

# Vaishnavi Shrivastava

---

CONTACT	Email: <a href="mailto:vaish.shrivastava@stanford.edu">vaish.shrivastava@stanford.edu</a>	Homepage: <a href="https://vshrivas.github.io/">https://vshrivas.github.io/</a>
KEYWORDS	LLM calibration, multi-hop reasoning, retrieval augmentation	
EDUCATION	<b>Stanford University</b> Master of Science, Computer Science Advisor: <a href="#">Percy Liang</a>	2022 - 2024 (projected)
	<b>California Institute of Technology (Caltech)</b> Bachelor of Science, Computer Science	2015 - 2019 <b>3.9/4.0</b>
PUBLICATIONS	<p>[1] <u>Llamas Know What GPTs Don't Show: Surrogate Models for Confidence Estimation.</u> <b>V. Shrivastava</b>, P. Liang, A. Kumar. 2023. <i>In submission</i></p> <p>[2] <u>Benchmarking and Improving Generator-Validator Consistency of Language Models.</u> X. Lisa Li, <b>V. Shrivastava</b>, S. Li, T. Hashimoto, P. Liang. 2023. <i>In submission</i> <a href="#">[arxiv]</a></p> <p>[3] <u>Bias Runs Deep: Implicit Reasoning Biases in Persona-Assigned LLMs.</u> S. Gupta, <b>V. Shrivastava</b>, A. Deshpande, A. Kalyan, P. Clark, A. Saharwal, T. Khot. 2023. <i>In submission</i></p> <p>[4] <u>UserIdentifier: Implicit User Representations for Simple and Effective Personalized Sentiment Analysis.</u> F. Mireshghallah, <b>V. Shrivastava</b>, M. Shokouhi, T. Berg-Kirkpatrick, R. Sim, D. Dimitriadis. 2021. <i>North American Chapter of the Association for Computational Linguistics (NAACL) 2022</i> <a href="#">[arxiv]</a></p> <p>[5] <u>Exploring Low-Cost Transformer Model Compression for Large-Scale Commercial Reply Suggestions.</u> <b>V. Shrivastava*</b>, R. Gaonkar*, S. Gupta*, A. Jha. 2021. <i>Preprint</i> <a href="#">[arxiv]</a></p>	
RESEARCH EXPERIENCE	<b>Research Assistant:</b> <ul style="list-style-type: none"><li><b>Stanford University:</b> Advised by Percy Liang <span style="float: right;"><i>(Sep'22 - Current)</i></span> <i>Themes: LLMs, Calibration, Reasoning</i></li><li><b>Allen Institute for AI:</b> Advised by Tushar Khot, Peter Clark <span style="float: right;"><i>(Jun'23 - Current)</i></span> <i>Themes: Reasoning, Persona-guided LLMs, Calibration</i></li></ul>	
WORK EXPERIENCE	<b>Applied Scientist:</b> <ul style="list-style-type: none"><li><b>Microsoft AI:</b> Suggested Replies &amp; Summarization <span style="float: right;"><i>(Sep'19 - Aug'22)</i></span> <i>Themes: Dialog Systems, Model Compression, Personalization, Summarization</i></li></ul> <b>Software Engineering Intern:</b> <ul style="list-style-type: none"><li><b>Microsoft AI:</b> Knowledge Mining and Graphs Group <span style="float: right;"><i>(Jul'18 - Sep'18)</i></span> <i>Themes: Key-Phrase Extraction, Part-of-Speech Tagging, Email Search</i></li><li><b>Microsoft:</b> Substrate Data Store Group <span style="float: right;"><i>(Jun'17 - Sep'17)</i></span> <i>Themes: Multi-threading, Backend, Thread-Safe Caching</i></li><li><b>Dell-EMC:</b> <span style="float: right;"><i>(Jun'16 - Sep'16)</i></span> <i>Themes: Distributed Computing Algorithms, Concurrent Services</i></li></ul>	
TEACHING EXPERIENCE	<b>Teaching Assistant:</b> <ul style="list-style-type: none"><li><b>Caltech:</b> Machine Learning &amp; Data Mining, CS 155 <span style="float: right;"><i>(Jan'19 - Mar'19)</i></span></li><li><b>Caltech:</b> Database System Implementation, CS 122 <span style="float: right;"><i>(Jan'18 - Mar'18)</i></span></li></ul>	

SELECTED  
RESEARCH  
PROJECTS

**Surrogate Models for Confidence Estimation**

(Jul'23 - Sep'23)

*Advisor: Percy Liang, Ananya Kumar - Stanford University*

- SoTA models like GPT-4 and Claude don't provide access to their probabilities making it difficult to assess their confidences in their outputs. Prompting for confidences doesn't work well.
- We introduce surrogate model calibration - using a white-box surrogate like Llama-2 to approximate the internal confidences of a black-box model like GPT-4.
- Mixing surrogate probabilities and prompted confidences leads to further gains.

**Implicit Reasoning Biases in Persona-Assigned LLMs**

(Jun'23 - Sep'23)

*Advisor: Tushar Khot, Ashish Sabarwal - Allen Institute for AI*

- LLMs have deep-rooted biases which can be surfaced through personas.
- Performance on 24 reasoning tasks shows that models assigned personas of certain demographic groups may abstain or make more implicit reasoning errors, conforming with social stereotypes.

**Improving Generator-Validator Consistency in LLMs**

(Apr'23 - Jun'23)

*Advisor: Percy Liang, Lisa Li - Stanford University*

[\[arxiv\]](#)

- LLM are inconsistent with generative (What is 7+8?) vs validation queries (7+8=15, True/False?).
- We propose a fine-tuning scheme to improve generator-validator consistency and show accuracy improvements and performance transfer between generators and validators.

**Belief Aggregation for Factually Correct Reasoning**

(Sep'22 - Mar'23)

*Advisor: Percy Liang - Stanford University*

- We propose sampling chains-of-thought and extracting LLM's 'beliefs' from those chains.
- Beliefs can be composed and used to verify LLM's world model for factually correct reasoning.

**Implicit Personalized User Representations**

(Jul - Dec'21)

*Microsoft Research*

[\[arxiv\]](#)

- We investigate using uniformly distributed, non-trainable, user-specific prompts for user-personalization, instead of trainable embeddings, to circumvent periodically training embeddings per user.
- We demonstrate that we can outperform SOTA prefix-tuning based results on a suite of sentiment analysis by up to 13%, resulting in a paper.

**Low-Cost Transformer Model Compression**

(Jul - Nov'20)

*Microsoft Search, Assistant and Intelligence*

[\[arxiv\]](#)

- We experiment with low-cost methods to compress Transformer bi-encoder based reply suggestion system, reducing training and inference times by 42% and 35% respectively.
- Investigate how dataset size, pre-trained model use, and domain adaptation of the pre-trained model affected the performance of compression techniques.
- We discover that large-data settings allow low-cost techniques to be very effective in compressing pre-trained model based architectures. Insights led to a paper and a talk.

TALKS

*"Supercharging Reply Suggestions: Model Compression Solutions and Insights from a Real-World Setting". Microsoft Machine Learning, AI and Data Science Conference (MLADS) 2021*

SELECTED  
LEADERSHIP  
POSITIONS

- Corporate Vice President, *Caltech IEEE*
- Treasurer, *Caltech Society of Women Engineers*
- Secretary, *Caltech Robogals*

REFERENCES

**Percy Liang**, Associate Professor, Stanford University  
**Milad Shokouhi**, Partner Applied Scientist, Microsoft  
**Dan Schwartz**, Principal Applied Scientist, Microsoft  
**Donnie Pinkston**, Lecturer, Caltech