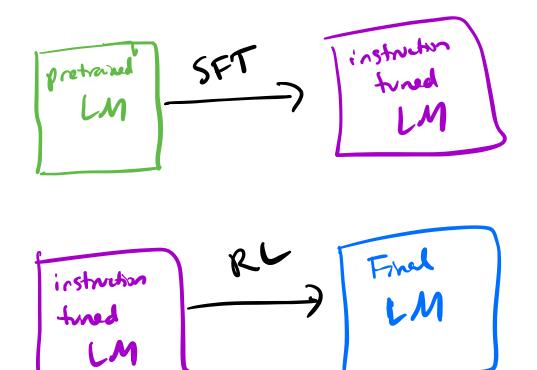
SFT to RL:



Ly instruction tuning requires a diserse dataset of instruction linear lowers types

$$L = \frac{1}{N} \sum_{i=1}^{N} -\log P(Y_i) \times_i$$

by ground touth output y is generally created by a human

La expensive to create HQ instruction following datasets L) Synthetic data Lyone acceptable y per inpot x 1) no learning from negative examples La distribution mismatch. 6 training time: expert HQ Y Li test time: model-gon y D exposure bias G doesn't directly involve human prefs