

# Vaishnavi Shrivastava

---

BASIC INFORMATION	Email: <a href="mailto:vashri@microsoft.com">vashri@microsoft.com</a> Homepage: <a href="https://vshrivas.github.io/">https://vshrivas.github.io/</a>	Pronouns: <i>she/her/hers</i> Phone Number: (+1) 408-477-5322
EDUCATION	<b>California Institute of Technology (Caltech)</b> Bachelor of Science, Computer Science	<b>Sep'15 – Jun'19</b> <b>3.9/4.0</b>
RESEARCH INTERESTS	<b>Natural Language Processing:</b> Transfer Learning & Language Models, Question Answering, Commonsense Reasoning, Abstractive Summarization, Dialog Systems, Grounded Language Learning <b>Machine Learning:</b> Few-shot Learning, Federated Learning, Deep Reinforcement Learning, Model Interpretability, Multi-modal Learning	
TECHNICAL SKILLS	<b>Languages:</b> <i>Proficient:</i> Python, Java, C, C++   <i>Basic:</i> C#, SQL <b>Toolkits:</b> PyTorch, Keras, Tensorflow	
PUBLICATIONS	[1] (Preprint) <b>V. Shrivastava*</b> , R. Gaonkar*, S. Gupta*, A. Jha. 2021. Exploring Low-Cost Transformer Model Compression for Large-Scale Commercial Reply Suggestions. <a href="#">arXiv: 2111.13999</a>  [2] (Under review) F. Miresghallah, <b>V. Shrivastava</b> , M. Shokouhi, T. Berg-Kirkpatrick, R. Sim, D. Dimitriadis. 2021. UserIdentifier: Implicit User Representations for Simple and Effective Personalized Sentiment Analysis. <a href="#">arXiv: 2110.00135</a>	
WORK AND RESEARCH EXPERIENCE	<b>Applied Scientist:</b> <ul style="list-style-type: none"><li>● <b>Microsoft AI:</b> Suggested Replies &amp; Summarization (Sep'19 - Present) <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 (Jul'18 - Sep'18) <i>Themes: Key-Phrase Extraction, Part-of-Speech Tagging, Email Search</i></li><li>● <b>Microsoft:</b> Substrate Data Store Group (Jun'17 - Sep'17) <i>Themes: Multi-threading, Backend, Thread-Safe Caching</i></li><li>● <b>Dell-EMC:</b> (Jun'16 - Sep'16) <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 (Jan'19 - Mar'19)</li><li>● <b>Caltech:</b> Database System Implementation, CS 122 (Jan'18 - Mar'18)</li></ul>	
RECENT PROJECTS	<b>Personalized Language Models</b> (Jul'21 - Present) <ul style="list-style-type: none"><li>– Aim is building user-level personalized generative reply suggestion dialog systems with GPT-2.</li><li>– Developed a modified <i>Prefix-Tuning</i> based approach to learn user-embeddings to condition GPT-2 model for personalization, improving validation perplexity by 9% over vanilla prefix-tuning.</li><li>– Jointly trained a network with GPT-2 to generate embeddings as a function of user n-gram language signals, to solve the cold-start problem of personalizing responses for unseen users.</li><li>– Using <i>LoRA: Low-Rank Adaptation of Large Language Models</i> technique for more fine-grained personalization by directly personalizing weight updates to GPT-2's attention matrices.</li></ul> <b>Implicit Personalized User Representations</b> (Jul - Sep'21) <a href="#">Paper</a> <ul style="list-style-type: none"><li>– Investigated using uniformly distributed, non-trainable, user-specific prompts for user-personalization, instead of trainable embeddings, to circumvent periodically training embeddings per user.</li></ul>	

- Demonstrated that we can outperform SOTA prefix-tuning based results on a suite of sentiment analysis by up to 13%, resulting in a paper.

### Federating Adapters

(Jul - Aug'21)

- To reduce communication overhead for large language models (LMs) during federated learning, proposed inserting bottleneck adapter layers and sharing client-server updates only on those layers.
- Improved communication costs by 121x on sentiment analysis, without significant accuracy drops.
- Proposed a user clustering mechanism to leverage *AdapterFusion* and further improve accuracy.

### Factual Consistency for Abstractive Summarization

(Mar - Jun'21)

- Developed an automated metric for evaluating factual consistency of summaries by few-shot tuning GPT-3 for question generation (QG) and question answering (QA).
- Generated questions on the summary using QG model, and answers to those questions first based on the source and then based on the summary using the QA model.
- Evaluated answer similarity between source and summary using an F1 score.

### Multi-turn Conversation Modeling

(Nov'20 - Feb'21)

- Modeled multi-turn conversations for contextualized response suggestions in dialog systems.
- Implemented shared-weight Hierarchical Transformers to encode prior utterances separately and aggregate them using a self-attention layer to form contextualized input representations.
- Saw substantial gains in offline metrics compared to previous single-turn model and baseline concatenating previous utterances as new input.

### Low-Cost Transformer Model Compression

(Jul - Nov'20)

[Paper](#)

- Experimented with low-cost methods to compress Transformer bi-encoder based reply suggestion system, reducing training and inference times by 42% and 35% respectively.
- Investigated how dataset size, pre-trained model use, and domain adaptation of the pre-trained model affected the performance of compression techniques.
- Discovered 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.

### Dialog System Triggering

(Feb - Jun'20)

- Trained a light-weight biLSTM classifier to decide which messages to trigger core reply suggestion system on, to prevent suggesting irrelevant responses to open-ended questions.
- Shipping the model led to reductions in latency and improved user engagement metrics, since fewer suggestions were shown to users for open-ended questions that the model struggled with.

SELECTED  
PREVIOUS  
PROJECTS

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

**Milad Shokouhi**, *Partner Applied Scientist, Microsoft*  
**Dan Schwartz**, *Principal Applied Scientist, Microsoft*  
**Abhishek Jha**, *ML Engineering Manager, Stripe*  
**Donnie Pinkston**, *Lecturer, Caltech*