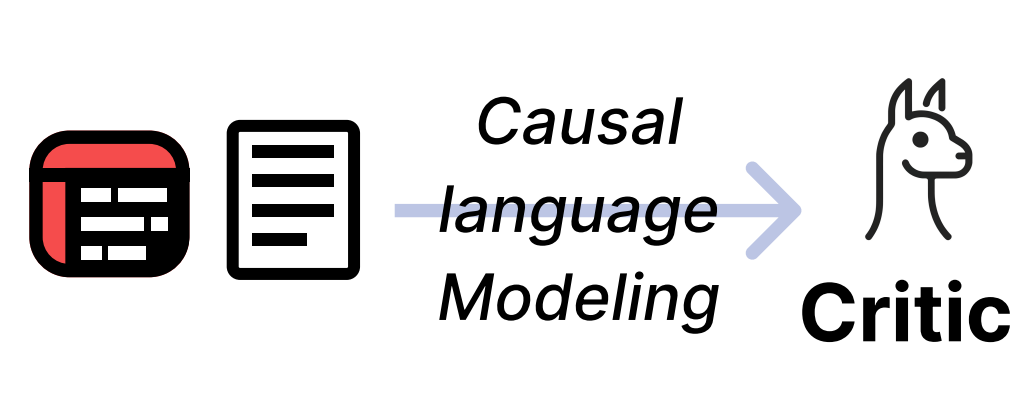
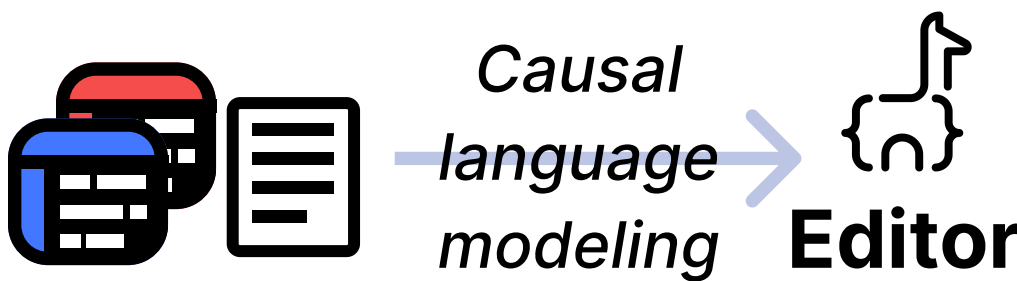


Phase I: Supervised Fine-tuning with COFFEE

Training Critic

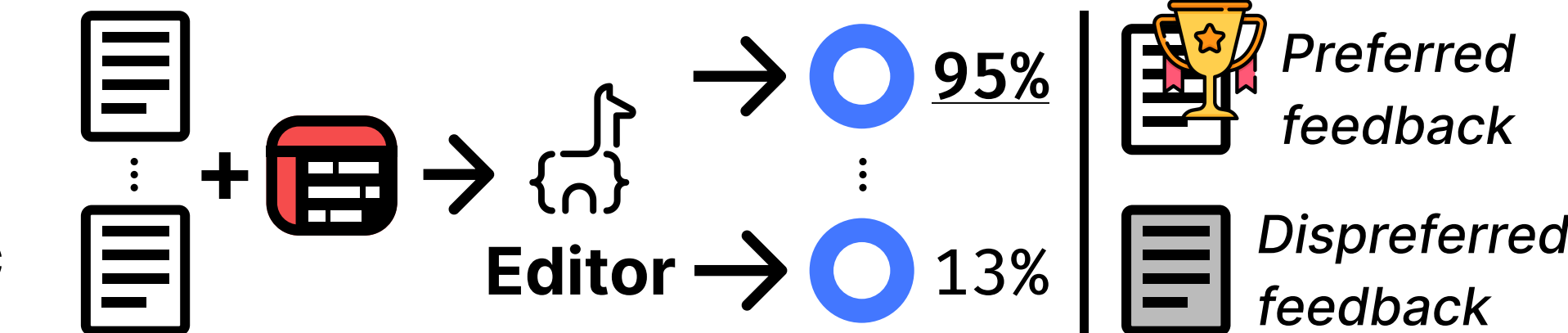


Training the Editor

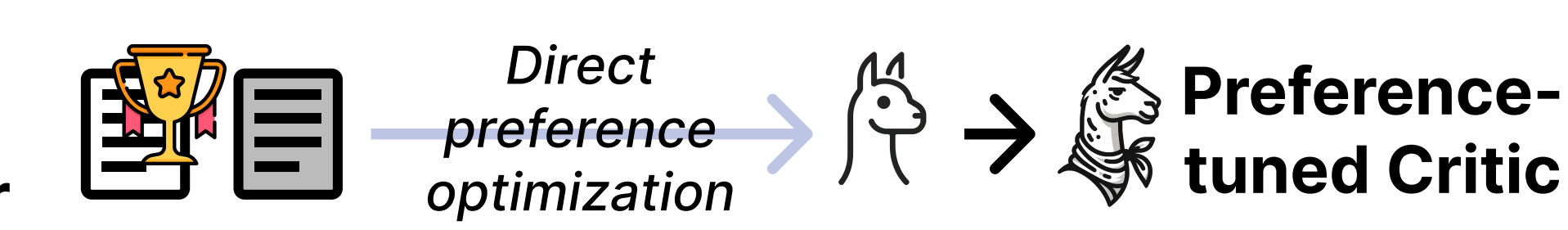


Phase II: Preference Tuning on Critic

Collecting Preference Pairs via Hidden Test Cases

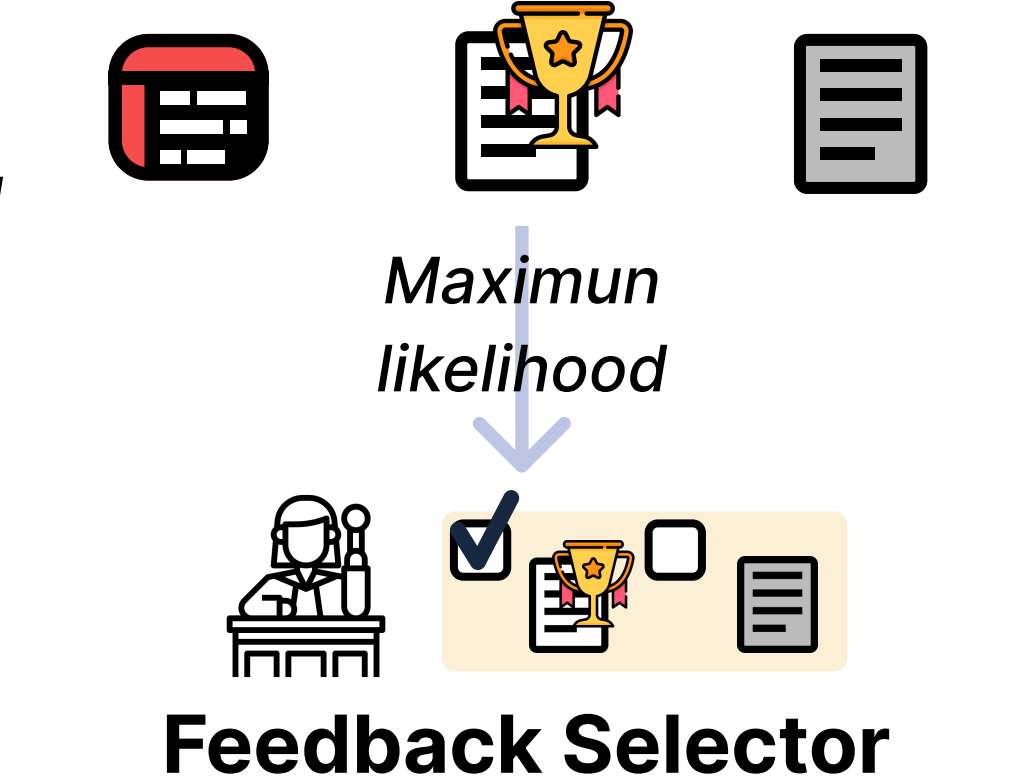


Preference-tuning the Critic Model



Phase III: Preference Selection via the Feedback Selector

Training the Feedback Selector



The Overall Workflow



Legends

