# 1 What is this

This is the final iteration of the training where we are using PPO to train our already supervised-trained LLM.