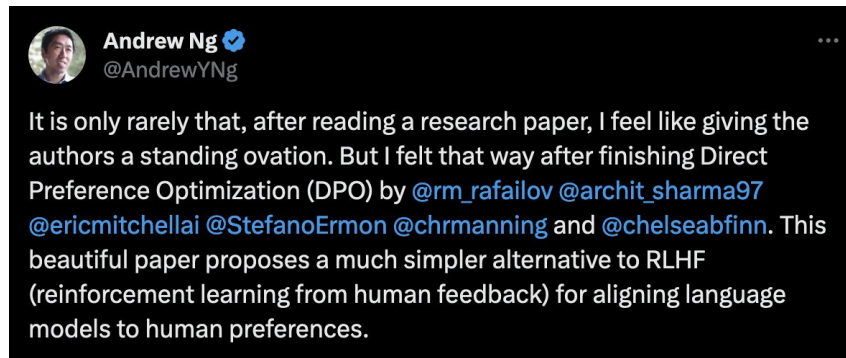# Direct Preference Optimization Explained, really

Roger Lam
Iamroger.com

# What's the hype?

Andrew Ng, Co-Founder of Coursera, Stanford CS adjunct faculty and former head of Baidu AI Group/Google Brain, wanted to give the authors a "standing ovation".

"This beautiful paper proposes a much simpler alternative to RLHF…"

Andrew Ng ✔
@AndrewYNg

It is only rarely that, after reading a research paper, I feel like giving the authors a standing ovation. But I felt that way after finishing Direct Preference Optimization (DPO) by @rm_rafailov @archit_sharma97 @ericmitchellai @StefanoErmon @chrmanning and @chelseabfinn. This beautiful paper proposes a much simpler alternative to RLHF (reinforcement learning from human feedback) for aligning language models to human preferences.
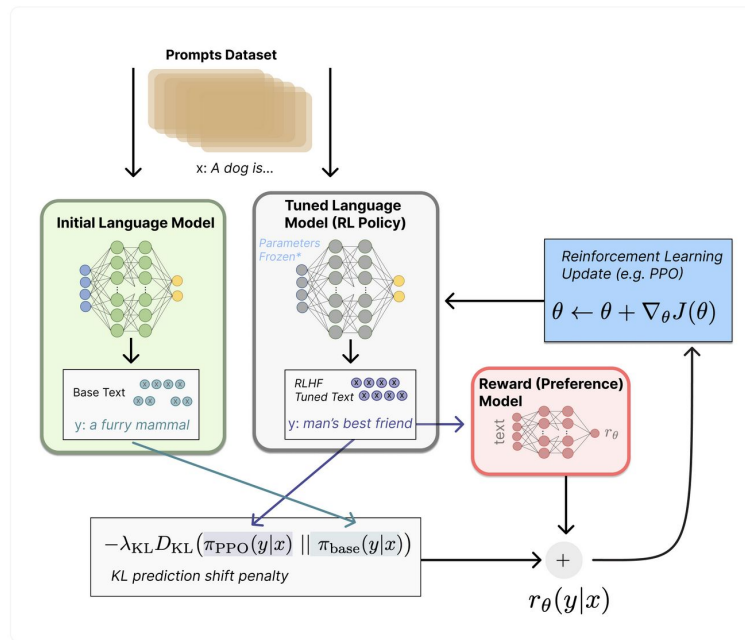
# Reinforcement Learning from Human Feedback

RLHF is performed after the initial LLM is trained to align its responses.

It uses two candidate outputs, asks a reward model what it thinks, and tweaks accordingly.

What's tough is that RLHF requires creating a reward model (red) to fine-tune based on human preference.

That makes it pretty heavy in complexity and compute.



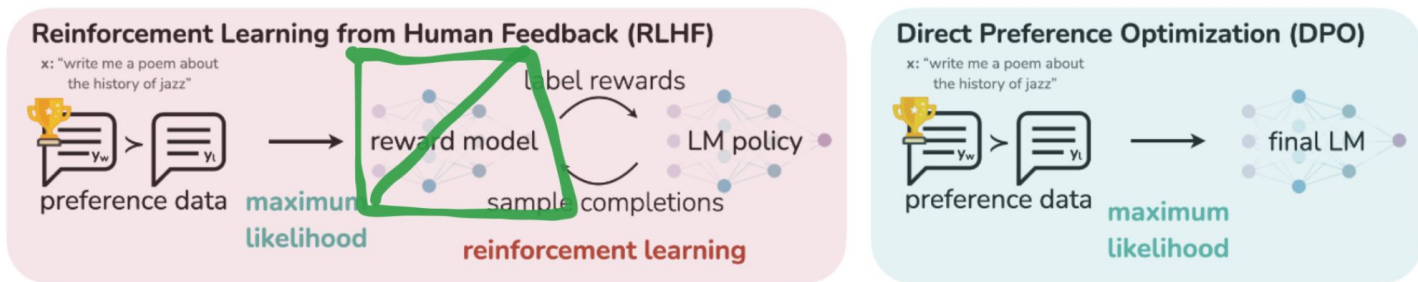https://huggingface.co/blog/rlhf

# What is Direct Preference Optimization?

DPO is a faster and more elegant way to tune your LLM than the traditional RLHF way.

It drops the reward model and instead uses the preference data directly to fine-tune the model.



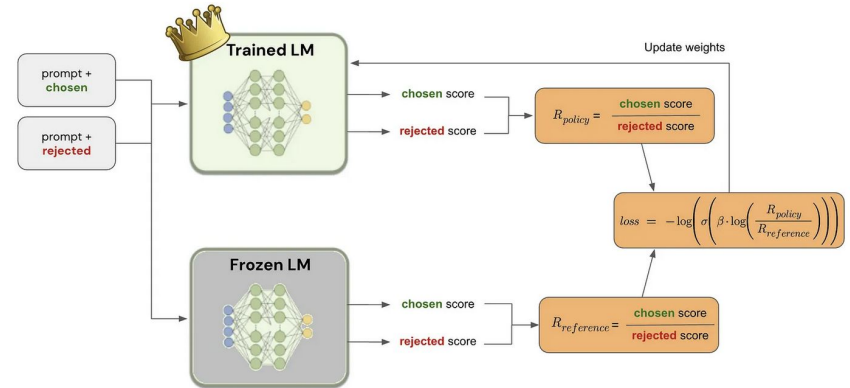https://arxiv.org/pdf/2305.18290.pdf

# DPO

Through some clever math, the authors discovered you could derive the optimal policy model and therefore the loss function mathematically.

Together with asking what itself thinks the score of each response should be, the fine-tuned LLM adjusts based on the results.

# Results

At relatively small scale (6B parameters), DPO works maybe even a little better than traditional PPO / RLHF.
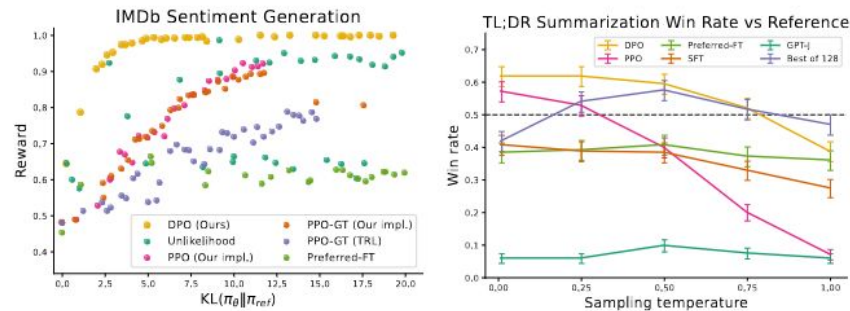
Very promising! Free gains!



Figure 2: **Left.** The frontier of expected reward vs KL to the reference policy. DPO provides the highest expected reward for all KL values, demonstrating the quality of the optimization. **Right.** TL;DR summarization win rates vs. human-written summaries, using GPT-4 as evaluator. DPO exceeds PPO's best-case performance on summarization, while being more robust to changes in the sampling temperature.
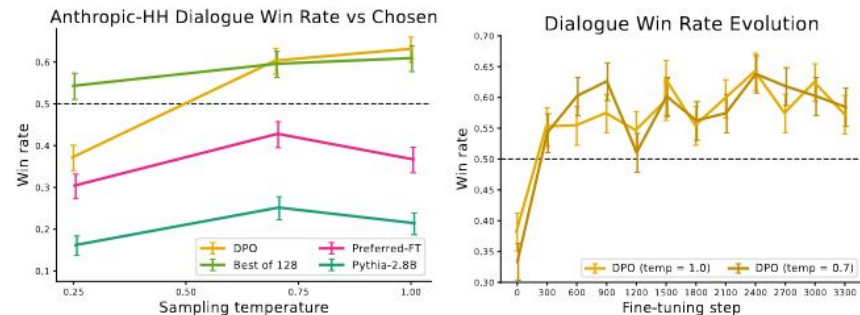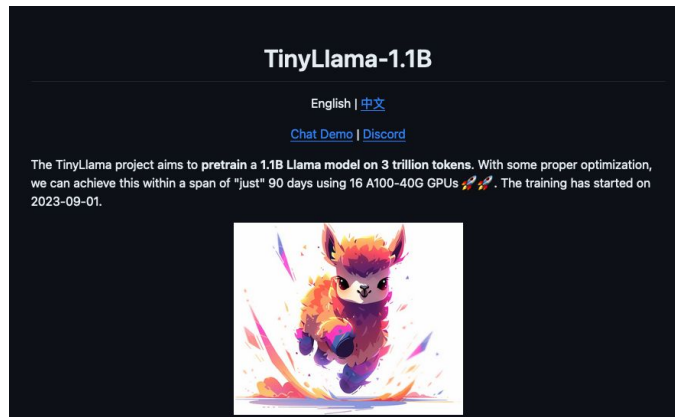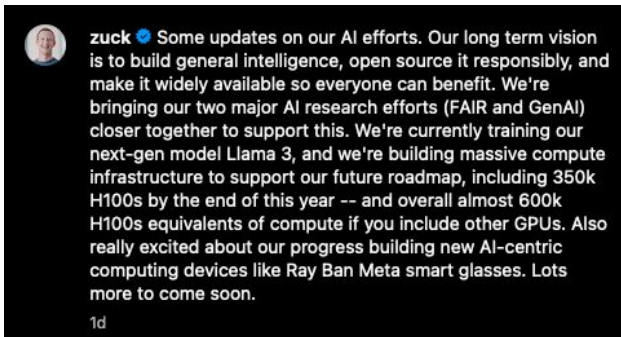


Figure 3: **Left.** Win rates computed by GPT-4 for Anthropic-HH one-step dialogue; DPO is the only method that improves over chosen summaries in the Anthropic-HH test set. **Right.** Win rates for different sampling temperatures over the course of training. DPO's improvement over the dataset labels is fairly stable over the course of training for different sampling temperatures.

# Why so exciting?

- We still have so much to discover!
- It's beautifully succinct.
- Optimizations like these lower the barrier to LLMs.
  - GPU-rich vs GPU-poor
  - What else can we do to make AI/ML accessible?

# Resources

- [DPO Paper](#)
- Andrew Ng explains DPO much [more elegantly](#)
- [Practical Simplified DPO](#) by [João Lages](#)
- [Mathematical Explanation of DPO](#) by [Pakhapoom Sarapat](#)
- [RLHF video](#) and [blog post](#) by Huggingface

# Thank You

If this was helpful, please share with your network, follow, subscribe, all the good things.

Roger Lam

https://www.linkedin.com/in/lam-roger/

https://www.lamroger.com/