







Idiosyncrasies in Large Language Models

Mingjie Sun*, Yida Yin*, Zhiqiu Xu, J Zico Kolter, Zhuang Liu

Carnegie Mellon University, UC Berkeley, University of Pennsylvania, Princeton University

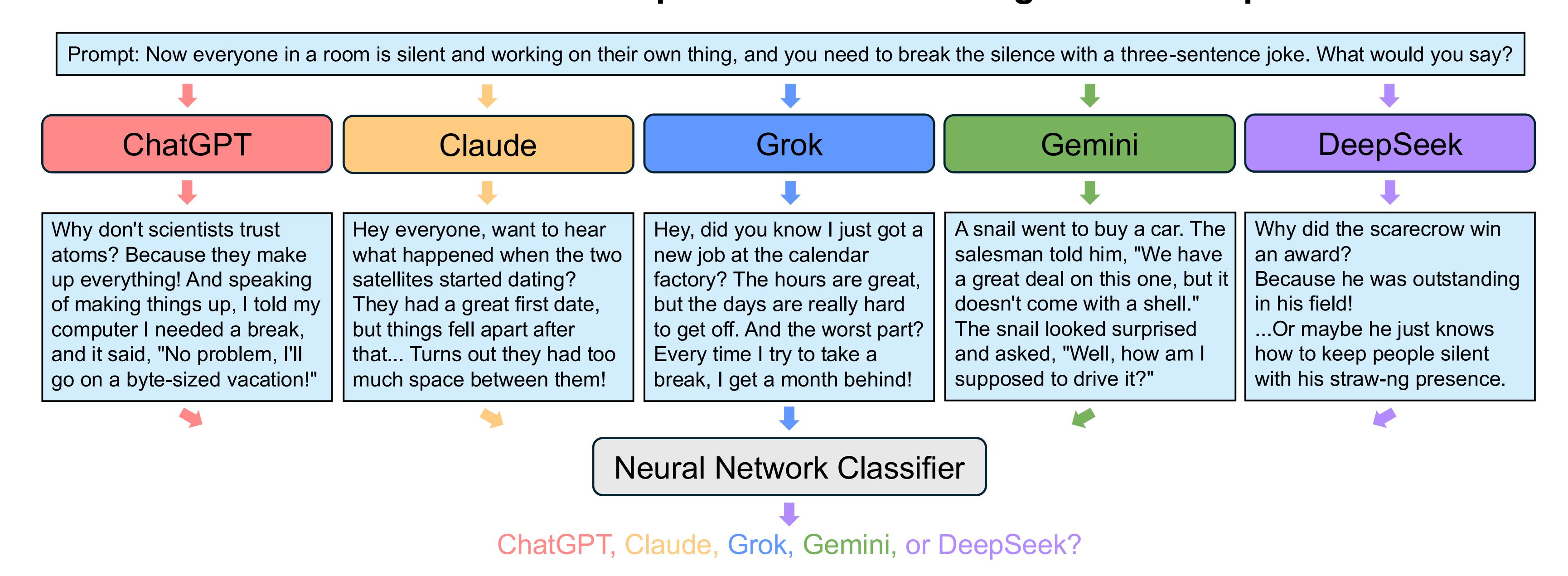
Background

With new LLMs being released every day, what distinguishes one model from another?



Method and Models

We train a neural network classifier to predict which model generates a specific text.



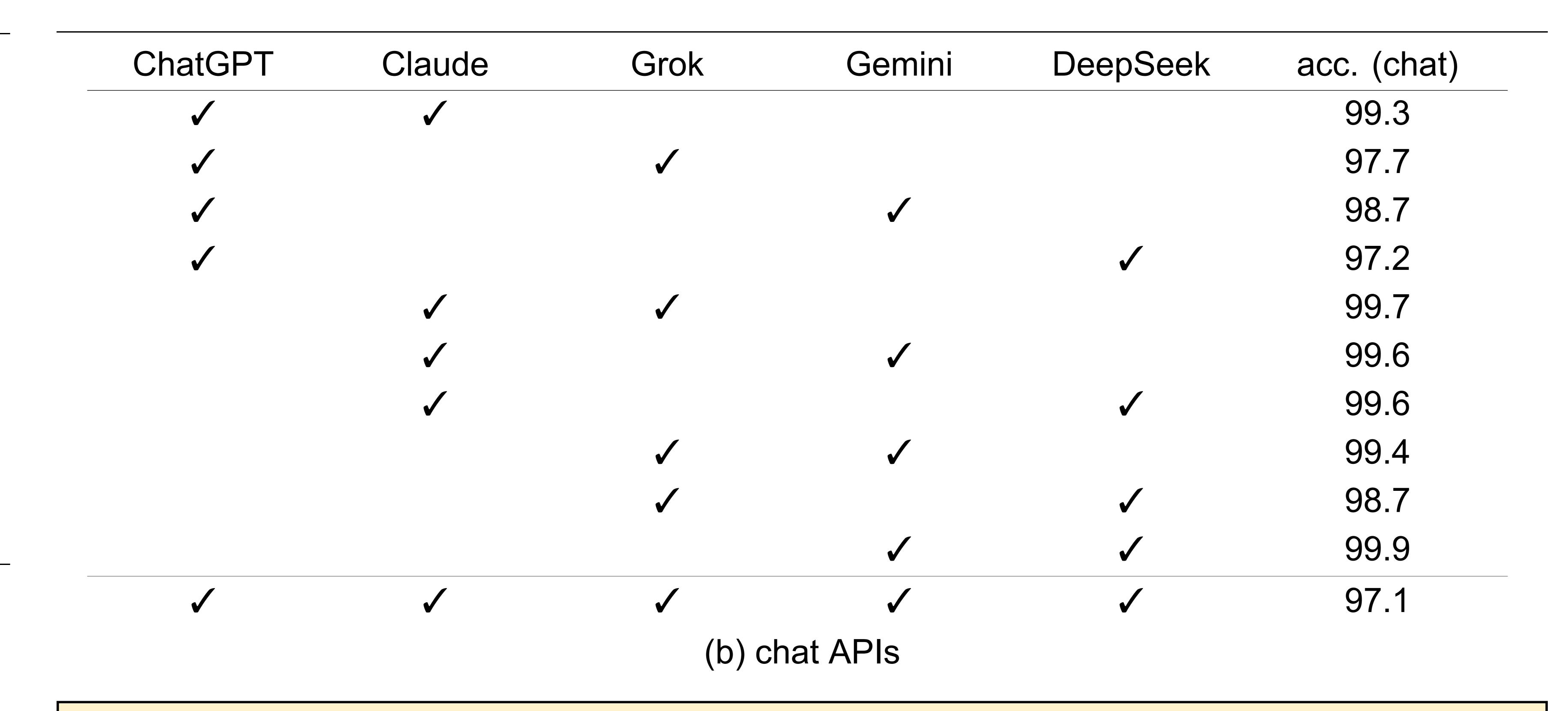
We study the above classification task across three groups of LLMs:

- 1. Chat APIs ("chat"): ChatGPT, Claude, Grok, Gemini, and DeepSeek
- 2. Instruct LLMs ("instruct"): Llama, Gemma, Qwen, and Mistral
- 3. Base LLMs ("base"): base versions of instruct LLMs

Evaluating Idiosyncrasies in LLMs

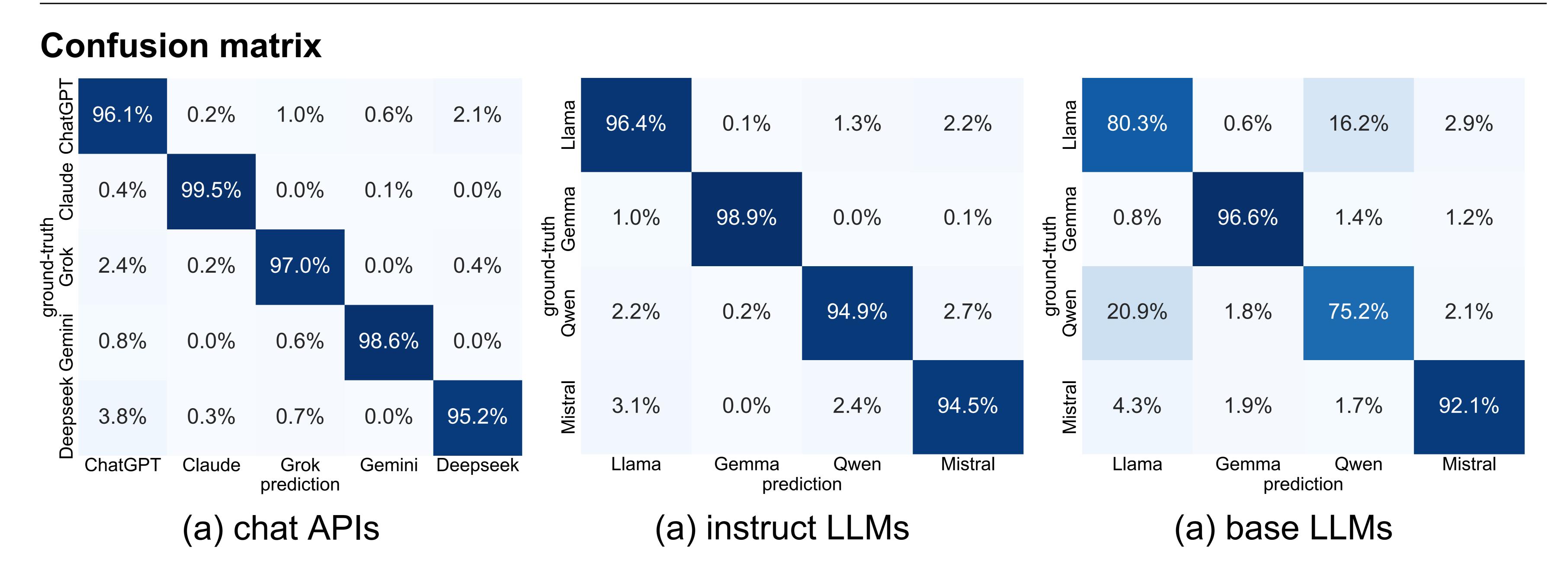
High accuracies are observed across different model families.

Llama	Gemma	Qwen	Mistral	acc. (instruct)	acc. (base)			
				99.9	98.3			
				97.8	81.7			
				97.0	96.3			
				99.9	98.3			
				99.9	98.4			
				96.1	95.7			
				96.3	87.3			
(a) instruct and base LLMs								



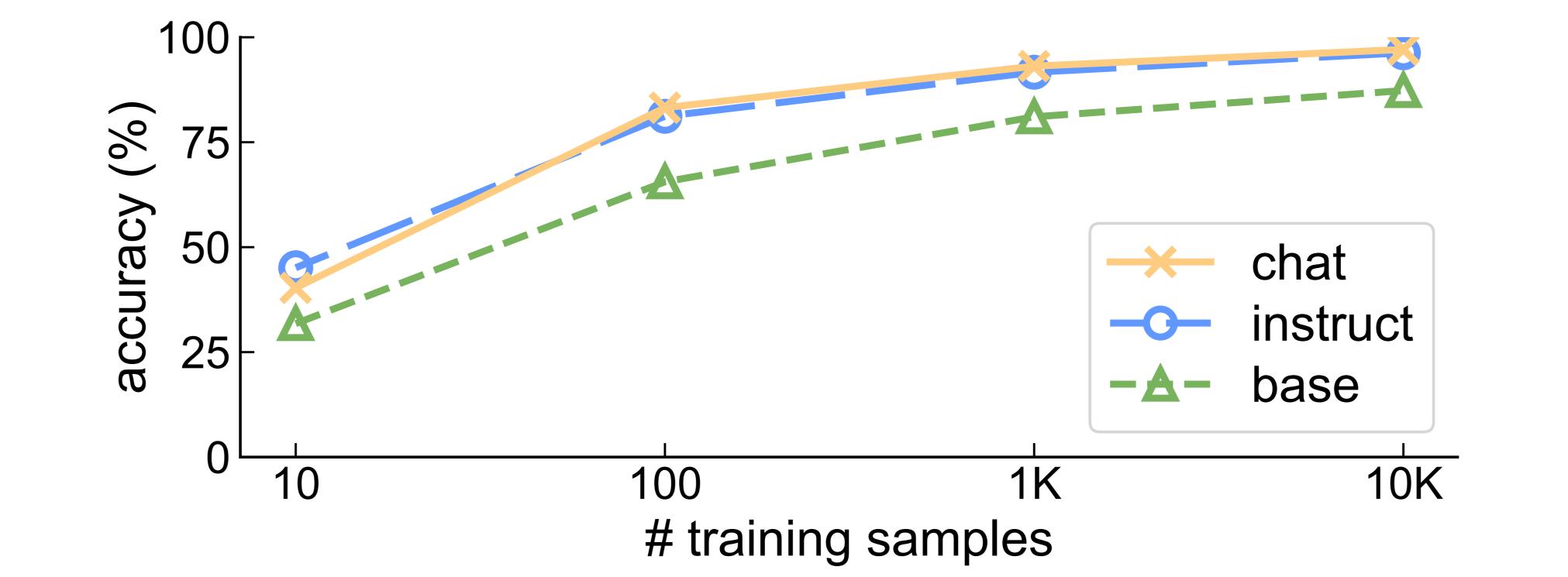
Each LLM has linguistic fingerprints that can identify the model behind a given text!

Analysis



Similar behaviors to text classification

method	chat	instruct	base
ELMo	90.8	91.0	69.8
BERT	91.1	91.5	66.0
T5	90.5	89.8	67.9
GPT-2	92.1	92.3	80.2
LLM2vec	97.1	96.3	87.3



PRINCETON PRINCETON

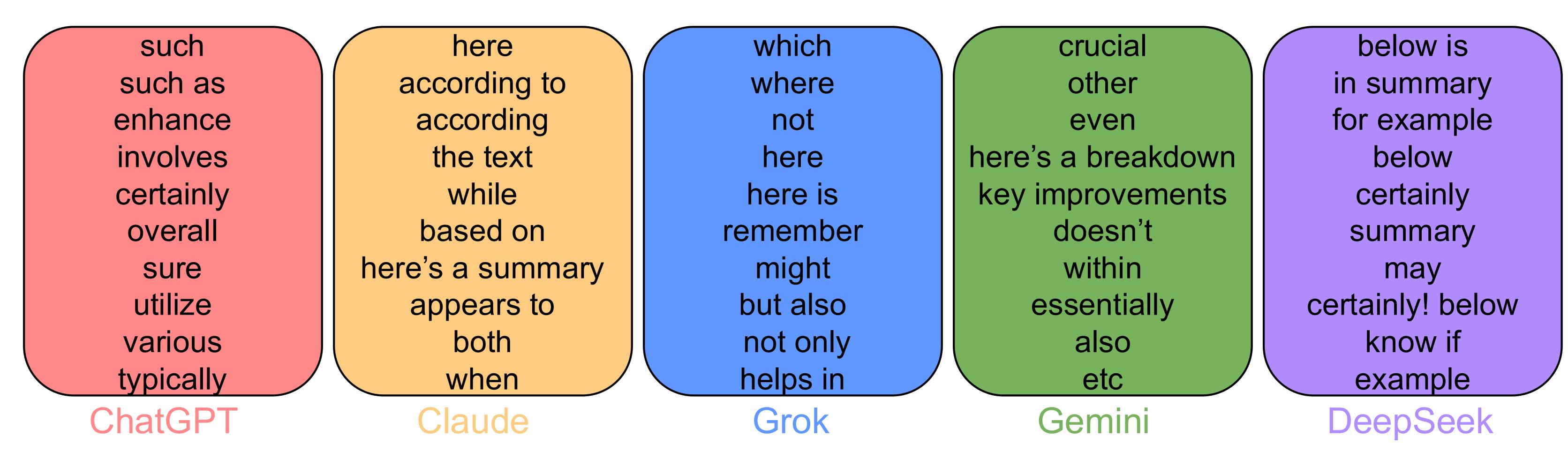




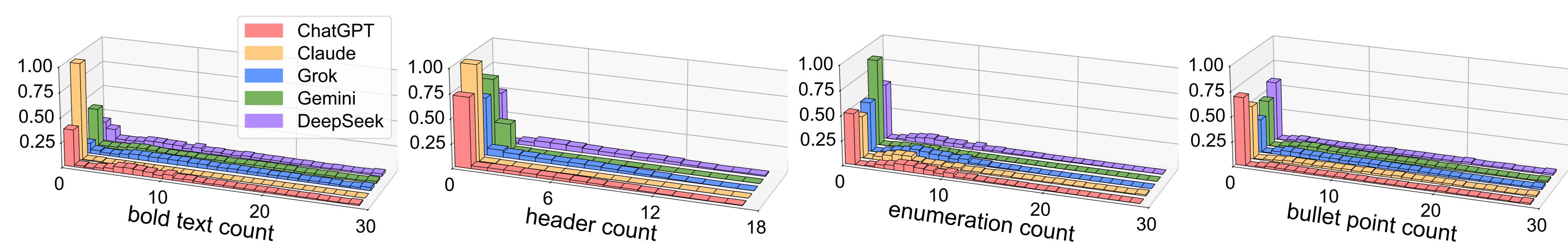


Concrete Idiosyncrasies in LLMs

Characteristic phrases



Unique markdown formatting elements



Open-ended language analysis

- 1. Descriptive and Detailed Tone: Often uses narrative styles with an informative, engaging, or vivid tone.
- 2. Specific and Technical Word: Employs descriptive and technical vocabulary, enhancing depth and specificity.
- 3. Structured and Contextual Opening Lines: Typically begins with context-setting or narrative introductions.
- 4. Markdown Formatting for Organization: Utilizes various markdown elements like headings, lists, and bold text for clarity.
- explanations, background, and broader topics.
 - 5. Comprehensive and In-Depth Content: Offers rich detail, focusing on
- factual, or succinct tone. 2. Functional and Clear Word Choices: Prefers simple or action-oriented
 - language prioritizing clarity and practicality. 3. Immediate and Direct Opening Lines: Often starts with a

1. Concise and Straightforward Tone: Generally adopts a more direct,

- straightforward statement or summary without extended context. 4. Minimal Markdown or List Use: Relies on plain lists or simple
- formatting for quick reference.
- 5. Focused and Summarized Content: Concentrates on essential points and specific phenomena, avoiding extensive detail.

Example responses from (and

Our products feature innovative	According to the text, Kai Fusser	1. Deliver Exceptional Service: The	1. Deliver Exceptional Service
sustainable materials, such as	believes that traditional cardio	foundation of word-of-mouth	Consistently exceed customer
Certainly! If you're looking for cheese	Based on the text provided, here are	marketing is consistent excellence.	expectations
alternatives to replace Brie in your	the key details about Armon Binns'	Providing top-notch services or	 Focus on quality and attention
Overall, while there are challenges,	While many winter sports in the	Ingredients:	Ingredients:
Tanzania is making progress	Pyrenees are similar to those found	• 2 (3 oz) packages of orange-	• 2 boxes orange-flavored Jello
Sure! Here's a simple guide to cooking	This appears to be a fragment of	flavored Jello	 1 can evaporated milk
	poetry that creates a pastoral	 1 cup tonic water (this is what 	Tonic water
ChatGPT	Claude	ChatGPT	Claude

Implication

Inferring model similarity

