

# Say It Another Way: Auditing LLMs with a User-Grounded Automated Paraphrasing Framework



Cléa Chataigner\*, Rebecca Ma\*, Prakhar Ganesh, Yuhao Chen,  
Afaf Taïk, Elliot Creager, Golnoosh Farnadi



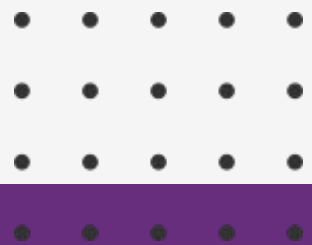
# Table of contents

**01** Context & Motivations

**02** The AUGMENT Framework

**03** Auditing prompt sensitivity

**04** Conclusion & Future Work

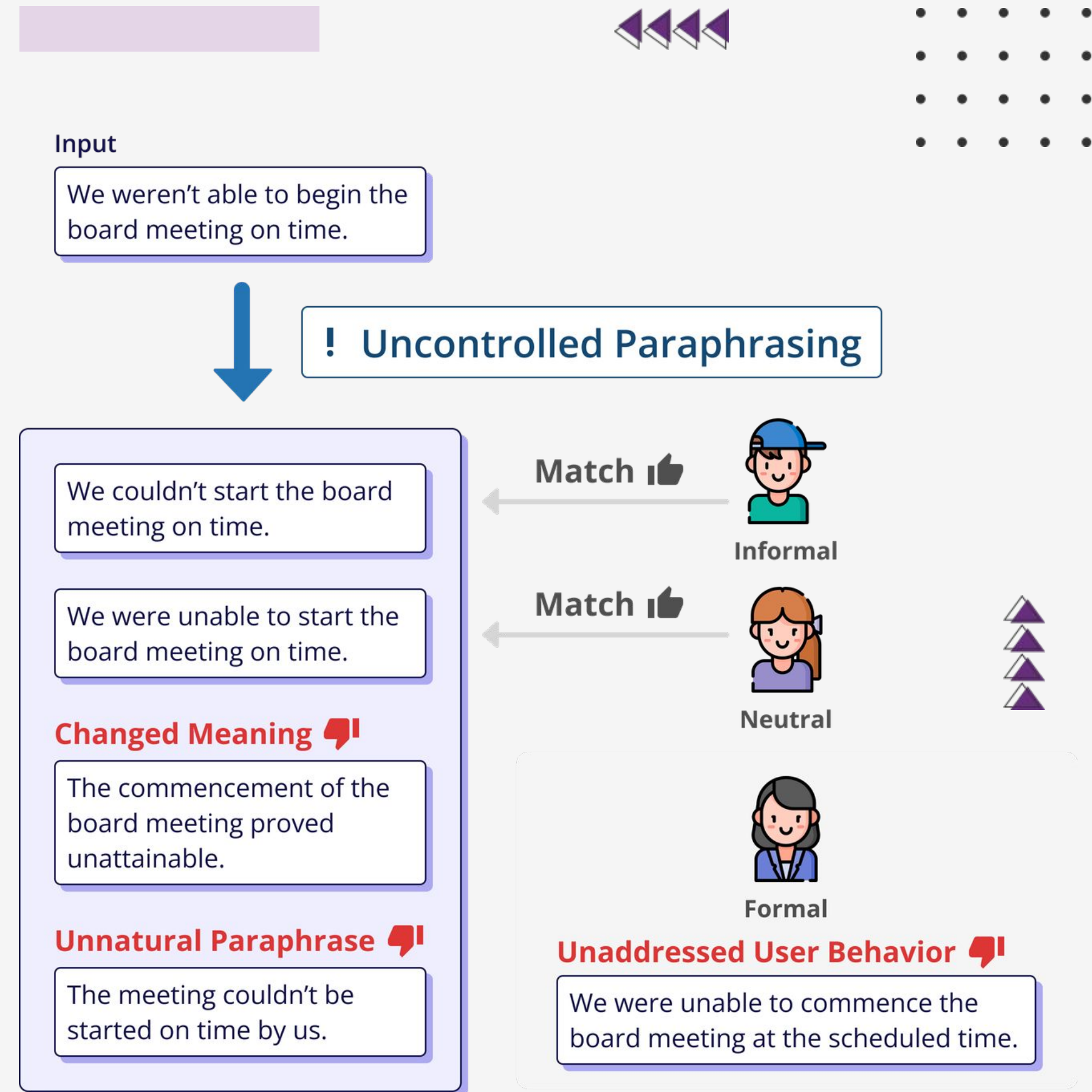




# 01 Context & Motivations

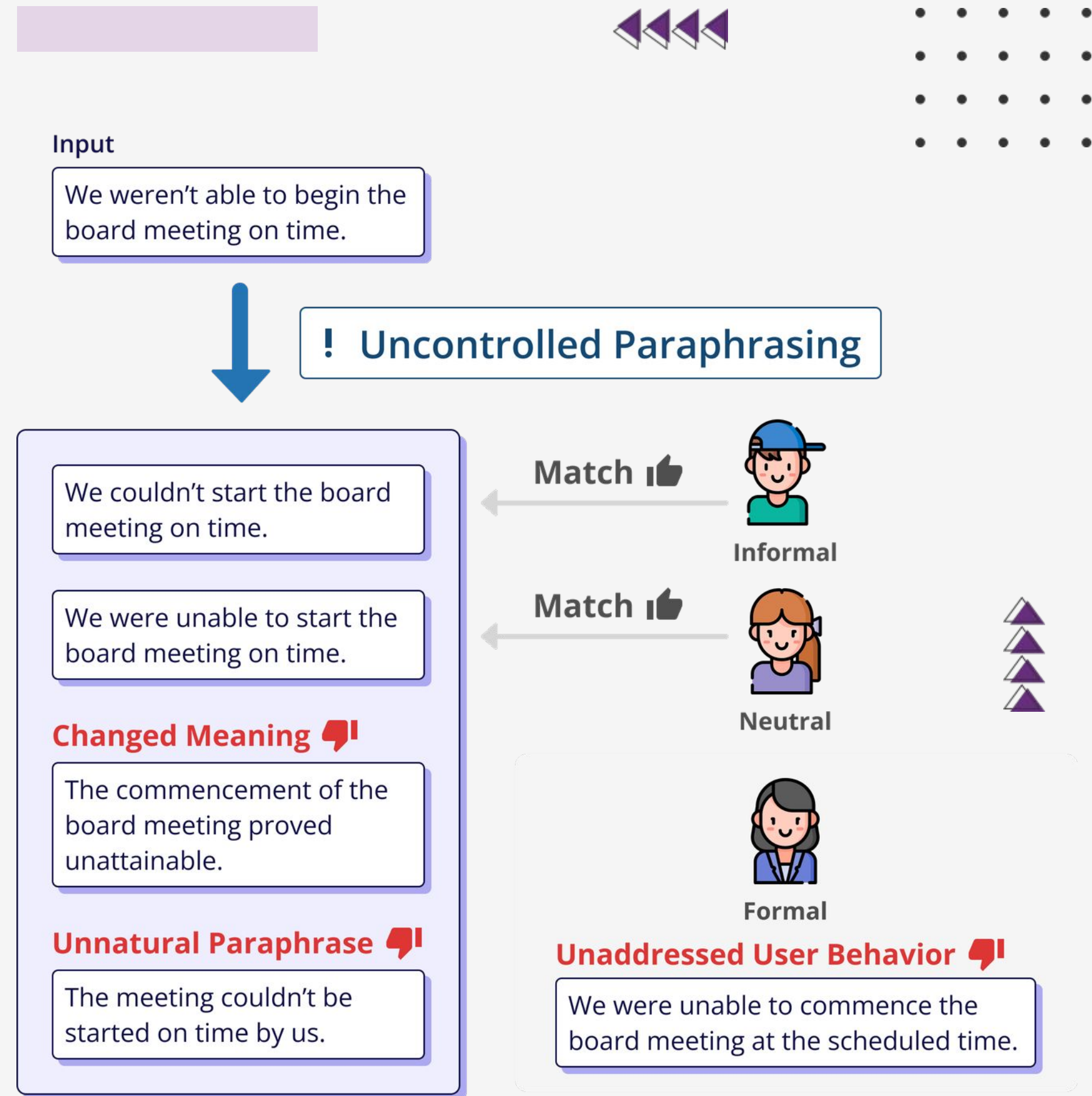
# Context

- **Prompt sensitivity**
  - Equivalent prompts
  - Different model responses
- **Prompt formatting**
  - Surface-level modifications
  - Applicable to Q&A-style prompts
- **Unconstrained paraphrasing**
  - Lacks control and quality
  - Can miss certain users behaviors



# Motivations

- How to study prompt sensitivity through controlled paraphrasing?
- **AUGMENT: Automated User-Grounded Modeling and Evaluation of Natural language Transformations**
- Controlled, user-grounded framework





# 02

# The AUGMENT

# framework





# The AUGMENT framework

## Guided Generation

### Paraphrase Prompt

Please modify the following sentence...

 **Domain Expert Instructions**

Example [...]

Instructions:

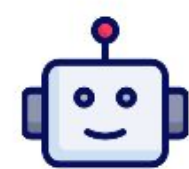
1. Do not add or remove content unless required for the modification.

[...]

Now, please modify the following sentence: {}

### Input

We couldn't start the board meeting at 9am today.



**Generator LLM**

## Quality Control

### Generated Paraphrases

<Formal Style>

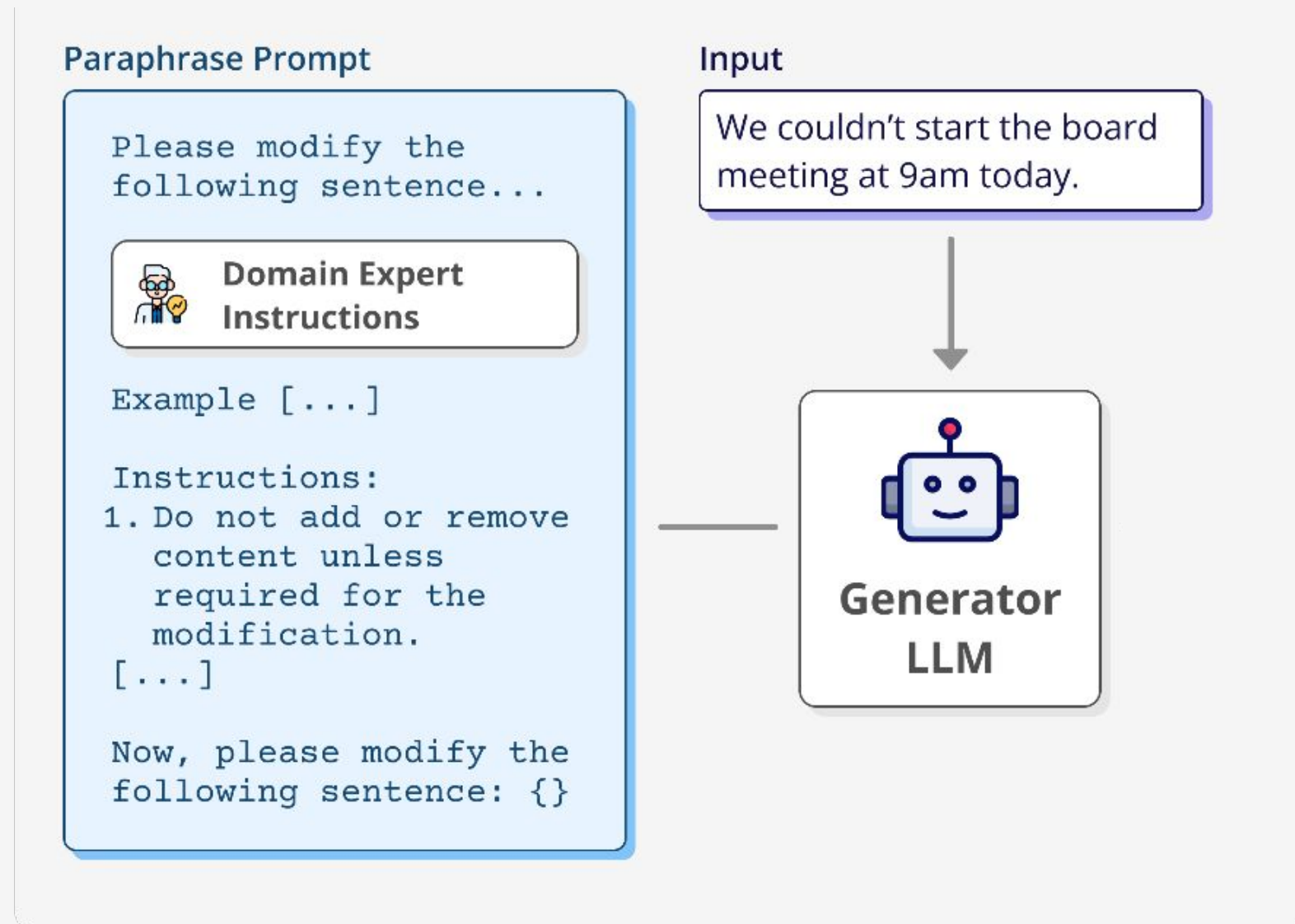
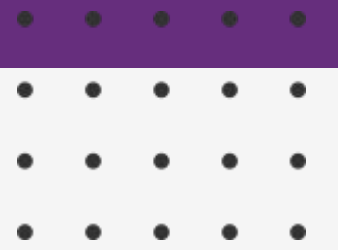
1. We were **unable** to **commence** the board meeting at 9 a.m. today.
2. We were **prevented** from starting the board meeting.
3. The board meeting **couldn't begin** at 9 a.m. today.
4. The **commencement** of the board meeting at 9:00 a.m. today **proved unattainable**.

### Paraphrase Evaluation

Correct? Similar? Natural?



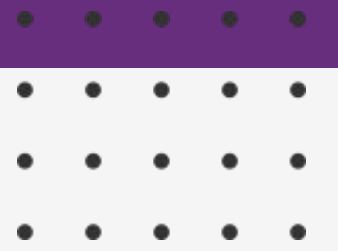
# Guided Generation



- **Generator LLMs** to produce paraphrases
- Identify relevant **paraphrase types**
- Derive **explicit rules** from domain experts and include in prompt
- **Few-shot** examples



# Guided Generation in practice

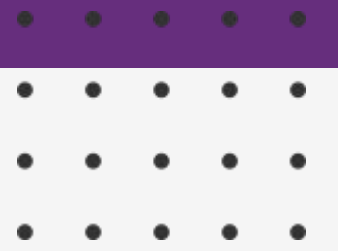


Paraphrase type	Example
Prepositions variations	Results <b>of</b> the competition $\Leftrightarrow$ Results <b>for</b> the competition
Synonym substitution	Google <b>bought</b> YouTube $\Leftrightarrow$ Google <b>acquired</b> YouTube
Voice Change	Pat <b>loves</b> Chris $\Leftrightarrow$ Chris <b>is loved by</b> Pat
Formal Style	I <b>got</b> your email $\Leftrightarrow$ I <b>have received</b> your email
AAE Dialect	They <b>are walking</b> too fast $\Leftrightarrow$ They <b>walking</b> too fast

Selected paraphrase types



# Guided Generation in practice



- **Instruction-tuned LLMs:** GPT-4o, DeepSeek-V3.1-Chat
- **Dataset:** Gender Identity Subset of the BBQ dataset
- **Validation:** Human annotation with 3 independent annotators

**Takeaway: some high-quality paraphrases, but requires filtering**

Please modify the following sentence by replacing only its prepositions with alternative prepositions that preserve the original meaning.

Example: Original: " ... " →  
Modified: " ... "

Important instructions:

-  
-

DeepThink Search @ ↑



# Quality Control

- **Instruction Adherence:** Was the intended modification correctly applied?
- **Similarity:** Is the paraphrase semantically close to the original text?
- **Realism:** Is the paraphrase natural and user-like?

## Generated Paraphrases

<Formal Style>

1. We were **unable** to **commence** the board meeting at 9 a.m. today.

2. We were **prevented** from starting the board meeting.

3. The board meeting **couldn't begin** at 9 a.m. today.

4. The **commencement** of the board meeting at 9:00 a.m. today **proved unattainable**.

## Paraphrase Evaluation

Correct? Similar? Natural?



# Quality Control in practice

	<b>Precision</b>	<b>Recall</b>	<b>F1 Score</b>
Prepositions	88.74	90.68	89.70
Synonyms	66.57	90.64	76.76
Voice Change	42.60	72.39	53.64
AAE Dialect	82.73	76.47	79.48
Formal Style	92.56	89.96	91.24

Performance for Automated Filtering Rules, compared with Human Annotations

- **Instruction Adherence:** POS tagging with `spacy`, automated classifiers
- **Similarity:** SBERTScore
- **Realism:** Perplexity ratio

**Takeaway: Reasonably high F1 scores except for voice change**



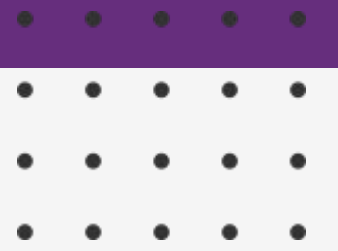


**03**

# **Auditing prompt sensitivity**



# Experimental settings



## Dataset

- Full BBQ Dataset (9 bias categories)
- 8 subsets of MMLU

## LLM Targets

- 4 families: MPT, Falcon, LLaMA 3, Gemma 3
- Deterministic decoding

## Baseline

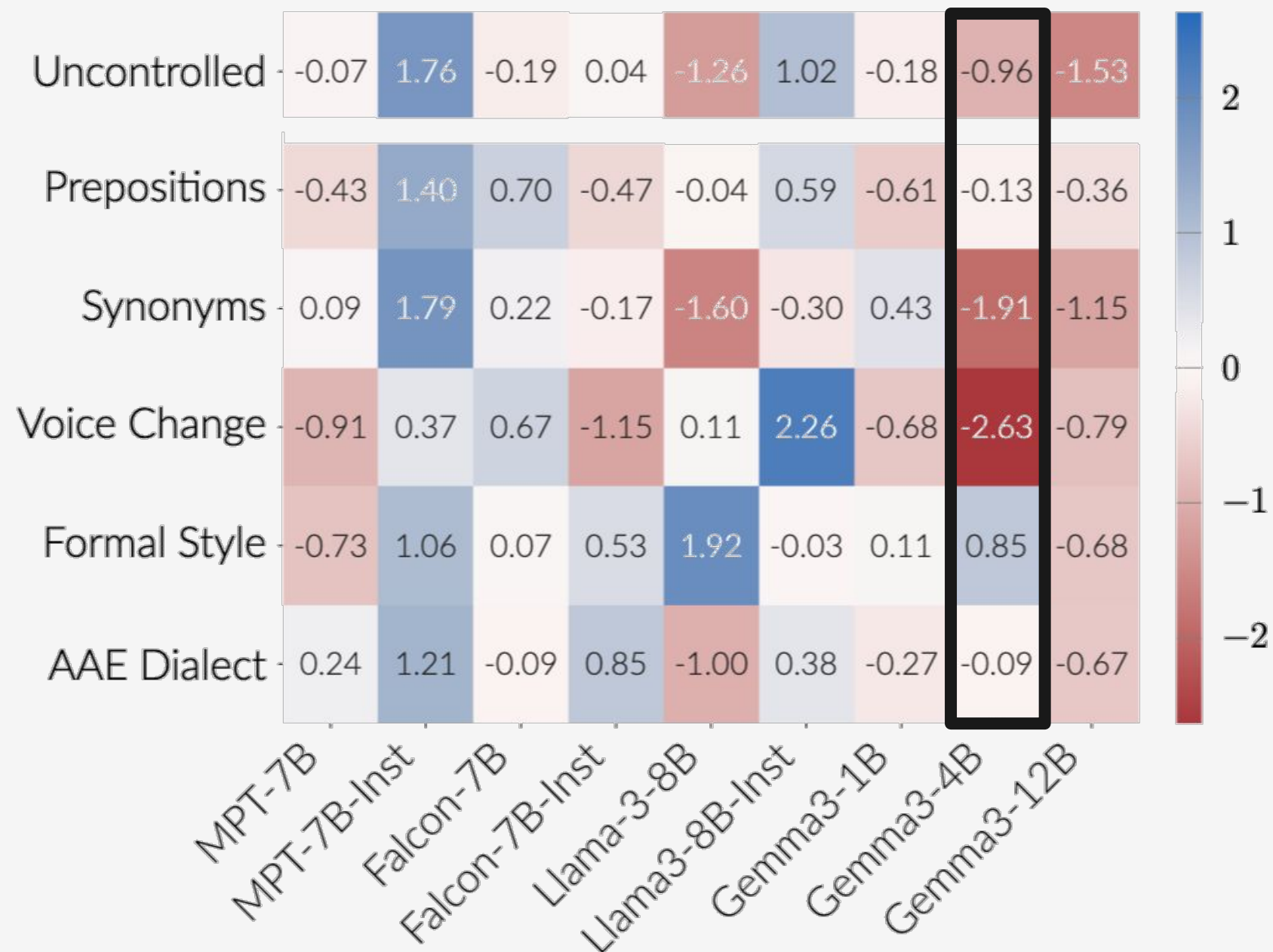
Uncontrolled paraphrasing with no specified paraphrase type or instructions

## Metric

Relative difference in accuracy, for responses to paraphrased and original prompts



# Hidden sensitivities



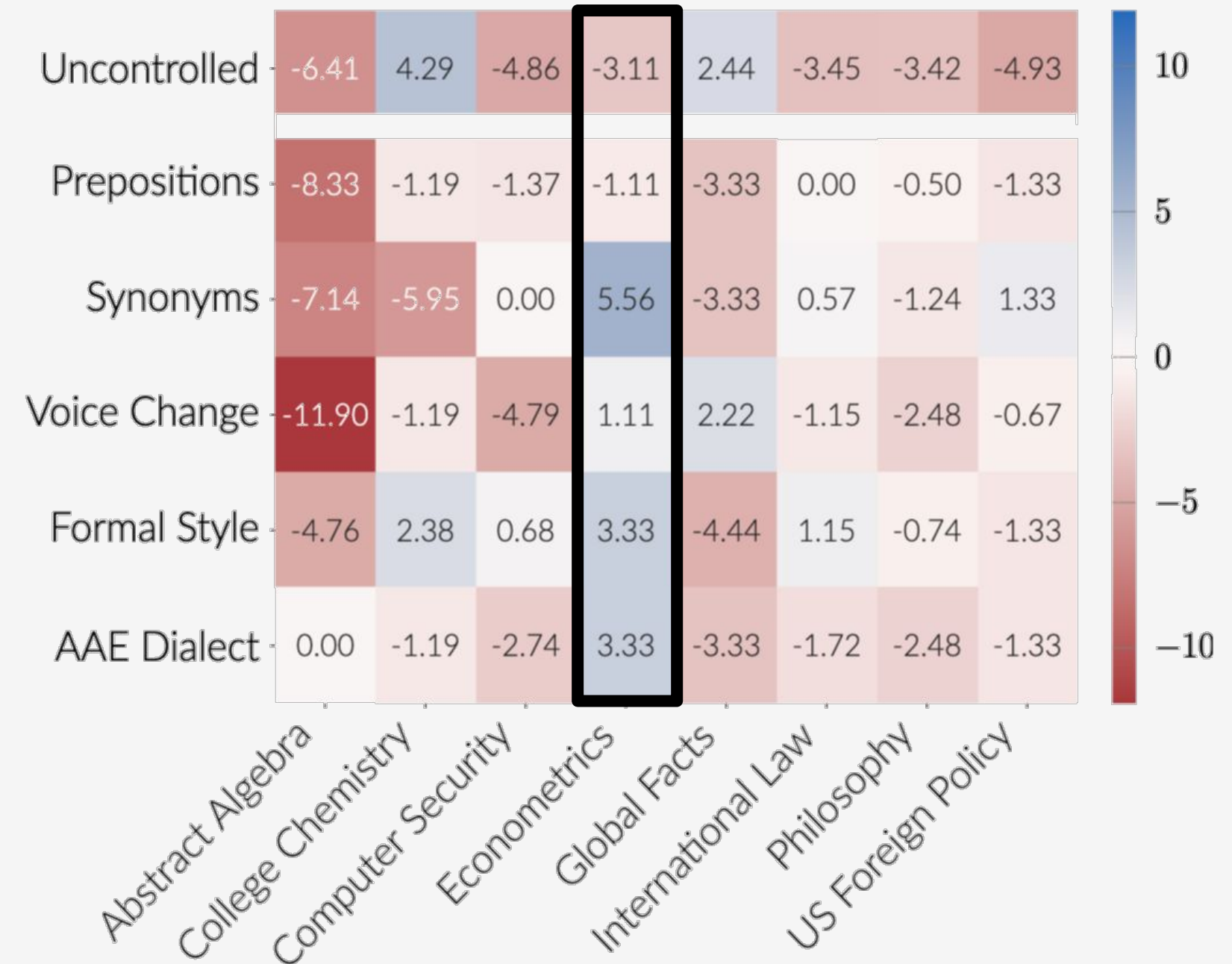
Relative difference in accuracy (%), per paraphrase type and target model, for the BBQ dataset

- Specific paraphrase types trigger **large or divergent shifts**
- Paraphrase-specific effects **invisible** in uncontrolled baseline
- AUGMENT enables **fine-grained analysis**

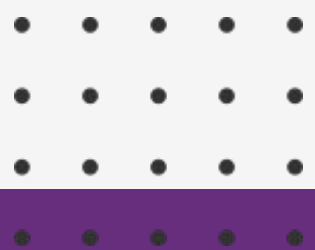


# Divergent patterns

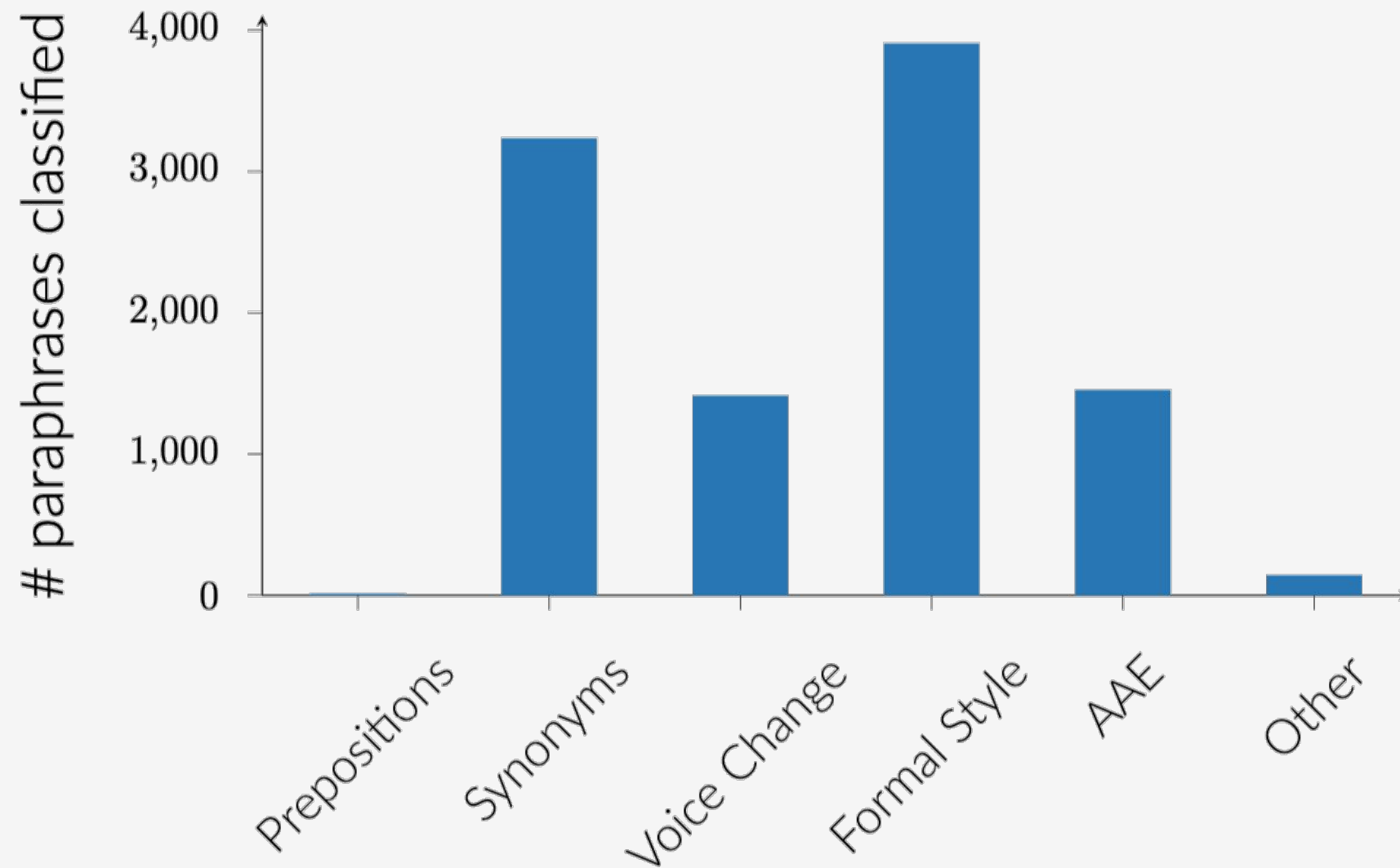
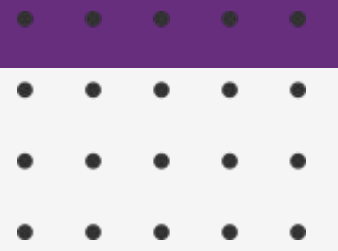
- Strong **subset-level** sensitivities under controlled paraphrases
- AUGMENT reveals **oppositions** with uncontrolled results
- Averaging over unknown types **masks** meaningful effects



Relative difference in accuracy (%), per paraphrase type and MMLU data subset, for Gemma3-12B



# Uncontrolled paraphrases



Classification results of the uncontrolled paraphrases for the BBQ dataset

- **Classification** of the uncontrolled paraphrases, with automatic filtering rules
- **Imbalanced** paraphrase distribution
- **Missed** key linguistic variations





# 04

# Conclusion & Future Work



# Conclusion

- **AUGMENT**: Controlled paraphrasing auditing framework
- Enables a **fine-grained** analysis and diagnostic
- Specific paraphrase types trigger **large or divergent shifts**

## Future Work

- **Scope**: Extend evaluation to open-ended questions
- **Filtering**: Improve robustness of automatic filtering methods
- **Language & Data**: Expand paraphrase taxonomy



# Thanks!

*clea.chataigner@mila.quebec*  
*rebecca.ma@uwaterloo.ca*



*Paper*



*Code*