

# **Modern LLMs & “post” training**

*Some slides and content today derived from materials by Mohit Iyyer (CS685 @ UMass) and Anna Rogers (“A Primer on BERTology”, TACL 2020)*

# Why isn't pre-training enough?

What we *want*: Generally useful models

What we *get*: Models capable of producing text capably. This is what we asked for! Given a big dataset of unlabeled data  $D$ :

$$\min_{\theta} - \frac{1}{|D|} \sum_{x \in D} \sum_{t=1}^{T-1} \log \pi_{\theta}(x_{t+1} \mid x_t, \dots, x_1)$$

# Why isn't pre-training enough?

What we *want*: Generally useful models

What we *get*: Models capable of producing text capably. This is what we asked for! Given a big dataset of unlabeled data  $D$ :

$$\min_{\theta} - \frac{1}{|D|} \sum_{x \in D} \sum_{t=1}^{T-1} \log \pi_{\theta}(x_{t+1} \mid x_t, \dots, x_1)$$

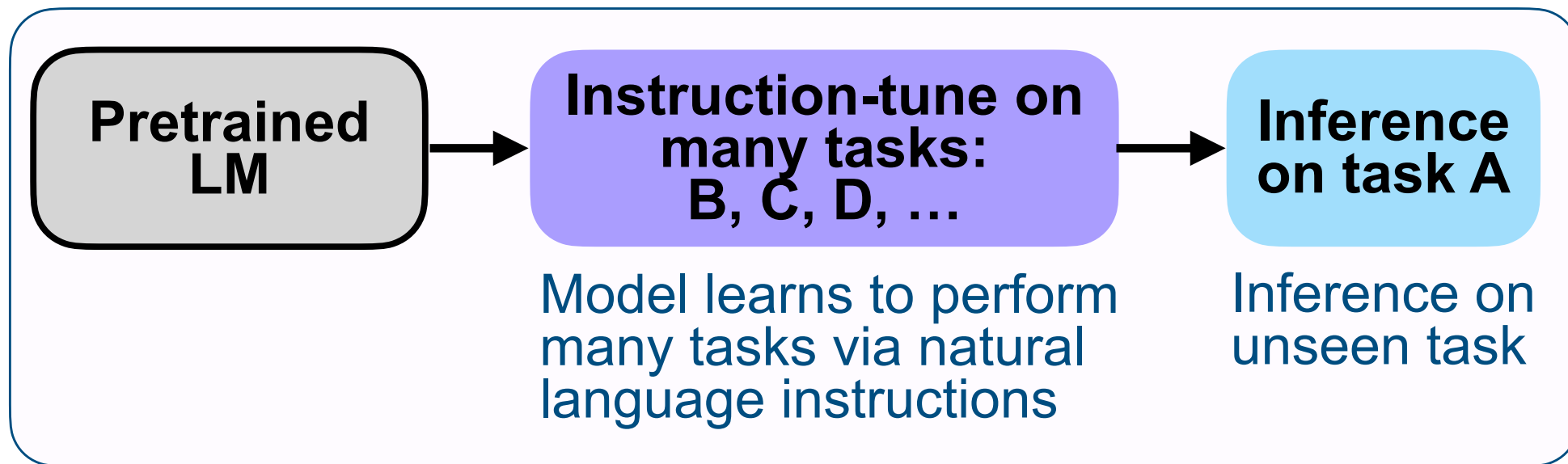
*Post*-training is the idea of “**aligning**” the model with what we want. This requires some sort of *supervision*.

# How to “align”

An active area of research

Two main strategies we'll discuss: *Instruction fine-tuning* and *Reinforcement Learning from Human Feedback (RLHF)*

# Instruction fine-tuning



FINETUNED LANGUAGE MODELS ARE ZERO-SHOT LEARNERS

Jason Wei\*, Maarten Bosma\*, Vincent Y. Zhao\*, Kelvin Guu\*, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le

Google Research

# Instruction fine-tuning

Finetune on many tasks (“instruction-tuning”)

**Input (Commonsense Reasoning)**

Here is a goal: Get a cool sleep on summer days.  
How would you accomplish this goal?  
OPTIONS:  
 -Keep stack of pillow cases in fridge.  
 -Keep stack of pillow cases in oven.

**Target**

keep stack of pillow cases in fridge

**Input (Translation)**

Translate this sentence to Spanish:  
The new office building was built in less than three months.

**Target**

El nuevo edificio de oficinas se construyó en tres meses.

Sentiment analysis tasks

Coreference resolution tasks

...

Inference on unseen task type

**Input (Natural Language Inference)**

Premise: At my age you will probably have learnt one lesson.  
Hypothesis: It's not certain how many lessons you'll learn by your thirties.  
Does the premise entail the hypothesis?  
OPTIONS:  
 -yes  -it is not possible to tell  -no

**FLAN Response**

It is not possible to tell

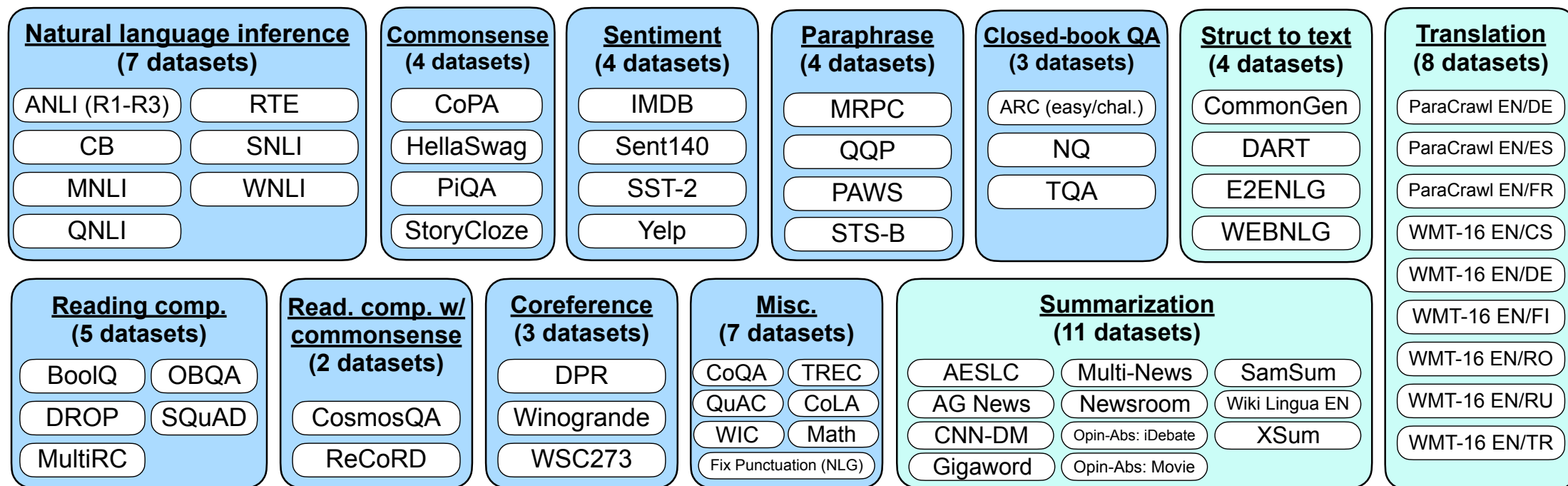


FINETUNED LANGUAGE MODELS ARE ZERO-SHOT LEARNERS

Jason Wei\*, Maarten Bosma\*, Vincent Y. Zhao\*, Kelvin Guu\*, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le

Google Research

# Instruction fine-tuning



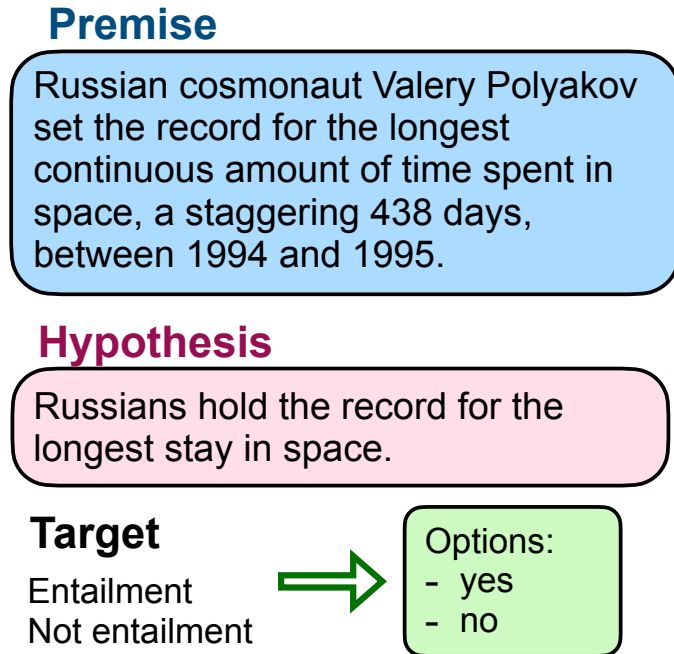
# FLAN

FINETUNED LANGUAGE MODELS ARE ZERO-SHOT LEARNERS

Jason Wei\*, Maarten Bosma\*, Vincent Y. Zhao\*, Kelvin Guu\*, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le

Google Research

# Multiple instruction templates per task



## Template 1

<premise>  
Based on the paragraph above, can we conclude that <hypothesis>?  
<options>

## Template 2

<premise>  
Can we infer the following?  
<hypothesis>  
<options>

## Template 3

Read the following and determine if the hypothesis can be inferred from the premise:  
Premise: <premise>  
Hypothesis: <hypothesis>  
<options>

## Template 4, ...

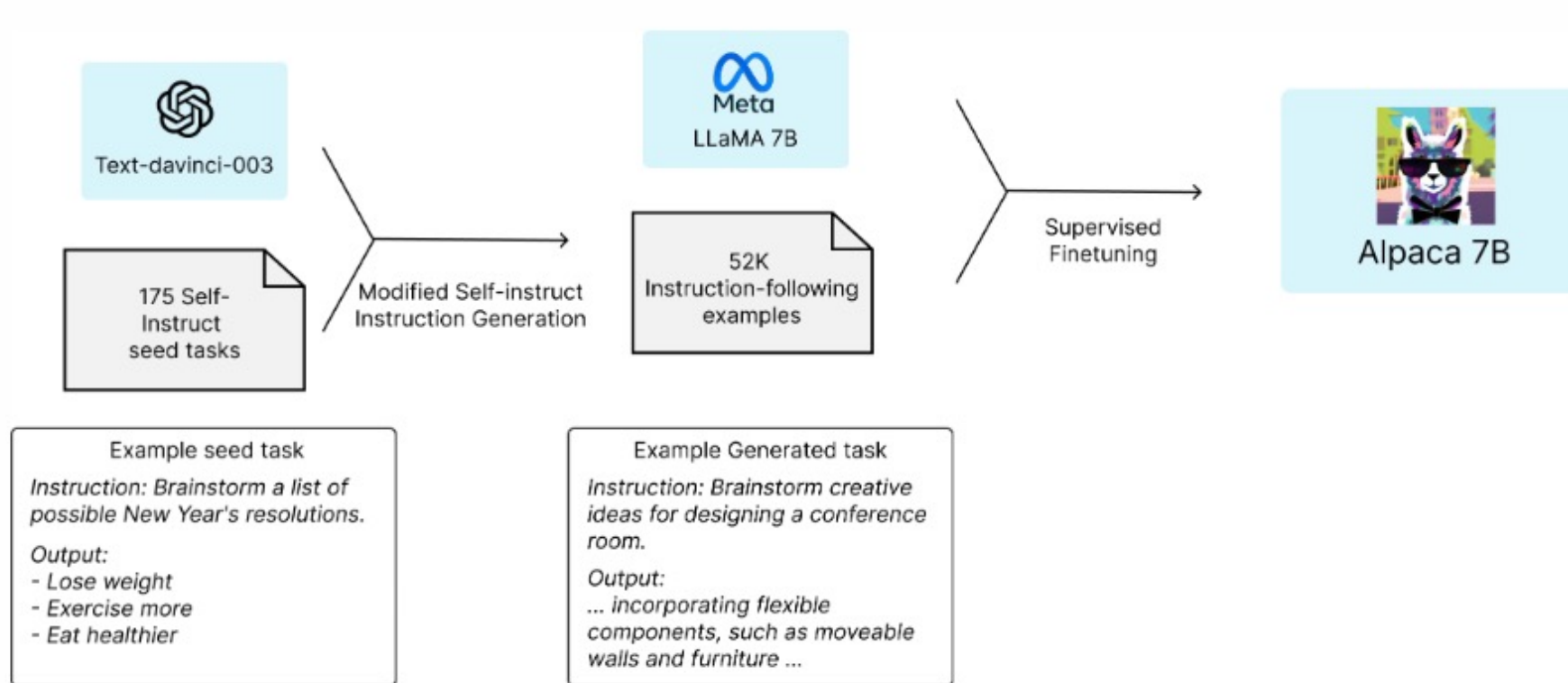
FINETUNED LANGUAGE MODELS ARE ZERO-SHOT LEARNERS

Jason Wei\*, Maarten Bosma\*, Vincent Y. Zhao\*, Kelvin Guu\*, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le

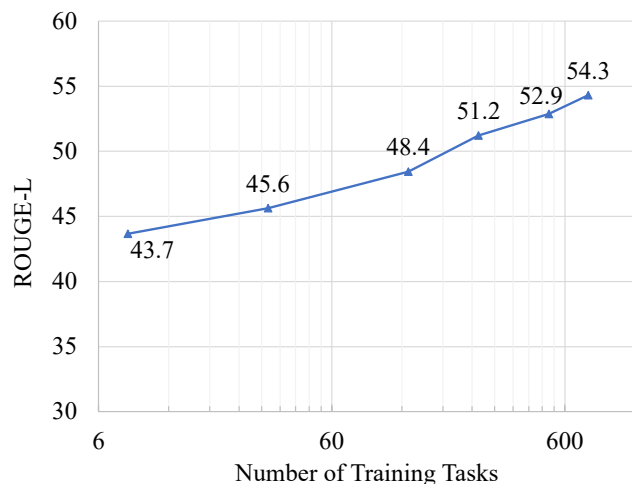
Google Research



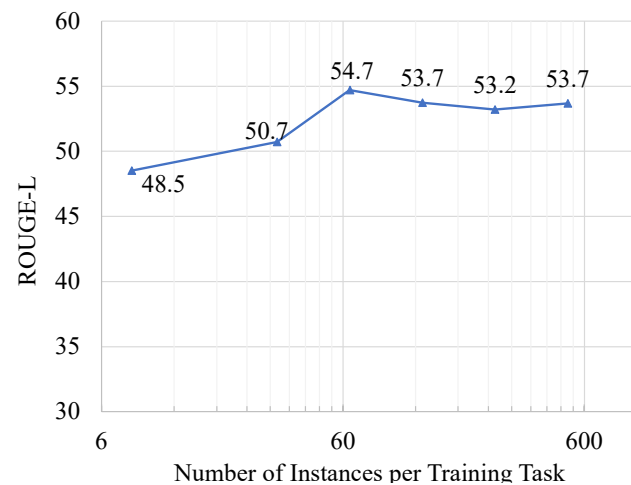
# Alpaca: Deriving training data from LLMs



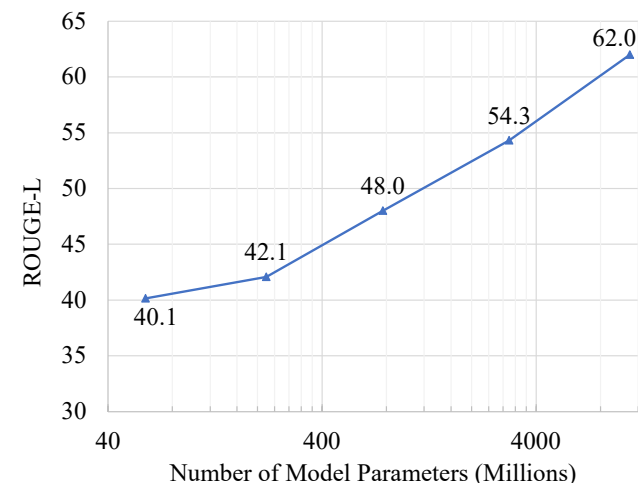




(a)



(b)



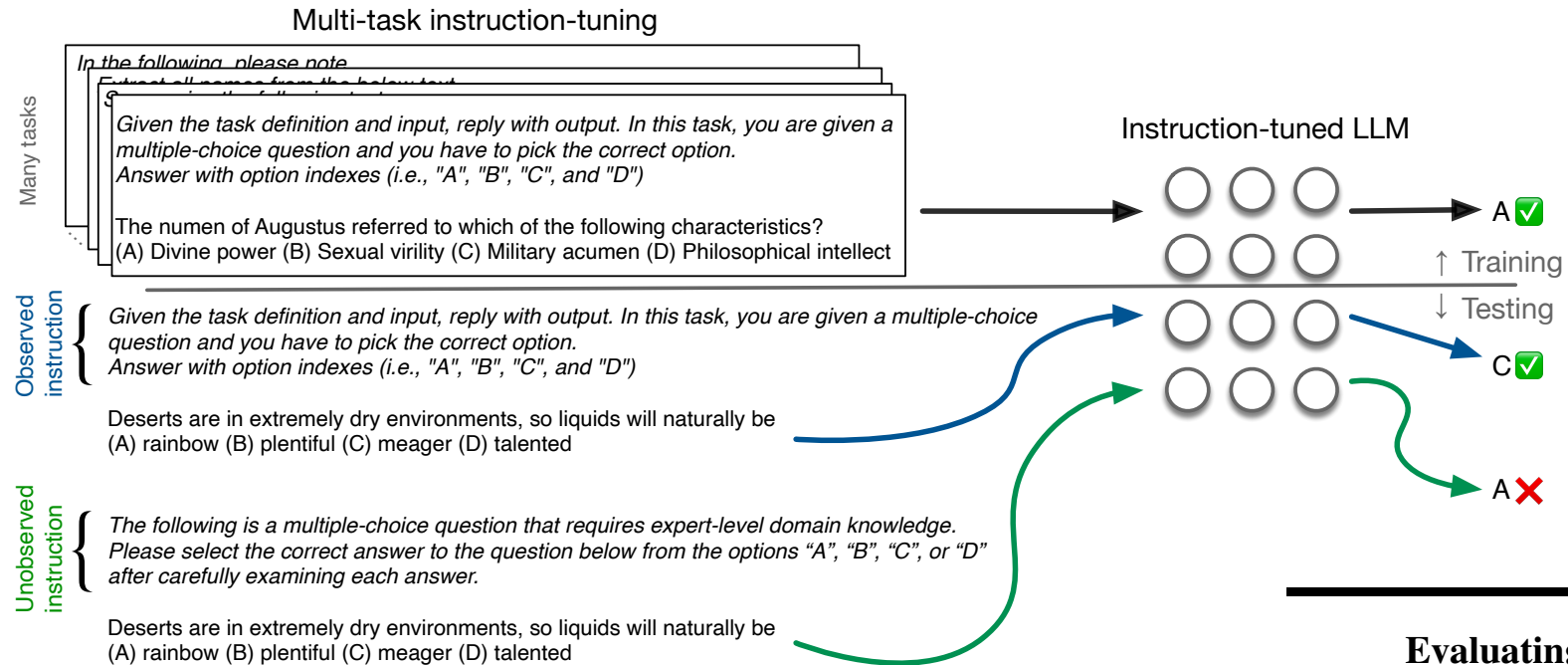
(c)

Figure 5: Scaling trends of models performance (§7.1) as a function of (a) the number of training tasks; (b) the number of instances per training task; (c) model sizes.  $x$ -axes are in log scale. The **linear growth of model performance with exponential increase in observed tasks and model size** is a promising trend. Evidently, the performance gain from more instances is limited.

### SUPER-NATURALINSTRUCTIONS: Generalization via Declarative Instructions on 1600+ NLP Tasks

◇Yizhong Wang<sup>2</sup> ◇Swaroop Mishra<sup>3</sup> \*Pegah Alipoormolabashi<sup>4</sup> \*Yeganeh Kordi<sup>5</sup>  
 Amirreza Mirzaei<sup>4</sup> Anjana Arunkumar<sup>3</sup> Arjun Ashok<sup>6</sup> Arut Selvan Dhanasekaran<sup>3</sup>  
 Atharva Naik<sup>7</sup> David Stap<sup>8</sup> Eshaan Pathak<sup>9</sup> Giannis Karamanolakis<sup>10</sup> Haizhi Gary Lai<sup>11</sup>  
 Ishan Purohit<sup>12</sup> Ishani Mondal<sup>13</sup> Jacob Anderson<sup>3</sup> Kirby Kuznia<sup>3</sup> Krima Doshi<sup>3</sup> Maitreya Patel<sup>3</sup>  
 Kuntal Kumar Pal<sup>3</sup> Mehrad Moradshahi<sup>14</sup> Mihir Parmar<sup>3</sup> Mirali Purohit<sup>15</sup> Neeraj Varshney<sup>3</sup>  
 Phani Rohitha Kaza<sup>3</sup> Pulkit Verma<sup>3</sup> Ravsehaj Singh Puri<sup>3</sup> Rushang Karia<sup>3</sup> Shailaja Keyur Sampat<sup>3</sup>  
 Savan Doshi<sup>3</sup> Siddhartha Mishra<sup>16</sup> Sujan Reddy<sup>17</sup> Sumanta Patro<sup>18</sup> Tanay Dixit<sup>19</sup> Xudong Shen<sup>20</sup>  
 Chitta Baral<sup>3</sup> Yejin Choi<sup>1,2</sup> Noah A. Smith<sup>1,2</sup> Hannaneh Hajishirzi<sup>1,2</sup> Daniel Khoshdel<sup>21</sup>

# But!



## Evaluating the Zero-shot Robustness of Instruction-tuned Language Models

**Jiuding Sun**  
Khoury College of Computer Sciences  
Northeastern University  
sun.jiu@northeastern.edu

**Chantal Shaib**  
Khoury College of Computer Sciences  
Northeastern University  
shaib.c@northeastern.edu

**Byron C. Wallace**  
Khoury College of Computer Sciences  
Northeastern University  
b.wallace@northeastern.edu

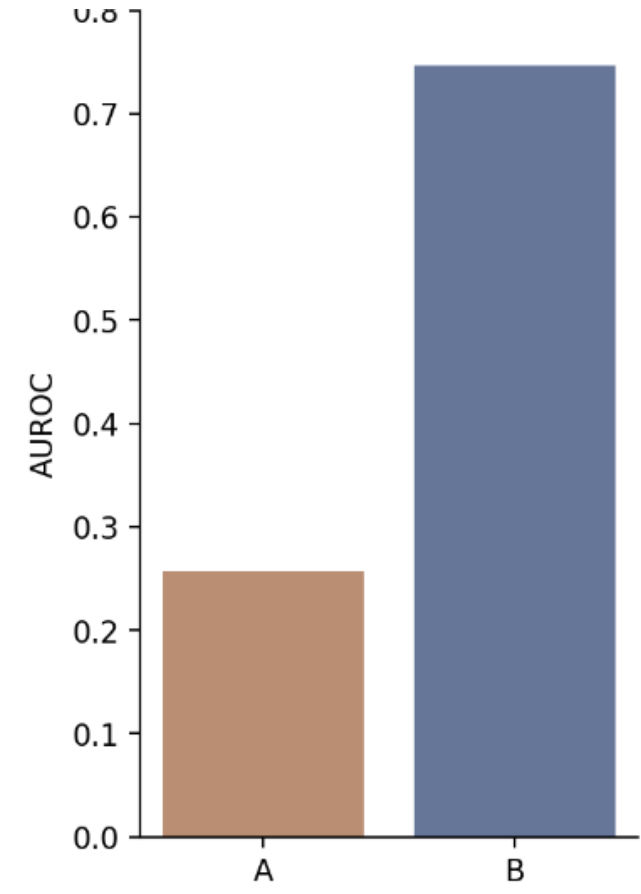


A

Analyze the patient's longitudinal ... include 'Yes' if the patient meets the definition of alcohol abuse, and 'No' if they do not.

B

In the above notes, please ... client's AUDIT score (input example- Alcohol Use Disorders Identification Test). Using ... yes or no answer.



### Open (Clinical) LLMs are Sensitive to Instruction Phrasings

Alberto Mario Ceballos Arroyo\*<sup>γ</sup> Monica Munnangi\*<sup>γ</sup> Jiuding Sun<sup>γ</sup>  
Karen Y.C. Zhang<sup>γ</sup> Denis Jered McInerney<sup>γ</sup>◇ Denis Jered McInerney<sup>γ</sup>◇ Byron C. Wallace<sup>γ</sup> Silvio Amir<sup>γ</sup>  
<sup>γ</sup>Northeastern University ◇Codametrix

{ceballosarroyo.a, munnangi.m, sun.jiu, zhang.yuchen, b.wallace,s.amir}@northeastern.edu

jmcinerney@codametrix.com

# Human preferences

Often more natural to elicit *preferences* between pairs of outputs than to provide explicit examples

For instance, if we want LLMs to generate “more polite” or less biased outputs, difficult to write a bunch of examples explicitly demonstrating these things:  
Easier to show two examples and ask which is “more polite”

# (Reinforcement) Learning from Human Feedback

## 1 Collect human feedback

A Reddit post is sampled from the Reddit TL;DR dataset.



Various policies are used to sample a set of summaries.



Two summaries are selected for evaluation.



A human judges which is a better summary of the post.



"j is better than k"

## 2 Train reward model

One post with two summaries judged by a human are fed to the reward model.



The reward model calculates a reward  $r$  for each summary.



The loss is calculated based on the rewards and human label, and is used to update the reward model.

$$r_j \quad r_k$$

$$\text{loss} = \log(\sigma(r_j - r_k))$$

"j is better than k"

## 3 Train policy with PPO

A new post is sampled from the dataset.



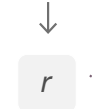
The policy  $\pi$  generates a summary for the post.



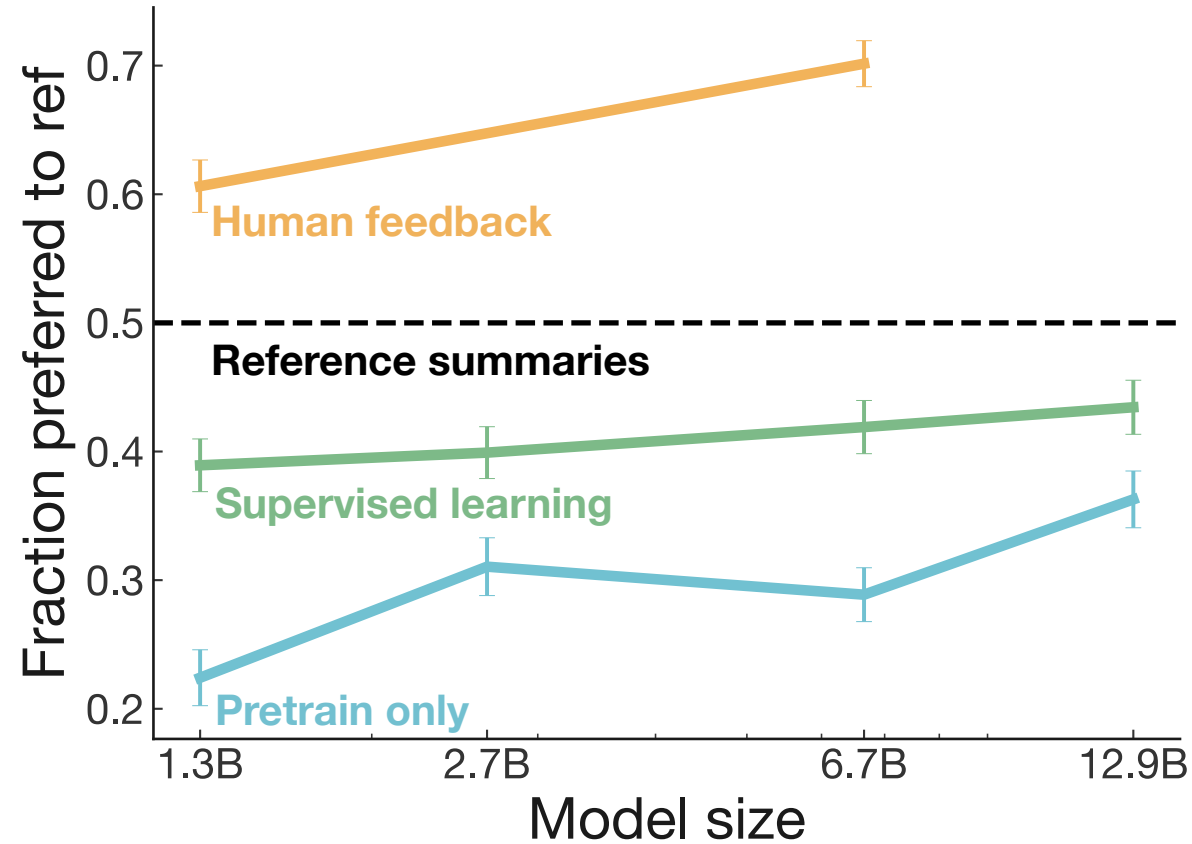
The reward model calculates a reward for the summary.



The reward is used to update the policy via PPO.



# (Reinforcement) Learning from Human Feedback



---

Learning to summarize from human feedback

---

Nisan Stiennon\* Long Ouyang\* Jeff Wu\* Daniel M. Ziegler\* Ryan Lowe\*  
Chelsea Voss\* Alec Radford Dario Amodei Paul Christiano\*



**Let's talk RL & PPO**  
**[see notes]**

# But who wants to deal w/RL?

Direct Preference Optimization (DPO) says: Oh, we can just use supervised learning to directly optimize for preference feedback labels

---

**Direct Preference Optimization:  
Your Language Model is Secretly a Reward Model**

---

Rafael Rafailov\*†

Archit Sharma\*†

Eric Mitchell\*†

Stefano Ermon†‡

Christopher D. Manning†

Chelsea Finn†

# The objective

$$\mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_{\theta}(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi_{\theta}(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right]$$

(w preferred to l)

### TL;DR Summarization Win Rate vs Reference

