# ACL 2024 Sharing
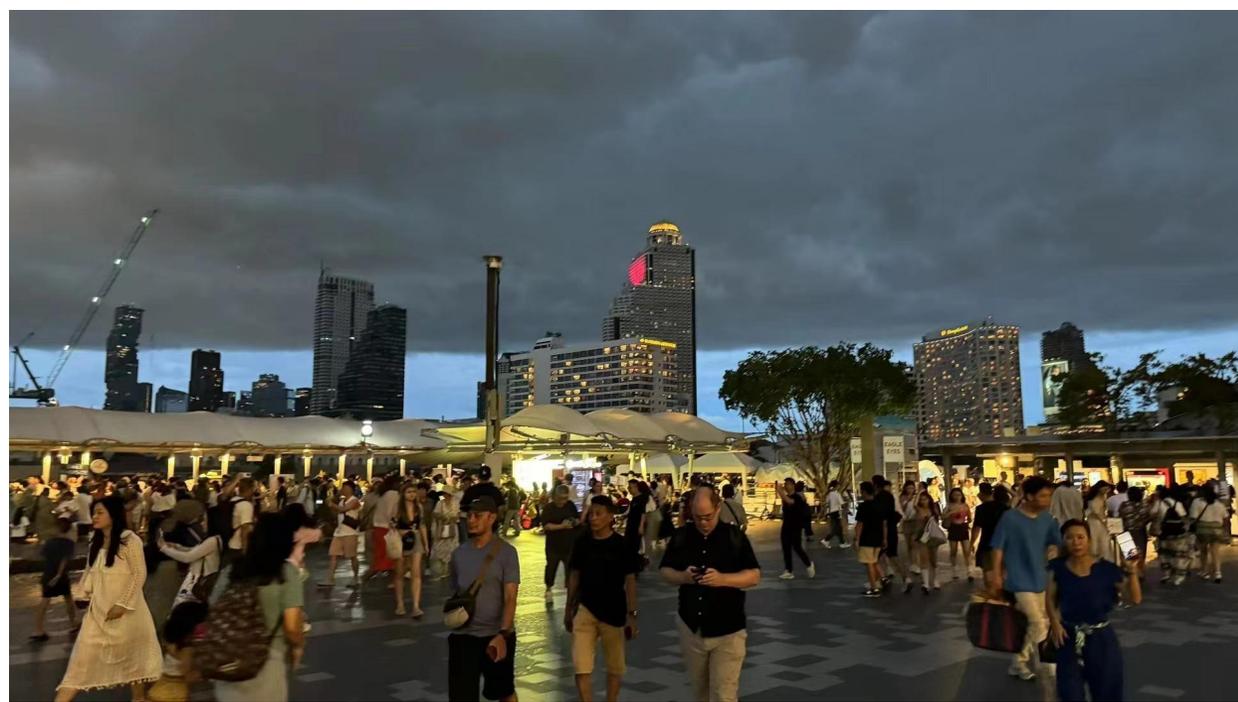
Zhu Liu

2024.09.26

# ACL 2024

- 会议时间：8月11日到16日（17,18 还有2天workshop）

- 会议地点：泰国曼谷的 Centara Grand and Bangkok Convention Centre at CentralWorld

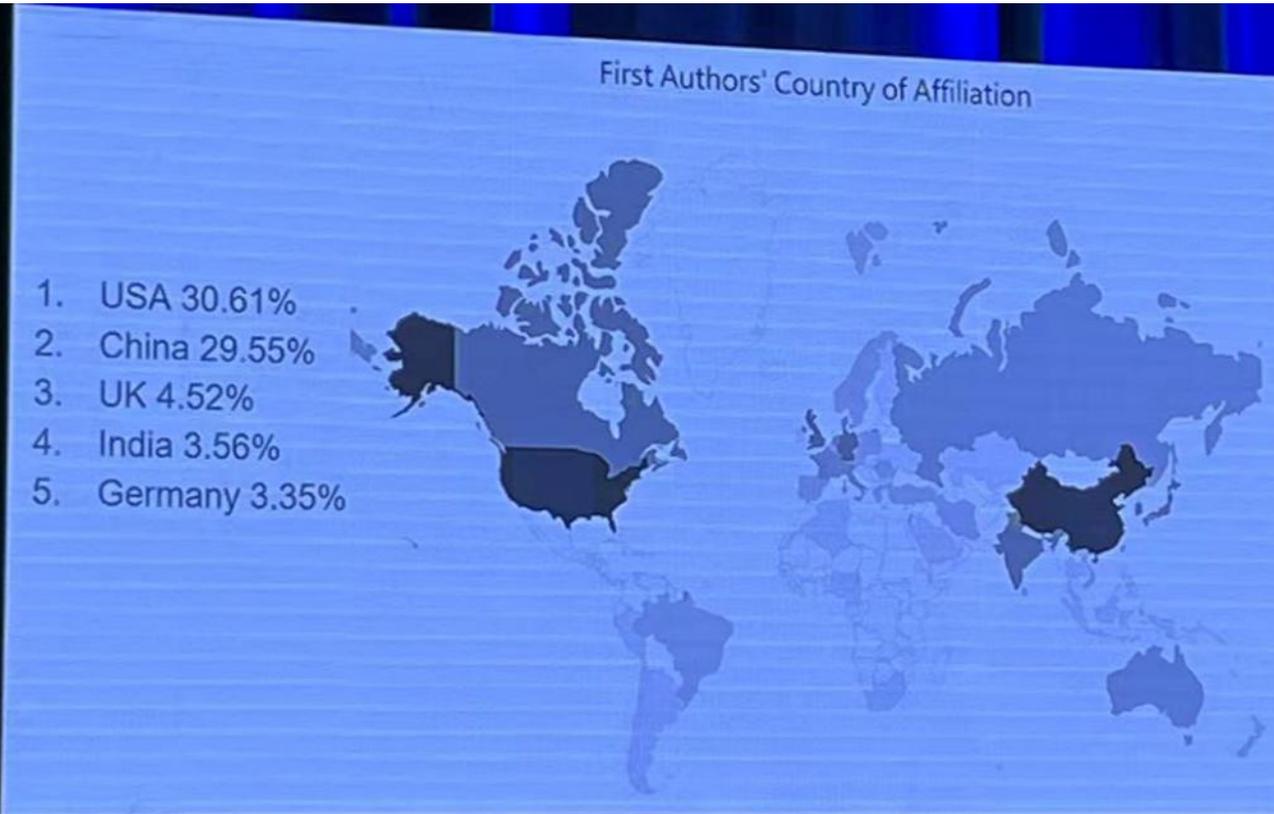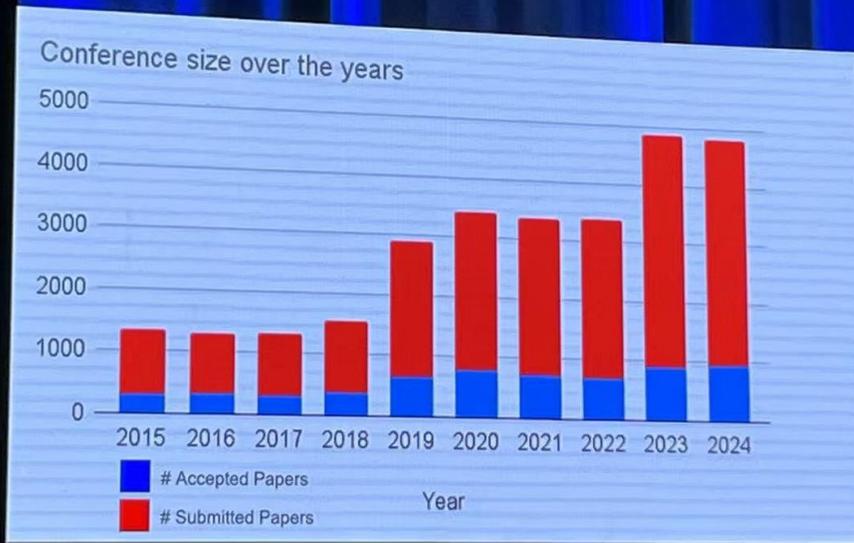# 曼谷整体印象

- 传统和现代的交融
  - 传统：很多佛教元素、保留王室、神秘的建筑
  - 现代：高楼大厦、便捷的购物体验...
- 天气湿热
- 饮食偏辛辣
- 一些可能吸引办会的点：国际化程度高、交通便捷、签证方便、英语普及度较高、很安全

# 议程

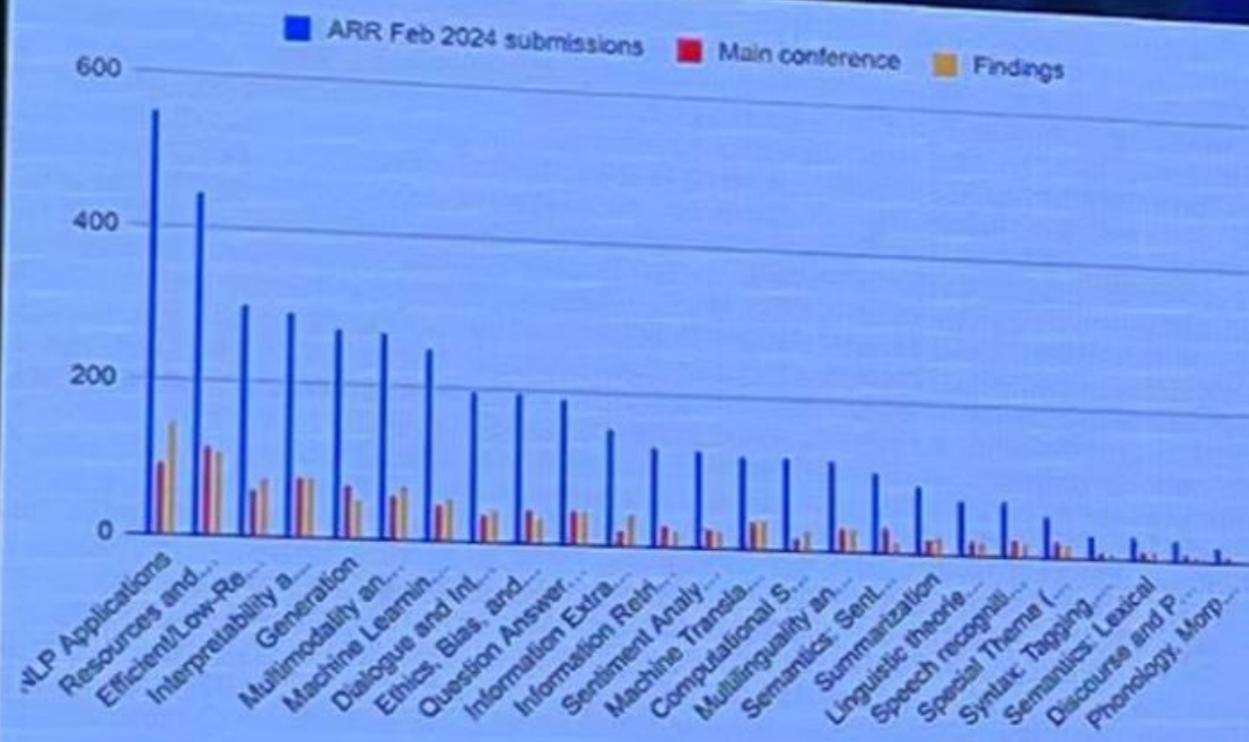| 日期 | 时间 | 活动内容 |
| --- | --- | --- |
| 8月10日 | 14:00-19:00 | 注册 |
| 8月11日 | 07:30-20:30 | 注册、教程、欢迎招待会 |
| 8月12日 | 07:30-17:30 | 注册、主题演讲、口头报告/海报/演示 |
| 8月13日 | 08:30-16:30 | 注册、主题演讲、圆桌讨论、社交晚宴 |
| 8月14日 | 08:30-16:30 | 注册、主题演讲、奖项与闭幕 |
| 8月15-16日 | 08:00-16:30 | 注册、研讨会 |

# ACL 2024 介绍

- 主会
  - 约5000篇投稿，72个SAC，716个AC，4208个审稿人
  - 940篇main（21.3%），975篇findings（22.1%），6 JCL，31 TACL
  - 3个主旨演讲，1个圆桌讨论
- 其他
  - 18个workshop，6个教程，38个demos和60篇SRW的论文

What's New This Year **NEW!**

- **All papers** presented as posters, some papers also as orals
- All findings papers can present (as posters)
- Virtual day **after** the conference (next week: Thursday, August 22)
- Non-publicized paper award
- Theme: Open science, open data, and open models for reproducible NLP

# 主旨演讲

- Does In-Context-Learning Offer the Best Tradeoff in Accuracy, Robustness, and Efficiency for Model Adaptation?

- Can LLMs Reason and Plan?

- Are LLMs Narrowing Our Horizon? Let's Embrace Variation in NLP!

# 主旨演讲

- Does In-Context-Learning Offer the Best Tradeoff in Accuracy, Robustness, and Efficiency for Model Adaptation?
- 主讲人：Sunita Sarawagi
  - Professor, IIT Bombay
  - sequence models for text and time-series, domain adaptation, effective human intervention in learning, graphical models and structured learning

# Background

- The model adaptation problem: from one domain (training set) to another (test set)

- Challenges:
  - Accuracy (overfit vs. no change)
  - Robustness
  - Efficiency

- Related topics: domain adaptation/transfer learning/few-shot learning

- Methods: Fine-tuning and its variants; Mixture of experts; Task-vectors; Matching-based methods

# Evaluation

| (small T) | Fine-tuning | MoE | Task vectors | Matching |
|---|---|---|---|---|
| Accuracy | ✓✓ | ✓✓✓ | ✓ | ✓ |
| Robustness | ✓ | ✓✓ | ✓✓✓ | ✓✓ |
| Efficiency — Adapt | ✓ | ✓ | Online ✓✓✓ | Online ✓✓✓ |
| Test | ✓✓✓ | ✓✓✓ | ✓✓✓ | ✓ |
| Others | | ✓ Difficult to train | | |

# LLM-era

- In-context learning (ICL): Adaptation is just a forward pass
- Why does it work?
  - H1: Transformers implement gradient descent algorithm over IC examples
  - H2: IC examples recognize tasks in pre-training, e.g. via **task vectors**
  - H3: Self-attention implements **matching-based** adaptation via induction heads

# Re-evaluation

| (small T) | Fine-tuning | | ICL |
|---|---|---|---|
| Accuracy | ✓✓ | | ✓✓ |
| Robustness | ✓ | | ✓ |
| Efficiency | | | Online |
| Adapt | ✓ | | ✓✓✓ |
| Test | ✓✓✓ | | ✓✓ |
| Others | | | ✓✓✓ |
| | | | Ease of use in multi-tenancy models |

Improvement:
- Structuring of prompts
- Pre-training strategies
- Fine-tuning
- Example retrieval
- Model architectures
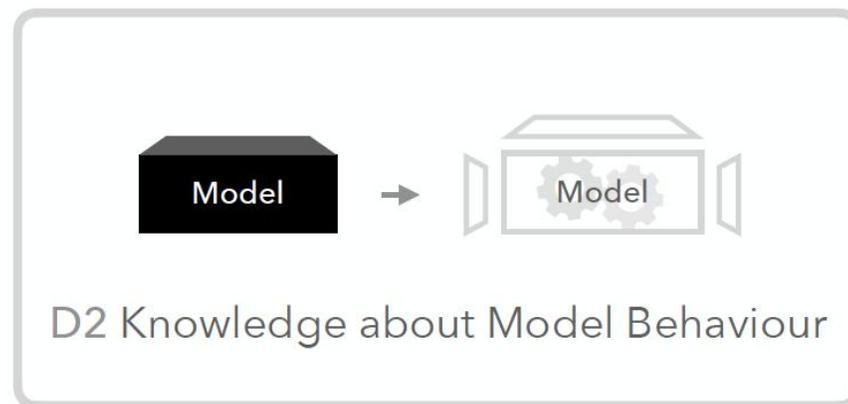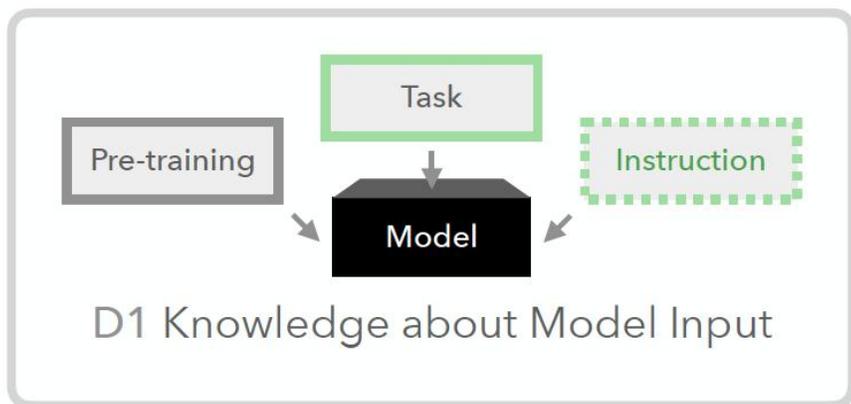
# Can LLMs Really Reason and Plan?

- Speaker: Subbarao Kambhampati
- Arizona State University
- 关注的问题：LLMs在推理过程中会出现幻觉等问题
- 作者对LLM常用的微调和提示语进行了很多测试
- 结论：LLMs的训练和使用方式并没有显示出它们能够进行原则性推理的迹象。

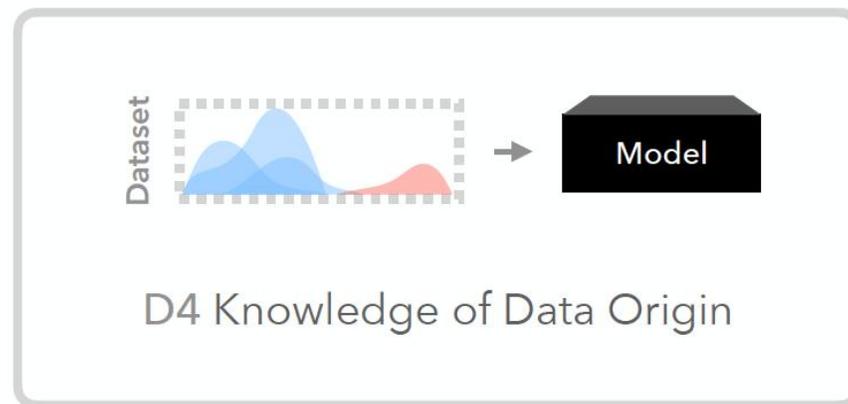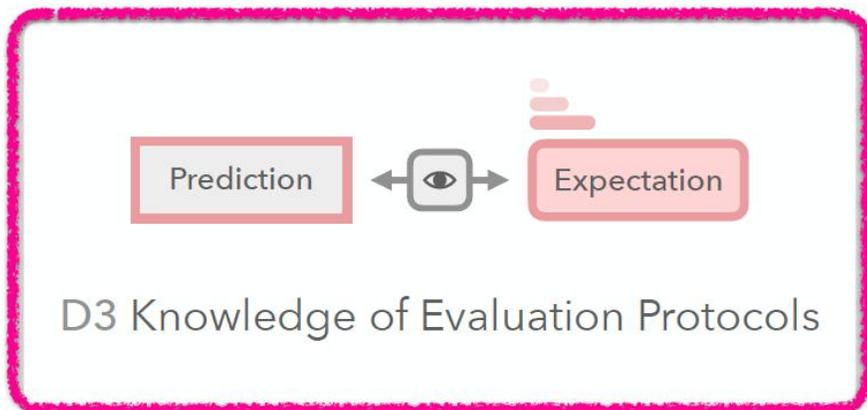# Are LLMs Narrowing Our Horizon? Let's Embrace Variation in NLP!

- Speaker: Barbara Plank LMU Munich & IT University of Copenhagen

- LLMs & Trust: A Short Look at AI history

- Regrain trust from four aspects:

D1 Knowledge about Model Input

D2 Knowledge about Model Behaviour

Trust arises from **knowledge of origin** as well as from **knowledge of functional capacity.**

*Trustworthiness - Working Definition by David G. Hays, 1979*

D3 Knowledge of Evaluation Protocols

D4 Knowledge of Data Origin

# We need to embrace variation holistically

- Inputs: linguistic variation, low-resource languages & dialects
- Outputs: human label variation as signal (not error) [Uncertainty]
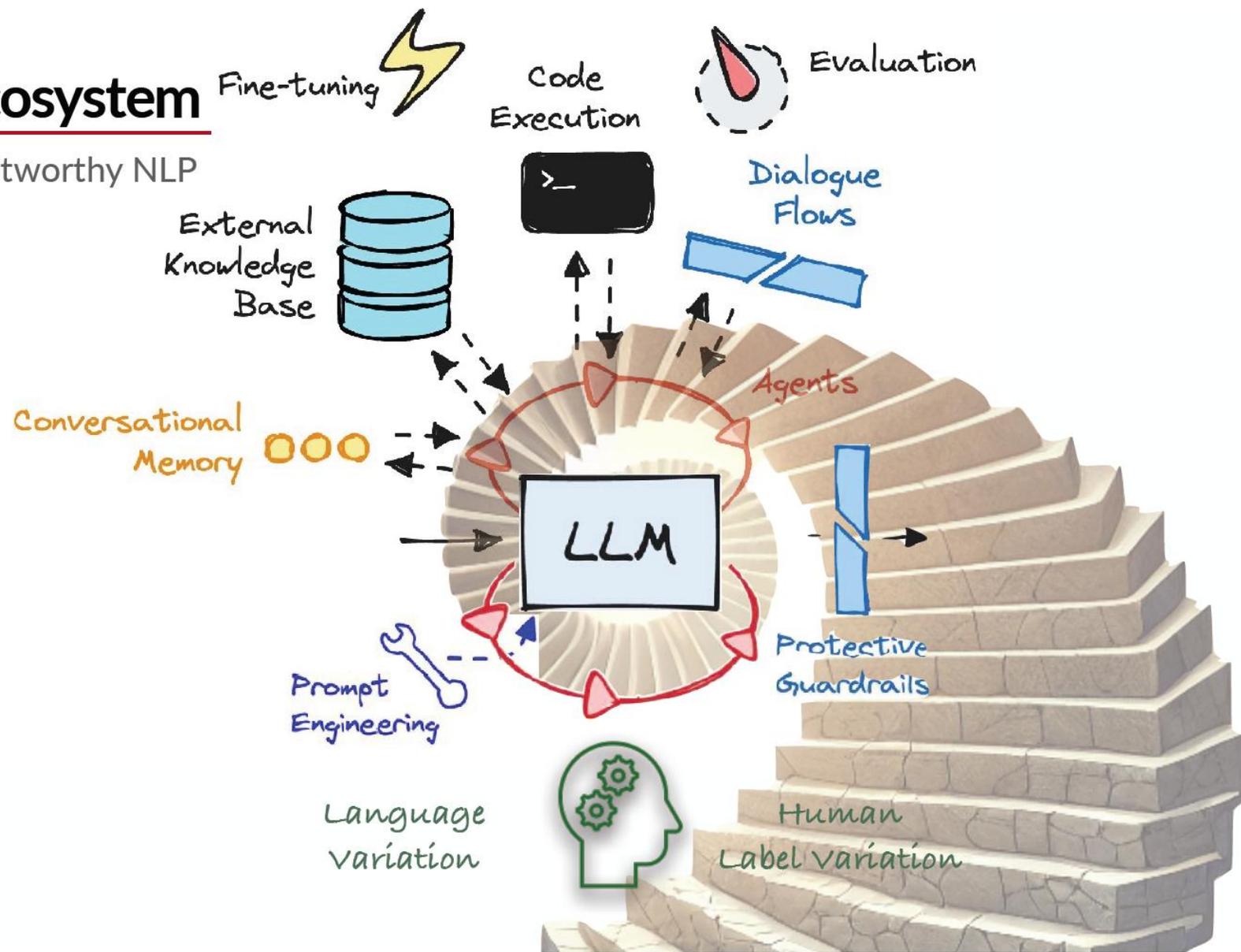- Research: Language as the Bridge

# Missing Ingredients:

Fine-tuning

Code Execution

Evaluation

External Knowledge Base

Dialogue Flows

Agents

## Social Sciences
e.g. Hovy & Yang, 2021 NAACL; Cornitzer et al., 2024.

Conversational Memory

## Survey Science
e.g. Eckman, Plank, Kreuter. ICML 2024.

LLM

## Sociolinguistics & Variation Linguistics
e.g. Grieve et al., 2024; Purschke et al., 2024 IClave 12

Prompt Engineering

Protective Guardrails

## HCI / label variation
e.g. Gordon et al., 2022 CHI, Plank 2022 EMNLP

## Interpretability & Actionable Insights
e.g. Marius Mosbach, NAACL 2024 workshop invited talk

69

# Trust LLM Ecosystem

Human-facing, Trustworthy NLP

Fine-tuning

Code Execution

Evaluation

External Knowledge Base

Dialogue Flows

Conversational Memory

Agents

LLM

Protective Guardrails

Prompt Engineering

Language Variation

Human Label Variation

# Tutorials

- **Tutorial 1– Computational Linguistics for Brain Encoding and Decoding: Principles, Practices and Beyond**

  Room : Lotus 1-4 (Level 22)

  *Jingyuan Sun, and Shaonan Wang, and Zijiao Chen, and Jixing Li, and Marie-Francine Moens*

- **Tutorial 2 - Automatic and Human-AI Interactive Text Generation (with a focus on Text Simplification and Revision)**

  Room : Lotus 5-7 (Level 22)

  *Yao Dou, and Philippe Laban, and Claire Gardent, and Wei Xu*

- **Tutorial 3 - Vulnerabilities of Large Language Models to Adversarial Attacks**

  Room : World Ballroom B (Level 23)

  *Yu Fu, Erfan Shayegan, and Md. Mamun Al Abdullah, and Pedram Zaree, and Nael Abu-Ghazaleh, and Yue Dong*

# Tutorials

**14:00 - 17:30**         Tutorial 4 - 6

- **Tutorial 4 – Computational Expressivity of Neural Language Models**

  Room: Lotus 1-4 (Level 22)

  *Alexandra Butoi and Ryan Cotterell and Anej Svete*

- **Tutorial 5– Watermarking for Large Language Model**

  Room: Lotus 5-7 (Level 22)

  *Xuandong Zhao, and Yu-Xiang Wang, and Lei Li*

- **Tutorial 6; Presentation Matters: How to Communicate Science in the NLP Venues and in the Wild?**

  Room: World Ballroom B (Level 23)

  *Sarvnaz Karimi, and Cecile Paris, and Gholamreza Haffari*

# T4: Representational Capacity of Neural Language Models

- A group of ETH Zürich led by Ryan Cotterell
- formal language theory to understand the representational capcacity of LLMs
- RNN vs. Finite-State Automata
- LLMs vs. finite-state automata and Turing machines.

# Semantics主题 (oral)

- - NounAtlas: Filling the Gap in Nominal Semantic Role Labeling

  Roberto Navigli, Marco Lo Pinto, Pasquale Silvestri, Dennis Rotondi, Simone Ciciliano, Alessandro Scirè

- - UG-schematic Annotation for Event Nominals: A Case Study in Mandarin Chinese

  Wenxi Li | Yutong Zhang | Guy Emerson | Weiwei Sun

  Computational Linguistics, Volume 50, Issue 2 - June 2023

- Distributional Inclusion Hypothesis and Quantifications: Probing for Hypernymy in Functional Distributional Semantics

  Chun Hei Lo, Wai Lam, Hong Cheng, Guy Emerson

# NounAtlas: Filling the Gap in Nominal Semantic Role Labeling



**Nominal SRL is under studied**

Existing SRL has been primarily focused on **verbal predicates**. But **nominal predicates are frequent in real-world settings!**

Dr. Jones's astounding discovery in his laboratory yesterday night
AGENT — nominal predicate — LOCATION — TIME

Newspaper Headlines     Dialogues

Short Messages     Social Media Posts

# Inventory

# Evaluation



Wordnet-based synset-to-frame mapping: Evaluation

1. **Unambiguous links**, evaluated on 100 items: **82%** in **agreement** with **manual annotations**, 18% made up of equally-valid selections
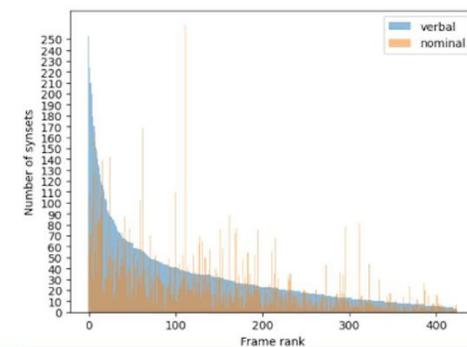
2. **Ambiguous links**, evaluated on 100 items: tie broken through **manual curation**, **Cohen's Kappa = 0.57 (moderate)**

3. **Missing links:** We use a Cross-Encoder trained on the other types of link to score predicates and, aggregated on ≤10 predicates, frames. We consider the top-5 frames and manually link.

- As a result of manual annotation, **99.6% of Unlinked synsets were added to a frame**
- **88.2% of correct frames** contained in our **model's top-5 ranking**, 49.3% in the top-1

# Corpus + model

## Creating a nominal SRL corpus

### Predicate nominalization: Generation

**Idea** → start from **SemCor** (Miller et al. 1993):
**Verbal predicates** in SemCor annotated with WordNet synsets **are directly linked** to VerbAtlas frames!

- We retain all SemCor sentences featuring verbal predicates
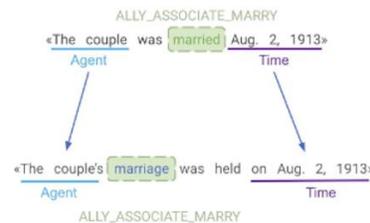- We **generate** the corresponding **nominalized version** of each sentence using **Gemini-Pro**

Change the sentence by nominalizing the verb "married" indicated by **. Use exactly one of these deverbal nouns: "marriage ceremony", "marriage", "wedding". Indicate the chosen deverbal noun with **:
"The couple was **married** Aug. 2, 1913 .

✦ The couple's **"marriage"** was held on Aug. 2, 1913.

- We **validate** the generation by keeping the sentence if:

1. The deverbal noun is **identified**
2. The deverbal noun is a **noun**
3. Its **lemma is among the candidates** provided

### Verbal-to-nominal role propagation

We annotate the verbal sentences with **InVeRo-XL** (Conia et al. 2021)

ALLY_ASSOCIATE_MARRY

«The couple was [married] Aug. 2, 1913»
        Agent                    Time

«The couple's [marriage] was held on Aug. 2, 1913»
        Agent                              Time

ALLY_ASSOCIATE_MARRY

We validate ~500 random sentences to generate a **gold test set for Nominal SRL**:

unchanged sentences: **92.11%**
unchanged frames: **95.56%**
unchanged role spans: **77.11%**
unchanged role labels: **95.27%**

Robust approach

## A unified model for SRL

Setup
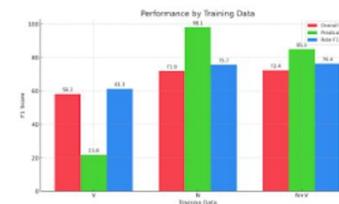**Backbone model**: RoBERTa-based model from Conia and Navigli (2020)
**Training data**:
- Nominal SRL: Our dataset
- Verbal SRL: OntoNotes 5.0 (V)
- The combination of the two datasets (N+V)
**Test data**:
- Our ~500 manually-curated nominal sentences (nouns)
- The OntoNotes test set (verbs) converted to VerbAtlas annotations

Gold-standard nominal test set

OntoNotes verbal test set

Our model, **jointly-trained on nominal and verbal data (N+V)**, achieves competitive performance on nominal SRL! 🏆
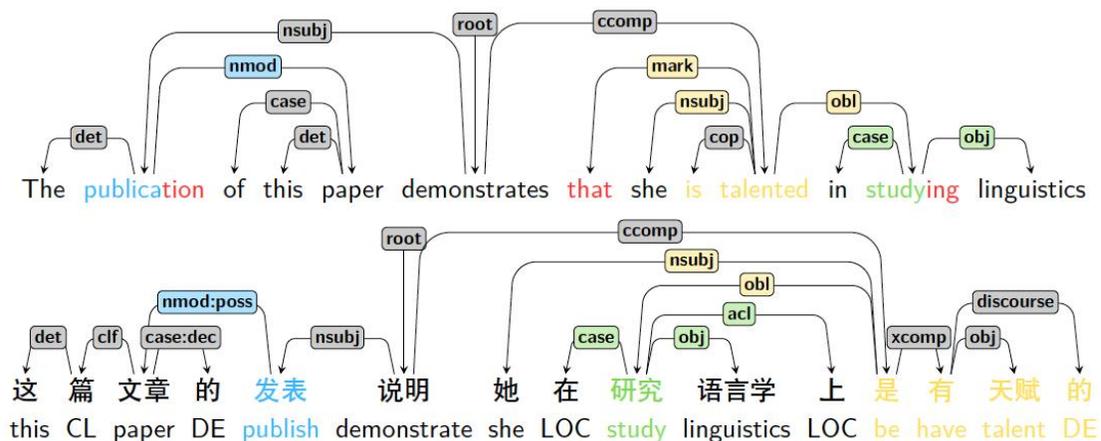
Also robust on the OntoNotes test set (verb frames).

Thanks to NounAtlas, we devised a unified approach for nominal and verbal SRL!

# Background



Multilingual Heterogeneity of Event Nominals

- **Identification**: Languages differ in their use of functional morphemes for nominalizations and the types of morphemes they use.
- **Nominal SRL**: Languages differ in how they realize participants of event nominals through syntactic constructions

# Research Question

## Multilingual Heterogeneity of Event Nominals

| | Mandarin | English | German |
|---|---|---|---|
| Subject with nominative case | − | ? | + |
| Accusative case on object | − | ? | + |
| Projection of outer Aspect | − | + | + |
| Modal or auxiliary verb | + | + | + |
| Complementizer | − | + | + |
| Verb suffix | − | + | + |
| Genitive/PP-subject | − | + | + |
| Genitive/PP-object | − | + | + |
| Gender | − | ? | + |
| Quantity | − | + | + |
| Determiner | − | + | + |
| Noun suffix | − | + | + |

Table: Verb- and noun-related features of event nominals

### Research Question

How can we accommodate multilingually heterogeneous phenomena in a unified way, or more theoretically, uncover universals of the world's languages beneath their surface-level variations?

# Distributional Inclusion Hypothesis and Quantifications: Probing for Hypernymy in Functional Distributional Semantics

## DIH and Quantifications

**DIH.** $r_2$ is a hypernym of $r_1$ iff $r_1$'s characteristic contexts $\subseteq r_2$'s.

**Quantifications.** A corpus with only universally quantified statements results in the reverse of DIH (rDIH).

$$animal \, \{\xleftarrow{\text{ARG1}} eats\}$$

$$dog \, \{\xleftarrow{\text{ARG1}} bark\} \quad bat \, \{\xleftarrow{\text{ARG1}} fly\}$$

Figure 1. A taxonomic hierarchy of nouns. Next to each noun are the contexts applicable to it and its descendants.

| Corpus 1 (DIH) | Corpus 2 (rDIH) |
|---|---|
| some dog barks | every dog barks |
| some animal barks | every dog eats |
| some bat flies | every bat flies |
| some animal flies | every bat eats |
| some animal eats | every animal eats |

Table 1. Corpora generated from the hierarchy in Fig. 1.

## FDS

**Entity Vectors.** $z \in \mathbb{R}^d$

**Truth-Conditional Semantic Functions.**

$$t^{(dog)}(z) = P\left(dog(z) = \top \mid z\right)$$
$$= \text{sigmoid}\left(v^{(dog)\top} z + b^{(dog)}\right)$$

**Representing Hypernymy.**

$$\forall z \text{ s.t. } \|z\|_2 \leq 1: t^{(dog)}(z) < t^{(animal)}(z)$$

which is true iff $s(dog, animal) > 0$, where

$$s(r_h, r_H) = b^{(r_H)} - b^{(r_h)} - \left\| v^{(r_H)} - v^{(r_h)} \right\|_2$$

**Model Training.** By Lo et al. (2023), given a DMRS graph:

$$some \xrightarrow{\text{RSTR}} dog \xleftarrow{\text{ARG1}} bark$$

Variational Inference: $q_\phi(z \mid \xleftarrow{\text{ARG1}} bark)$

Reconstruction: $\max \ln \mathbb{E}_{q_\phi(z \mid \xleftarrow{\text{ARG1}} bark)} \left[t^{(dog)}(z)\right] + \ldots$

**New objective (FDS$_\forall$).** Optimizes over regions of the entity s for handling universal quantifications

# Posters

# What does Kiki look like?

## Cross-modal associations between speech sounds and visual shapes in vision-and-language models

Tessa Verhoef, Kiana Shahrasbi and Tom Kouwenhoven

Leiden Institute of Advanced Computer Science, Leiden University, The Netherlands
Correspondence: t.verhoef@liacs.leidenuniv.nl, t.kouwenhoven@liacs.leidenuniv.nl

### INTRODUCTION

Humans have clear cross-modal preferences when matching novel words to visual shapes.

Example: bouba-kiki effect!

The development of multimodal models, including VLMs, can potentially revolutionise how machines understand and interact with humans.

Do VLMs associate non-words and visual stimuli in a human-like way?

### BACKGROUND

Non-arbitrariness as a **general property of language** [8]. Affects language **learning** and shapes language **emergence** [4,9,10].

Possible sources:
- Orthography [2]
- Acoustics and articulation [6,9,11]
- Affective–semantic properties of vocal communication [7]
- Physical properties relating to audiovisual regularities in the environment [3]

Alper and Averbuch-Elor[1] reported strong evidence for a bouba-kiki effect in CLIP and Stable Diffusion > surprising given model training and the absence of relevant data sources such as auditory information and experience with physical object properties.

Bouba or Kiki? From [5,9]

### METHODS

| Models | Architecture | Attention | #Params | #img,caps (M) |
|---|---|---|---|---|
| CLIP | Dual-stream | Modality-specific | 151.3M | 400, 400 |
| ViLT | Single-stream | Merged | 87.4M | 4.10, 9.85 |
| BLIP2 | Dual-stream | Q-Former | ~3.8B | 129, 258 |
| GPT-4o | Unknown | Unknown | Unknown | Unknown |

**Experiments:** Probed VLMs using a cognitive science paradigm to test for the bouba-kiki effect.

**Images:** Originals ⇑ from [5,9] as well as other image pairs from prior human experiments and entirely novel generated images, following [7].

4 pairs from [6]
4 pairs from [11]
8 pairs newly generated [7]
Examples ⇔

**Pseudowords:** from [7]

Sonorant consonants: m, n, l
Plosive consonants: t, k, p
Rounded vowels: oo, oh, ah
Non-rounded vowels: ee, uh, ay
Examples: moonch, lohmah, nahmoo
Examples: teekay, kuhpee, keepay

4 separate tests: Pseudoword generation (PWG) and probability score comparison (PSC), each with single syllables (1) or two-syllable words (2).

### RESULTS

**Fig 1:** Curved or Jagged images matched to Sonorant consonants (left bars) or Rounded vowels (right) in PWG1 test. The expected pattern would show higher bars for the Curved than for the Jagged shapes in both sets. Only true for CLIP and GPT-4o.

**Fig 2:** Probability scores in PSC1 for two pairs of original pseudowords (bouba & kiki, takete & maluma) + four generated syllable types: Sonorant-Rounded (S-R, expected to be most Curved), Sonorant-Non-Rounded (S-NR), Plosive-Rounded (P-R) and Plosive-Non-Rounded (P-NR, expected to be most Jagged), paired with Jagged or Curved shapes. **No effect found.**

| Model | PWG1 | PSC1 | PWG2 | PSC2 |
|---|---|---|---|---|
| CLIP | | | | |
| GPT-4o | | | | |
| ViLT | | | | |
| BLIP2 | | | | |

**Fig 3:** Evidence for a bouba-kiki effect found in only three tests (green), for CLIP and GPT-4o. ViLT and BLIP2 never passed.

### DISCUSSION

It is too early to conclude that VLMs understand sound-symbolism or map visio-linguistic representations in a human-like way. Results depend heavily on which specific model is tested and how the task is formulated.

However, results tentatively suggest that cross-modal preferences can, to some extent, be learned from statistical regularities in data.

Model features such as architecture design, training objective, number of parameters, and input data seem to affect the results.

These findings inform discussions on the origins of the bouba-kiki effect in human cognition and future developments of VLMs that align well with human cross-modal associations.

### REFERENCES

---

# Aligning to Adults Is Easy, Aligning to Children Is Hard: A Study of Linguistic Alignment in Dialogue Systems

Dorothea French, Sidney D'Mello, Katharina von der Wense
University of Colorado Boulder

### Introduction

Over the course of a conversation people start to align with each other. They use similar words, syntax, and talk about similar ideas. This phenomenon, called linguistic alignment, helps conversational participants understand each other and reduces the effort needed to do so, amongst other positive effects.

Linguistic alignment is particularly important when one of the participants is a child, or a non-fluent speaker of the language. Levels of parent's alignment have been shown to correspond to a child's vocabulary and language development. Repetition and correction also help language learners improve naturally.

In this era, characterized by widespread use of dialogue system, it is important to ask how well they can communicate with humans. Alignment of dialogue systems can reduce frustration and perceived task load. It is also a critical component of communication. Thus we must ask, how well do our systems do? Do they reach human-like levels, or is this an area of friction?

In this work we ask two questions:

To what degree do two state-of-the-art dialogue systems - ChatGPT and Llama2 - align to users, and how does their alignment compare to that typically seen between humans?

To what degree do ChatGPT and Llama2 align to children (i.e. non-fluent speakers), and how does this level of alignment compare to a parent's?

### Approach and Data

We use the Switchboard Dialogue Acts Corpus (SWDA) of a series of phone conversations on a variety of topics to assess model alignment with adults.

We use the CHILDES dataset of conversations between caretakers and children to assess the models' alignment to children, i.e. language learners.

**SWDA Context**

| | |
|---|---|
| A | Any jury's not going to disregard the evidence, you know. |
| B | Uh, that's true |
| B | I, I, I think our judicial system is attorney welfare myself. |
| A | That may very well be. |

Response

| | |
|---|---|
| B | I, I hold it in the utmost contempt. |
| GPT | It's definitely a possibility that needs to be looked into. |
| Llama2 | Yeah, it's like, you know, they're just trying |

Table 1: Last lines of SWDA transcript plus true and generated responses

**CHILDES Context**

| | |
|---|---|
| MOT | Hm? |
| CHI | Where Mommy go? |
| MOT | Mommy went to the university this morning to get some books. |
| CHI | Where's Mommy's books? |

Response

| | |
|---|---|
| MOT | They're in the hallway in a big bag. |
| GPT | Mommy will bring the books home this evening |
| Llama2 | Mommy left her books in the car. |

Table 2: Last lines of CHILDES transcript plus true and generated responses

From the datasets we select transcript excerpts of 36 turns - 35 turns of context, and to following utterance.

Our two models, GPT-3.5 and Llama2, are prompted with and respond to the 35 turns of context as a participant or parent respectively

We look at three types of alignment (syntactic, lexical, and semantic) between the last turn of the context and the response and compare these metrics between the models and gold responses

### Results

**Table 3: Average alignment for responses to SWDA**

| Response | Syntactic | Lexical | Semantic |
|---|---|---|---|
| True | 0.444 | 0.170 | 0.308 |
| ChatGPT | 0.443 | 0.151 | 0.340 |
| Llama2 13b | 0.472 | 0.207 | 0.350 |
| Llama2 7b | 0.475 | 0.213 | 0.374 |

**Table 4: Average alignment for responses to CHILDES**

| Response | Syntactic | Lexical | Semantic |
|---|---|---|---|
| True | 0.490 | 0.278 | 0.411 |
| ChatGPT | 0.436 | 0.190 | 0.347 |
| Llama2 13b | 0.464 | 0.227 | 0.345 |
| Llama2 7b | 0.473 | 0.251 | 0.370 |

Looking at the SWDA responses, we see ChatGPT aligns at human-like levels. Whereas Llama2 exhibits elevated levels of alignment. Upon manual inspection, we see differences in the responses. ChatGPT responses are generally better, scoring an average of 4.37 out of 5 for quality, compared to Llama2's 3.5 out of 5. The ChatGPT responses are more likely to contain novel or information that drives the conversation forward - confirmed by a lower semantic and lexical alignment.

Llama2 is more repetitive, and it does a very good job mimicking the stylistic elements of the conversation, but often makes less sense.

When responding to children, both models have lower than human levels of alignment. Like with adults, ChatGPT has slightly better responses, and a much lower proportion of very poor (1-2 out of 5) responses. Once again, Llama2 has more repetitive responses, but higher alignment. In both response sets we see the smaller Llama2 model consistently has higher alignment.

### Conclusion

Dialogue systems show great potential to assist humans across a variety of tasks. The success of these interactions correlates with linguistic alignment. We find that, when responding to adult speakers, ChatGPT shows approximately human-level alignments and provides constructive responses. Llama2, however, overly mimics the conversation. This could be positive when talking with children or language learners as it results in elevated alignment. We conclude that SOTA dialogue systems can improve in regards to tailoring levels of alignment to match various circumstances, without reducing quality.

In the future, we plan to investigate alignment to adult learners, non-typical speakers, and explore methods to create dialogue systems with a closer to human levels of alignment.

# Social Events

# 颁奖环节

- The 2024 winner of the 1999 Test-of-Time Paper Award is:
- Lillian Lee. 1999. Measures of Distributional Similarity.
- In Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics (ACL), pages 25–32, College Park, Maryland, USA.

- The 2024 winner of the 2014 Test-of-Time Paper Award is:
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. GloVe: Global Vectors for Word Representation.
- In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 1532–1543, Doha, Qatar.

# Best Paper Award

- Mission: Impossible Language Models
- Why are Sensitive Functions Hard for Transformers?
- Deciphering Oracle Bone Language with Diffusion Models
- Causal Estimation of Memorisation Profiles
- Aya Model: An Instruction Finetuned Open-Access Multilingual Language Model
- Semisupervised Neural Proto-Language Reconstruction
- Natural Language Satisfiability: Exploring the Problem Distribution and Evaluating Transformer-based Language Models

# One more thing... from EMB

**ACL is Not an AI Conference**

- I am not using "AI" as a synonym for "ML".
- Machine learning (including deep learning) provides many techniques that are useful for language technology and computational linguistics
  - (Though both of those terms are problematic.)
- The problems of CL/NLP can also be illuminating for questions about ML

- The issues that I am concerned with arise when the focus shifts to "AI"

**AI as a research & commercial field**

Asks questions like:

- How do we build "thinking machines" that can do "human-like" reasoning?
- How do we build "thinking machines" that can "surpass" humans in cognitive work?
  - (and cure cancer, solve the climate crisis, make end-of-life decisions, etc)
- How do we automate the scientific method?
- How do we automate away such creative work as painting and writing?
  - Or: How do we steal artwork at scale and try to convince people this is "for the common good"?

# Compling/NLP asks questions such as

- How are languages similar/different?
- How is information represented in languages?
- How can we build technology that assists with: transcription, translation, summarization, information access … in different languages?
- How can we evaluate such technology?
- What kinds of intermediate representations are useful for such technology?
- How well do different ML techniques work for different tasks?
- How do language technologies interact with existing systems of power and oppression?

# Workshops

- Rep4NLP (Probing -> causal inference)
  - - Latent Space Exploration for Safe and Trustworthy AI by Hassan Sajjad
  - The gap between what language models say and what they know
  - Generalisation in LLMs – and beyond
  - Efficiency as an Inductive Bias: Towards Tokenizer-free and Dynamically Sparse Language Models

- Heng Ji: AI Plays Medicinal Chemist

- Anna Rogers: Ai for reasearch workflow

# Q & A

THANK YOU

# Note