



CAspER in the Machine

Insights into Character Variety in LLM-Generated Stories



Anneliese Brei¹, Abhishek Sharma², Nicholas Sanaie¹, Lu Wang³, Snigdha Chaturvedi¹

(1) University of North Carolina at Chapel Hill, Department of Computer Science;
(3) University of Michigan, Computer Science and Engineering Division

(2) Georgia Institute of Technology, School of Computer Science;

Contact: abrei@cs.unc.edu

Overview

- Adapt narratological methods to automatically understand characters at level of portrayal;
- Construct high-quality labeled dataset of human-written & LLM-generated short stories;
- Compare and analyze LLM-generated & human-written characters.

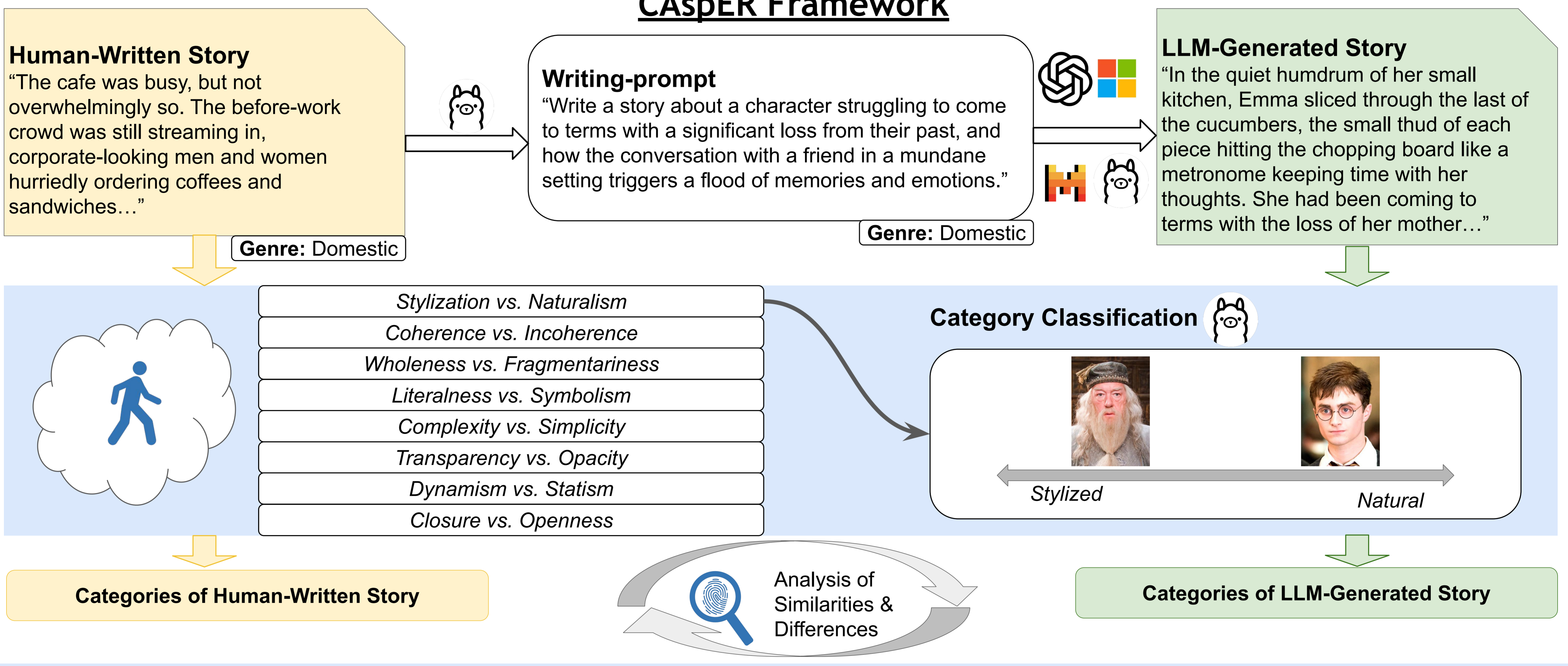
The Task

- CAspER: method for analyzing 8 axes of character portrayal;
- Using CAspER, address:
 - Do LLM-generated stories and human-written stories have similar characters?
 - Do LLMs generate stories with a variety of characters?

Experiments

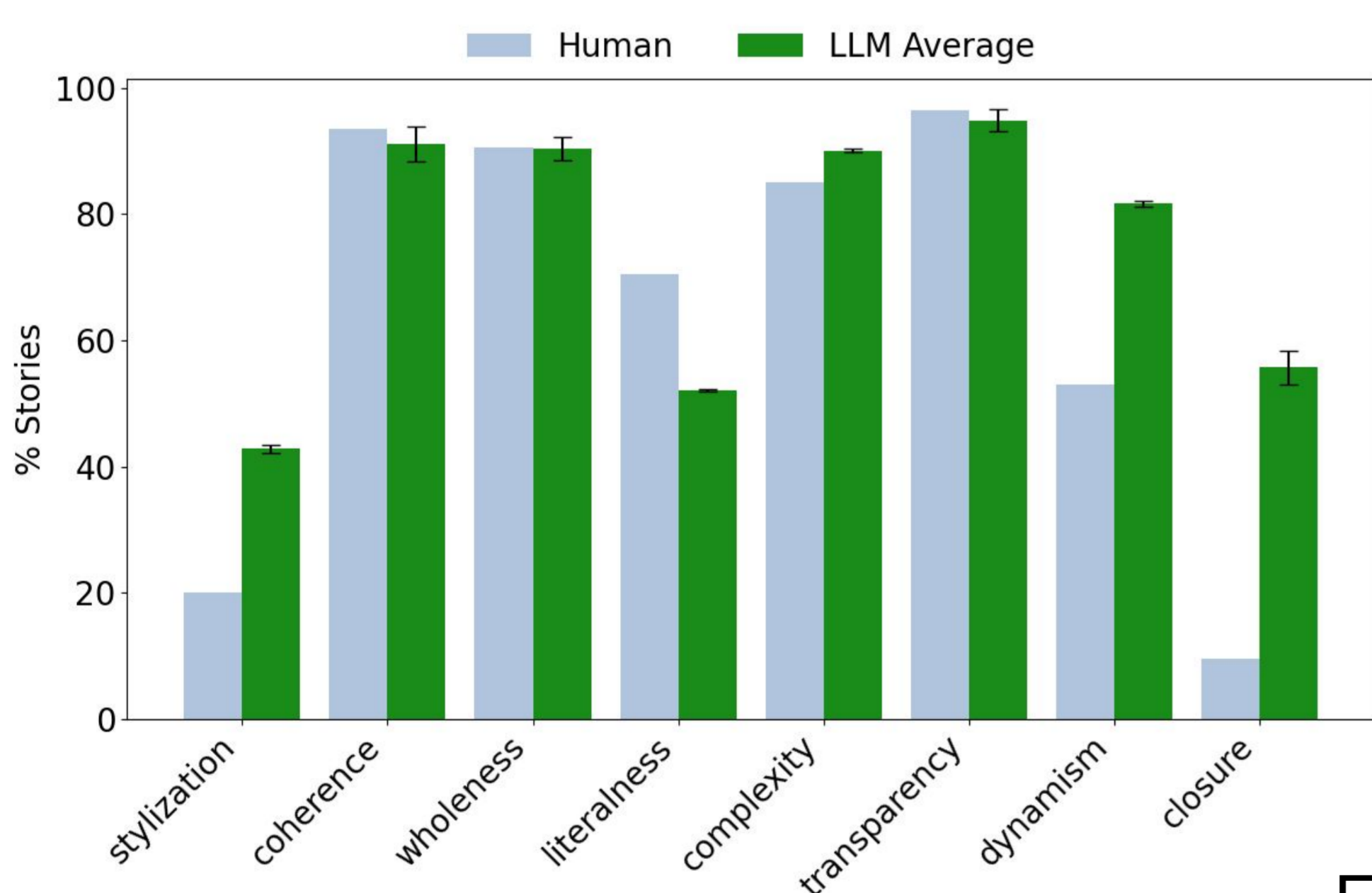
- Corpus creation:**
 - 200 human-written from *r/shortstories*
 - 4400 LLM-generated from 8 models (*Llama, Phi, Mistral, GPT4o-mini*)
- Character classification:** using LLM-judge evaluated on hand-labeled testset of 100 stories with average F1-macro score=84.64.

CAspER Framework

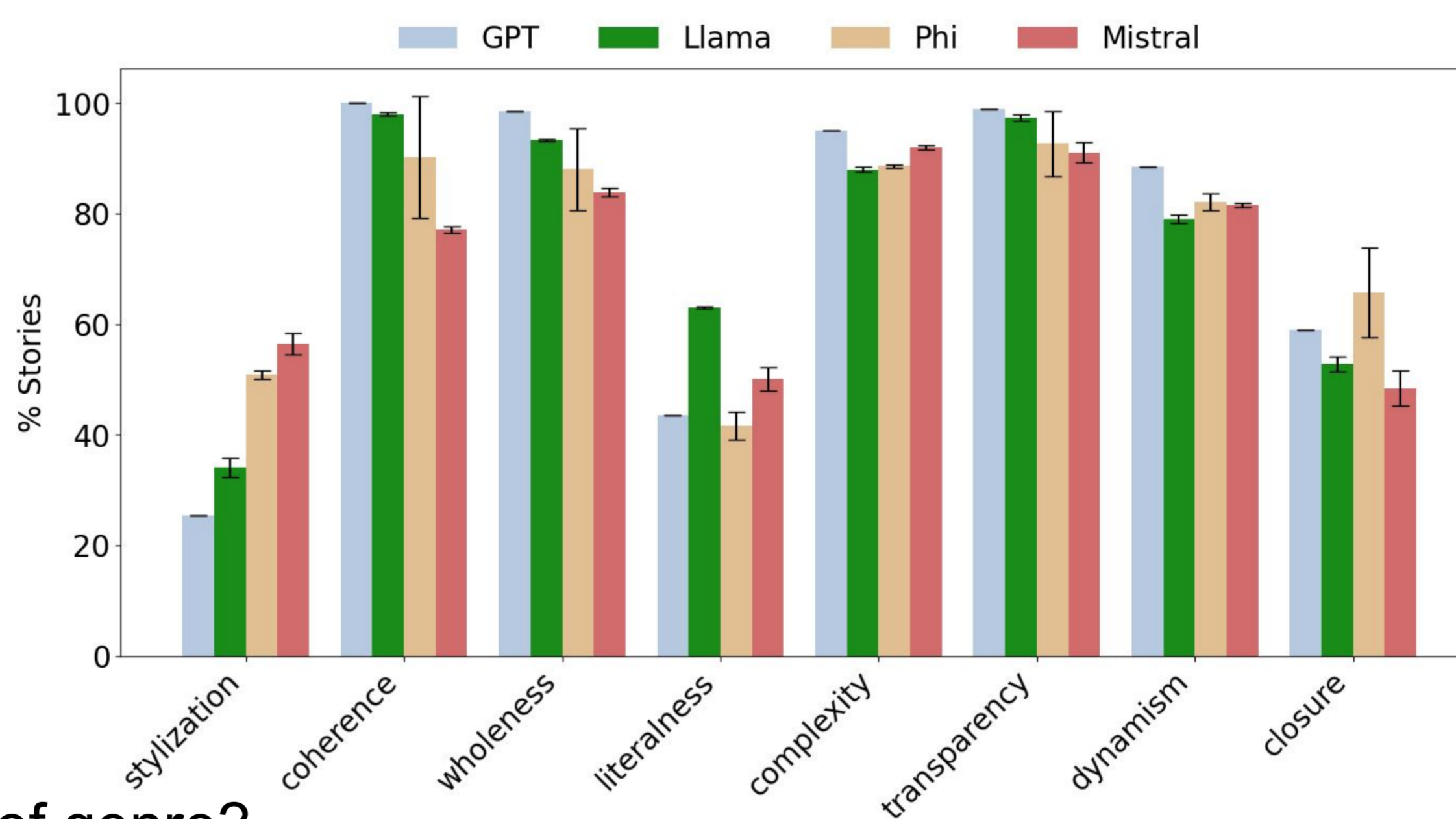


Analysis (Highlights)

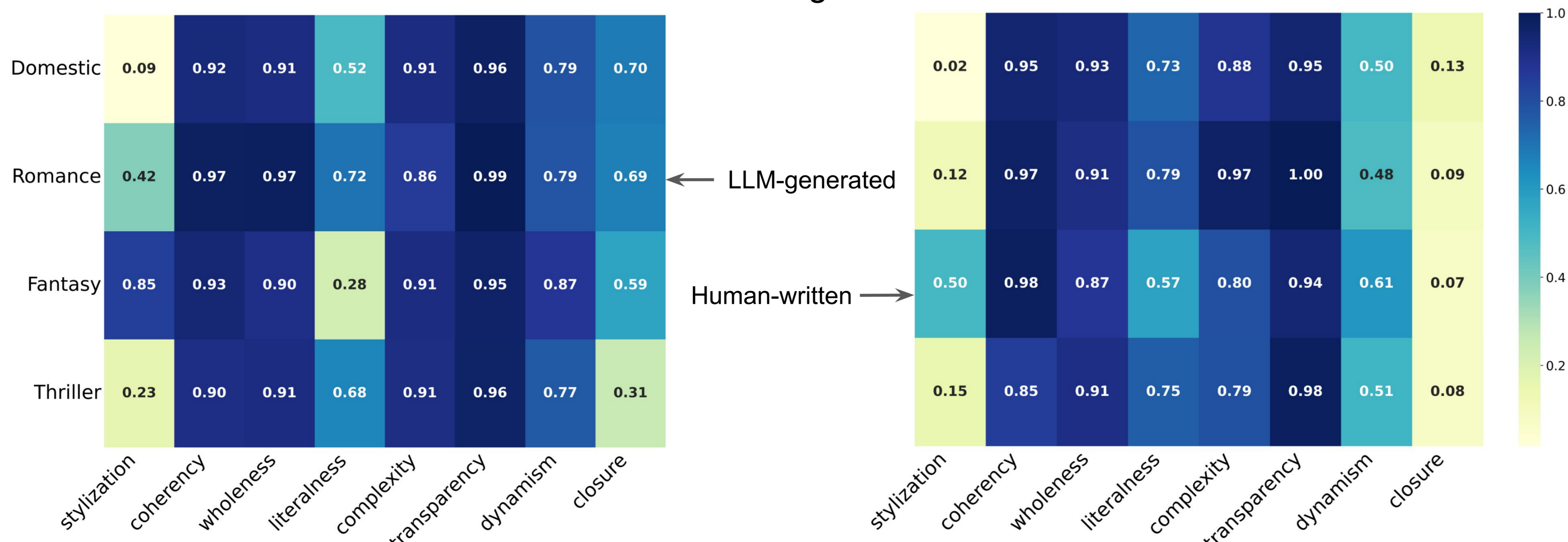
LLM vs. human-written characters?



Characters by different LLM families?



Effect of genre?



Key Takeaways

- LLMs-generated stories are more likely to "play it safe."
- Within family, size does not change types of characters generated.
- Phi family generates most diverse; Llama family generates least diverse characters.
- Genre plays major role in types of character
- When re-generating stories from a prompt, *literalness* & *closure* vary most.