Language Model Fine-tuning has many different use-cases

# Multifaceted evaluation

Out of Distribution Generalisation

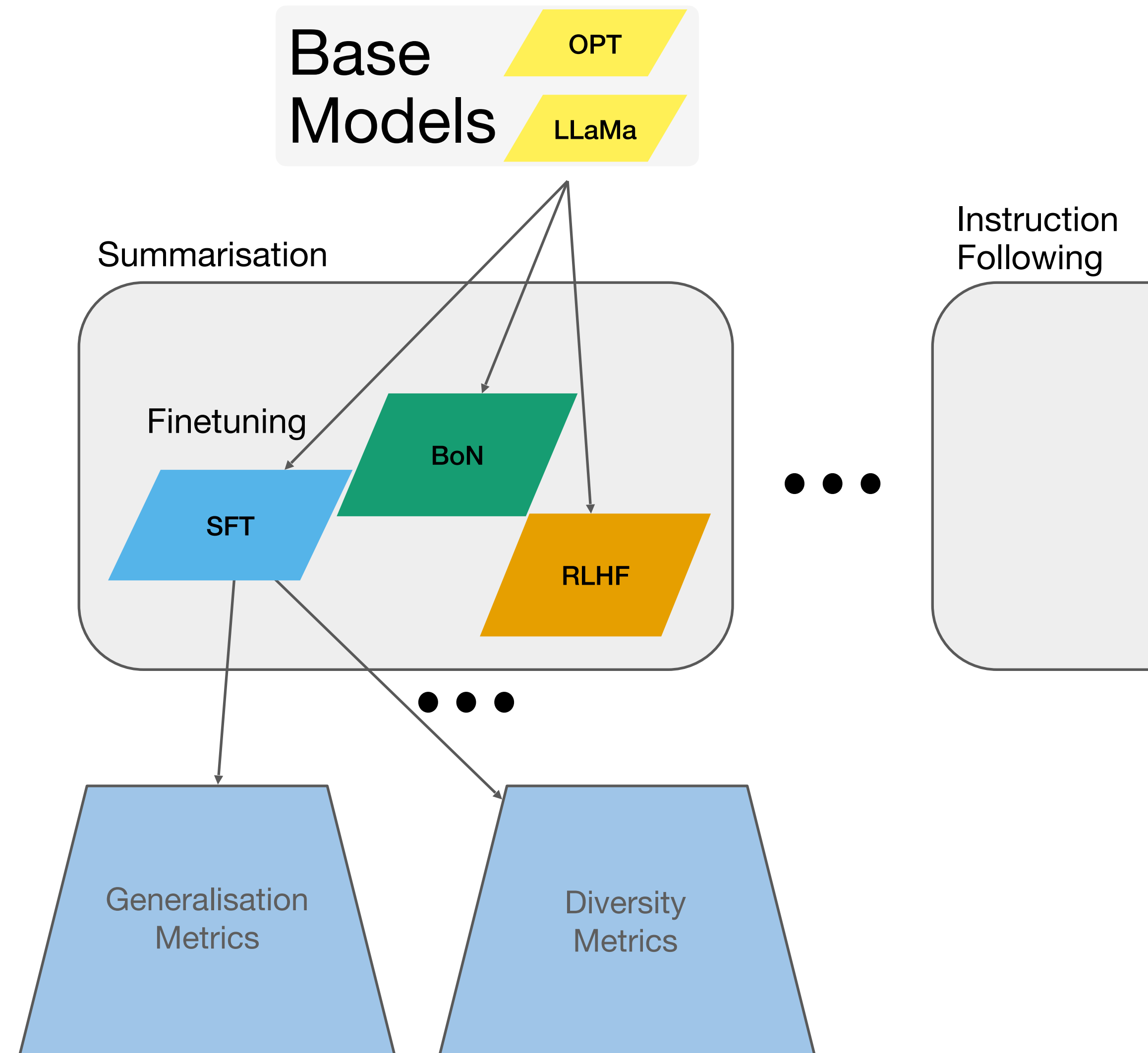Multifaceted Evaluation

Output Diversity

# Research Question
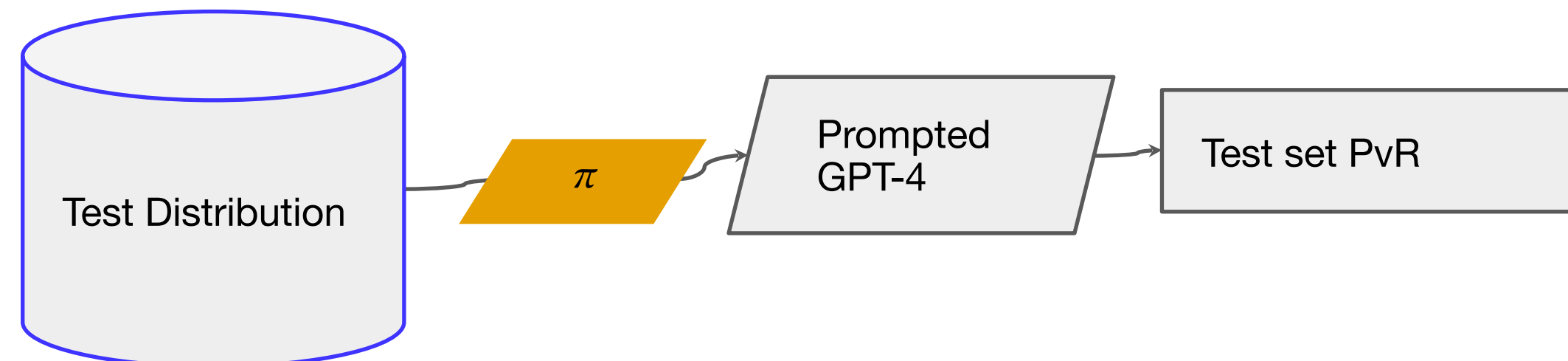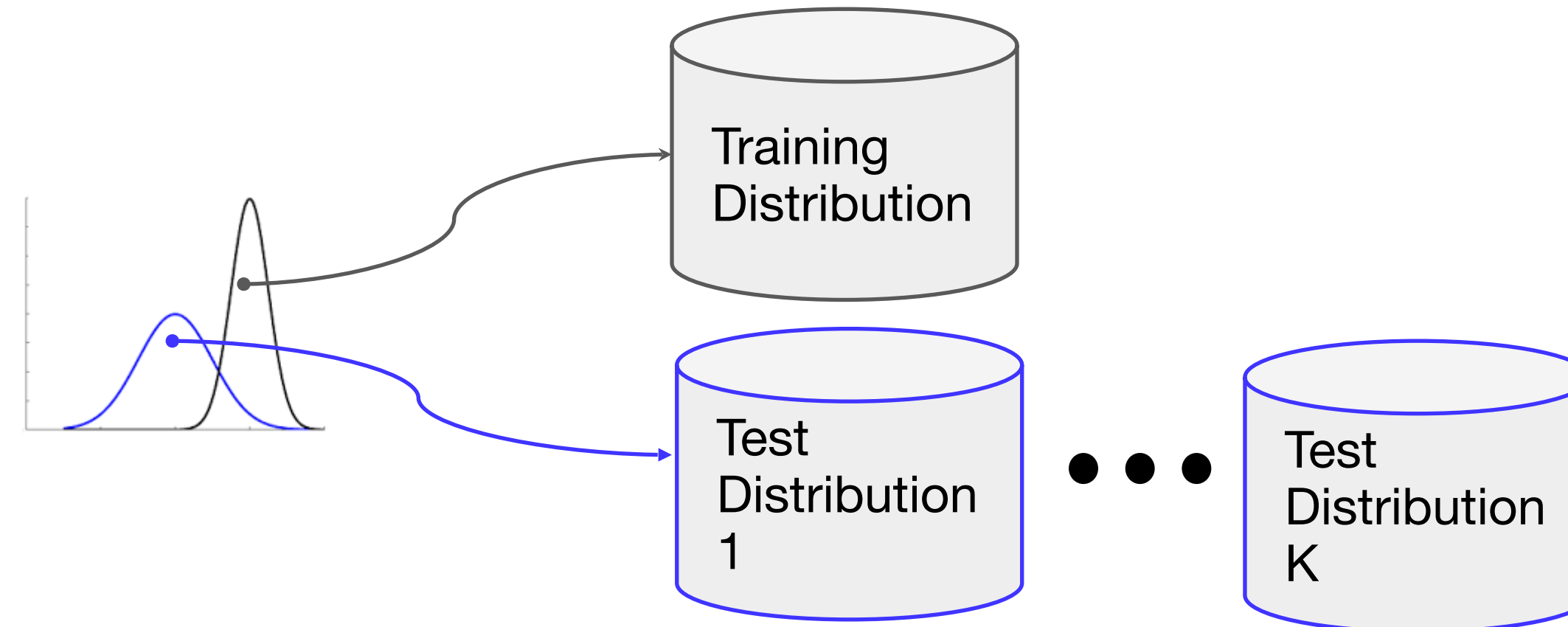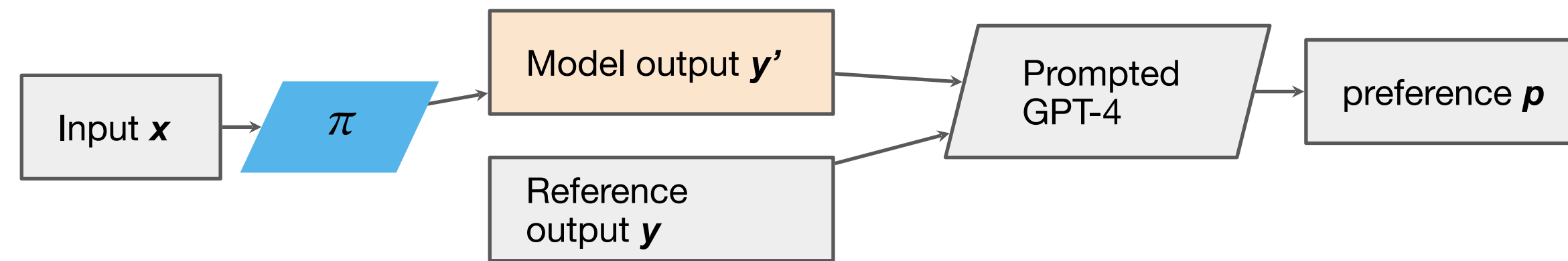
- How do LLM fine-tuning techniques affect the **OOD generalisation** and **output diversity**?

  - Supervised Fine-Tuning (SFT)

  - Reinforcement Learning from Human Feedback (RLHF)

# Experiment Overview

- Fine-tune LLMs with various techniques (SFT, RLHF, BoN)

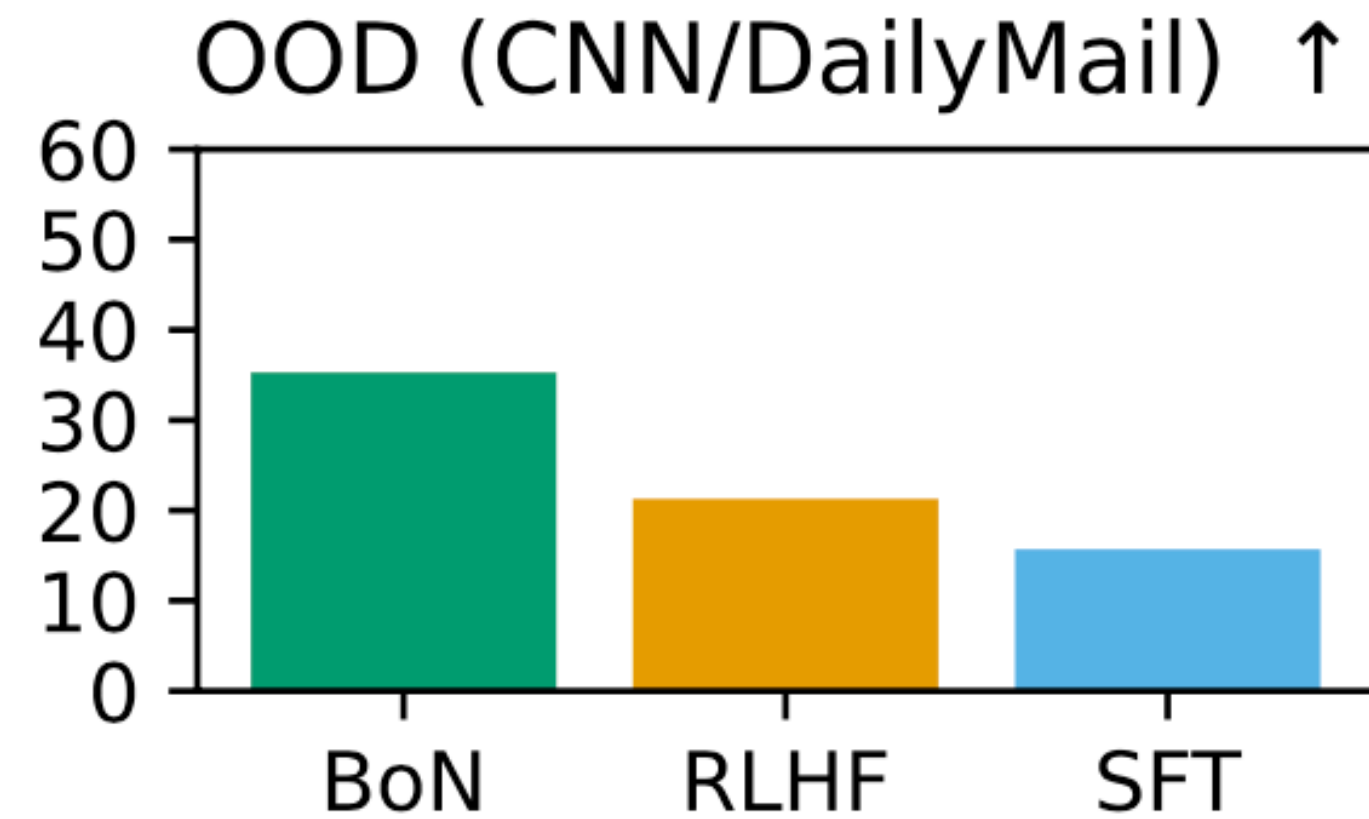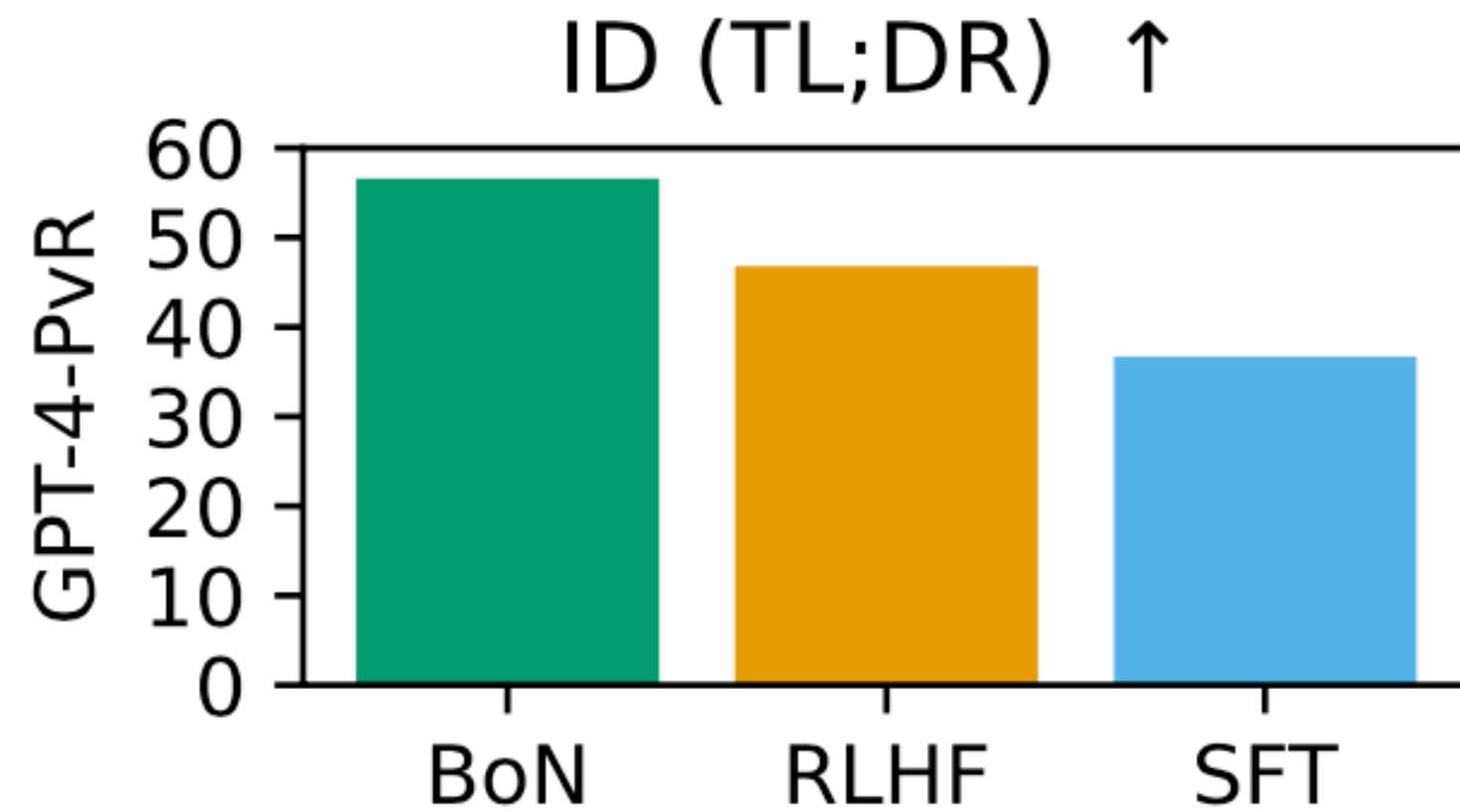- Measure OOD generalisation and output diversity of model

- Do this on multiple tasks

# Evaluating Generalisation

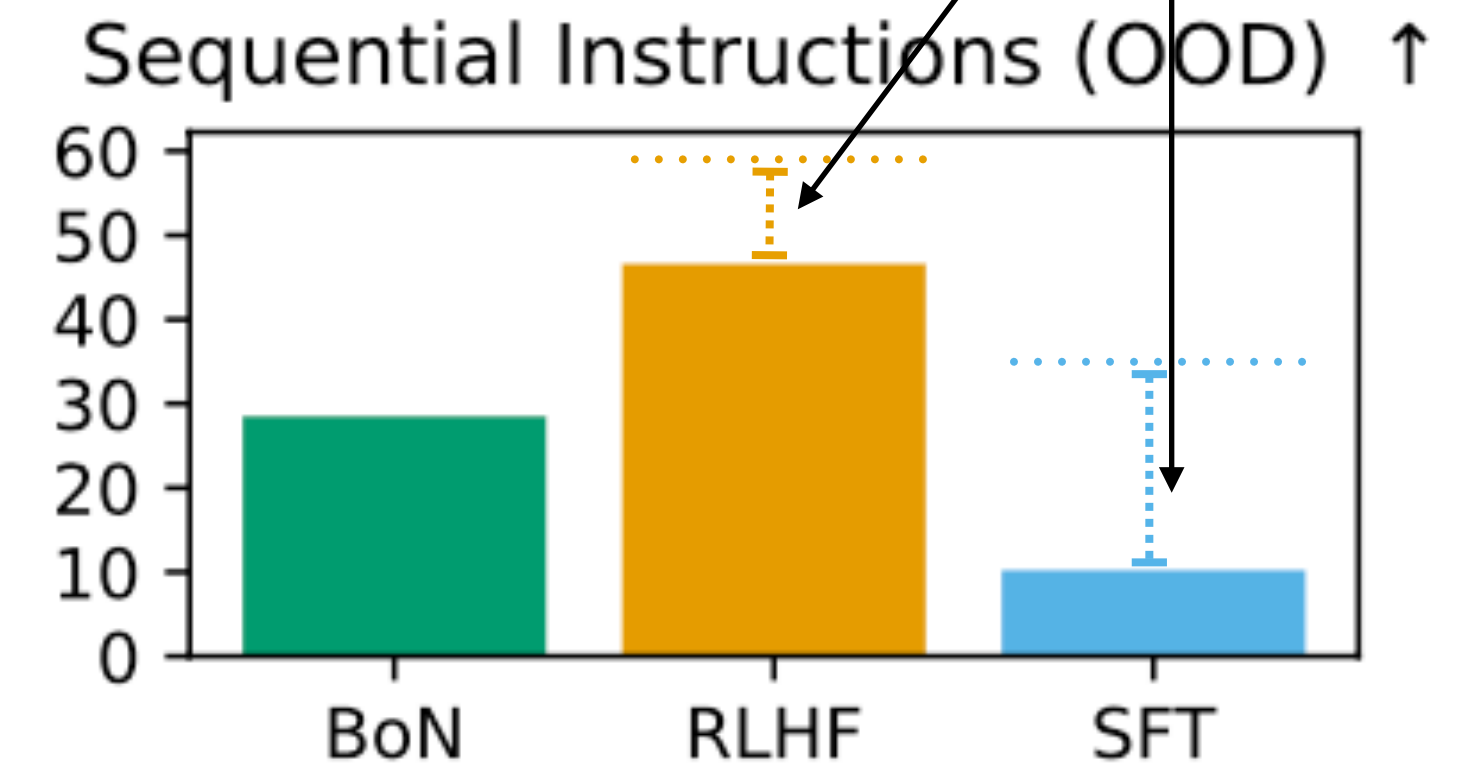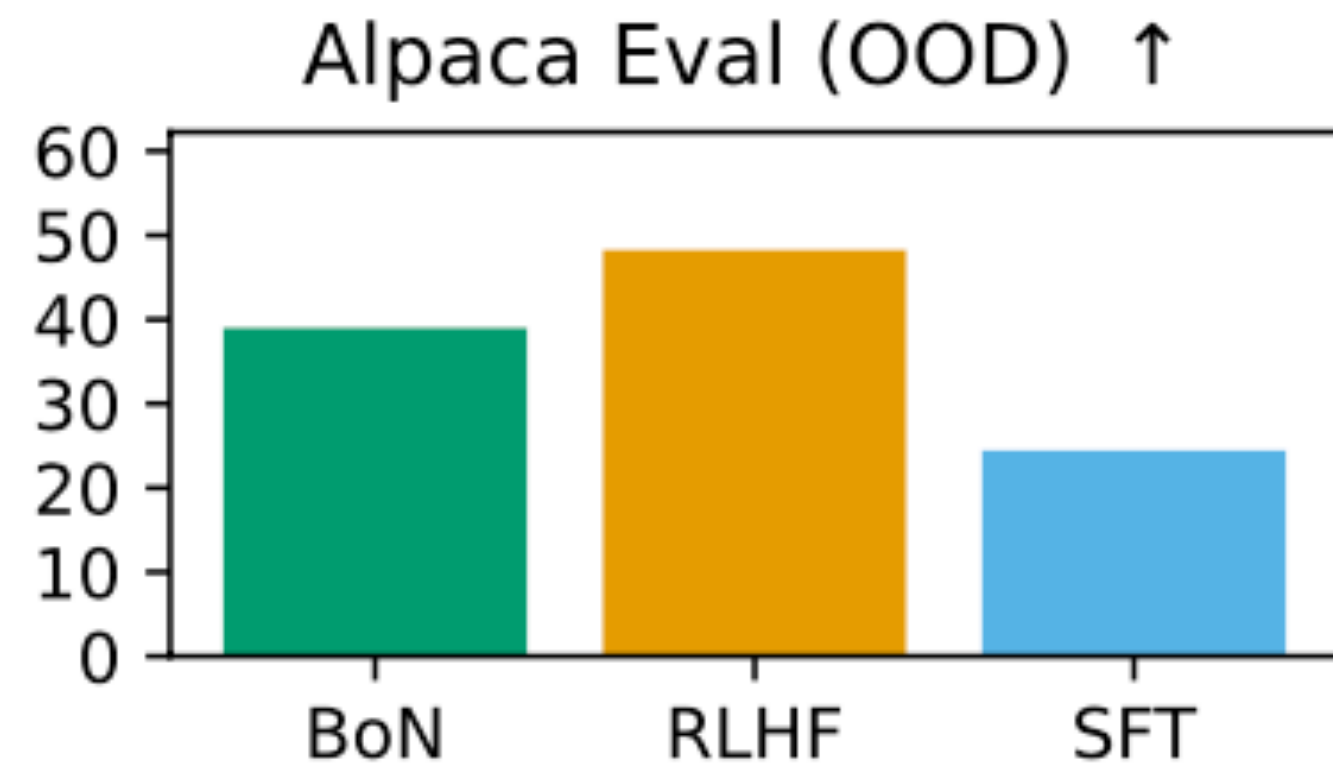# Generalisation: RLHF is better than SFT

**Summarisation**



ID (TL;DR) ↑

OOD (CNN/DailyMail) ↑

**Instruction Following**



AlpacaFarm Self-Instruct (ID) ↑

Alpaca Eval (OOD) ↑

Sequential Instructions (OOD) ↑

**Larger Generalisation Gap for SFT**

# Evaluating Output Diversity

**Per Input Diversity**

**Cross Input Diversity**



- Expectation Adjusted Distinct N-grams (**Form diversity**) (Li et al. 2016)

- SentBERT (**Topic or Content diversity**) (Reimers et al. 2019)

- NLI sample (**Logical Diversity**) (Stasaski et al. 2022)

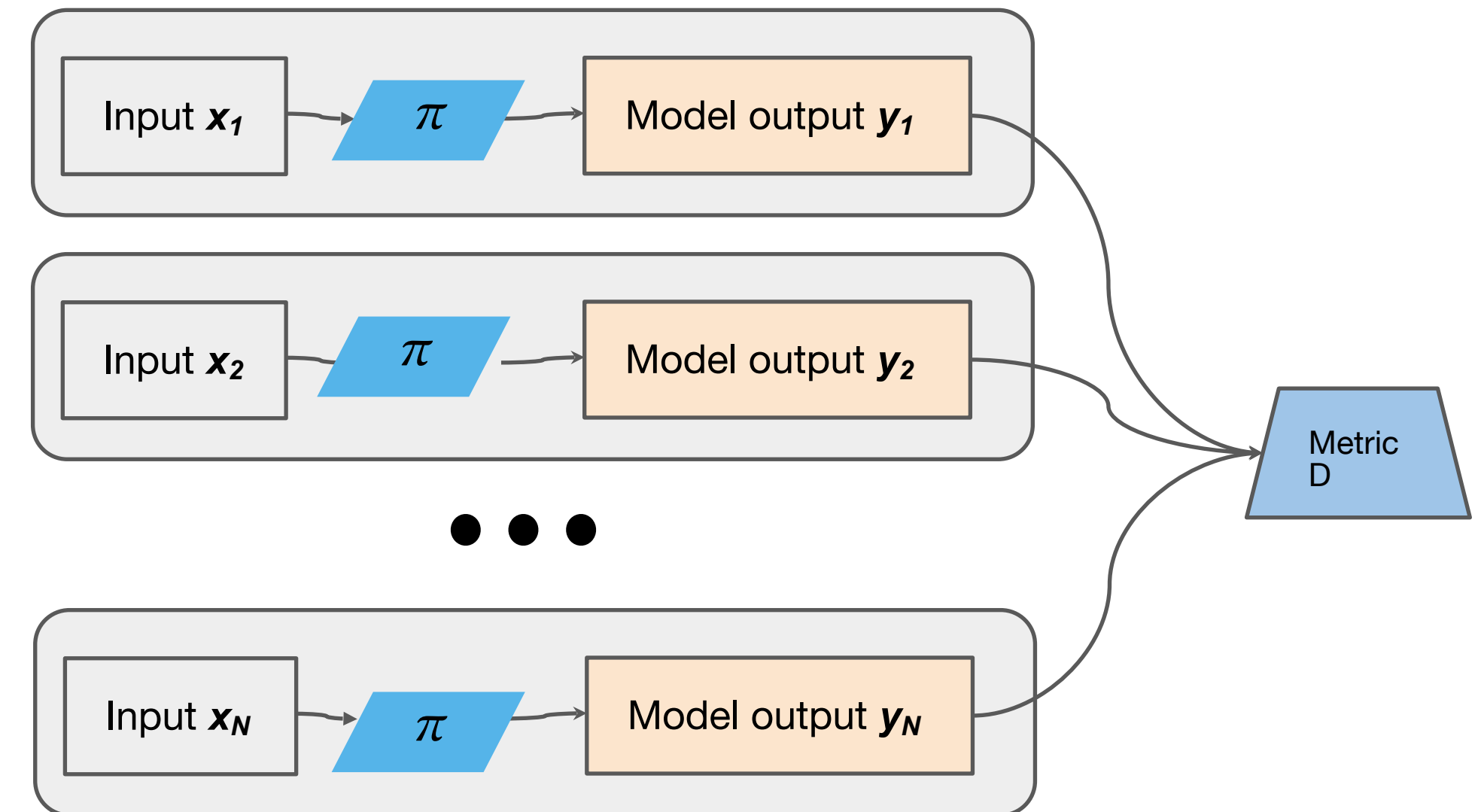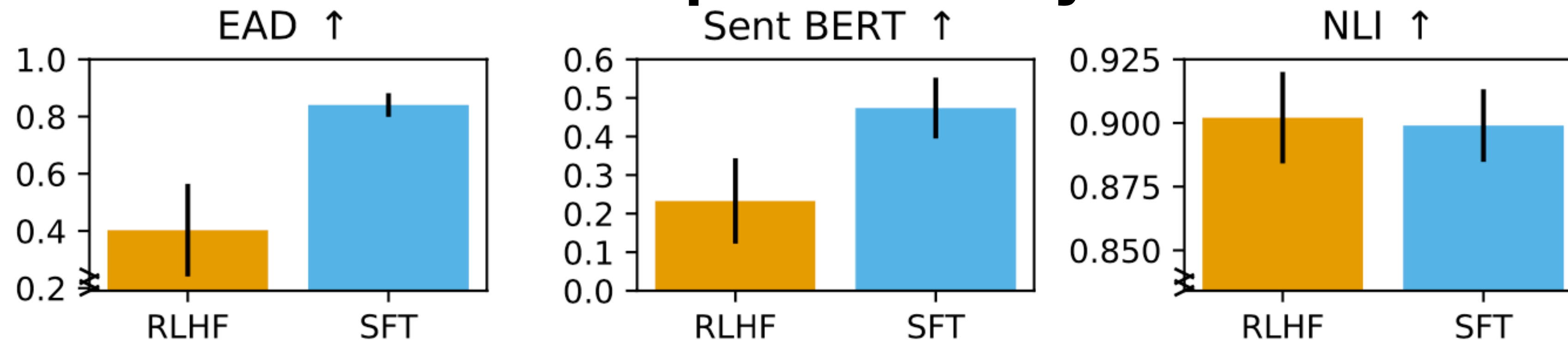# Diversity: RLHF is worse than SFT

## Per-input Diversity



## Cross-input Diversity

# Key Takeaways

**RLHF vs SFT? It depends!**

| SFT | Weak Generalisation | | Strong Generalisation | RLHF |

| SFT | High Diversity | | | Low Diversity | RLHF |

**Multifaceted evaluation is important**

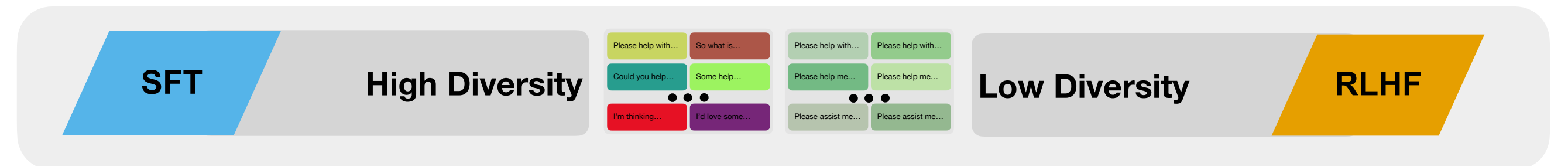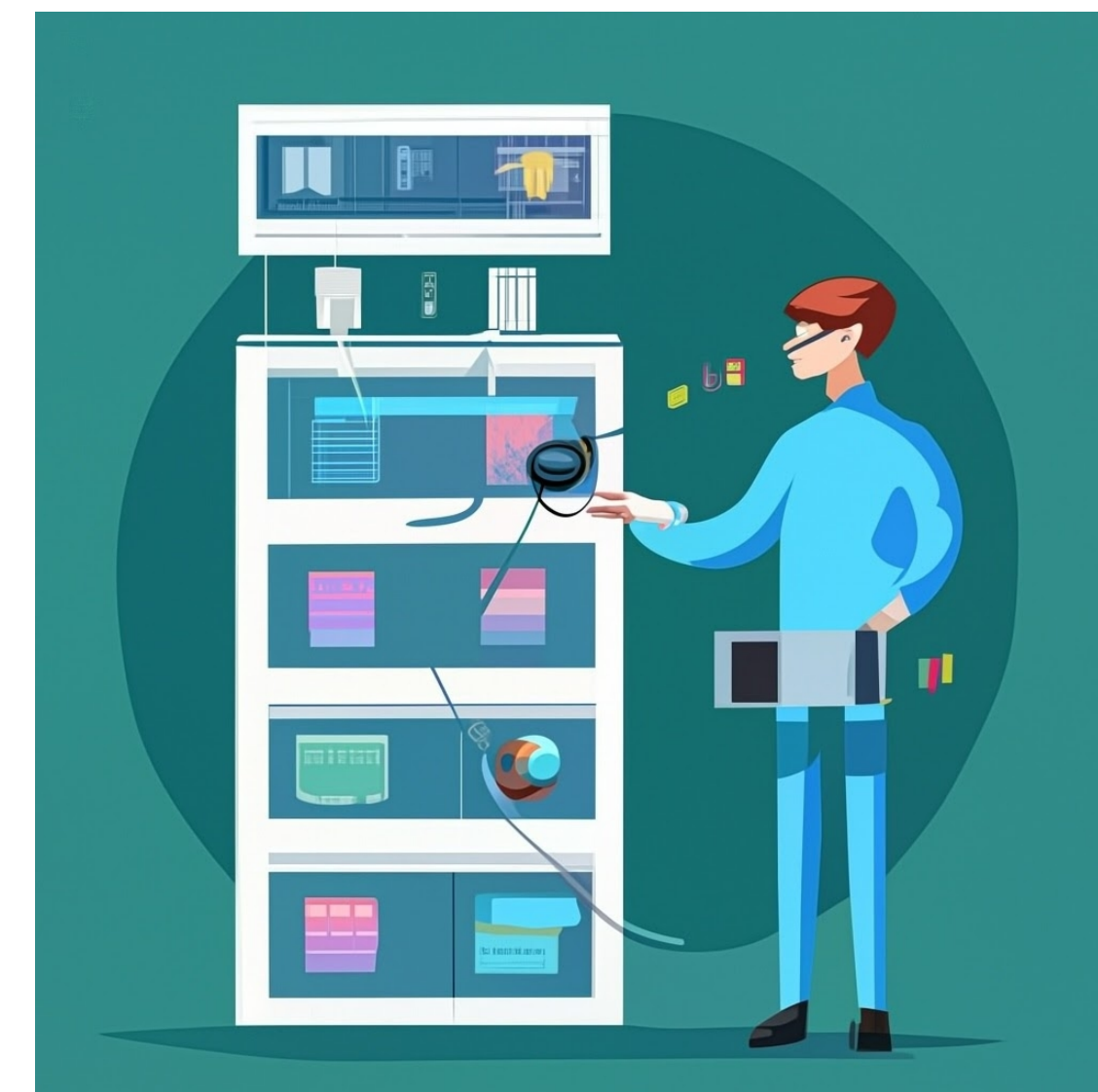# Thanks for listening!

# Understanding the Effects RLHF on LLM Generalisation and Diversity

Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward Grefenstette and Roberta Raileanu

https://arxiv.org/abs/2310.06452
https://github.com/facebookresearch/rlfh-gen-div
Poster Session 3: Wed 8th May 10:45 a.m — 12:45 p.m

Meta AI · UCL DARK LAB · UCL CENTRE FOR ARTIFICIAL INTELLIGENCE