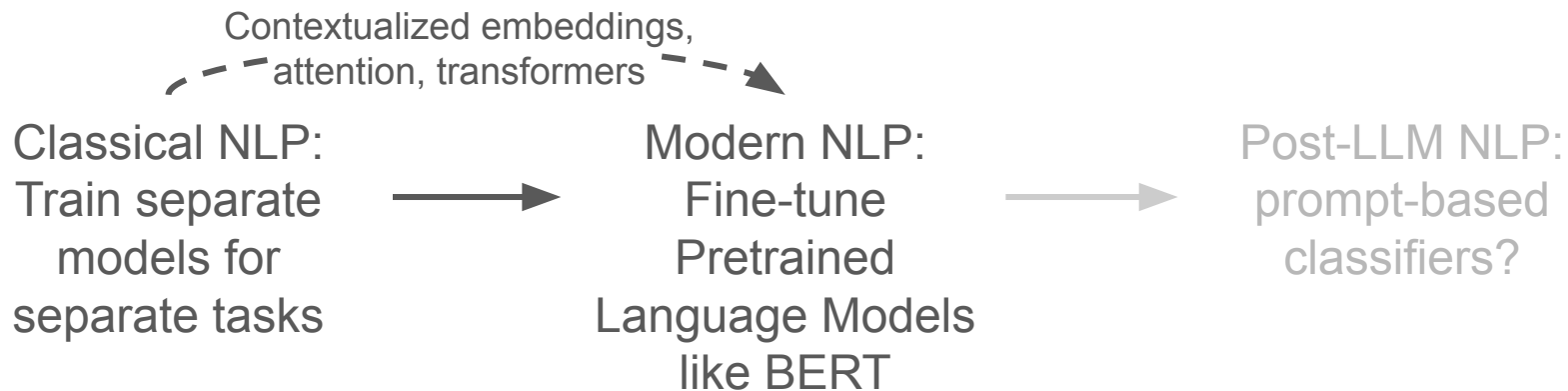


NLP Crash Course III

Indira Sen and David Garcia

The timeline so far: Supervised NLP



The timeline so far: Supervised NLP

Classical NLP:
Train separate
models for
separate tasks



Modern NLP:
Fine-tune
Pretrained
Language Models
like BERT



Post-LLM NLP:
prompt-based
classifiers?

Natural Language Generation (NLG)

NLG focuses on systems that produce fluent, coherent and useful language output for human consumption

Natural Language Generation (NLG)

NLG focuses on systems that produce fluent, coherent and useful language output for human consumption

Open-ended generation: the output distribution still has high freedom

Non-open-ended generation: the input mostly determines the output generation.

Natural Language Generation (NLG)

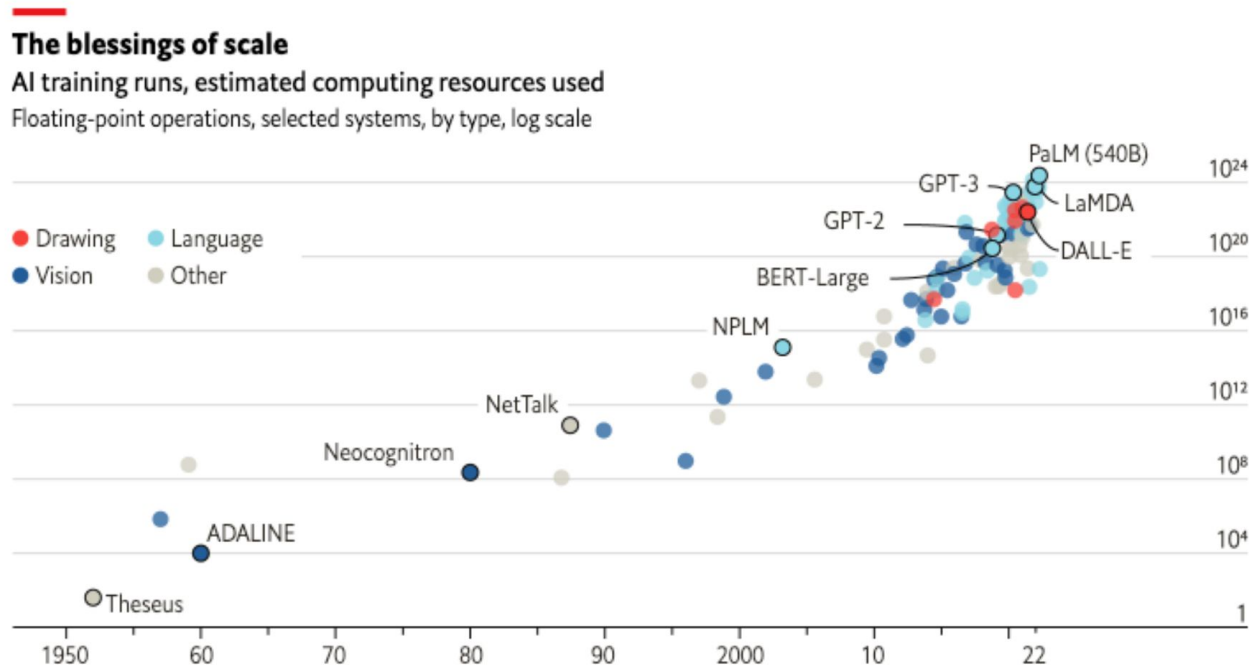
NLG focuses on systems that produce fluent, coherent and useful language output for human consumption

Open-ended generation: the output distribution still has high freedom

Non-open-ended generation: the input mostly determines the output generation.



Emergent Properties of LMs



Sources: "Compute trends across three eras of machine learning", by J. Sevilla et al., arXiv, 2022; Our World in Data

Emergent Properties of LMs

Let's revisit the Generative Pretrained Transformer (GPT) models from OpenAI as an example:

GPT-2 (1.5B parameters; Radford et al., 2019)

- Same architecture as GPT, just bigger (117M -> 1.5B)
- But trained on much more data: 4GB -> 40GB of internet text data (WebText)
- Scrape links posted on Reddit w/ at least 3 upvotes (rough proxy of human quality)

[Language models are unsupervised multitask learners](#)

Language Models are Unsupervised Multitask Learners

Alec Radford^{*1} Jeffrey Wu^{*1} Rewon Child¹ David Luan¹ Dario Amodei^{**1} Ilya Sutskever¹

Abstract

Natural language processing tasks, such as question answering, machine translation, reading comprehension, and summarization, are typically approached with supervised learning on task-specific datasets. We demonstrate that language models begin to learn these tasks without any explicit supervision when trained on a new dataset of millions of webpages called WebText. When conditioned on a document plus questions, the answers generated by the language model reach 55 F1 on the CoQA dataset - matching or exceeding the performance of 3 out of 4 baseline systems without using the 127,000+ training examples. The capacity of the language model is essential to the success of zero-shot task transfer and increasing it improves performance in a log-linear fashion across tasks. Our largest model, GPT-2,

competent generalists. We would like to make general systems which can perform many tasks without the need to manually create and curate a dataset for each one.

The dominant approach to creating ML systems is to select a dataset of training examples demonstrating the desired behavior for a desired task, train a system to learn these behaviors, and then test its performance on a new, identically distributed (IID) held-out dataset. This has served well to make progress on narrow tasks, but the often erratic behavior of captioning models (Ji et al., 2017), reading comprehension systems (Jawahar et al., 2018), and image classifiers (Alcorn et al., 2018) on a wide variety of possible inputs highlights shortcomings of this approach.

Our suspicion is that the prevalence of such behavior on single domain datasets is a major consequence of generalization observed in current systems.

Emergent Properties of LMs

One key emergent ability in GPT-2 is **zero-shot learning**:

the ability to do many tasks with no examples, and no gradient updates, by simply:

- Specifying the right sequence prediction problem (e.g. question answering):

Passage: Tom Brady... Q: Where was Tom Brady born? A: ...

[Language models are unsupervised multitask learners](#)

Language Models are Unsupervised Multitask Learners

Alec Radford^{*1} Jeffrey Wu^{*1} Rewon Child¹ David Luan¹ Dario Amodei^{**1} Ilya Sutskever¹

Abstract

Natural language processing tasks, such as question answering, machine translation, reading comprehension, and summarization, are typically approached with supervised learning on task-specific datasets. We demonstrate that language models begin to learn these tasks without any explicit supervision when trained on a new dataset of millions of webpages called WebText. When conditioned on a document plus questions, the answers generated by the language model reach 55 F1 on the CoQA dataset - matching or exceeding the performance of 3 out of 4 baseline systems without using the 127,000+ training examples. The capacity of the language model is essential to the success of zero-shot task transfer and increasing it improves performance in a log-linear fashion across tasks. Our largest model, GPT-2,

competent generalists. We would like to see more general systems which can perform many tasks without the need to manually create and train a dataset for each one.

The dominant approach to creating ML systems is to select a dataset of training examples demonstrating the desired behavior for a desired task, train a system to learn that behavior, and then test its performance on a new, identically distributed (IID) held-out dataset. This has served well to make progress on narrow tasks, but the often erratic behavior of captioning models (Ji et al., 2017), reading comprehension systems (Jawahar et al., 2018), and image classifiers (Alcorn et al., 2018) on a wide variety of possible inputs highlights the limitations of this approach.

Our suspicion is that the prevalence of such behavior on single domain datasets is a major contributor to the lack of generalization observed in current systems.

Emergent zero-shot and few-shot learning

Specify a task by simply prepending examples of the task before your example

- Also called **in-context learning**, to stress that no gradient updates are performed when learning a new task

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

Abstract

We demonstrate that scaling up language models greatly improves task-agnostic, few-shot performance, sometimes even becoming competitive with prior state-of-the-art fine-tuning approaches. Specifically, we train GPT-3, an autoregressive language model with 175 billion parameters, 10x more than any previous non-sparse language model, and test its performance in the few-shot setting. For all tasks, GPT-3 is applied without any gradient updates or fine-tuning, with tasks and few-shot demonstrations specified purely via text interaction with the model. GPT-3 achieves strong performance on many NLP datasets, including translation, question-answering, and cloze tasks. We also identify some datasets where GPT-3's few-shot learning still struggles, as well as some datasets where GPT-3 faces methodological issues related to training on large web corpora.

Instruction Finetuning

Before: sentence completion, not answering a prompt.

Language models are not aligned with user intent

PROMPT *Explain the moon landing to a 6 year old in a few sentences.*

COMPLETION GPT-3

Explain the theory of gravity to a 6 year old.

Explain the theory of relativity to a 6 year old in a few sentences.

Explain the big bang theory to a 6 year old.

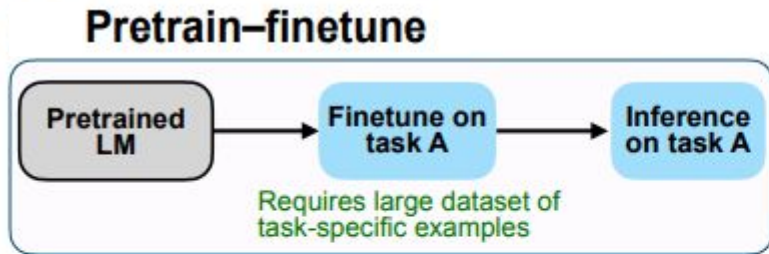
Explain evolution to a 6 year old.

Instruction Finetuning

Before: sentence completion, not answering a prompt.

Language models are not aligned with user intent

Instruction Finetuning!



Finetuned Language Models Are Zero-Shot Learners

Training language models to follow instructions with human feedback

Long Ouyang* Jeff Wu* Xu Jiang* Diogo Almeida* Carroll L. Wainwright*

Pamela Mishkin* Chong Zhang Sandhini Agarwal Katarina Slama Alex Ray

John Schulman Jacob Hilton Fraser Kelton Luke Miller Maddie Simens

Amanda Askell†

Peter Welinder

Paul Christiano*†

Jan Leike*

Ryan Lowe*

OpenAI

Abstract

Making language models bigger does not inherently make them better at following a user's intent. For example, large language models can generate outputs that are untruthful, toxic, or simply not helpful to the user. In other words, these models are not *aligned* with their users. In this paper, we show an avenue for aligning language models with user intent on a wide range of tasks by fine-tuning with human feedback. Starting with a set of labeler-written prompts and prompts submitted through a language model API, we collect a dataset of labeler demonstrations of the desired model behavior, which we use to fine-tune GPT-3 using supervised learning. We then collect a dataset of rankings of model outputs, which we use to further fine-tune this supervised model using reinforcement learning from human feedback. We call the resulting models *InstructGPT*. In human evaluations on our prompt distribution, outputs from the 1.3B parameter InstructGPT model are preferred to outputs from the 175B GPT-3, despite having 100x fewer parameters. Moreover, InstructGPT models show improvements in truthfulness and reductions in toxic content generation while having minimal performance regression on public

Training language models to follow instructions with human feedback

at fine-tuning with human feedback is a promising direction for aligning language models with human intent.

Instruction Finetuning

Before: sentence completion, not answering a prompt.

Language models are not aligned with user intent

Instruction Finetuning!



Finetuned Language Models Are Zero-Shot Learners

Training language models to follow instructions with human feedback

Long Ouyang* Jeff Wu* Xu Jiang* Diogo Almeida* Carroll L. Wainwright*

Pamela Mishkin* Chong Zhang Sandhini Agarwal Katarina Slama Alex Ray

John Schulman Jacob Hilton Fraser Kelton Luke Miller Maddie Simens

Amanda Askell†

Peter Welinder

Paul Christiano*†

Jan Leike*

Ryan Lowe*

OpenAI

Abstract

Making language models bigger does not inherently make them better at following a user's intent. For example, large language models can generate outputs that are untruthful, toxic, or simply not helpful to the user. In other words, these models are not *aligned* with their users. In this paper, we show an avenue for aligning language models with user intent on a wide range of tasks by fine-tuning with human feedback. Starting with a set of labeler-written prompts and prompts submitted through a language model API, we collect a dataset of labeler demonstrations of the desired model behavior, which we use to fine-tune GPT-3 using supervised learning. We then collect a dataset of rankings of model outputs, which we use to further fine-tune this supervised model using reinforcement learning from human feedback. We call the resulting models *InstructGPT*. In human evaluations on our prompt distribution, outputs from the 1.3B parameter InstructGPT model are preferred to outputs from the 175B GPT-3, despite having 100x fewer parameters. Moreover, InstructGPT models show improvements in truthfulness and reductions in toxic content generation while having minimal performance regression on public

Training language models to follow instructions with human feedback

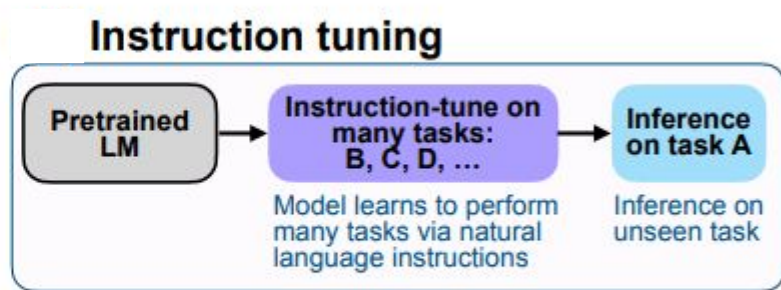
at fine-tuning with human feedback is a promising direction for aligning models with human intent.

Instruction Finetuning

Before: sentence completion, not answering a prompt.

Language models are not aligned with user intent

Instruction Finetuning!



Finetuned Language Models Are Zero-Shot Learners

Training language models to follow instructions with human feedback

Long Ouyang* Jeff Wu* Xu Jiang* Diogo Almeida* Carroll L. Wainwright*

Pamela Mishkin* Chong Zhang Sandhini Agarwal Katarina Slama Alex Ray

John Schulman Jacob Hilton Fraser Kelton Luke Miller Maddie Simens

Amanda Askell†

Peter Welinder

Paul Christiano*†

Jan Leike*

Ryan Lowe*

OpenAI

Abstract

Making language models bigger does not inherently make them better at following a user's intent. For example, large language models can generate outputs that are untruthful, toxic, or simply not helpful to the user. In other words, these models are not *aligned* with their users. In this paper, we show an avenue for aligning language models with user intent on a wide range of tasks by fine-tuning with human feedback. Starting with a set of labeler-written prompts and prompts submitted through a language model API, we collect a dataset of labeler demonstrations of the desired model behavior, which we use to fine-tune GPT-3 using supervised learning. We then collect a dataset of rankings of model outputs, which we use to further fine-tune this supervised model using reinforcement learning from human feedback. We call the resulting models *InstructGPT*. In human evaluations on our prompt distribution, outputs from the 1.3B parameter InstructGPT model are preferred to outputs from the 175B GPT-3, despite having 100x fewer parameters. Moreover, InstructGPT models show improvements in truthfulness and reductions in toxic content generation while having minimal performance regression on public

Training language models to follow instructions with human feedback

at fine-tuning with human feedback is a promising direction for aligning language models with human intent.

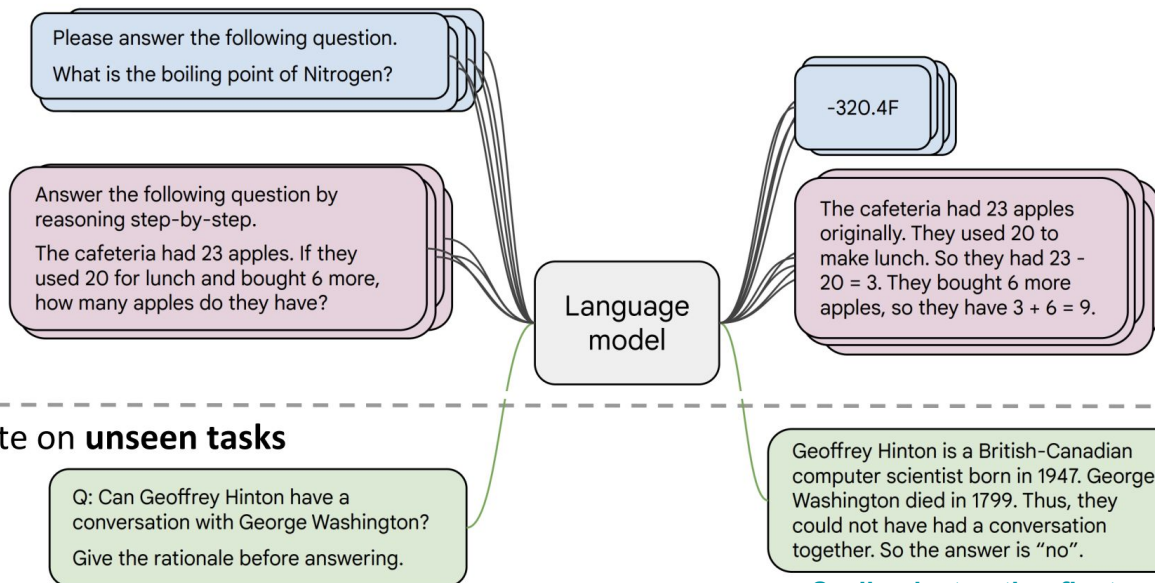
Instruction Finetuning

Also in Flan-T5

Scaling Instruction-Finetuned Language Models

Hyung Won Chung* Le Hou* Shayne Longpre* Barret Zoph† Yi Tay†
William Fedus† Yunxuan Li Xuezhi Wang Mostafa Dehghani Siddhartha Brahma
Albert Webson Shixiang Shane Gu Zhuyun Dai Mirac Suzgun Xinyun Chen
Aakanksha Chowdhery Alex Castro-Ros Marie Pellat Kevin Robinson
Dasha Valter Sharan Narang Gaurav Mishra Adams Yu Vincent Zhao
Yanning Huang Andrew Dai Hongkun Yu Slav Petrov Ed H. Chi
Jenny Zhou Quoc V. Le

- Collect examples of (instruction, output) pairs across many tasks and finetune an LM

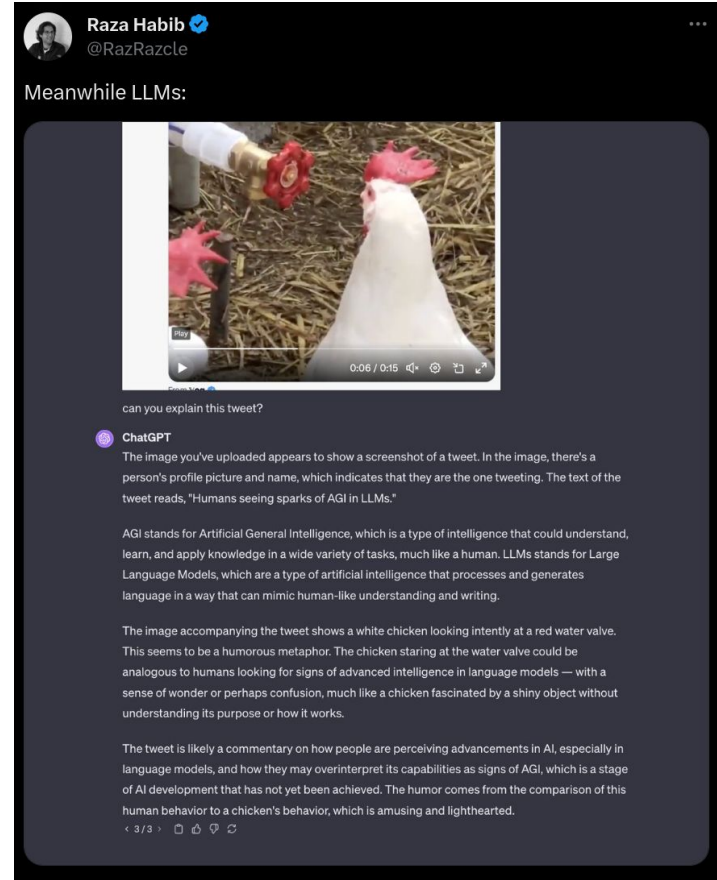


Instructions have been shown to improve performance as we explore instruction finetuning across different model sizes, and (3) finetuning on these aspects dramatically improves performance on various prompting setups (zero-shot, few-shot, CoT), text-to-text generation, RealToxicityPrompts). The model performs PaLM 540B by a large margin on several benchmarks, such as MMLU, which achieve strong few-shot performance. Overall, instruction finetuning is a promising approach for training language models.

- Evaluate on **unseen tasks**

[Scaling instruction-finetuned language models](#)

Sparks of AGI



<https://twitter.com/deliprao/status/1724163062830153814>

<https://twitter.com/RazRazcle/status/1724174924367974827>

Reinforcement Learning from Human Feedback (RLHF)

Instruction finetuning is great but has some limitations

Even with instruction finetuning, there a mismatch between the LM objective and the objective of “satisfy human preferences”!

Can we explicitly attempt to satisfy human preferences?

Reinforcement Learning from Human Feedback (RLHF)

Instruction finetuning is great but has some limitations

Even with instruction finetuning, there a mismatch between the LM objective and the objective of “satisfy human preferences”!

Can we explicitly attempt to satisfy human preferences?

Reinforcement Learning: area of machine learning and optimal control concerned with how an intelligent agent ought to take actions in a dynamic environment in order to maximize the cumulative reward.

In the InstructGPT paper, they also had RLHF

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-3.5 with supervised learning.



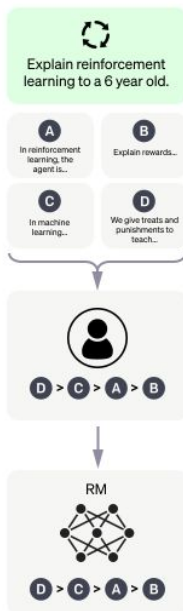
Step 2

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.



Training language models to follow instructions with human feedback

Training language models to follow instructions with human feedback

Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

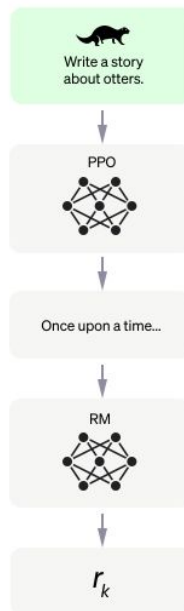
A new prompt is sampled from the dataset.

The PPO model is initialized from the supervised policy.

The policy generates an output.

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.



Imeida* Carroll L. Wainwright*

wal Katarina Slama Alex Ray

Luke Miller Maddie Simens

Paul Christiano*†

Ryan Lowe*

make them better at following instructions. In other words, these models are an avenue for aligning language models by fine-tuning with human prompts and prompts submitted by a set of labeler demonstrations of GPT-3 using supervised model outputs, which we use to train GPT-3 using reinforcement learning from human feedback (RLHF). In human evaluations on the parameter InstructGPT model are able to have 100x fewer parameters. In terms of truthfulness and reductions in performance regressions on public benchmarks, our results show promising direction for aligning

Tasks collected from labelers

- **Plain:** We simply ask the labelers to come up with an arbitrary task, while ensuring the tasks had sufficient diversity.
- **Few-shot:** We ask the labelers to come up with an instruction, and multiple query/response pairs for that instruction.
- **User-based:** We had a number of use-cases stated in waitlist applications to the OpenAI API. We asked labelers to come up with prompts corresponding to these use cases.

Use-case	Prompt
Brainstorming	List five ideas for how to regain enthusiasm for my career
Generation	Write a short story where a bear goes to the beach, makes friends with a seal, and then returns home.

The results

PROMPT *Explain the moon landing to a 6 year old in a few sentences.*

COMPLETION GPT-3

Explain the theory of gravity to a 6 year old.

Explain the theory of relativity to a 6 year old in a few sentences.

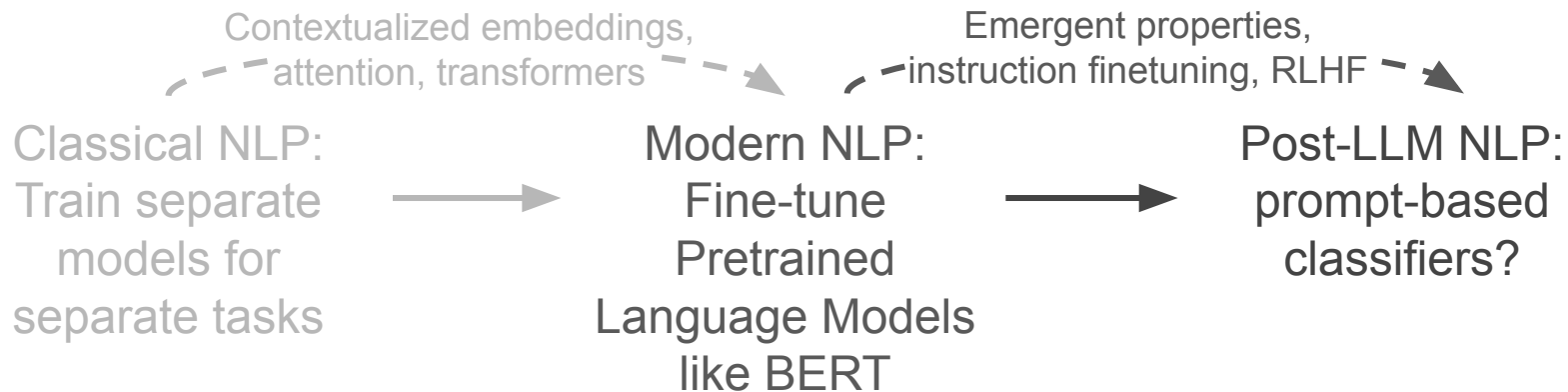
Explain the big bang theory to a 6 year old.

Explain evolution to a 6 year old.

InstructGPT

People went to the moon, and they took pictures of what they saw, and sent them back to the earth so we could all see them.

The timeline so far: Supervised NLP



Next week: debating content analysis with LLM

PNAS

BRIEF REPORT

POLITICAL SCIENCES

OPEN ACCESS



ChatGPT outperforms crowd workers for text-annotation tasks

Fabrizio Gilaridi^{A1}, Meysam Alizadeh^{A2}, and Maël Kubli^{A3}

Edited by Mary Waters, Harvard University, Cambridge, MA; received March 27, 2023; accepted June 2, 2023

Many NLP applications require manual text annotations for a variety of tasks, notably to train classifiers or evaluate the performance of unsupervised models. Depending on the size and degree of complexity, the tasks may be conducted by crowd workers on platforms such as MTurk as well as trained annotators, such as research assistants. Using four samples of tweets and news articles ($n = 6,183$), we show that ChatGPT outperforms crowd workers for several annotation tasks, including relevance, stance, topics, and frame detection. Across the four datasets, the zero-shot accuracy of ChatGPT exceeds that of crowd workers by about 25 percentage points on average, while ChatGPT's intercoder agreement exceeds that of both crowd workers and trained annotators for all tasks. Moreover, the per-annotation cost of ChatGPT is less than \$0.003—about thirty times cheaper than MTurk. These results demonstrate the potential of large language models to drastically increase the efficiency of text classification.

ChatGPT | text classification | large language models | human annotations | text as data

VS

Chatbots Are Not Reliable Text Annotators

Ross Deans Kristensen-McLachlan^{*ab}, Miceal Canavan^c, Márton Kardos^a,
Mia Jacobsen^a, and Lene Aarøe^{cd}

^aCenter for Humanities Computing

^bDepartment of Linguistics, Cognitive Science, and Semiotics

^cDepartment of Political Science

^dAarhus Institute of Advances Studies

November 13, 2023