# Learning from Feedback

Haw-Shiuan Chang

# Deadlines

- **https://people.cs.umass.edu/~hschang/cs685/schedule.html**

- **3/14**: HW 1 due
- **3/17:** Quiz 3
- **4/11**: HW 2 due
  - Will be released before the spring break

- **5/9**: Last day to submit extra credit
  - Please check the announcement at Piazza for the recording link

# Fine-tuning mostly Changes the Style

**Query:** How does actor critic improve over REINFORCE?

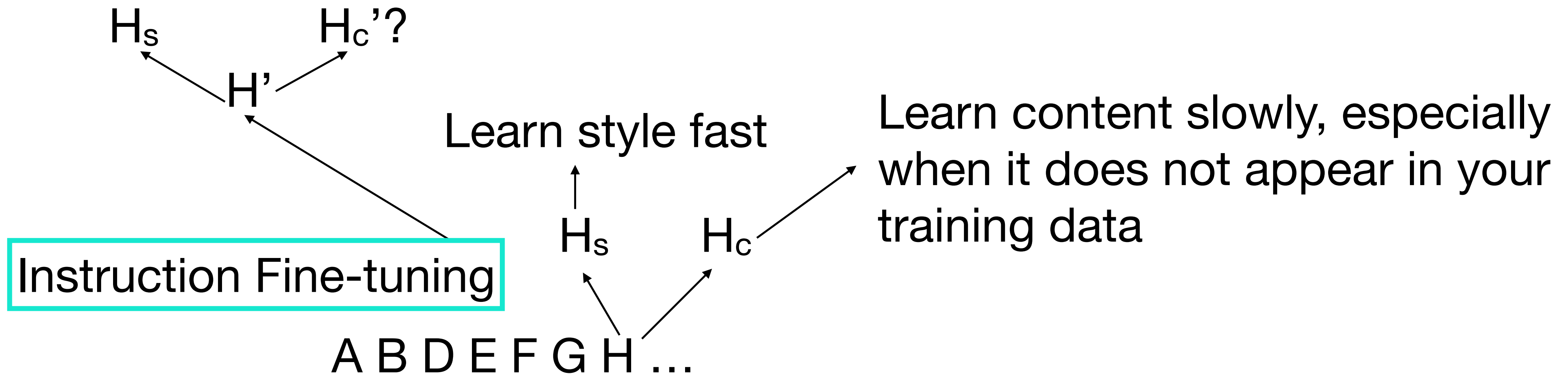| ChatGPT Response | Imitation Model Response |
|---|---|
| Actor-critic algorithms are an extension of the REINFORCE algorithm that combine both policy-based and value-based methods. Some of the changes that actor-critic algorithms make over standard REINFORCE: | Actor-critic algorithms are a type of reinforcement learning algorithm that improves the REINFORCE algorithm by combining the actor (policy) and critic (value) components. The main changes that actor-critic algorithms make over standard REINFORCE are: |
| 1. Function Approximation: Actor-critic algorithms typically use function approximation methods to represent the value and policy functions. This allows them to learn in high-dimensional state spaces. | 1. Actor-critic algorithms use a single critic to evaluate the value function of the policy, whereas standard REINFORCE uses a separate critic for each action. |
| 2. Critic Network: In actor-critic algorithms, there is an additional critic network that estimates the value function. The critic provides an estimate of the value of a state, which is then used to update the policy. This reduces the variance of the update and stabilizes learning. | 2. Actor-critic algorithms use a single actor to determine the actions to take, whereas standard REINFORCE uses a separate actor for each action. |
| 3. Advantage Estimation: Actor-critic algorithms use the advantage function to estimate the quality of an action taken by the policy. The advantage function is the difference between the return of an action and the estimated value of a state. This reduces the variance of the update and allows for more efficient learning. | 3. Actor-critic algorithms learn the policy and actor simultaneously, whereas standard REINFORCE learns the policy and actor separately. |
| | 4. Actor-critic algorithms use an optimistic estimate of the value function to determine the actor's actions, whereas standard REINFORCE uses the critic's estimate. |
| Overall, actor-critic algorithms provide a more efficient and stable way of learning in reinforcement learning tasks, particularly in high-dimensional state spaces. | Overall, actor-critic algorithms combine the critic's and actor's estimates in order to learn both the policy and actor simultaneously, while standard REINFORCE learns them separately. This can lead to faster convergence and potentially better generalization, but the choice between the two algorithms depends on the specific problem and the desired behavior of the agent. |

The False Promise of Imitating Proprietary LLMs (https://arxiv.org/pdf/2305.15717)

3

# Limitations of Fine-tuning

Cannot find it in the pretraining data -> Hallucination

$H_s$     $H_c$'?

$H'$

Learn style fast     Learn content slowly, especially when it does not appear in your training data

Instruction Fine-tuning

$H_s$     $H_c$

A B D E F G H …

+

Memorized pretraining corpus:

A*03$AFJ@*(!c()@kflm!@!cnvaodi

# Old Fine-Tuning

- Main purpose
  - Learn to do the task

- Data Quality Requirement
  - Usually medium

- Fine-tuning compared to few-shot?
  - Better

# Instruction Tuning

- Main purpose
  - Learn to do the task / Learn the output format

- Data Quality Requirement
  - Usually medium

- Fine-tuning compared to few-shot?
  - Better

# SFT

- Main purpose
  - Often inducing the high-quality pretraining data
    - (The mechanisms are not fully clear yet)

- Data Quality Requirement
  - High
- Fine-tuning compared to few-shot?
  - Could be worse
    - (contradict to the traditional ML common sense)

# Midterm Example Question

- We have an LLM base model (only pretrained). Which of the following fine-tuning data might boost the LLM's helpfulness as a chatbot?

- (A) 200 expert-written summaries

- (B) 1k expert-written summaries

- (C) 10k stories from Reddit

- (D) 200 most upvoted posts in Stack Exchange

- (E) 100k News in the domains LLM are not familiar with

- (D)

# Midterm Example Question

- We have an LLM base model (only pretrained). Which of the following fine-tuning data might boost the LLM's few-shot performance in the corresponding task the most?

- (A) 200 expert-written summaries for summarization

- (B) 1k expert-written summaries for summarization

- (C) 10k stories from Reddit for story generation

- (D) 200 most upvoted posts in StackOverflow for QA in code domain

- (E) 100k News in the domains LLM are not familiar with for news generation

- (D) > (B) > (A) > (E) >? 0 > (C)

# Limitations of SFT

- Too expensive

  - Your quality needs to be close to the best response on the Internet

  - Hiring experts is expensive

- Fine-tuning on unfamiliar materials could cause hallucination

- Could easily affect the different tasks

- Do not have negative examples

  - LLM only knows what it should say. It does not know what it should not say

  - Hard to prevent generating harmful/toxic responses

# LLM Development

Internet low-quality text (e.g., from trolls or haters)



Internet high-quality text

Post-training stage
(Filtering process)

- Architectures
  - MLP
  - RNN
  - Transformer

- Training Stages
  - Pretraining
  - Supervised Fine-tuning (SFT)
  - Alignment
    - **Learning from Human Feedback (LHF)**
  - Reasoning

# LLM Evaluation



"You insist that there is something a machine cannot do. If you tell me precisely what it is a machine cannot do, then I can always make a machine which will do just that."

- John von Neumann, 1948

https://www.reddit.com/r/singularity/comments/18t02br/john_von_neumann_was_the_first_who_used_the/

# Evaluation is Hard, Why?

Human as the Judge

LLM → Response 1

Score 4.5

Score 3.6

Score 1.9

Score 4.3

Noisy.

Every person has their preference.

Scores are easily affected by the context

# Evaluation -> Loss

Human as the Judge /
LLM as the Judge

**LLM**

Response 1 ✅  We Should Encourage the LLM to output more of this

Response 2 ❌

# Differentiable Review

- Gradient descent is the easiest and most stable way

- If you change a little, the output cannot change a little.

  - That is not differentiable

# Last Year Note

Typos



$$\max_{\pi_\theta} \mathbb{E}_{x\sim\mathcal{D},y\sim\pi_\theta(y|x)}\big[r_\phi(x,y)\big] - \beta\mathbb{D}_{\mathrm{KL}}\big[\pi_\theta(y\mid x)\,||\,\pi_{\mathrm{ref}}(y\mid x)\big]$$

# Helpfulness vs Harmlessness



Normal Prompt

Adverserial Prompt

Helpfulness Scores

Harmlessness Scores (52B)

**Professional Writers**
**Context Distilled**
**Static HH RLHF**
**Online HH RLHF (52B)**
**Online Helpful RLHF (52B)**

Elo Scores

Parameters

Crowdworker Preference Frequency

**Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback (https://arxiv.org/abs/2204.05862)**

# RLAIF



Constitutional AI: Harmlessness from AI Feedback (https://arxiv.org/abs/2212.08073)

# Why can LLM be a Judge?

- Some websites have higher quality

- Different quality of text inducing different kinds of responses



My guessing. Not know a paper

# Best of N

Reward Model

LLM → Response 1    0.6    ❌

Response 2    0.9    ✔️

......

Response N    0.3    ❌

# SFT Rejection Sampling / RAFT



Reward Model

SFT Finetuning

LLM

Response 1     0.6     ❌

Response 2     0.9     ✅

……

Response N     0.3     ❌

Llama 2: Open Foundation and Fine-Tuned Chat Models (https://arxiv.org/pdf/2307.09288)

# Iterative RLHF is Better, why?



RLHF Workflow: From Reward Modeling to Online RLHF (https://arxiv.org/abs/2405.07863)

# Multiple Rounds



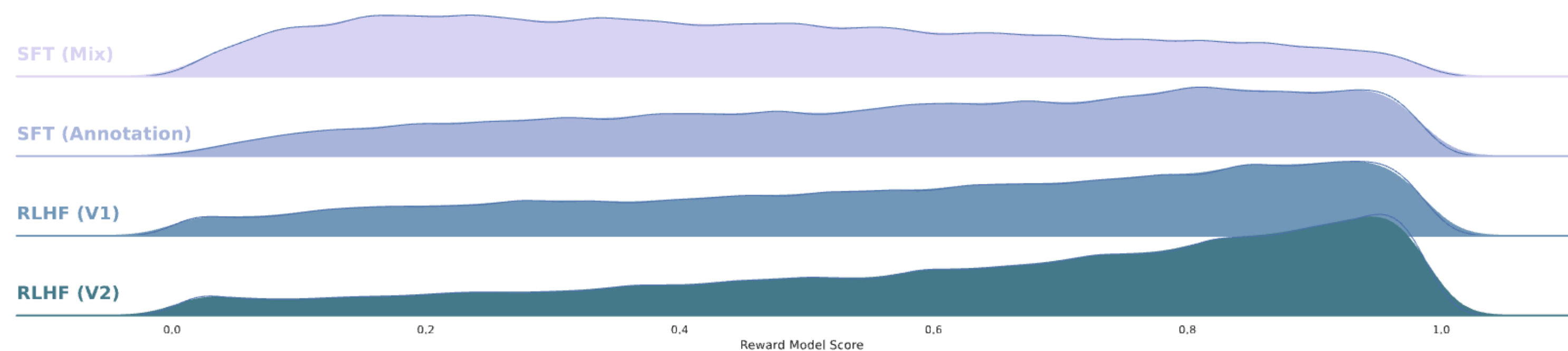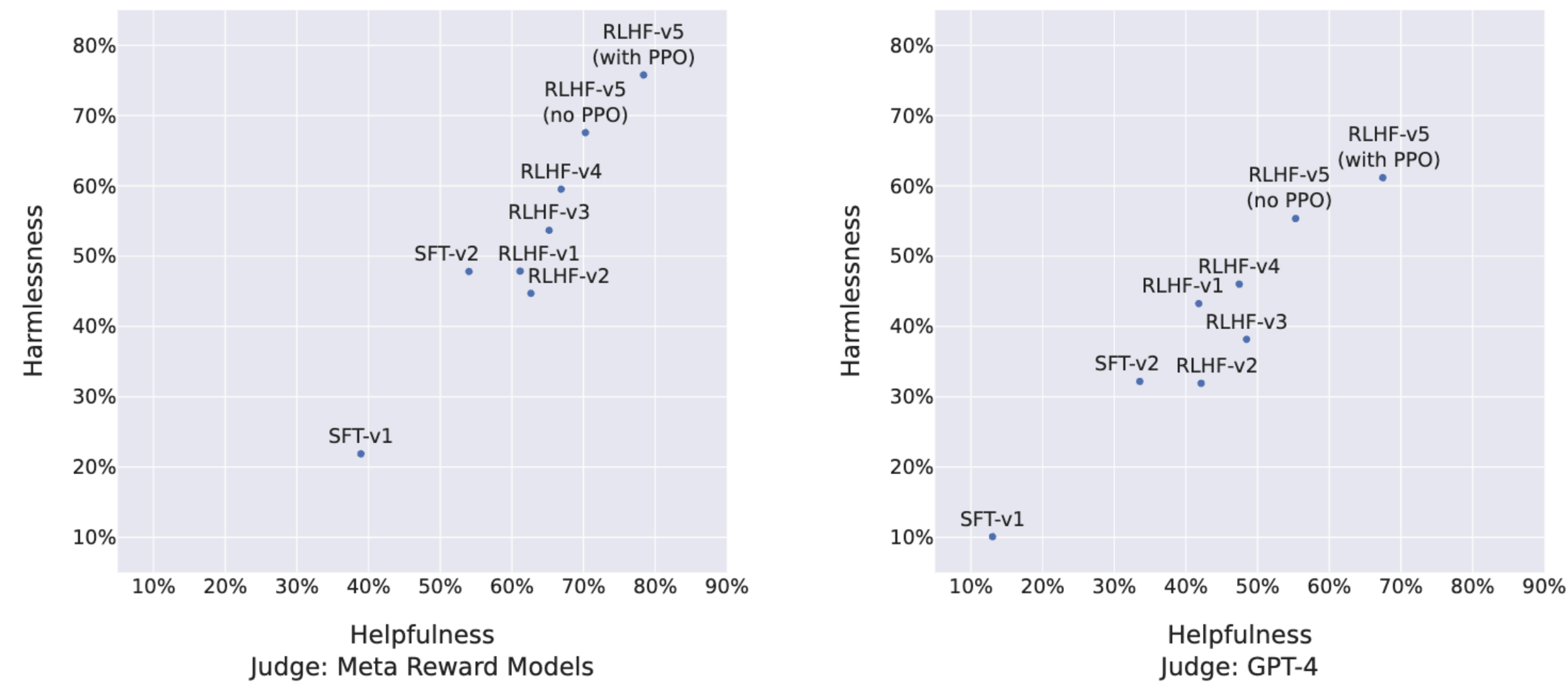**Figure 20: Distribution shift for progressive versions of LLAMA 2-CHAT,** from SFT models towards RLHF.

Llama 2: Open Foundation and Fine-Tuned Chat Models (https://arxiv.org/pdf/2307.09288)