

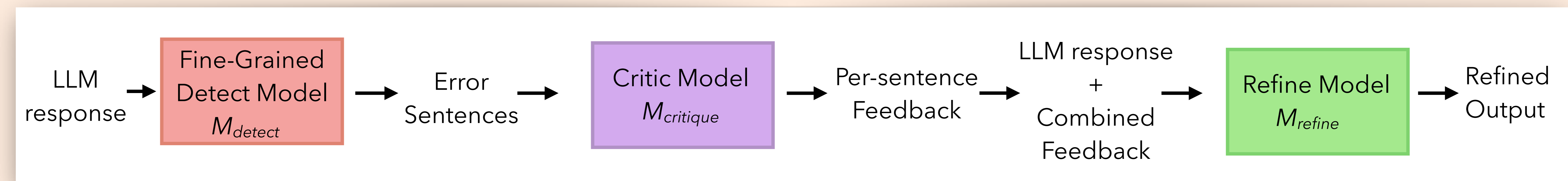
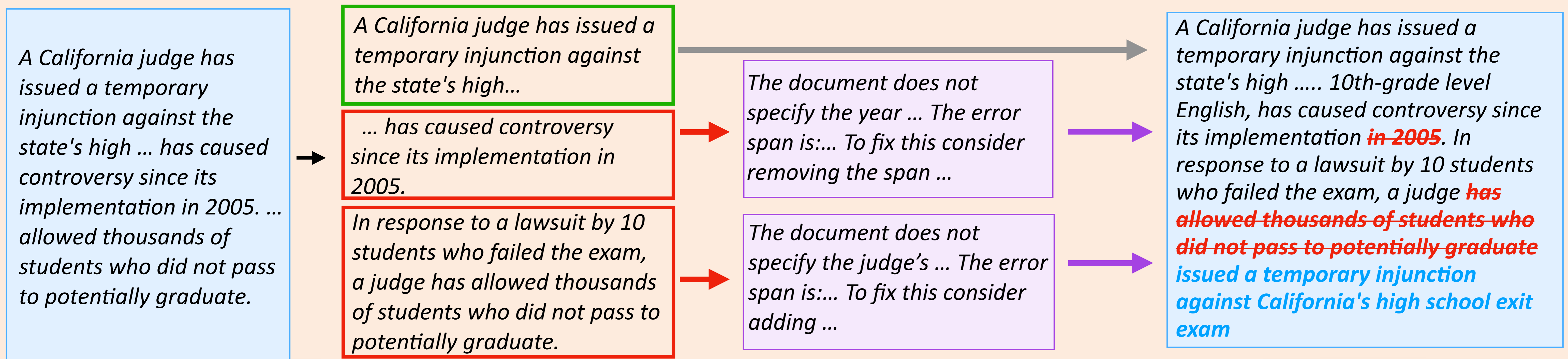
Learning to Refine with Fine-Grained Natural Language Feedback

Manya Wadhwa, Xinyu Zhao, Jessy Li, Greg Durrett
manya.wadhwa@utexas.edu



Detect, Critique, and Refine (DCR) enables better post-hoc refinement

Input (LLM response) **Step 1: Detect** M_{detect} marks incorrect sentences **Step 2: Critique** $M_{critique}$ generates feedback per error sentence **Step 3: Refine** M_{refine} edits the input using combined per-sentence feedback



Experiments

Task: Improve factual consistency of document grounded summaries

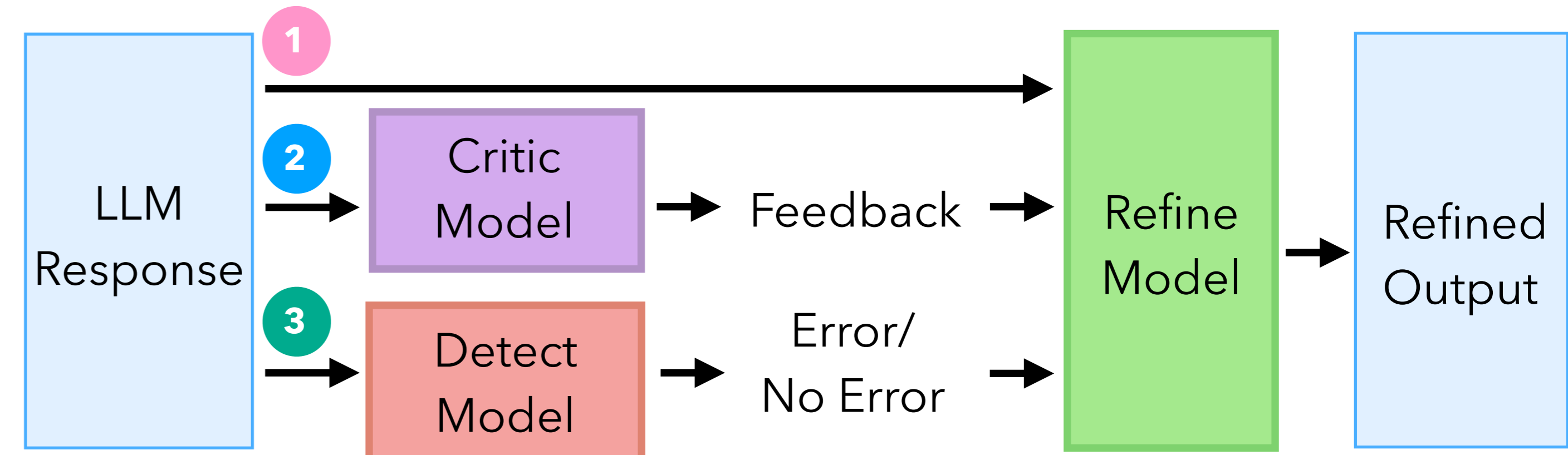
Datasets: UltraChat (Ding et al, 2023), TofuEval (Tang et al, 2024)

Evaluation Metrics: Δ AlignScore, GPT-4 Win Rate on Factuality (more metrics in the paper!)

Detect Model	MiniCheck (Tang et al, 2024)
Critic Model	Llama-2-7b-Chat (FT) Llama-3-8b-Instruct (FT) GPT-4 (prompted)
Refine Model	Llama-2-7b-Chat (FT) Llama-3-8b-Instruct (FT) GPT-4 (prompted)

Baseline Refinement Methods

- 1 Direct Refine
- 2 Feed + Refine
- 3 Detect + Refine

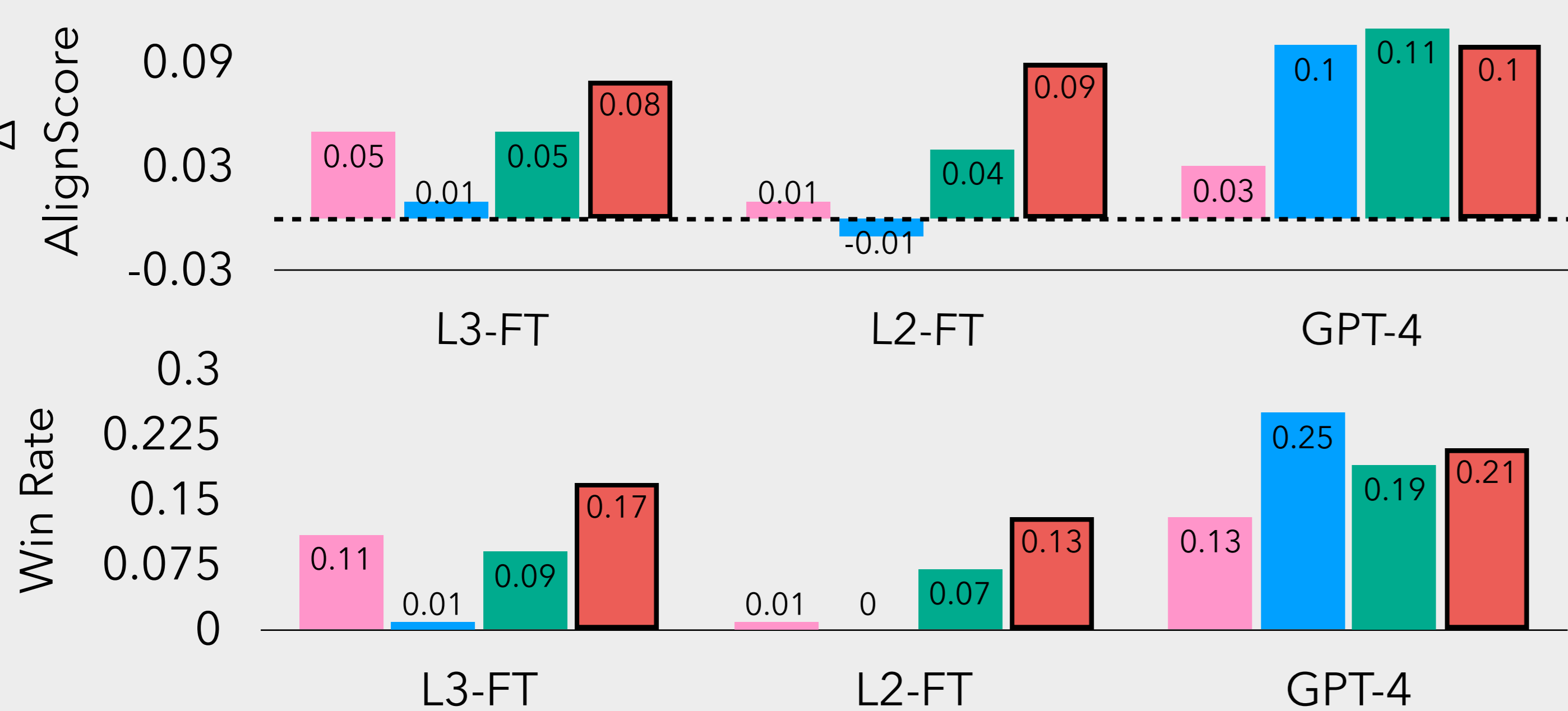


Results

Result 1: DCR refinement outperforms standard refinement strategies

- Direct Refine
- Feed+Refine
- Detect + Refine
- DCR (proposed)

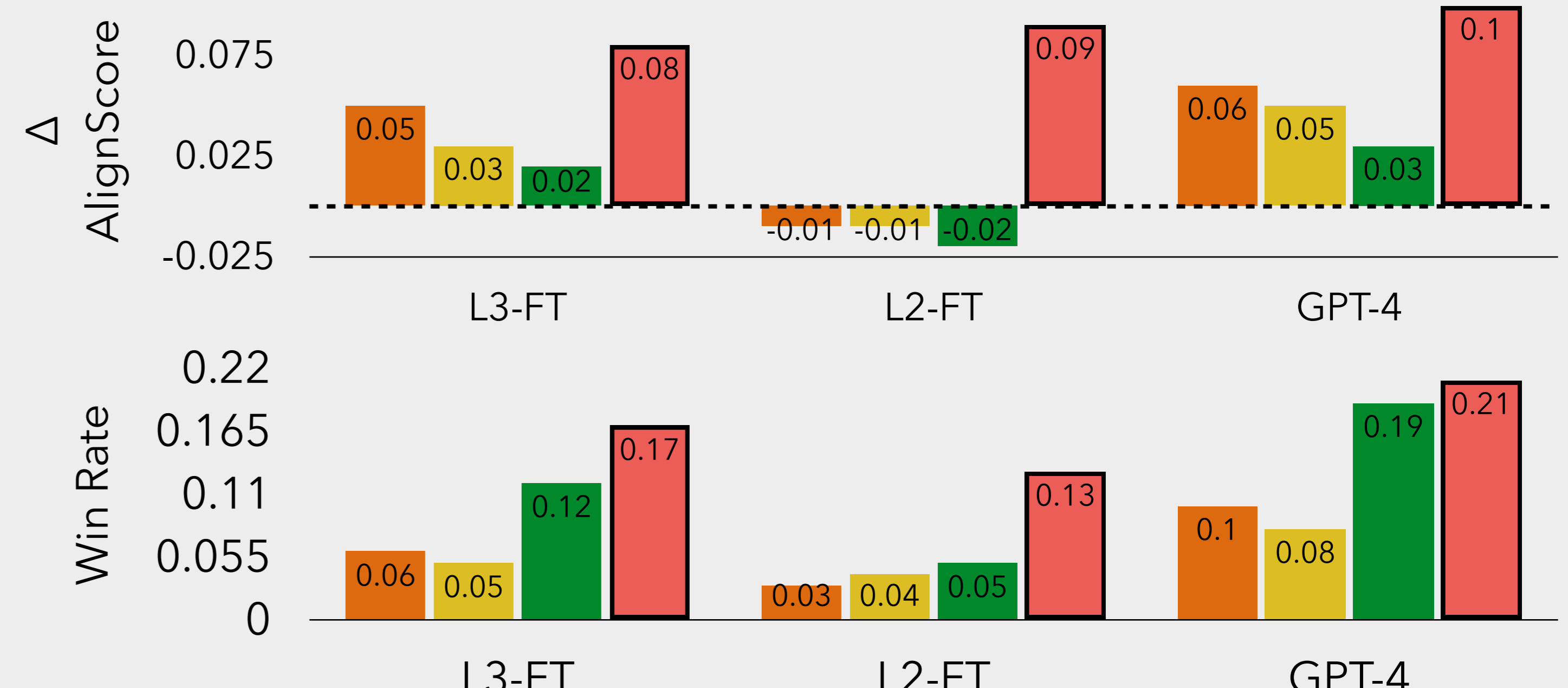
Downstream evaluation of refinement on TofuEval



Result 2: Proposed M_{critic} in DCR outperforms refining with off-the-shelf critic models

- Shepherd
- UltraCM
- Selfee-13b
- DCR (proposed)

Downstream evaluation of refinement on TofuEval



Takeaways

- DCR does better than baseline refinement techniques
- Decomposing Critique into Detect and Critique helps train targeted and more accurate critic models
- Fine-tuned models can perform at-par with stronger models using DCR

More results on critique correctness and faithful refinements in the paper!