

# Collaborative Creative Writing

Vishakh Padmakumar



ML<sup>2</sup>

Machine Learning  
for Language

# Acknowledgements

Coauthors: He He, Tuhin Chakrabarty



Feedback and Discussions: Richard Pang, Nitish Joshi, Alicia Parrish, Angelica Chen, Nicholas Lourie, Chen Zhao, Jason Phang and the rest of the ML<sup>2</sup> group :)

VI

What are the important and timely research directions in NLP?



There are several important and timely research directions in natural language processing (NLP), some of which include:




1. **Pre-training and fine-tuning:** Developing techniques for pre-training large language models on large amounts of data, and then fine-tuning them on smaller, task-specific datasets. This has been very successful in a number of NLP tasks, and researchers are continuing to work on improving the performance and efficiency of these approaches.
2. **Transfer learning:** Developing techniques that allow NLP models to transfer knowledge from one language or domain to another, which can be particularly useful in situations where there is limited data available in a particular language or domain.
3. **Dialogue systems:** Developing NLP models that can effectively handle complex, open-ended conversations, such as those that might be used in virtual assistants or chatbots. This involves developing models that can understand and generate natural language, as well as models that can reason about and respond appropriately to different types of input.
4. **Multilingual and cross-lingual NLP:** Developing NLP models that can handle multiple languages and translate between languages, which is important for applications such as machine translation and information retrieval in a global context.
5. **Sentiment analysis:** Developing NLP models that can accurately identify and classify the sentiment (e.g., positive, negative, or neutral) of text, which has applications in a



There are several important and timely research directions in natural language processing (NLP), some of which include: 👍 🗨️

1. **Pre-training and fine-tuning:** Developing techniques for pre-training large language models on large amounts of data, and then fine-tuning them on smaller, task-specific datasets. This has been very successful in a number of NLP tasks, and researchers continue to work on improving the performance and efficiency of these approaches.
2. **Transfer learning:** Developing techniques that allow NLP models to transfer knowledge from one language or domain to another, which can be particularly useful in situations where there is limited data available in a particular language or domain.
3. **Dialogue systems:** Developing NLP models that can effectively handle complex, open-ended conversations, such as those that might be used in virtual assistants or chatbots. This involves developing models that can understand and generate natural language, as well as models that can reason about and respond appropriately to different types of input.
4. **Multilingual and cross-lingual NLP:** Developing NLP models that can handle multiple languages and translate between languages, which is important for applications such as machine translation and information retrieval in a global context.
5. **Sentiment analysis:** Developing NLP models that can accurately identify and classify the sentiment (e.g., positive, negative, or neutral) of text, which has applications in a



**Maybe a little  
too generic**

VI

Can you write a paragraph about why we should work on human-AI collaboration



There are several reasons why it is important to work on human-AI collaboration.



First, as AI continues to advance and become more prevalent in a wide range of industries, it is important to find ways to effectively integrate AI into human workflows and decision-making processes. By working on human-AI collaboration, we can help to ensure that AI is used in a way that is beneficial to both humans and AI, and that it is aligned with ethical and societal norms.

Second, human-AI collaboration has the potential to bring many benefits, including increased efficiency, productivity, and effectiveness. By combining the strengths of both humans and AI, it is possible to achieve outcomes that would not be possible with either humans or AI alone.

Third, as AI becomes more advanced and capable of handling more complex tasks, it is important to consider the potential negative impacts of AI, such as the displacement of human workers or the potential for AI to make biased or unfair decisions. By involving humans in the decision-making process, it is possible to incorporate human judgment and values into AI systems, helping to mitigate these negative impacts.



There are several reasons why it is important to work on human-AI collaboration.



First, as AI continues to advance and become more prevalent in a wide range of industries, it is important to find ways to effectively integrate AI into human workflows and decision-making processes. By working on human-AI collaboration, we can help to ensure that AI is used in a way that is beneficial to both humans and AI, and that it is aligned with ethical and societal norms.

Second, human-AI collaboration has the potential to bring many benefits, including increased efficiency, productivity, and effectiveness. By combining the strengths of both humans and AI, it is possible to achieve outcomes that would not be possible with either humans or AI alone.

Third, as AI becomes more advanced and capable of handling more complex tasks, it is important to consider the potential negative impacts of AI, such as the displacement of human workers or the potential for AI to make biased or unfair decisions. By involving humans in the decision-making process, it is possible to incorporate human judgment and values into AI systems, helping to mitigate these negative impacts.



# Overview

- Broad Direction: How can we assist writers at creative writing tasks?

# Overview

- Broad Direction: How can we assist writers at creative writing tasks?
- Content generation is exploding and there is a lot of demand for high-quality prose

# Overview

- Broad Direction: How can we assist writers at creative writing tasks?
- Content generation is exploding and there is a lot of demand for high-quality prose
- LLMs generate fluent text and can be finetuned for various purposes

# Overview

- Broad Direction: How can we assist writers at creative writing tasks?
- Content generation is exploding and there is a lot of demand for high-quality prose
- LLMs generate fluent text and can be finetuned for various purposes
- Opportune time to work on writing assistants<sup>[1,2,3]</sup>
  - [Verse-by-Verse](#), [HemingwayApp](#), [Adept](#)

[1] Andy Coenen, Luke Davis, Daphne Ippolito, Emily Reif, and Ann Yuan. 2021. Wordcraft: a human-ai collaborative editor for story writing. CoRR. abs/2107.07430

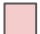









[2] Lee, Mina, Percy Liang, and Qian Yang. "CoAuthor: Designing a Human-AI Collaborative Writing Dataset for Exploring Language Model Capabilities." arXiv preprint arXiv:2201.06796 (2022).

[3] Du, Wanyu, et al. "Read, Revise, Repeat: A System Demonstration for Human-in-the-loop Iterative Text Revision." arXiv preprint arXiv:2204.03685 (2022).

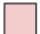




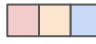

# Overview

- Broad Direction: How can we assist writers at creative writing tasks?
- Many open questions **bridging NLP and HCI**:
  - How do we design the ***most effective collaboration setup*** to help human users?
  - How do we train models to generate ***helpful*** suggestions?
  - What is the ***best form of user feedback*** and how do we ***incorporate it in model training***?
  - How do we assist users in ***content planning*** for long form creative writing?
  - How do we ensure ***equitable creative writing assistance*** to all users?

# Overview

- Broad Direction: How can we assist writers at creative writing tasks?
- Many open questions **bridging NLP and HCI**:
  - How do we design the **most effective collaboration setup** to help human users? 
  - How do we train models to generate **helpful** suggestions? 
  - What is the **best form of user feedback** and how do we **incorporate it in model training**? 
  - How do we assist users in **content planning** for long form creative writing? 
  - How do we ensure **equitable creative writing assistance** to all users? 
- Our work:
  - Machine-in-the-Loop Rewriting for Creative Image Captioning   
  - Collaborative Poetry writing with Instruction Tuning  

# Overview

- Broad Direction: How can we assist writers at creative writing tasks?
- Many open questions **bridging NLP and HCI**:
  - How do we design the **most effective collaboration setup** to help human users? 
  - How do we train models to generate **helpful** suggestions? 
  - What is the **best form of user feedback** and how do we **incorporate it in model training**? 
  - How do we assist users in **content planning** for long form creative writing? 
  - How do we ensure **equitable creative writing assistance** to all users? 
- Our Work:
  - **Machine-in-the-Loop Rewriting for Creative Image Captioning** 
  - Collaborative Poetry writing with Instruction Tuning 

# Machine-in-the-Loop Rewriting for Creative Image Captioning

Vishakh Padmakumar, He He

NAACL 2022



ML<sup>2</sup>

Machine Learning  
for Language



# Machine-In-The-Loop Creative Writing

- Creative writing tasks can be challenging for both humans and machines.

# Machine-In-The-Loop Creative Writing

- Creative writing tasks can be challenging for both humans and machines.
  - Humans would benefit from suggestions on wording and framing their ideas<sup>[1]</sup>
  - Models are able to rewrite spans of text<sup>[4]</sup> but struggle with global coherence<sup>[2,3]</sup>

[1] Monica J Garfield. 2008. Creativity support systems. In *Handbook on Decision Support Systems 2*, pages 745–758. Springer

[2] Elizabeth Clark, Anne Spencer Ross, Chenhao Tan, Yangfeng Ji, and Noah A Smith. 2018. Creative writing with a machine in the loop: Case studies on slogans and stories. In *23rd International Conference on Intelligent User Interfaces*, pages 329–340.

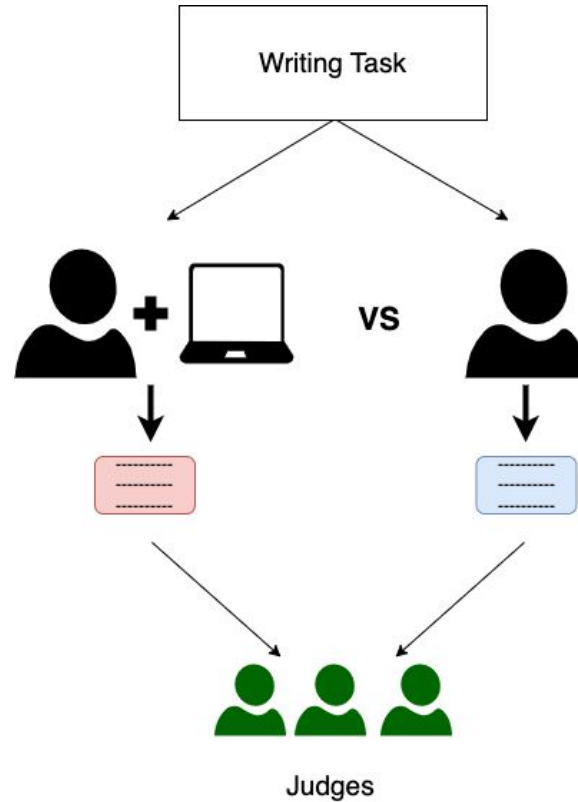
[3] Nader Akoury, Shufan Wang, Josh Whiting, Stephen Hood, Nanyun Peng, and Mohit Iyyer. 2020. STORIUM: A Dataset and Evaluation Platform for Machine-in-the-Loop Story Generation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6470–6484

[4] Chris Donahue, Mina Lee, and Percy Liang. 2020. Enabling language models to fill in the blanks. In *Association for Computational Linguistics (ACL)*

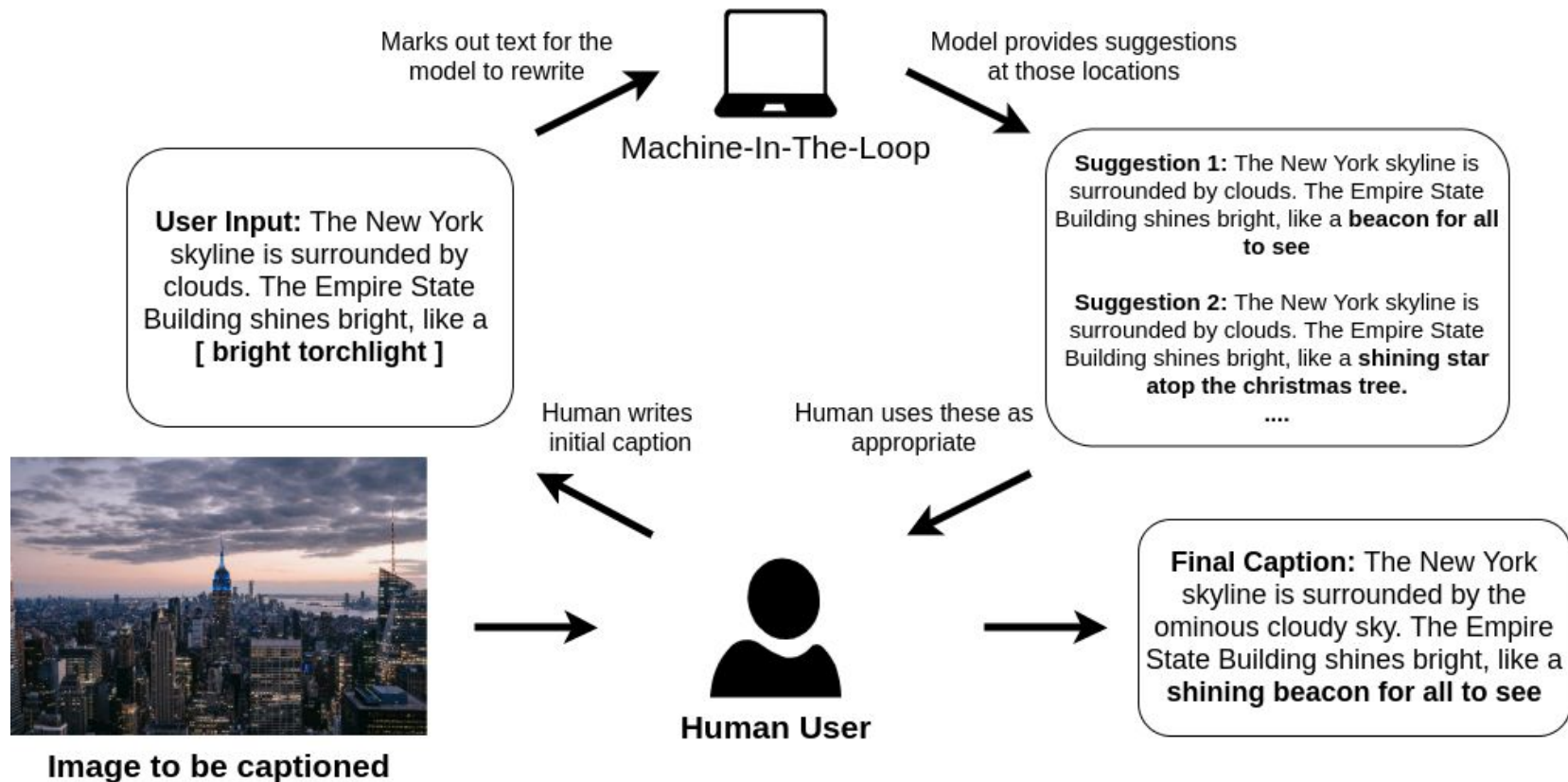
# Machine-In-The-Loop Creative Writing

- Creative writing tasks can be challenging for both humans and machines.
  - Humans would benefit from suggestions on wording and framing their ideas
  - Models are able to rewrite spans of text but struggle with global coherence
- Motivates a cooperative setting: Can a model help the author improve their creative output?

# Machine-In-The-Loop Creative Writing

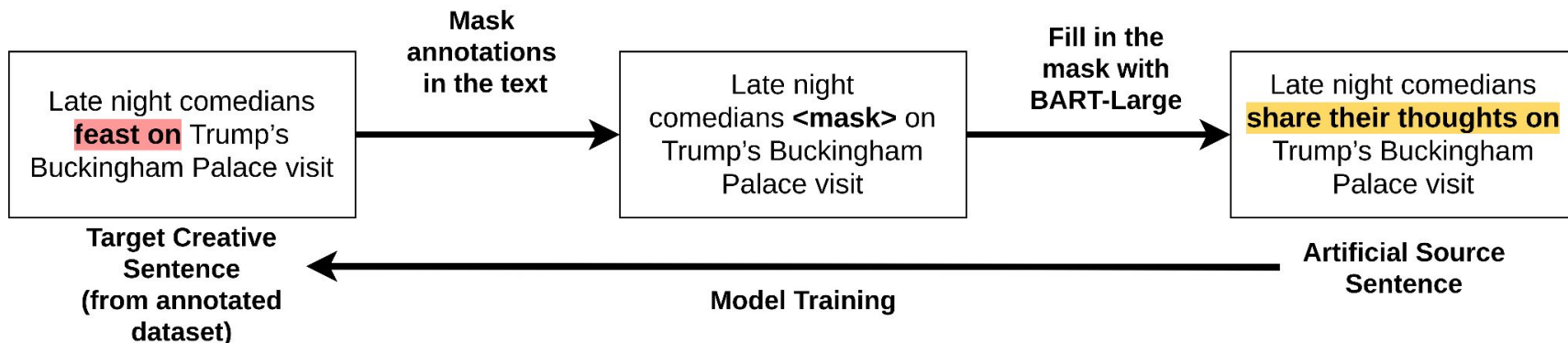


# Task Setup - Creative Image Captioning



# Training the Creative Rewriting Assistant (CRA) Model

- Fine-tuning Data: We create a pseudo-parallel corpus of creative sentences (annotated for literary devices) and corresponding generic sentences
- **CRA** is a fine-tuned BART-Large model



# Demo



Enter your text

A young man wearing jeans and a green shirt is kneeling down by a mud puddle, while sticking his finger in the water. An owl of beautiful contemplation is a troubled youth.

Suggest

Finish

Our result: 111

Select the suggestion that you like best ( ... ) (Original text)



**Suggestion 1**

A young man wearing jeans and a green shirt is kneeling down by a mud puddle, while sticking his finger in the water. An owl of beautiful contemplation is a **troubled** youth.



**Suggestion 2**

A young man wearing jeans and a green shirt is kneeling down by a mud puddle, while sticking his finger in the water. An owl of beautiful contemplation is a **troubled** youth.



**Suggestion 3**

A young man wearing jeans and a green shirt is kneeling down by a mud puddle, while sticking his finger in the water. An owl of beautiful contemplation is a **troubled** youth.

**SUCCESS!** Well for you, what you can do is write a recommendation after the input and well use the area where you would like suggestions (marked by a start tag, **T**) and an end tag, **T** and then hit **MOCKIT**. You can also hover **\_** to indicate a blank line (should be **\_**), for example

**Input:** The New York skyline (in the background surrounded by the intense dark sky, but the Empire State Building shines bright, like a [purple highlight])

**Output:** The New York skyline (in the background surrounded by the intense dark sky, but the Empire State Building shines bright, like a **blaze of blue** in the darkness)

3. Once you hit **SUCCESS** you will receive 3 suggestions on the bottom-right pane. Select one which you feel is appropriate and use **with the chosen suggestion** again and use the tool repeatedly.

4. Continue until you are happy with the description with a **minimum of 2 iterations**.

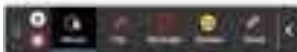
**History**

All interactions are logged here.



**Older Description 1**

A young man wearing jeans and a green shirt is kneeling down by a mud puddle, while sticking his finger in the water. An owl of beautiful contemplation is a **troubled** youth.



# Project Roadmap

- Do users find CRA model suggestions helpful?
- Are users more effective at creative image captioning with model help?
- How does collaboration with the model impact different users?
- Can the model be adapted to learn from observed user interactions to provide better assistance?



# Project Roadmap

- **Do users find CRA model suggestions helpful?**
- Are users more effective at creative image captioning with model help?
- How does collaboration with the model impact different users?
- Can the model be adapted to learn from observed user interactions to provide better assistance?

# Do Users Find Model Suggestions Helpful?

We compare the **CRA model** to a **baseline BART model** with an A/B user study (n=50)

# Do Users Find Model Suggestions Helpful?

We compare the **CRA model** to a **baseline BART model** with an A/B user study (n=50)

- On average, users find the **CRA model** to be more helpful than **BART** by a statistically significant margin

Results from Post Completion Survey

	<b>BART</b>	<b>CRA</b>
<b>Model Helpfulness</b>	2.23	<b>3.06</b>

# Do Users Find Model Suggestions Helpful?

We compare the **CRA model** to a **baseline BART model** with an A/B user study (n=50)

- On average, users find the **CRA model** to be more helpful than **BART** by a statistically significant margin
- Users accept larger fraction of suggestions from the **CRA model**

Results from Post Completion Survey

	<b>BART</b>	<b>CRA</b>
<b>Model Helpfulness</b>	2.23	<b>3.06</b>

Model Comparison via Interaction Logs

	<b>BART</b>	<b>CRA</b>
<b>Avg # Requests</b>	<b>3.02</b>	2.82
<b>% Acceptance</b>	24.5%	<b>31.9%</b>
<b>Rouge-L Retention</b>	0.744	<b>0.824</b>

# Project Roadmap

- **Do users find CRA model suggestions helpful?**
- Are users more effective at creative image captioning with model help?
- How does collaboration with the model impact different users?
- Can the model be adapted to learn from observed user interactions to provide better assistance?

# Project Roadmap

- Users find **CRA suggestions more helpful** than a baseline model
- **Are users more effective at creative image captioning with model help?**
- How does collaboration with the model impact different users?
- Can the model be adapted to learn from observed user interactions to provide better assistance?

# Are Users More Effective With Model Help?

- We collected captions for 100 images from **solo writers**, from users collaborating with the **CRA model** and the **BART model**.

# Are Users More Effective With Model Help?

- We collected captions for 100 images from **solo writers**, from users collaborating with the **CRA model** and the **BART model**.
- For each caption pair, we collect 3 annotations for which is better and take a majority vote



# Are Users More Effective With Model Help?

- We collected captions for 100 images from **solo writers**, from users collaborating with the **CRA model** and the **BART model**.

	# Majority Vote Wins		
<b>Human + CRA</b>	<b>57</b>	43	<b>Solo Writers</b>
<b>Human + CRA</b>	<b>54</b>	46	<b>Human + BART</b>
<b>Human + BART</b>	<b>55</b>	45	<b>Solo Writers</b>

# Project Roadmap

- Users find **CRA suggestions more helpful** than a baseline model
- **Are users more effective at creative image captioning with model help?**
- How does collaboration with the model impact different users?
- Can the model be adapted to learn from observed user interactions to provide better assistance?

# Project Roadmap

- Users find **CRA suggestions more helpful** than a baseline model
- **Collaborative users are more effective** at the creative writing task
- **How does collaboration with the model impact different users?**
- Can the model be adapted to learn from observed user interactions to provide better assistance?

# How Does **CRA** Impact Different Users?

**Effect of User Skill Level:** We divide users into two groups, **novice** and **skilled**, based on their self-rated writing skill.

# How Does **CRA** Impact Different Users?

**Effect of User Skill Level:** We divide users into two groups, **novice** and **skilled**, based on their self-rated writing skill.

**Takeaway:** **Skilled** users find the CRA model to be significantly more helpful

Results from Post Completion Survey

	<b>Novice (n=22)</b>	<b>Skilled (n=28)</b>
<b>Helpfulness</b>	2.27	<b>3.23</b>
<b># Requests</b>	<b>3.04</b>	2.64
<b>% Accepted</b>	29.8%	<b>33.7%</b>

# When is the Model Effective?

## Profile of Suggestion:

- The model performs best when rewriting shorter spans of larger texts

### Figurative Rewriting - Accepted Suggestion



A solemn woman place her mother's diary on a stepping stone her late father laid in the garden. The [ **surrounding pale grass gently sway in the cold breeze** ] while the woman ponders times of the past.

A solemn woman place her mother's diary on a stepping stone her late father laid in the garden. The **pale grass gently danced and teased in the wind** while the woman pondered times of the past.



# When is the Model Effective?

## Profile of Suggestion:

- The model performs best when rewriting shorter spans of larger texts
- Skilled writers tend to request this profile of suggestion



A solemn woman place her mother's diary on a stepping stone her late father laid in the garden. The [ **surrounding pale grass gently sway in the cold breeze** ] while the woman ponders times of the past.

A solemn woman place her mother's diary on a stepping stone her late father laid in the garden. The **pale grass gently danced and teased in the wind** while the woman pondered times of the past.



A child stands tall [ **in a wave on the beach.** ]

A child stands tall **by the waves on the beach.**



# Project Roadmap

- Users find **CRA suggestions more helpful** than a baseline model
- **Collaborative users are more effective** at the creative writing task
- **How does collaboration with the model impact different users?**
- Can the model be adapted to learn from observed user interactions to provide better assistance?



# Project Roadmap

- Users find **CRA suggestions more helpful** than a baseline model
- **Collaborative users are more effective** at the creative writing task
- **Model helps skilled writers more** potentially widening the gap in performance
- **Can the model be adapted to learn from observed user interactions to provide better assistance?**

# Can We Learn from User Feedback?

- We create a dataset from 50 sets of observed interactions.
- Sentence Pairs:
  - Original Text  $\mapsto$  Accepted Suggestion
  - Rejected Suggestion  $\mapsto$  Original Text
- Fine-tune the **initial CRA model** to **User-adapted CRA Model**

# Can We Learn from User Feedback?

- Fine-tune the **initial CRA model** to **User-adapted CRA Model**
- Compare the two models with an A/B user study (n=50)

# Can We Learn from User Feedback?

- Fine-tune the **initial CRA model** to **User-adapted CRA Model**
- Compare the two models with an A/B user study (n=50)
- On average, users find the **User-adapted CRA model** to be **more helpful** than **CRA model**, but not by a statistically significant margin

Results from Post Completion Survey

	<b>Initial CRA</b>	<b>User-adapted CRA</b>
<b>Helpfulness</b>	2.81	<b>3.05</b>
<b>Satisfaction</b>	3.67	<b>3.78</b>

Model Comparison via Interaction Logs

	<b>CRA</b>	<b>User-adapted CRA</b>
<b># Requests</b>	2.88	2.76
<b>% Acceptance</b>	<b>31.9%</b>	31.8%

# Project Roadmap

- Users find **CRA suggestions more helpful** than a baseline model
- **Collaborative users are more effective** at the creative writing task
- **Model helps skilled writers more** potentially widening the gap in performance
- **Can the model be adapted to learn from observed user interactions to provide better assistance?**

# Takeaways

- Users find **CRA suggestions more helpful** than a baseline model
- **Collaborative users are more effective** at the creative writing task
- **Model helps skilled writers more** potentially widening the gap in performance
- The **model becomes more helpful after updating it** from user interactions, **but not by much**

# What Next?

- **Model helps skilled writers more** potentially widening the gap in performance
  - How to better assist writers who aren't as comfortable with the language?
- The **model becomes more helpful after updating it** from user feedback, **but not by much**
  - How to learn more effectively from aggregated observed interactions?

# What Next?

- **Model helps skilled writers more** potentially widening the gap in performance
  - How to better assist writers who aren't as comfortable with the language?
- The model becomes more helpful after updating it from user feedback, but not by much
  - How to learn more effectively from aggregated observed interactions?



How can we design more accessible interactions?

# Cross-Task Generalization via Natural Language Crowdsourcing Instructions

Swaroop Mishra<sup>3\*</sup>, Daniel Khashabi<sup>1</sup>, Chitta Baral<sup>3</sup>, Hannaneh Hajishirzi<sup>1,2</sup>  
<sup>1</sup>University of Washington, <sup>2</sup>University of Washington, <sup>3</sup>Arizona State University

Input: She chose to make a salad for lunch on Sunday.  
How long did it take for her to make a salad?  
Output:

# InstructionNER: A Multi-Task Instruction-Based Generative Framework for Few-shot NER

Liwen Wang<sup>1\*</sup>, Rumei Li<sup>2\*</sup>, Yang Yan<sup>1</sup>, Yuanmeng Yan<sup>1</sup>, Sirui Wu<sup>1</sup>  
<sup>1</sup>Beijing University of Posts and Telecommunications, <sup>2</sup>Meituan Inc.

# How can instructions be accessible

# Instructions

tasks (e.g., text classification) do not solve classification tasks. The challenge in AI is to build models that learn a new task by understanding the readable instructions that define it. To study

# LANGUAGE MODELS ARE ZERO-SHOT

Jacob Brian, Lantao Yu, Karan Vasudevan, Martin Bosma\*, Vincent Y. Zhao\*, Kelvin Guu\*, Adams Wei Yu, Nan Du, Andrew M. Dai, and Quoc V. Le  
Google Research

## ABSTRACT

This paper explores a simple method for improving the performance of language models. We show that instruction tuning on a collection of datasets described via instructions improves zero-shot performance on unseen tasks. We take a 137B parameter model and evaluate this instruction-tuned model on 1000 tasks from the FLAN suite.

# Benchmarking Generalization via In-Context Instructions on 1,600+ Language Tasks

Yizhong Wang<sup>2</sup>, Swaroop Mishra<sup>3</sup>, Pegah Alipoormolabashi<sup>4</sup>, Yeganeh Kordi<sup>5</sup>, Amirreza Mirzaei<sup>7</sup>, Anjana Arunkumar<sup>3</sup>, Arjun Ashok<sup>6</sup>, Arut Selvan Dhanasekaran<sup>3</sup>, Atharva Naik<sup>7</sup>, David Stap<sup>8</sup>, Eshaan Pathak<sup>9</sup>, Giannis Karamanolakis<sup>10</sup>, Haizhi Gary Lai<sup>11</sup>, Ishan Purohit<sup>12</sup>, Ishani Mondal<sup>13</sup>, Jacob Anderson<sup>3</sup>, Kirby Kuznia<sup>3</sup>, Krima Doshi<sup>3</sup>, Maitreya Patej<sup>3</sup>, Kuntal Kumar Pal<sup>3</sup>, Mehrad Moradshahi<sup>14</sup>, Mihir Parmar<sup>3</sup>, Mirali Purohit<sup>15</sup>, Neeraj Varshney<sup>3</sup>, Ananya Kaza<sup>3</sup>, Pulkit Verma<sup>3</sup>, Ravsehaj Singh Puri<sup>3</sup>, Rushang Karia<sup>3</sup>, Shailaja Keyur Sampat<sup>3</sup>, Anurag Mishra<sup>16</sup>, Sujan Reddy<sup>17</sup>, Sumanta Patro<sup>18</sup>, Tanay Dixit<sup>19</sup>, Xudong Shen<sup>20</sup>, Noah A. Smith<sup>1,2</sup>, Hannaneh Hajishirzi<sup>1,2</sup>, Daniel Khashabi<sup>1</sup>  
<sup>1</sup>Stanford Univ., <sup>2</sup>Sharif Univ. of Tech., <sup>3</sup>Tehran Polytechnic, <sup>4</sup>PSG College of Tech., <sup>5</sup>IIT Kharagpur, <sup>6</sup>Govt. Polytechnic Rajkot, <sup>7</sup>Govt. Polytechnic Rajkot, <sup>8</sup>Stanford Univ., <sup>9</sup>Karnataka, <sup>10</sup>TCS Research, <sup>11</sup>IIT Madras, <sup>12</sup>National Univ. of Singapore, <sup>13</sup>Microsoft Research, <sup>14</sup>Stanford Univ., <sup>15</sup>IIT Madras, <sup>16</sup>National Univ. of Singapore, <sup>17</sup>IIT Kharagpur, <sup>18</sup>TCS Research, <sup>19</sup>IIT Madras, <sup>20</sup>National Univ. of Singapore

# Collaborative Poetry Writing with Instructions

{Tuhin Chakrabarty, Vishakh Padmakumar}, He He

EMNLP 2022



ML<sup>2</sup>

Machine Learning  
for Language

# Project Roadmap

- Can we train LLMs to satisfy creative writing instructions for poetry writing tasks?
- Can LLMs compose instructions seen at train time in unseen combinations?
- Can we help users complete creative writing tasks (poetry writing) using natural language instructions?

# Project Roadmap

- **Can we train LLMs to satisfy creative writing instructions for poetry writing tasks?**
- Can LLMs compose instructions seen at train time in unseen combinations?
- Can we help users complete creative writing tasks (poetry writing) using natural language instructions?

# Dataset of Instructions

- Staying on Subject:

Write a poetic verse that ends in a word which rhymes with 'late'

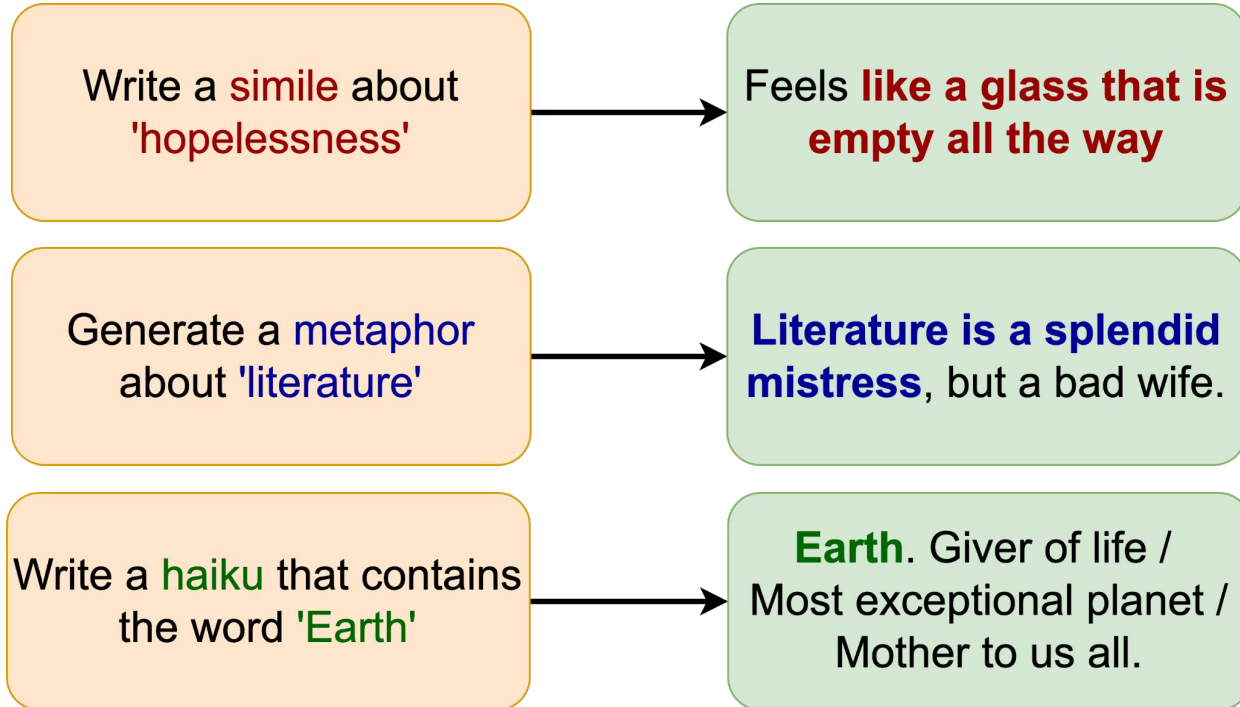
And homeward he turned,  
his path now **straight**.

Write a poetic sentence about 'god' and ending in 'eyes'

An all-powerful **God**, no  
escaping his **eyes**

# Dataset of Instructions

- Control on Literary Devices:



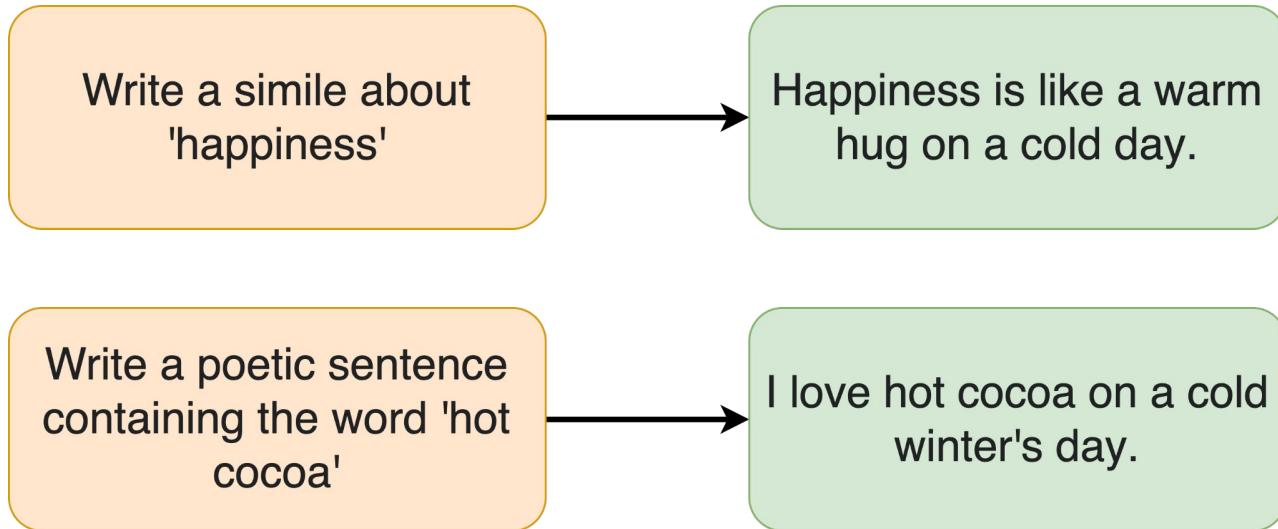
# Evaluation

- Hand crafted test sets of instructions for different kinds of capabilities



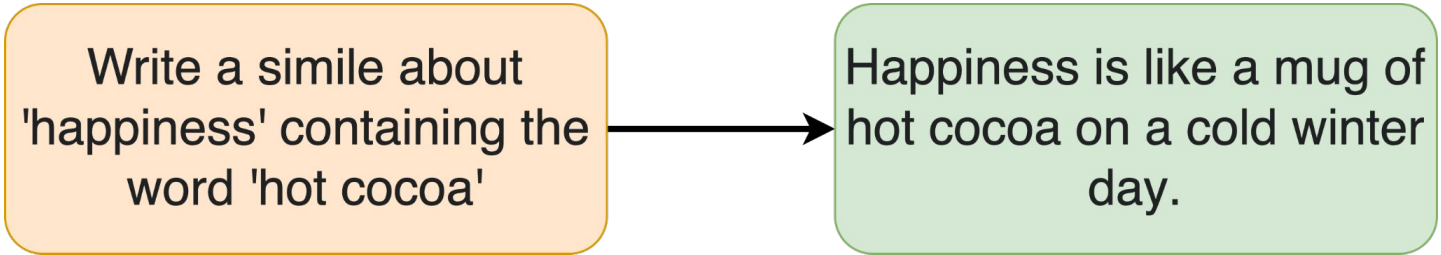
# Evaluation

- Hand crafted test sets of instructions for different kinds of capabilities
  - **Known Instruction Templates**



# Evaluation

- Hand crafted test sets of instructions for different kinds of capabilities
  - Known Instruction Templates
  - **Compositional Instruction Templates**



# Evaluation

- Hand crafted test sets of instructions for different kinds of capabilities
  - Known Instruction Templates
  - Compositional Instruction Templates
- Baselines
  - T0 models 🙌
    - **T0 - 3B** Finetuned + **T0pp - 11B** Few-Shot
  - **InstructGPT - 175B** - Zero Shot + Few-Shot



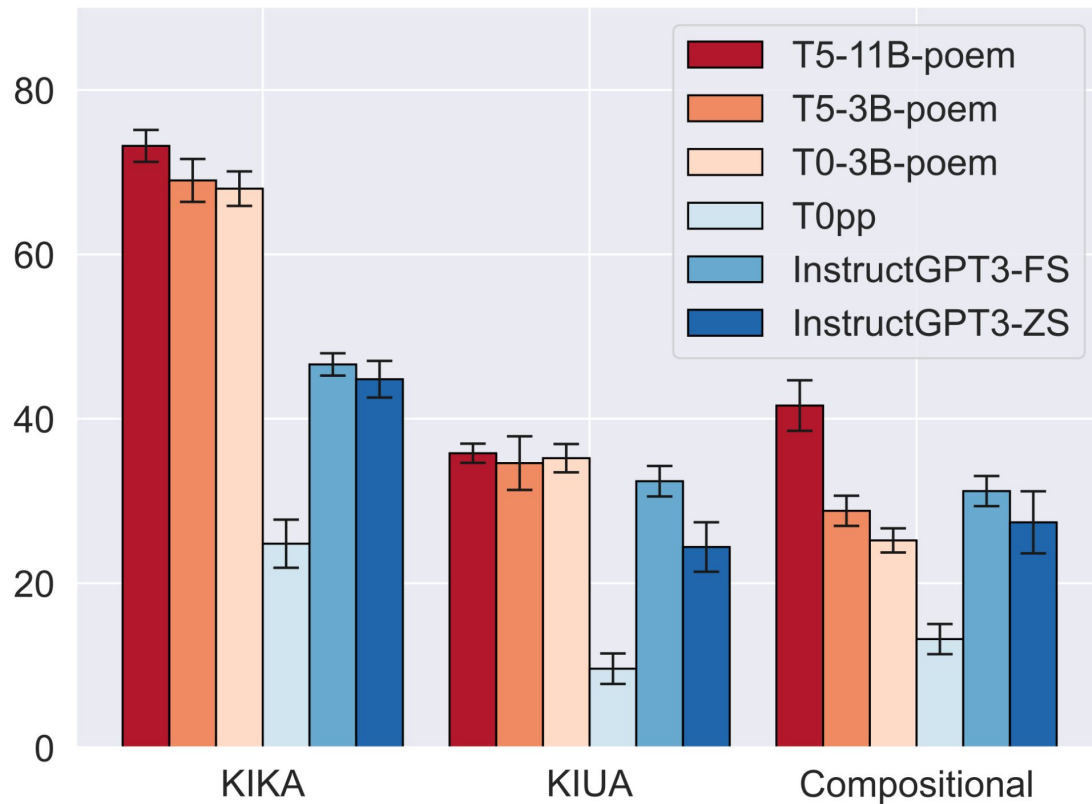
# Research Questions

- Can we train LLMs to satisfy creative writing instructions for poetry writing tasks?
- Can models compose instructions seen at train time in unseen combinations?
- Can we help users complete creative writing tasks using natural language instructions?

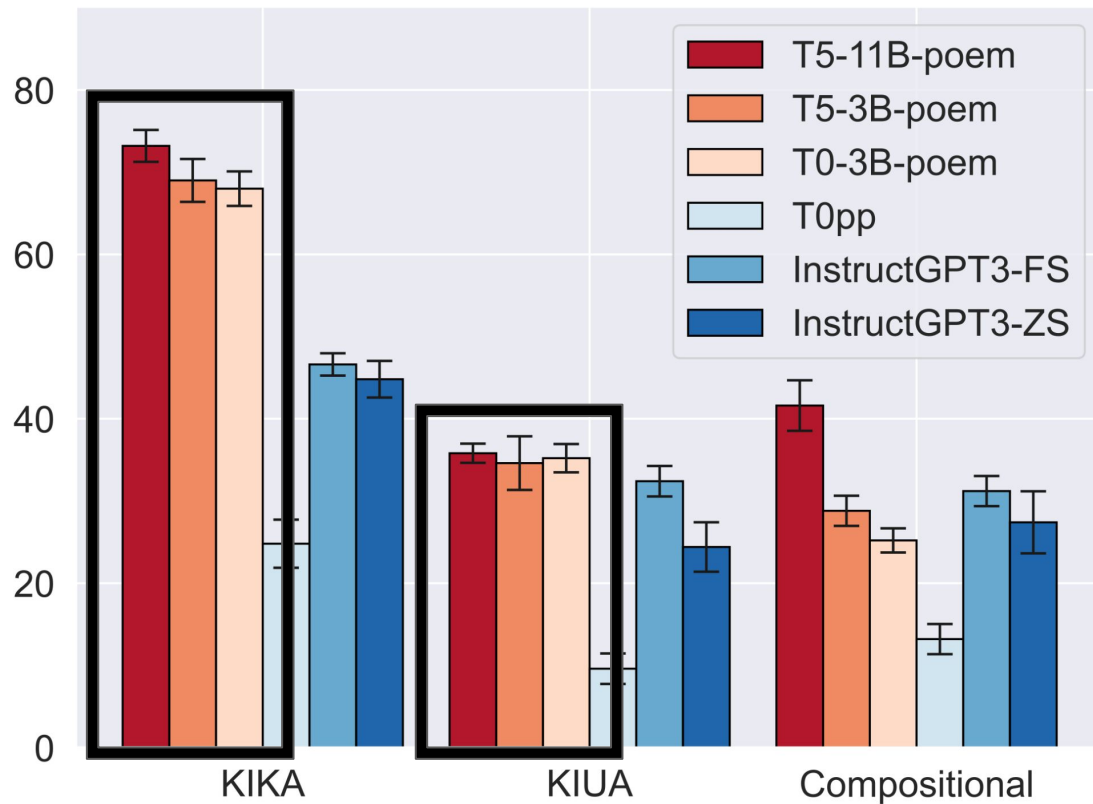
# Research Questions

- **Can we train LLMs to satisfy creative writing instructions for poetry writing tasks?**
- **Can models compose instructions seen at train time in unseen combinations?**
- Can we help users complete creative writing tasks using natural language instructions?

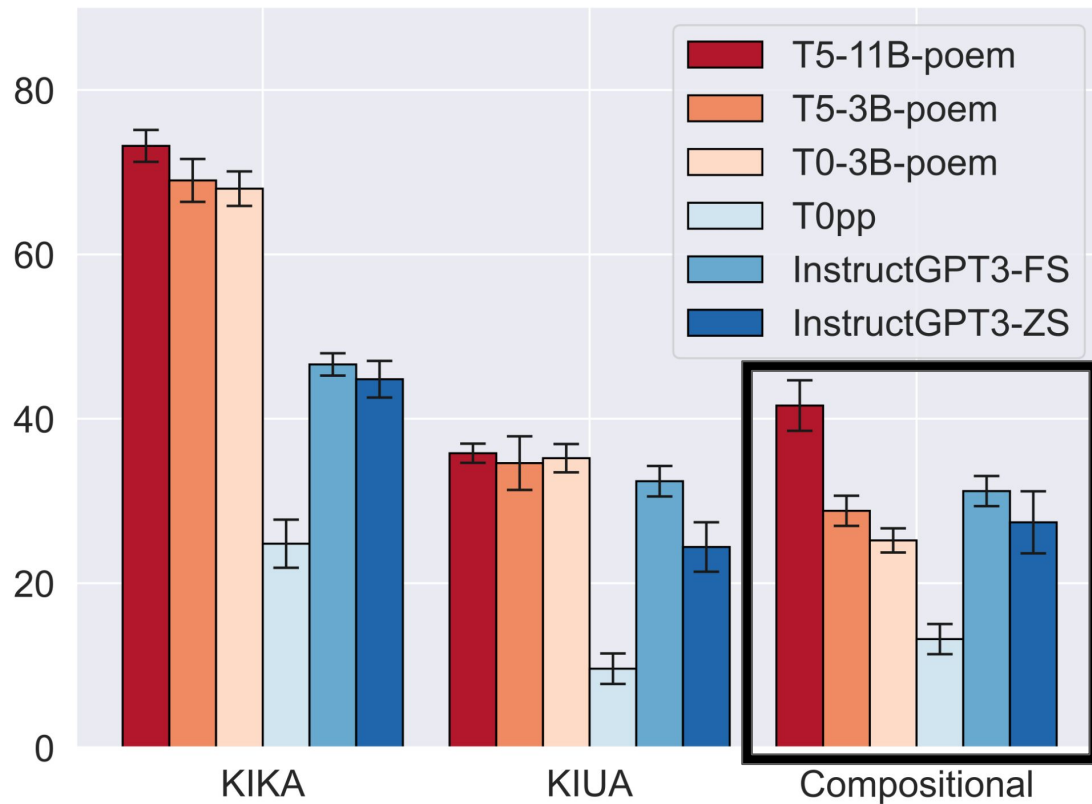
# Instruction Tuning - Evaluation



# Finetuned Models Are Strong In-Domain But Drop on Out-Of-Domain Data

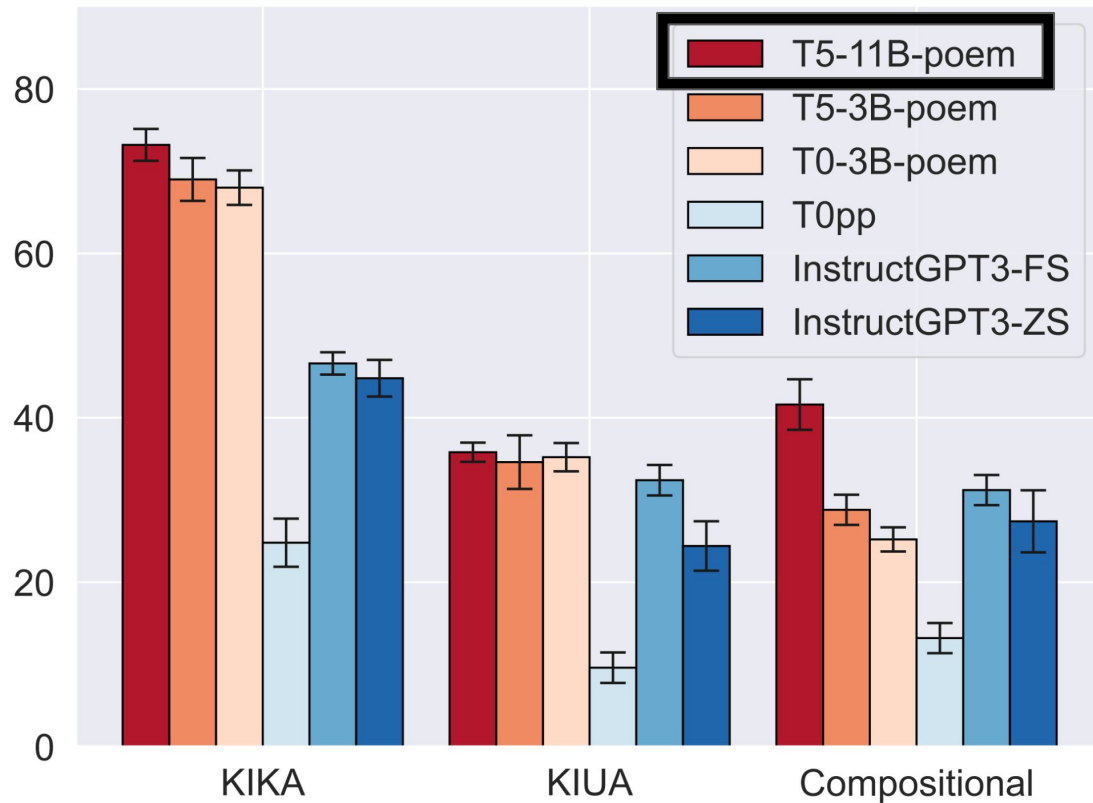


# Larger Models Compose Instructions Better



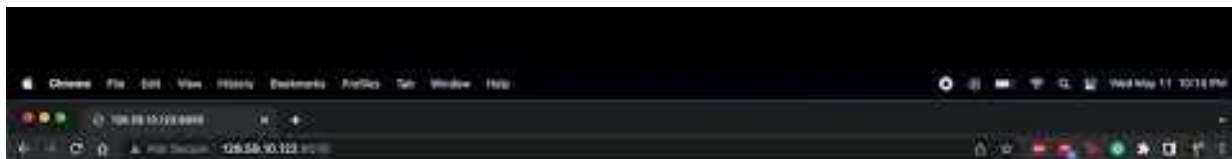


We use T5-11B for the User Study.



# Research Questions

- Can we train LLMs to satisfy creative writing instructions for poetry writing tasks?
- Can models compose instructions seen at train time in unseen combinations?
- **Can we help users complete creative writing tasks using natural language instructions?**



## CoPoet: Collaborative Poetry Writing with Instructions

### Poem:

The sunshine sways in the trees.

### Choose from below options:

- The birds sing a song of remembrance.
- The wind plays with the leaves.
- The chirping of the birds.
- The birds are singing and the breeze is humming.
- The breeze caresses the leaves.
- None of the above.

### Write your poem here:

#### Poem so far

The sunshine sways in the trees.  
The wind plays with the leaves.

#### Poem Title

Refresh

### Tools:

Choose an instruction template or write one below:

#### Suggest sentences about a topic:

- Suggest a sentence about a topic.
- Suggest a sentence ending with a custom word.
- Suggest a sentence starting with a custom word.
- Suggest a sentence about a specific topic and ending in a custom word.
- Suggest a sentence starting with a custom word and ending in a custom word.

#### Suggesting the next sentence:

- Suggest next sentence given what you've written so far.
- Suggest a topic for next sentence given what you've written so far.

#### Suggest a rhyming sentence:

- Suggest a sentence with a specific topic and rhyming with previous sentence.
- Suggest a sentence with a specific topic and rhyming with a custom word.

#### Suggest a simile or metaphor:

- Suggest a metaphor about a specific topic.
- Suggest a simile about a specific topic.

#### Your Instruction

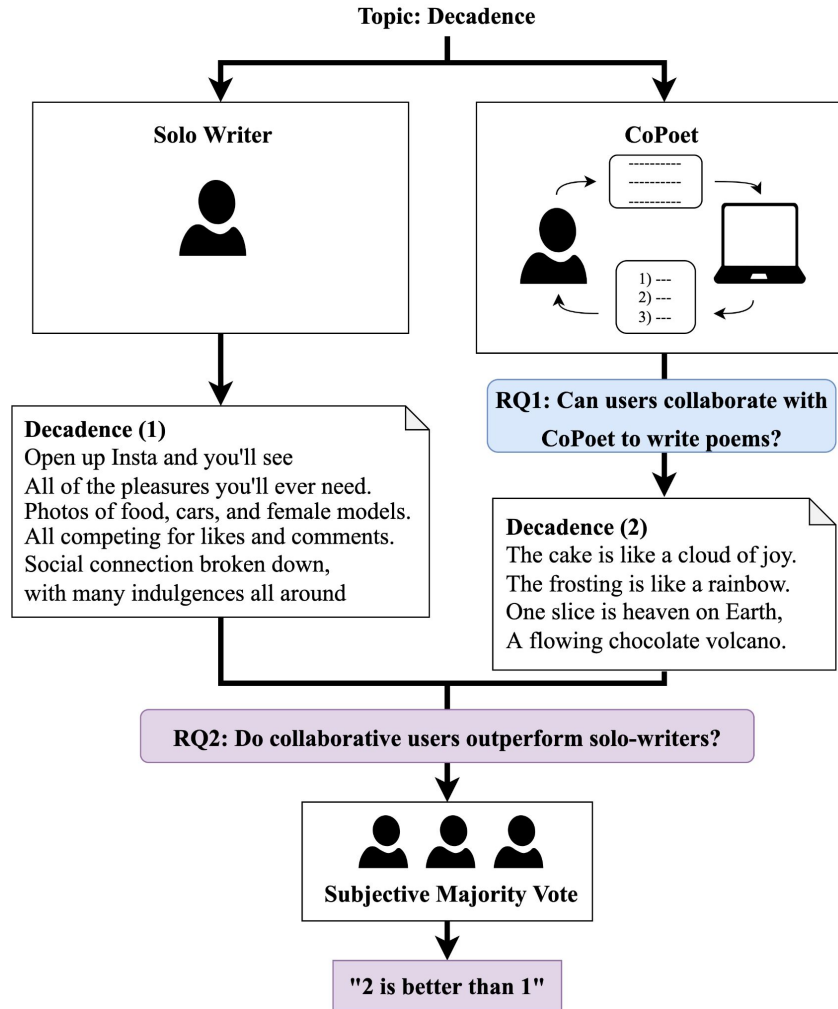
Submit

#### Find a rhyming word for:

Find Rhyming Words



# Overview of User Study



# Collaborative Writing Process

## Poem Draft



My dreams shatter like a glass ceiling

My dreams shatter like a glass ceiling  
As you walk away, dragging a broken  
shard

My dreams shatter like a glass ceiling  
As you walk away, dragging a broken  
shard  
I watch you leave, you no longer have  
meaning

My dreams shatter like a glass ceiling  
As you walk away, dragging a broken  
shard  
I watch you leave, you no longer have  
meaning  
As I lay there in the dark, my heart turns  
hard



## Interactions



Write a poetic sentence that contains  
the word 'dreams'!

- And he who, in the dreams of the night.
- And in your dreams, I am with you.



Write a simile about 'shattered  
dreams'!

- My dreams shattered like a glass bottle.
- My dreams shatter like a glass ceiling



Write a next sentence in a poetry given the previous  
sentence 'My dreams shatter like a glass ceiling

- As you walk away, dragging a broken shard.
- When I see your empty eyes.

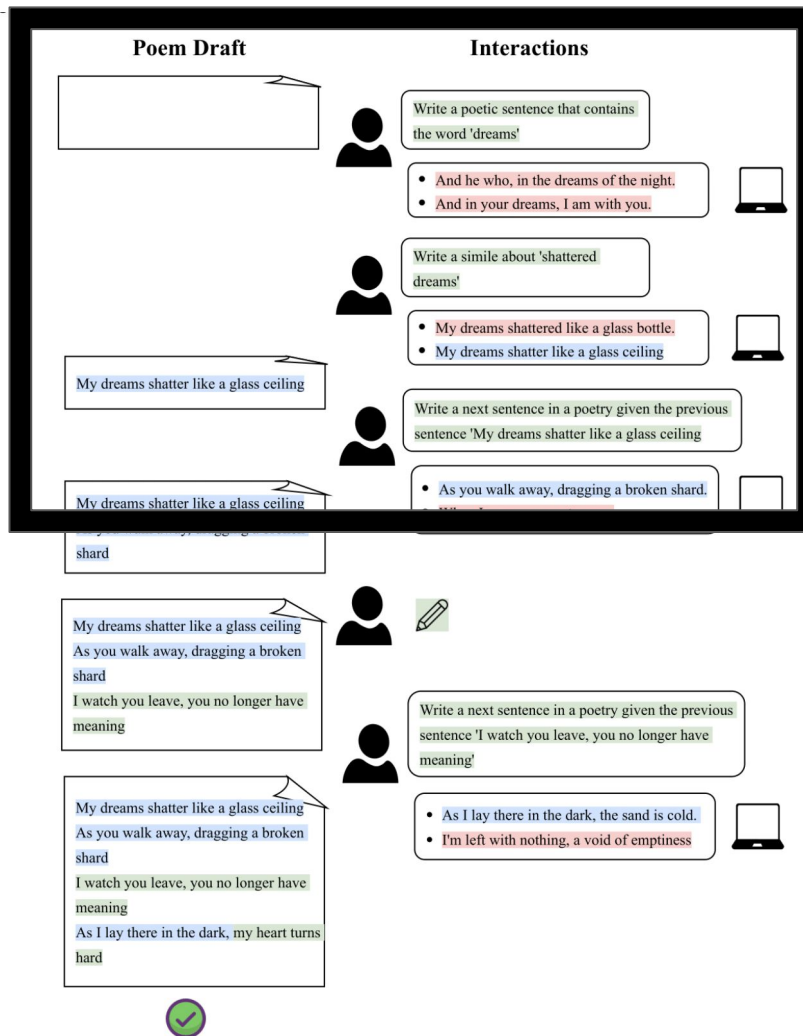


Write a next sentence in a poetry given the previous  
sentence 'I watch you leave, you no longer have  
meaning'!

- As I lay there in the dark, the sand is cold.
- I'm left with nothing, a void of emptiness



# Collaborative Writing Process



# Collaborative Writing Process

## Poem Draft



My dreams shatter like a glass ceiling

My dreams shatter like a glass ceiling  
As you walk away, dragging a broken  
shard

My dreams shatter like a glass ceiling  
As you walk away, dragging a broken  
shard  
I watch you leave, you no longer have  
meaning

My dreams shatter like a glass ceiling  
As you walk away, dragging a broken  
shard  
I watch you leave, you no longer have  
meaning  
As I lay there in the dark, my heart turns  
hard



## Interactions



Write a poetic sentence that contains  
the word 'dreams'!

- And he who, in the dreams of the night.
- And in your dreams, I am with you.



Write a simile about 'shattered  
dreams'!

- My dreams shattered like a glass bottle.
- My dreams shatter like a glass ceiling



Write a next sentence in a poetry given the previous  
sentence 'My dreams shatter like a glass ceiling

- As you walk away, dragging a broken shard.
- When I see your empty eyes.



Write a next sentence in a poetry given the previous  
sentence 'I watch you leave, you no longer have  
meaning'

- As I lay there in the dark, the sand is cold.
- I'm left with nothing, a void of emptiness



# Collaborative Poem Writing

## Another Day

The world has not yet awakened.

Darkness still creeps, but the day is not far.

Oh wait! there's the sun, and thus a solitary regret.

I still can't believe I haven't been to bed yet.

## Instructions:

- Write a poetic sentence that contains the word 'Morning'
- Write a simile about 'Night'
- **Write a poetic sentence that contains the word 'sun' and ending in a rhyme for 'yet'**
- Write a poetic sentence that contains the word 'Darkness' and ending in a rhyme for 'awakened'



# Collaborative Poem Writing

## Another Day

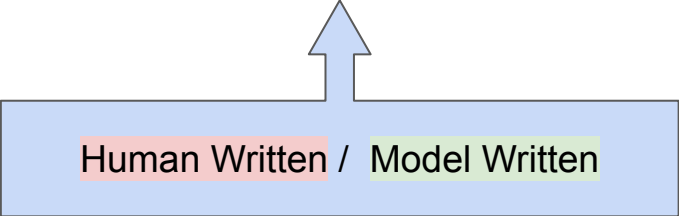
The world has not yet awakened.

Darkness still creeps, but the day is not far.

Oh wait! there's the sun, and thus a solitary regret.

I still can't believe I haven't been to bed yet.

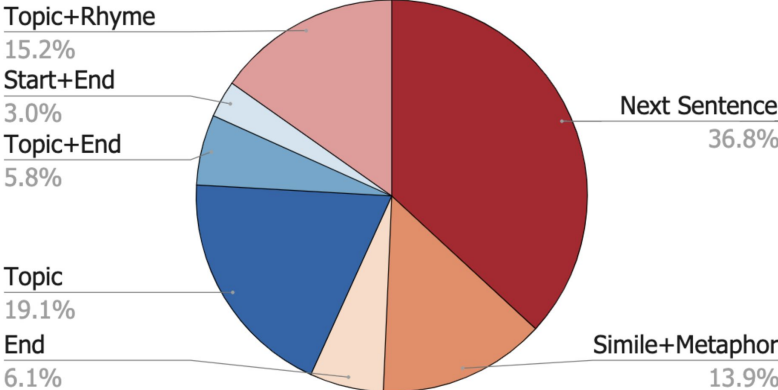
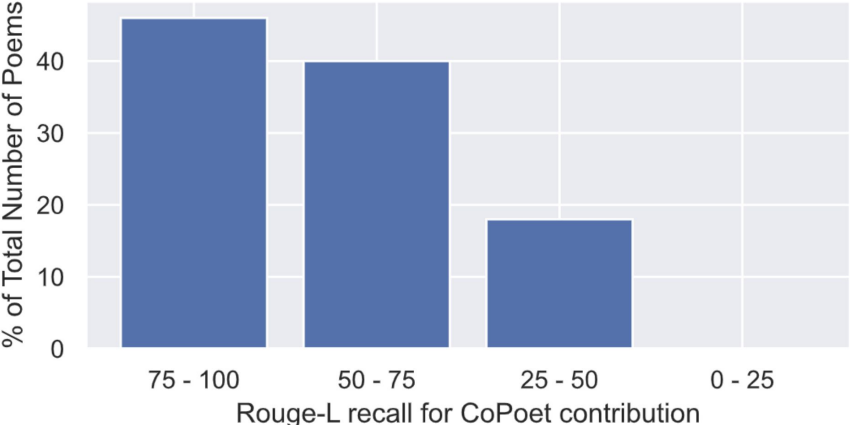
Human Written / Model Written



## Instructions:

- Write a poetic sentence that contains the word 'Morning'
- Write a simile about 'Night'
- **Write a poetic sentence that contains the word 'sun' and ending in a rhyme for 'yet'**
- Write a poetic sentence that contains the word 'Darkness' and ending in a rhyme for 'awakened'

# Do Users Find The Model Helpful?



# Do Users Write *Better* Poems With Model Help?

---

	Relevant %	Preferred %
Solo	96	43
Collaborative	98	57

---

---

	Preferred %	Not Preferred %
Diversity	63.0	37.0
Rhyme	72.5	27.5

---

# More Examples + Model Contributions



# Takeaways

- **Instruction tuning** can be an **effective way to help users write poems**

# Takeaways

- **Instruction tuning** can be an **effective way to help users write poems**
- InstructGPT3 is **pretty good at staying on subject** but still has **difficulty with more challenging generation instructions**

Write a poetic sentence  
that ends in 'happiness'









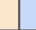



The key to happiness is to  
find what you love and  
hold onto it tight.



# Takeaways

- **Instruction tuning** can be an **effective way to help users write poems**
- InstructGPT3 is **pretty good at staying on subject** but still has **difficulty with more challenging generation instructions**

# Big Picture

- Many open questions bridging NLP and HCI:
  - How do we design the **most effective collaboration setup** to help human users? 
  - How do we train models to generate **helpful** suggestions? 
  - What is the **best form of user feedback** and how do we **incorporate it in model training**? 
  - How do we assist users in **content planning** for long form creative writing? 
  - How do we ensure **equitable creative writing assistance** to all users? 
- My Projects:
  - Machine-in-the-Loop Rewriting for Creative Image Captioning   
  - Collaborative Poetry writing with Instruction Tuning  



# Big Picture

- How do we train models to generate *helpful* suggestions?
  - Controlling stylistic attributes of text such as sentiment (hopefully going to ICML)
    - Also works on non-textual sequences like proteins

I had an amazing evening here! The food was yummy, the ambience was great and the service was so fast.  
Would recommend 10/10

I had a nice time. The food was good, and the ambience and service were fine. A decent night out!

The food was okay, the ambience and service were nothing to write home about. All in all it's fine for a one time visit but I don't think I'll be going back

# Big Picture

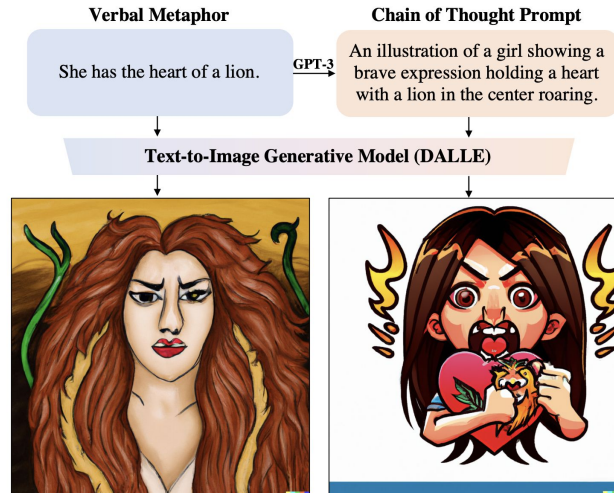
- What is the *best form of user feedback* and how do we *incorporate it in model training*?
  - Few-shot LLM personalization as an alternative to aggregation of feedback
  - Machine Teaching
    - Can we train models that can perform the content selection + presentation task needed to help human students?

# Big Picture

- How do we design the ***most effective collaboration setup*** to help human users?
  - Providing assistance in more specialised domains (medical texts)

# Big Picture

- How do we design the **most effective collaboration setup** to help human users?
  - Multimodal Creativity - Check out Tuhin's work on [Visual Metaphors](#) :)



# Thank You

# Backup Slides

# How Well Do Models Compose Instructions?

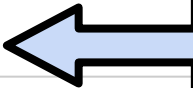
		<b>T5 - 11B</b>	<b>T5 - 3B</b>	<b>T0 - 3B</b>	<b>InstructGPT- ZS (175B)</b>	<b>InstructGPT - FS (175B)</b>	<b>T0pp</b>
<b>Subject (55)</b>	% - Match	76.36%	60%	54.54%	72.72%	<b>87.87%</b>	65.45%
	% - Match w/ Ending (34)	<b>47.05%</b>	41.17%	38.23%	26.31%	29.41%	29.41%
<b>Rhyme (16)</b>	% - Match -w & Rhyme Success Rate	<b>43.75%</b>	25.00%	37.50%	26.31%	37.50%	0.00%
<b>Simile (4)</b>	% Subject + Comparator	25.00%	50.00%	0.00%	66.66%	<b>75%</b>	25.00%
	% Comparator	<b>100.00%</b>	<b>100.00%</b>	<b>100.00%</b>	83.33%	<b>100%</b>	25.00%
<b>Metaphor (4)</b>	% Subject + Comparator	<b>50.00%</b>	<b>50.00%</b>	25%	25%	25%	25.00%
<b>Haiku (5)</b>	% Subject + [15-19] Syllables	<b>60.00%</b>	20.00%	<b>60%</b>	0%	20%	0.00%

# How Well Do Models Compose Instructions?

		T5 - 11B	T5 - 3B	T0 - 3B	InstructGPT- ZS (175B)	InstructGPT - FS (175B)	T0pp
<b>Subject (55)</b>	% - Match	76.36%	60%	54.54%	72.72%		65.45%
	% - Match w/ Ending (34)	<b>47.05%</b>	41.17%	38.23%	26.31%		29.41%
<b>Rhyme (16)</b>	% - Match -w & Rhyme Success Rate	<b>43.75%</b>	25.00%	37.50%			0.00%
<b>Simile (4)</b>	% Subject + Comparator	25.00%	50.00%	0.00%	66.66%		25.00%
	% Comparator	<b>100.00%</b>	<b>100.00%</b>	<b>100.00%</b>	83.33%		25.00%
<b>Metaphor (4)</b>	% Subject + Comparator	<b>50.00%</b>	<b>50.00%</b>	25%	25%	25%	25.00%
<b>Haiku (5)</b>	% Subject + [15-19] Syllables	<b>60.00%</b>	20.00%	<b>60%</b>	0%	20%	0.00%

Composition improves with model size

T5 outperforms T0 again





# How Well Do Models Compose Instructions?

		<b>T5 - 11B</b>	<b>T5 - 3B</b>	<b>T0 - 3B</b>	<b>InstructGPT- ZS (175B)</b>	<b>InstructGPT - FS (175B)</b>	<b>T0pp</b>
<b>Subject (55)</b>	% - Match	76.36%	60%	54.54%	72.72%	87.87%	65.45%
	% - Match w/ Ending (34)	47.05%	41.17%		26.31%	29.41%	29.41%
<b>Rhyme (16)</b>	% - Match -w & Rhyme Success Rate	43.75%	25.00%		26.19%	37.50%	0.00%
<b>Simile (4)</b>	% Subject + Comparator	25.00%	50.00%		66.66%	75%	25.00%
	% Comparator	100.00%	100.00%		83.33%	100%	25.00%
<b>Metaphor (4)</b>	% Subject + Comparator	50.00%	50.00%		25%	25%	25.00%
<b>Haiku (5)</b>	% Subject + [15-19] Syllables	60.00%	20.00%	60%	0%	20%	0.00%

Comparable performance with T5 edging it on challenging instructions

# Scalar Controlled Text Generation

w/ Richard Pang, He He and Ankur Parikh



ML<sup>2</sup>

Machine Learning  
for Language

# Motivation

- As opposed to controlling the output of text with literary devices or instructions, writers might want to control scalar attributes
  - Sentiment Control, Toxicity etc.

# Motivation

- As opposed to controlling the output of text with literary devices or instructions, writers might want to control scalar attributes
- An example of sentiment control between a **positive**, **neutral** and **slightly negative** version of the same sentence

I had an amazing evening here! The food was yummy, the ambience was great and the service was so fast.  
Would recommend 10/10

I had a nice time. The food was good, and the ambience and service were fine. A decent night out!

The food was okay, the ambience and service were nothing to write home about. All in all it's fine for a one time visit but I don't think I'll be going back

# Problem Setup

- Assume we have an oracle scorer  $f_s$ 
  - Maps from an input sequence to the range of the score
- Given an input text  $x$  and a target score  $s_t$
- **Goal: Generate  $x'$  s.t.  $f_s(x') = s_t$**

# Approach

- **Generate  $x'$  s.t.  $f_s(x') = s_t$  iteratively**
  - First generate  $x'_i = x_{i-1} + \partial f_s / \partial x$
  - Increase the number of iterations “ $i$ ” in order to achieve higher/lower target scores
  - Allows generalization to OOD target scores

# Approach

- Generate  $x'$  s.t.  $f_s(x') = s_t$  iteratively
  - First generate  $x'_i = x_{i-1} + \partial f_s / \partial x$
  - Increase the number of iterations “ $i$ ” in order to achieve higher/lower target scores
  - Allows generalization to OOD target scores
  - **How do we train a model to perform this?**

# Approach

- Generate  $x'$  s.t.  $f_s(x') = s_t$  iteratively
  - Update  $x'_i = x_{i-1} + \partial f_s / \partial x$
- **Learn to approximate  $\partial f_s / \partial x$  from Perturbations**
  - Create paired data using masking + infilling
  - Learn Control Tags with LLMs to edit the text and move up and down this scale



# Learning from Perturbations

Example:

- **Source:** "<dec> The desserts come with the **easy to-believe claim that they contain just** under 200 calories.", **Score: 2.607**
- **Target:** "The desserts come with the **hard-to-believe claim that they are all** under 200 calories.", **Score: 2.132**

# Example

Trained T5-base on the created dataset

X = “Top notch doctor in a top notch practise”,  $s(x) = 4.904$ ,  $s_t = 3.0$

Iteration	Text	Sentiment Score
1	<dec> Top notch doctor in a top notch practice.	4.904
2	<dec> A top notch doctor is in the practice.	4.701
3	<dec> The practice is staffed by a top notch <b>hygenist</b> .	4.003
4	<dec> The practice is managed by a top notch hygenist.	3.647
5	The practice is run by a very good hygenist.	2.955

# Experiments

- For sentiment analysis, scorer model is linear regression classifier
  - Score range: [0, 4]
- Select only training data in [1, 3] range

# Experiments

- For sentiment analysis, scorer model is linear regression classifier
  - Score range: [0, 4]
- Select only training data in [1, 3] range
- Evaluation:
  - Given a source sentence, alter it to two separate target scores
  - Report success rate of reaching the target
    - Very In-Domain - Source score +/- 0.3
    - In-Domain - {1.5, 2.5}
    - Out-of-Domain - {0.5, 3.5}

# Results

Test set size = 1831		VID	ID	OOD
<b>With Scorer</b>	Our Model (small, $n_{\text{iter}} = 10$ , $n_{\text{seq}} = 5$ )	0.965	<b>0.971</b>	0.429
	Our Model (large, $n_{\text{iter}} = 10$ , $n_{\text{seq}} = 5$ )	0.938	0.930	<b>0.623</b>
	Genhance ( $d_z = 0.15$ , $n = 1$ )	0.377	0.287	0.04375
	Genhance ( $d_z = 0.15$ , $n = 50$ )	0.9515	0.9075	0.3865
	Genhance ( $d_z = 0.15$ , $n = 100$ )	<b>0.9775</b>	0.9535	0.54965

# Metrics of Evaluation

Test set size = 1831		VID	ID	OOD
With Scorer	Output (d_z = 0.15, n_seq = 10, n = 1)	0.965	<b>0.971</b>	0.429
	Output (d_z = 0.15, n_seq = 10, n = 5)	0.938	0.930	<b>0.623</b>
	Output (d_z = 0.15, n_seq = 10, n = 10)	0.377	0.287	0.04375
	Output (d_z = 0.15, n_seq = 10, n = 50)	0.9515	0.9075	0.3865
	Genhance (d_z = 0.15, n = 100)	<b>0.9775</b>	0.9535	0.54965

Each value in these columns is the success rate at achieving various target attribute values for different models



# Comparison to Baseline Genhance Model

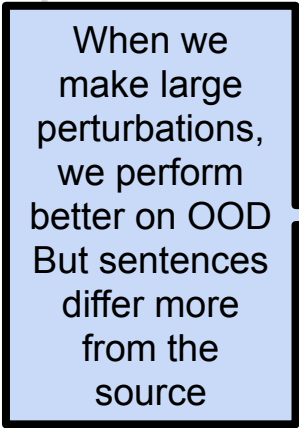
Test set size = 1831		VID	ID	OOD
	Our Model (small, $n_{\text{iter}} = 10$ , $n_{\text{seq}} = 5$ )	0.965	<b>0.971</b>	0.429
	Our Model (large, $n_{\text{iter}} = 10$ , $n_{\text{seq}} = 5$ )	0.938	0.930	<b>0.623</b>
	Genhance ( $d_z = 0.15$ , $n = 1$ )	0.377	0.287	0.04375
	Genhance ( $d_z = 0.15$ , $n = 50$ )	0.9515	0.9075	0.3865
	Genhance ( $d_z = 0.15$ , $n = 100$ )	<b>0.9775</b>	0.9535	0.54965

We outperform a controlled generation baseline on ID and OOD



With Score

# Trade-off based on Perturbation Size

	1831	VID	ID	OOD
<p>When we make large perturbations, we perform better on OOD. But sentences differ more from the source.</p> 	Our Model (small, n_iter = 10, n_seq = 5)	0.965	<b>0.971</b>	0.429
	Our Model (large, n_iter = 10, n_seq = 5)	0.938	0.930	<b>0.623</b>
	Genhance (d_z = 0.15, n = 1)	0.377	0.287	0.04375
	Genhance (d_z = 0.15, n = 50)	0.9515	0.9075	0.3865
	<b>With Scorer</b> Genhance (d_z = 0.15, n = 100)	<b>0.9775</b>	0.9535	0.54965



# Fine-Grained Comparison

The success rate is a coarse metric. We want to examine the average case.

Given all the source examples, how much are we able to change the score of each?

# Fine-Grained Comparison

The success rate is a coarse metric. We want to examine the average case.

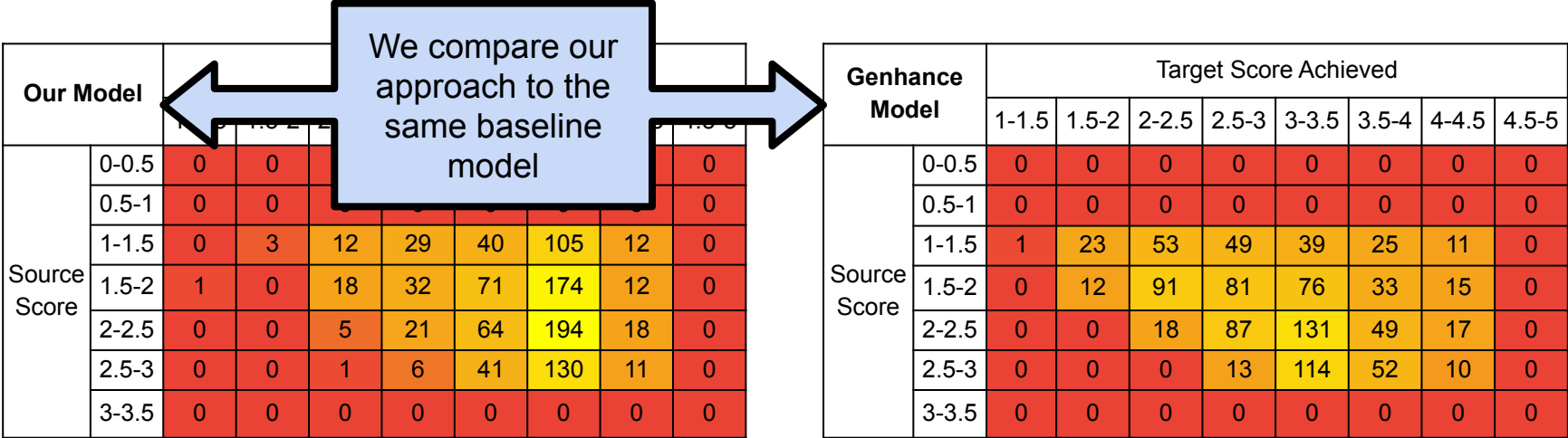
Given all the source examples, how much are we able to change the score of each?

Our Model		Target Score Achieved							
		1-1.5	1.5-2	2-2.5	2.5-3	3-3.5	3.5-4	4-4.5	4.5-5
Source Score	0-0.5	0	0	0	0	0	0	0	0
	0.5-1	0	0	0	0	0	0	0	0
	1-1.5	0	3	12	29	40	105	12	0
	1.5-2	1	0	18	32	71	174	12	0
	2-2.5	0	0	5	21	64	194	18	0
	2.5-3	0	0	1	6	41	130	11	0
	3-3.5	0	0	0	0	0	0	0	0

Genhance Model		Target Score Achieved							
		1-1.5	1.5-2	2-2.5	2.5-3	3-3.5	3.5-4	4-4.5	4.5-5
Source Score	0-0.5	0	0	0	0	0	0	0	0
	0.5-1	0	0	0	0	0	0	0	0
	1-1.5	1	23	53	49	39	25	11	0
	1.5-2	0	12	91	81	76	33	15	0
	2-2.5	0	0	18	87	131	49	17	0
	2.5-3	0	0	0	13	114	52	10	0
	3-3.5	0	0	0	0	0	0	0	0

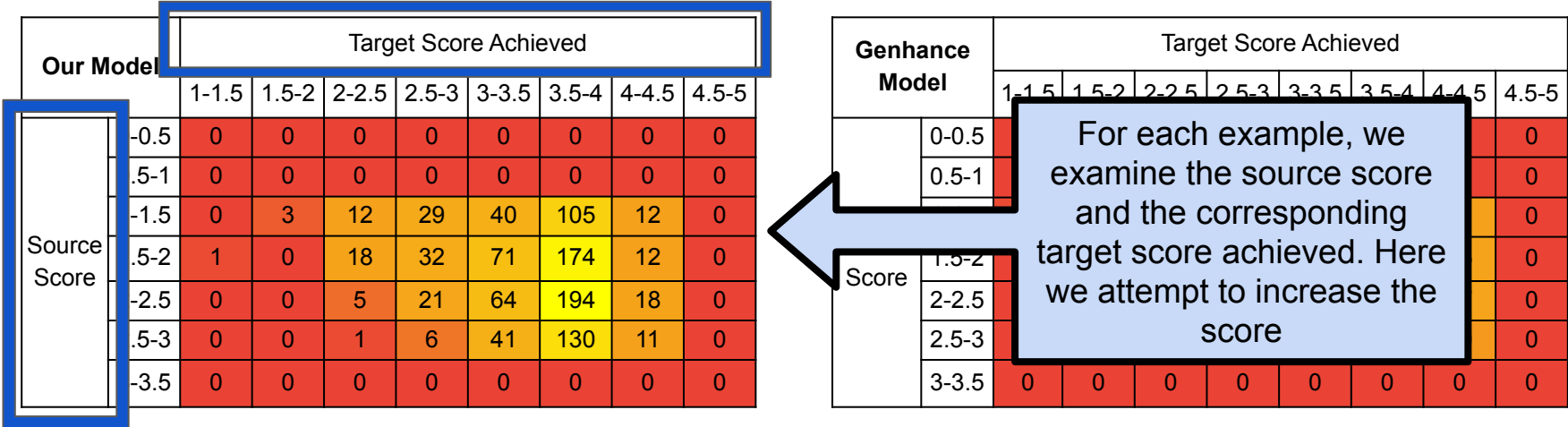
# Fine-Grained Comparison

Given all the source examples, how much are we able to change the score of each?



# Fine-Grained Comparison

Given all the source examples, how much are we able to change the score of each?



# Fine-Grained Comparison

Given all the source examples, how much are we able to change the score of each?

Our Model		Target Score Achieved							
		1-1.5	1.5-2	2-2.5	2.5-3	3-3.5	3.5-4	4-4.5	4.5-5
Source Score	0-0.5	0	0	0	0	0	0	0	0
	0.5-1	0	0	0	0	0	0	0	0
	1-1.5	0	3	12	29	40	105	18	0
	1.5-2	1	0	18	32	71	174	12	0
	2-2.5	0	0	5	21	64	194	18	0
	2.5-3	0	0	1	6	41	130	11	0
	3-3.5	0	0	0	0	0	0	0	0

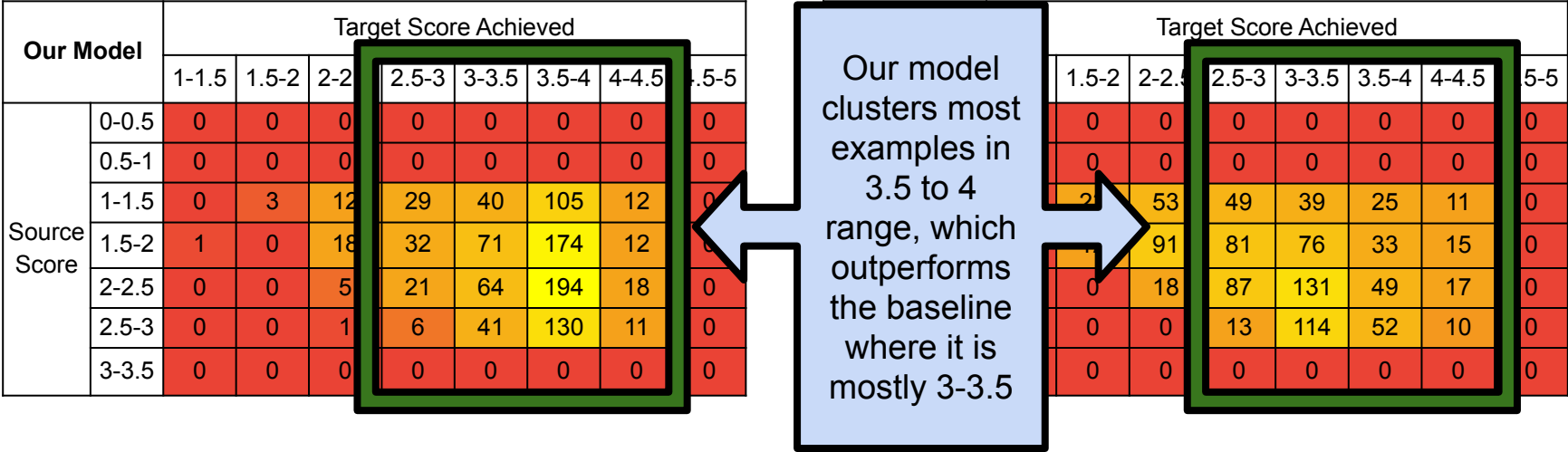
  

Genhance Model		Target Score Achieved							
		1-1.5	1.5-2	2-2.5	2.5-3	3-3.5	3.5-4	4-4.5	4.5-5
Source Score	1-1.5	0	0	0	0	0	0	0	0
	1.5-2	0	0	0	0	0	0	0	0
	2-2.5	0	0	0	13	114	52	10	0
	3-3.5	0	0	0	0	0	0	0	0

This cell means that there were **105 examples** with **source score in 1-1.5** that achieved **target score between 3.5 and 4**

# Fine-Grained Comparison

Given all the source examples, how much are we able to change the score of each?



# Fine-Grained Comparison

Our model is able to better shift the distribution of scores of examples in the desired direction.

Our Model		Target Score Achieved							
		1-1.5	1.5-2	2-2.5	2.5-3	3-3.5	3.5-4	4-4.5	4.5-5
Source Score	0-0.5	0	0	0	0	0	0	0	0
	0.5-1	0	0	0	0	0	0	0	0
	1-1.5	0	3	12	29	40	105	12	0
	1.5-2	1	0	18	32	71	174	12	0
	2-2.5	0	0	5	21	64	194	18	0
	2.5-3	0	0	1	6	41	130	11	0
	3-3.5	0	0	0	0	0	0	0	0

Genhance Model		Target Score Achieved							
		1-1.5	1.5-2	2-2.5	2.5-3	3-3.5	3.5-4	4-4.5	4.5-5
Source Score	0-0.5	0	0	0	0	0	0	0	0
	0.5-1	0	0	0	0	0	0	0	0
	1-1.5	1	23	53	49	39	25	11	0
	1.5-2	0	12	91	81	76	33	15	0
	2-2.5	0	0	18	87	131	49	17	0
	2.5-3	0	0	0	13	114	52	10	0
	3-3.5	0	0	0	0	0	0	0	0

# Takeaways

- We're able to achieve better controlled generation than baselines, particularly towards OOD target attribute values
- Our approach learns the distribution of the data and trains a model to move up and down the scale of an attribute along the data distribution.



# Next Steps

- Nothing that we do is specific to text!
- We're trying to generate protein sequences where we control the attribute of the stability of the molecule.
  - Not possible in other approaches that require a differentiable scoring function

# Collaborative Poem Writing

## Sad Reality

No one prepares you for growing older

**Aging is a symptom of your dreams  
foreclosure**

You're not a musician, but you smoke like them

Hurry up, your lunch is over in 10

## Instructions:

- Write a poetic sentence that contains the word 'Growing up'
- **Write a metaphor about 'Aging'**
- Write a next sentence in a poetry given the previous sentence 'Aging is a symptom of your dreams fading'
- Write a next sentence in a poetry given the previous sentence 'You're not a musician, but you smoke like them'

# User Interface

Describe the image below



Enter your text

The majestic lion lies poised and ready to strike, an [awesome sight to see].

Suggest

Finish

Char count: 77

Select the suggestion that you like best: 1  2  3  Original Text



### Suggestion 1

The majestic lion lies poised and ready to strike, an **alert** predator eyeing his **prey**



### Suggestion 2

The majestic lion lies poised and ready to strike, an **avenger**.



### Suggestion 3

The majestic lion lies poised and ready to strike, an **abyss waiting**.

# Error Analysis

- Most common error case is content drift
  - Model changes the meaning of the sentence when rewriting
- Copying of the source text verbatim
- Repetition in generated text



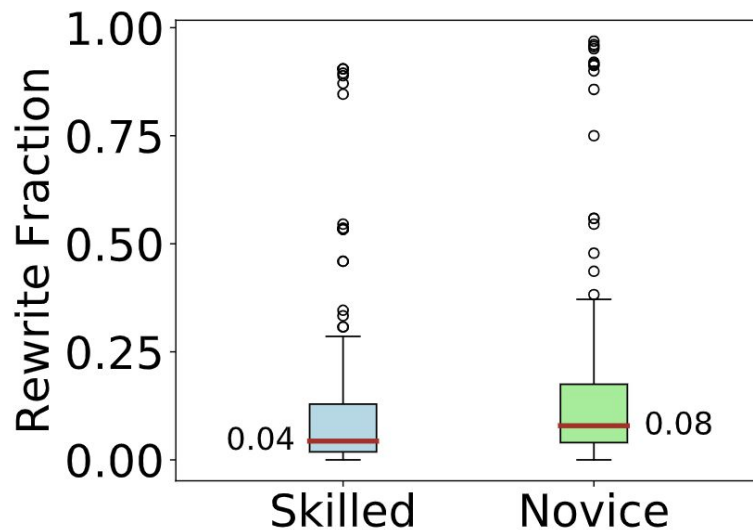
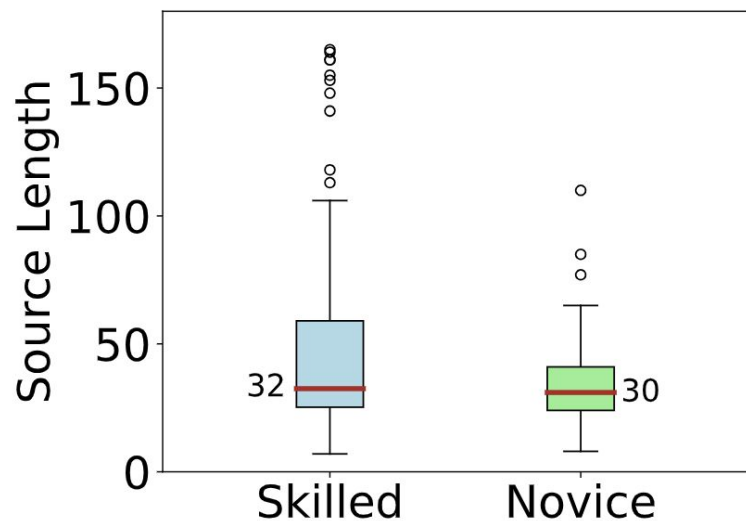
In front of a wall, a girl with blonde hair is on her hands who seems to be **[ coming out of a magical door ]**

In front of a wall, a girl with blonde hair is on her hands who seems to be **laughing out loud**



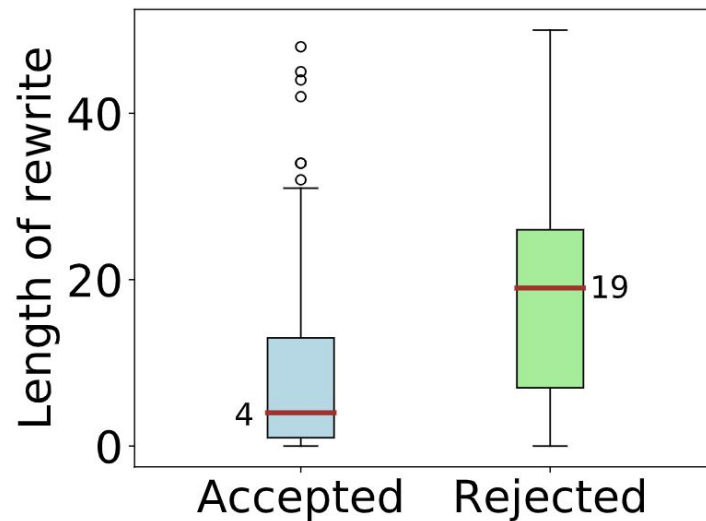
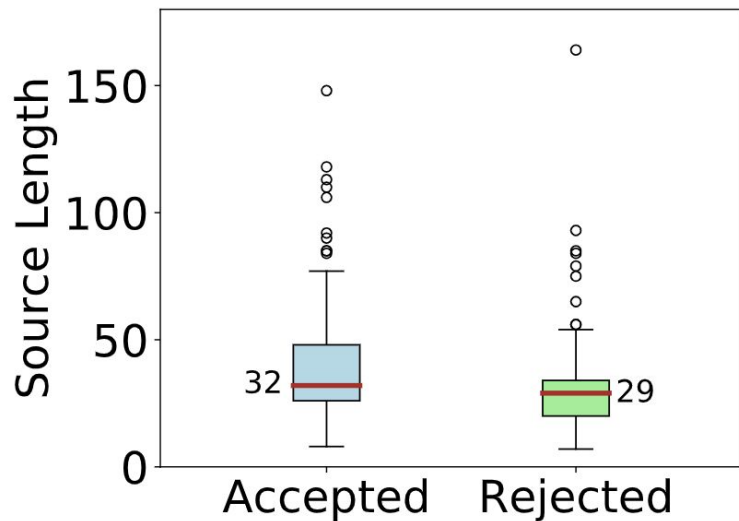
# When is the Model Effective?

Skilled writers tend to write longer sentences and request shorter fractions to be rewritten



# When is the Model Effective?

Shorter rewrites in longer sentences tend to be accepted more



# Automatic Evaluation - Unknown Entities

		<b>T5 - 11B</b>	<b>T5 - 3B</b>	<b>T0 - 3B</b>	<b>InstructGPT- ZS (175B)</b>	<b>InstructGPT - FS (175B)</b>	<b>T0pp</b>
<b>Subject (31)</b>	% - Match	<b>80.64%</b>	74.19%	77.41%	69.23%	74.19%	51.61%
	% - Match w/ Ending (2)	<b>100%</b>	<b>100%</b>	<b>100%</b>	0%	0.00%	0.00%
<b>Rhyme (11)</b>	Success Rate	<b>36.36%</b>	36.36%	9.09%	9.09%	18.18%	9.09%
<b>Simile (9)</b>	% - Subject + Comparator	33.33%	44.44%	33.33%	11.11%	<b>55.55%</b>	11.11%
	% - Comparator	<b>100.00%</b>	<b>100.00%</b>	<b>100.00%</b>	66.66%	88.88%	22.22%
<b>Metaphor (7)</b>	% - Subject + Comparator	<b>100.00%</b>	<b>100.00%</b>	<b>100.00%</b>	85.71%	<b>100.00%</b>	28.57%
<b>Haiku (7)</b>	% Subject + (15-19) Syllables	<b>71.42%</b>	<b>71.42%</b>	57.14%	57.14%	42.85%	0.00%

# Automatic Evaluation - Unknown Entities

		T5 - 11B	T5 - 3B	T0 - 3B	InstructGPT- ZS (175B)	InstructGPT - FS (175B)	T0pp
<b>Subject (31)</b>	% - Match	<b>80.64%</b>	74.19%	77.41%	69.23%	74.19%	51.61%
	% - Match w/ Ending (2)	<b>100%</b>	<b>100%</b>	66.67%	66.67%	0.00%	0.00%
<b>Rhyme (11)</b>	Success Rate	<b>36.36%</b>	36.36%	36.36%	36.36%	18.18%	9.09%
<b>Simile (9)</b>	% - Subject + Comparator	33.33%	33.33%	33.33%	33.33%	<b>55.55%</b>	11.11%
	% - Comparator	<b>100.00%</b>	<b>100.00%</b>	100.00%	66.67%	88.88%	22.22%
<b>Metaphor (7)</b>	% - Subject + Comparator	<b>100.00%</b>	<b>100.00%</b>	100.00%	100.00%	<b>100.00%</b>	28.57%
<b>Haiku (7)</b>	% Subject + (15-19) Syllables	<b>71.42%</b>	<b>71.42%</b>	57.14%	57.14%	42.85%	0.00%

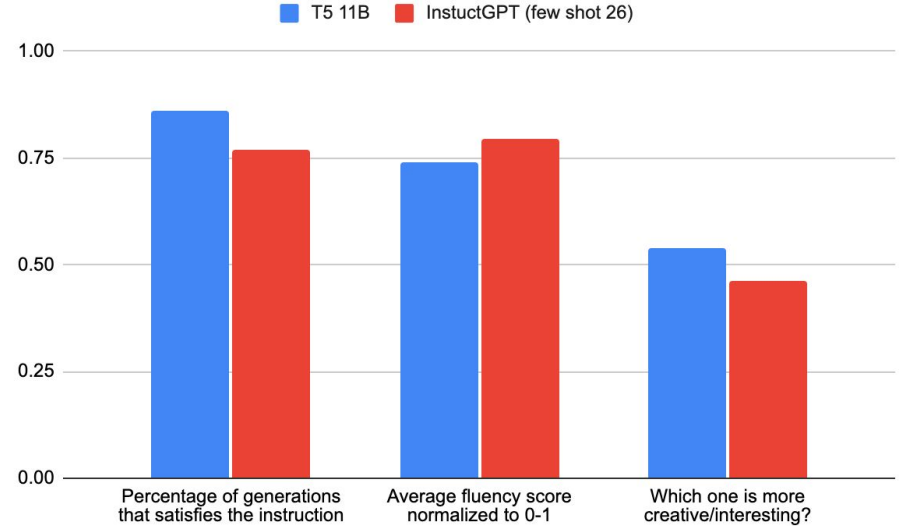
Similar trends as with known entities with lower overall performance



# Human Evaluation

## Known Entities

	T5 11B	InstructGPT
<b>Percentage that satisfies the instruction</b>	<b>0.862</b>	0.769
Fluency, on a scale of 1-5	3.697	<b>3.969</b>
Which one is more creative/interesting?	<b>0.538</b>	0.462



# Human Evaluation

## Known Entities

	T5 11B	InstructGPT
Percentage that satisfies the instruction	0.862	0.769
Fluency, on a scale of 1-5	3.697	<b>3.969</b>
Which one is more creative/interesting?	<b>0.538</b>	0.462

## Unknown Entities

	T5 11B	InstructGPT
Percentage that satisfies the instruction	0.925	0.865
Fluency, on a scale of 1-5	3.865	<b>3.905</b>
Which one is more creative/interesting?	<b>0.567</b>	0.433

InstructGPT is more fluent

# Human Evaluation - Compositional Instructions

## Compositional Test Set

	T5 11B	InstructGPT
<b>Percentage that satisfies the instruction</b>	<b>0.776</b>	0.552
Fluency, on a scale of 1-5	3.483	<b>3.756</b>
Which one is more creative/interesting?	0.477	<b>0.523</b>

# Collaborative Poem Writing

## The harshness of time.

Time is a very harsh mistress all bitter and cold.

Time is never ending everlasting and very bold.

Time will elapse you in moments that matter.

**Time will pass and I'll get older and fatter.**

## Instructions:

- Write a poetic sentence that contains the word 'Time'
- Write a next sentence in a poetry given the previous sentence 'Time is my horse that stays always with me.'
- Write a poetic sentence that contains the word 'time' and ending in a rhyme for 'me'
- Write a poetic sentence that contains the word 'flow' and ending in a rhyme for 'me'
- **Write a poetic sentence that contains the word 'Time' and ending in a rhyme for 'matter'**
- Write a poetic sentence that contains the word 'time' and ending in a rhyme for 'cold'

# How Well Do Models Follow Instructions?

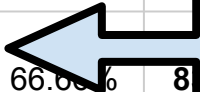
Known Templates		T5 - 11B	T5 - 3B	T0 - 3B	InstructGPT- Zero Shot (175B)	InstructGPT - Few Shot (175B)	T0pp (11B)
<b>Subject (51)</b>	% - Match	92.15%	92.15%	86.27%	86.27%	<b>96.07%</b>	84.31%
	% - Match w/ Ending (22)	<b>95.45%</b>	<b>95.45%</b>	86.36%	13.63%	18.18%	40.90%
<b>Rhyme (14)</b>	Success Rate	78.57%	<b>85.71%</b>	<b>85.71%</b>	57.14%	71.42%	0.00%
<b>Simile (6)</b>	% Subject + Comparator	66.66%	<b>83.33%</b>	66.66%	83.33%	66.66%	16.66%
	% Comparator	<b>100.00%</b>	<b>100%</b>	<b>100.00%</b>	100%	83.33%	16.66%
<b>Metaphor (5)</b>	% Subject + Comparator	<b>100.00%</b>	<b>100.00%</b>	<b>100.00%</b>	80%	<b>100%</b>	0.00%
<b>Haiku (5)</b>	% Subject + (15-19) Syllables	20.00%	40.00%	0.00%	<b>80%</b>	40%	0.00%

# How Well Do Models Follow Instructions?

Known Templates		T5 - 11B	T5 - 3B	T0 - 3B	InstructGPT- Zero Shot (175B)	InstructGPT - Few Shot (175B)	T0pp (11B)
Subject (51)	% - Match	92.15%	90.00%	90.00%	86.27%	<b>96.07%</b>	84.31%
	% - Match w/ Ending (22)	<b>95.45%</b>	90.00%	90.00%	13.63%	18.18%	40.90%
Rhyme (14)	Success Rate	78.57%	80.00%	80.00%	57.14%	71.42%	0.00%
Simile (6)	% Subject + Comparator	66.66%	80.00%	80.00%	33.33%	66.66%	16.66%
	% Comparator	<b>100.00%</b>	100.00%	100.00%	100%	83.33%	16.66%
Metaphor (5)	% Subject + Comparator	<b>100.00%</b>	100.00%	100.00%	80%	<b>100%</b>	0.00%
Haiku (5)	% Subject + (15-19) Syllables	20.00%	40.00%	0.00%	<b>80%</b>	40%	0.00%

**Hand Crafted Metrics** for each kind of instruction

These are **soft metrics!**



# How Well Do Models Follow Instructions?

Known Templates		T5 - 11B	T5 - 3B	T0 - 3B	InstructGPT- Zero Shot (175B)	InstructGPT - Few Shot (175B)	T0pp (11B)
Subject (51)	% - Match	92.15%	92.15%	86.27%	86.27%		
	% - Match w/ Ending (22)	<b>95.45%</b>	<b>95.45%</b>	86.36%	13.63%		
Rhyme (14)	Success Rate	78.57%	<b>85.71%</b>	<b>85.71%</b>	57.14%		
Simile (6)	% Subject + Comparator	66.66%	<b>83.33%</b>	66.66%	83.33%		
	% Comparator	<b>100.00%</b>	<b>100%</b>	<b>100.00%</b>	100%		
Metaphor (5)	% Subject + Comparator	<b>100.00%</b>	<b>100.00%</b>	<b>100.00%</b>	80%		
Haiku (5)	% Subject + (15-19) Syllables	20.00%	40.00%	0.00%	<b>80%</b>	40%	0.00%

**T5 largely outperforms T0**  
 ↓  
 Transfer of Instruction Tuning to new/unrelated tasks



# How Well Do Models Follow Instructions?

Known Templates		T5 - 11B	T5 - 3B	T0 - 3B	InstructGPT- Zero Shot (175B)	InstructGPT - Few Shot (175B)	T0pp (11B)
			15%	86.27%	86.27%	<b>96.07%</b>	84.31%
Subject (51)			45%	86.36%	13.63%	18.18%	40.90%
Rhyme (14)			71%	<b>85.71%</b>	57.14%	71.42%	0.00%
			33%	66.66%	83.33%	66.66%	16.66%
Simile (6)			0%	<b>100.00%</b>	100%	83.33%	16.66%
Metaphor (5)			0%	<b>100.00%</b>	80%	<b>100%</b>	0.00%
Haiku (5)	% Subject (15-19) Syllables	20.00%	40.00%	0.00%	<b>80%</b>	40%	0.00%

**InstructGPT performance improves in few-shot setting.**

**Both outperform T0pp**

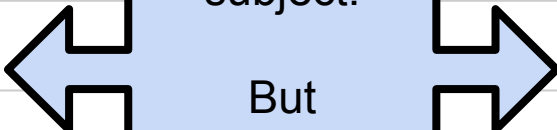


# How Well Do Models Follow Instructions?

Known Templates		T5 - 11B	T5 - 3B	InstructGPT - 175B	InstructGPT - Few Shot (175B)	T0pp (11B)	
Subject (51)	% - Match	92.15%	92.15%	97%	96.07%	84.31%	
	% - Match w/ Ending (22)	95.45%	95.45%	93%	18.18%	40.90%	
Rhyme (14)	Success Rate	78.57%			71.42%	0.00%	
Simile (6)	% Subject + Comparator	66.66%	83.33%	63%	66.66%	16.66%	
	% Comparator	100.00%	100%	67%	83.33%	16.66%	
Metaphor (5)	% Subject + Comparator	100.00%	100.00%	60%	100%	0.00%	
Haiku (5)	% Subject + (15-19) Syllables	20.00%	40.00%	0.00%	80%	40%	0.00%

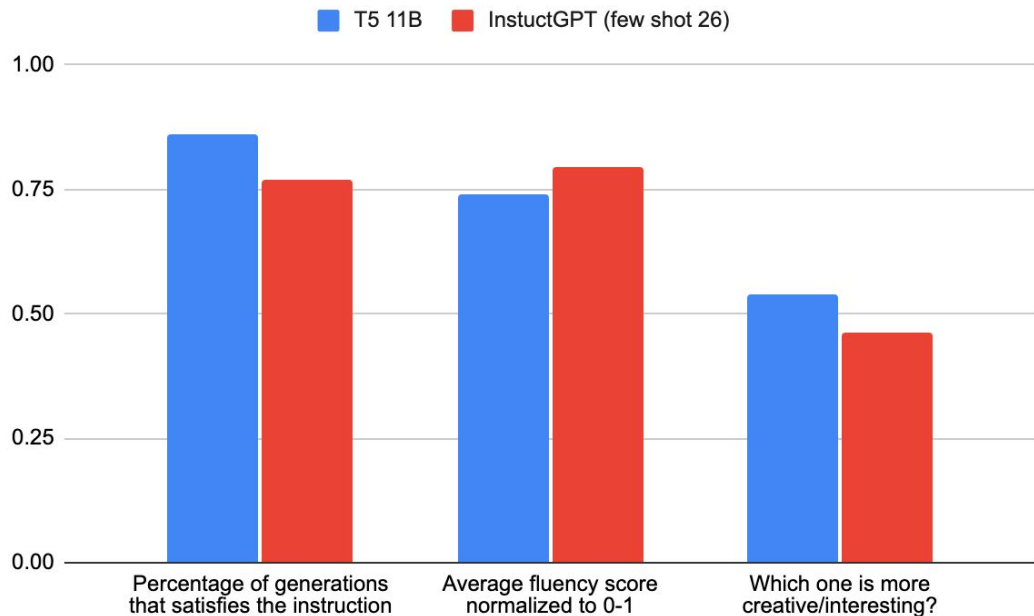
InstructGPT is good at staying on subject.

But performance drops on harder instructions



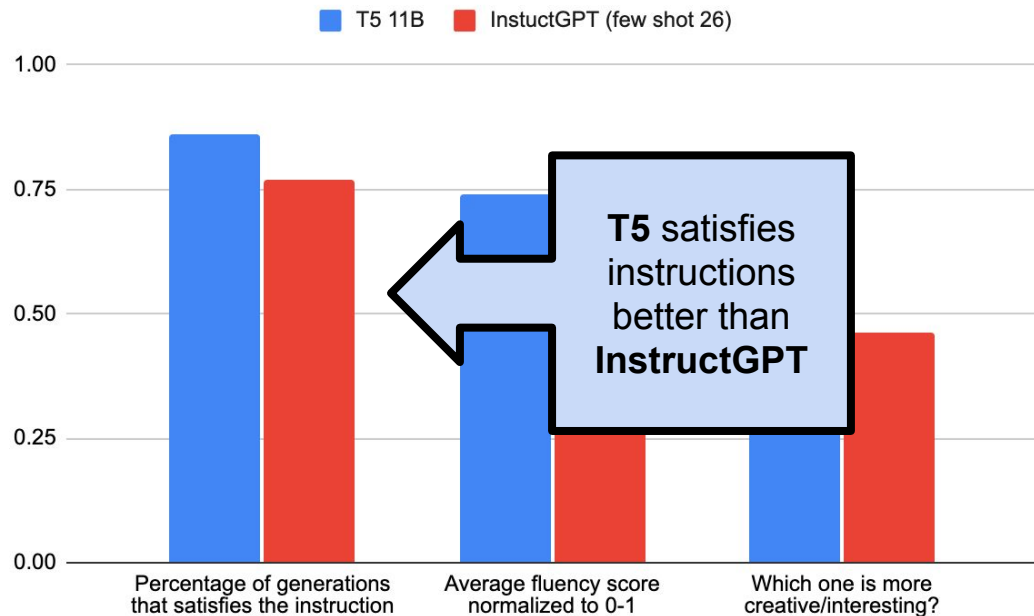
# How Well Do Models Follow Instructions?

- Sample model generations for each instruction
- Majority vote from 3 human annotators on each of the following axes



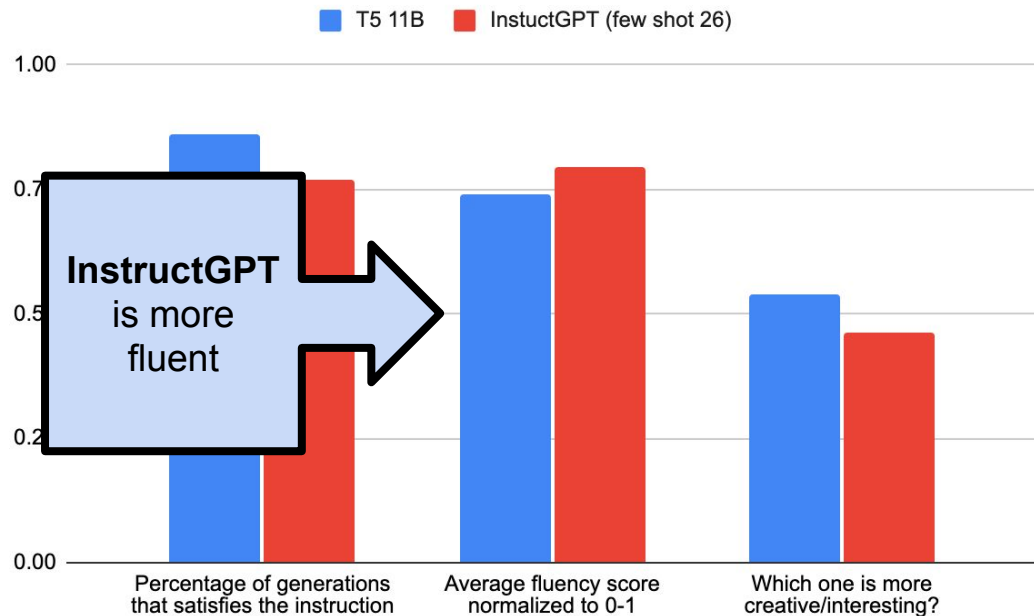
# How Well Do Models Follow Instructions?

- Sample model generations for each instruction
- Majority vote from 3 human annotators on each of the following axes



# How Well Do Models Follow Instructions?

- Sample model generations for each instruction
- Majority vote from 3 human annotators on each of the following axes



# Research Questions

- Can we train LLMs to satisfy creative writing instructions for poetry writing tasks?
- **Can models compose instructions seen at train time in unseen combinations?**
- Can we help users complete creative writing tasks using natural language instructions?

# How Well Do Models Compose Instructions?

