Sunnyvale, CA
yash−s20.github.io

# YASH SHARMA

(845) 290-4694
yash.sharma200999@gmail.com
github.com/yash−s20

## EDUCATION

**Cornell University** *Ithaca, NY* **Aug 2022 – May 2024**
MS in **Computer Science** | Minor in **Cognitive Science** GPA: 3.91 / 4
Computational Sustainability, Advanced Language Technologies, Advanced Programming Languages

**Indian Institute of Technology Bombay** *Mumbai, India* **Aug 2017 – May 2021**
B.Tech in **Computer Science & Engineering** (Honors) | Minor in **Artificial Intelligence** GPA: 9.68 / 10
Deep Learning for NLP, Advanced Machine Learning, Analysis of Concurrent Programs

## SOFTWARE SKILLS

**Systems & Programming** | python, C/C++, bash, Rust, Haskell, Java, Javascript, SQL, AVX, Git, Perforce, Docker, KVM
**Machine Learning** | PyTorch, TensorFlow, TensorRT & onnxruntime

## WORK EXPERIENCE

**ML Research Engineer, Matic Robots** *Mountain View, CA* **Jun 2024 – present**
- Part of the **Neural Networks** team building robust, secure and autonomous **perception and understanding**
- Building, training and evaluating **transformer-based** 3D reconstruction networks that run **real-time** on edge devices

**Software Engineer, Samsung Electronics** *Suwon, South Korea* **Sep 2021 – Aug 2022**
- Developed high-performance **5G-NR** virtual L1 layer as part of **P**hysical **U**plink **S**hared **Ch**annel team
- Utilized Intel®Intrinsics (**AVX-512**) for efficient parallel processing of data, focusing on cache bottleneck optimization
- Reduced bottlenecks in uplink signal processing pipeline to achieve upto **20% speedup**

**Network Engineer Intern, Samsung Electronics** *remote* **Jun 2020 – July 2020**
- Built an automated network load testing framework to evaluate performance of in-production load balancing services

**Summer Research Intern, TU Braunschweig** *Braunschweig, Germany* **May 2019 – July 2019**
- Designed and built **WeLineation**, a full-stack app using **Expectation Maximization** for medical image segmentation

## RESEARCH EXPERIENCE

**Master's Thesis - Prof. Sanjiban Choudhury** *Cornell University* **Feb 2023 – Apr 2024**
- Built a learning system using **Vision-Language transformers** to allow the transfer of human skills to household robots
- Collaborated on a **speech-interactive task planner** for human-robot collaborative cooking, and a web-based evaluator

**Undergraduate Research - Prof. Preethi Jyothi** *IIT Bombay & Microsoft*
*Improving code-switched Automatic Speech Recognition (ASR) using Transformers* **Aug 2020 – Jun 2021**
- Built a new bilingual **speech recognition** model conditioned on language using CUDA accelerated dynamic programming
*Improving Low Resource Code-switched ASR using Augmented Code-switched TTS* **Dec 2019 – Jun 2020**
- Used end-to-end ASR models trained on Hindi and English monolingual corpi and code-switched synthetic data to improve performance in low-resource settings

## PUBLICATIONS

- **Demo2Code:** From Summarizing Demonstrations to Synthesizing Code via Extended Chain-of-Thought *[NeurIPS 2023]*
- Improving **low resource code-switched ASR** using augmented code-switched TTS *[INTERSPEECH 2020]*
- **WeLineation:** crowdsourcing delineations for reliable ground truth estimation *[SPIE Medical Imaging 2020]*

## TEACHING ASSISTANTSHIPS

**Cornell University**

| | |
|---|---|
| **Intro. to Machine Learning** *Spring 2024* | **Computer System Organization & Programming** *Fall 2022, 2023* |
| **Intro. to Analysis of Algorithm** *Summer 2023* | **Computational Sustainability** *Spring 2023* |

**IIT Bombay**

| | |
|---|---|
| **Software Systems Lab** *Fall 2019, 2020* | **Calculus** *Fall 2018* |

## KEY PROJECTS

| | | | |
|---|---|---|---|
| **Psychological analysis of ChatGPT** | *Prof. Valerie Reyna* | **Cornell** | **Fall 2023** |
| Research course exploring decision making of LLMs in risky and ethically ambiguous situations | | | |
| **Modeling misinformation in organizations** | *Prof. Jon Kleinberg* | **Cornell** | **Spring 2023** |
| Formalize the effect of corruption in hierarchical organizations using information networks | | | |
| **Few-shot action recognition on egocentric data** | *Prof. Kilian Weinberger* | **Cornell** | **Fall 2022** |
| Building a two-head action recognition system for EPIC-Kitchens tackling long-tail labels | | | |
| **Morphological Inflection through Deep Learning** | *Prof. Pushpak Bhattacharyya* | **IITB** | **2021** |
| **Maze Solving with Evolutionary RL** | *Prof. S. Kalyanakrishnan* | **IITB** | **2020** |