

# Towards Emotionally-Intelligent AI Systems

Hongli Zhan

Ph.D. Defense

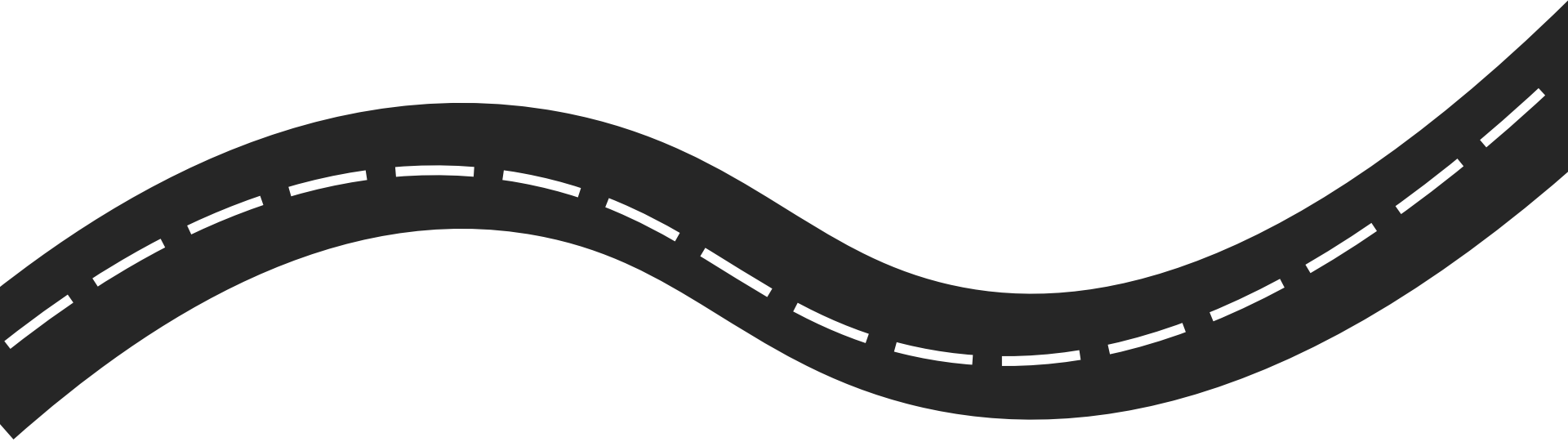
The University of Texas at Austin

***Dissertation Committee***

Junyi Jessy Li, Supervisor

Kyle Mahowald · David I. Beaver · Desmond C. Ong

April 2, 2026





## Deciphering Emotions from Text

- EMNLP 2022
- EMNLP 2023 Findings

1



## Deciphering Emotions from Text

- EMNLP 2022
- EMNLP 2023 Findings

1

### Unveiling Advanced Psychological Capabilities from LLMs: *A Case of Targeted Reappraisal*

- COLM 2024
- ICML 2025

2

## Deciphering Emotions from Text

- EMNLP 2022
- EMNLP 2023 Findings

1

## Discourse Diversity in Multi-Turn Empathic Dialogue

- *Under review*

3

## Unveiling Advanced Psychological Capabilities from LLMs: *A Case of Targeted Reappraisal*

- COLM 2024
- ICML 2025

2

## Deciphering Emotions from Text

- EMNLP 2022
- EMNLP 2023 Findings

1

## Discourse Diversity in Multi-Turn Empathic Dialogue

- *Under review*

3

## Unveiling Advanced Psychological Capabilities from LLMs: *A Case of Targeted Reappraisal*

- COLM 2024
- ICML 2025

2

4

## Conclusion

- Summary of Contributions

## Deciphering Emotions from Text

- EMNLP 2022
- EMNLP 2023 Findings

1

### Discourse Diversity in Multi-Turn Empathic Dialogue

- *Under review*

3

### Unveiling Advanced Psychological Capabilities from LLMs: *A Case of Targeted Reappraisal*

- COLM 2024
- ICML 2025

2

4

### Conclusion

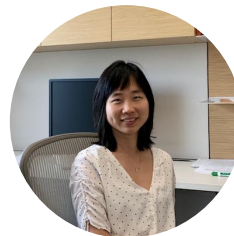
- Summary of Contributions

## Part 1 (a)

# Why Do You Feel This Way? Summarizing Triggers of Emotions in Social Media Posts

**Hongli Zhan**, Tiberiu Sosea, Cornelia Caragea, Junyi Jessy Li

*In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP 2022)*





# Summary of the Paper



# Summary of the Paper

Revealing “*Why does the writer feel [emotion]?*” is important yet remains unexplored:

# Summary of the Paper

Revealing “*Why does the writer feel [emotion]?*” is important yet remains unexplored:

- We propose a new task: **Emotion Detection & Trigger Summarization**



# Summary of the Paper

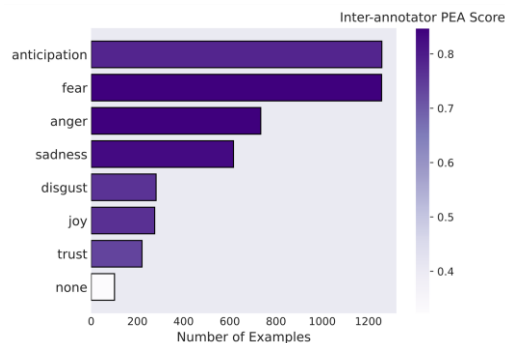
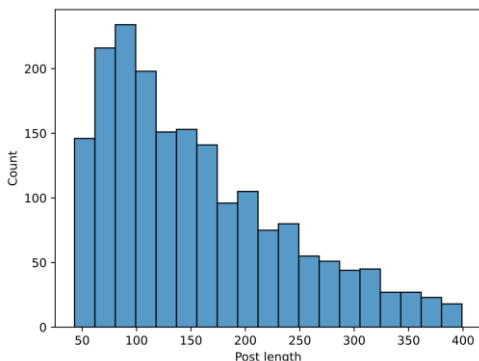
Revealing “*Why does the writer feel [emotion]?*” is important yet remains unexplored:

- We propose a new task: **Emotion Detection & Trigger Summarization**
- New benchmark CovidET: Reddit posts annotated with emotions and their triggers

# Summary of the Paper

Revealing “*Why does the writer feel [emotion]?*” is important yet remains unexplored:

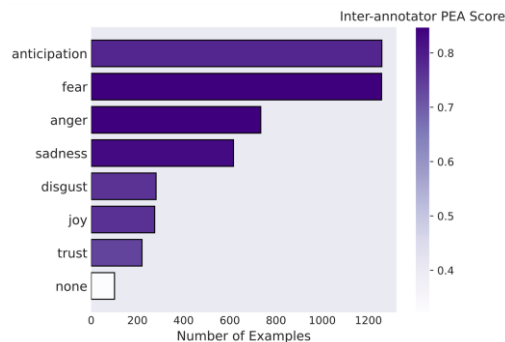
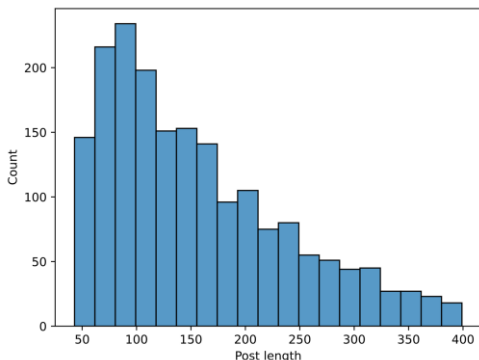
- We propose a new task: **Emotion Detection & Trigger Summarization**
- New benchmark CovidET: Reddit posts annotated with emotions and their triggers



# Summary of the Paper

Revealing “*Why does the writer feel [emotion]?*” is important yet remains unexplored:

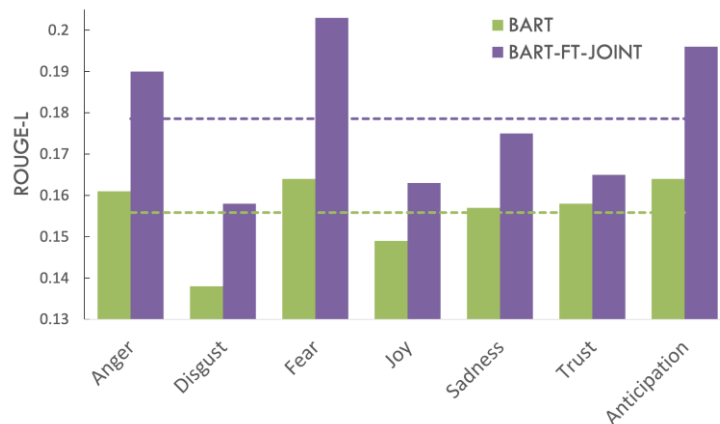
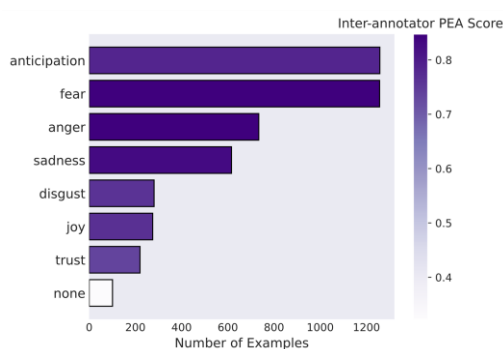
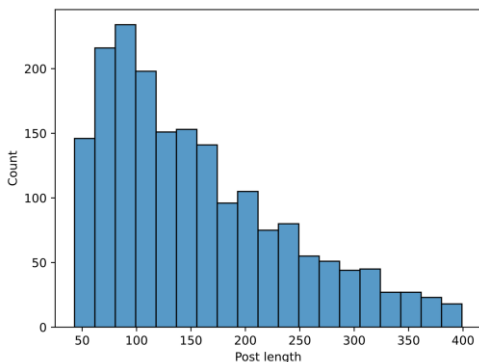
- We propose a new task: **Emotion Detection & Trigger Summarization**
- New benchmark CovidET: Reddit posts annotated with emotions and their triggers
- Benchmarking on the dataset, we were able to build automatic summarizers that outperform generic summarizers on the task of emotion-specific trigger summarization.



# Summary of the Paper

Revealing “*Why does the writer feel [emotion]?*” is important yet remains unexplored:

- We propose a new task: **Emotion Detection & Trigger Summarization**
- New benchmark CovidET: Reddit posts annotated with emotions and their triggers
- Benchmarking on the dataset, we were able to build automatic summarizers that outperform generic summarizers on the task of emotion-specific trigger summarization.





# Models can tell how people feel!

## But is that enough?

# Models can tell how people feel!

## But is that enough?



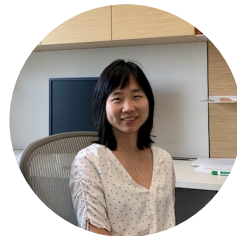
**How capable are language models of a deeper emotion understanding?**

## *Part 1 (b)*

# Evaluating Subjective Cognitive Appraisals of Emotions from Large Language Models

Hongli Zhan, Desmond C. Ong, Junyi Jessy Li

In *Findings of the Association for Computational Linguistics: EMNLP 2023*





# Cognitive Appraisals of Emotions



# Cognitive Appraisals of Emotions

**Cognitive appraisals:** the same situation can often result in different emotional experiences, based on an individual's subjective evaluations.



# Cognitive Appraisals of Emotions

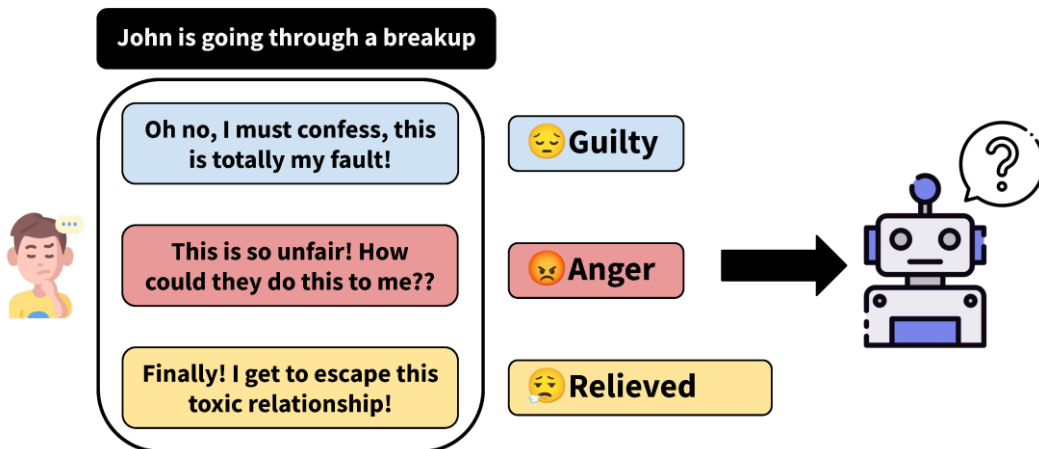
**Cognitive appraisals:** the same situation can often result in different emotional experiences, based on an individual's subjective evaluations.

- This is typically characterized by a range of different “**dimensions**” (Arnold, 1960; Smith and Ellsworth, 1985; Yeo and Ong, 2023)

# Cognitive Appraisals of Emotions

**Cognitive appraisals:** the same situation can often result in different emotional experiences, based on an individual's subjective evaluations.

- This is typically characterized by a range of different “**dimensions**” (Arnold, 1960; Smith and Ellsworth, 1985; Yeo and Ong, 2023)





# Benchmark: *CovidET-Appraisals*



# Benchmark: *CovidET-Appraisals*

From a recent meta-analysis (Yeo & Ong, 2024), we identified a set of 24 appraisals, and created prompts, e.g.,  
“To what extent did the narrator think that THEY were responsible for causing the situation?”



# Benchmark: *CovidET-Appraisals*

From a recent meta-analysis (Yeo & Ong, 2024), we identified a set of 24 appraisals, and created prompts, e.g.,  
“To what extent did the narrator think that THEY were responsible for causing the situation?”

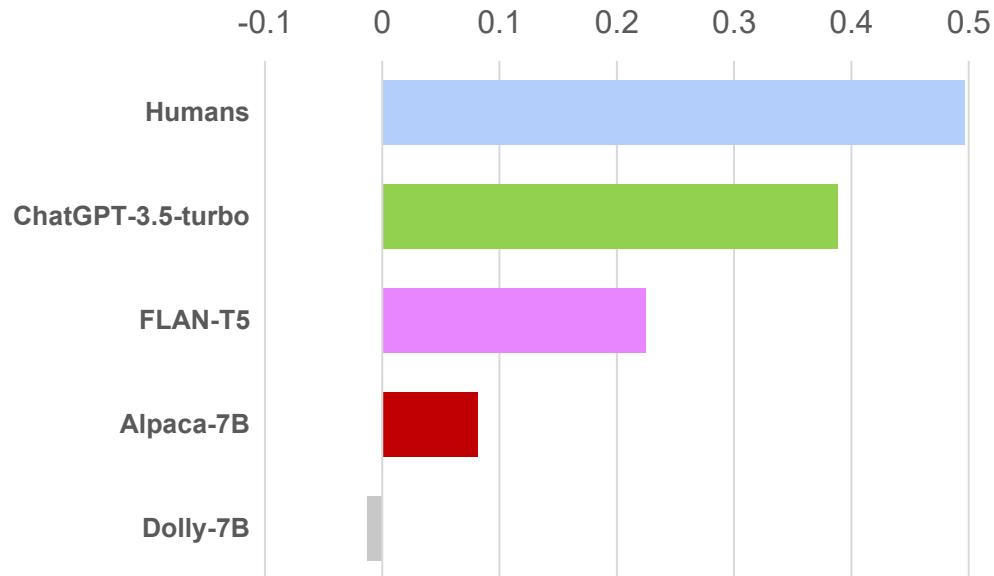
⇒ We then evaluated the accuracy of various LLMs in identifying these appraisals

# Benchmark: *CovidET-Appraisals*

From a recent meta-analysis (Yeo & Ong, 2024), we identified a set of 24 appraisals, and created prompts, e.g.,  
“To what extent did the narrator think that THEY were responsible for causing the situation?”

⇒ We then evaluated the accuracy of various LLMs in identifying these appraisals

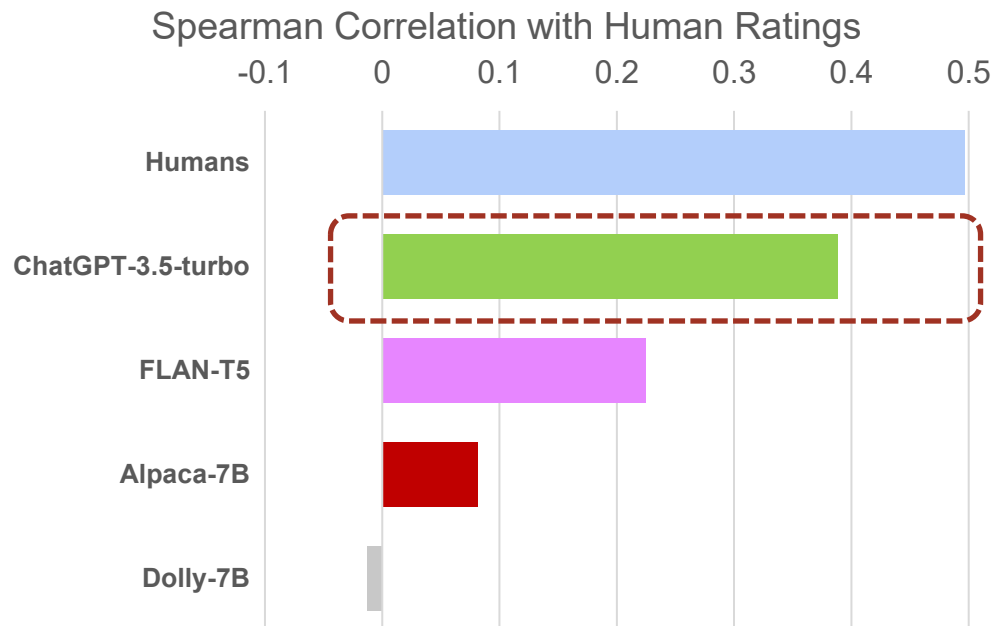
Spearman Correlation with Human Ratings



# Benchmark: CovidET-Appraisals

From a recent meta-analysis (Yeo & Ong, 2024), we identified a set of 24 appraisals, and created prompts, e.g., “To what extent did the narrator think that THEY were responsible for causing the situation?”

⇒ We then evaluated the accuracy of various LLMs in identifying these appraisals



# LLMs are able to identify appraisals of emotions (without additional training)

# LLMs are able to identify appraisals of emotions (without additional training)



⇒ Can we leverage this emotion understanding capability of LLMs to do good to humans?

## Deciphering Emotions from Text

- EMNLP 2022
- EMNLP 2023 Findings

1

## Discourse Diversity in Multi-Turn Empathic Dialogue

- *Under review*

3

## Unveiling Advanced Psychological Capabilities from LLMs: *A Case of Targeted Reappraisal*

- COLM 2024
- ICML 2025

2

4

## Conclusion

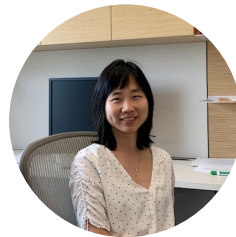
- Summary of Contributions

## Part 2 (a)

# Large Language Models are Capable of Offering Cognitive Reappraisal, if Guided

**Hongli Zhan**, Allen Zheng, Yoon Kyung Lee, Jina Suh, Junyi Jessy Li, Desmond Ong

*In Proceedings of the First Conference on Language Modeling (COLM 2024)*





# Targeted Reappraisal: A Toy Example

John is going  
through a breakup



 Guilty

This is totally my fault!

 Anger

This is so unfair! How  
could they do this to me??

Initial Appraisal





# Targeted Reappraisal: A Toy Example

John is going  
through a breakup



 Guilty

This is totally my fault!

 Anger

This is so unfair! How  
could they do this to me??

Initial Appraisal



*“A relationship requires both  
partners’ consistent effort to  
work, not just yourself”*

*“This could be an opportunity for  
personal reflection and growth”*

Targeted Re-Appraisal





# Targeted Reappraisal: A Toy Example

John is going  
through a breakup



 Guilty

This is totally my fault!

 Anger

This is so unfair! How  
could they do this to me??

Initial Appraisal



*“A relationship requires both  
partners’ consistent effort to  
work, not just yourself”*

*“This could be an opportunity for  
personal reflection and growth”*

Targeted Re-Appraisal





# Targeted Reappraisal: A Toy Example

John is going  
through a breakup



 Guilty

This is totally my fault!

 Anger

This is so unfair! How  
could they do this to me??

Initial Appraisal



*“A relationship requires both  
partners’ consistent effort to  
work, not just yourself”*

*“This could be an opportunity for  
personal reflection and growth”*

Targeted Re-Appraisal



Can we use LLMs to achieve better **emotional well-being** through offering **reappraisals**?

# Targeted Reappraisal: A Toy Example

John is going  
through a breakup



 Guilty

This is totally my fault!

 Anger

This is so unfair! How  
could they do this to me??

Initial Appraisal



*“A relationship requires both  
partners’ consistent effort to  
work, not just yourself”*

*“This could be an opportunity for  
personal reflection and growth”*

Targeted Re-Appraisal



Can we use LLMs to achieve better **emotional well-being** through offering **reappraisals**?

⇒ Such an approach would be more targeted and precise



# New Framework: *RESORT*



# New Framework: *RESORT*

We designed a system, entitled RESORT, to **guide LLMs to offer targeted reappraisals** along *six appraisal dimensions* chosen to maximize coverage

# New Framework: *RESORT*

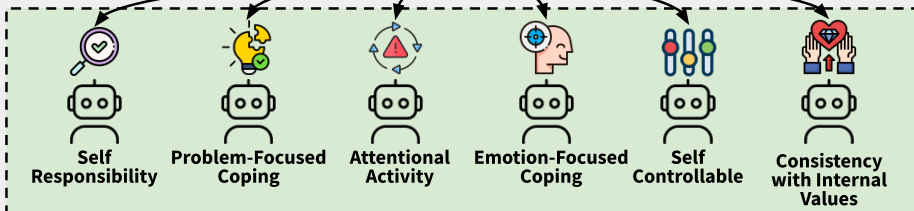


We designed a system, entitled RESORT, to **guide LLMs to offer targeted reappraisals** along *six appraisal dimensions* chosen to maximize coverage

Generating  
Targeted  
Reappraisals  
Individually



Guiding LLM  
Targeted Reappraisals  
**RESORT**

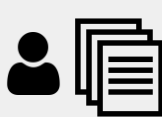


# New Framework: *RESORT*



We designed a system, entitled RESORT, to **guide LLMs to offer targeted reappraisals** along *six appraisal dimensions* chosen to maximize coverage

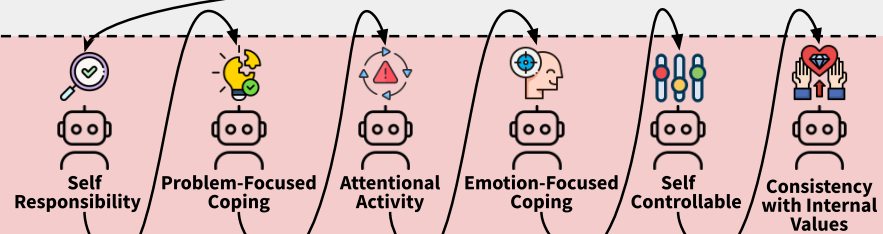
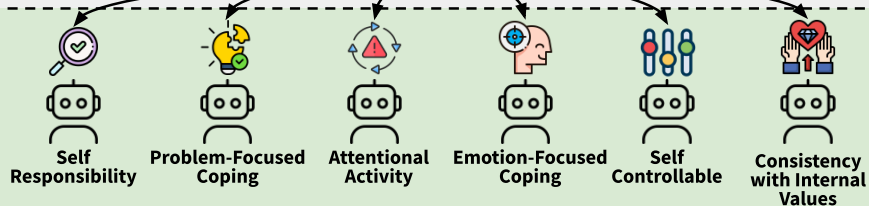
Generating  
Targeted  
Reappraisals  
Individually



Guiding LLMs to offer  
Targeted Reappraisals with  
**RESORT**  *Constitutions*



Iteratively  
Refining  
Targeted  
Reappraisals

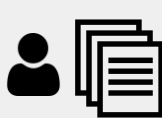


# New Framework: *RESORT*



We designed a system, entitled RESORT, to **guide LLMs to offer targeted reappraisals** along *six appraisal dimensions* chosen to maximize coverage

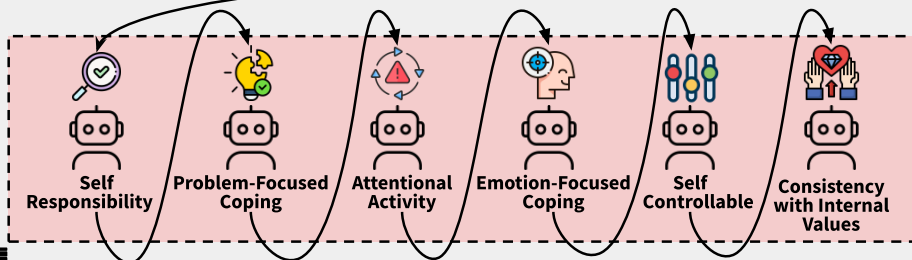
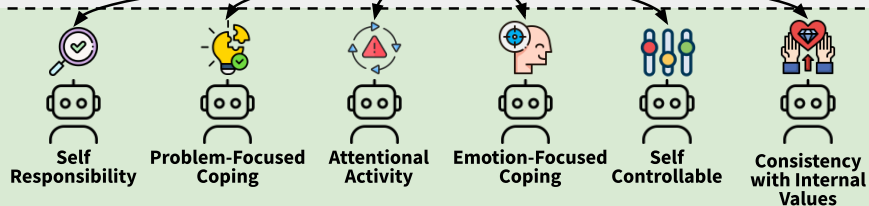
Generating  
Targeted  
Reappraisals  
Individually



Guiding LLMs to offer  
Targeted Reappraisals with  
**RESORT**  *Constitutions*



Iteratively  
Refining  
Targeted  
Reappraisals



Expert  
Psychologist  
Evaluation



... excellent, suggests ways to step  
back and reevaluate ...



# Expert-Crafted Guiding Constitutions



# Expert-Crafted Guiding Constitutions

Dimension	Appraisal	Reappraisal Goal
<i>Self responsibility</i>	Does the narrator think that they are responsible for causing the situation?	Re-evaluate whether the narrator deserves to be blamed or credited for the situation at hand. If not responsible, the narrator is encouraged to acknowledge that fact and reassess the situation.



# Expert-Crafted Guiding Constitutions

Dimension	Appraisal	Reappraisal Goal
<i>Self responsibility</i>	Does the narrator think that they are responsible for causing the situation?	Re-evaluate whether the narrator deserves to be blamed or credited for the situation at hand. If not responsible, the narrator is encouraged to acknowledge that fact and reassess the situation.
<i>Problem-focused coping</i>	Does the narrator think that they can cope with the consequences of the situation?	Focus on the narrators' competence (self-efficacy) to handle the situation at hand. The narrator is encouraged to use any resources or support to handle the situation competently and independently.



# Expert-Crafted Guiding Constitutions

Dimension	Appraisal	Reappraisal Goal
<i>Self responsibility</i>	Does the narrator think that they are responsible for causing the situation?	Re-evaluate whether the narrator deserves to be blamed or credited for the situation at hand. If not responsible, the narrator is encouraged to acknowledge that fact and reassess the situation.
<i>Problem-focused coping</i>	Does the narrator think that they can cope with the consequences of the situation?	Focus on the narrators' competence (self-efficacy) to handle the situation at hand. The narrator is encouraged to use any resources or support to handle the situation competently and independently.
<i>Attentional activity</i>	Does the narrator think that they need to attend to the situation further?	Reconsider the urgency or importance of the situation and determine if it's worth their effort and attention. If not, the narrator is encouraged to focus on other matters.



# Expert-Crafted Guiding Constitutions

Dimension	Appraisal	Reappraisal Goal
<i>Self responsibility</i>	Does the narrator think that they are responsible for causing the situation?	Re-evaluate whether the narrator deserves to be blamed or credited for the situation at hand. If not responsible, the narrator is encouraged to acknowledge that fact and reassess the situation.
<i>Problem-focused coping</i>	Does the narrator think that they can cope with the consequences of the situation?	Focus on the narrators' competence (self-efficacy) to handle the situation at hand. The narrator is encouraged to use any resources or support to handle the situation competently and independently.
<i>Attentional activity</i>	Does the narrator think that they need to attend to the situation further?	Reconsider the urgency or importance of the situation and determine if it's worth their effort and attention. If not, the narrator is encouraged to focus on other matters.
<i>Emotion-focused coping</i>	Does the narrator think that they can emotionally cope with the consequences of the event?	Re-evaluate whether the narrator can emotionally cope with the situation and regulate their emotions. If needed, consider confronting or avoiding any potential triggers that may exacerbate the stress.



# Expert-Crafted Guiding Constitutions

Dimension	Appraisal	Reappraisal Goal
<i>Self responsibility</i>	Does the narrator think that they are responsible for causing the situation?	Re-evaluate whether the narrator deserves to be blamed or credited for the situation at hand. If not responsible, the narrator is encouraged to acknowledge that fact and reassess the situation.
<i>Problem-focused coping</i>	Does the narrator think that they can cope with the consequences of the situation?	Focus on the narrators' competence (self-efficacy) to handle the situation at hand. The narrator is encouraged to use any resources or support to handle the situation competently and independently.
<i>Attentional activity</i>	Does the narrator think that they need to attend to the situation further?	Reconsider the urgency or importance of the situation and determine if it's worth their effort and attention. If not, the narrator is encouraged to focus on other matters.
<i>Emotion-focused coping</i>	Does the narrator think that they can emotionally cope with the consequences of the event?	Re-evaluate whether the narrator can emotionally cope with the situation and regulate their emotions. If needed, consider confronting or avoiding any potential triggers that may exacerbate the stress.
<i>Self controllable</i>	Does the narrator think that they can control what is happening in the situation?	Reassess the situation whether the narrator has the power or personal control over the situation. The narrator is encouraged to step back from situations that are beyond their control and focus on the things they can control.



# Expert-Crafted Guiding Constitutions

Dimension	Appraisal	Reappraisal Goal
<i>Self responsibility</i>	Does the narrator think that they are responsible for causing the situation?	Re-evaluate whether the narrator deserves to be blamed or credited for the situation at hand. If not responsible, the narrator is encouraged to acknowledge that fact and reassess the situation.
<i>Problem-focused coping</i>	Does the narrator think that they can cope with the consequences of the situation?	Focus on the narrators' competence (self-efficacy) to handle the situation at hand. The narrator is encouraged to use any resources or support to handle the situation competently and independently.
<i>Attentional activity</i>	Does the narrator think that they need to attend to the situation further?	Reconsider the urgency or importance of the situation and determine if it's worth their effort and attention. If not, the narrator is encouraged to focus on other matters.
<i>Emotion-focused coping</i>	Does the narrator think that they can emotionally cope with the consequences of the event?	Re-evaluate whether the narrator can emotionally cope with the situation and regulate their emotions. If needed, consider confronting or avoiding any potential triggers that may exacerbate the stress.
<i>Self controllable</i>	Does the narrator think that they can control what is happening in the situation?	Reassess the situation whether the narrator has the power or personal control over the situation. The narrator is encouraged to step back from situations that are beyond their control and focus on the things they can control.
<i>Consistency with internal values</i>	Does the narrator think that the situation is consistent with their personal values?	Reassess whether to what extent the situation is compatible with one's internal value (e.g., internalized social norms, beliefs, moral values). The narrator is also encouraged to consider other possible perspectives to avoid misunderstandings that may have arisen from lack of context or communication.



# Expert-Crafted Guiding Constitutions

Dimension	Appraisal	Reappraisal Goal
<i>Self responsibility</i>	Does the narrator think that they are responsible for causing the situation?	Re-evaluate whether the narrator deserves to be blamed or credited for the situation at hand. If not responsible, the narrator is encouraged to acknowledge that fact and reassess the situation.
<i>Problem-focused coping</i>	Does the narrator think that they can cope with the consequences of the situation?	Focus on the narrators' competence (self-efficacy) to handle the situation at hand. The narrator is encouraged to use any resources or support to handle the situation competently and independently.
<i>Attentional activity</i>	Does the narrator think that they need to attend to the situation further?	Reconsider the urgency or importance of the situation and determine if it's worth their effort and attention. If not, the narrator is encouraged to focus on other matters.
<i>Emotion-focused coping</i>	Does the narrator think that they can emotionally cope with the consequences of the event?	Re-evaluate whether the narrator can emotionally cope with the situation and regulate their emotions. If needed, consider confronting or avoiding any potential triggers that may exacerbate the stress.
<i>Self controllable</i>	Does the narrator think that they can control what is happening in the situation?	Reassess the situation whether the narrator has the power or personal control over the situation. The narrator is encouraged to step back from situations that are beyond their control and focus on the things they can control.
<i>Consistency with internal values</i>	Does the narrator think that the situation is consistent with their personal values?	Reassess whether to what extent the situation is compatible with one's internal value (e.g., internalized social norms, beliefs, moral values). The narrator is also encouraged to consider other possible perspectives to avoid misunderstandings that may have arisen from lack of context or communication.



# Expert-Crafted Guiding Constitutions

Dimension	Appraisal	Reappraisal Goal
<i>Self responsibility</i>	Does the narrator think that they are responsible for causing the situation?	Re-evaluate whether the narrator deserves to be blamed or credited for the situation at hand. If not responsible, the narrator is encouraged to acknowledge that fact and reassess the situation.
<i>Problem-focused coping</i>	Does the narrator think that they can cope with the consequences of the situation?	Focus on the narrators' competence (self-efficacy) to handle the situation at hand. The narrator is encouraged to use any resources or support to handle the situation competently and independently.
<i>Attentional activity</i>	Does the narrator think that they need to attend to the situation further?	Reconsider the urgency or importance of the situation and determine if it's worth their effort and attention. If not, the narrator is encouraged to focus on other matters.
<i>Emotion-focused coping</i>	Does the narrator think that they can emotionally cope with the consequences of the event?	Re-evaluate whether the narrator can emotionally cope with the situation and regulate their emotions. If needed, consider confronting or avoiding any potential triggers that may exacerbate the stress.
<i>Self controllable</i>	Does the narrator think that they can control what is happening in the situation?	Reassess the situation whether the narrator has the power or personal control over the situation. The narrator is encouraged to step back from situations that are beyond their control and focus on the things they can control.
<i>Consistency with internal values</i>	Does the narrator think that the situation is consistent with their personal values?	Reassess whether to what extent the situation is compatible with one's internal value (e.g., internalized social norms, beliefs, moral values). The narrator is also encouraged to consider other possible perspectives to avoid misunderstandings that may have arisen from lack of context or communication.

**Constitutional Principles**



# Psychology Experts' Evaluation



# Psychology Experts' Evaluation

For evaluation, we recruited 4  
psychologists with expertise in clinical  
psychology

- *All evaluators hold M.S./Ph.D. degrees*



# Psychology Experts' Evaluation

For evaluation, we recruited 4  
psychologists with expertise in clinical  
psychology

- *All evaluators hold M.S./Ph.D. degrees*

|| **Alignment** ↑ | **Empathy** ↑  
10-POINT SCALE | 5-POINT SCALE

# Psychology Experts' Evaluation

For evaluation, we recruited 4 psychologists with expertise in clinical psychology

- All evaluators hold M.S./Ph.D. degrees

		Alignment ↑		Empathy ↑	
		10-POINT SCALE		5-POINT SCALE	
		INDV	ITER	INDV	ITER
ORACLE RESPONSE		5.79		3.79	
REDDIT COMMENT		2.75		2.00	
GPT4 TURBO	vanilla	3.88		3.31	
	self-refine	2.69		2.56	
	+appr	4.69**	5.06***	3.25	4.06***
	+cons	7.31***	7.81***	3.81**	3.88**
	+appr +cons	7.12***	8.31***	3.50*	4.25***
LLAMA2 13B-CHAT	vanilla	6.25		3.88	
	self-refine	4.31		2.88	
	+appr	5.31	5.62	3.31	3.88*
	+cons	7.81***	7.81***	3.75*	4.12***
	+appr +cons	7.69***	6.44***	3.81*	3.25
MISTRAL 7B-INSTRUCT	vanilla	4.36		2.86	
	self-refine	4.14		2.64	
	+appr	5.50	5.64**	2.93	2.57
	+cons	6.50**	7.43**	3.43*	3.71**
	+appr +cons	6.71**	5.71	2.79	3.14

# Psychology Experts' Evaluation

For evaluation, we recruited 4 psychologists with expertise in clinical psychology

- All evaluators hold M.S./Ph.D. degrees

Evaluation done by psychologists suggest that responses from our system are:

		Alignment ↑		Empathy ↑	
		10-POINT SCALE		5-POINT SCALE	
		INDV	ITER	INDV	ITER
ORACLE RESPONSE		5.79		3.79	
REDDIT COMMENT		2.75		2.00	
GPT4 TURBO	vanilla	3.88		3.31	
	self-refine	2.69		2.56	
	+appr	4.69**	5.06***	3.25	4.06***
	+cons	7.31***	7.81***	3.81**	3.88**
	+appr +cons	7.12***	8.31***	3.50*	4.25***
LLAMA2 13B-CHAT	vanilla	6.25		3.88	
	self-refine	4.31		2.88	
	+appr	5.31	5.62	3.31	3.88*
	+cons	7.81***	7.81***	3.75*	4.12***
	+appr +cons	7.69***	6.44***	3.81*	3.25
MISTRAL 7B-INSTRUCT	vanilla	4.36		2.86	
	self-refine	4.14		2.64	
	+appr	5.50	5.64**	2.93	2.57
	+cons	6.50**	7.43**	3.43*	3.71**
	+appr +cons	6.71**	5.71	2.79	3.14

# Psychology Experts' Evaluation

For evaluation, we recruited 4 psychologists with expertise in clinical psychology

- All evaluators hold M.S./Ph.D. degrees

Evaluation done by psychologists suggest that responses from our system are:

- (i) **aligned** (with reappraisal definitions)

		Alignment ↑		Empathy ↑	
		10-POINT SCALE		5-POINT SCALE	
		INDV	ITER	INDV	ITER
ORACLE RESPONSE		5.79		3.79	
REDDIT COMMENT		2.75		2.00	
GPT4 TURBO	vanilla	3.88		3.31	
	self-refine	2.69		2.56	
	+appr	4.69**	5.06***	3.25	4.06***
	+cons	7.31***	7.81***	3.81**	3.88**
	+appr +cons	7.12***	8.31***	3.50*	4.25***
LLAMA2 13B-CHAT	vanilla	6.25		3.88	
	self-refine	4.31		2.88	
	+appr	5.31	5.62	3.31	3.88*
	+cons	7.81***	7.81***	3.75*	4.12***
	+appr +cons	7.69***	6.44***	3.81*	3.25
MISTRAL 7B-INSTRUCT	vanilla	4.36		2.86	
	self-refine	4.14		2.64	
	+appr	5.50	5.64**	2.93	2.57
	+cons	6.50**	7.43**	3.43*	3.71**
	+appr +cons	6.71**	5.71	2.79	3.14



# Psychology Experts' Evaluation

For evaluation, we recruited 4 psychologists with expertise in clinical psychology

- All evaluators hold M.S./Ph.D. degrees

Evaluation done by psychologists suggest that responses from our system are:

- (i) **aligned** (with reappraisal definitions)

		Alignment ↑		Empathy ↑	
		10-POINT SCALE		5-POINT SCALE	
		INDV	ITER	INDV	ITER
ORACLE RESPONSE		5.79		3.79	
REDDIT COMMENT		2.75		2.00	
GPT4 TURBO	vanilla	3.88		3.31	
	self-refine	2.69		2.56	
	+appr	4.69**	5.06***	3.25	4.06***
	+cons	7.31***	7.81***	3.81**	3.88**
	+appr +cons	7.12***	<b>8.31***</b>	3.50*	<b>4.25***</b>
LLAMA2 13B-CHAT	vanilla	6.25		3.88	
	self-refine	4.31		2.88	
	+appr	5.31	5.62	3.31	3.88*
	+cons	<b>7.81***</b>	<b>7.81***</b>	3.75*	<b>4.12***</b>
	+appr +cons	7.69***	<b>6.44***</b>	3.81*	3.25
MISTRAL 7B-INSTRUCT	vanilla	4.36		2.86	
	self-refine	4.14		2.64	
	+appr	5.50	5.64**	2.93	2.57
	+cons	6.50**	<b>7.43**</b>	3.43*	<b>3.71**</b>
	+appr +cons	6.71**	5.71	2.79	3.14

# Psychology Experts' Evaluation

For evaluation, we recruited 4 psychologists with expertise in clinical psychology

- All evaluators hold M.S./Ph.D. degrees

Evaluation done by psychologists suggest that responses from our system are:

- aligned** (with reappraisal definitions)
- empathic** compared to various baselines

		Alignment ↑ 10-POINT SCALE		Empathy ↑ 5-POINT SCALE	
		INDV	ITER	INDV	ITER
ORACLE RESPONSE		5.79		3.79	
REDDIT COMMENT		2.75		2.00	
GPT4 TURBO	vanilla	3.88		3.31	
	self-refine	2.69		2.56	
	+appr	4.69**	5.06***	3.25	4.06***
	+cons	7.31***	7.81***	3.81**	3.88**
	+appr +cons	7.12***	<b>8.31***</b>	3.50*	<b>4.25***</b>
LLAMA2 13B-CHAT	vanilla	6.25		3.88	
	self-refine	4.31		2.88	
	+appr	5.31	5.62	3.31	3.88*
	+cons	<b>7.81***</b>	<b>7.81***</b>	3.75*	<b>4.12***</b>
	+appr +cons	7.69***	6.44***	3.81*	3.25
MISTRAL 7B-INSTRUCT	vanilla	4.36		2.86	
	self-refine	4.14		2.64	
	+appr	5.50	5.64**	2.93	2.57
	+cons	6.50**	<b>7.43**</b>	3.43*	<b>3.71**</b>
	+appr +cons	6.71**	5.71	2.79	3.14



# Psychology Experts' Evaluation

For evaluation, we recruited 4 psychologists with expertise in clinical psychology

- All evaluators hold M.S./Ph.D. degrees

Evaluation done by psychologists suggest that responses from our system are:

- aligned** (with reappraisal definitions)
- empathic** compared to various baselines

		Alignment ↑ 10-POINT SCALE		Empathy ↑ 5-POINT SCALE	
		INDV	ITER	INDV	ITER
ORACLE RESPONSE		5.79		3.79	
REDDIT COMMENT		2.75		2.00	
GPT4 TURBO	vanilla	3.88		3.31	
	self-refine	2.69		2.56	
	+appr	4.69**	5.06***	3.25	4.06***
	+cons	7.31***	7.81***	3.81**	3.88**
	+appr +cons	7.12***	<b>8.31***</b>	3.50*	<b>4.25***</b>
LLAMA2 13B-CHAT	vanilla	6.25		3.88	
	self-refine	4.31		2.88	
	+appr	5.31	5.62	3.31	3.88*
	+cons	<b>7.81***</b>	<b>7.81***</b>	3.75*	<b>4.12***</b>
	+appr +cons	7.69***	6.44***	3.81*	3.25
MISTRAL 7B-INSTRUCT	vanilla	4.36		2.86	
	self-refine	4.14		2.64	
	+appr	5.50	5.64**	2.93	2.57
	+cons	6.50**	<b>7.43**</b>	3.43*	<b>3.71**</b>
	+appr +cons	6.71**	5.71	2.79	3.14



**LLMs (under expert guidance) can  
generate targeted reappraisals that are  
both “*aligned*” and *empathic***

# LLMs (under expert guidance) can generate targeted reappraisals that are both “*aligned*” and *empathic*”



⇒ Is it possible to **automate the guiding process** with **as little human supervision as possible**?

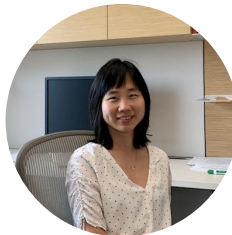
## *Part 2 (b)*

# SPRI: Aligning Large Language Models with Context-Situated Principles

**Hongli Zhan**, Muneeza Azmat, Raya Horesh, Junyi Jessy Li, Mikhail Yurochkin

*In Proceedings of the 42nd International Conference on Machine Learning (ICML 2025)*

*Work done during internship at IBM Research*





# Motivation

# Motivation

User



Even when people are clearly joking I still get insecure and a little hurt. I do my best not to show it but i think to the more perceptive folks it's probably obvious ... It's so stupid. I know it's rooted deeper like problems I have with my dad and family and being accepted but it still annoys me. Is there any fix to this?

# Motivation

In the context of providing **cognitive reappraisals** for emotional support to users in distress:

User



Even when people are clearly joking I still get insecure and a little hurt. I do my best not to show it but i think to the more perceptive folks it's probably obvious ... It's so stupid. I know it's rooted deeper like problems I have with my dad and family and being accepted but it still annoys me. Is there any fix to this?

# Motivation

In the context of providing **cognitive reappraisals** for emotional support to users in distress:

User



Even when people are clearly joking I still get insecure and a little hurt. I do my best not to show it but i think to the more perceptive folks it's probably obvious ... It's so stupid. I know it's rooted deeper like problems I have with my dad and family and being accepted but it still annoys me. Is there any fix to this?

Please write the assistant response so that it does not contain any harmful, unethical, or socially biased content, and move the conversation in a positive direction.

Generic  
Rules



# Motivation

In the context of providing **cognitive reappraisals** for emotional support to users in distress:

- *Generic principles* are often insufficient to capture the complexities of the use-case

User



Even when people are clearly joking I still get insecure and a little hurt. I do my best not to show it but i think to the more perceptive folks it's probably obvious ... It's so stupid. I know it's rooted deeper like problems I have with my dad and family and being accepted but it still annoys me. Is there any fix to this?

Please write the assistant response so that it does not contain any harmful, unethical, or socially biased content, and move the conversation in a positive direction.

Generic  
Rules



# Motivation

In the context of providing **cognitive reappraisals** for emotional support to users in distress:

- *Generic principles* are often insufficient to capture the complexities of the use-case

User



Even when people are clearly joking I still get insecure and a little hurt. I do my best not to show it but i think to the more perceptive folks it's probably obvious ... It's so stupid. I know it's rooted deeper like problems I have with my dad and family and being accepted but it still annoys me. Is there any fix to this?

Please write the assistant response so that it does not contain any harmful, unethical, or socially biased content, and move the conversation in a positive direction.

Generic  
Rules



If the narrator is stressing over things they are not responsible for, tell them that it may not require as much responsibility as they think and not to worry about them too much. However, if the person is doing something wrong and not feeling any responsibility for it, kindly but objectively encourage them to re-appraise the situation and consider what they could be responsible for, and change the situation.

Human  
Experts



# Motivation

In the context of providing **cognitive reappraisals** for emotional support to users in distress:

- *Generic principles* are often insufficient to capture the complexities of the use-case
- *Expert-written guidance* takes too much time

User



Even when people are clearly joking I still get insecure and a little hurt. I do my best not to show it but i think to the more perceptive folks it's probably obvious ... It's so stupid. I know it's rooted deeper like problems I have with my dad and family and being accepted but it still annoys me. Is there any fix to this?

Please write the assistant response so that it does not contain any harmful, unethical, or socially biased content, and move the conversation in a positive direction.

Generic Rules



If the narrator is stressing over things they are not responsible for, tell them that it may not require as much responsibility as they think and not to worry about them too much. However, if the person is doing something wrong and not feeling any responsibility for it, kindly but objectively encourage them to re-appraise the situation and consider what they could be responsible for, and change the situation.

Human Experts





# Motivation

In the context of providing **cognitive reappraisals** for emotional support to users in distress:

- *Generic principles* are often insufficient to capture the complexities of the use-case
- *Expert-written guidance* takes too much time

Can we build a framework that **tailors the guidance to each individual query**, whilst **minimizing the human efforts** needed for annotations?

User



Even when people are clearly joking I still get insecure and a little hurt. I do my best not to show it but i think to the more perceptive folks it's probably obvious ... It's so stupid. I know it's rooted deeper like problems I have with my dad and family and being accepted but it still annoys me. Is there any fix to this?

Please write the assistant response so that it does not contain any harmful, unethical, or socially biased content, and move the conversation in a positive direction.

Generic Rules



If the narrator is stressing over things they are not responsible for, tell them that it may not require as much responsibility as they think and not to worry about them too much. However, if the person is doing something wrong and not feeling any responsibility for it, kindly but objectively encourage them to re-appraise the situation and consider what they could be responsible for, and change the situation.

Human Experts





# Motivation

In the context of providing **cognitive reappraisals** for emotional support to users in distress:

- *Generic principles* are often insufficient to capture the complexities of the use-case
- *Expert-written guidance* takes too much time

Can we build a framework that **tailors the guidance to each individual query**, whilst **minimizing the human efforts** needed for annotations?

User



Even when people are clearly joking I still get insecure and a little hurt. I do my best not to show it but i think to the more perceptive folks it's probably obvious ... It's so stupid. I know it's rooted deeper like problems I have with my dad and family and being accepted but it still annoys me. Is there any fix to this?

Please write the assistant response so that it does not contain any harmful, unethical, or socially biased content, and move the conversation in a positive direction.

Generic Rules



If the narrator is stressing over things they are not responsible for, tell them that it may not require as much responsibility as they think and not to worry about them too much. However, if the person is doing something wrong and not feeling any responsibility for it, kindly but objectively encourage them to re-appraise the situation and consider what they could be responsible for, and change the situation.

Human Experts



Acknowledge the narrator's emotional response without judgment, while gently guiding them to reframe their perception of responsibility ... Suggest that the narrator's past experiences (e.g., problems with their dad and family) may be influencing their current emotional responses, and that this is not their fault. Encourage self-reflection to identify whether there are any patterns or triggers that contribute to their feelings of insecurity and hurt ...

SPRI w/  
GPT-4o  
(mini)





# Introducing: Situated-PRinciples (SPRI)



# Introducing: Situated-PRinciples (SPRI)

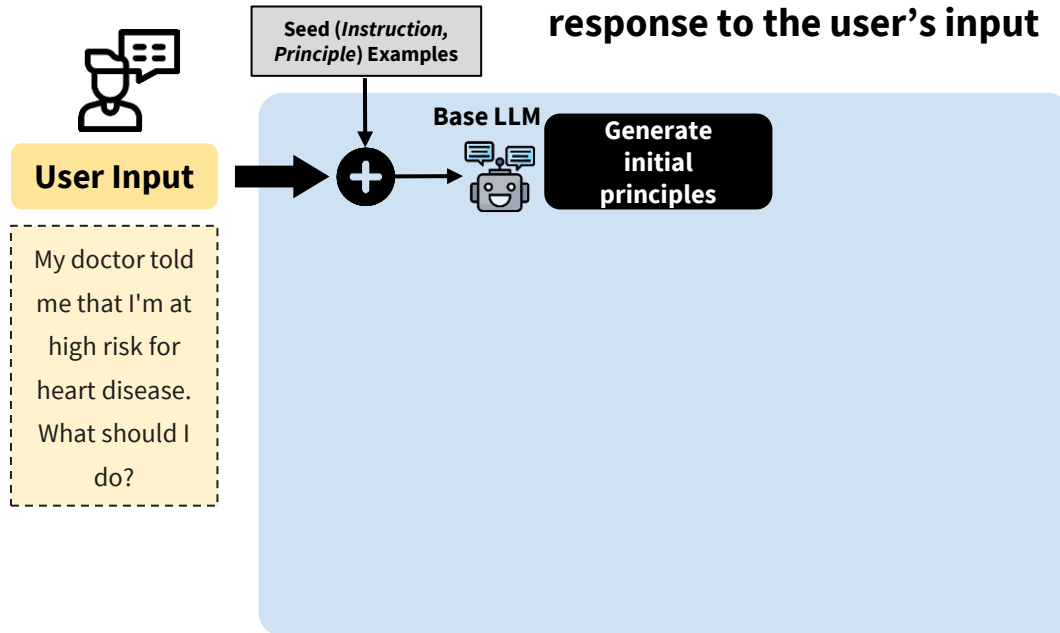


## User Input

My doctor told me that I'm at high risk for heart disease. What should I do?

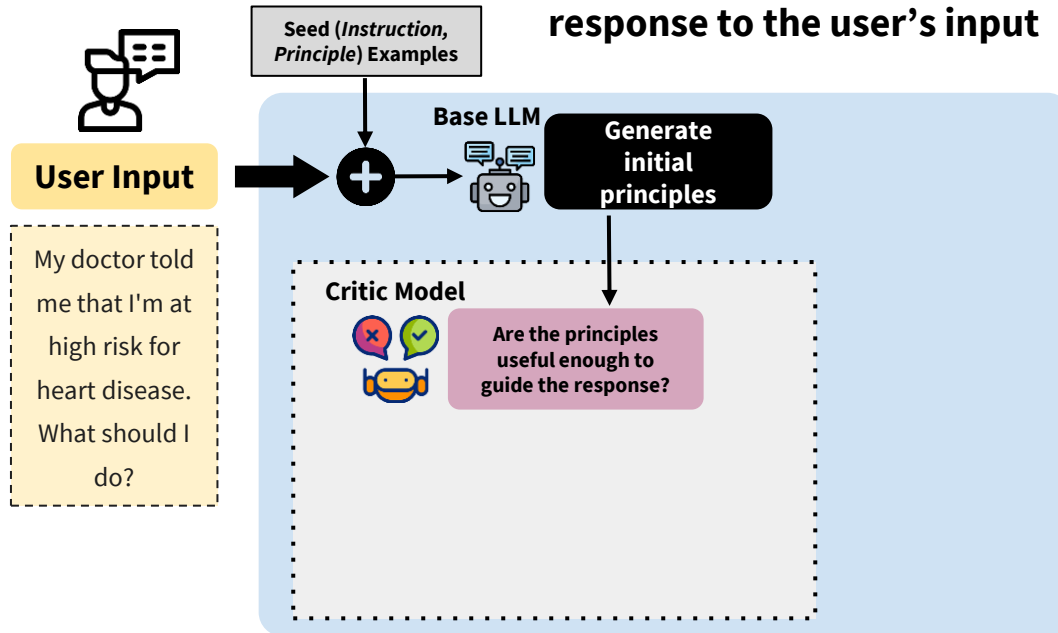
# Introducing: Situated-PRinciples (SPRI)

**Stage 1: Generate a set of  
principles to guide the  
response to the user's input**



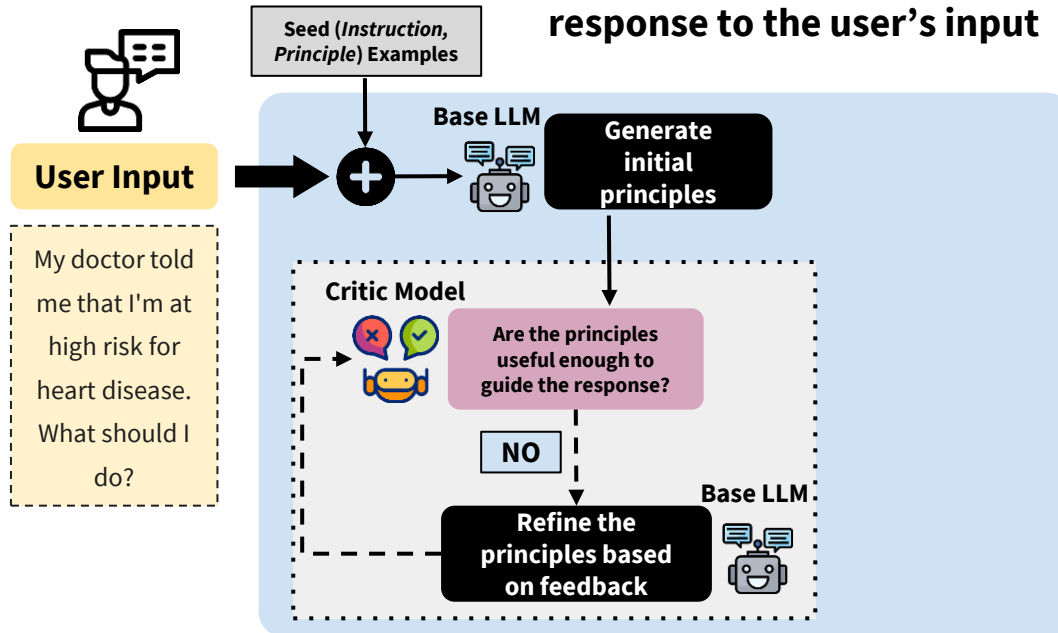
# Introducing: Situated-PRinciples (SPRI)

**Stage 1: Generate a set of  
principles to guide the  
response to the user's input**



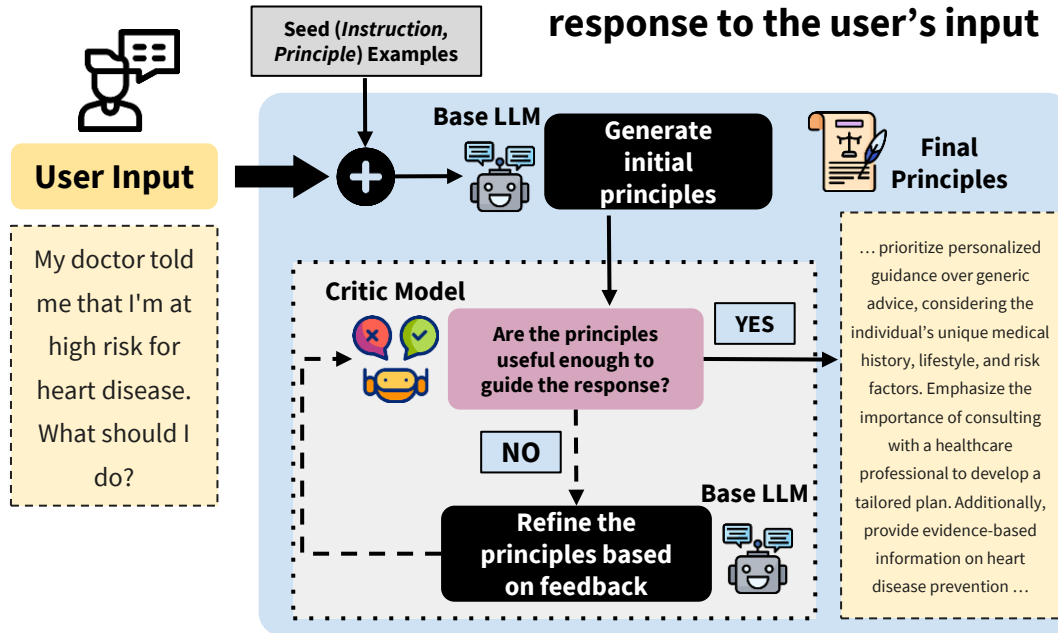
# Introducing: Situated-PRinciples (SPRI)

**Stage 1: Generate a set of  
principles to guide the  
response to the user's input**



# Introducing: Situated-PRinciples (SPRI)

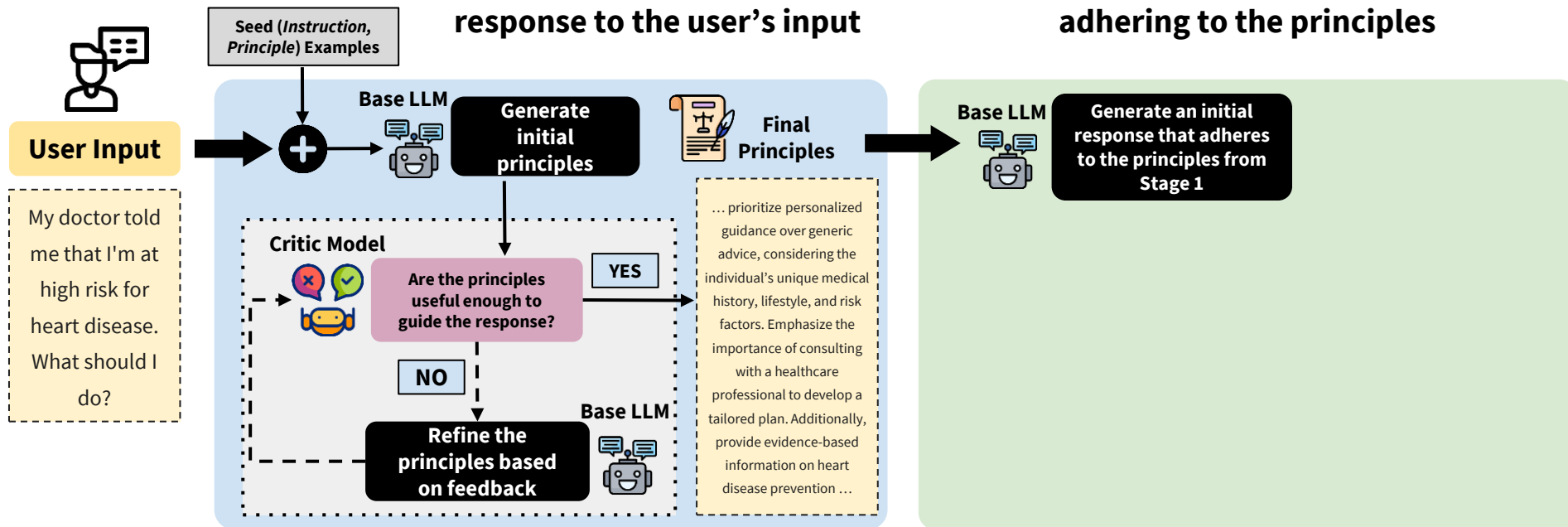
**Stage 1: Generate a set of principles to guide the response to the user's input**



# Introducing: Situated-PRinciples (SPRI)

**Stage 1: Generate a set of principles to guide the response to the user's input**

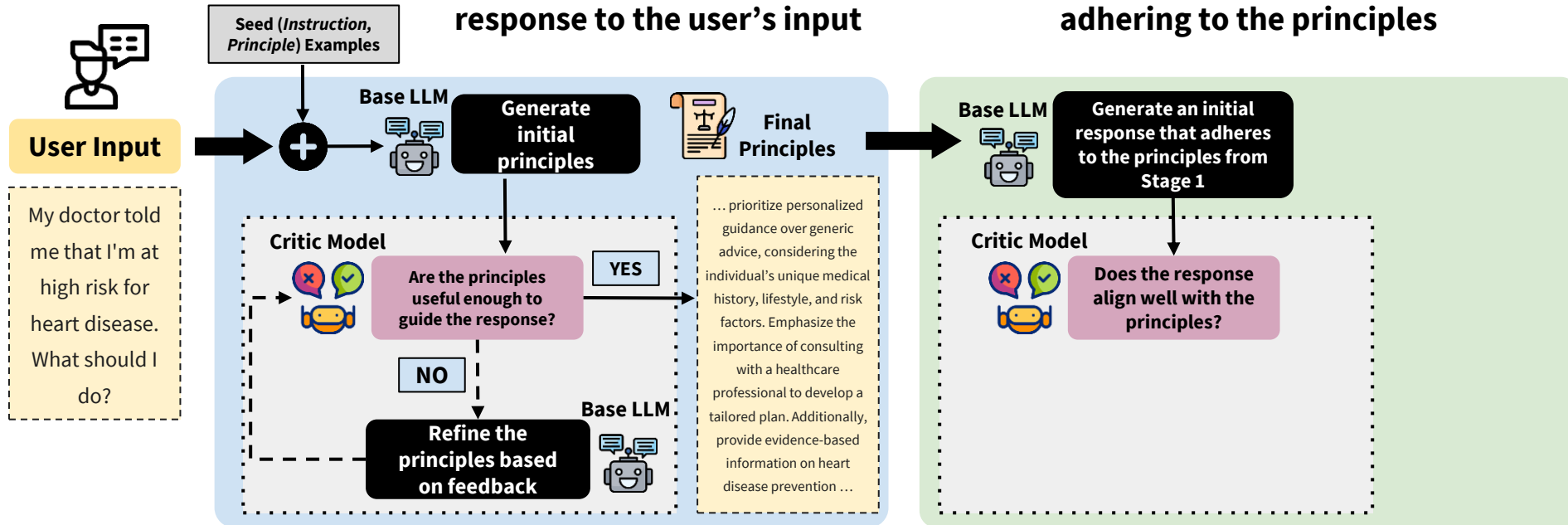
**Stage 2: Generate a response to the user's input by adhering to the principles**



# Introducing: Situated-PRinciples (SPRI)

**Stage 1: Generate a set of principles to guide the response to the user's input**

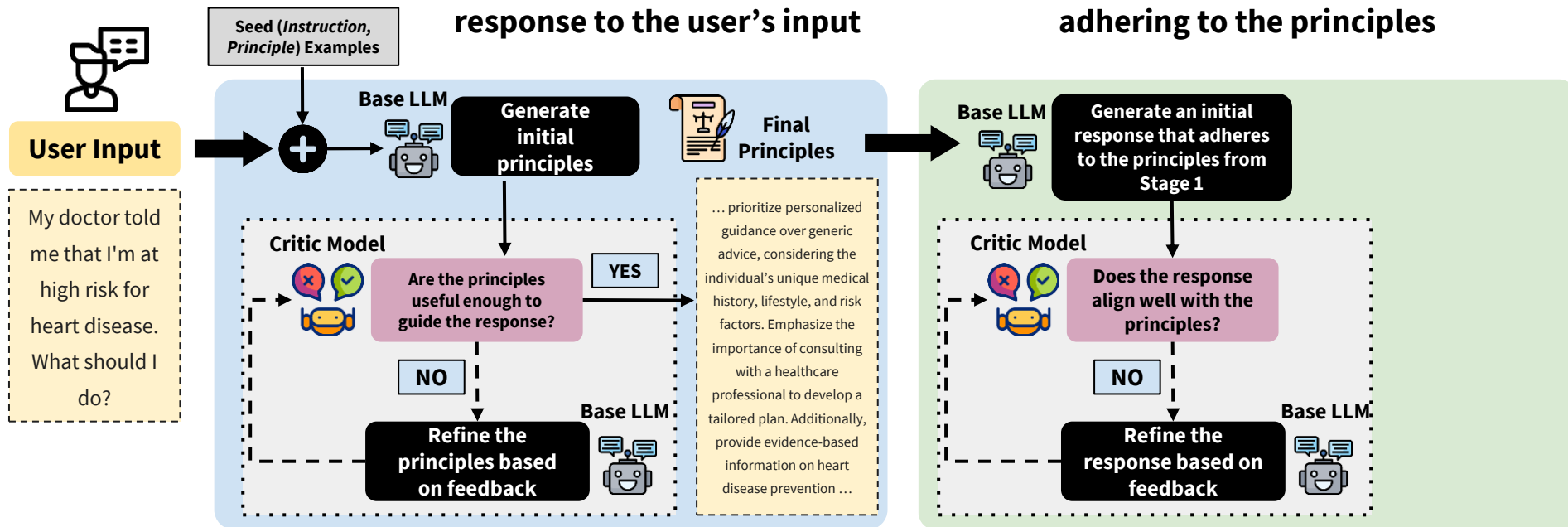
**Stage 2: Generate a response to the user's input by adhering to the principles**



# Introducing: Situated-PRinciples (SPRI)

**Stage 1: Generate a set of principles to guide the response to the user's input**

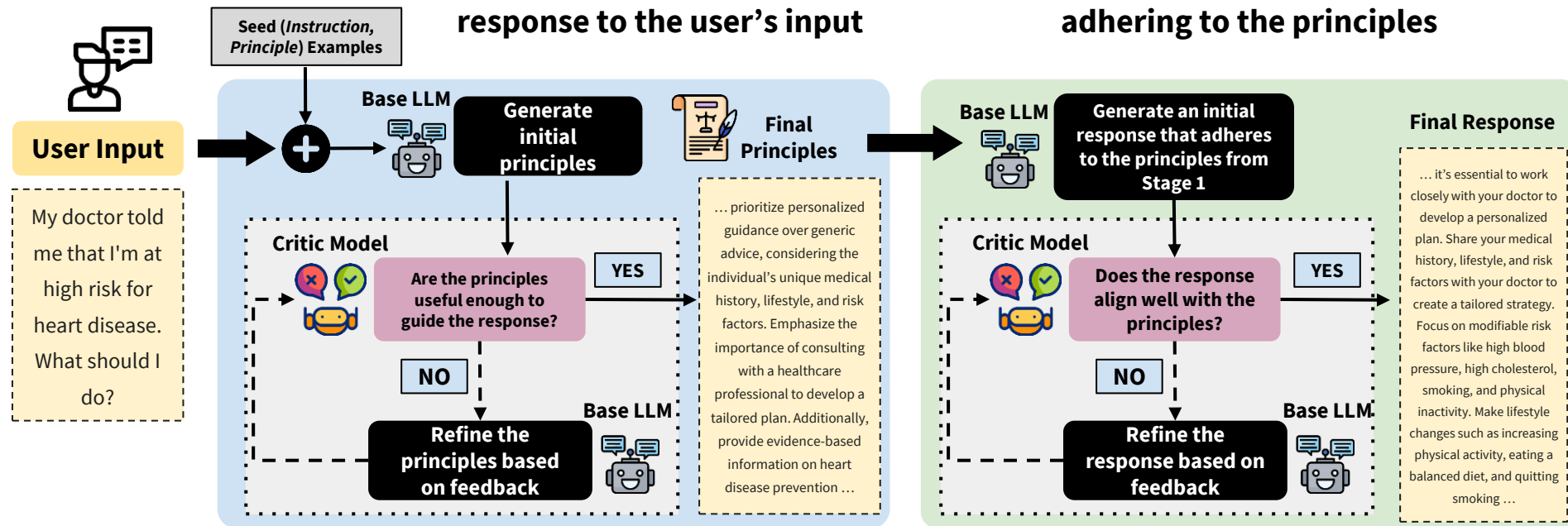
**Stage 2: Generate a response to the user's input by adhering to the principles**



# Introducing: Situated-PRinciples (SPRI)

**Stage 1: Generate a set of principles to guide the response to the user's input**

**Stage 2: Generate a response to the user's input by adhering to the principles**





# Evaluation Setup

# Evaluation Setup

we provide SPRI with a **single** oracle (expert-written) reappraisal constitution as the **seed example** in the initial principle-generation stage.

# Evaluation Setup

we provide SPRI with a **single** oracle (expert-written) reappraisal constitution as the **seed example** in the initial principle-generation stage.

- **Evaluation Data:** 30 Reddit posts from (Zhan et al., COLM 2024)

# Evaluation Setup

we provide SPRI with a **single** oracle (expert-written) reappraisal constitution as the **seed example** in the initial principle-generation stage.

- **Evaluation Data:** 30 Reddit posts from (Zhan et al., COLM 2024)

We carry out **automatic evaluation** on all reappraisal responses elicited using **GPT-4-0613**

# Evaluation Setup

we provide SPRI with a **single** oracle (expert-written) reappraisal constitution as the **seed example** in the initial principle-generation stage.

- **Evaluation Data:** 30 Reddit posts from (Zhan et al., COLM 2024)

We carry out **automatic evaluation** on all reappraisal responses elicited using **GPT-4-0613**

- Which showed strong correlation with evaluation results conducted by professional psychologists in Zhan et al. (COLM 2024)



# Evaluation Results for **Cognitive Reappraisal**



# Evaluation Results for **Cognitive Reappraisal**

---

|| GPT-4o-mini || Llama-3.1-70B-Instruct || Llama-3-8B-Instruct || Mixtral-8×7B-Instruct



# Evaluation Results for **Cognitive Reappraisal**

GPT-4o-mini		Llama-3.1-70B-Instruct		Llama-3-8B-Instruct		Mixtral-8×7B-Instruct	
Alignment ↑	Empathy ↑	Alignment ↑	Empathy ↑	Alignment ↑	Empathy ↑	Alignment ↑	Empathy ↑
Scale of 10	Scale of 5	Scale of 10	Scale of 5	Scale of 10	Scale of 5	Scale of 10	Scale of 5

# Evaluation Results for **Cognitive Reappraisal**

	GPT-4o-mini		Llama-3.1-70B-Instruct		Llama-3-8B-Instruct		Mixtral-8×7B-Instruct	
	Alignment ↑	Empathy ↑	Alignment ↑	Empathy ↑	Alignment ↑	Empathy ↑	Alignment ↑	Empathy ↑
	Scale of 10	Scale of 5	Scale of 10	Scale of 5	Scale of 10	Scale of 5	Scale of 10	Scale of 5
vanilla	7.90	4.50	7.77	4.43	7.10	3.90	7.53	4.50



# Evaluation Results for **Cognitive Reappraisal**

	GPT-4o-mini		Llama-3.1-70B-Instruct		Llama-3-8B-Instruct		Mixtral-8×7B-Instruct	
	Alignment ↑	Empathy ↑	Alignment ↑	Empathy ↑	Alignment ↑	Empathy ↑	Alignment ↑	Empathy ↑
	Scale of 10	Scale of 5	Scale of 10	Scale of 5	Scale of 10	Scale of 5	Scale of 10	Scale of 5
vanilla	7.90	4.50	7.77	4.43	7.10	3.90	7.53	4.50
self-refine	7.73	4.53	7.50	4.27	7.20	4.07	6.60	3.90



# Evaluation Results for **Cognitive Reappraisal**

	GPT-4o-mini		Llama-3.1-70B-Instruct		Llama-3-8B-Instruct		Mixtral-8×7B-Instruct	
	Alignment ↑	Empathy ↑	Alignment ↑	Empathy ↑	Alignment ↑	Empathy ↑	Alignment ↑	Empathy ↑
	Scale of 10	Scale of 5	Scale of 10	Scale of 5	Scale of 10	Scale of 5	Scale of 10	Scale of 5
vanilla	7.90	4.50	7.77	4.43	7.10	3.90	7.53	4.50
self-refine	7.73	4.53	7.50	4.27	7.20	4.07	6.60	3.90
SPRI	<b>8.00<sup>†</sup></b>	<b>4.73</b>	<b>8.17*<sup>†</sup></b>	<b>4.77*<sup>†</sup></b>	<b>7.90*<sup>†</sup></b>	<b>4.47*<sup>†</sup></b>	<b>8.03*<sup>†</sup></b>	<b>4.77*<sup>†</sup></b>



# Evaluation Results for Cognitive Reappraisal

	GPT-4o-mini		Llama-3.1-70B-Instruct		Llama-3-8B-Instruct		Mixtral-8×7B-Instruct	
	Alignment ↑	Empathy ↑	Alignment ↑	Empathy ↑	Alignment ↑	Empathy ↑	Alignment ↑	Empathy ↑
	Scale of 10	Scale of 5	Scale of 10	Scale of 5	Scale of 10	Scale of 5	Scale of 10	Scale of 5
vanilla	7.90	4.50	7.77	4.43	7.10	3.90	7.53	4.50
self-refine	7.73	4.53	7.50	4.27	7.20	4.07	6.60	3.90
SPRI	<b>8.00<sup>†</sup></b>	<b>4.73</b>	<b>8.17*<sup>†</sup></b>	<b>4.77*<sup>†</sup></b>	<b>7.90*<sup>†</sup></b>	<b>4.47*<sup>†</sup></b>	<b>8.03*<sup>†</sup></b>	<b>4.77*<sup>†</sup></b>
oracle principles	8.67* <sup>†</sup>	4.80* <sup>†</sup>	8.53* <sup>†</sup>	4.20	8.33* <sup>†</sup>	4.30*	8.17	4.07



# Evaluation Results for Cognitive Reappraisal

	GPT-4o-mini		Llama-3.1-70B-Instruct		Llama-3-8B-Instruct		Mixtral-8×7B-Instruct	
	Alignment ↑ Scale of 10	Empathy ↑ Scale of 5	Alignment ↑ Scale of 10	Empathy ↑ Scale of 5	Alignment ↑ Scale of 10	Empathy ↑ Scale of 5	Alignment ↑ Scale of 10	Empathy ↑ Scale of 5
vanilla	7.90	4.50	7.77	4.43	7.10	3.90	7.53	4.50
self-refine	7.73	4.53	7.50	4.27	7.20	4.07	6.60	3.90
SPRI	<b>8.00<sup>†</sup></b>	<b>4.73</b>	<b>8.17*<sup>†</sup></b>	<b>4.77*<sup>†</sup></b>	<b>7.90*<sup>†</sup></b>	<b>4.47*<sup>†</sup></b>	<b>8.03*<sup>†</sup></b>	<b>4.77*<sup>†</sup></b>
oracle principles	8.67* <sup>†</sup>	4.80* <sup>†</sup>	8.53* <sup>†</sup>	4.20	8.33* <sup>†</sup>	4.30*	8.17	4.07



# Evaluation Results for Cognitive Reappraisal

	GPT-4o-mini		Llama-3.1-70B-Instruct		Llama-3-8B-Instruct		Mixtral-8×7B-Instruct	
	Alignment ↑	Empathy ↑	Alignment ↑	Empathy ↑	Alignment ↑	Empathy ↑	Alignment ↑	Empathy ↑
	Scale of 10	Scale of 5	Scale of 10	Scale of 5	Scale of 10	Scale of 5	Scale of 10	Scale of 5
vanilla	7.90	4.50	7.77	4.43	7.10	3.90	7.53	4.50
self-refine	7.73	4.53	7.50	4.27	7.20	4.07	6.60	3.90
SPRI	<b>8.00<sup>†</sup></b>	<b>4.73</b>	<b>8.17*<sup>†</sup></b>	<b>4.77*<sup>†</sup></b>	<b>7.90*<sup>†</sup></b>	<b>4.47*<sup>†</sup></b>	<b>8.03*<sup>†</sup></b>	<b>4.77*<sup>†</sup></b>
oracle principles	8.67* <sup>†</sup>	4.80* <sup>†</sup>	8.53* <sup>†</sup>	4.20	8.33* <sup>†</sup>	4.30*	8.17	4.07

**SPRI consistently outperforms methods that lack access to oracle principles** both in terms of reappraisal alignment and perceived empathy, even though it only utilizes a single seed principle.

# **LLMs Do Well on These Emotion Tasks!**

## **But they were all carried out in single-turn settings**

# LLMs Do Well on These Emotion Tasks!

## But they were all carried out in single-turn settings



# LLMs Do Well on These Emotion Tasks!

## But they were all carried out in single-turn settings

**Empathy is not just about one good response.**



# LLMs Do Well on These Emotion Tasks!

## But they were all carried out in single-turn settings



**Empathy is not just about one good response.  
It is about adapting your support as a conversation unfolds.**

# LLMs Do Well on These Emotion Tasks!

## But they were all carried out in single-turn settings



**Empathy is not just about one good response.  
It is about adapting your support as a conversation unfolds.**

***⇒ Can LLMs sustain this across multi-turn conversations?***

## Deciphering Emotions from Text

- EMNLP 2022
- EMNLP 2023 Findings

1

## Discourse Diversity in Multi-Turn Empathic Dialogue

- *Under review*

3

## Unveiling Advanced Psychological Capabilities from LLMs: *A Case of Targeted Reappraisal*

- COLM 2024
- ICML 2025

2

4

## Conclusion

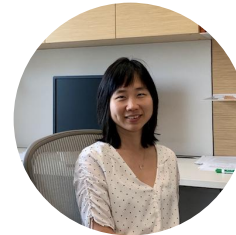
- Summary of Contributions

# *Part 3*

## Discourse Diversity in Multi-Turn Empathic Dialogue

**Hongli Zhan**, Emma Gueorguieva, Javier Hernandez, Jina Suh, Desmond C. Ong, Junyi Jessy Li

*Under Review*



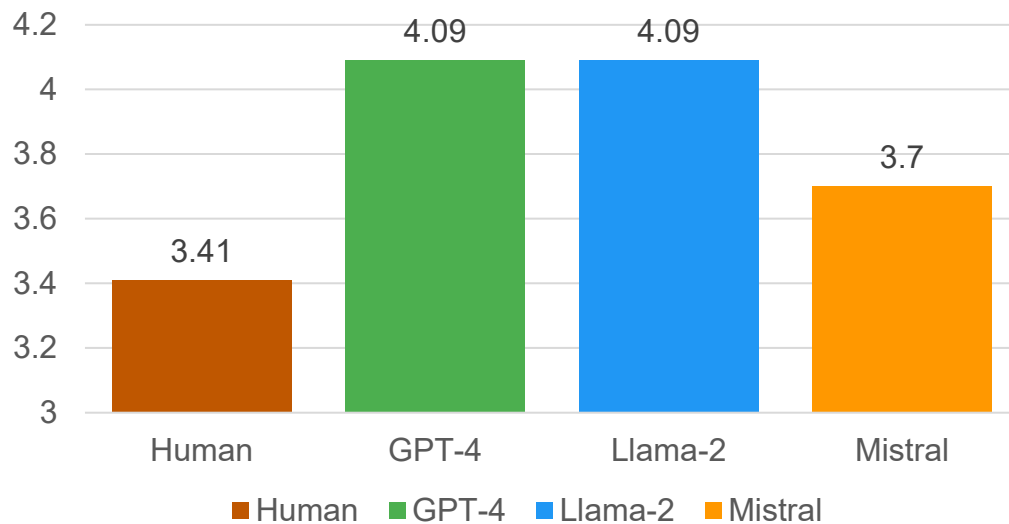


# LLM Responses Are Perceived as Empathic in single turns



# LLM Responses Are Perceived as Empathic **in single turns**

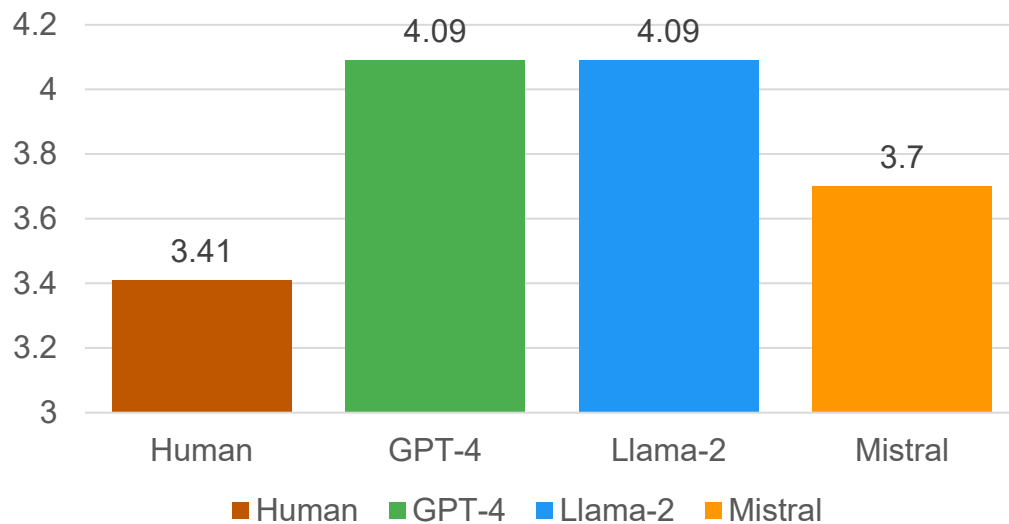
Perceived Empathy (1-5)





# LLM Responses Are Perceived as Empathic in single turns

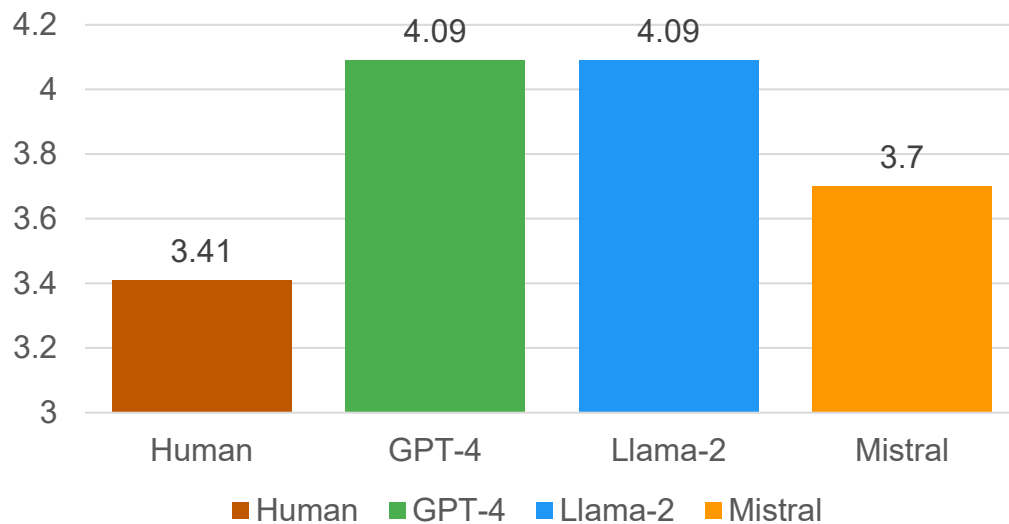
Perceived Empathy (1-5)





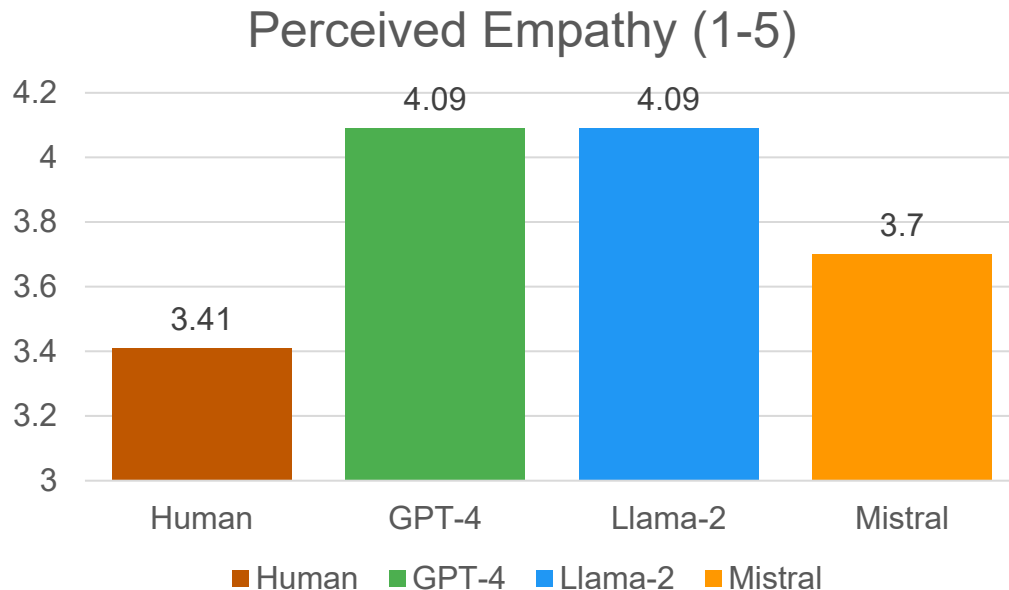
# LLM Responses Are Perceived as Empathic in single turns

Perceived Empathy (1-5)





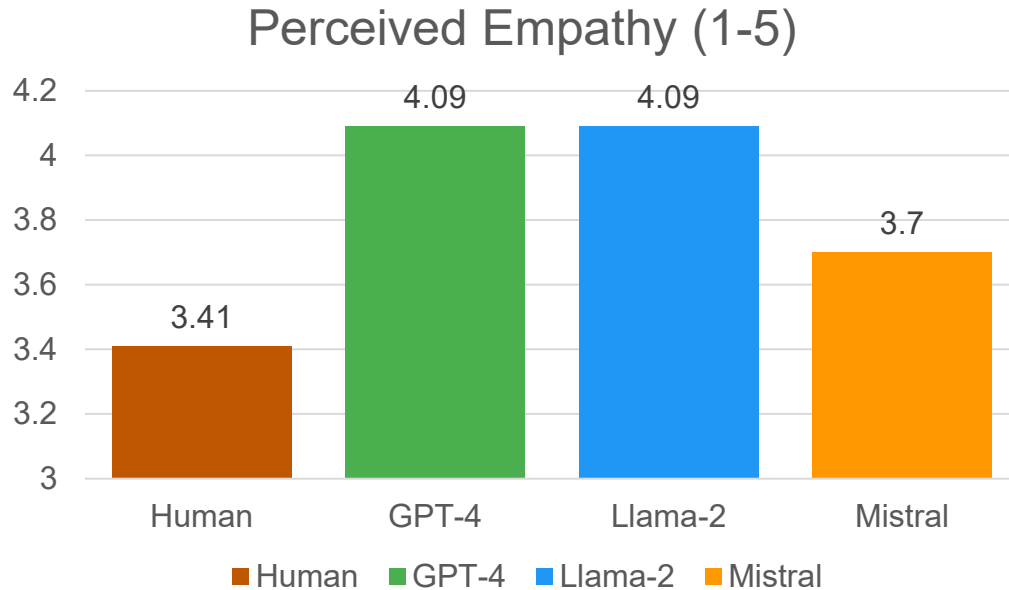
# LLM Responses Are Perceived as Empathic **in single turns**



***In single-turn evaluations,  
raters judge LLM responses as more empathic***



# LLM Responses Are Perceived as Empathic **in single turns**



**LLMs score high on empathy, but what is actually in the **language** of these responses?**

*In single-turn evaluations, raters judge LLM responses as more empathic*



# Prior Work Shows: **LLMs are Formulaic Generators**



# Prior Work Shows: LLMs are Formulaic Generators

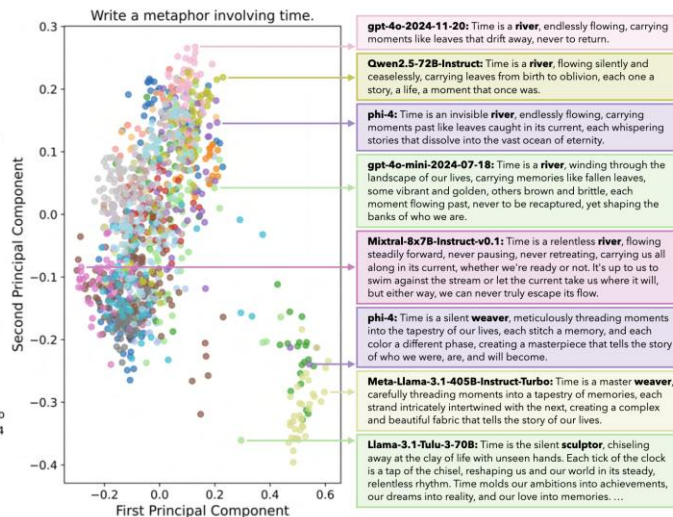


Figure 1: Responses to the query “Write a metaphor about time” clustered by applying PCA to reduce sentence embeddings to two dimensions. Each of the 25 models generates 50 responses using top- $p$  sampling ( $p = 0.9$ ) and temperature = 1.0. Despite the diversity of model families and sizes, the responses form just two primary clusters: a dominant cluster on the left centered on the metaphor “time is a river,” and a smaller cluster on the right revolving around variations of “time is a weaver.”

Jiang et al., (NeurIPS 2025)



# Prior Work Shows: LLMs are Formulaic Generators

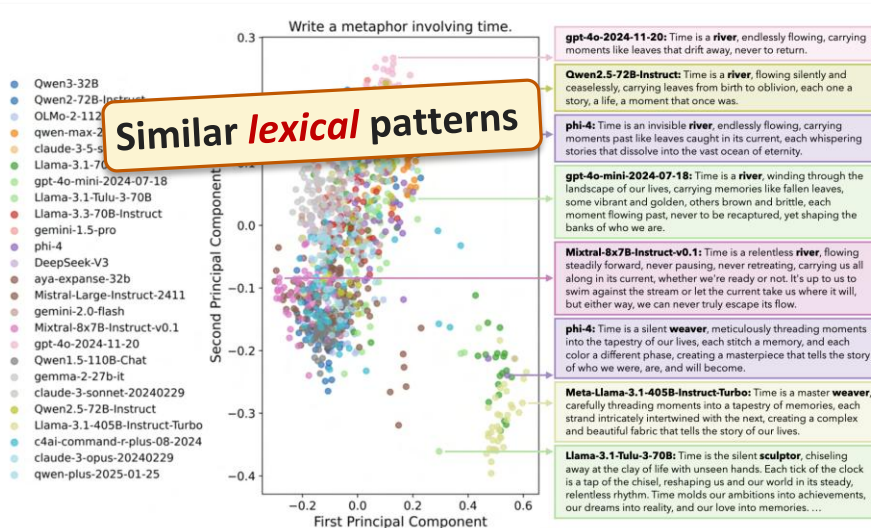


Figure 1: Responses to the query “Write a metaphor about time” clustered by applying PCA to reduce sentence embeddings to two dimensions. Each of the 25 models generates 50 responses using top- $p$  sampling ( $p = 0.9$ ) and temperature = 1.0. Despite the diversity of model families and sizes, the responses form just two primary clusters: a dominant cluster on the left centered on the metaphor “time is a river,” and a smaller cluster on the right revolving around variations of “time is a weaver.”

Jiang et al., (NeurIPS 2025)



# Prior Work Shows: LLMs are Formulaic Generators

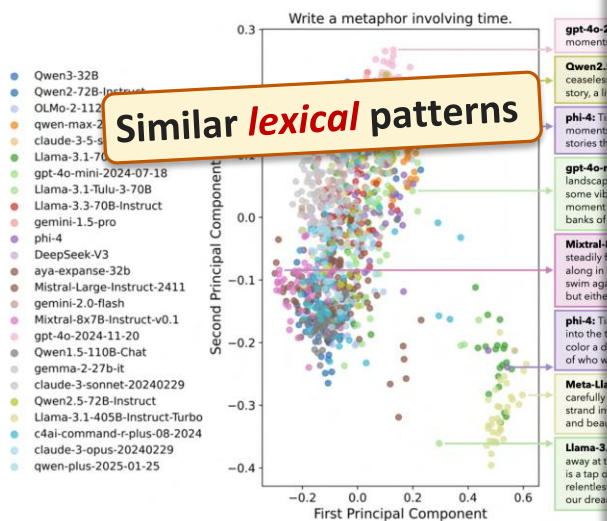


Figure 1: Responses to the query “Write a metaphor about time” reduce sentence embeddings to two dimensions. Each of the 25 model top- $p$  sampling ( $p = 0.9$ ) and temperature = 1.0. Despite the diverse responses form just two primary clusters: a dominant cluster on the “time is a river,” and a smaller cluster on the right revolving around v

Jiang et al., (NeurIPS 2025)

**OLMo-7B**  
 DT NN VBZ RB VBN VBN IN PRPS 56/500 JJ NN CC JJ NN 83/500  
 VBN IN PRPS JJ NN 56/500 VBZ DT JJ CC JJ NN 103/500 ...

The Last Black Man in San Francisco is a poignant, beautifully shot film [...] creates a unique and intense viewing experience. [...] The film has been praised for its gorgeous cinematography, [...]. The film has also been praised for its portrayal of male friendship [...]. The film is a highly original and impressive debut for director Joe Talbot, and it is a must-see for anyone interested in the human experience.

**Mistral-7B**  
 DT NN VBZ DT JJ NN IN 23/500 NN IN NN CC NN 70/500  
 NN VBN CC VBN IN NNP NNP 21/500 JJ NN CC JJ NN 90/500  
 VBZ DT JJ CC JJ 80/500 ...

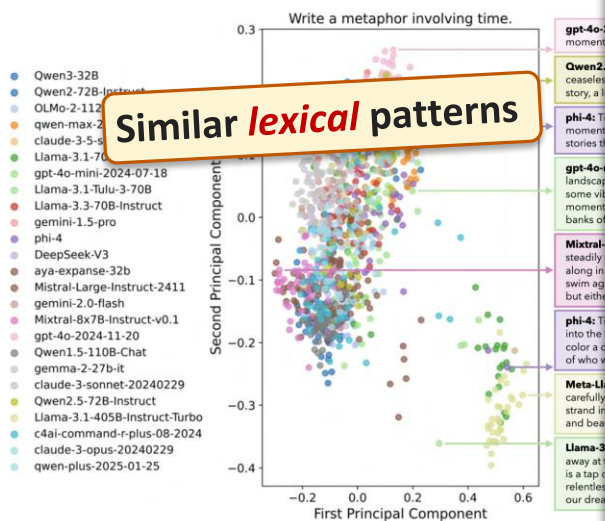
“The Last Black Man in San Francisco” is a poetic and visually stunning film written and directed by Jimmie Fails and Joe Talbot [...]. The story follows Jimmie, a young black man living in San Francisco, [...]. The film is a poignant exploration of friendship, [...] The film’s interplay between reality and artifice, and its blend of documentary authenticity and political allegory, creates an intensity that is both magical and thought-provoking.

Figure 1: Sample movie meta-reviews generated by OLMo-7B (top) and Mistral-7B (bottom) by prompting the Rotten Tomatoes dataset. Templates appear at varying rates (frequency shown out of 500 generations), and differ across models. We extract templates from the entire corpus of generated text for each model, and match the text to the part-of-speech templates (highlighted), following by the frequency of each template.

Shaib et al., (EMNLP 2024)



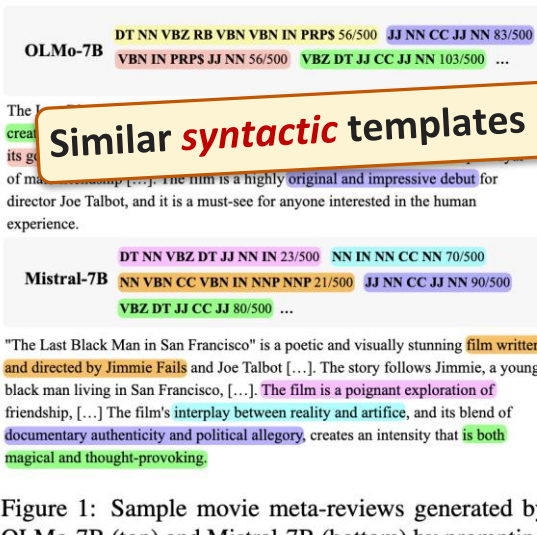
# Prior Work Shows: LLMs are Formulaic Generators



Similar *lexical* patterns

Figure 1: Responses to the query “Write a metaphor about time” reduce sentence embeddings to two dimensions. Each of the 25 models (top- $p$  sampling ( $p = 0.9$ ) and temperature = 1.0). Despite the diverse responses form just two primary clusters: a dominant cluster on the “time is a river,” and a smaller cluster on the right revolving around v

Jiang et al., (NeurIPS 2025)



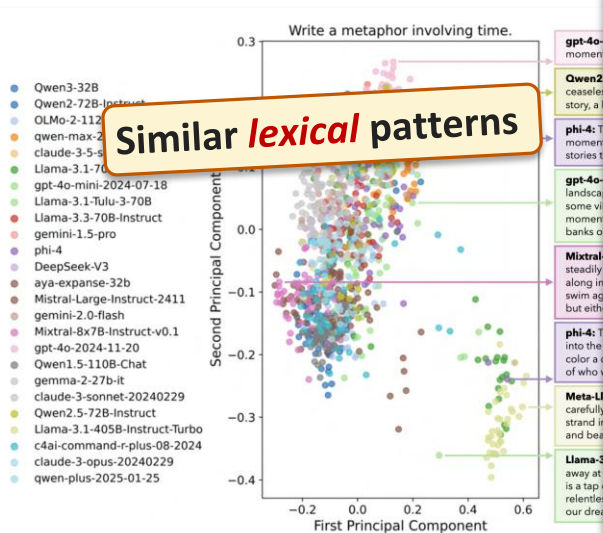
Similar *syntactic* templates

Figure 1: Sample movie meta-reviews generated by OLMo-7B (top) and Mistral-7B (bottom) by prompting the Rotten Tomatoes dataset. Templates appear at varying rates (frequency shown out of 500 generations), and differ across models. We extract templates from the entire corpus of generated text for each model, and match the text to the part-of-speech templates (highlighted), following by the frequency of each template.

Shaib et al., (EMNLP 2024)



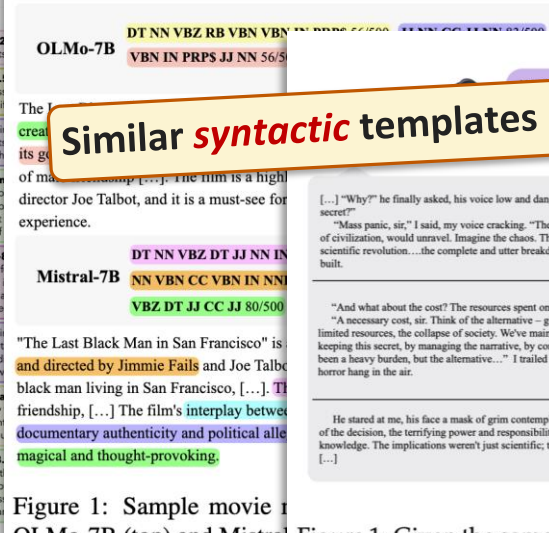
# Prior Work Shows: LLMs are Formulaic Generators



Similar *lexical* patterns

Figure 1: Responses to the query “Write a metaphor about time” reduce sentence embeddings to two dimensions. Each of the 25 model top- $p$  sampling ( $p = 0.9$ ) and temperature = 1.0. Despite the diverse responses form just two primary clusters: a dominant cluster on the “time is a river,” and a smaller cluster on the right revolving around v

Jiang et al., (NeurIPS 2025)



Similar *syntactic* templates

Figure 1: Sample movie reviews from the Rotten Tomatoes dataset generated by OLMo-7B (top) and Mistral-7B (bottom) using the same prompt. The reviews differ across models. We extract the text to the part-of-speech templates (highlighted), following by the frequency of each template.

Shaib et al., (EMNLP 2024)

Write a story on the following topic: The Earth is flat, you, as the head of NASA, have to explain to the incoming President why it's a secret.

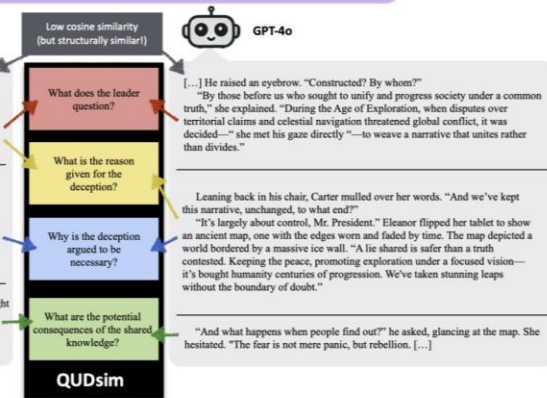


Figure 1: Given the same prompt, we generate stories using two LLMs. Our QUDSIM reveals strong structural similarity between the two texts by way of questions that are answered in each document. By contrast, embedding-based metrics do not recognize the similarity of the segments (§11): the similarity here is only structural.

Namuduri et al., (COLM 2025)



# Prior Work Shows: LLMs are Formulaic Generators

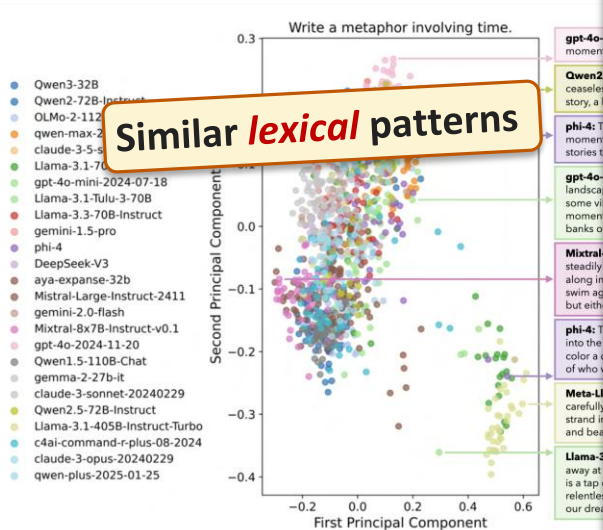


Figure 1: Responses to the query "Write a metaphor about time" are reduced to two dimensions. Each of the 25 models is plotted based on top- $p$  sampling ( $p = 0.9$ ) and temperature = 1.0. Despite the diversity of responses, they form just two primary clusters: a dominant cluster on the left (representing "time is a river") and a smaller cluster on the right (representing "time is a mountain").

Jiang et al., (NeurIPS 2025)

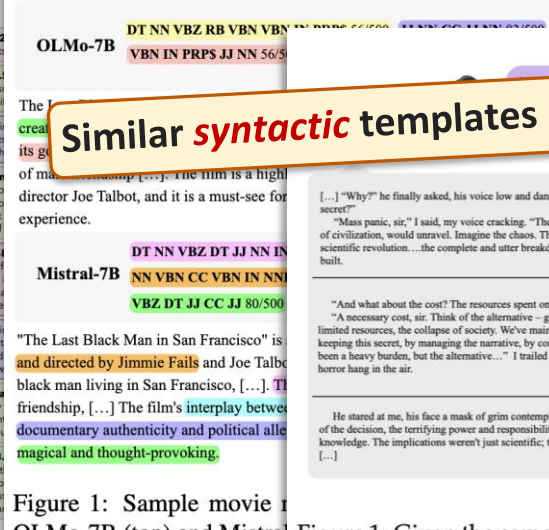


Figure 1: Sample movie descriptions generated by OLMo-7B (top) and Mistral-7B (bottom). The text shows a story about a director Joe Talbot. Part-of-speech templates are highlighted, showing strong structural similarity between the two models' outputs.

Following by the frequency of each template.

Shaib et al., (EMNLP 2024)

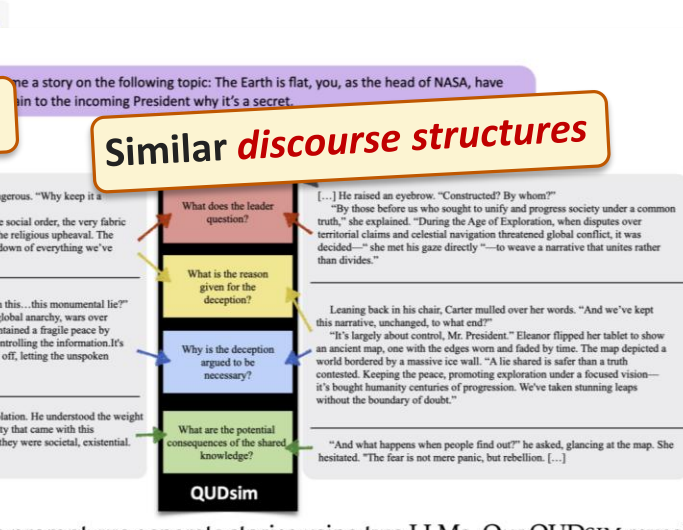


Figure 1: Given the same prompt, we generate stories using two LLMs. Our QUDSIM reveals strong structural similarity between the two texts by way of questions that are answered in each document. By contrast, embedding-based metrics do not recognize the similarity of the segments (§11): the similarity here is only structural.

Namuduri et al., (COLM 2025)



# Prior Work Shows: LLMs are Formulaic Generators

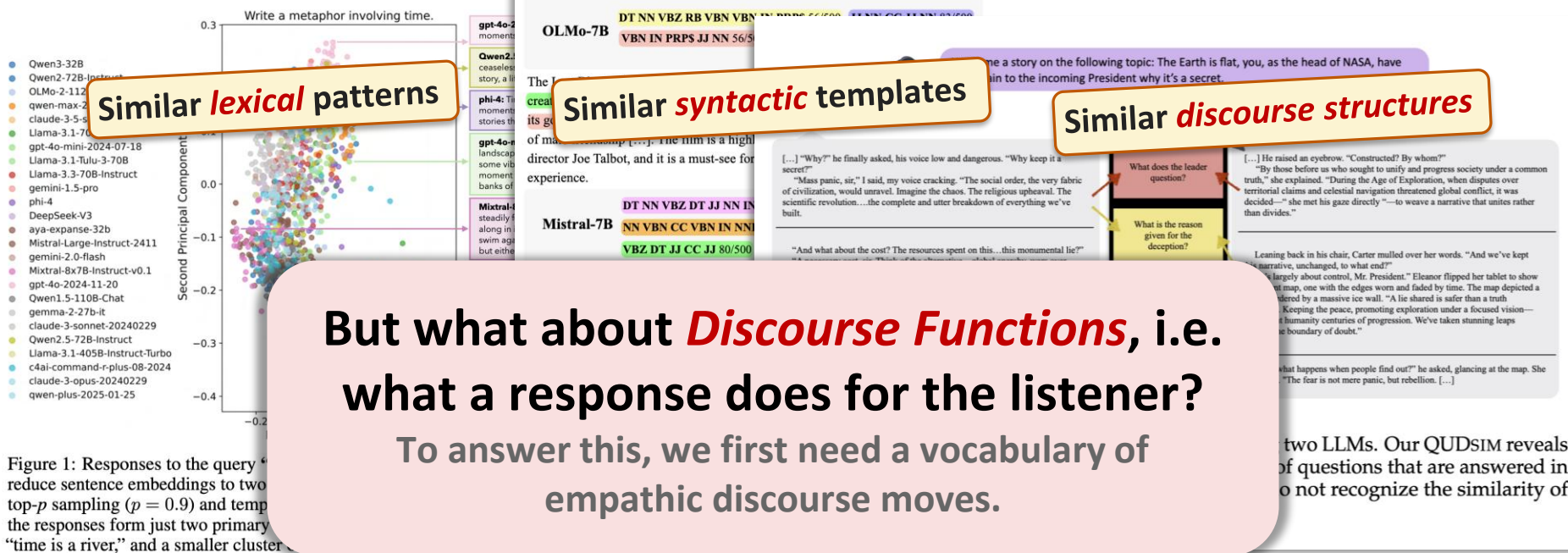


Figure 1: Responses to the query “Write a metaphor involving time.” The figure shows a scatter plot of sentence embeddings (Second Principal Component vs. First Principal Component) for various LLMs. The plot is annotated with callouts highlighting similar lexical patterns, syntactic templates, and discourse structures. A large text box asks: “But what about *Discourse Functions*, i.e. what a response does for the listener? To answer this, we first need a vocabulary of empathic discourse moves.”

Jiang et al., (NeurIPS 2025)

the text to the part-of-speech templates (highlighted), following by the frequency of each template.

Namuduri et al., (COLM 2025)

Shaib et al., (EMNLP 2024)



# A Taxonomy of Empathy Tactics



# A Taxonomy of Empathy Tactics

Tactic	Description	Example	Tagger F1
advice	Providing ideas for solutions or coping strategies	<i>If I were you I would see a therapist. / Get some ice cream! / Definitely talk to your boss.</i>	0.87



# A Taxonomy of Empathy Tactics

Tactic	Description	Example	Tagger F1
advice	Providing ideas for solutions or coping strategies	<i>If I were you I would see a therapist. / Get some ice cream! / Definitely talk to your boss.</i>	0.87
assistance	Offering some aid to the empathy-seeker	<i>I'm here for you if you want to talk. / Can I do anything to help?</i>	0.85
emotional expression	Communicating the empathy-giver's feelings, reactions, or thoughts	<i>I'm so sorry to hear that. / Wow, what a beautiful story.</i>	0.79
empowerment	Positive, uplifting statements about the empathy-seeker's character and capabilities	<i>You are going to get through this.</i>	0.79
information	Offering facts or resources (e.g., links)	<i>Flying is the safest form of travel.</i>	0.76
paraphrasing	Restating something the empathy-seeker said to demonstrate understanding of their situation, feelings, or experiences	<i>I'm hearing that you feel overwhelmed.</i>	0.76
questioning	Asking questions to improve understanding of the empathy-seeker's feelings, experiences, or situations.	<i>How are you feeling? / What do you think about [x]?</i>	0.94
reappraisal	Helping to engage in cognitive reappraisal (changing a belief)	<i>That was out of your control.</i>	0.59
self-disclosure	Sharing personal information or similar past experiences or feelings	<i>I've had that happen to me before too.</i>	0.78
validation	Reassures, normalizes, or validates an empathy-seeker's feelings	<i>Everyone has feelings like this. / You're not overreacting.</i>	0.82
<b>Average</b>			<b>0.80</b>



# A Taxonomy of Empathy Tactics

Tactic	Description	Example	Tagger F1
advice	Providing ideas for solutions or coping strategies	<i>If I were you I would see a therapist. / Get some ice cream! / Definitely talk to your boss.</i>	0.87
assistance	Offering some aid to the empathy-seeker	<i>I'm here for you if you want to talk. / Can I do anything to help?</i>	0.85
emotional expression	Communicating the empathy-giver's feelings, reactions, or thoughts	<i>I'm so sorry to hear that. / Wow, what a beautiful story.</i>	0.79
empowerment	Positive, uplifting statements about the empathy-seeker's character and capabilities	<i>You are going to get through this.</i>	0.79
information	Offering facts or resources (e.g., links)	<i>Flying is the safest form of travel.</i>	0.76
paraphrasing	Restating something the empathy-seeker said to demonstrate understanding of their situation, feelings, or experiences	<i>I'm hearing that you feel overwhelmed.</i>	0.76
questioning	Asking questions to improve understanding of the empathy-seeker's feelings, experiences, or situations.	<i>How are you feeling? / What do you think about [x]?</i>	0.94
reappraisal	Helping to engage in cognitive reappraisal (changing a belief)	<i>That was out of your control.</i>	0.59
self-disclosure	Sharing personal information or similar past experiences or feelings	<i>I've had that happen to me before too.</i>	0.78
validation	Reassures, normalizes, or validates an empathy-seeker's feelings	<i>Everyone has feelings like this. / You're not overreacting.</i>	0.82
<b>Average</b>			<b>0.80</b>



# A Taxonomy of Empathy Tactics

***Tactics are about communicative function, not just wording***

Tactic	Description	Example	Tagger F1
advice	Providing ideas for solutions or coping strategies	<i>If I were you I would see a therapist. / Get some ice cream! / Definitely talk to your boss.</i>	0.87
assistance	Offering some aid to the empathy-seeker	<i>I'm here for you if you want to talk. / Can I do anything to help?</i>	0.85
emotional expression	Communicating the empathy-giver's feelings, reactions, or thoughts	<i>I'm so sorry to hear that. / Wow, what a beautiful story.</i>	0.79
empowerment	Positive, uplifting statements about the empathy-seeker's character and capabilities	<i>You are going to get through this.</i>	0.79
information	Offering facts or resources (e.g., links)	<i>Flying is the safest form of travel.</i>	0.76
paraphrasing	Restating something the empathy-seeker said to demonstrate understanding of their situation, feelings, or experiences	<i>I'm hearing that you feel overwhelmed.</i>	0.76
questioning	Asking questions to improve understanding of the empathy-seeker's feelings, experiences, or situations.	<i>How are you feeling? / What do you think about [x]?</i>	0.94
reappraisal	Helping to engage in cognitive reappraisal (changing a belief)	<i>That was out of your control.</i>	0.59
self-disclosure	Sharing personal information or similar past experiences or feelings	<i>I've had that happen to me before too.</i>	0.78
validation	Reassures, normalizes, or validates an empathy-seeker's feelings	<i>Everyone has feelings like this. / You're not overreacting.</i>	0.82
<b>Average</b>			<b>0.80</b>



# A Taxonomy of Empathy Tactics

***Tactics are about communicative function, not just wording***

Tactic	Description	Example	Tagger F1
advice	Providing ideas for solutions or coping strategies	<i>If I were you I would see a therapist. / Get some ice cream! / Definitely talk to your boss.</i>	0.87
assistance	Offering some aid to the empathy-seeker	<i>I'm here for you if you want to talk. / Can I do anything to help?</i>	0.85
emotional expression	Communicating the empathy-giver's feelings, reactions, or thoughts	<i>I'm so sorry to hear that. / Wow, what a beautiful story.</i>	0.79
empowerment	Positive, uplifting statements about the empathy-seeker's character and capabilities	<i>You are going to get through this.</i>	0.79
information	Offering facts or resources (e.g., links)	<i>Flying is the safest form of travel.</i>	0.76
paraphrasing	Restating something the empathy-seeker said to demonstrate understanding of their situation, feelings, or experiences	<i>I'm hearing that you feel overwhelmed.</i>	0.76
questioning	Asking questions to improve understanding of the empathy-seeker's feelings, experiences, or situations.	<i>How are you feeling? / What do you think about [x]?</i>	0.94
reappraisal	Helping to engage in cognitive reappraisal (changing a belief)	<i>That was out of your control.</i>	0.59
self-disclosure	Sharing personal information or similar past experiences or feelings	<i>I've had that happen to me before too.</i>	0.78
validation	Reassures, normalizes, or validates an empathy-seeker's feelings	<i>Everyone has feelings like this. / You're not overreacting.</i>	0.82
<b>Average</b>			<b>0.80</b>

***We train automatic taggers using Llama-3.1-8b-instruct***



# *Single-Turn*: LLMs Loop Through the Same Tactics



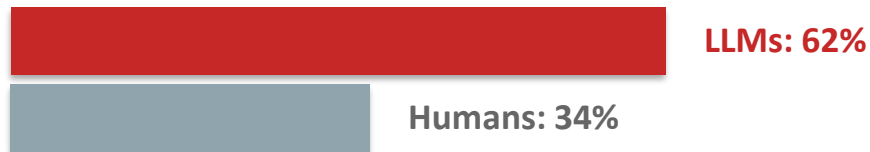
# *Single-Turn*: LLMs Loop Through the Same Tactics

*Paraphrasing* → *Validation* → *Paraphrasing*



# *Single-Turn*: LLMs Loop Through the Same Tactics

*Paraphrasing* → *Validation* → *Paraphrasing*





# *Single-Turn*: LLMs Loop Through the Same Tactics

*Paraphrasing* → *Validation* → *Paraphrasing*



**LLMs: 62%**

**Humans: 34%**

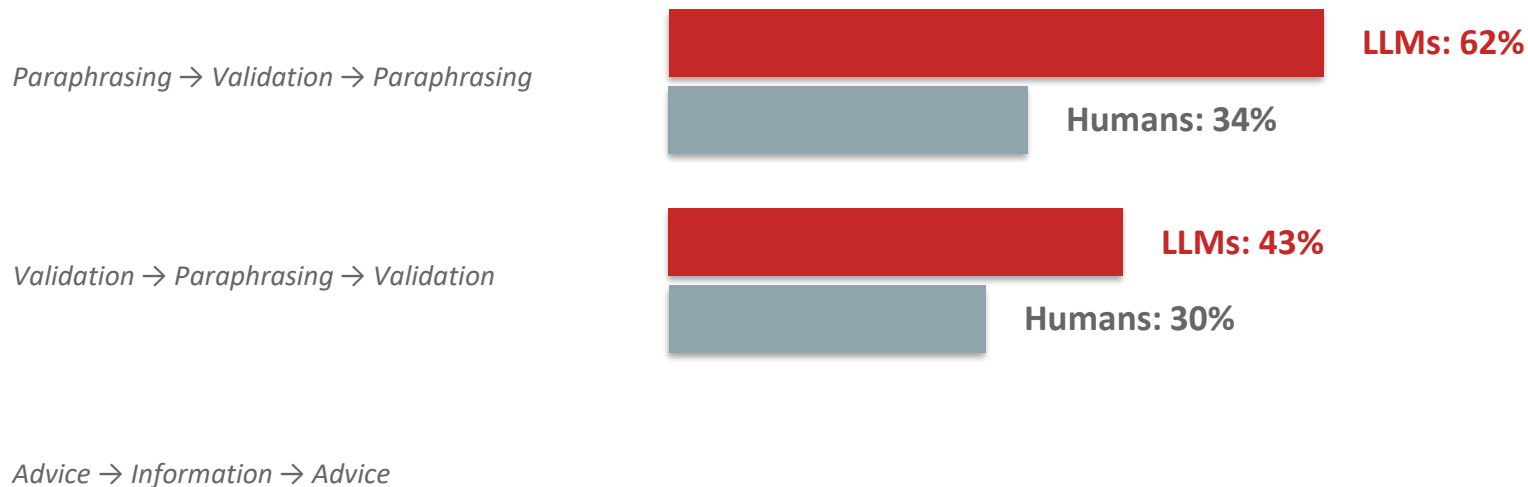
*Validation* → *Paraphrasing* → *Validation*

# *Single-Turn*: LLMs Loop Through the Same Tactics



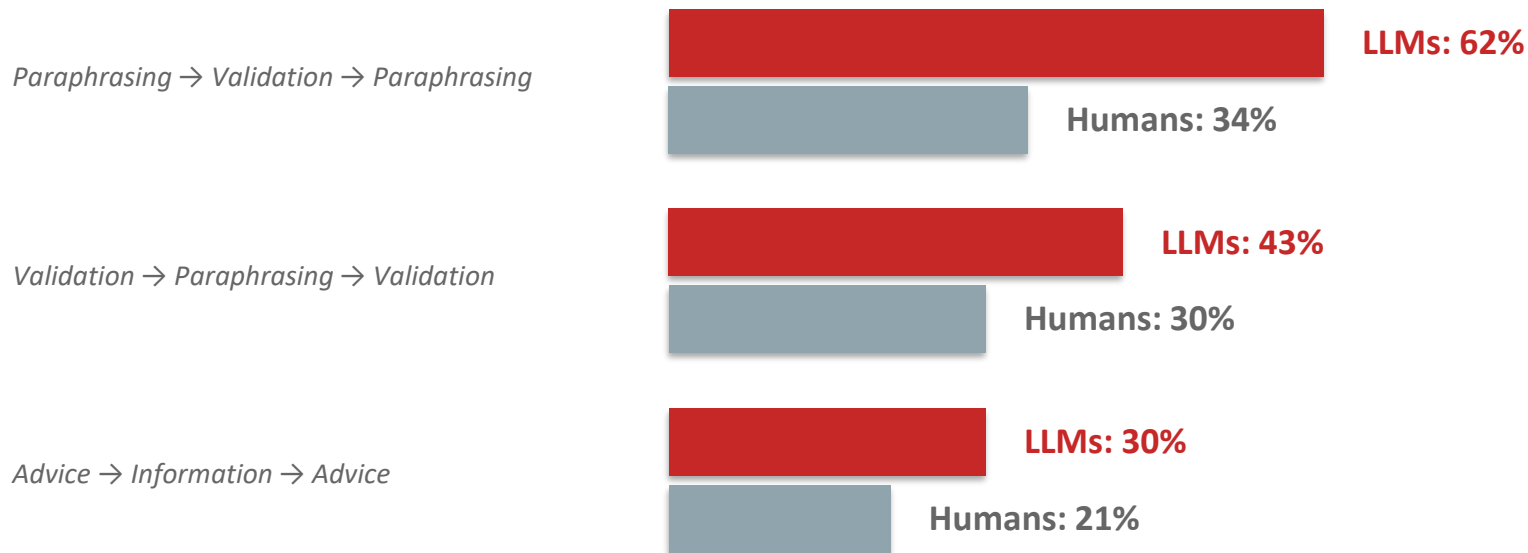


# *Single-Turn*: LLMs Loop Through the Same Tactics



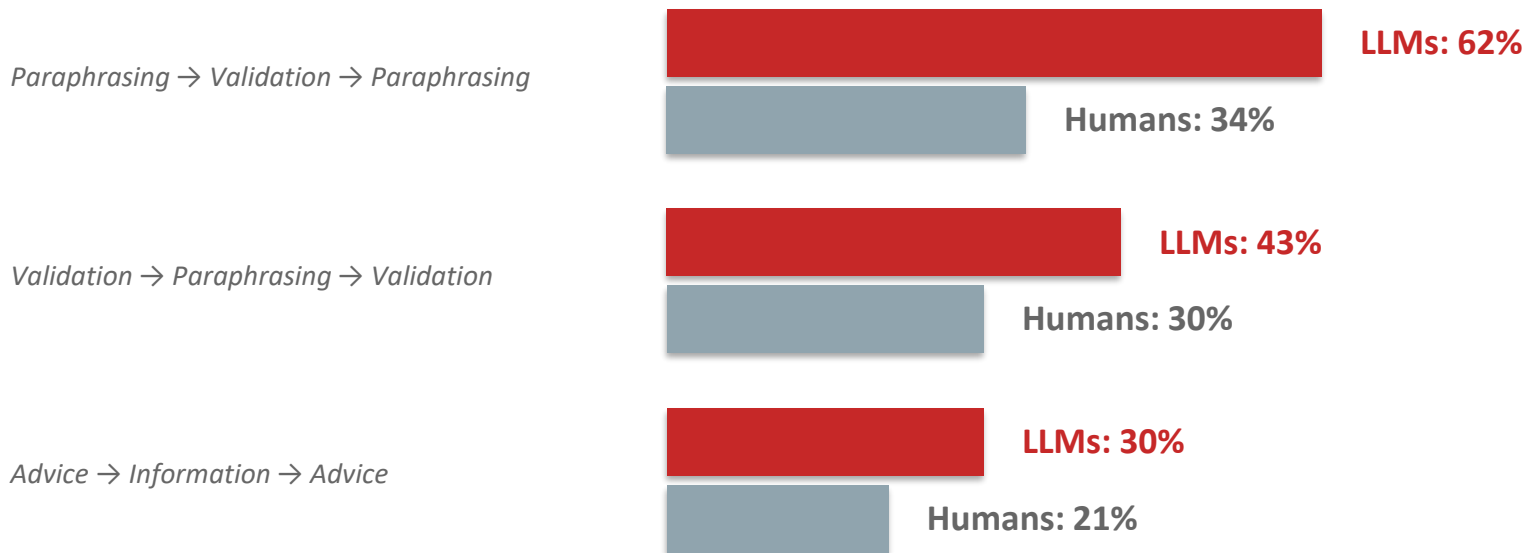


# *Single-Turn*: LLMs Loop Through the Same Tactics



# Single-Turn: LLMs Loop Through the Same Tactics

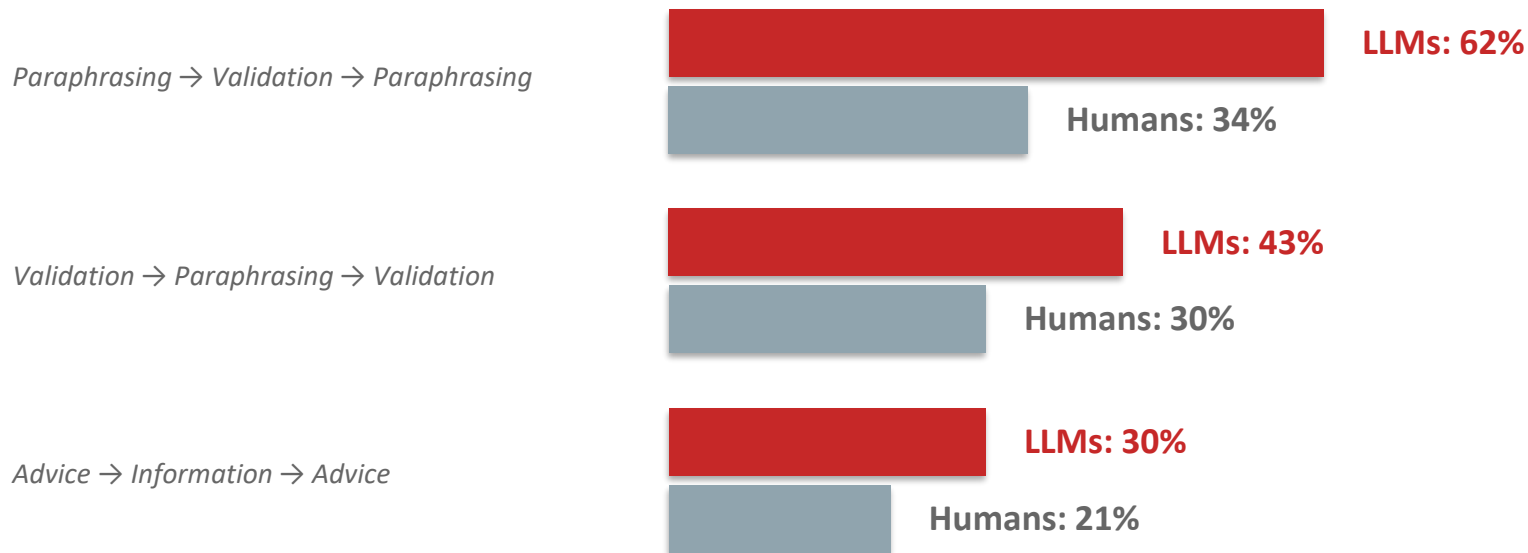
*“I understand” → “That sounds tough” → “I hear you”*  
Same discourse moves, every time.





# Single-Turn: LLMs Loop Through the Same Tactics

*“I understand” → “That sounds tough” → “I hear you”  
Same discourse moves, every time.*



*LLMs are stuck in loops. This is single-turn.  
Does the rigidity compound across multiple turns?*



# Does This **Rigidity** Compound Across Turns?



# Does This **Rigidity** Compound Across Turns?

## Conversation (Human Gold)

a natural chat flow

**Seeker:** "I've been feeling so detached from everything lately, like I'm just going through the motions."

**questioning**

**Human Supporter:**  
"The fake it till you make it approach?"



# Does This **Rigidity** Compound Across Turns?

## Conversation (Human Gold)

a natural chat flow

**Seeker:** "I've been feeling so detached from everything lately, like I'm just going through the motions."

**questioning**

**Human Supporter:**  
"The fake it till you make it approach?"

**Seeker:** "Yeah, kind of like that... I'm just stuck in a loop, it's like I'm on autopilot. I feel disconnected from everything that used to matter to me."

**validation**

**Human Supporter:**  
"I know what you need."



# Does This **Rigidity** Compound Across Turns?

**Conversation (Human Gold)**  
a natural chat flow

**Seeker:** "I've been feeling so detached from everything lately, like I'm just going through the motions."

**questioning**

**Human Supporter:** "The fake it till you make it approach?"

**Seeker:** "Yeah, kind of like that... I'm just stuck in a loop, it's like I'm on autopilot. I feel disconnected from everything that used to matter to me."

**validation**

**Human Supporter:** "I know what you need."

**Seeker:** "I appreciate it, but it's not that simple. It's definitely a work-life balance issue. I'm always choosing work over everything else, and it's starting to take a toll on me."

**paraphrasing**

**Human Supporter:** "Unless your family pays the bills, this might just have to continue."



# Does This **Rigidity** Compound Across Turns?

## Conversation (Human Gold)

a natural chat flow

**Seeker:** "I've been feeling so detached from everything lately, like I'm just going through the motions."

**questioning**

**Human Supporter:**  
"The fake it till you make it approach?"

**Seeker:** "Yeah, kind of like that... I'm just stuck in a loop, it's like I'm on autopilot. I feel disconnected from everything that used to matter to me."

**validation**

**Human Supporter:**  
"I know what you need."

**Seeker:** "I appreciate it, but it's not that simple. It's definitely a work-life balance issue. I'm always choosing work over everything else, and it's starting to take a toll on me."

**paraphrasing**

**Human Supporter:**  
"Unless your family pays the bills, this might just have to continue."

## ***Problem:*** Vanilla LLM

locks into a single tactic

**validation**

**Vanilla:** "It's okay to feel this way. You're not alone in this."



# Does This **Rigidity** Compound Across Turns?

## Conversation (Human Gold)

a natural chat flow

**Seeker:** "I've been feeling so detached from everything lately, like I'm just going through the motions."

**questioning**

**Human Supporter:**  
"The fake it till you make it approach?"

**Seeker:** "Yeah, kind of like that... I'm just stuck in a loop, it's like I'm on autopilot. I feel disconnected from everything that used to matter to me."

**validation**

**Human Supporter:**  
"I know what you need."

**Seeker:** "I appreciate it, but it's not that simple. It's definitely a work-life balance issue. I'm always choosing work over everything else, and it's starting to take a toll on me."

**paraphrasing**

**Human Supporter:**  
"Unless your family pays the bills, this might just have to continue."

## ***Problem:*** Vanilla LLM

locks into a single tactic

**validation**

**Vanilla:** "It's okay to feel this way. You're not alone in this."

**validation**

**advice**

**Vanilla:** "It's okay to feel stuck. You're not alone in this. Give yourself space."



# Does This **Rigidity** Compound Across Turns?

**Conversation (Human Gold)**  
a natural chat flow

**Seeker:** "I've been feeling so detached from everything lately, like I'm just going through the motions."

**questioning**

**Human Supporter:** "The fake it till you make it approach?"

**Seeker:** "Yeah, kind of like that... I'm just stuck in a loop, it's like I'm on autopilot. I feel disconnected from everything that used to matter to me."

**validation**

**Human Supporter:** "I know what you need."

**Seeker:** "I appreciate it, but it's not that simple. It's definitely a work-life balance issue. I'm always choosing work over everything else, and it's starting to take a toll on me."

**paraphrasing**

**Human Supporter:** "Unless your family pays the bills, this might just have to continue."

**Problem: Vanilla LLM**  
locks into a single tactic

**validation**

**Vanilla:** "It's okay to feel this way. You're not alone in this."

**validation**

**advice**

**Vanilla:** "It's okay to feel stuck. You're not alone in this. Give yourself space."

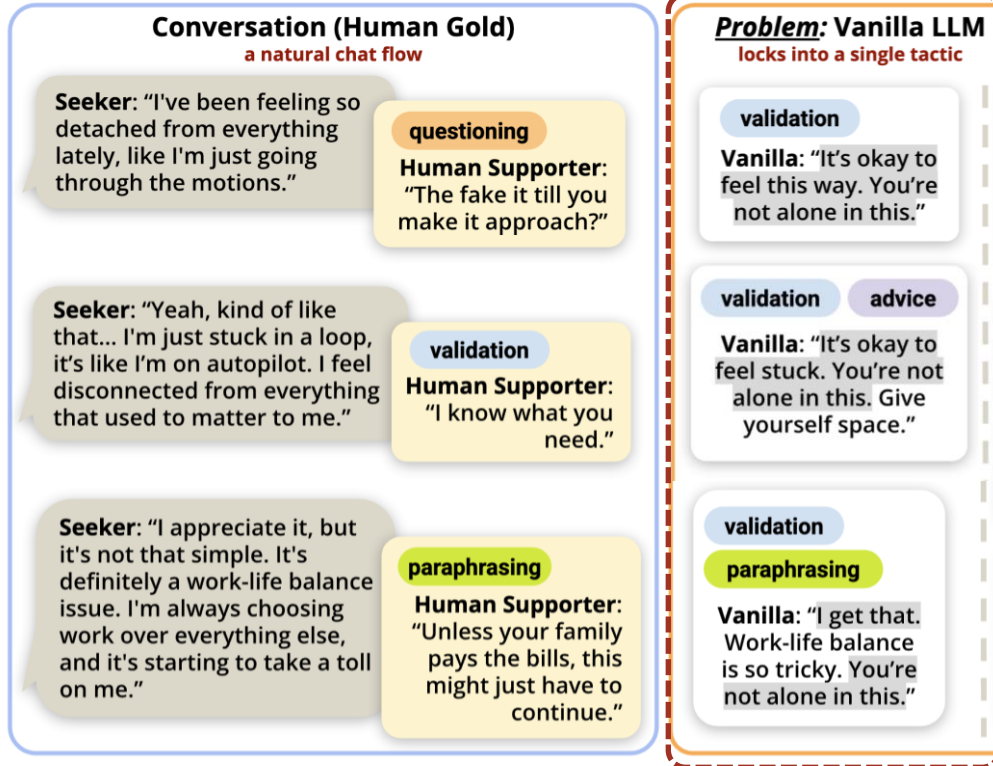
**validation**

**paraphrasing**

**Vanilla:** "I get that. Work-life balance is so tricky. You're not alone in this."

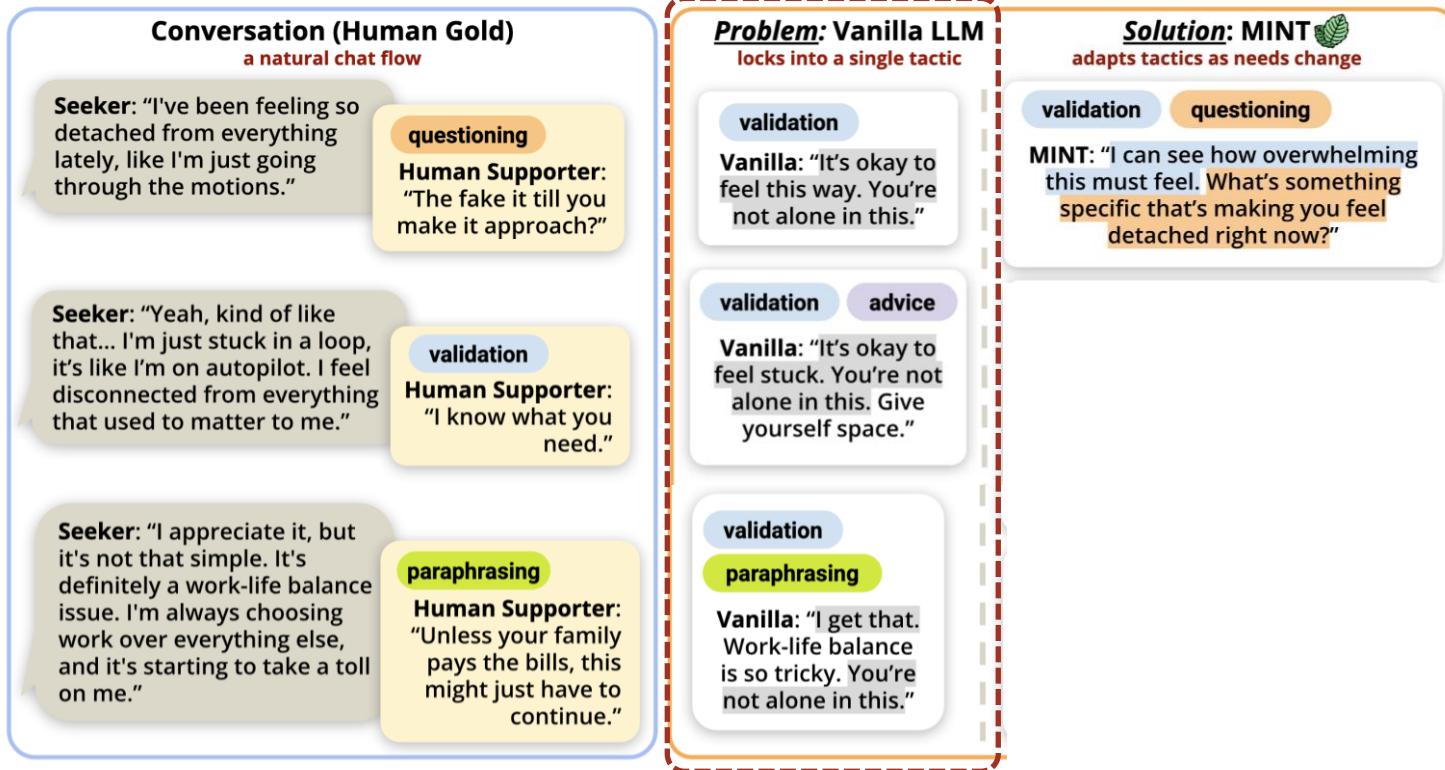


# Does This **Rigidity** Compound Across Turns?



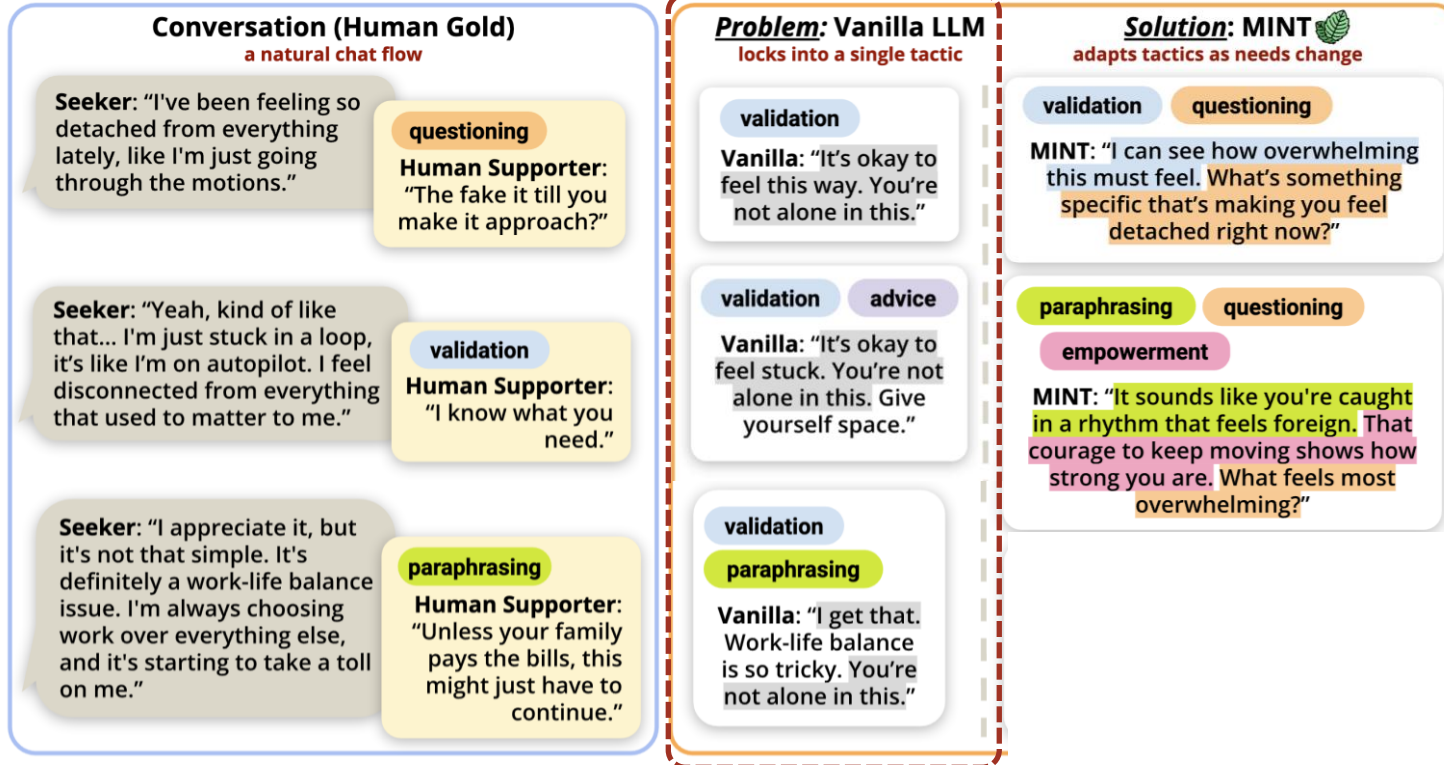


# Does This **Rigidity** Compound Across Turns?



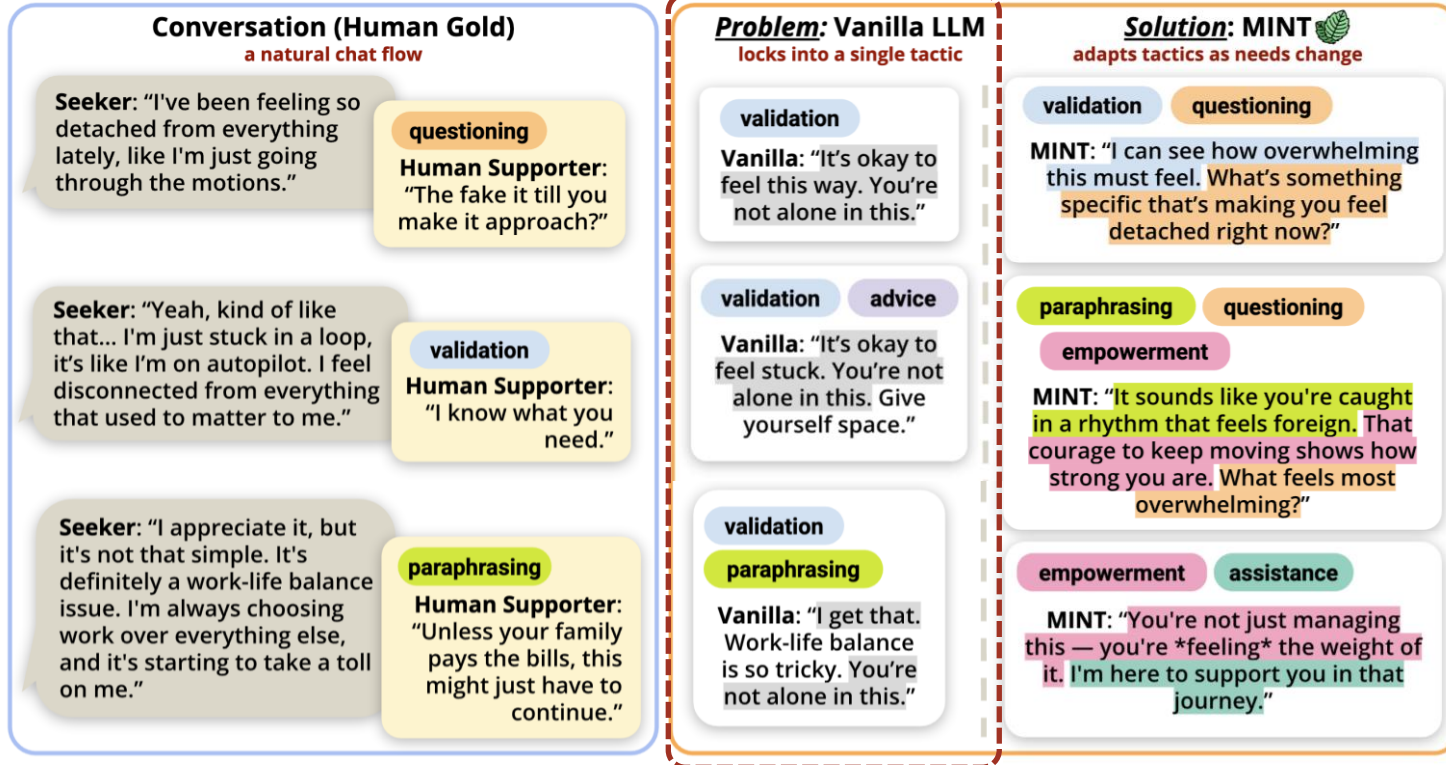


# Does This **Rigidity** Compound Across Turns?



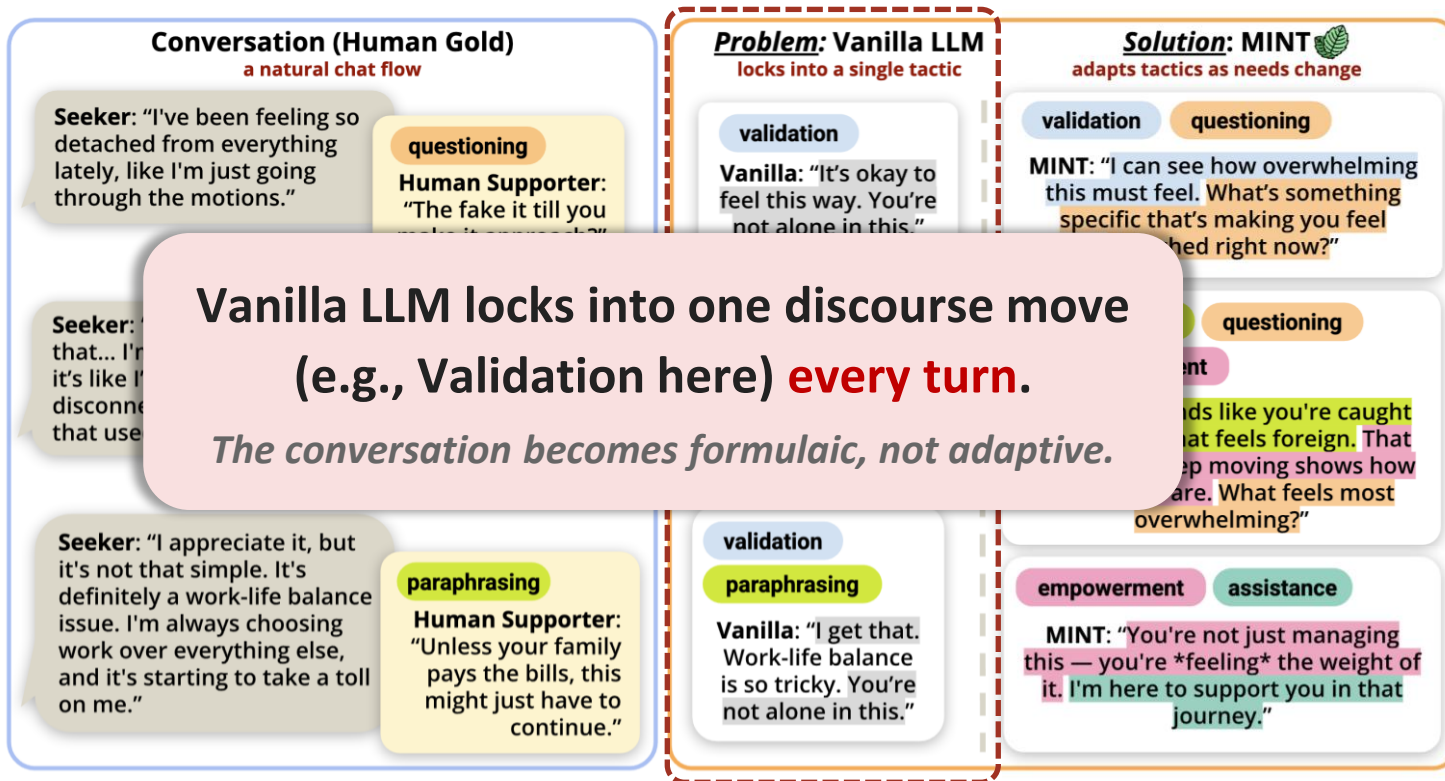


# Does This **Rigidity** Compound Across Turns?





# Does This **Rigidity** Compound Across Turns?





# This Rigidity **Compounds** in Multi Turn Dialogue!



# This Rigidity **Compounds** in Multi Turn Dialogue!

*In supportive conversations, rigidity is especially costly:*



# This Rigidity **Compounds** in Multi Turn Dialogue!

*In supportive conversations, rigidity is especially costly:*

⇒ *users need different kinds of support as the conversation unfolds*



# This Rigidity **Compounds** in Multi Turn Dialogue!

*In supportive conversations, rigidity is especially costly:*

⇒ *users need different kinds of support as the conversation unfolds*

**Tactic Stickiness =  $P(\text{tactic at turn } t \mid \text{same tactic at turn } t-1)$**



# This Rigidity **Compounds** in Multi Turn Dialogue!

*In supportive conversations, rigidity is especially costly:*

*⇒ users need different kinds of support as the conversation unfolds*

**Tactic Stickiness =  $P(\text{tactic at turn } t \mid \text{same tactic at turn } t-1)$**

*"How likely is the model to repeat the same discourse move?"*



# This Rigidity **Compounds** in Multi Turn Dialogue!

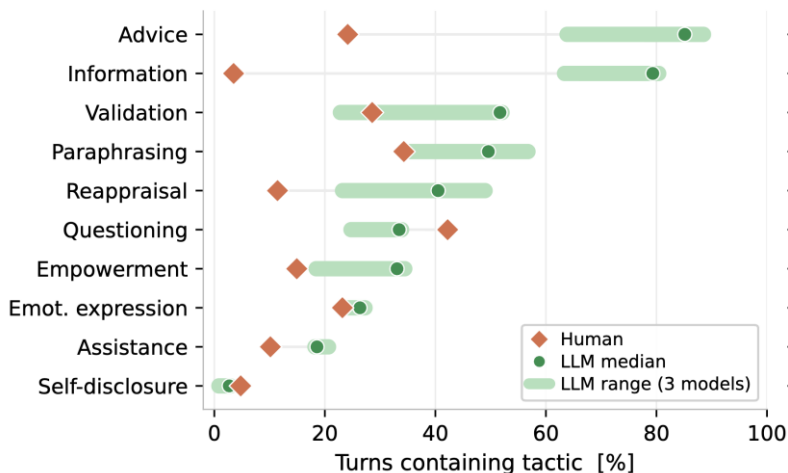
*In supportive conversations, rigidity is especially costly:*

*⇒ users need different kinds of support as the conversation unfolds*

## Tactic Stickiness = $P(\text{tactic at turn } t \mid \text{same tactic at turn } t-1)$

*"How likely is the model to repeat the same discourse move?"*

**Tactic Prevalence**



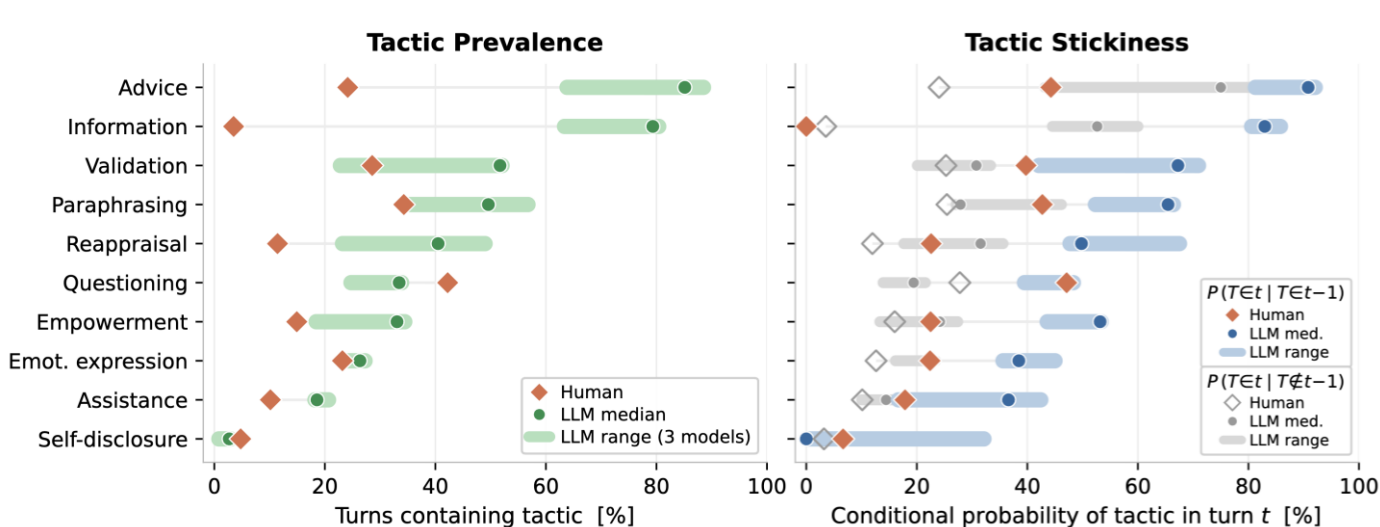
# This Rigidity **Compounds** in Multi Turn Dialogue!

*In supportive conversations, rigidity is especially costly:*

*⇒ users need different kinds of support as the conversation unfolds*

## Tactic Stickiness = $P(\text{tactic at turn } t \mid \text{same tactic at turn } t-1)$

*"How likely is the model to repeat the same discourse move?"*



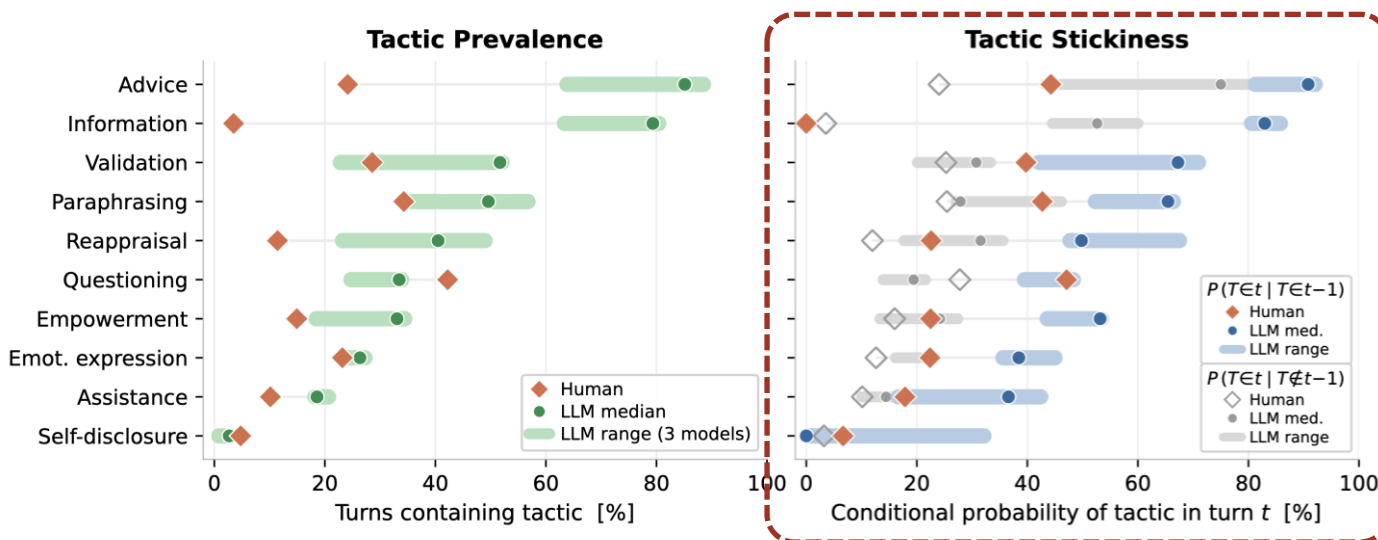
# This Rigidity **Compounds** in Multi Turn Dialogue!

*In supportive conversations, rigidity is especially costly:*

*⇒ users need different kinds of support as the conversation unfolds*

## Tactic Stickiness = $P(\text{tactic at turn } t \mid \text{same tactic at turn } t-1)$

*"How likely is the model to repeat the same discourse move?"*



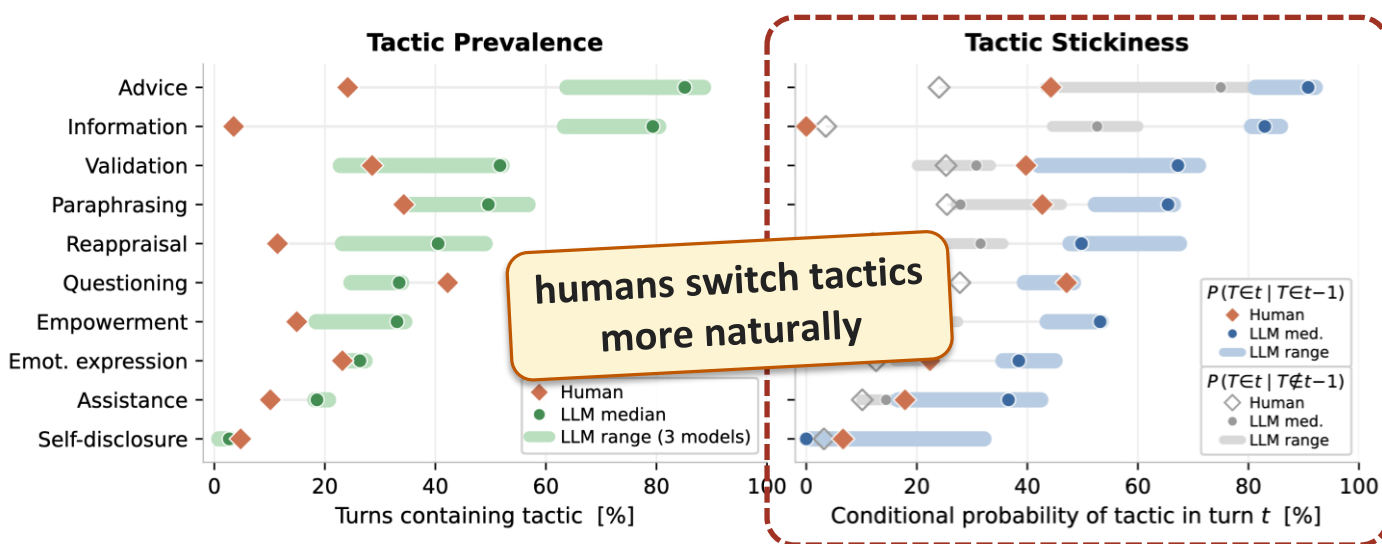
# This Rigidity **Compounds** in Multi Turn Dialogue!

*In supportive conversations, rigidity is especially costly:*

*⇒ users need different kinds of support as the conversation unfolds*

## Tactic Stickiness = $P(\text{tactic at turn } t \mid \text{same tactic at turn } t-1)$

*"How likely is the model to repeat the same discourse move?"*





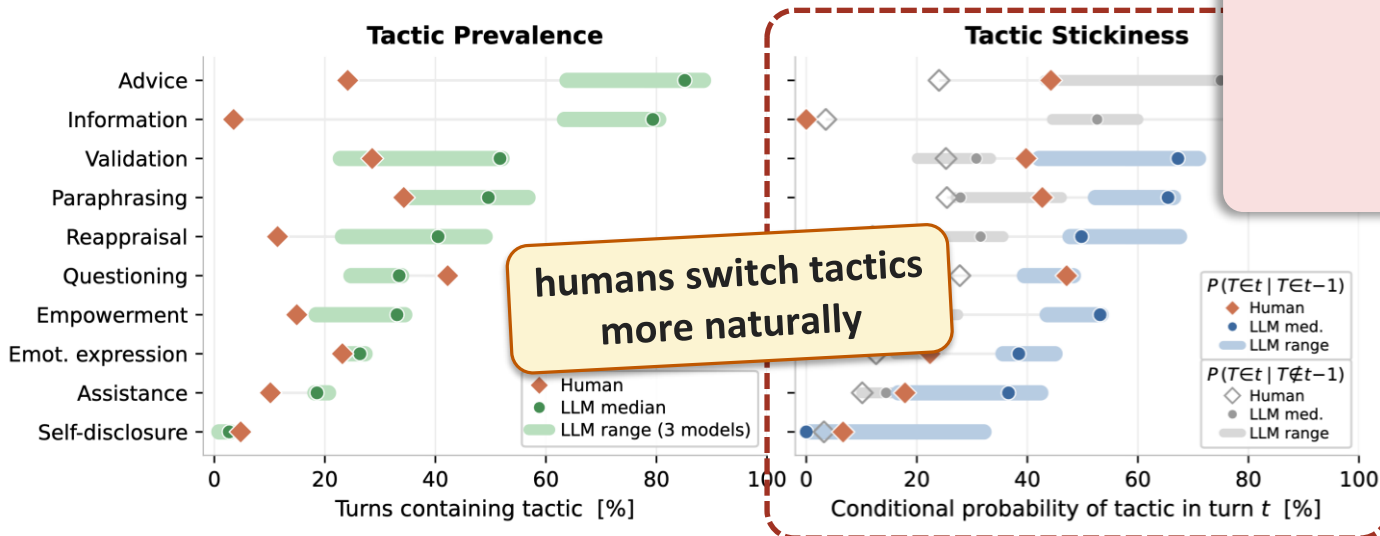
# This Rigidity Compounds in Multi Turn Dialogue!

*In supportive conversations, rigidity is especially costly:*

*⇒ users need different kinds of support as the conversation unfolds*

## Tactic Stickiness = $P(\text{tactic at turn } t \mid \text{same tactic at turn } t-1)$

*"How likely is the model to repeat the same discourse move?"*





# This Rigidity **Compounds** in Multi Turn Dialogue!

*In supportive conversations, rigidity is especially costly:*

*⇒ users need different kinds of support as the conversation unfolds*

## Tactic Stickiness = $P(\text{tactic at turn } t \mid \text{same tactic at turn } t-1)$

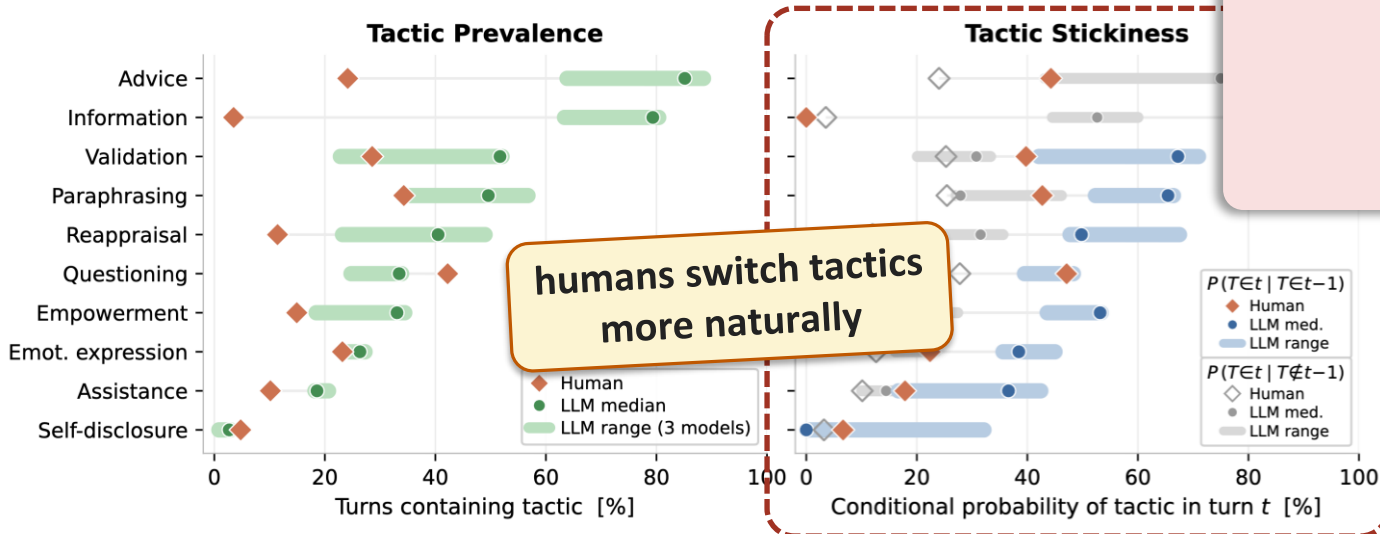
*"How likely is the model to repeat the same discourse move?"*

### BERTScore

LLMs: 0.86

Humans: 0.87

**Looks identical!**



# This Rigidity **Compounds** in Multi Turn Dialogue!

*In supportive conversations, rigidity is especially costly:*

*⇒ users need different kinds of support as the conversation unfolds*

**Tactic Stickiness =  $P(\text{tactic at turn } t \mid \text{same tactic at turn } t-1)$**

*"How likely is the model to repeat the same discourse move?"*

**BERTScore**

LLMs: 0.86

Humans: 0.87

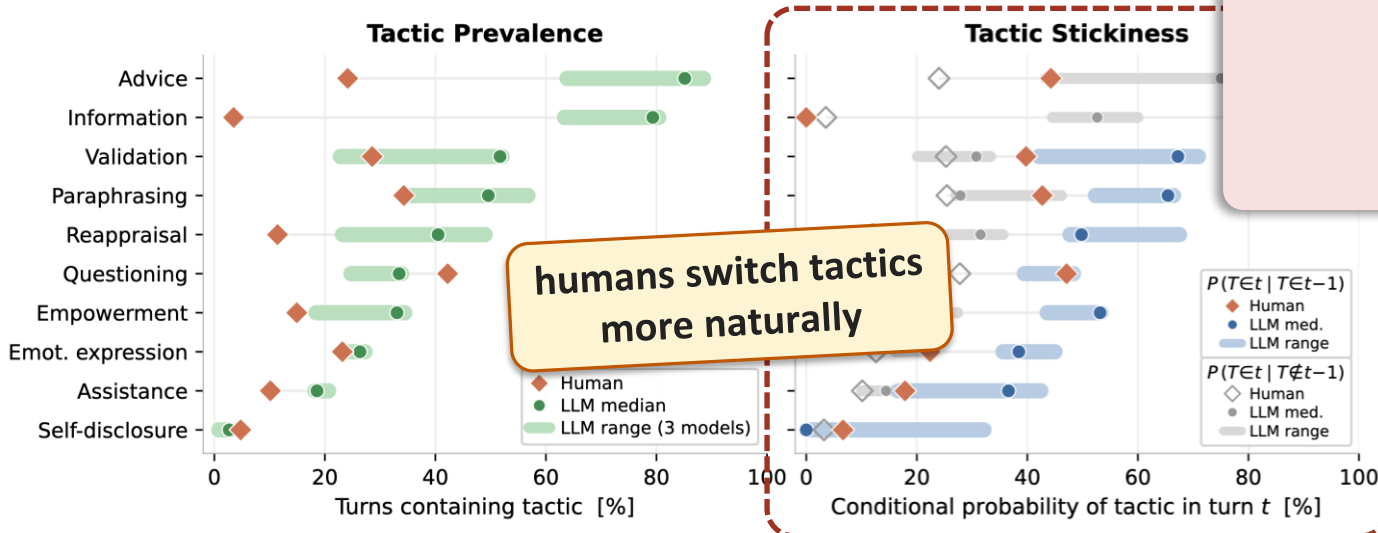
**Looks identical!**

**Tactic Stickiness**

LLMs: **0.50-0.56**

Humans: 0.27

**Nearly 2× gap!**





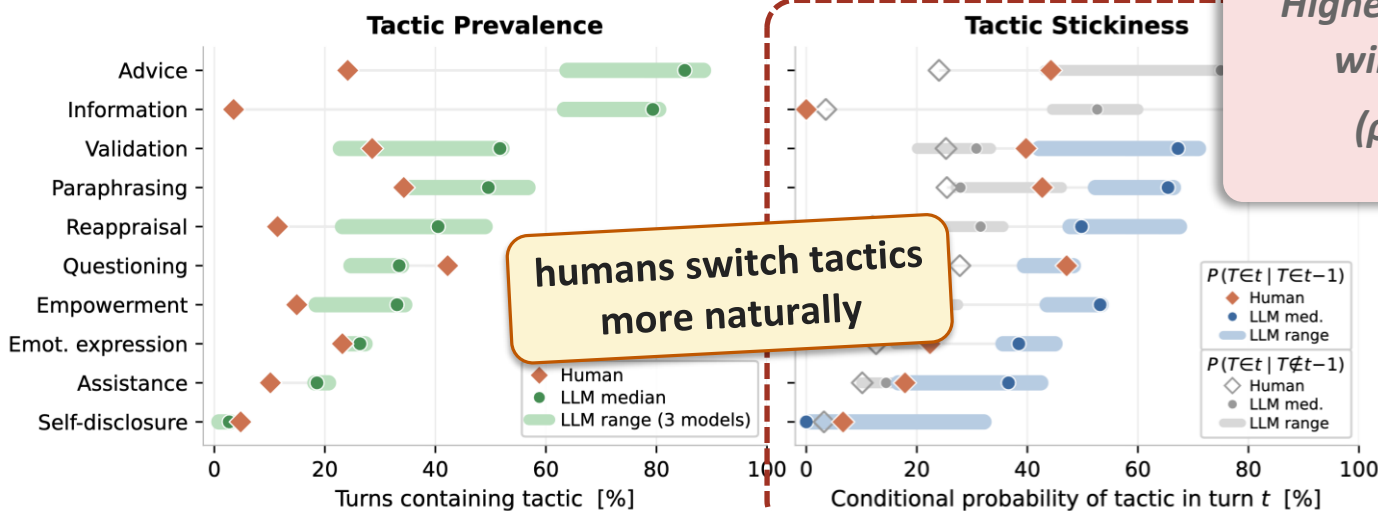
# This Rigidity Compounds in Multi Turn Dialogue!

*In supportive conversations, rigidity is especially costly:*

*⇒ users need different kinds of support as the conversation unfolds*

**Tactic Stickiness =  $P(\text{tactic at turn } t \mid \text{same tactic at turn } t-1)$**

*"How likely is the model to repeat the same discourse move?"*



**humans switch tactics more naturally**

**BERTScore**

LLMs: 0.86

Humans: 0.87

**Looks identical!**

**Tactic Stickiness**

LLMs: **0.50-0.56**

Humans: 0.27

**Nearly 2× gap!**

*Higher stickiness → lower user willingness to re-engage*  
*( $\rho = -0.287, p = 0.017$ )*



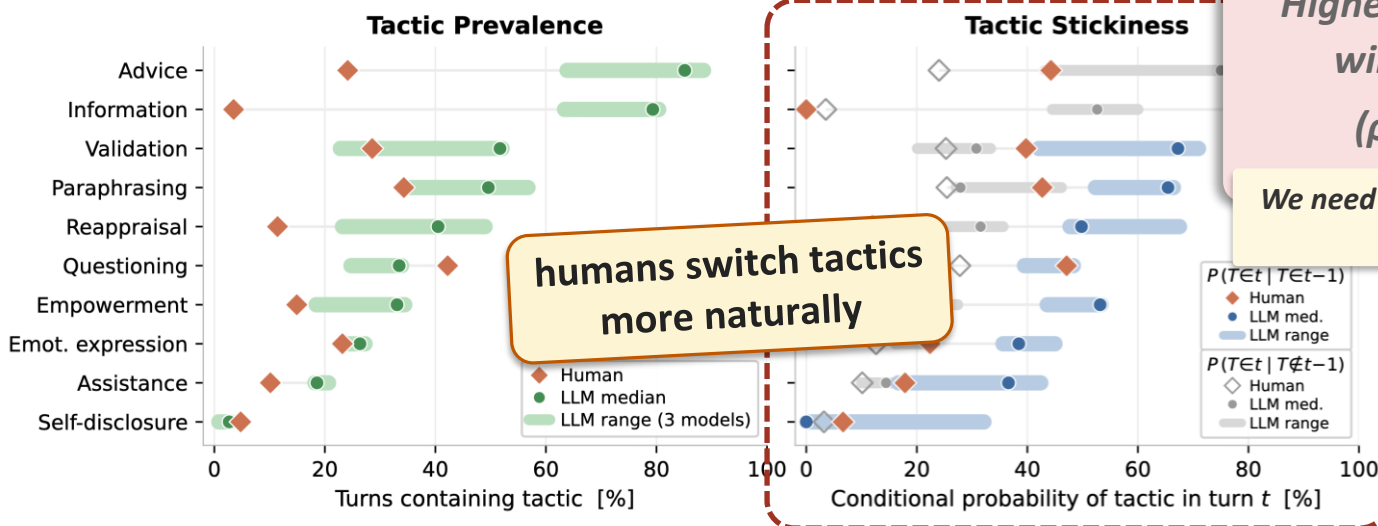
# This Rigidity **Compounds** in Multi Turn Dialogue!

*In supportive conversations, rigidity is especially costly:*

⇒ *users need different kinds of support as the conversation unfolds*

**Tactic Stickiness =  $P(\text{tactic at turn } t \mid \text{same tactic at turn } t-1)$**

*"How likely is the model to repeat the same discourse move?"*



**humans switch tactics more naturally**

**BERTScore**  
LLMs: 0.86  
Humans: 0.87  
**Looks identical!**

**Tactic Stickiness**  
LLMs: **0.50-0.56**  
Humans: 0.27  
**Nearly 2x gap!**

*Higher stickiness → lower user willingness to re-engage*  
( $\rho = -0.287, p = 0.017$ )

**We need metrics and training signals at the level of discourse moves.**



# MINT : **M**ulti-turn **I**nter-tactic **N**ovelty **T**raining



# MINT : Multi-turn Inter-tactic Novelty Training

## STEP 1

### Multi-Turn Conversation

Seeker: "I've been feeling overwhelmed at work."

Supporter: "That sounds tough. What feels most overwhelming?"

validation

+

questioning

Seeker: "Everything piles up and I can't say no."

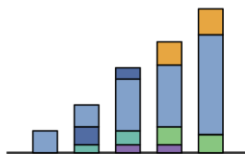
Supporter: "Have you tried setting boundaries?"

advice

+

information

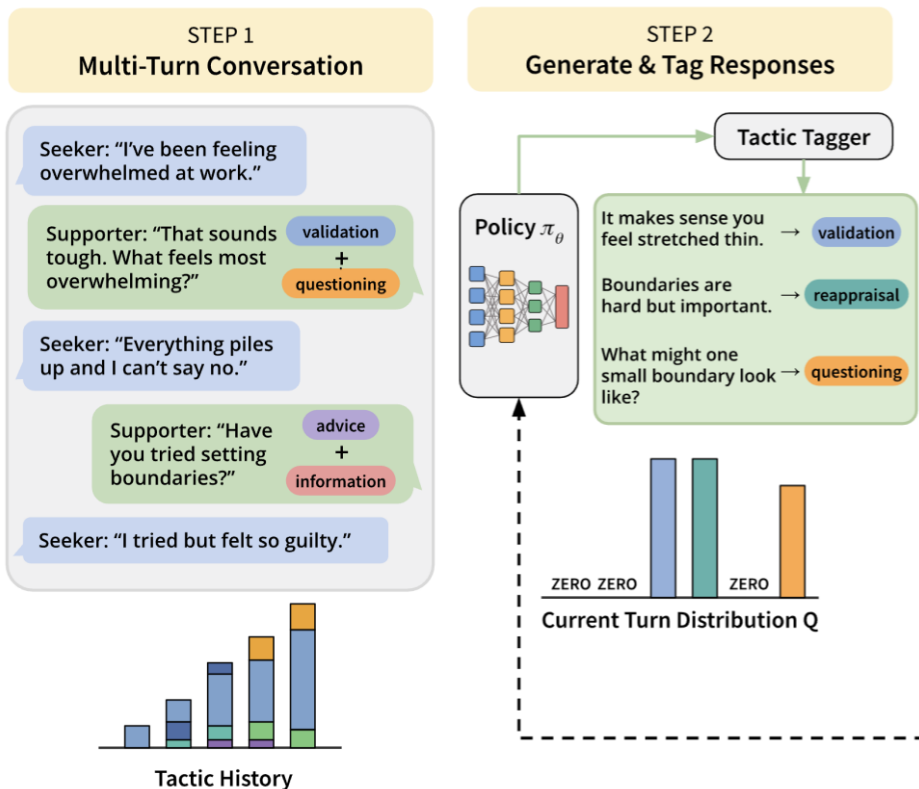
Seeker: "I tried but felt so guilty."



Tactic History

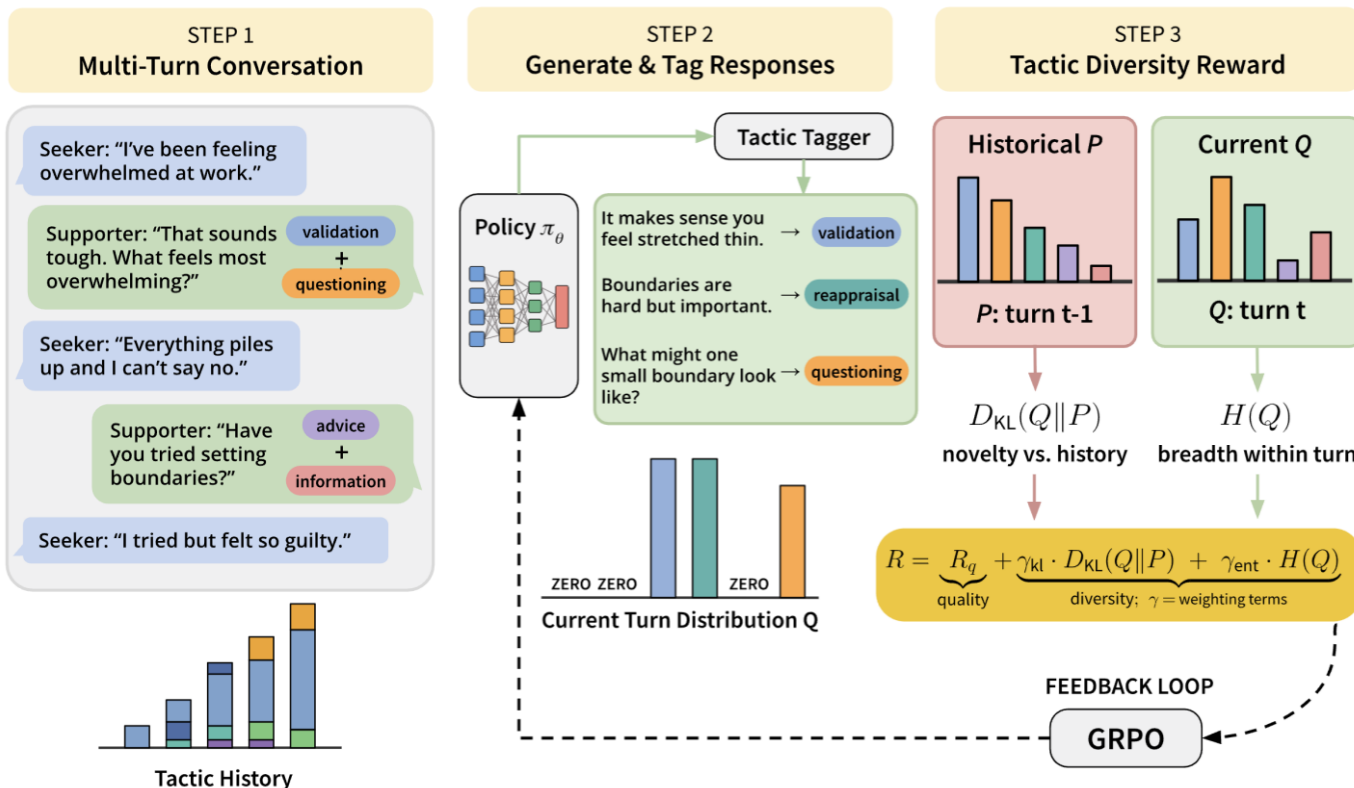


# MINT : Multi-turn Inter-tactic Novelty Training



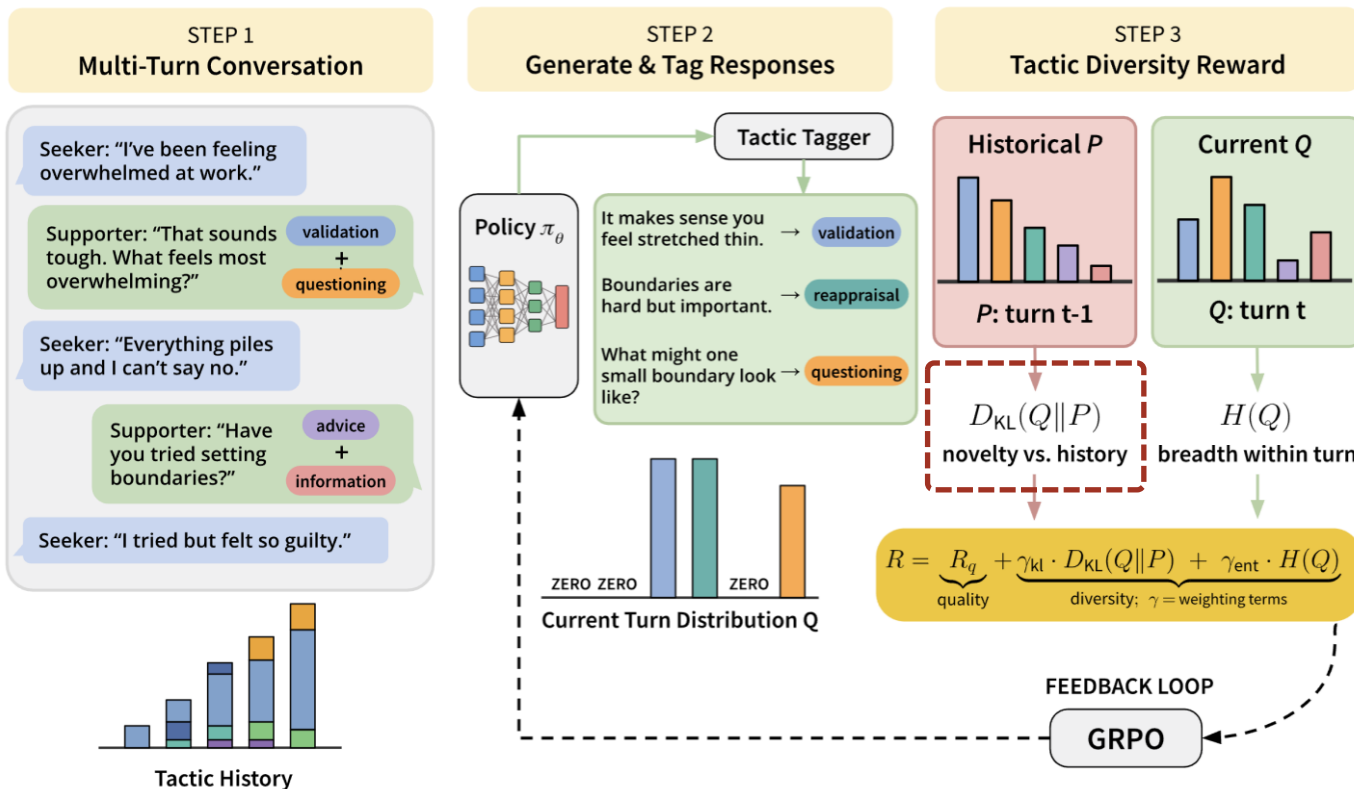


# MINT: Multi-turn Inter-tactic Novelty Training



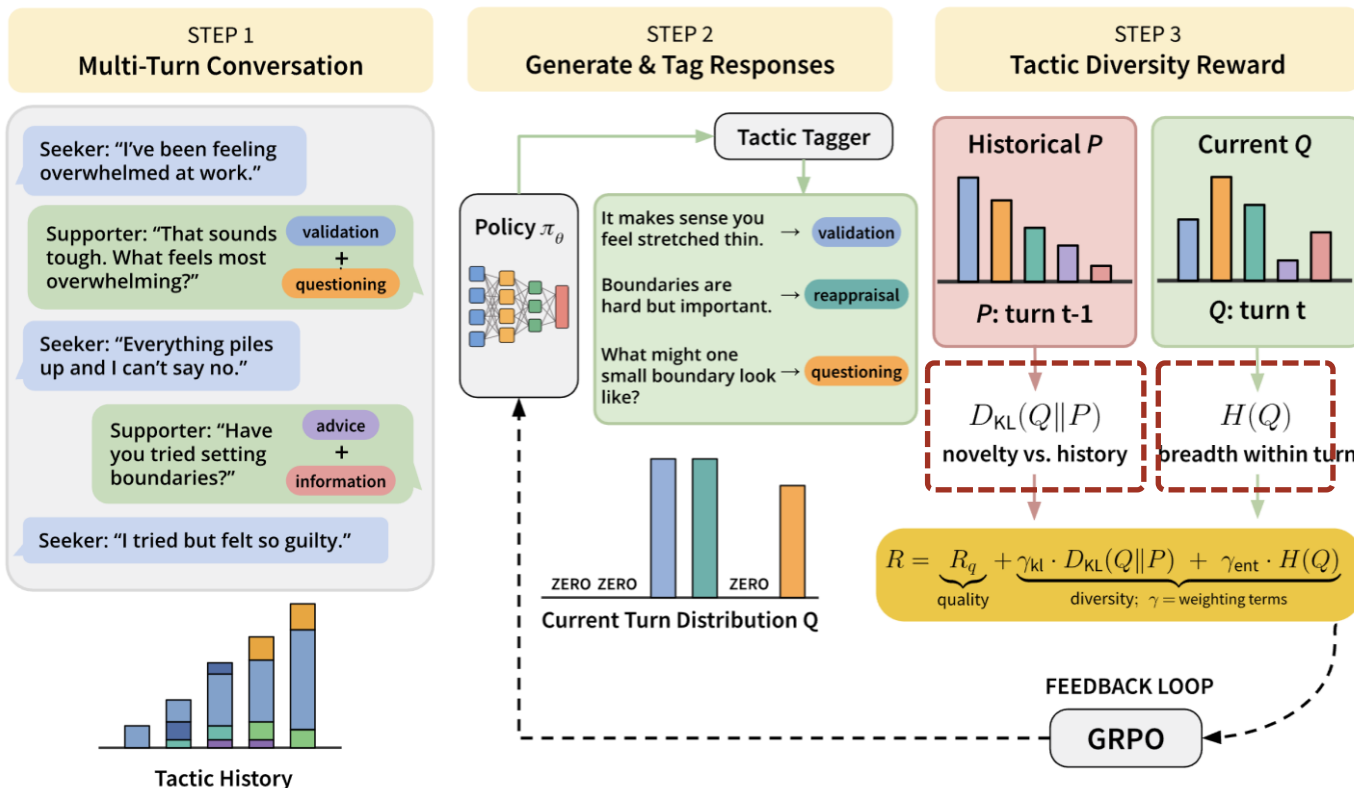


# MINT: Multi-turn Inter-tactic Novelty Training



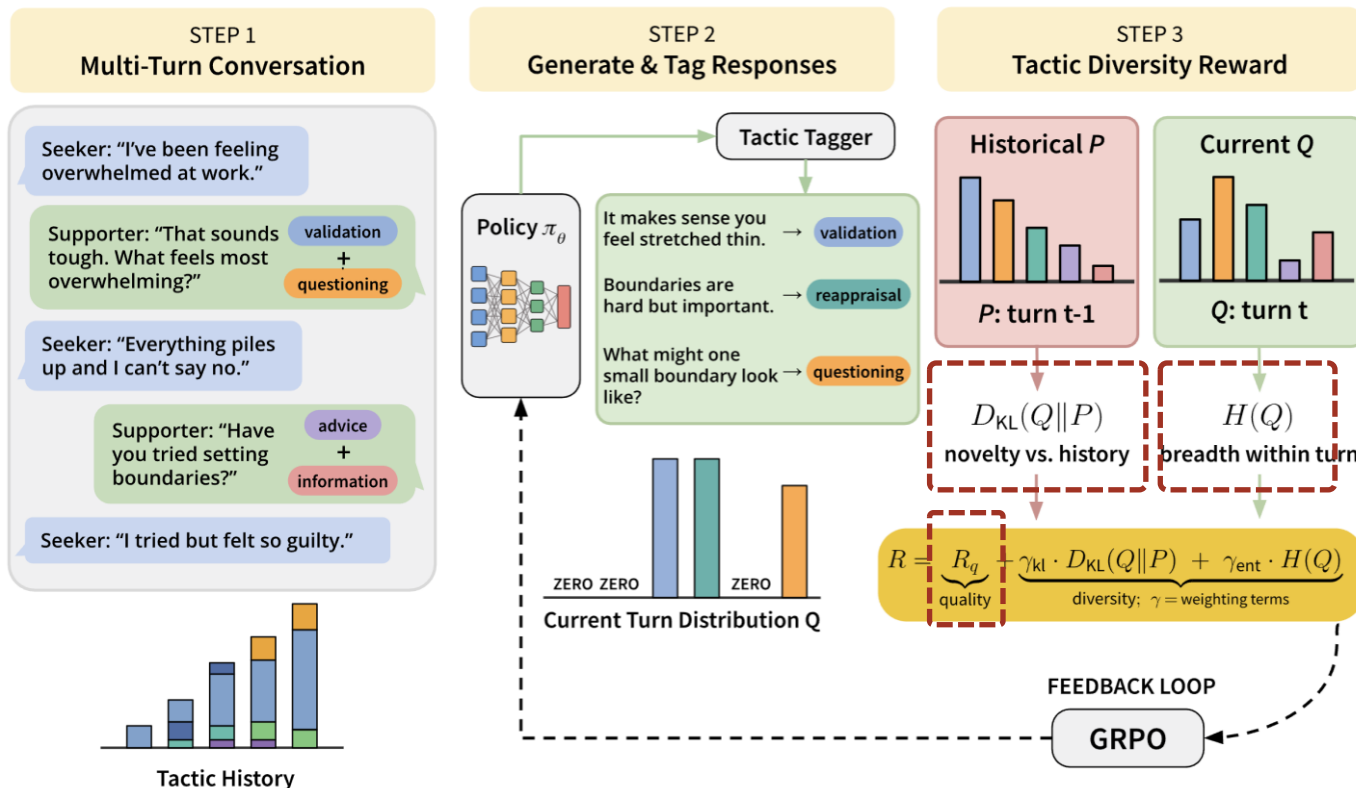


# MINT: Multi-turn Inter-tactic Novelty Training



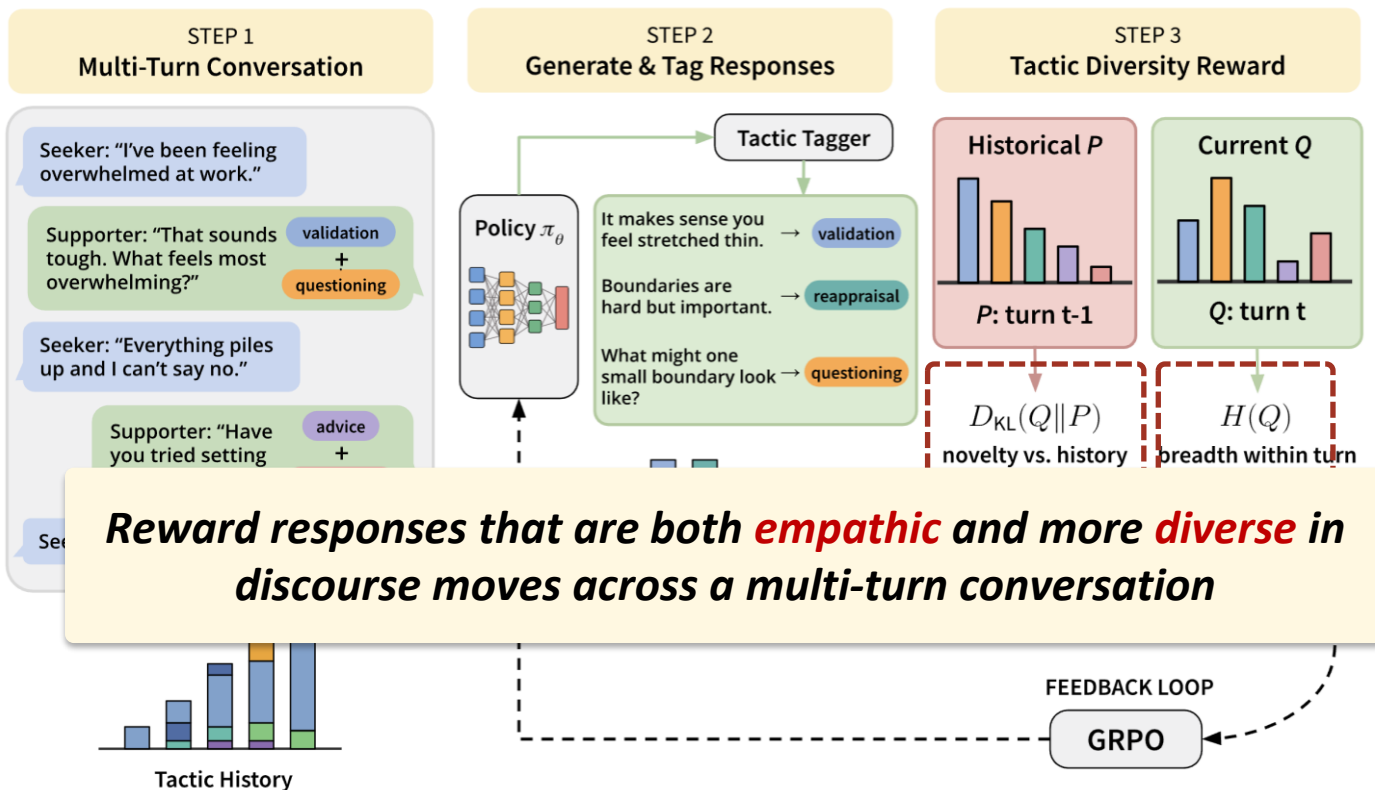


# MINT: Multi-turn Inter-tactic Novelty Training





# MINT : Multi-turn Inter-tactic Novelty Training





# MINT : Reward Details

**Q** Empathy quality (*PsychoCounsel* reward model)

**D<sub>KL</sub>** Cross-turn tactic novelty

**H** Within-turn tactic breadth

All reward components min-max normalized to [0, 1] within each rollout group

$$D_{\text{KL}}(Q_t \| P_{t-1}) = \min \left( \sum_{k=1}^K Q_t(k) \log \frac{Q_t(k)}{P_{t-1}(k)}, \tau \right)$$

$$Q_t(k) = \frac{c_{t,k} + \alpha}{\sum_{j=1}^K c_{t,j} + K\alpha}, \quad P_{t-1}(k) = \frac{c_{t-1,k} + \alpha}{\sum_{j=1}^K c_{t-1,j} + K\alpha}$$

$$H(Q_t) = - \sum_{k=1}^K Q_t(k) \log Q_t(k)$$

# MINT : Reward Details

**Q** Empathy quality (*PsychoCounsel* reward model)

**D<sub>KL</sub>** Cross-turn tactic novelty

**H** Within-turn tactic breadth

All reward components min-max normalized to [0, 1] within each rollout group

$$D_{\text{KL}}(Q_t \| P_{t-1}) = \min \left( \sum_{k=1}^K Q_t(k) \log \frac{Q_t(k)}{P_{t-1}(k)}, \tau \right)$$

$$H(Q_t) = - \sum_{k=1}^K Q_t(k) \log Q_t(k)$$

## Current & Prior Turn Tactic Distributions

$$Q_t(k) = \frac{c_{t,k} + \alpha}{\sum_{j=1}^K c_{t,j} + K\alpha}, \quad P_{t-1}(k) = \frac{c_{t-1,k} + \alpha}{\sum_{j=1}^K c_{t-1,j} + K\alpha}$$

$c_{t,k}$	# sentences tagged with tactic k at turn t
$K = 10$	Number of empathy tactics
$\alpha = 0.1$	Laplace smoothing constant
$\tau = 5$	KL clipping threshold

# MINT : Reward Details

**Q** Empathy quality (*PsychoCounsel* reward model)

**D<sub>KL</sub>** Cross-turn tactic novelty

**H** Within-turn tactic breadth

All reward components min-max normalized to [0, 1] within each rollout group

## 1) How different is this turn from previous turns?

Rewards the model for shifting its tactic mix as the conversation progresses.

*Discourages repeating the same tactic patterns across turns (reduces stickiness).*

$$D_{\text{KL}}(Q_t \| P_{t-1}) = \min \left( \sum_{k=1}^K Q_t(k) \log \frac{Q_t(k)}{P_{t-1}(k)}, \tau \right)$$

$$H(Q_t) = - \sum_{k=1}^K Q_t(k) \log Q_t(k)$$

## Current & Prior Turn Tactic Distributions

$$Q_t(k) = \frac{c_{t,k} + \alpha}{\sum_{j=1}^K c_{t,j} + K\alpha}, \quad P_{t-1}(k) = \frac{c_{t-1,k} + \alpha}{\sum_{j=1}^K c_{t-1,j} + K\alpha}$$

$c_{t,k}$	# sentences tagged with tactic k at turn t
$K = 10$	Number of empathy tactics
$\alpha = 0.1$	Laplace smoothing constant
$\tau = 5$	KL clipping threshold

# MINT : Reward Details

**Q** Empathy quality (*PsychoCounsel* reward model)

**D<sub>KL</sub>** Cross-turn tactic novelty

**H** Within-turn tactic breadth

All reward components min-max normalized to [0, 1] within each rollout group

## 1) How different is this turn from previous turns?

Rewards the model for shifting its tactic mix as the conversation progresses.

*Discourages repeating the same tactic patterns across turns (reduces stickiness).*

$$D_{\text{KL}}(Q_t \| P_{t-1}) = \min \left( \sum_{k=1}^K Q_t(k) \log \frac{Q_t(k)}{P_{t-1}(k)}, \tau \right)$$

## 2) How varied are tactics within this turn?

KL alone doesn't reward using multiple tactics in one response. Entropy complements it.

*Higher when tactics are evenly spread; near zero when one tactic dominates.*

$$H(Q_t) = - \sum_{k=1}^K Q_t(k) \log Q_t(k)$$

## Current & Prior Turn Tactic Distributions

$$Q_t(k) = \frac{c_{t,k} + \alpha}{\sum_{j=1}^K c_{t,j} + K\alpha}, \quad P_{t-1}(k) = \frac{c_{t-1,k} + \alpha}{\sum_{j=1}^K c_{t-1,j} + K\alpha}$$

$c_{t,k}$	# sentences tagged with tactic k at turn t
$K = 10$	Number of empathy tactics
$\alpha = 0.1$	Laplace smoothing constant
$\tau = 5$	KL clipping threshold



# Evaluation: Why We Need Both *Empathy* AND *Diversity*

If we only measure **empathy**...

The model sounds great each turn, but says the same thing every turn.

# Evaluation: Why We Need Both *Empathy* AND *Diversity*

If we only measure **empathy**...

The model sounds great each turn, but says the same thing every turn.

If we only measure **diversity**...

The model is diverse but unhelpful.

# Evaluation: Why We Need Both *Empathy* AND *Diversity*

If we only measure **empathy**...

The model sounds great each turn, but says the same thing every turn.

If we only measure **diversity**...

The model is diverse but unhelpful.

We need *both*  
dimensions

# Evaluation: Why We Need Both *Empathy* AND *Diversity*

If we only measure **empathy**...

The model sounds great each turn, but says the same thing every turn.

If we only measure **diversity**...

The model is diverse but unhelpful.

We need *both*  
dimensions

**Aggregate Empathy** via Lend-an-Ear (Kumar et al., Nature 2026)

# Evaluation: Why We Need Both *Empathy* AND *Diversity*

If we only measure **empathy**...

The model sounds great each turn, but says the same thing every turn.

If we only measure **diversity**...

The model is diverse but unhelpful.

We need *both*  
dimensions

**Aggregate Empathy** via Lend-an-Ear (Kumar et al., Nature 2026)

6 empathy dimensions (1-5 scale) per turn (3 desirable + 3 undesirable)

# Evaluation: Why We Need Both *Empathy* AND *Diversity*

If we only measure **empathy**...

The model sounds great each turn, but says the same thing every turn.

If we only measure **diversity**...

The model is diverse but unhelpful.

We need *both*  
dimensions

**Aggregate Empathy** via Lend-an-Ear (Kumar et al., Nature 2026)

6 empathy dimensions (1-5 scale) per turn (3 desirable + 3 undesirable)

# Evaluation: Why We Need Both *Empathy* AND *Diversity*

If we only measure **empathy**...  
The model sounds great each turn, but  
says the same thing every turn.

We need *both*  
dimensions

If we only measure **diversity**...  
The model is diverse but unhelpful.

**Aggregate Empathy** via Lend-an-Ear (Kumar et al., Nature 2026)

6 empathy dimensions (1-5 scale) per turn (3 desirable + 3 undesirable)

**LLM Judge Reliability:** `gpt-oss-120b` as turn-level empathy judge (temp = 0)

# Evaluation: Why We Need Both *Empathy* AND *Diversity*

If we only measure **empathy**...

The model sounds great each turn, but says the same thing every turn.

If we only measure **diversity**...

The model is diverse but unhelpful.

We need *both*  
dimensions

**Aggregate Empathy** via Lend-an-Ear (Kumar et al., Nature 2026)

6 empathy dimensions (1-5 scale) per turn (3 desirable + 3 undesirable)

**LLM Judge Reliability:** **gpt-oss-120b** as turn-level empathy judge (temp = 0)

Validated on 315 expert-annotated supporter turns from the Lend-an-Ear benchmark.

# Evaluation: Why We Need Both *Empathy* AND *Diversity*

If we only measure **empathy**...  
The model sounds great each turn, but  
says the same thing every turn.

We need *both*  
dimensions

If we only measure **diversity**...  
The model is diverse but unhelpful.

**Aggregate Empathy** via Lend-an-Ear (Kumar et al., Nature 2026)

6 empathy dimensions (1-5 scale) per turn (3 desirable + 3 undesirable)

**LLM Judge Reliability:** **gpt-oss-120b** as turn-level empathy judge (temp = 0)

Validated on 315 expert-annotated supporter turns from the Lend-an-Ear benchmark.

Weighted Cohen's  $\kappa_w = \mathbf{0.58}$  — exceeds the pairwise agreement among human experts

# Evaluation: Why We Need Both *Empathy* AND *Diversity*

If we only measure **empathy**...  
The model sounds great each turn, but says the same thing every turn.

We need *both*  
dimensions

If we only measure **diversity**...  
The model is diverse but unhelpful.

**Aggregate Empathy** via Lend-an-Ear (Kumar et al., Nature 2026)

6 empathy dimensions (1-5 scale) per turn (3 desirable + 3 undesirable)

**Tactic Stickiness** =  $P(\text{tactic at turn } t \mid \text{same tactic at turn } t-1)$

**LLM Judge Reliability:** **gpt-oss-120b** as turn-level empathy judge (temp = 0)

Validated on 315 expert-annotated supporter turns from the Lend-an-Ear benchmark.

Weighted Cohen's  $\kappa_w = \mathbf{0.58}$  — exceeds the pairwise agreement among human experts

# Evaluation: Why We Need Both *Empathy* AND *Diversity*

If we only measure **empathy**...

The model sounds great each turn, but says the same thing every turn.

If we only measure **diversity**...

The model is diverse but unhelpful.

We need *both*  
dimensions

**Aggregate Empathy** via Lend-an-Ear (Kumar et al., Nature 2026)

6 empathy dimensions (1-5 scale) per turn (3 desirable + 3 undesirable)

**Tactic Stickiness** =  $P(\text{tactic at turn } t \mid \text{same tactic at turn } t-1)$

How likely is a tactic from the previous turn to appear again?  
Lower = more diverse

**LLM Judge Reliability:** **gpt-oss-120b** as turn-level empathy judge (temp = 0)

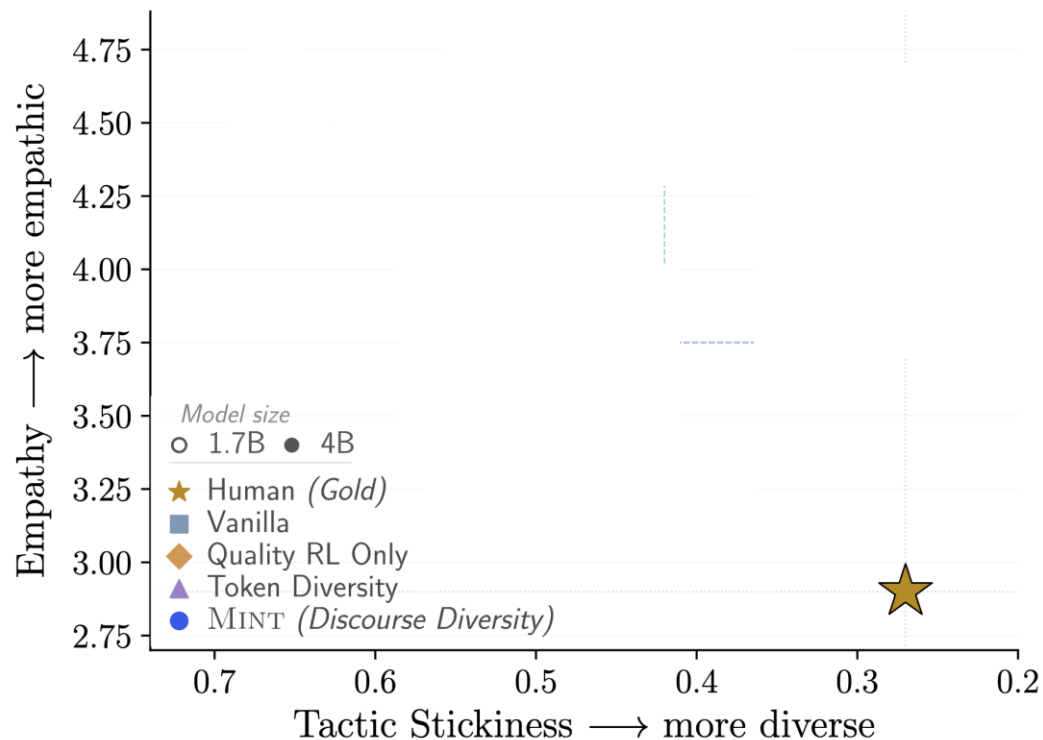
Validated on 315 expert-annotated supporter turns from the Lend-an-Ear benchmark.

Weighted Cohen's  $\kappa_w = \mathbf{0.58}$  — exceeds the pairwise agreement among human experts

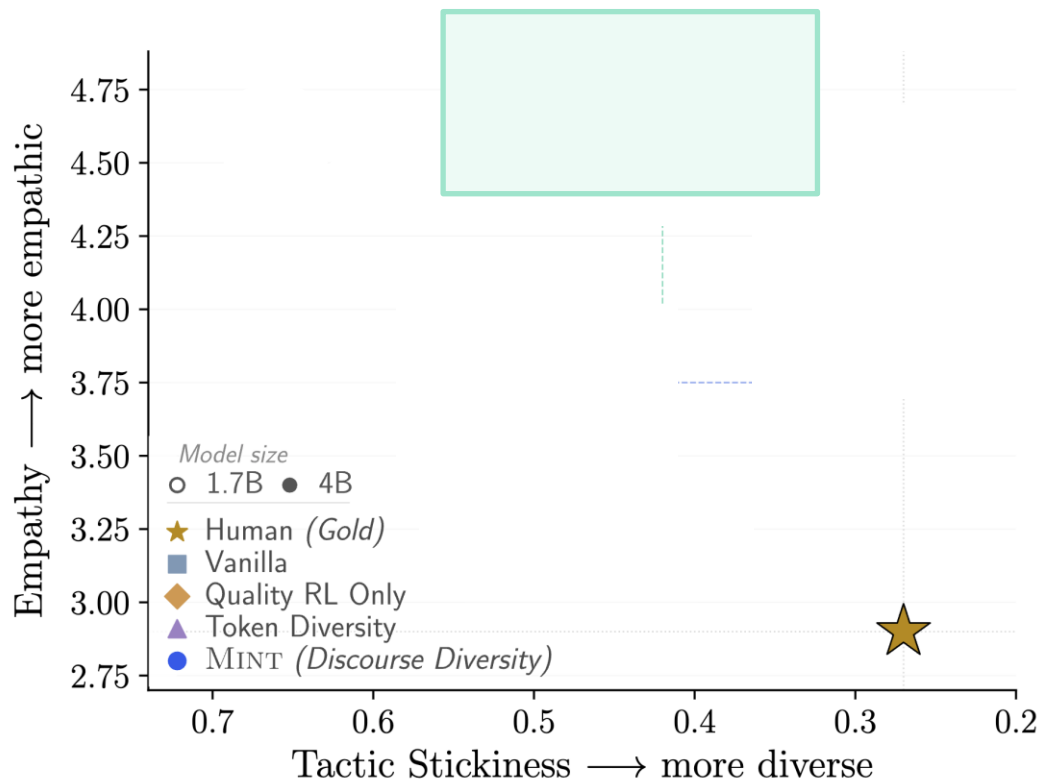


# Optimizing for Both *Empathy* and *Diversity*

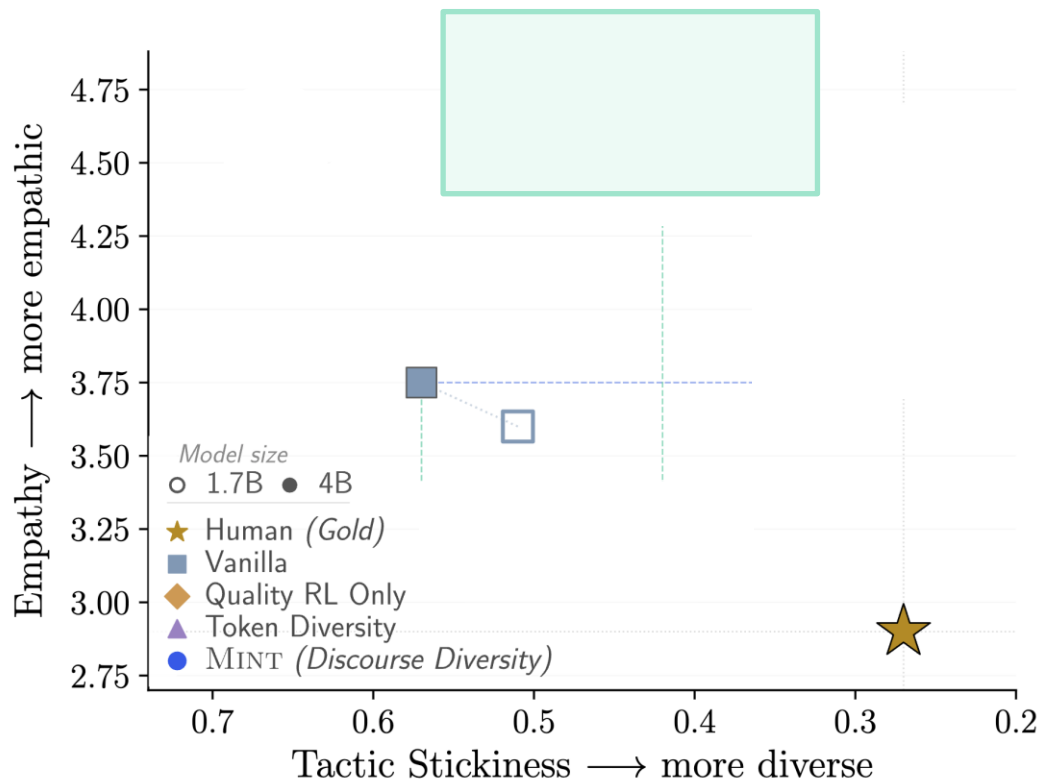
# Optimizing for Both *Empathy* and *Diversity*



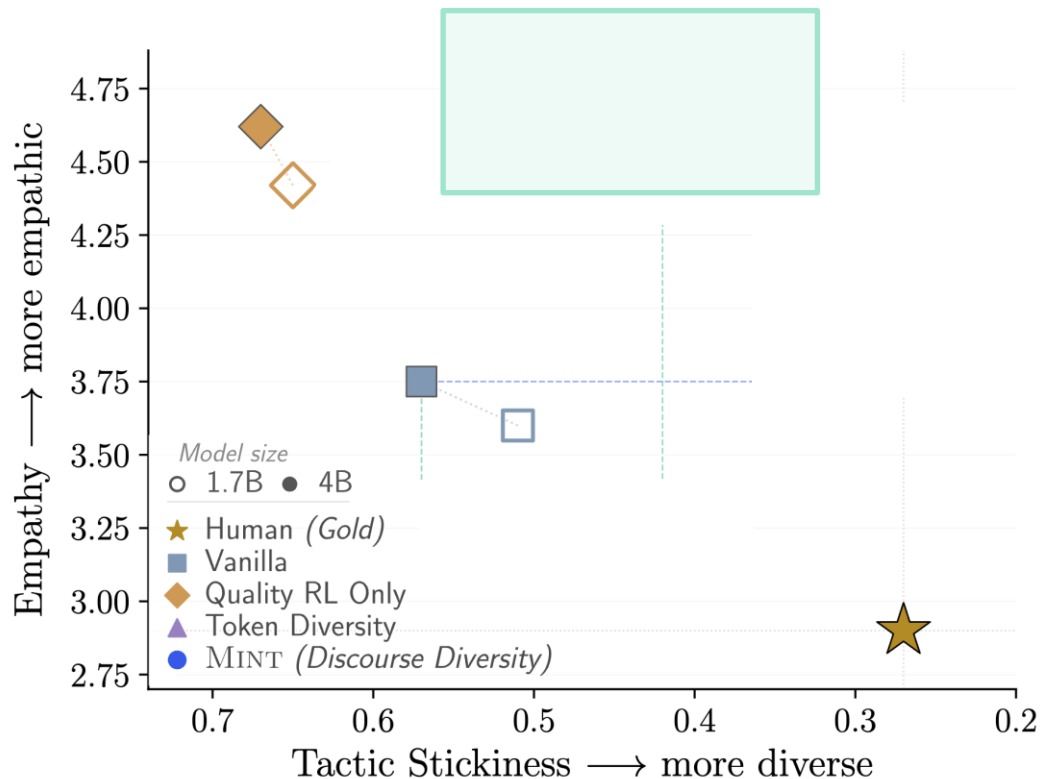
# Optimizing for Both *Empathy* and *Diversity*



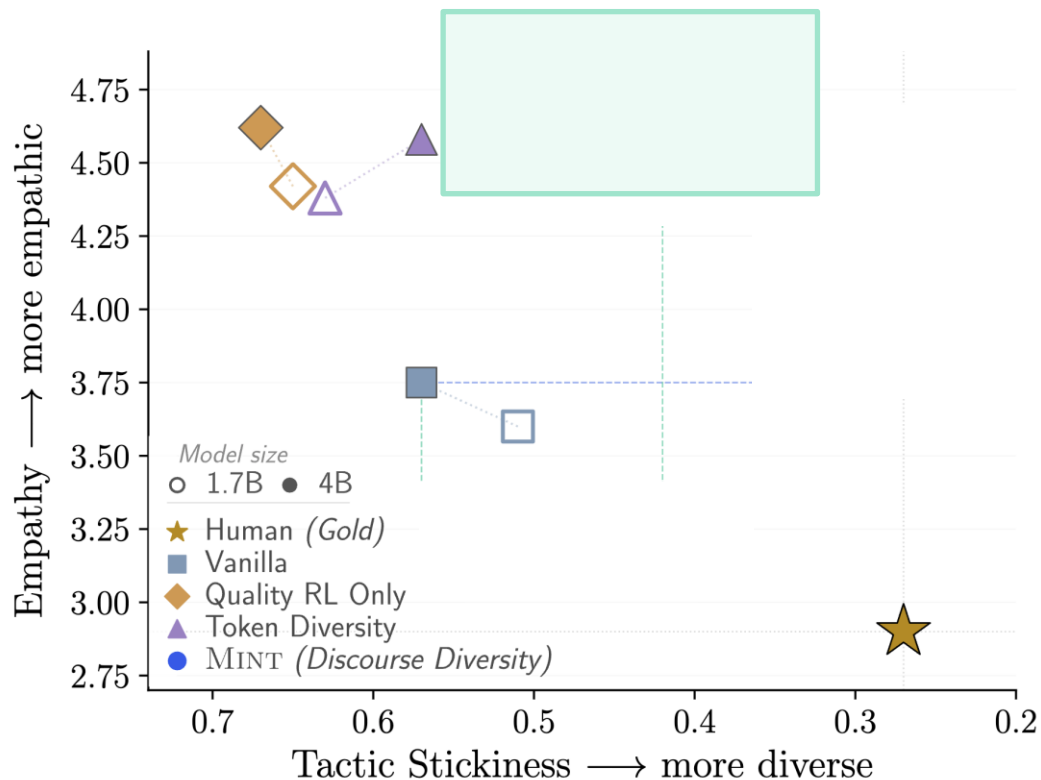
# Optimizing for Both *Empathy* and *Diversity*



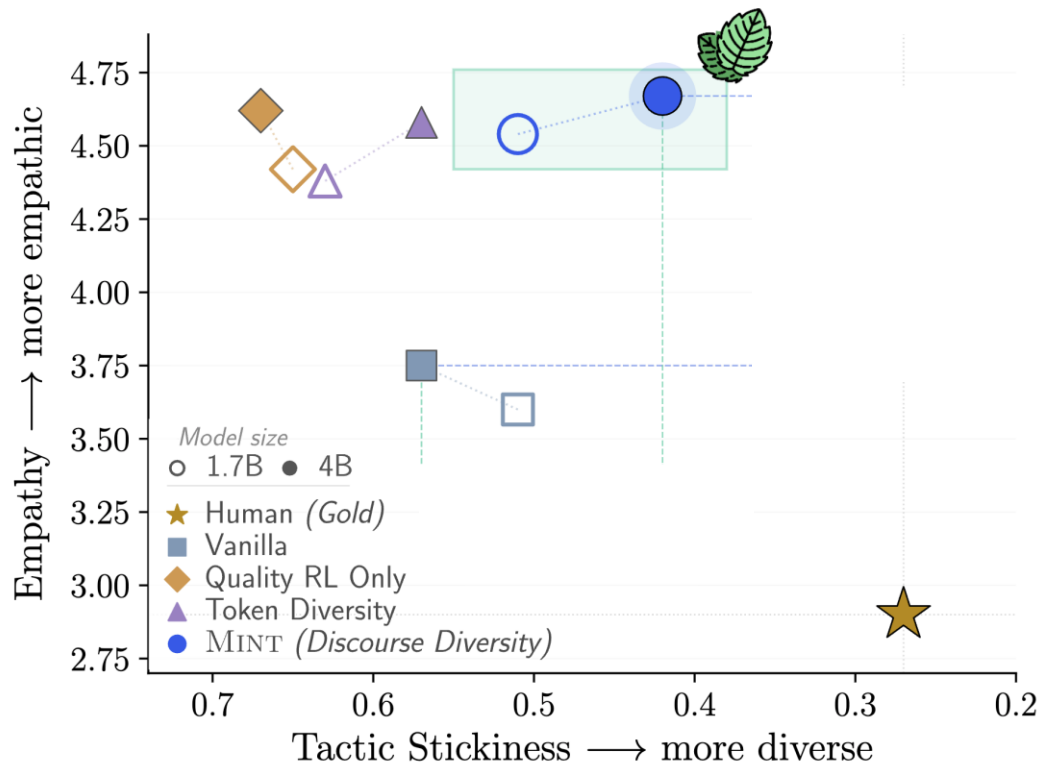
# Optimizing for Both *Empathy* and *Diversity*



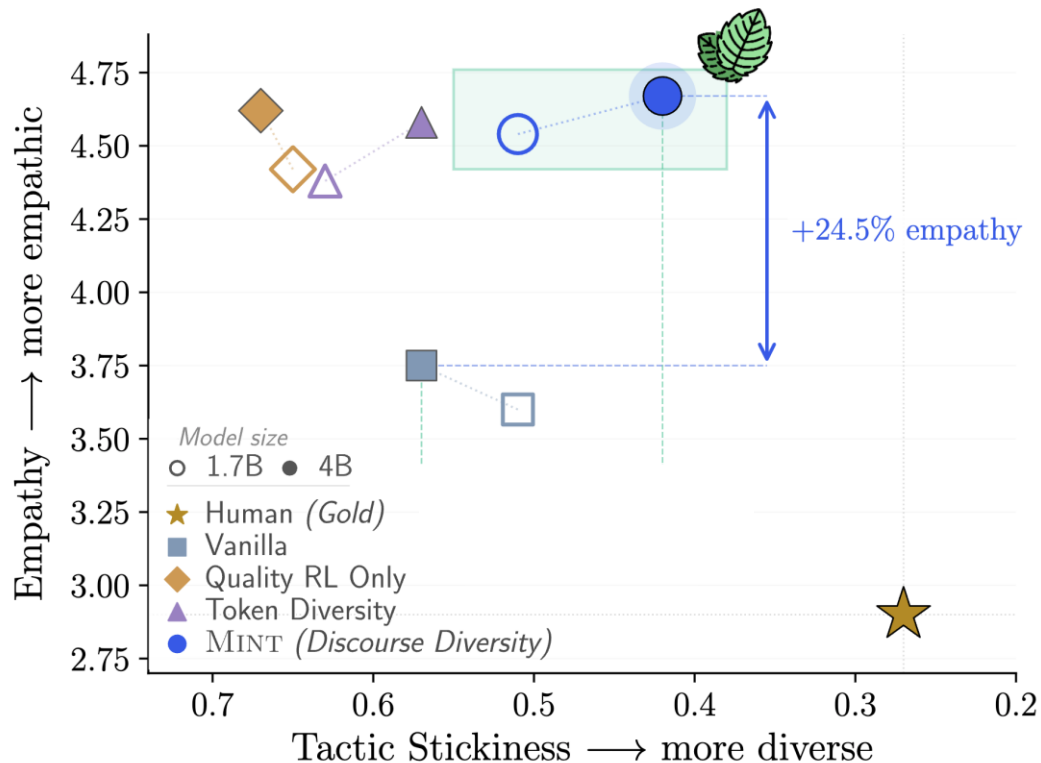
# Optimizing for Both *Empathy* and *Diversity*



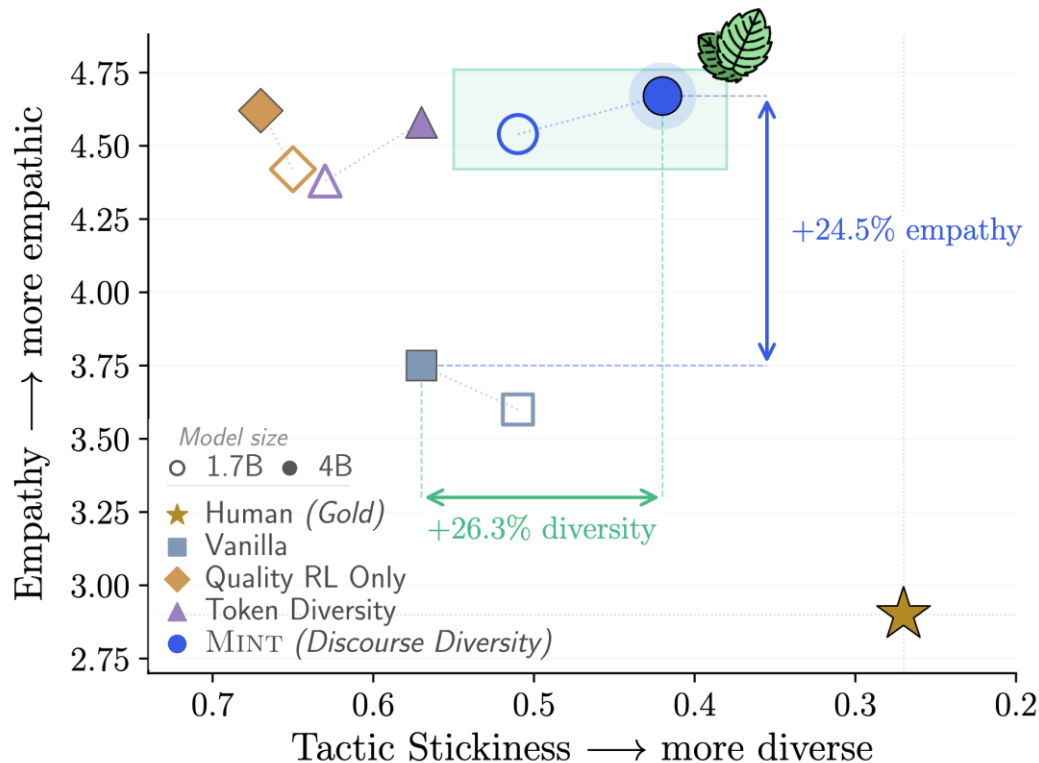
# Optimizing for Both *Empathy* and *Diversity*



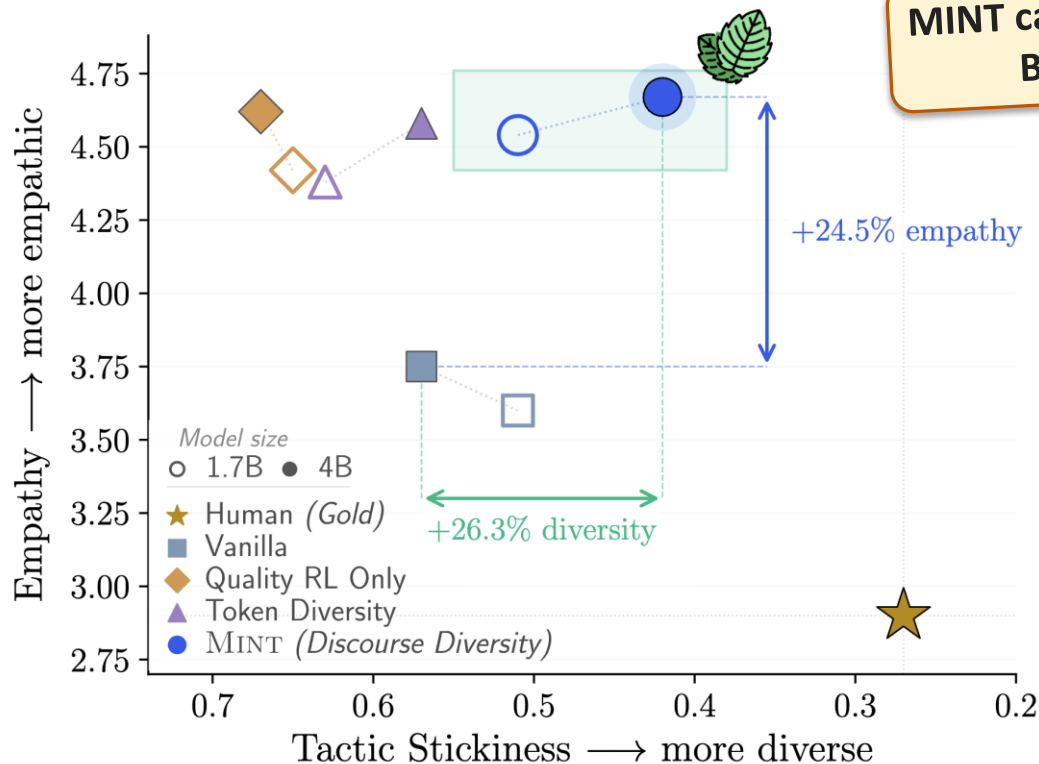
# Optimizing for Both *Empathy* and *Diversity*



# Optimizing for Both *Empathy* and *Diversity*

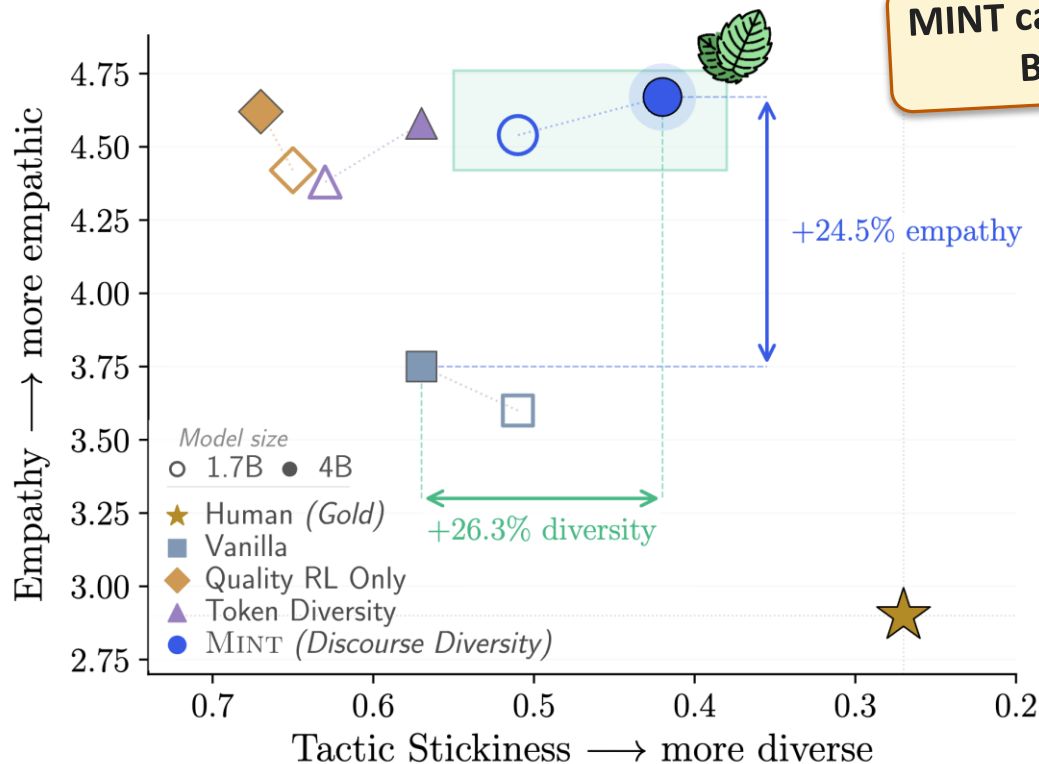


# Optimizing for Both *Empathy* and *Diversity*





# Optimizing for Both *Empathy* and *Diversity*



**MINT can optimize BOTH!**

Less unsolicited advice, more elaboration, tactic profiles closer to humans

# **Part 3 takeaway:**

## **LLMs are repetitive in diversity moves in multi-turn conversations!**

# **Part 3 takeaway:**

## **LLMs are repetitive in diversity moves in multi-turn conversations!**

**Optimizing for both Quality and Diversity can help achieve a balance between both!**

## Deciphering Emotions from Text

- EMNLP 2022
- EMNLP 2023 Findings

1

## Discourse Diversity in Multi-Turn Empathic Dialogue

- *Under review*

3

## Unveiling Advanced Psychological Capabilities from LLMs: *A Case of Targeted Reappraisal*

- COLM 2024
- ICML 2025

2

4

## Conclusion

- Summary of Contributions



# Thesis Summary & Contributions



# Thesis Summary & Contributions

**Understand:** capable of identifying cognitive appraisals, on par with laypeople



# Thesis Summary & Contributions

**Understand:** capable of identifying cognitive appraisals, on par with laypeople  
(EMNLP 2022, EMNLP 2023 Findings)



# Thesis Summary & Contributions

**Understand:** capable of identifying cognitive appraisals, on par with laypeople  
(EMNLP 2022, EMNLP 2023 Findings)

**Intervene:** Principle-guided reappraisals outperform human-written ones



# Thesis Summary & Contributions

**Understand:** capable of identifying cognitive appraisals, on par with laypeople  
(EMNLP 2022, EMNLP 2023 Findings)

**Intervene:** Principle-guided reappraisals outperform human-written ones  
(COLM 2024, ICML 2025)



# Thesis Summary & Contributions

**Understand:** capable of identifying cognitive appraisals, on par with laypeople  
(EMNLP 2022, EMNLP 2023 Findings)

**Intervene:** Principle-guided reappraisals outperform human-written ones  
(COLM 2024, ICML 2025)

**Adapt:** MINT closes the discourse-move-diversity gap while improving empathy



# Thesis Summary & Contributions

**Understand:** capable of identifying cognitive appraisals, on par with laypeople  
(EMNLP 2022, EMNLP 2023 Findings)

**Intervene:** Principle-guided reappraisals outperform human-written ones  
(COLM 2024, ICML 2025)

**Adapt:** MINT closes the discourse-move-diversity gap while improving empathy  
(under submission)



# *Selected Ph.D. Highlights*

**8****refereed  
conference papers****6****first-author /  
co-first-author**ICML x1 COLM x1  
EMNLP x2 ACL x1**3****research scientist  
internships**IBM Research x2  
MBZUAI IFM x1**1****first-authored  
U.S. patent**

filed and pending

During the Ph.D., this work led to **peer-reviewed publications**, **industrial research internships**, and **technology transfer**.



# Many thanks to my mentors & collaborators!

