

# Huong Ngo

Email: hvn2002@uw.edu Phone: 5084887263 GitHub: github.com/huongngo-8 Website: huongngo-8.github.io

## Education

---

**University of Washington, Seattle**

**Seattle, WA**

*Applied Computational Mathematical Sciences: Data Science & Statistics B.Sc.* 3.71 GPA

Sep 2020–Dec 2024

## Publications

---

**Objaverse-XL: A Universe of 10M+ 3D Objects**

**NeurIPS Dataset and Benchmarks Track 2023**

Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, Eli VanderBilt, Aniruddha Kembhavi, Carl Vondrick, Georgia Gkioxari, Kiana Ehsani, Ludwig Schmidt, Ali Farhadi

## Work Experiences

---

**Allen Institute for AI**

**Seattle, WA**

*Research Intern*

Oct 2023–Dec 2023

- Working in the **Perceptual Reasoning and Interaction Research (PRIOR) Team**

**USAFacts**

**Seattle, WA**

*Data Engineer Intern*

Jun 2023–Sep 2023

- Built batch processing pipeline in **Azure Synapse** for data ingestion into data warehouse, saving over **10000 dollars** of cloud storage costs and **150 hours** of developer hours of operational costs annually
- Implemented **PySpark ETL** data pipeline tool in Azure Synapse to automate extraction and transformation **10M** data entries from over **10000** Excel tables (unstructured data) leading to a **97% reduction** in manual work hours

**Paul G. Allen Center for Computer Science and Engineering**

**Seattle, WA**

*Deep Learning Research Assistant*

Mar 2023–ongoing

- Conducting **computer vision** and **multimodal learning** research under guidance of **Matt Detike**
- Applied large-scale data processing pipelines to **120M** object images using CLIP to annotate object aesthetic scores and build quality tiers in dataset
- Developing open-source distributed training of **OpenAI's Whisper** model on **1224 hours** of multilingual speech data with **PyTorch, PyTorch Lightning, Slurm, Weights and Biases**
- Implementing and designing modifications to Whisper model to expand multilingual speech transcription capabilities

## Teaching Experiences

---

**Paul G. Allen Center for Computer Science and Engineering**

**Seattle, WA**

*Machine Learning and Database Teaching Assistant*

Sep 2022–June 2023

**University of Washington, Department of Statistics**

**Seattle, WA**

*Statistics Tutor*

Sept 2021–June 2022

## Relevant Projects

---

**Text2Midi - Generating Symbolic Music Representation From Text**

- Architected a novel multimodal generative model that generates symbolic music representation from text descriptions by leveraging language modeling, contrastive language-music learning and pre-trained models
- Developed data processing pipeline to ingest, label and transform dataset of over **22000** songs for training
- Trained model that is a two-tower parallel Transformer-based encoder (text and music) using Music-BERT (RoBERTa) and BERT, Transformer-based decoder, and a joint embedding space

**Gehirn - Automated Generation of Symbolic Music Representation Datasets**

- Co-authored a paper that introduces a novel system for generating datasets with transcriptions, audio and text captions for music generation tasks
- Designed system that is a pipeline connecting a Python data mining script, a Transformer-based automatic music transcription model to obtain transcriptions, and GPT-3.5 text completion to produce semantic descriptions

## Skills

---

**Languages:** Python, SQL, R, Java

**Technologies:** NumPy, pandas, matplotlib, PyTorch, PyTorch Lightning, Weights and Biases, scikit-learn, PySpark, SparkSQL, OpenCV, AWS, BeautifulSoup

**Developer Tools:** Jupyter, GitHub/Git, Slurm

## Relevant Coursework

---

**University of Washington, Seattle**

**Seattle, WA**

Machine Learning Systems, Machine Learning for Big Data, Machine Learning, Artificial Intelligence, Databases, Data Structures & Algorithms, Linear Algebra, Statistics & Probability