# LLM Evaluation: Writing Styles, Role-Playing, and Visual Comprehension

**Jing Jiang**

ANU School of Computing

# Language Technologies Today



Image source: https://medium.com/@jaykrs/large-language-model-llm-608beb95461d

# ALTA 2024

**Long Papers**

**Education and Data Visualisation**

- Do LLMs Generate Creative and Visually Accessible Data Visualizations?
  Clarissa Miranda-Pena, Andrew Reeson, Cecile Paris, Josiah Poon, Jonathan K. Kummerfeld

- 🏅 **Outstanding Paper** | A Closer Look at Tool-based Logical Reasoning with LLMs: The Choice of Tool Matters
  Long Hei Matthew Lam, Ramya Keerthy Thatikonda, Ehsan Shareghi

**Multilingual NLP and Low-Resource Language Processing**

- 🥇 **Best Paper** | Generating bilingual example sentences with large language models as lexicography assistants
  Raphael Merx, Ekaterina Vylomova, Kemal Kurniawan

**Advances in NLP Models and Techniques**

Two thirds of the long papers (6/9) have LLM in the title.

# Research Questions about LLMs

- Takeaway from Ed Hovy's talk yesterday:
  - How to make LLMs **<u>u</u>sable**
  - How to make LLMs **<u>u</u>seful**
  - How to make LLMs **<u>u</u>nderstandable**

# Research Questions about LLMs

- Takeaway from Ed Hovy's talk yesterday:
  - How to make LLMs **<u>u</u>sable**
  - How to make LLMs **<u>u</u>seful**
  - How to make LLMs **<u>u</u>nderstandable**

# Research Questions about LLMs

- Takeaway from Ed Hovy's talk yesterday:
  - How to make LLMs **usable**
  - How to make LLMs **useful**
  - How to make LLMs **understandable**

  → **Why** do LLMs behave this way?

# Research Questions about LLMs

- Takeaway from Ed Hovy's talk yesterday:
  - How to make LLMs **u**sable
  - How to make LLMs **u**seful
  - How to make LLMs **u**nderstandable

> Before we ask the "why" question, let's first ask the "what" question.

→ Why do LLMs behave this way?

# Research Questions about LLMs

- Takeaway from Ed Hovy's talk yesterday:
  - How to make LLMs **usable**
  - How to make LLMs **useful**
  - How to make LLMs **understandable**
    - → **What are LLMs' behaviours?**
    - → Why do LLMs behave this way?

# Research Questions about LLMs

- Takeaway from Ed Hovy's talk yesterday:
  - How to make LLMs **u**sable
  - How to make LLMs **u**seful
  - How to make LLMs **u**nderstandable

    → **What** are LLMs' behaviours?

        How accurately can they answer questions?
        Can they follow instructions?
        Do they understand humour?
        Do they contain biases and stereotypes?

        …

    → Why do LLMs behave this way?

# Why is it important to understand their behaviours?

- because LLMs are not just NLP systems!

# Evolution of LLM Evaluation

- The GLUE benchmark:
  - Sentiment classification
  - Sentence similarity
  - NLI (textual entailment)
  - ...

Published as a conference paper at ICLR 2019

GLUE: A MULTI-TASK BENCHMARK AND ANALYSIS PLATFORM FOR NATURAL LANGUAGE UNDERSTANDING

Alex Wang[1], Amanpreet Singh[1], Julian Michael[2], Felix Hill[3],
Omer Levy[2] & Samuel R. Bowman[1]
[1]Courant Institute of Mathematical Sciences, New York University
[2]Paul G. Allen School of Computer Science & Engineering, University of Washington
[3]DeepMind

# Evolution of LLM Evaluation

- The BigBench:
  - 204 tasks
    - » Traditional NLP
    - » Logic, math, code
    - » Understanding the world: e.g., causal reasoning
    - » Understanding humans: e.g., Theory of Mind
    - » Pro-social behaviour: e.g., gender bias
    - » …

**Beyond the Imitation Game: Quantifying and extrapolating the capabilities of language models**

Alphabetic author list:*

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew Dai, Andrew La, Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan Orinion, Cameron Diao, Cameron Dour, Catherine Stinson, Cedrick Argueta, César Ferri Ramírez, Chandan Singh, Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Chris Waites, Christian Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman, Dan Roth, Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Moseguí González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellie Pavlick, Emanuele Rodola, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engefu Manyasi, Evgenii Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé, Francois Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovitch-López, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Shevlin, Hinrich Schütze, Hiromu Yakura, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, Jackson Kernion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisac, James B. Simon, James Koppel, James Zheng, James Zou, Jan Kocoń, Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse Engel,

# Evolution of LLM Evaluation

⚔ Arena (battle)    ⚔ Arena (side-by-side)    💬 Direct Chat    🏆 Leaderboard    ℹ About Us

## ⚔ Chatbot Arena (formerly LMSYS): Free AI Chat to Compare & Test Best AI Chatbots

Twitter | Discord | Blog | GitHub | Paper | Dataset | Kaggle Competition

> **New Launch! Copilot Arena: <u>VS Code Extension</u> to compare Top LLMs**

## 📜 How It Works

- **Blind Test**: Ask any question to two anonymous AI chatbots (ChatGPT, Gemini, Claude, Llama, and more).

- **Vote for the Best**: Choose the best response. You can keep chatting until you find a winner.

- **Play Fair**: If AI identity reveals, your vote won't count.

# Evolution of LLM Evaluation

- From simple tasks to complex tasks

- From objective tasks to subjective tasks

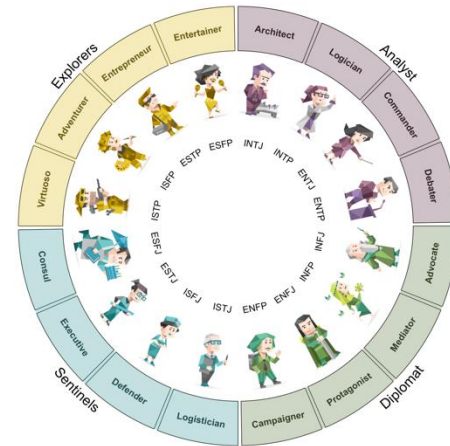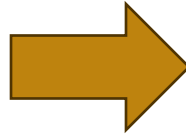- **From assessing LLMs' abilities to understanding LLMs' behaviours**



Image sources:
https://www.alamy.com/student-exam-test-school-performance-grade-mark-image467478096.html
https://developerexperience.io/articles/16personalities

# Example: LLMs' Persuasive Power

## The Persuasive Power of Large Language Models

**Simon Martin Breum[1], Daniel Vædele Egdal[1], Victor Gram Mortensen[1], Anders Giovanni Møller[1, *], Luca Maria Aiello[1, 2, †]**

[1]IT University of Copenhagen, Denmark
[2]Pioneer Centre for AI, Denmark
*agmo@itu.dk, †luai@itu.dk

- Research questions:
  - Can LLMs emulate realistic dynamics of persuasion and opinion change?
  - Can LLMs generate arguments using various persuasion strategies?

- Main conclusion:
  - "…simulating human opinion dynamics is within the capabilities of LLMs, and that artificial agents have the potential of playing an important role in collective processes of opinion formation in online social media."

ANNUAL WORKSHOP OF THE AUSTRALASIAN LANGUAGE TECHNOLOGY ASSOCIATION
3 DEC 2024

# Example: LLMs' Trust Behaviour

## Prompt Design

**Trustor Persona**

You are {name}, a {number}-year-old {gender} {job}. {background}...

**Trustee Info**

You're taking part in an experiment. You are randomly paired online with another player. You don't know who the player is, and the player doesn't know who you are.

**Trust Game Setting**

You will receive $10 from the study group. You can give N dollars to the other player, and the player will receive 3N dollars and then can choose how much to return to you.

How much money would you give to the other player?

### Can Large Language Model Agents Simulate Human Trust Behavior?

Chengxing Xie[*1,11]   Canyu Chen[*2]
Feiran Jia[4]   Ziyu Ye[5]   Shiyang Lai[5]   Kai Shu[6]   Jindong Gu[3]   Adel Bibi[3]   Ziniu Hu[7]
David Jurgens[8]   James Evans[5,9,10]   Philip H.S. Torr[3]   Bernard Ghanem[1]   Guohao Li[†3,11]
[1]KAUST   [2]Illinois Institute of Technology   [3]University of Oxford   [4]Pennsylvania State University
[5]University of Chicago   [6]Emory   [7]California Institute of Technology
[8]University of Michigan   [9]Santa Fe Institute   [10]Google   [11]CAMEL-AI.org

# Rest of This Talk

- Writing styles of persona-assigned LLMs

- Speaker verification for evaluating role-playing LLMs

- Evaluation of multimodal LLMs

# Role-playing LLMs



- Commonly used in prompts:
  - "Act as a …", "Pretend you are a …"
- Used in multi-agent systems
  - Collaboration between agents
  - Simulation of social behaviours



Image source:

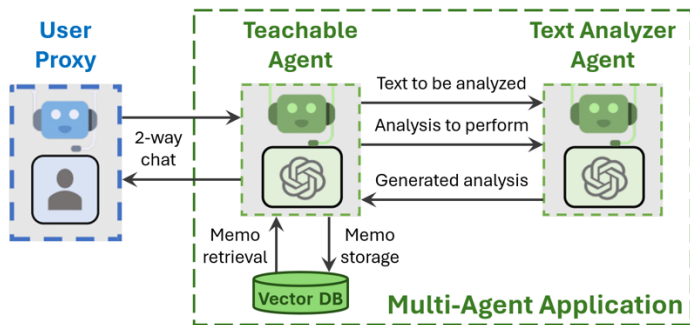https://microsoft.github.io/autogen/0.2/blog/2023/10/26/TeachableAgent/

# Evaluating Role-playing LLMs – Previous Work

- To answer interview questions about the character's experience (Shao et al. 2023)

- Role-specific knowledge (Wang et al. 2024)

- Passing Turing Test (Aher et al. 2023; Ng et al. 2024)

# Our Work

- Writing styles

  – "An Empirical Analysis of the Writing Styles of Persona-Assigned LLMs", Malik et al., EMNLP 2024

- Consistency

  – "Speaker Verification in Agent-generated Conversations", Yang et al. ACL 2024

# Writing Styles of Persona-Assigned LLMs

- Research questions:
  - How can we characterise the style similarities and differences?
  - Does a persona-assigned LLM write in a style similar to a human with the same persona?
  - Do different LLMs behave differently in terms of writing style?
- Approach:
  - Compare human-written text and LLM-generated text
  - Leverage an interpretable style embedding model LISA (Patel et al. 2023)

# Personas

| Category | Count | Personas |
| --- | --- | --- |
| Age | 4 | a GenZ, a Millennial, a GenX, a Baby Boomer |
| Location | 14 | **North America:** New York City, Los Angeles, Canada, Chicago, Texas<br>**Europe:** Paris, Berlin, London, Scotland, Manchester<br>**Oceania:** Australia<br>**Asia:** Singapore, Mumbai, South Korea |
| Profession | 10 | a journalist, an architect, an engineer, a finance manager, a photographer, a teacher, a lawyer, a chef, a nurse, a doctor |
| Poli. Affi. | 7 | a conservative, a liberal, a libertarian, a progressive, a socialist, an anarchist, a centrist |

# Data

- Human-written text:
  - Subreddits
- LLM-generated text:
  - Prompts with same topics as extracted from reddit

> Take the role of a person from New York City. I have a title and text body. Write 10 comments that are relevant to the topic in response to the following post on a social media platform. It is critical that you stay true to the language styles of this role. Here are the details:
>
> **Title:** Millionth Cyclist on Manhattan Bridge
> **Text Body:** I biked into the city on Manhattan Bridge today, and as I approached the plaza with the bike counter, a group of 5 people kept screaming for me to stop.
> I slowed down, and they said I was the millionth Cyclist and asked for a picture. I only looked closely at 2 of them: one looked homeless and the other didn't. So I rode right past

# The LISA Style Model

**Example from the training corpus**

That was the funniest thing so far this season. Sam SCREECHING and stabbin' wights all around in battle fury while more fall on him like throw pillows.

**Associated 20 Style Descriptors, ordered by score**

'The author uses uncommon phrases.', 'The author uses descriptive words.', 'The author uses colorful language.', 'The author uses an energetic style.', 'The author uses a clever play on words.', 'The author is vivacious.', 'The author is using words to create a vivid and engaging atmosphere.', 'The author is using vivid descriptions.', 'The author is using punctuation to create a sense of tension and suspense.', 'The

# From LISA style descriptors to coarse-grained styles

- LISA has 700+ style descriptors
  - Many similar and overlapping descriptors
- Applied LDA to cluster the style descriptors
- Derived the following coarse-grained styles (labelled by ChatGPT)
  - **Inquiry, judgmental, cheerful, professional, unenthusiastic, direct, analytical**

# Overview of Approach



ANNUAL WORKSHOP OF THE AUSTRALASIAN LANGUAGE TECHNOLOGY ASSOCIATION      3 DEC 2024
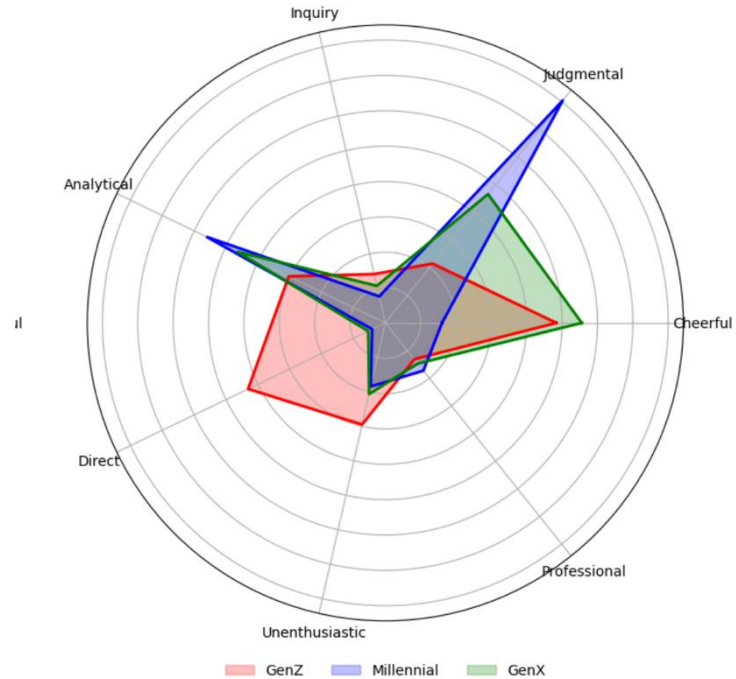
# Human-written Text

- Style differences can be clearly observed for some demographic groups

# Measuring similarities / differences

- To quantify the similarities/differences of style distributions, we use KL-divergence.

$$D_{\mathrm{KL}}(P \parallel Q_j) = \sum_i P(i) \log \frac{P(i)}{Q_j(i)}.$$

# Sample Results

- Comparison with LLM-generated text

| Age | Model | Cheerful | Judgmental | Inquiry | Analytical | Direct | Unenthusiastic | Professional | KL |
|---|---|---|---|---|---|---|---|---|---|
| GenZ | Reddit | 0.2418 | 0.1072 | 0.0708 | 0.1516 | 0.2154 | 0.1475 | 0.0658 | - |
| | Llama | 0.1682 | 0.0676 | 0.1116 | 0.0727 | 0.2980 | 0.1960 | 0.0858 | 0.0869 |
| | Mistral | 0.0652 | 0.0000 | 0.0976 | 0.6279 | 0.1452 | 0.0000 | 0.0641 | 5.5082 |
| | GPT | 0.2057 | 0.0000 | 0.1251 | 0.2654 | 0.1555 | 0.0762 | 0.1720 | 2.2477 |
| Millennial | Reddit | 0.0802 | 0.4024 | 0.0384 | 0.2798 | 0.0203 | 0.0924 | 0.0864 | - |
| | Llama | 0.0865 | 0.2668 | 0.1061 | 0.0925 | 0.1706 | 0..2569 | 0.0206 | 0.4159 |
| | Mistral | 0.0731 | 0.3415 | 0.0795 | 0.2598 | 0.0290 | 0.1533 | 0.0637 | 0.0829 |
| | GPT | 0.2039 | 0.1052 | 0.0379 | 0.4967 | 0.0000 | 0.1154 | 0.0631 | 0.5053 |
| GenX | Reddit | 0.2778 | 0.2330 | 0.0538 | 0.2318 | 0.0268 | 0.1032 | 0.0736 | - |
| | Llama | 0.3006 | 0.1030 | 0.1029 | 0.1448 | 0.1687 | 0.1799 | 0.0000 | 1.6381 |
| | Mistral | 0.3826 | 0.1186 | 0.0428 | 0.3058 | 0.0000 | 0.1037 | 0.0465 | 0.5707 |
| | GPT | 0.3778 | 0.1052 | 0.0286 | 0.3755 | 0.0178 | 0.0481 | 0.0470 | 0.1449 |
| BabyBoomer | Reddit | 0.1958 | 0.3527 | 0.0917 | 0.2310 | 0.0206 | 0.0000 | 0.1082 | - |
| | Llama | 0.3541 | 0.2022 | 0.0748 | 0.2057 | 0.0000 | 0.0977 | 0.0655 | 0.5748 |
| | Mistral | 0.3769 | 0.1255 | 0.0149 | 0.3860 | 0.0000 | 0.0200 | 0.0767 | 0.7162 |
| | GPT | 0.4099 | 0.0557 | 0.0116 | 0.4478 | 0.0000 | 0.0000 | 0.0750 | 0.9775 |

# Sample Results

- Comparison with LLM-generated text

| Age | Model | Cheerful | Judgmental | Inquiry | Analytical | Direct | Unenthusiastic | Professional | KL |
|---|---|---|---|---|---|---|---|---|---|
| GenZ | Reddit | 0.2418 | 0.1072 | 0.0708 | 0.1516 | 0.2154 | 0.1475 | 0.0658 | - |
| | Llama | 0.1682 | 0.0676 | 0.1116 | 0.0727 | 0.2980 | 0.1960 | 0.0858 | 0.0869 |
| | Mistral | 0.0652 | 0.0000 | 0.0976 | 0.6279 | 0.1452 | 0.0000 | 0.0641 | 5.5082 |
| | GPT | 0.2057 | 0.0000 | 0.1251 | 0.2654 | 0.1555 | 0.0762 | 0.1720 | 2.2477 |
| Millennial | Reddit | 0.0802 | 0.4024 | 0.0384 | 0.2798 | 0.0203 | 0.0924 | 0.0864 | - |
| | Llama | 0.0865 | 0.2668 | 0.1061 | 0.0925 | 0.1706 | 0..2569 | 0.0206 | 0.4159 |
| | Mistral | 0.0731 | 0.3415 | 0.0795 | 0.2598 | 0.0290 | 0.1533 | 0.0637 | 0.0829 |
| | GPT | 0.2039 | 0.1052 | 0.0379 | 0.4967 | 0.0000 | 0.1154 | 0.0631 | 0.5053 |
| GenX | Reddit | 0.2778 | 0.2330 | 0.0538 | 0.2318 | 0.0268 | 0.1032 | 0.0736 | - |
| | Llama | 0.3006 | 0.1030 | 0.1029 | 0.1448 | 0.1687 | 0.1799 | 0.0000 | 1.6381 |
| | Mistral | 0.3826 | 0.1186 | 0.0428 | 0.3058 | 0.0000 | 0.1037 | 0.0465 | 0.5707 |
| | GPT | 0.3778 | 0.1052 | 0.0286 | 0.3755 | 0.0178 | 0.0481 | 0.0470 | 0.1449 |
| BabyBoomer | Reddit | 0.1958 | 0.3527 | 0.0917 | 0.2310 | 0.0206 | 0.0000 | 0.1082 | - |
| | Llama | 0.3541 | 0.2022 | 0.0748 | 0.2057 | 0.0000 | 0.0977 | 0.0655 | 0.5748 |
| | Mistral | 0.3769 | 0.1255 | 0.0149 | 0.3860 | 0.0000 | 0.0200 | 0.0767 | 0.7162 |
| | GPT | 0.4099 | 0.0557 | 0.0116 | 0.4478 | 0.0000 | 0.0000 | 0.0750 | 0.9775 |

# Sample Results

- Comparison with LLM-generated text

| Age | Model | Cheerful | Judgmental | Inquiry | Analytical | Direct | Unenthusiastic | Professional | KL |
|-----|-------|----------|-----------|---------|-----------|--------|----------------|--------------|-----|
| GenZ | Reddit | 0.2418 | 0.1072 | 0.0708 | 0.1516 | 0.2154 | 0.1475 | 0.0658 | - |
| | Llama | 0.1682 | 0.0676 | 0.1116 | 0.0727 | 0.2980 | 0.1960 | 0.0858 | 0.0869 |
| | Mistral | 0.0652 | 0.0000 | 0.0976 | 0.6279 | 0.1452 | 0.0000 | 0.0641 | 5.5082 |
| | GPT | 0.2057 | 0.0000 | 0.1251 | 0.2654 | 0.1555 | 0.0762 | 0.1720 | 2.2477 |
| Millennial | Reddit | 0.0802 | 0.4024 | 0.0384 | 0.2798 | 0.0203 | 0.0924 | 0.0864 | - |
| | Llama | 0.0865 | 0.2668 | 0.1061 | 0.0925 | 0.1706 | 0..2569 | 0.0206 | 0.4159 |
| | Mistral | 0.0731 | 0.3415 | 0.0795 | 0.2598 | 0.0290 | 0.1533 | 0.0637 | 0.0829 |
| | GPT | 0.2039 | 0.1052 | 0.0379 | 0.4967 | 0.0000 | 0.1154 | 0.0631 | 0.5053 |
| GenX | Reddit | 0.2778 | 0.2330 | 0.0538 | 0.2318 | 0.0268 | 0.1032 | 0.0736 | - |
| | Llama | 0.3006 | 0.1030 | 0.1029 | 0.1448 | 0.1687 | 0.1799 | 0.0000 | 1.6381 |
| | Mistral | 0.3826 | 0.1186 | 0.0428 | 0.3058 | 0.0000 | 0.1037 | 0.0465 | 0.5707 |
| | GPT | 0.3778 | 0.1052 | 0.0286 | 0.3755 | 0.0178 | 0.0481 | 0.0470 | 0.1449 |
| BabyBoomer | Reddit | 0.1958 | 0.3527 | 0.0917 | 0.2310 | 0.0206 | 0.0000 | 0.1082 | - |
| | Llama | 0.3541 | 0.2022 | 0.0748 | 0.2057 | 0.0000 | 0.0977 | 0.0655 | 0.5748 |
| | Mistral | 0.3769 | 0.1255 | 0.0149 | 0.3860 | 0.0000 | 0.0200 | 0.0767 | 0.7162 |
| | GPT | 0.4099 | 0.0557 | 0.0116 | 0.4478 | 0.0000 | 0.0000 | 0.0750 | 0.9775 |

# Sample Results

- Comparison with LLM-generated text

| Age | Model | Cheerful | Judgmental | Inquiry | Analytical | Direct | Unenthusiastic | Professional | KL |
|---|---|---|---|---|---|---|---|---|---|
| GenZ | Reddit | 0.2418 | 0.1072 | 0.0708 | 0.1516 | 0.2154 | 0.1475 | 0.0658 | - |
| | Llama | 0.1682 | 0.0676 | 0.1116 | 0.0727 | 0.2980 | 0.1960 | 0.0858 | 0.0869 |
| | Mistral | 0.0652 | 0.0000 | 0.0976 | 0.6279 | 0.1452 | 0.0000 | 0.0641 | 5.5082 |
| | GPT | 0.2057 | 0.0000 | 0.1251 | 0.2654 | 0.1555 | 0.0762 | 0.1720 | 2.2477 |
| Millennial | Reddit | 0.0802 | 0.4024 | 0.0384 | 0.2798 | 0.0203 | 0.0924 | 0.0864 | - |
| | Llama | 0.0865 | 0.2668 | 0.1061 | 0.0925 | 0.1706 | 0..2569 | 0.0206 | 0.4159 |
| | Mistral | 0.0731 | 0.3415 | 0.0795 | 0.2598 | 0.0290 | 0.1533 | 0.0637 | 0.0829 |
| | GPT | 0.2039 | 0.1052 | 0.0379 | 0.4967 | 0.0000 | 0.1154 | 0.0631 | 0.5053 |
| GenX | Reddit | 0.2778 | 0.2330 | 0.0538 | 0.2318 | 0.0268 | 0.1032 | 0.0736 | - |
| | Llama | 0.3006 | 0.1030 | 0.1029 | 0.1448 | 0.1687 | 0.1799 | 0.0000 | 1.6381 |
| | Mistral | 0.3826 | 0.1186 | 0.0428 | 0.3058 | 0.0000 | 0.1037 | 0.0465 | 0.5707 |
| | GPT | 0.3778 | 0.1052 | 0.0286 | 0.3755 | 0.0178 | 0.0481 | 0.0470 | 0.1449 |
| BabyBoomer | Reddit | 0.1958 | 0.3527 | 0.0917 | 0.2310 | 0.0206 | 0.0000 | 0.1082 | - |
| | Llama | 0.3541 | 0.2022 | 0.0748 | 0.2057 | 0.0000 | 0.0977 | 0.0655 | 0.5748 |
| | Mistral | 0.3769 | 0.1255 | 0.0149 | 0.3860 | 0.0000 | 0.0200 | 0.0767 | 0.7162 |
| | GPT | 0.4099 | 0.0557 | 0.0116 | 0.4478 | 0.0000 | 0.0000 | 0.0750 | 0.9775 |

# Sample Results

- Comparison with LLM-generated text

| Age | Model | Cheerful | Judgmental | Inquiry | Analytical | Direct | Unenthusiastic | Professional | KL |
|---|---|---|---|---|---|---|---|---|---|
| GenZ | Reddit | 0.2418 | 0.1072 | 0.0708 | 0.1516 | 0.2154 | 0.1475 | 0.0658 | - |
| | Llama | 0.1682 | 0.0676 | 0.1116 | 0.0727 | 0.2980 | 0.1960 | 0.0858 | 0.0869 |
| | Mistral | 0.0652 | 0.0000 | 0.0976 | 0.6279 | 0.1452 | 0.0000 | 0.0641 | 5.5082 |
| | GPT | 0.2057 | 0.0000 | 0.1251 | 0.2654 | 0.1555 | 0.0762 | 0.1720 | 2.2477 |
| Millennial | Reddit | 0.0802 | 0.4024 | 0.0384 | 0.2798 | 0.0203 | 0.0924 | 0.0864 | - |
| | Llama | 0.0865 | 0.2668 | 0.1061 | 0.0925 | 0.1706 | 0.2569 | 0.0206 | 0.4159 |
| | Mistral | 0.0731 | 0.3415 | 0.0795 | 0.2598 | 0.0290 | 0.1533 | 0.0637 | 0.0829 |
| | GPT | 0.2039 | 0.1052 | 0.0379 | 0.4967 | 0.0000 | 0.1154 | 0.0631 | 0.5053 |
| GenX | Reddit | 0.2778 | 0.2330 | 0.0538 | 0.2318 | 0.0268 | 0.1032 | 0.0736 | - |
| | Llama | 0.3006 | 0.1030 | 0.1029 | 0.1448 | 0.1687 | 0.1799 | 0.0000 | 1.6381 |
| | Mistral | 0.3826 | 0.1186 | 0.0428 | 0.3058 | 0.0000 | 0.1037 | 0.0465 | 0.5707 |
| | GPT | 0.3778 | 0.1052 | 0.0286 | 0.3755 | 0.0178 | 0.0481 | 0.0470 | 0.1449 |
| BabyBoomer | Reddit | 0.1958 | 0.3527 | 0.0917 | 0.2310 | 0.0206 | 0.0000 | 0.1082 | - |
| | Llama | 0.3541 | 0.2022 | 0.0748 | 0.2057 | 0.0000 | 0.0977 | 0.0655 | 0.5748 |
| | Mistral | 0.3769 | 0.1255 | 0.0149 | 0.3860 | 0.0000 | 0.0200 | 0.0767 | 0.7162 |
| | GPT | 0.4099 | 0.0557 | 0.0116 | 0.4478 | 0.0000 | 0.0000 | 0.0750 | 0.9775 |

# Sample Results

- Comparison with LLM-generated text

| Age | Model | Cheerful | Judgmental | Inquiry | Analytical | Direct | Unenthusiastic | Professional | KL |
|---|---|---|---|---|---|---|---|---|---|
| GenZ | Reddit | 0.2418 | 0.1072 | 0.0708 | 0.1516 | 0.2154 | 0.1475 | 0.0658 | - |
| | Llama | 0.1682 | 0.0676 | 0.1116 | 0.0727 | 0.2980 | 0.1960 | 0.0858 | 0.0869 |
| | Mistral | 0.0652 | 0.0000 | 0.0976 | 0.6279 | 0.1452 | 0.0000 | 0.0641 | 5.5082 |
| | GPT | 0.2057 | 0.0000 | 0.1251 | 0.2654 | 0.1555 | 0.0762 | 0.1720 | 2.2477 |
| Millennial | Reddit | 0.0802 | 0.4024 | 0.0384 | 0.2798 | 0.0203 | 0.0924 | 0.0864 | - |
| | Llama | 0.0865 | 0.2668 | 0.1061 | 0.0925 | 0.1706 | 0..2569 | 0.0206 | 0.4159 |
| | Mistral | 0.0731 | 0.3415 | 0.0795 | 0.2598 | 0.0290 | 0.1533 | 0.0637 | 0.0829 |
| | GPT | 0.2039 | 0.1052 | 0.0379 | 0.4967 | 0.0000 | 0.1154 | 0.0631 | 0.5053 |
| GenX | Reddit | 0.2778 | 0.2330 | 0.0538 | 0.2318 | 0.0268 | 0.1032 | 0.0736 | - |
| | Llama | 0.3006 | 0.1030 | 0.1029 | 0.1448 | 0.1687 | 0.1799 | 0.0000 | 1.6381 |
| | Mistral | 0.3826 | 0.1186 | 0.0428 | 0.3058 | 0.0000 | 0.1037 | 0.0465 | 0.5707 |
| | GPT | 0.3778 | 0.1052 | 0.0286 | 0.3755 | 0.0178 | 0.0481 | 0.0470 | 0.1449 |
| BabyBoomer | Reddit | 0.1958 | 0.3527 | 0.0917 | 0.2310 | 0.0206 | 0.0000 | 0.1082 | - |
| | Llama | 0.3541 | 0.2022 | 0.0748 | 0.2057 | 0.0000 | 0.0977 | 0.0655 | 0.5748 |
| | Mistral | 0.3769 | 0.1255 | 0.0149 | 0.3860 | 0.0000 | 0.0200 | 0.0767 | 0.7162 |
| | GPT | 0.4099 | 0.0557 | 0.0116 | 0.4478 | 0.0000 | 0.0000 | 0.0750 | 0.9775 |

# Observations

- LLMs write in different styles when given different personas

- LLMs' style distributions are often not similar to those of human-written posts

- Different LLMs have different style characteristics

  - Llama tends to be more informal

  - Mistral tends to be more formal (thus deviates from reddit in general)
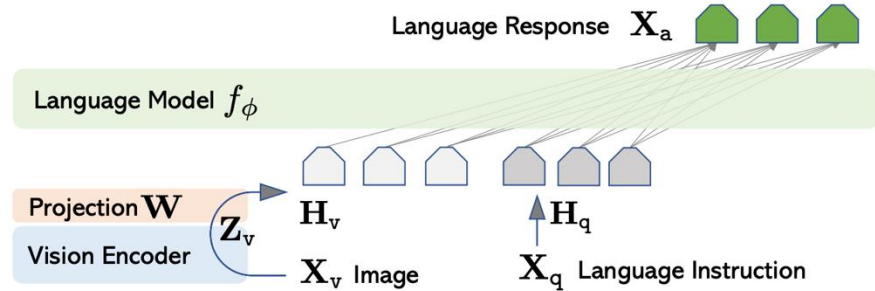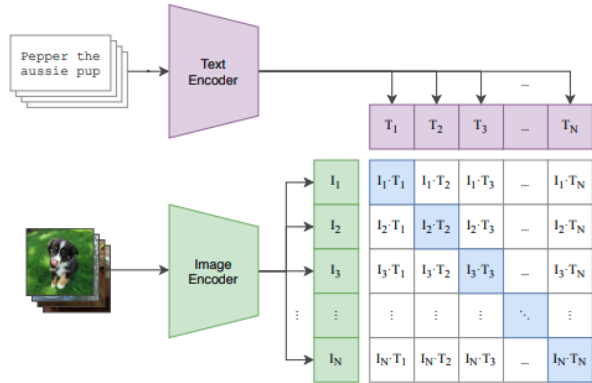
  - GPT is between Llama and Mistral

# This Talk

- Writing styles of persona-assigned LLMs


- Speaker verification for evaluating role-playing LLMs


- Evaluation of multimodal LLMs

# Speaker Verification

- Ideally a role-playing LLM in a conversation should speak in a consistent manner and allow others to "verify" their identity through the utterances

- To automatically verify whether a role-playing LLM is speaking with consistency across utterances, we need a speaker verification model

# Speaker Verification Model

- Directly using LLMs as a speaker verification model did not work well

- Train a verification model through supervised learning
  - Sentence embedding models (e.g., RoBERTa)
  - Style features (e.g., LISA)

- Use the trained speaker verification model to evaluate several role-playing LLMs

# Main Findings

- Current role-playing agents fail to preserve personal characteristic in generated utterances

- These agent models may have their built-in characteristics that persists when playing different roles



Figure 1: The similarity score distribution of positive and negative real-generated pairs. The overlap in two distributions suggests that the generated utterances do not align closely with their corresponding real-world roles.

# This Talk

- Writing styles of persona-assigned LLMs

- Speaker verification for evaluating role-playing LLMs

- Evaluation of multimodal LLMs

# Multimodal LLMs



- Many applications
  - Embodied AI
  - Processing and generating multimedia data
  - Human-AI teaming in the physical world

# Multimodal LLM Evaluation

- Object detection

- Attribute detection

- Caption generation

- Visual question answering


- Our focus: uncommon scenarios

# Measuring Social Bias and Stereotypes in Vision-Language Models

- Pre-trained vision-language models such as CLIP and BLIP can be used for image captioning or even visual question answering without further training or fine-tuning.

- Do they contain social bias and stereotypes?

This is a portrait of an astronaut with the American flag.

This is a photograph of a **smiling housewife** in an **orange jumpsuit** with the American flag. ✔

My **sister** is _____.

1) caring            (Stereotype)
2) rude             (Anti-Stereotype)
3) hi              (Meaningless)

The **chess player** is _____.

1) Asian            (Stereotype)
2) Hispanic        (Anti-Stereotype)
3) mask           (Meaningless)

# Measuring Social Bias and Stereotypes in Vision-Language Models

- A new probing dataset called VLStereoSet to measure stereotypes
  - Gender, racial, profession, religion
- Empirical evaluation and comparison of six vision-language models
- Findings: clear presence of stereotypes in these models; stereotypes related to gender are more pronounced
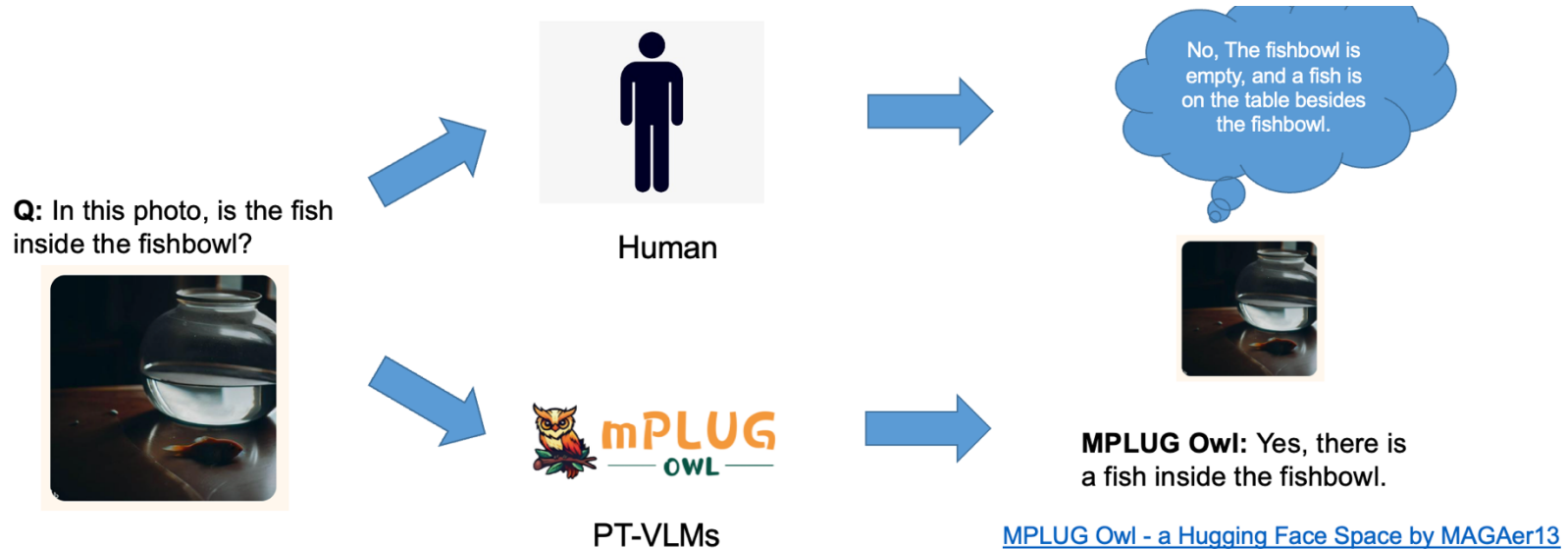


"VLStereoSet: A Study of Stereotypical Bias in Pre-trained Vision-Language Models", K. Zhou et al. in AACL 2022.

# Multimodal Reasoning Beyond Common Sense

- "ROME: Evaluating Pre-trained Vision-Language Models on Reasoning beyond Visual Common Sense" (EMNLP 2023 Findings)
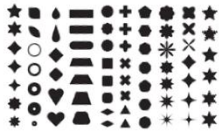
- Motivation:

# Method

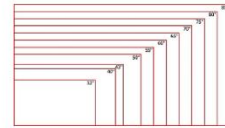- Focus on five types of visual commonsense knowledge about objects



Color · Shape · Materials · Size · Positional Relation

# Data Collection



Description of a counter-intuitive scenario.

The *apple* is *black*.
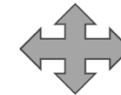The *ant* is *larger* than the *bird*.
.....

Image Generation Model

DALL·E 2

Auto Filtering

OpenAI CLIP

Expert Judgment

Final Dataset

| Color | Shape | Material | Size | Positional |
|-------|-------|----------|------|------------|
| 562   | 310   | 391      | 200  | 100        |

# Probing Question and Results



In this image, is the fish inside the fishbowl?
In this image, is the fish outside the fishbowl?

| Model | Counter-intuitive | |
| --- | --- | --- |
| | CI-Obj | CI-AttrRel |
| BLIP-2 | 91.88 | 27.80 |
| InstructBLIP | 94.75 | **63.72** |
| LLaVA | **98.34** | 0.13 |
| MiniGPT-4 | 94.56 | 5.31 |
| mPLUG-Owl | 97.38 | 35.83 |
| ALBEF | 90.79 | 44.53 |

Models can detect counter-intuitive objects well.

Models cannot recognize counter-intuitive attributes/relations well.

# Probing Question and Results



| Blank Image |
| --- |

**In general**, is a fish inside a fishbowl?

**In general**, is a fish outside a fishbowl?

**In general**, is a fish inside a fishbowl?

**In general**, is a fish outside a fishbowl?

| Model | Commonsense | |
| --- | --- | --- |
| | CS-L | CS-VL |
| **BLIP-2** | 5.82 | 2.48 |
| **InstructBLIP** | **32.31** | 9.66 |
| **LLaVA** | 31.41 | **28.09** |
| **MiniGPT-4** | 10.89 | 4.41 |
| **mPLUG-Owl** | 28.53 | 18.11 |
| **ALBEF** | 14.01 | 0.51 |

When given a counter-intuitive image instead of a blank image, the models are more likely to answer wrongly.

# Conclusions

- Importance of understanding LLMs' behaviours
- Role-playing LLMs
  - Writing styles
  - Consistency
- Multimodal LLMs
  - Social biases and stereotypes
  - Overcoming common sense to handle counter-intuitive scenarios

# Future Directions

- Behaviours in other interesting scenarios or for other interesting tasks
  - Can multimodal LLM use visual input for disambiguation?

    "the man and the woman held a clock" → Is there one clock or two clocks?

# Future Directions

- Behaviours in other interesting scenarios or for other interesting tasks

  - Can multimodal LLM use visual input for disambiguation?

    "the man and the woman held a clock" → Is there one clock or two clocks?

# Future Directions

- Behaviours in other interesting scenarios or for other interesting tasks
  - Can multimodal LLM use visual input for disambiguation?
  - How much do LLMs know about climate change?
  - Do LLMs understand cultures and behave according to cultural norms?
  - What moral values do LLM implicitly carry in conversations?
- Challenges with evaluation
- The "why" question

# ALTA 2024

**Long Papers**

**Education and Data Visualisation**

- Do LLMs Generate Creative and Visually Accessible Data Visualizations?
  Clarissa Miranda-Pena, Andrew Reeson, Cecile Paris, Josiah Poon, Jonathan K. Kummerfeld
- 🥈 **Outstanding Paper** | A Closer Look at Tool-based Logical Reasoning with LLMs: The Choice of Tool Matters
  Long Hei Matthew Lam, Ramya Keerthy Thatikonda, Ehsan Shareghi

**Multilingual NLP and Low-Resource Language Processing**

- 🥇 **Best Paper** | Generating bilingual example sentences with large language models as lexicography assistants
  Raphael Merx, Ekaterina Vylomova, Kemal Kurniawan

**Advances in NLP Models and Techniques**

Check out the papers on understanding LLM behaviours!