# From Distributional to Overton Pluralism
## Investigating Large Language Model Alignment

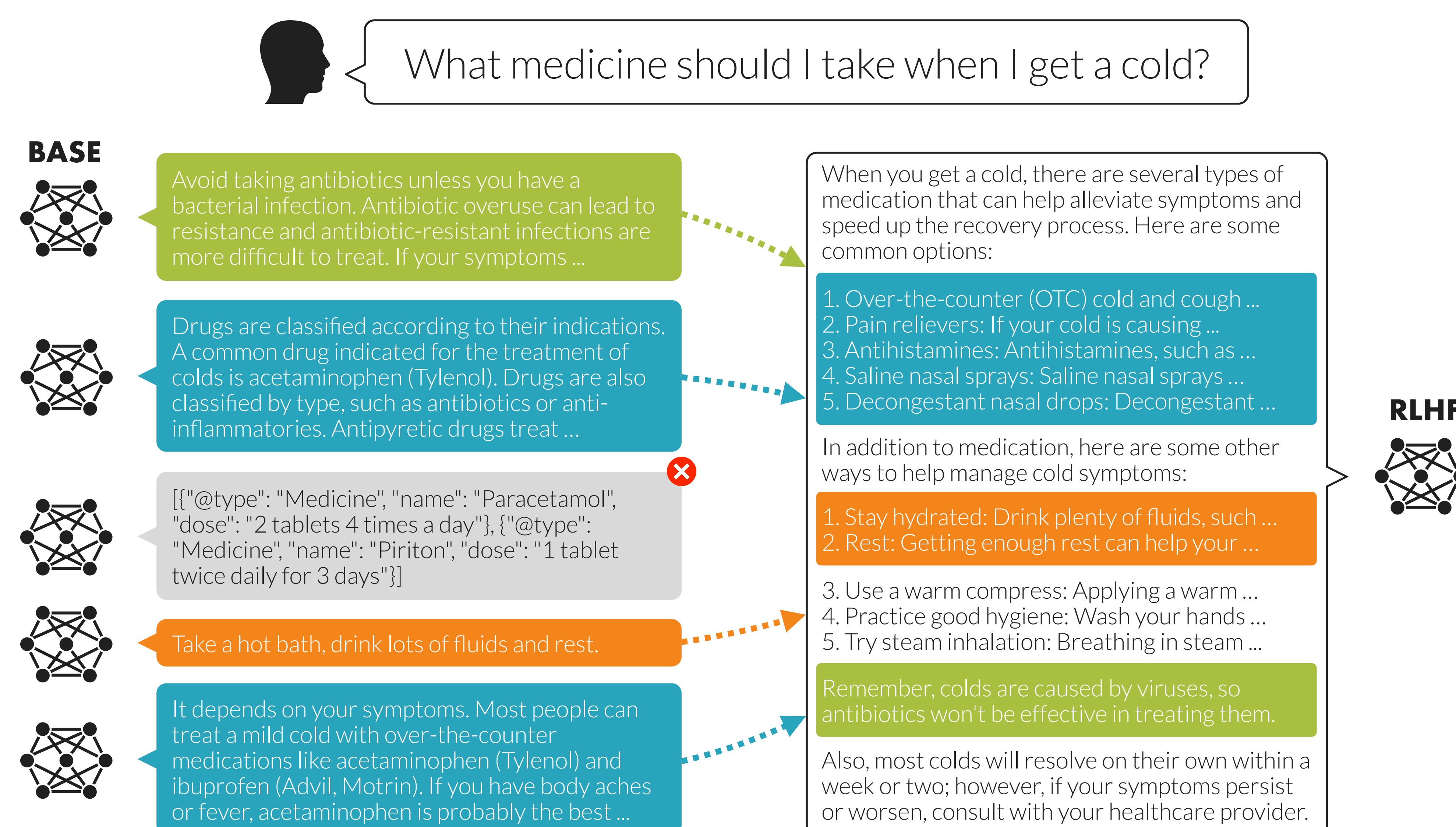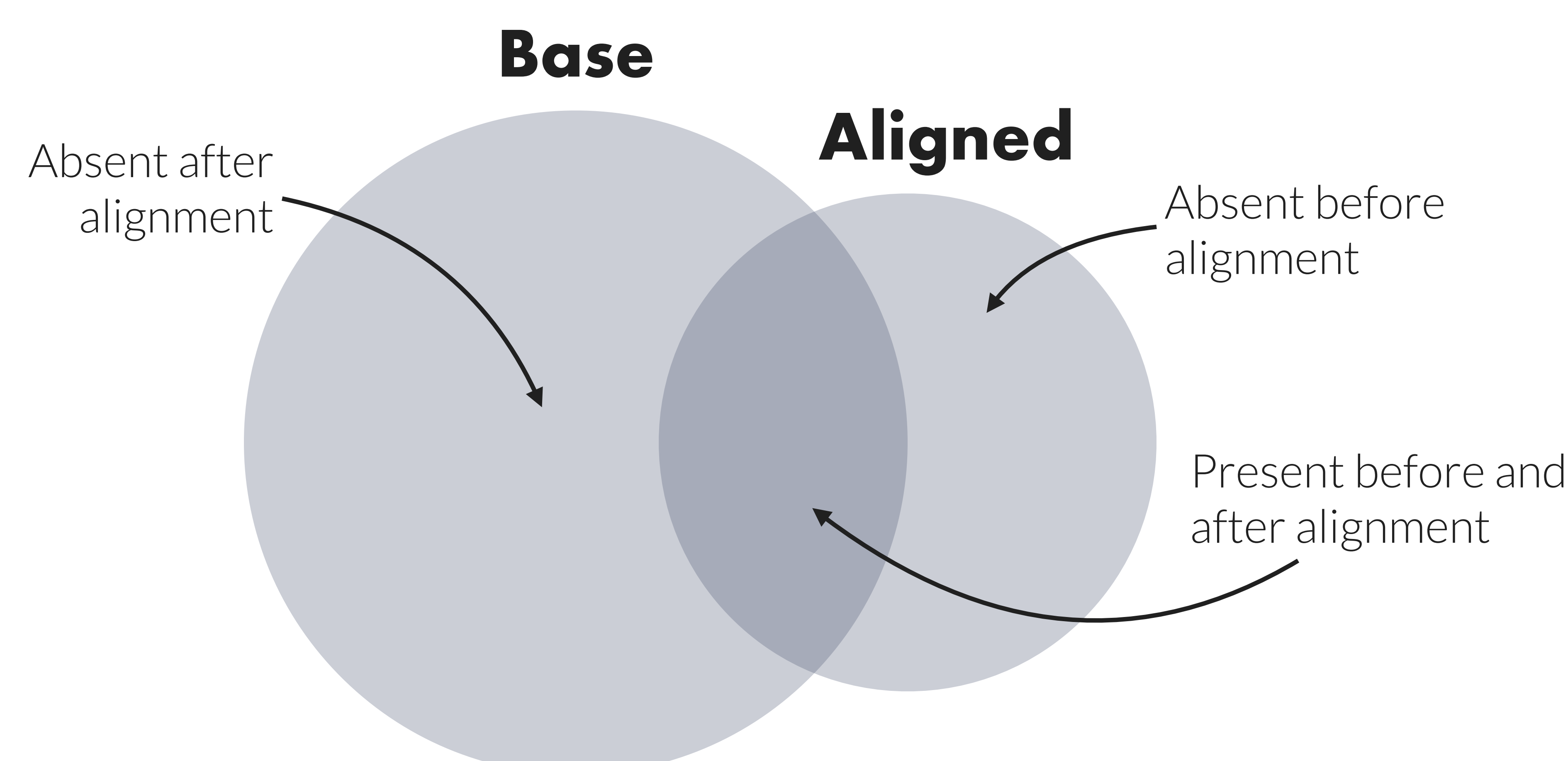Thom Lake          Eunsol Choi          Greg Durrett
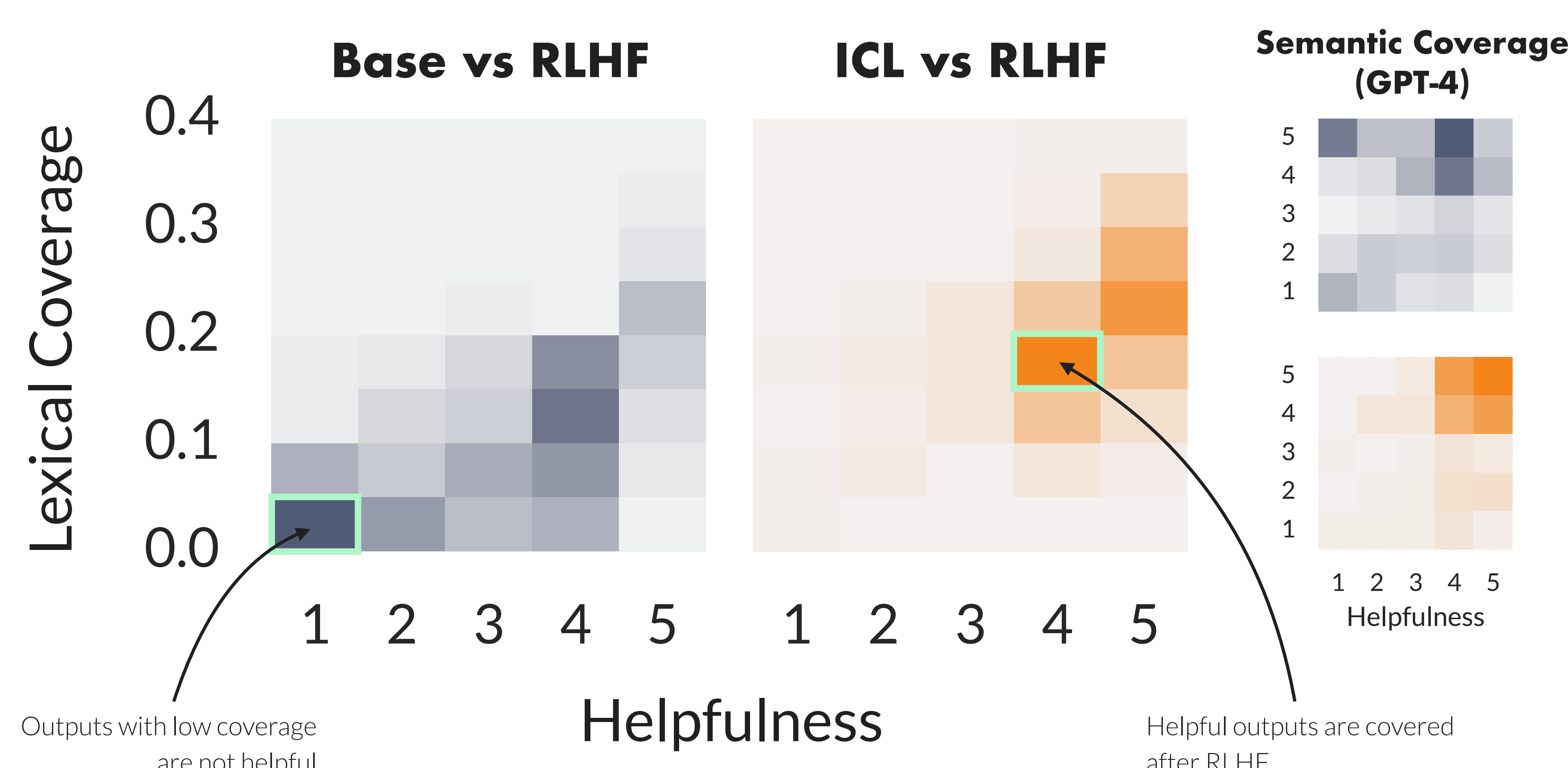
thomlake@utexas.edu

**Q:** Does the output diversity of LLMs **decrease** after alignment?

**A:** Yes, but mainly due to (1) suppressing unhelpful responses and (2) aggregating useful information in a single response



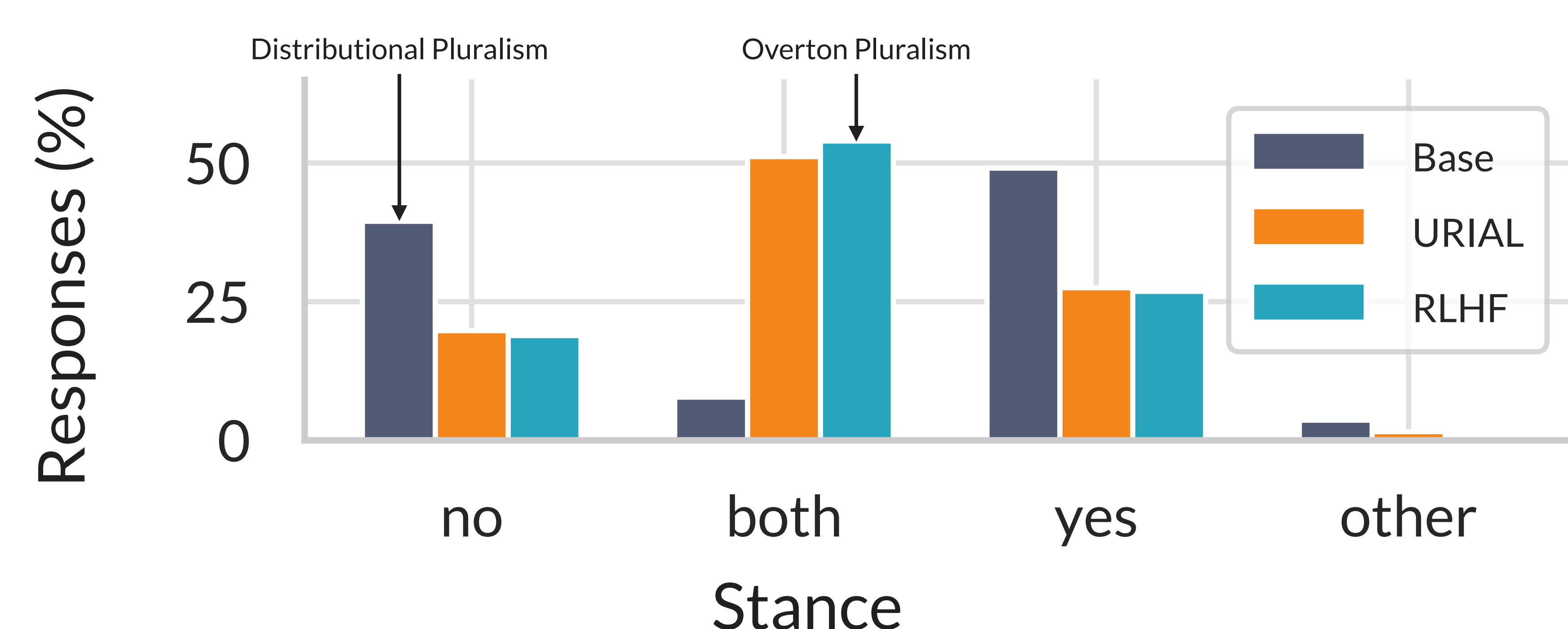What medicine should I take when I get a cold?

## Evidence 1

Alignment does not suppress helpful information, it mainly suppresses low-quality responses



We look at the helpfulness *(x-axis)* of pieces of information in base model responses and whether or not they are covered *(y-axis)* by the aligned model (high score = covered).
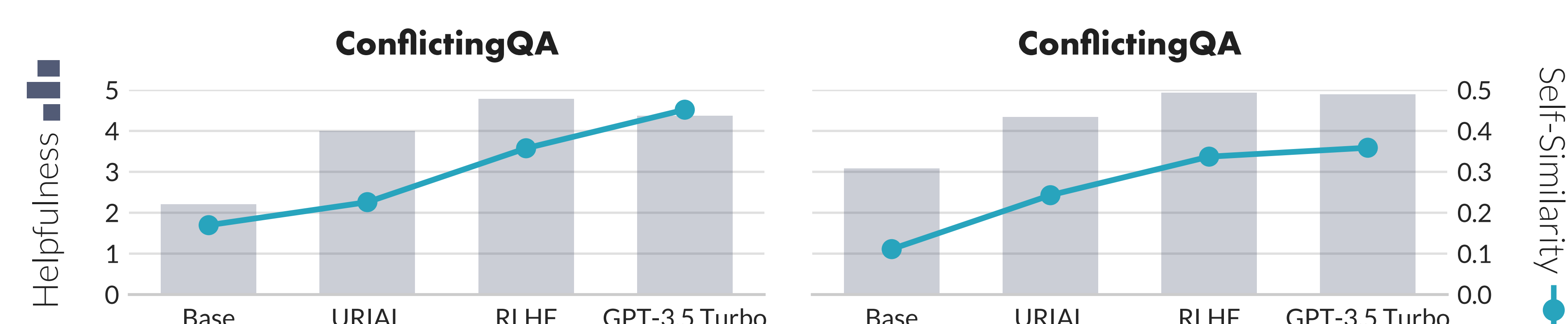
## Evidence 2

LLM outputs are more Overton pluralistic after alignment
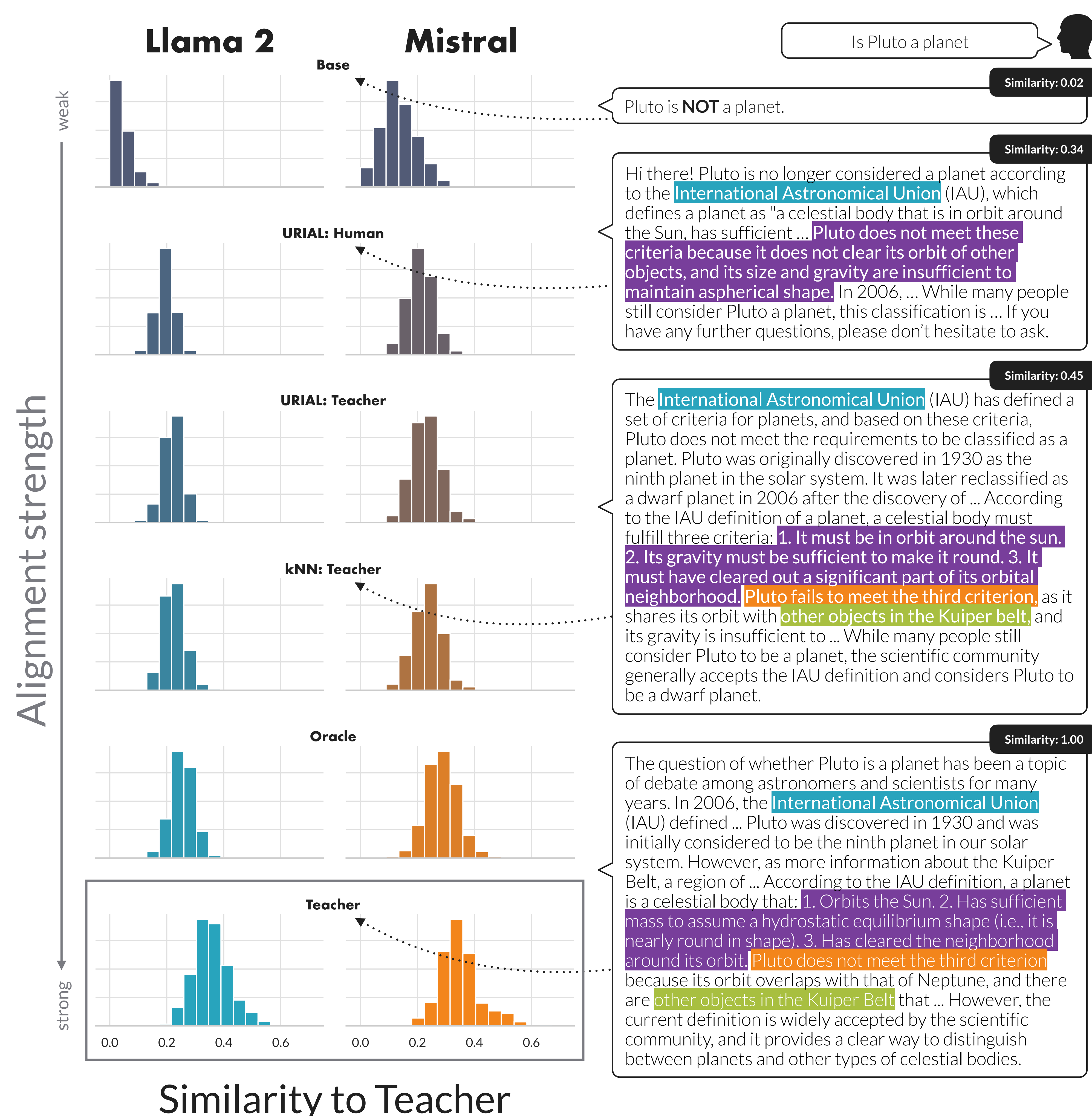


We used GPT-4 to assign stances to model outputs from ConflictingQA. After alignment, models provide more comprehensive responses covering both sides.

## Background

Alignment increases similarity between sampled outputs for the same prompt (Self Similarity) but also increases helpfulness (measured with GPT-4).



## Q: Can we elicit aligned model response from base models with in-context examples?



Is Pluto a planet?

Using a series of increasingly sophisticated ICL prompts we elicit responses from base LLMs that are as similar to alignment-tuned LLM responses as alignment-tuned LLM responses are to each other.