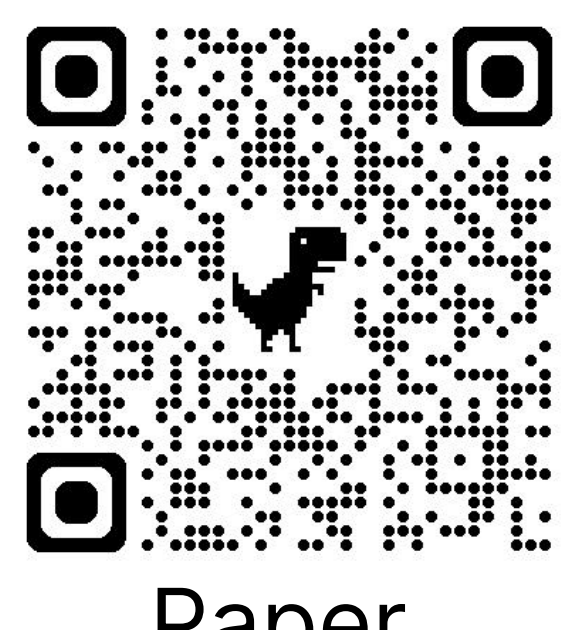


Contrastive Instruction Tuning (CoIN)

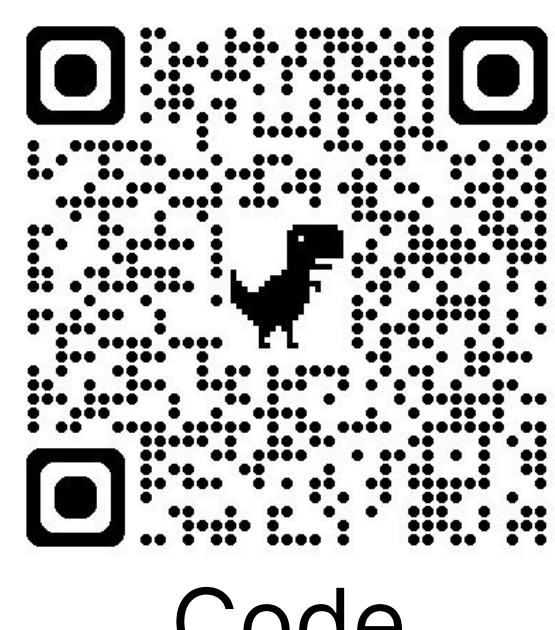
Motivation

LLMs are sensitive to variations in instructions

Tianyi Lorena Yan¹, Fei Wang¹, James Y. Huang¹, Wenxuan Zhou¹, Fan Yin²
Aram Galstyan¹, Wenpeng Yin³, Muhao Chen⁴



Paper



Code

Correspondence to: [✉ tianyiy@usc.edu](mailto:tianyiy@usc.edu) / [🐦 @LorenaYannnnn](https://twitter.com/LorenaYannnnn)

Instruction

Review the sentence below and identify whether its grammar is "acceptable" or "unacceptable".

Input The mechanical doll wriggled itself loose.

Acceptable. ✓

Instruction Variant

Please evaluate the grammar of the following sentences and mark them as "acceptable" or "unacceptable".

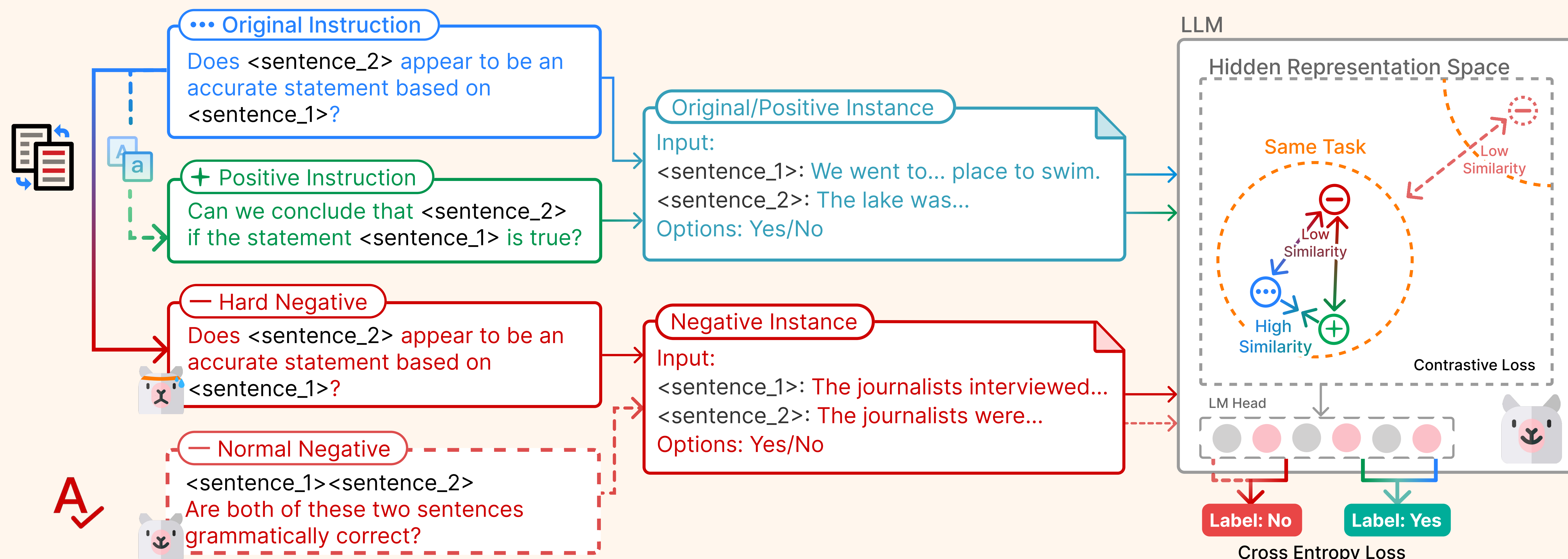
Input The mechanical doll wriggled itself loose.

Unacceptable. ✗

Method

In LLMs' hidden representation space:

Encourage **semantically equivalent** inputs to **stay close** to each other & **dissimilar** ones to be **far apart**



Training

- 25 datasets from FLAN collection (52k instruction-instance pairs)
- + **Positive sample:** Instruction paraphrases (Avoid making assumptions about types of variation in instructions)
- **Negative sample:** Randomly select one instance from the remaining dataset

Experiment Setup

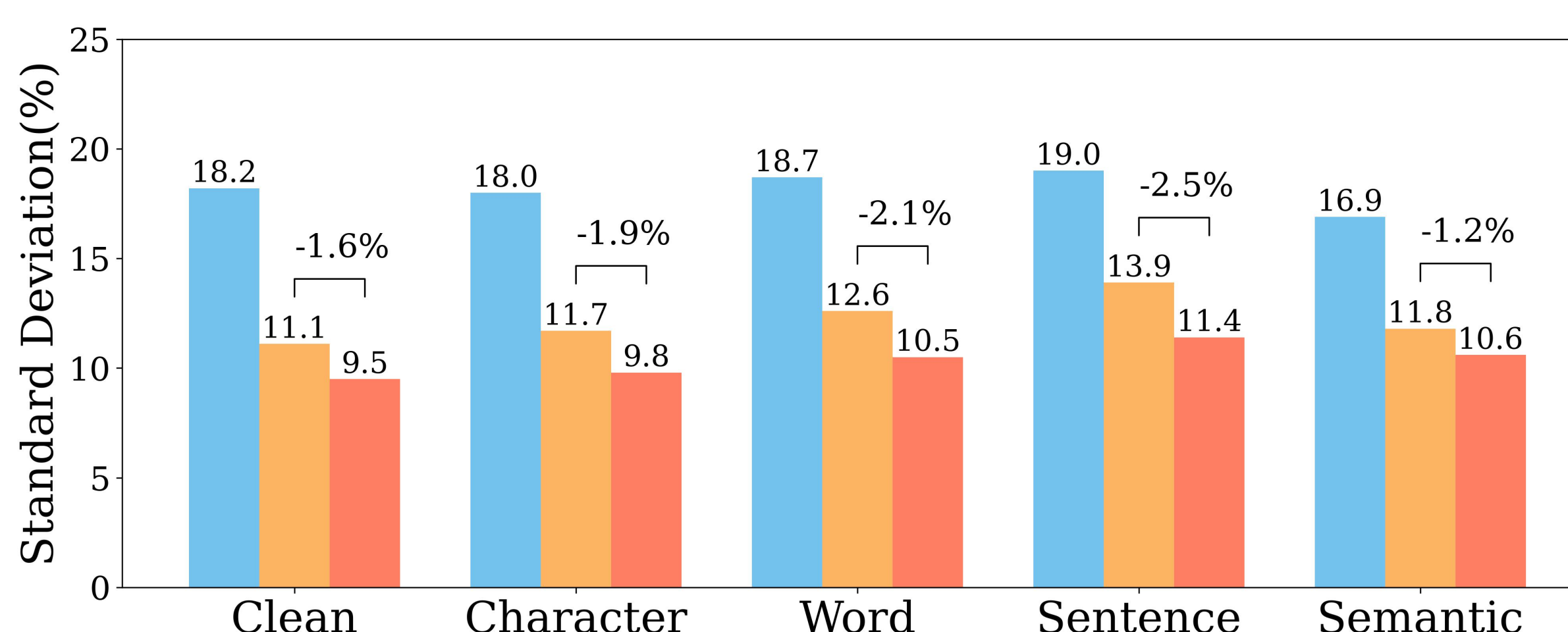
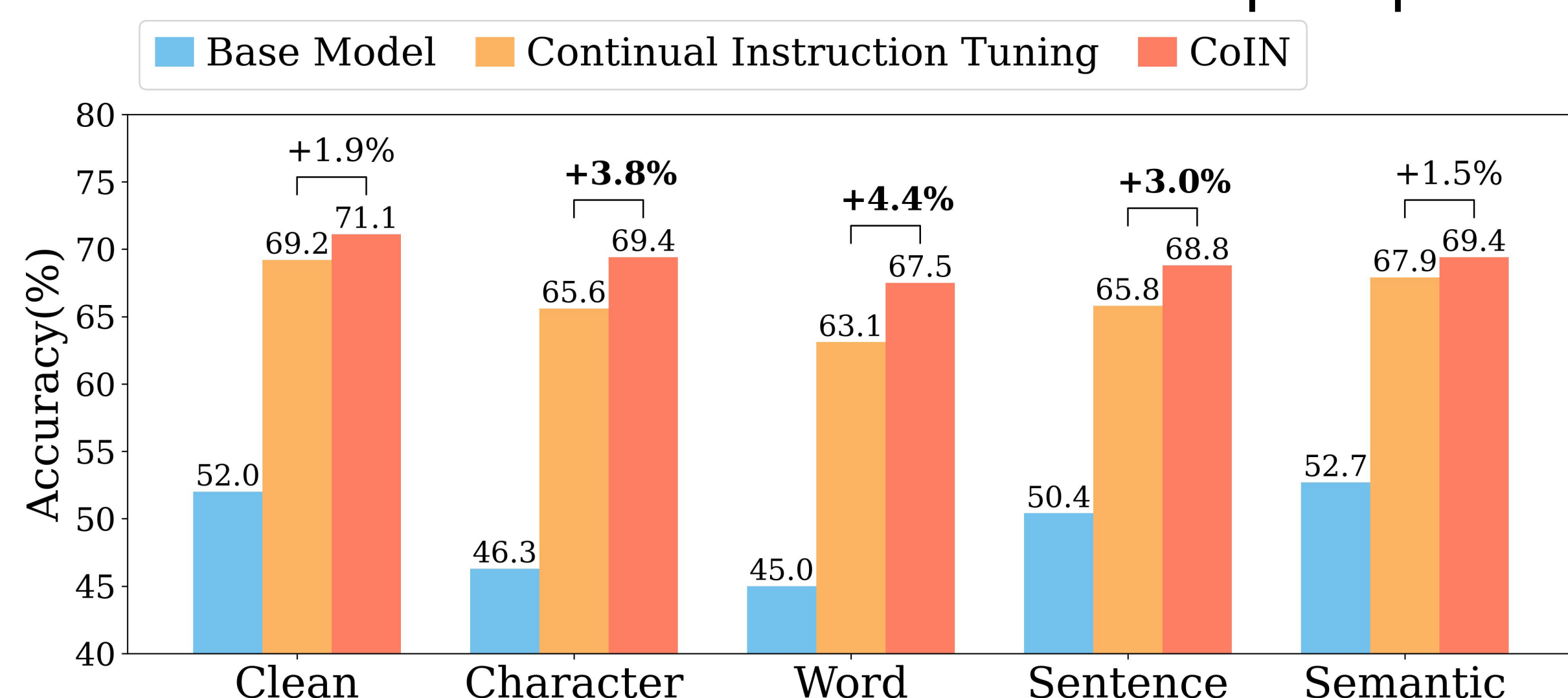
Evaluation

PromptBench + GLUE
Six clean instructions + perturbed versions

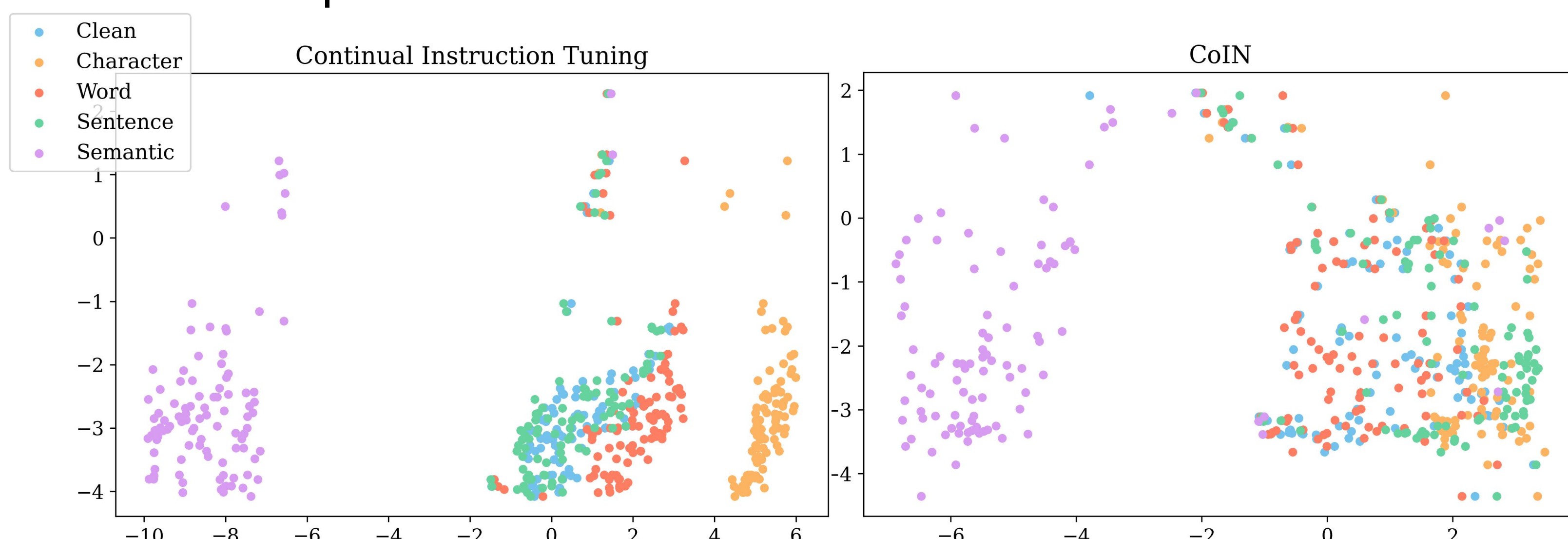
Clean: Review the sentence below and identify whether its grammar is 'Acceptable' or 'Unacceptable':
Character: Reiew the seVntence below and identify wheoher its gVammar is 'Acceptable' or 'Unacceptable':
Word: Analyzed the assertion below and ascertain whether its grammar is 'Acceptable' or 'Unacceptable':
Sentence: Review the sentence below and identify whether its grammar is 'Acceptable' or 'Unacceptable' LGOZMPXsPd:
Semantic: Evaluate the sentence below and determine if its grammar is 'Acceptable' or 'Unacceptable':
 * All instructions are **unseen** during training

Results & Analyses

1. Consistent improvement in performance w/o introducing any new data and training steps + Generalizes from paraphrases to all four types of variations in instructions



2. Closer Representations of Instruction Variations



3. Positive Impacts on Different Tasks

Task	Continual Instruction Tuning		COIN		Δ	
	Accuracy (%)	Std	Accuracy	Std	Accuracy	Std
Sentiment Analysis	89.0	4.1	90.4	3.1	+1.4	-1.1
Natural Language Inference	64.4	3.7	66.1	3.5	+1.7	-0.2
Paraphrase Identification	63.0	11.0	68.5	5.9	+5.4	-5.1
Grammar Correctness	62.0	9.2	68.4	3.9	+6.3	-5.3

Future Applications

CoIN can be applied to enhance models' robustness on other prompt component (e.g. system prompts, few-shot demonstration) and other modalities