

Contrastive Instruction Tuning



Tianyi Lorena Yan¹



Fei Wang¹



James Y. Huang¹



Wenxuan Zhou¹



Fan Yin²



Aram Galstyan¹



Wenpeng Yin³



Muhao Chen⁴



Are LLMs good human instruction followers yet?

Are LLMs good human instruction followers yet?

Instruction

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

Input

The mechanical doll wriggled itself loose.



Are LLMs good human instruction followers yet?

Instruction

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

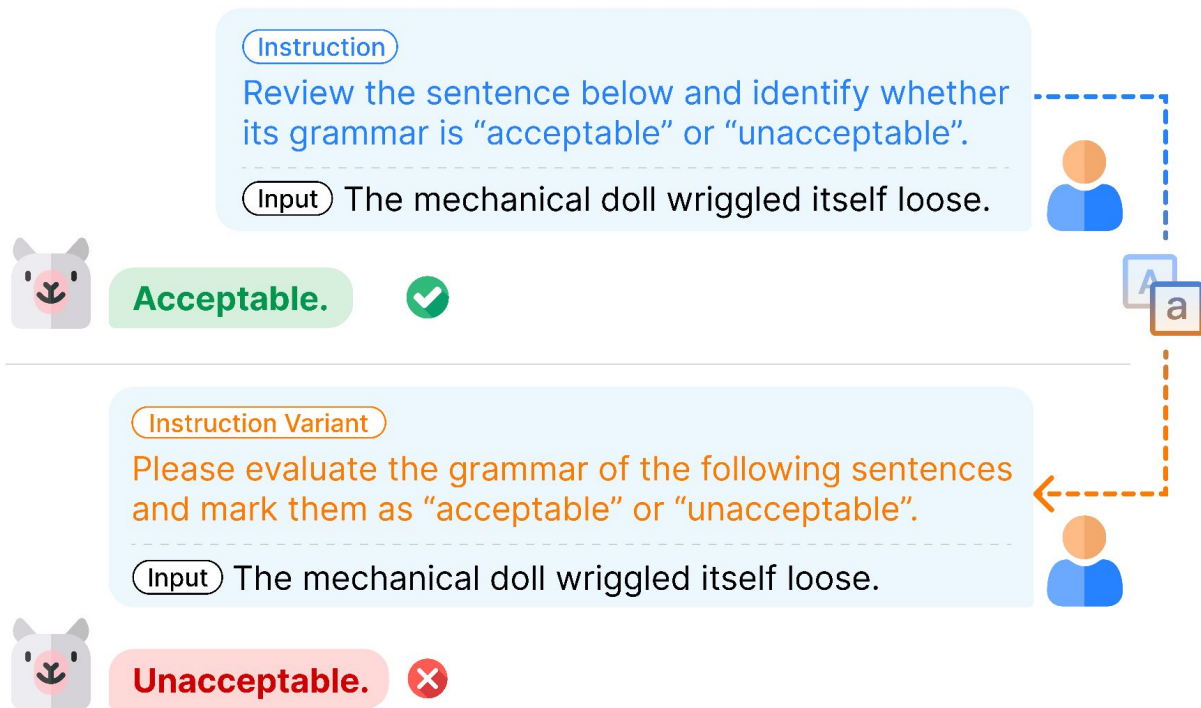
Input The mechanical doll wriggled itself loose.



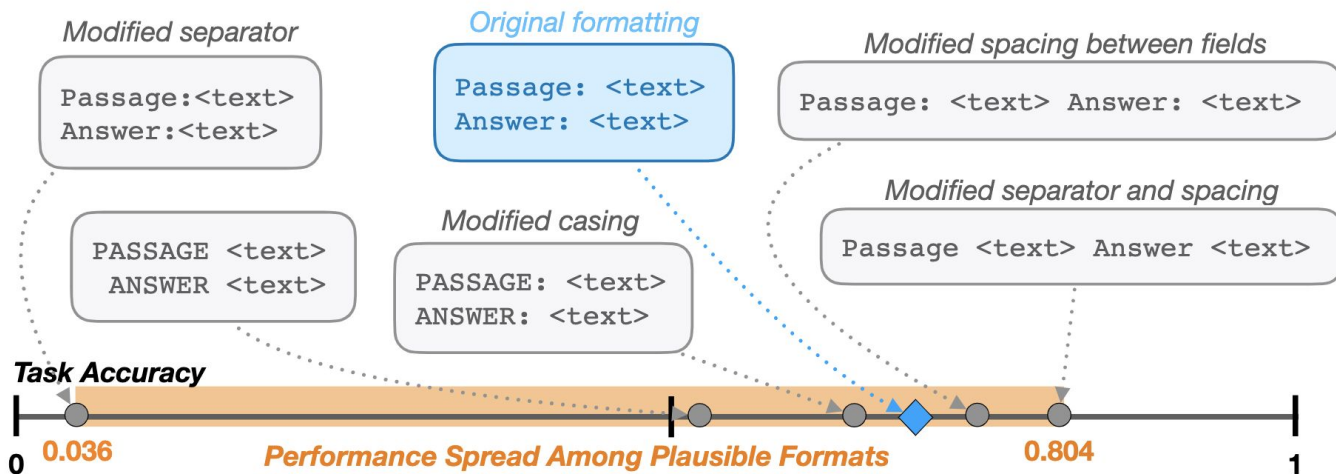
Acceptable.



Are LLMs good human instruction followers yet?





Issue: LLMs are sensitive to variations in instructions





76 accuracy points difference

Issue: LLMs are sensitive to variations in instructions

Prompt	Sample	
As a mathematics instructor, calculate the answer to the following problem related to if a number is a prime:	Question: Let $z(a) = -871*a + 415$. Is $z(-16)$ a composite number? Answer:	User 1
 Yes. ✓		
<hr/>		
As a mathematics instructor, calculate the answer to the following problem related to if a number is a prime:	Question: Let $z(a) = -871*a + 415$. Is $z(-16)$ a composite number? Answer:	User 2
 No. ✗		

(a) Typos lead to errors in math problems.

Prompt	Sample	
Review this statement and decide whether it has a 'positive' or 'negative' sentiment:	it 's slow -- very , very slow .	User 1
 Negative. ✓		
<hr/>		
Analyze this assertion and defining whether it is a 'positive' or 'negative' sentiment:	it 's slow -- very , very slow .	User 2
 Postive. ✗		

(b) Synonyms lead to errors in sentiment analysis problems.



ChatGPT gives inconsistent answers when facing variations in instructions.

Issue: LLMs are sensitive to variations in instructions

Prompt	Sample	User
As a mathematics instructor, calculate the answer to the following problem related to if a number is a prime:	Question: Let $z(a) = -871*a + 415$. Is $z(-16)$ a composite number? Answer:	User 1
As a mathematics instructor, calculate the answer to the following problem related to if a number is a prime:	Question: Let $z(a) = -871*a + 415$. Is $z(-16)$ a composite number? Answer:	User 2

Yes. ✓

No. ✗

(a) Typos lead to errors in math problems.

Prompt	Sample	User
Review this statement and decide whether it has a 'positive' or 'negative' sentiment:	it 's slow -- very , very slow .	User 1
Analyze this assertion and defining whether it is a 'positive' or 'negative' sentiment:	it 's slow -- very , very slow .	User 2

Negative. ✓

Postive. ✗

(b) Synonyms lead to errors in sentiment analysis problems.



ChatGPT gives inconsistent answers when facing variations in instructions.

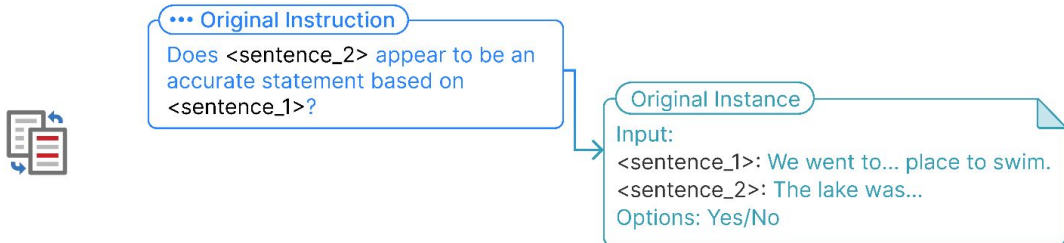


Our solution: Contrastive Instruction Tuning

Idea: Encourage semantically equivalent inputs to stay close to each other while dissimilar ones to be far apart in LLMs' hidden representation space

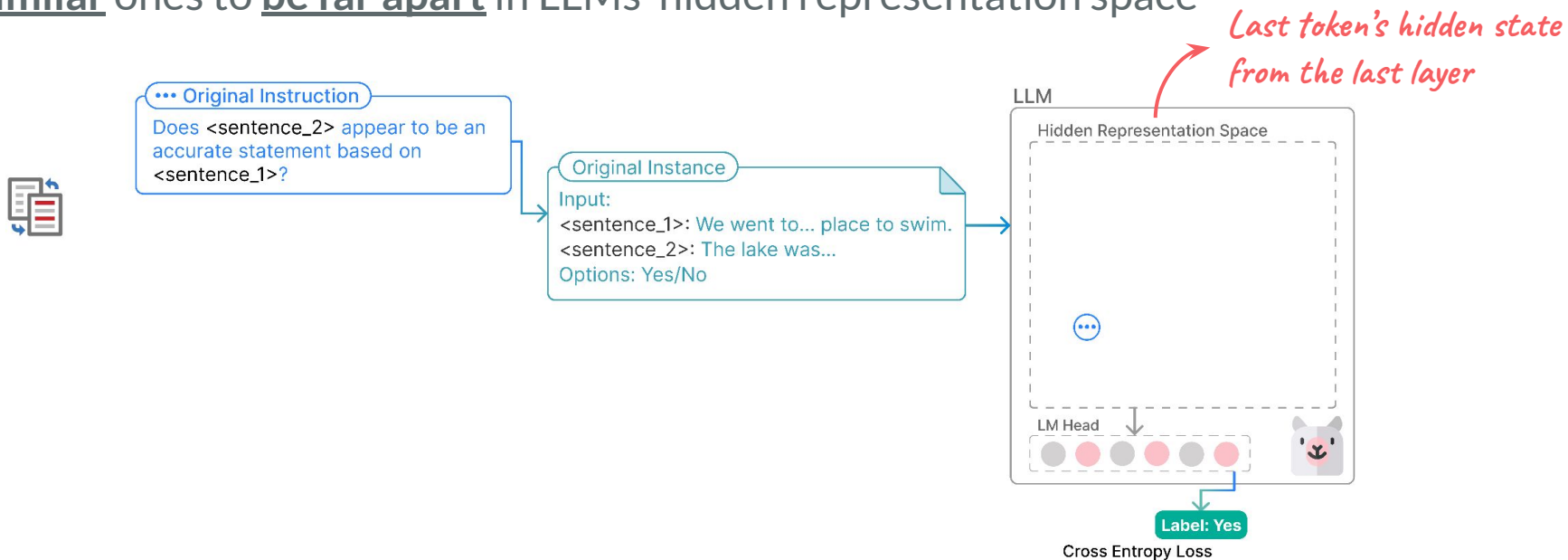
Our solution: Contrastive Instruction Tuning (CoIN)

Idea: Encourage semantically equivalent inputs to stay close to each other while dissimilar ones to be far apart in LLMs' hidden representation space



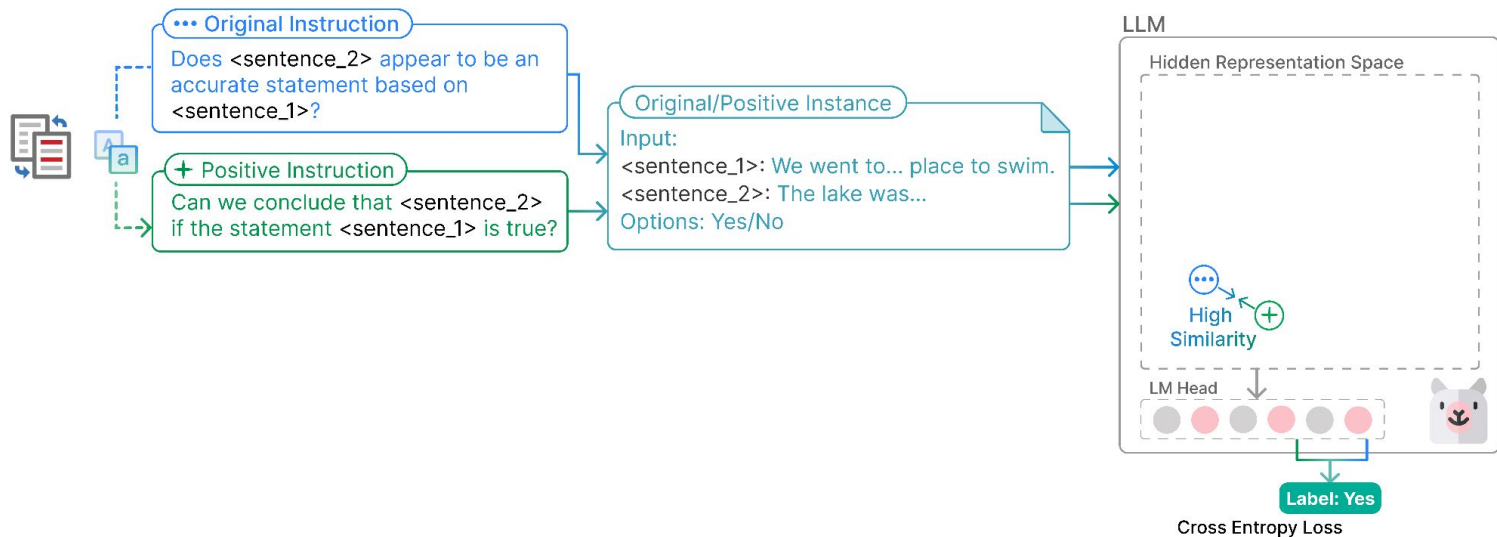
Our solution: Contrastive Instruction Tuning

Idea: Encourage semantically equivalent inputs to stay close to each other while dissimilar ones to be far apart in LLMs' hidden representation space



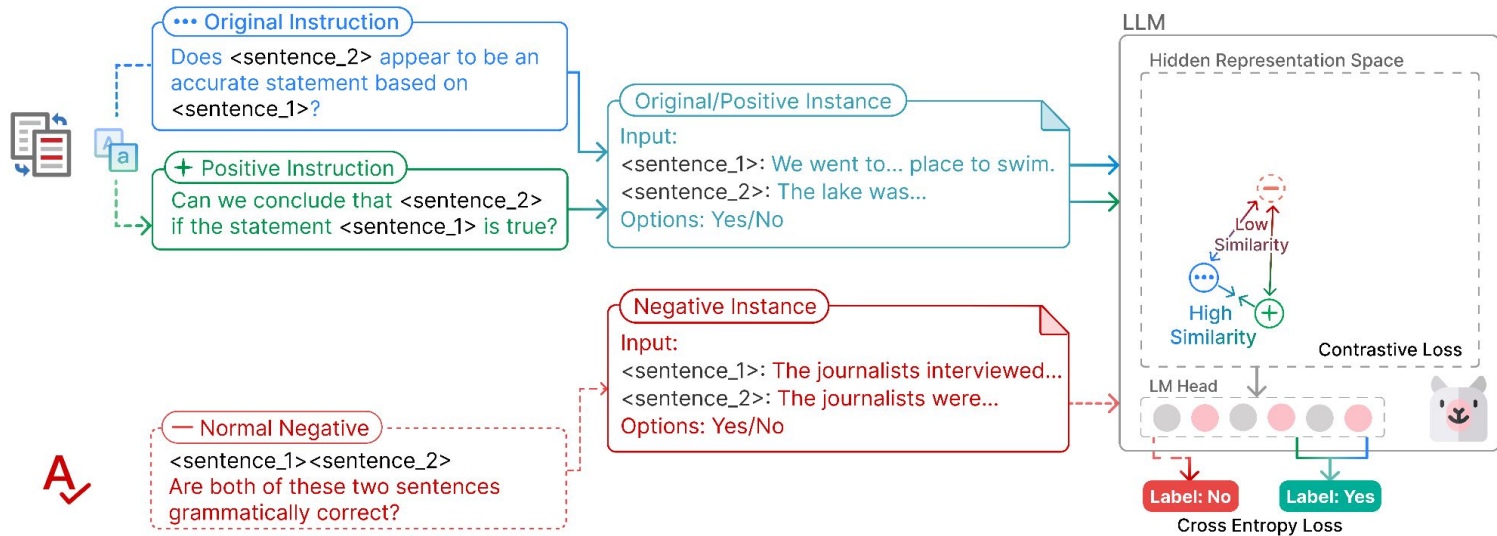
Our solution: Contrastive Instruction Tuning

Idea: Encourage semantically equivalent inputs to stay close to each other while dissimilar ones to be far apart in LLMs' hidden representation space



Our solution: Contrastive Instruction Tuning

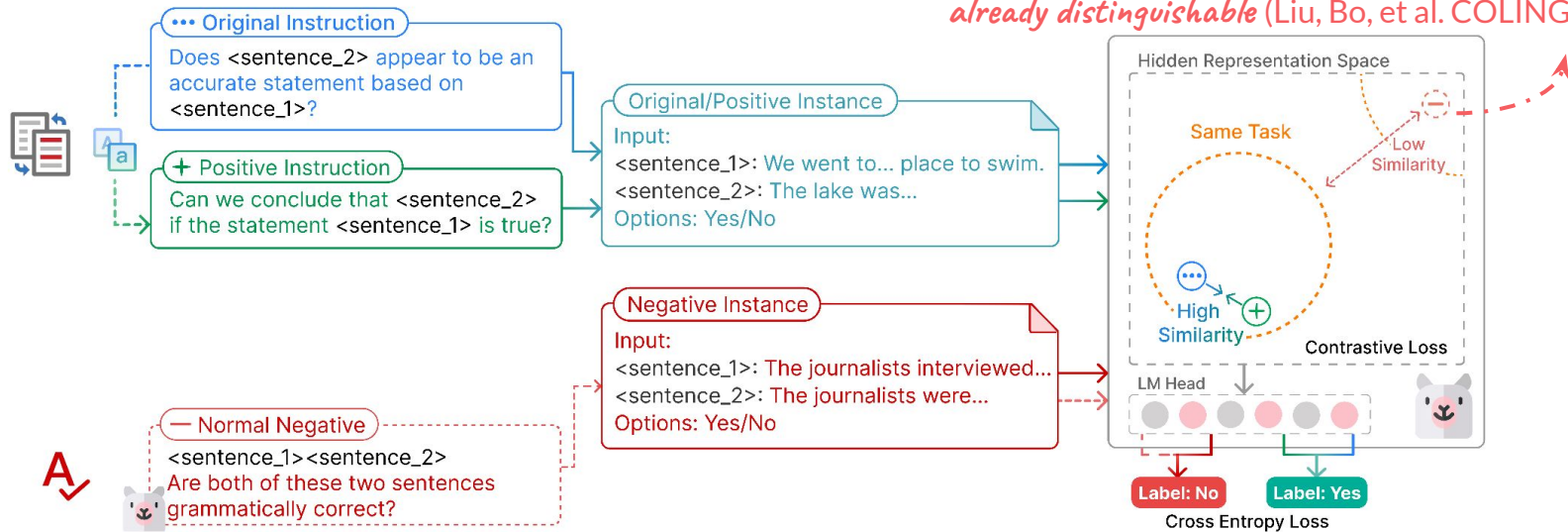
Idea: Encourage semantically equivalent inputs to stay close to each other while dissimilar ones to be far apart in LLMs' hidden representation space



Our solution: Contrastive Instruction Tuning

Idea: Encourage semantically equivalent inputs to stay close to each other while dissimilar ones to be far apart in LLMs' hidden representation space

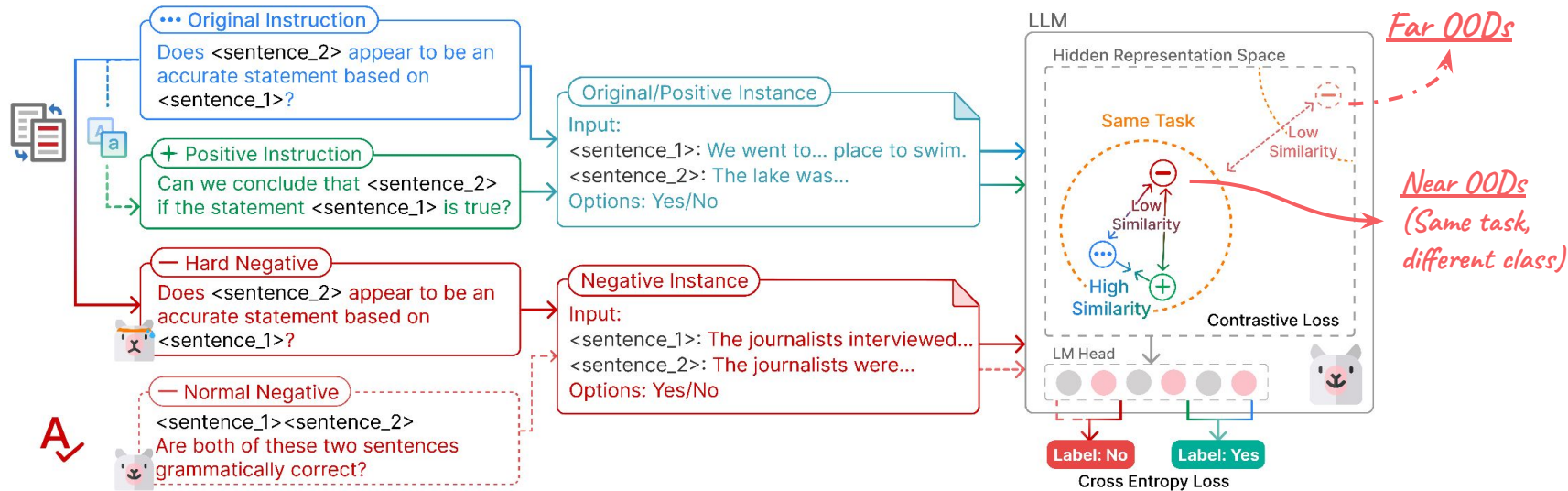
Data w/ instructions of different tasks (Far OODs) are already distinguishable (Liu, Bo, et al. COLING 2024)



Liu, B., Zhan, L., Lu, Z., Feng, Y., Xue, L., & Wu, X. M. How Good Are Large Language Models at Out-of-Distribution Detection?. COLING 2024

Our solution: Contrastive Instruction Tuning

Idea: Encourage semantically equivalent inputs to stay close to each other while dissimilar ones to be far apart in LLMs' hidden representation space



Experiment Setup: Training

- Data:
 - Datasets from the FLAN collection (52k instruction-instance pairs)
 - For every pair from a dataset
 - **Positive sample**: randomly select a predefined instruction template as paraphrases
(Avoid making assumptions about specific types of variations in instructions)
 - **Negative sample**: randomly select another pair from the remaining dataset
- Model: Alpaca LoRA

(Refer to paper for more experiment details)

Experiment Setup: Evaluation

- Sample 300 instruction-instance pairs from each of the 10 GLUE tasks
- Select six clean instructions predefined for each task & add perturbations at four levels following PromptBench

Experiment Setup: Evaluation

- Sample 300 instruction-instance pairs from each of the 10 GLUE tasks
- Select six clean instructions predefined for each task & add perturbations at four levels following PromptBench

Clean: Review the sentence below and identify whether its grammar is 'Acceptable' or 'Unacceptable':

Character: **Reiew** the **seVntence** below and **identifpy 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

Experiment Setup: Evaluation

- Sample 300 instruction-instance pairs from each of the 10 GLUE tasks
- Select six clean instructions predefined for each task & add perturbations at four levels following PromptBench

Clean: Review the sentence below and identify whether its grammar is 'Acceptable' or 'Unacceptable':

Character: **Reiew** the **seVntence** below and **identifpy wheoher** its **gVammar** is 'Acceptable' or 'Unacceptable':

Word: **Analyzed** the **assertion** below and **ascertain** whether its grammar is 'Acceptable' or 'Unacceptable':

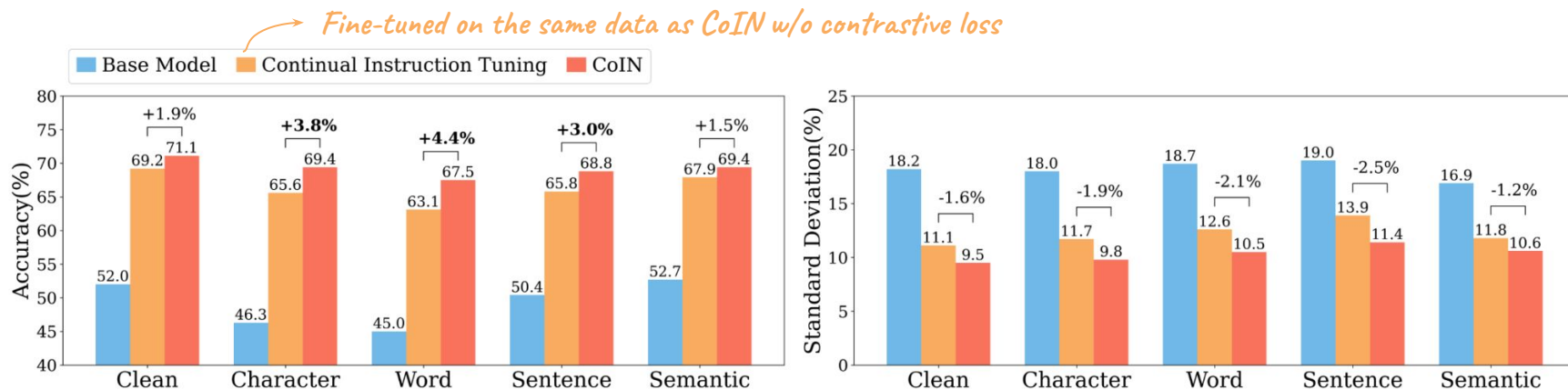
Sentence: Review the sentence below and identify whether its grammar is 'Acceptable' or 'Unacceptable' **LGOZMPXsPd**:

Semantic: **Evaluate the sentence below and determine if its grammar is 'Acceptable' or 'Unacceptable'**:

* All instructions are unseen during training

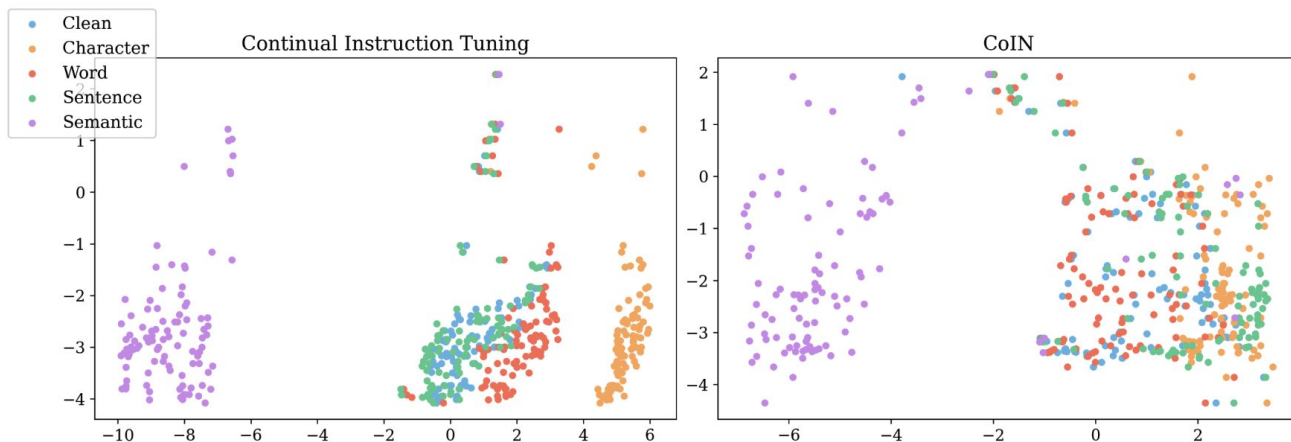
- Metric: Average accuracy (exact match) and standard deviation

Main Results



- 🤔 Consistent improvement in accuracy & decrease in standard deviation w/o introducing any new data & training steps
- 😊 Able to generalize from paraphrases to all types of variations in instructions

Analyses: Closer Representations of Instruction Variations



UMAP (McInnes et al., 2020) visualization of the hidden representations of decoder's last output token (300 datapoints from CoLA dataset)

- 🤔 Continual instruction tuning:
Instructions with different variations are clustered into distinct groups → Higher sensitivity
- 😊 CoIN:
Larger overlap between clean & perturbed instructions → More robust to instruction variations

Analyses: Impact on Different Tasks

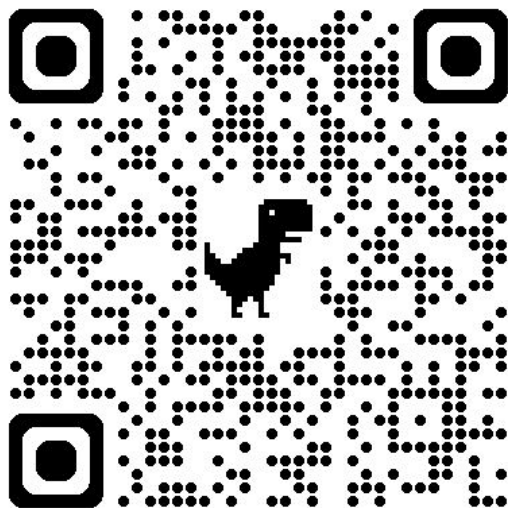
(%)	Continual Instruction Tuning		COIN		Δ	
Task	Accuracy	Std	Accuracy	Std	Accuracy	Std
Sentiment Analysis	89.0	4.1	90.4	3.1	+1.4	-1.1
Natural Language Inference	64.4	3.7	66.1	3.5	+1.7	-0.2
Paraphrase Identification	63.0	11.0	68.5	5.9	+5.4	-5.1
Grammar Correctness	62.0	9.2	68.4	3.9	+6.3	-5.3

- More evident improvement in paraphrase identification and grammar correctness
- Directly benefit from model's more refined ability to group textual inputs with similar semantic meanings

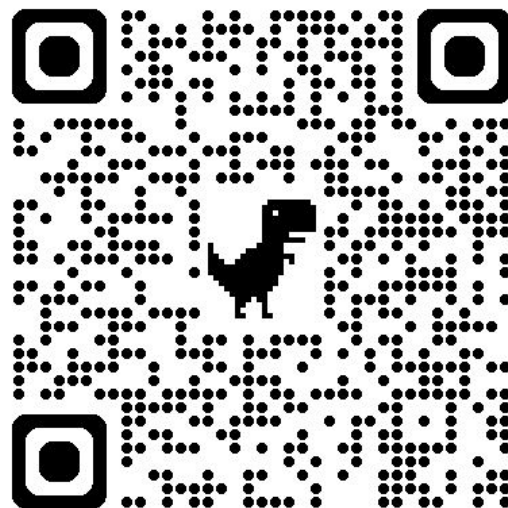
Conclusions

- We propose Contrastive Instruction Tuning (CoIN) that aligns hidden representations of semantically equivalent instruction-instance pairs
- Evaluation results on PromptBench w/ instruction variations at character, word, sentence, and semantic level demonstrate CoIN's effectiveness of enhancing LLMs' robustness to instruction variations
- CoIN can be applied to enhance models' robustness on other prompt component (e.g. system prompts, few-shot demonstration) and other modalities

QR Codes



Paper



Code



`tianyiy@usc.edu` /



`@LorenaYannnnn`

Special thanks to all my amazing collaborators!

