

# Fei Fang

## EDUCATION

<b>Stanford University</b> <i>Master of Science in Computer Science (Specialization in Artificial Intelligence)</i>	<i>Jan 2022</i> Stanford, CA
<b>Stanford University</b> <i>Bachelor of Science in Mathematics, with Distinction</i>	<i>Jan 2022</i> Stanford, CA

## EMPLOYMENT

<b>Cresta Intelligence</b> — Conversational AI for contact centers <i>Machine Learning Engineer</i>	<i>May 2024 - Present</i> Palo Alto, CA
<ul style="list-style-type: none"><li>• Finetuned 7B-parameter large language model (LLM) using knowledge distillation for a domain-specific question answering (QA) system backed by retrieval-augmented generation (RAG), outperforming GPT-4o in internal benchmarks.</li><li>• Deployed the optimized LLM to enhance real-time customer support efficiency at major enterprises such as United Airlines and Square.</li></ul>	
<b>Glean Technologies</b> — Enterprise-grade AI-driven search engine and assistant <i>Machine Learning Engineer</i>	<i>Jan 2022 - April 2024</i> Palo Alto, CA
<ul style="list-style-type: none"><li>• Led the development of a multilingual RAG-based LLM assistant for Japanese and English, deployed to 10+ major Japanese corporations including Toyota.</li><li>• Productionized the ColBERT model for reranking relevant contexts in RAG, demonstrating a 7-point increase in Mean Average Precision (MAP).</li><li>• Developed a novel algorithm inspired by ColBERT which efficiently probes a finetuned semantic bi-encoder for word-level measures of relevance at query time, making the retrieval system more interpretable and easier to debug.</li><li>• Defined quantitative metrics to measure the impact of semantic search on the hybrid retrieval system in production traffic, enabling A/B testing for all modeling experiments.</li><li>• Optimized serving efficiency across the vector search stack, resulting in a 14% reduction in median latency and an 11% reduction in 90th percentile latency.</li></ul>	

## PUBLICATIONS

- [1] E. Kreiss, **F. Fang**, N. D. Goodman, and C. Potts. "Concadia: Towards Image-Based Text Generation with a Purpose". In: *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. 2022, pp. 4667–4684.
- [2] **F. Fang**, K. Sinha, N. D. Goodman, C. Potts, and E. Kreiss. "Color Overmodification Emerges from Data-Driven Learning and Pragmatic Reasoning". In: *Proceedings of the Annual Meeting of the Cognitive Science Society*. Vol. 44. 2022.

## TEACHING

<b>CS221: Artificial Intelligence</b> , Stanford University <i>Co-Head Course Assistant, Student Liaison</i>	<i>Fall 2021</i>
<b>CS103: Mathematical Foundations of Computing</b> , Stanford University <i>Co-Lecturer</i>	<i>Summer 2021</i>
<b>CS103: Mathematical Foundations of Computing</b> , Stanford University <i>Tutorial Leader</i>	<i>Spring 2021</i>
<b>CS221: Artificial Intelligence</b> , Stanford University <i>Head Course Assistant</i>	<i>Winter 2021</i>
<b>CS103: Mathematical Foundations of Computing</b> , Stanford University <i>Course Assistant</i>	<i>Spring &amp; Fall 2019; Winter &amp; Spring 2020</i>

## SERVICE

<b>Lemontree (<a href="https://foodhelpline.org">foodhelpline.org</a>)</b> — SMS-based helpline connecting people to local food banks <i>AI Strategy Advisor</i>	<i>Oct 2024 - Present</i>
---	---------------------------

## AWARDS, GRANTS, FELLOWSHIPS

---

### Computer Science Teaching Fellow

*Appointed as summer session instructor for a core requirement course of the Computer Science department*

*Jan 2021*

Stanford University

### Kung-Yi Kao Prize for Outstanding Progress in the Study of Japanese

*Awarded to an outstanding undergraduate in an East Asian Language Study*

*Jun 2019*

Stanford University

### East Asia Undergraduate Summer Language Study Grant

*Received funding to undertake intensive Japanese study during the summer session*

*Jun 2017*

Stanford University

## TALKS AND POSTERS

---

### Wikimedia Research Showcase, virtual

*Talk on "Concadia: Towards Image-Based Text Generation with a Purpose"*

E. Kreiss (presenting), F. Fang, N. D. Goodman, C. Potts.

*Apr 2024*

### EMNLP 2022, virtual

*Poster on "Concadia: Towards Image-Based Text Generation with a Purpose"*

E. Kreiss (presenting), F. Fang, N. D. Goodman, C. Potts.

*Dec 2022*

### CogSci 2022, virtual

*Talk on "Color Overmodification Emerges from Data-Driven Learning and Pragmatic Reasoning"*

F. Fang & E. Kreiss (presenting jointly), K. Sinha, N. D. Goodman, C. Potts.

*Jul 2022*

## LANGUAGES

---

*Listed in chronological order of acquisition: **Mandarin** (native), **Cantonese** (native), **English** (native proficiency), **Spanish** (professional proficiency), **French** (elementary proficiency), **Japanese** (professional proficiency)*