, Not helpful to user

Solution



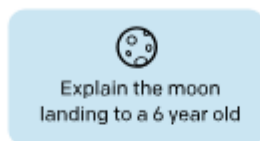
Reinforcement Learning from Human Feedback
(RLHF; Christiano et al., 2017; Stiennon et al., 2020)

Instruct GPT method

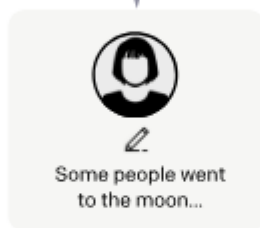
Step 1

Collect demonstration data, and train a supervised policy.

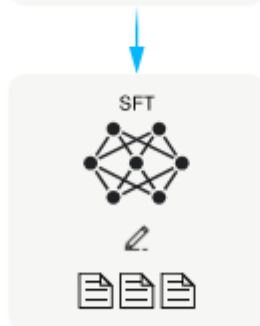
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



This data is used to fine-tune GPT-3 with supervised learning.



Step 2

Collect comparison data, and train a reward model.

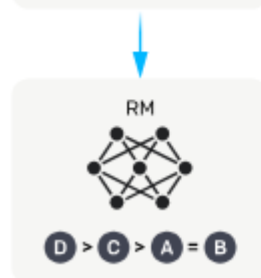
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



Step 3

Optimize a policy against the reward model using reinforcement learning.

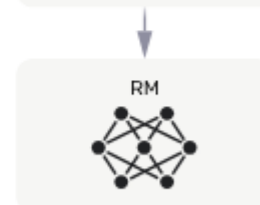
A new prompt is sampled from the dataset.



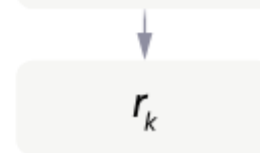
The policy generates an output.



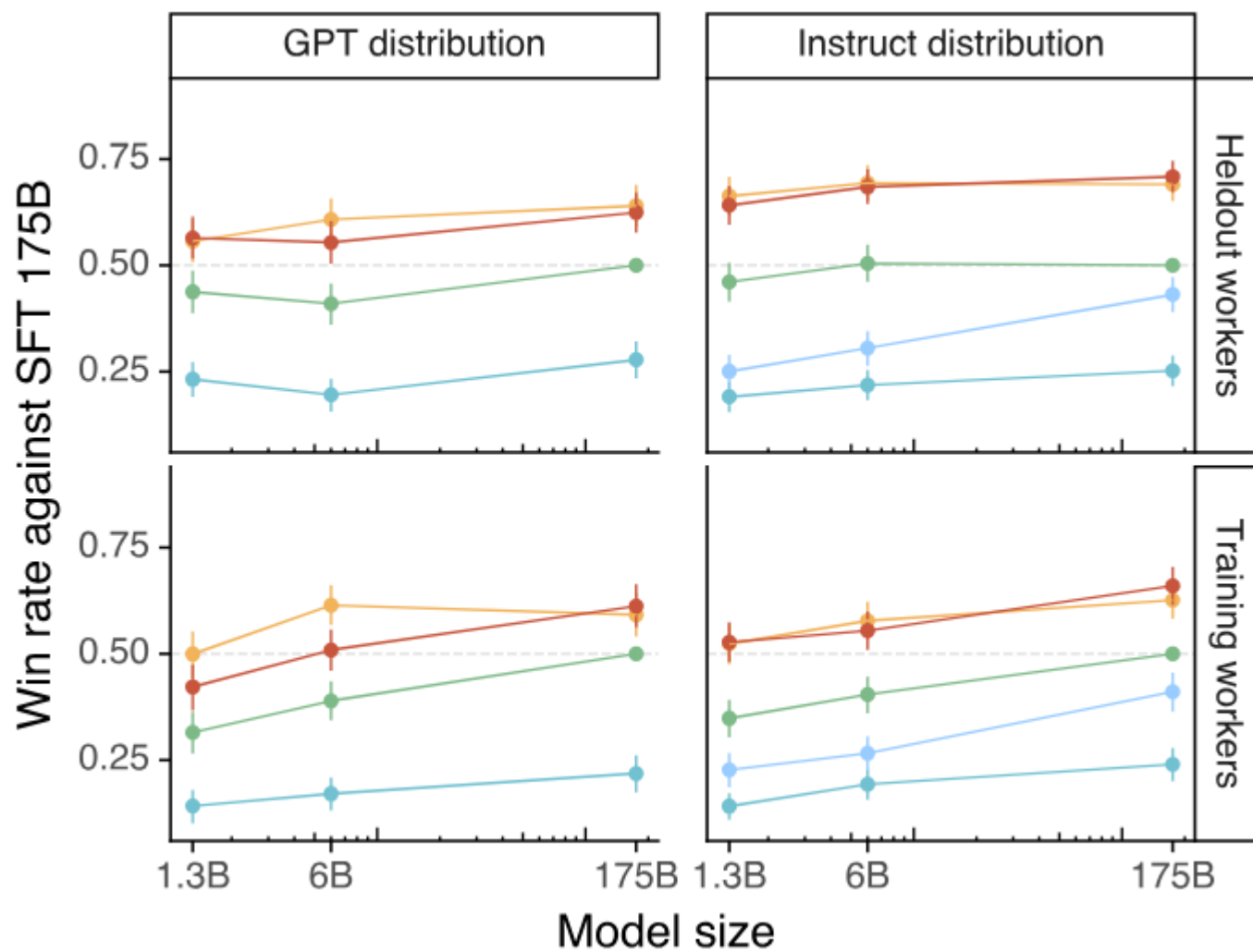
The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



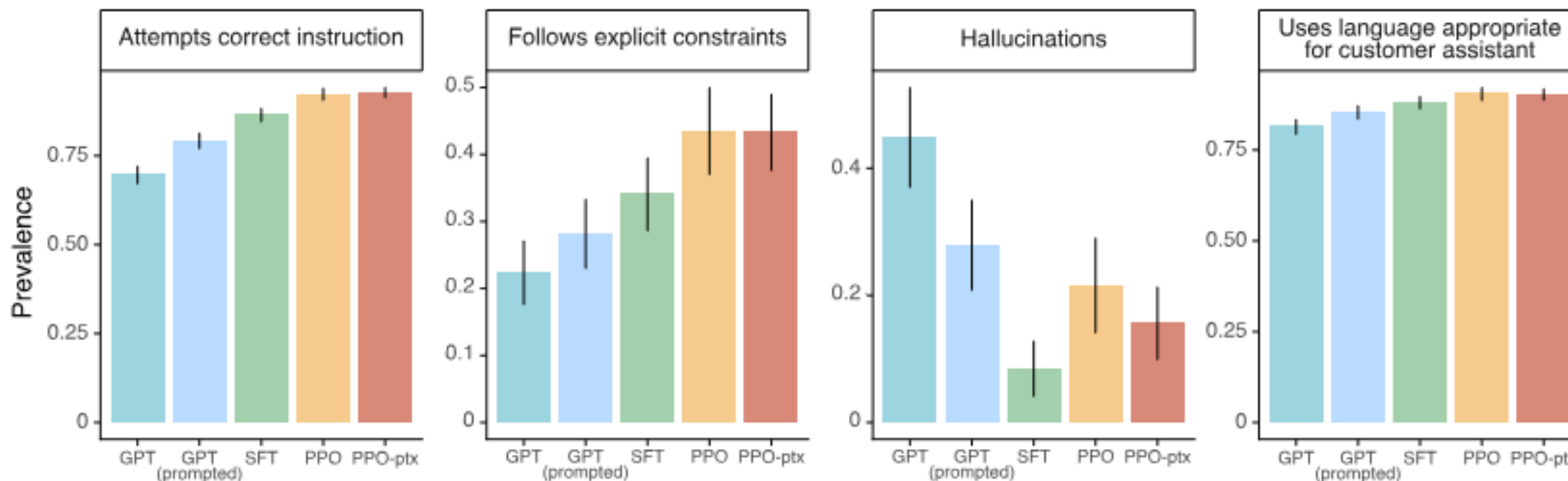
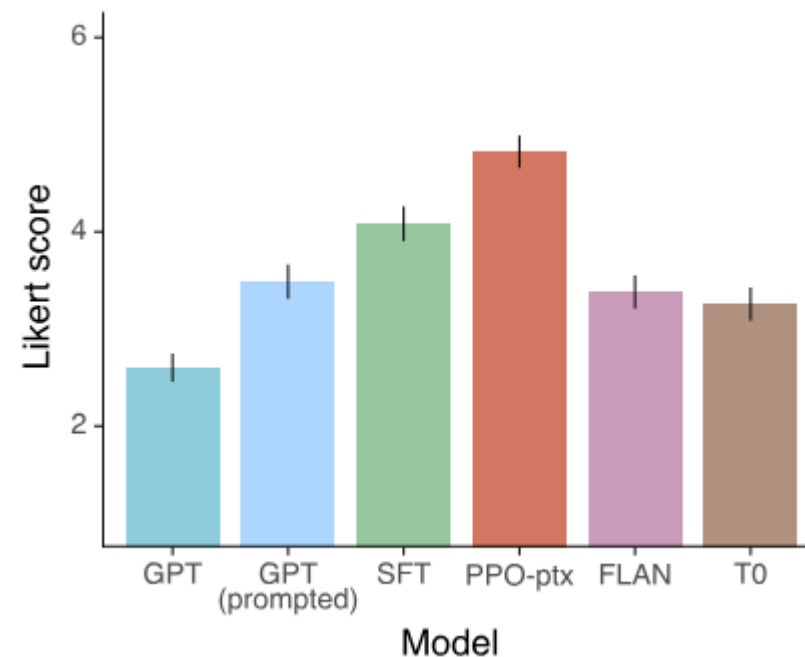
Instruct GPT Results



Instruct GPT Results

It is unclear how to measure honesty in purely generative models

Similarly to honesty, measuring the harms of language models also poses many challenges



prerequisite

Reinforcement Learning

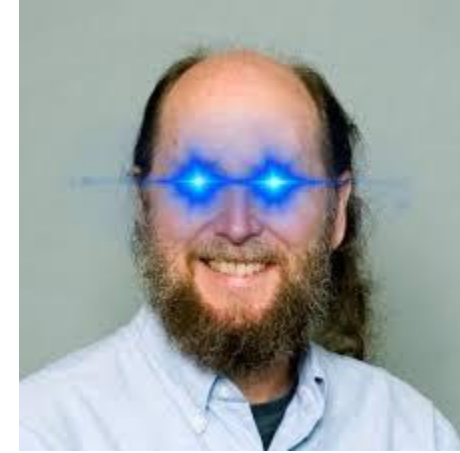
What is Reinforcement Learning?

prerequisite

Reinforcement Learning

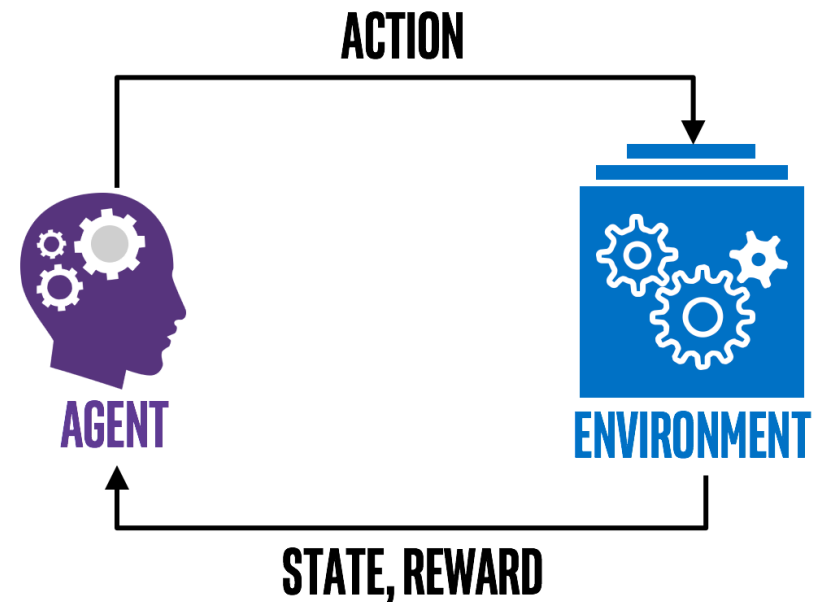
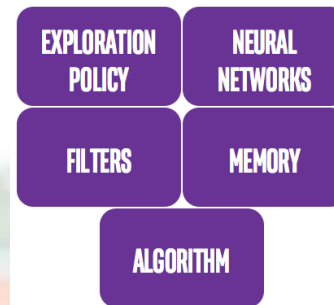
What is Reinforcement Learning?

MDP, PPO, Q-learning, DQN, actor-critic...



Richard S. Sutton

David Silver

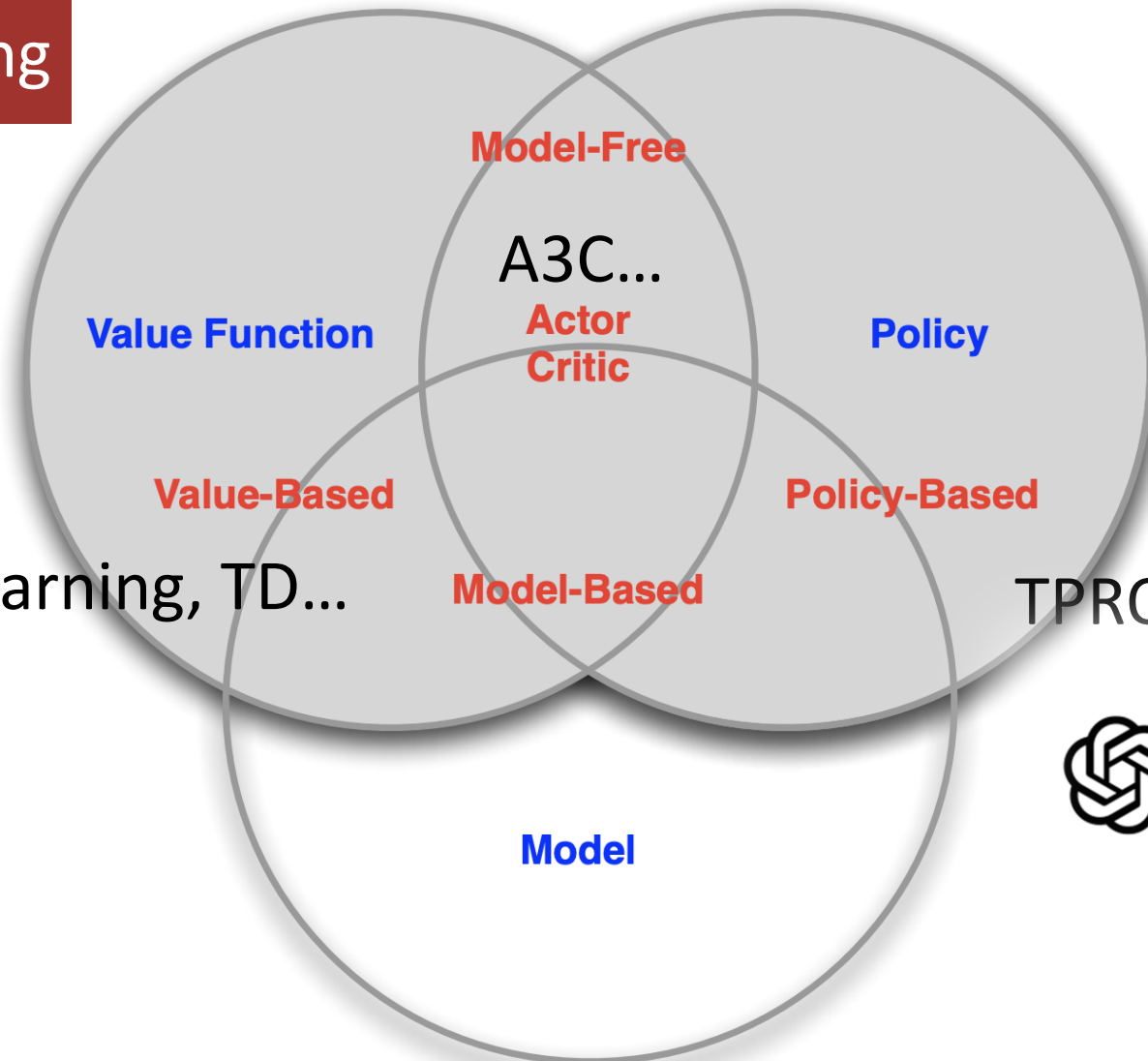


RL Overview

Reinforcement Learning

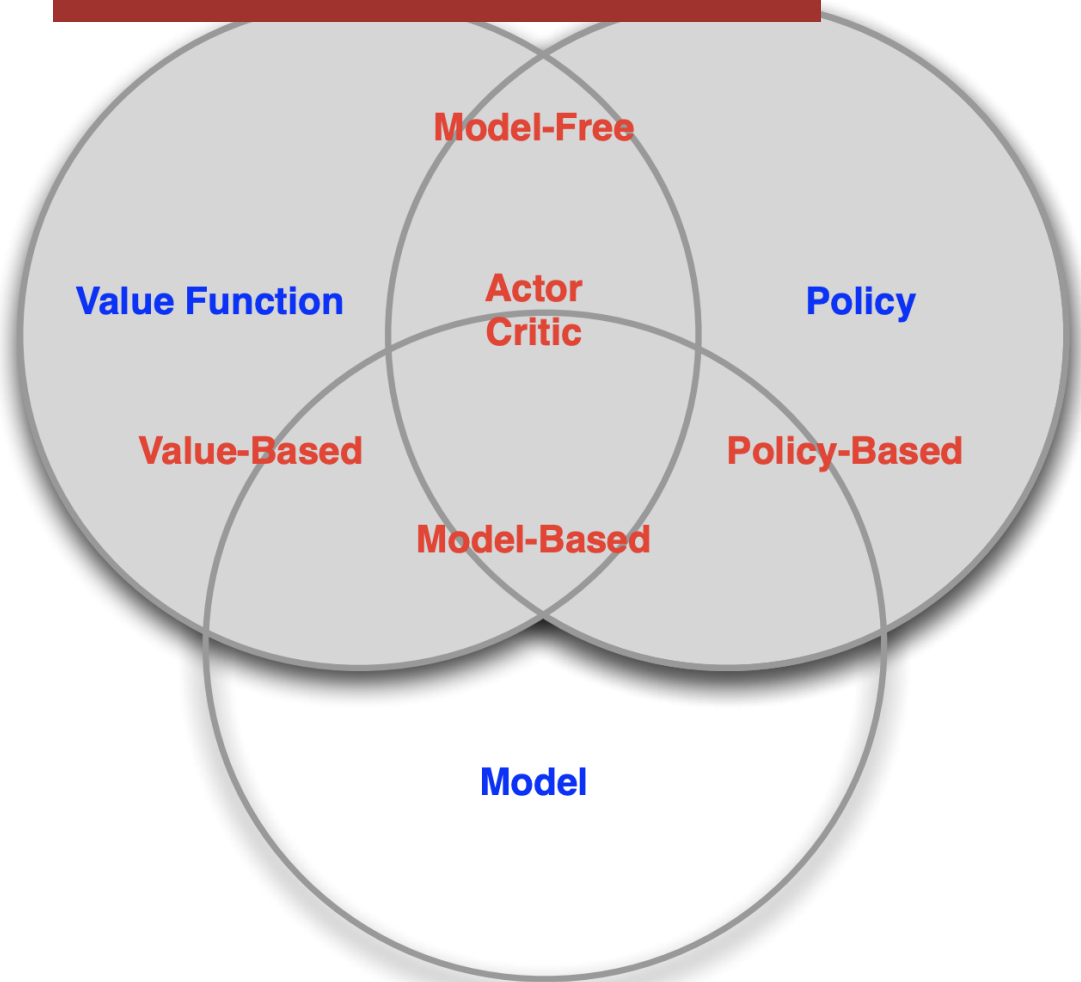


DQN, MC Learning, TD...



Contents

Reinforcement Learning



Markov Decision Process (MDP)

= Decision Process + Markov Property

$$P[S_{t+1} | S_t] = P[S_{t+1} | S_1, \dots, S_t]$$

A brown arrow points from the text 'Markov Property' to the state S_t in the equation above.

MDP = $\langle S, A, P, R, \gamma \rangle$ =

State, action, transition probability,
reward, discount factor

prerequisite

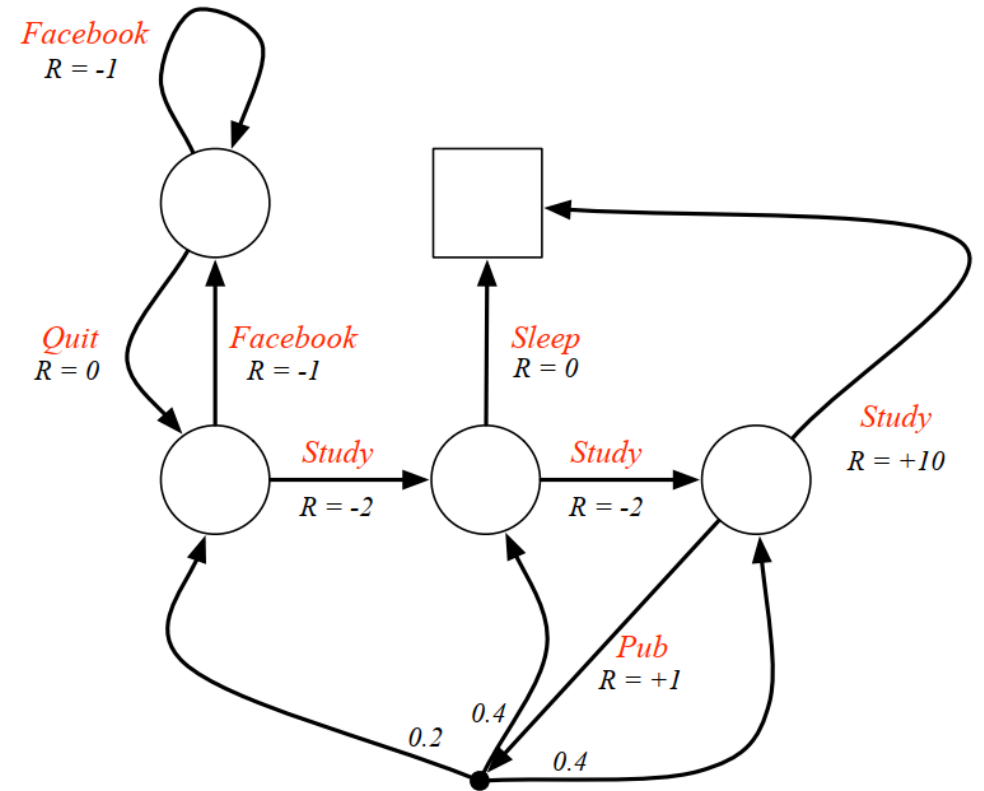
Bellman equation

Bellman equation, optimality...

$$v_*(s) = \max_{\pi} v_{\pi}(s)$$

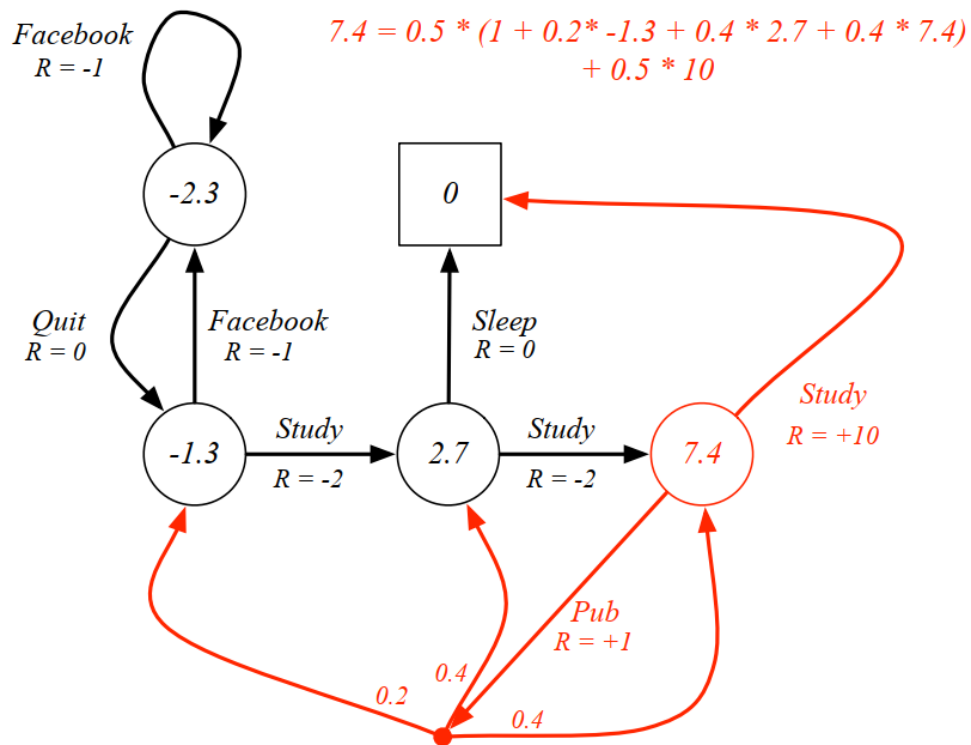
$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a)$$

$$\begin{aligned} v_{\pi}(s_t) &= \mathbb{E}_{\pi}[G_t] \\ &= \mathbb{E}_{\pi}[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots] \\ &= \mathbb{E}_{\pi}[r_{t+1} + \gamma(r_{t+2} + \gamma r_{t+3} + \dots)] \\ &= \mathbb{E}_{\pi}[r_{t+1} + \gamma G_{t+1}] \\ &= \mathbb{E}_{\pi}[r_{t+1} + \gamma v_{\pi}(s_{t+1})] \end{aligned}$$

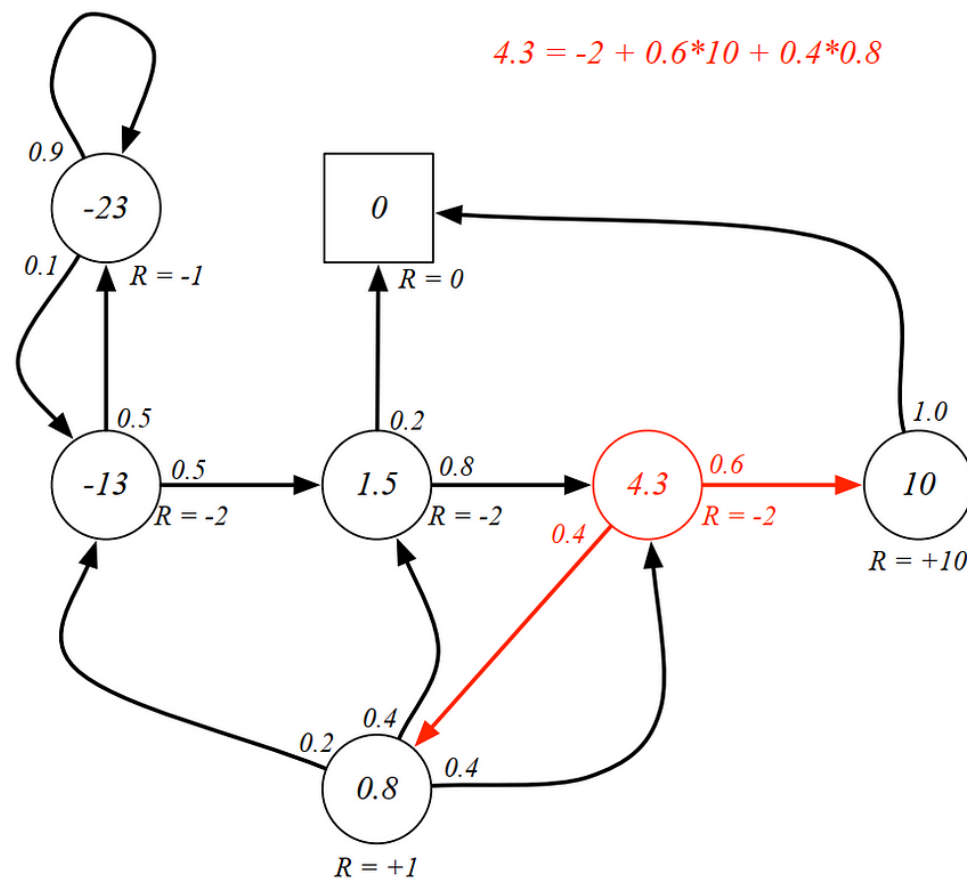


prerequisite

MDP (decision)



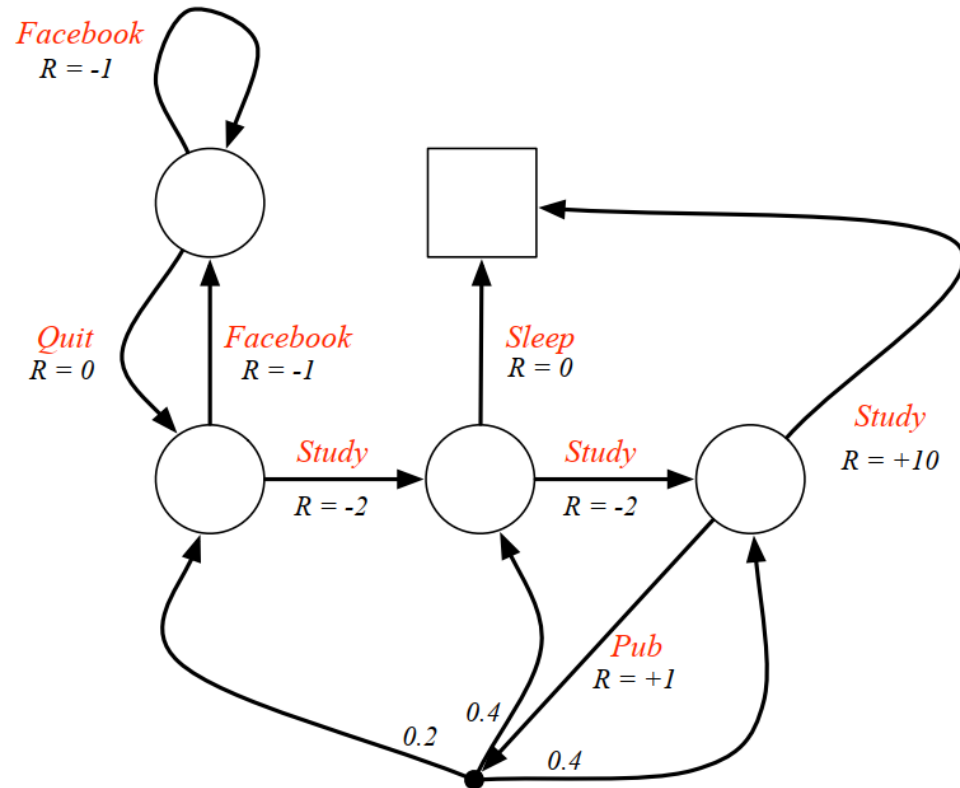
MRP (reward)



prerequisite

Reinforcement Learning

Value based? Policy based?
=> Different Object function



prerequisite

TRPO, PPO

Motivation: want to update in trust bound => How?



PPO (Proximal policy optimization)

Trust Region Policy Optimization (TRPO)

$$\begin{aligned} & \underset{\theta}{\text{maximize}} \mathbb{E}_{s \sim \rho_{\theta_{\text{old}}}, a \sim q} \left[\frac{\pi_{\theta}(a|s)}{q(a|s)} Q_{\theta_{\text{old}}}(s, a) \right] \quad (14) \\ & \text{subject to } \mathbb{E}_{s \sim \rho_{\theta_{\text{old}}}} [D_{\text{KL}}(\pi_{\theta_{\text{old}}}(\cdot|s) \parallel \pi_{\theta}(\cdot|s))] \leq \delta. \end{aligned}$$

Proximal policy optimization (PPO)

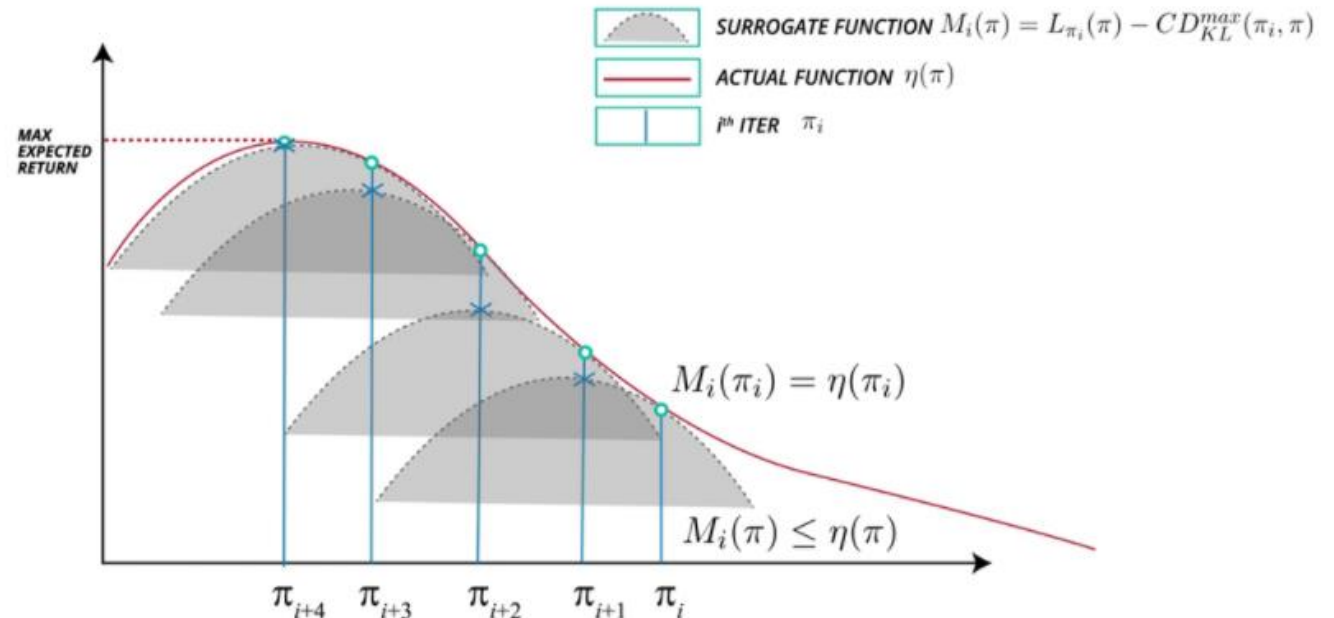
$$r_t(\theta) = \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)}$$

$$L^{CPI}(\theta) = \hat{\mathbb{E}}_t \left[\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t \right] = \hat{\mathbb{E}}_t [r_t(\theta) \hat{A}_t]$$

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t [\min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t)]$$

MM Algorithm

Minorization-Maximization

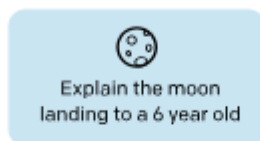


Instruct GPT method

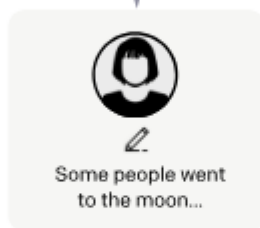
Step 1

Collect demonstration data, and train a supervised policy.

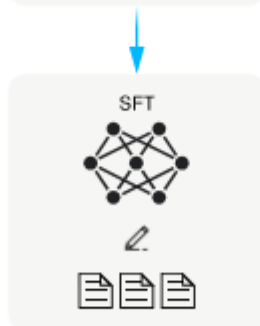
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



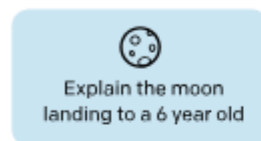
This data is used to fine-tune GPT-3 with supervised learning.



Step 2

Collect comparison data, and train a reward model.

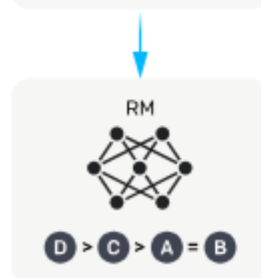
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



Step 3

Optimize a policy against the reward model using reinforcement learning.

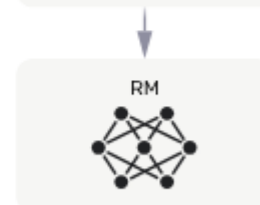
A new prompt is sampled from the dataset.



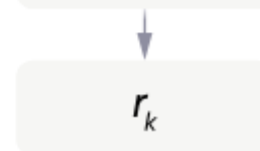
The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



Instruct GPT

Reward model

$$\text{loss}(\theta) = -\frac{1}{\binom{K}{2}} E_{(x, y_w, y_l) \sim D} [\log(\sigma(r_\theta(x, y_w) - r_\theta(x, y_l)))]$$

“We ran an experiment where we split our labelers into 5 groups, and train 5 RMs (with 3 different seeds)”

Step 2

Collect comparison data, and train a reward model.

A prompt and several model outputs are sampled.

🧠
Explain the moon landing to a 6 year old

A
Explain gravity...

B
Explain war...

C
Moon is natural satellite of...

D
People went to the moon...

A labeler ranks the outputs from best to worst.

👤
D > C > A = B

This data is used to train our reward model.

RM
D > C > A = B

Instruct GPT

PPO

$$\text{objective}(\phi) = E_{(x,y) \sim D_{\pi_{\phi}^{\text{RL}}}} \left[r_{\theta}(x, y) - \beta \log \left(\pi_{\phi}^{\text{RL}}(y | x) / \pi^{\text{SFT}}(y | x) \right) \right] + \gamma E_{x \sim D_{\text{pretrain}}} \left[\log(\pi_{\phi}^{\text{RL}}(x)) \right]$$

Step 3

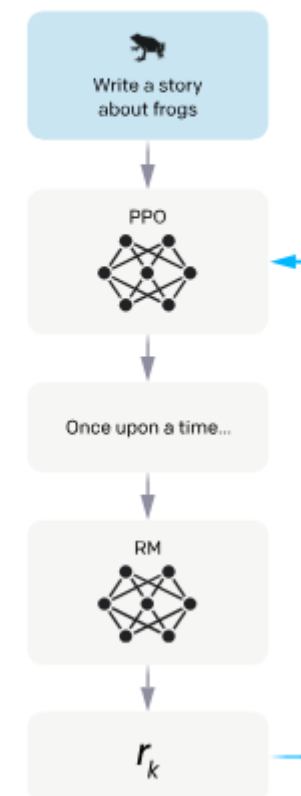
Optimize a policy against the reward model using reinforcement learning.

A new prompt is sampled from the dataset.

The policy generates an output.

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.



prerequisite

Reinforcement Learning

What is Reinforcement Learning?

We ran an experiment where we split our labelers into 5 groups,
and train 5 RMs (with 3 different seeds)

Limitation

Methodology

The behavior of our InstructGPT models is determined in part by the human feedback obtained from our contractors. ~~

However, this group is clearly not representative of the full spectrum of people who will use and be affected by our deployed models.

Models

Our models are neither fully aligned nor fully safe; they still generate toxic or biased outputs, make up facts, and generate sexual and violent content without explicit prompting. They can also fail to generate reasonable outputs on some inputs.

Open questions

Many methods could be tried to further decrease the models' propensity to generate toxic, biased, or otherwise harmful outputs

*While we mainly focus on RLHF, there are many other algorithms that could be used to train policies on our demonstration and comparison data to get even better results.
⇒ explore expert iteration, simpler behavior cloning methods...*

Comparisons are also not necessarily the most efficient way of providing an alignment signal.

...



TRAIN AND TEST