# Xuweiyi Chen

xuweic@email.virginia.edu | (206)532-9635 | xuweiyichen.github.io

## EDUCATIONAL BACKGROUND

**UNIVERSITY OF VIRGINIA**                                                                           **Charlottesville, VA**
*Ph.D. in Computer Science and Engineering*                                           *Aug. 2024 – May 2029 (Expected)*
Overall GPA: 4.0/4.0
Concentration: **3D Computer Vision and Multimodal Learning**

**UNIVERSITY OF MICHIGAN**                                                                                    **Ann Arbor, MI**
*M.S. in Computer Science and Engineering*                                                          *Aug. 2022 – May 2024*

**UNIVERSITY OF WASHINGTON**                                                                                    **Seattle, WA**
*B.S. in Applied and Computational Mathematical Sciences, CUM LAUDE*              *Sep. 2018 – June 2022*
Honors: $6000 CoMotion Mary Gates Innovation Scholarship
           $3000 Usha and S. Rao Varanassi SAFS Scholarship

## SELECTED INTERNSHIPS

**PixAI.art**                                                                                                                     **Remote**
*Machine Learning Engineering Intern.*                                                              *Jan. 2024 – May. 2024*
- Experience Large-scale Pretrained Image and Video Diffusion Models using 128 H100 GPUs
- Deployed Video Diffusion Models to generate personalized 2D cartoon-based natural videos for user applications.
- Led the integration of 3D computer vision with diffusion models, driving research and development of innovative user-facing products.

## SELECTED FIRST-AUTHOR PUBLICATIONS

**Learning 3D Representations from Procedural 3D Programs**
*UVA CV LAB supervised Prof. Zezhou Cheng*                                                                   *Nov. 2024*
- Procedurally generated shapes offer a scalable, copyright-free, and geometrically diverse alternative to labor-intensive human-designed 3D datasets like ShapeNet.
- We use procedurally generated 3D shapes to achieve strong results in object classification, part segmentation, and few-shot learning.
- Point-MAE-Zero can perform masked point cloud completion without fine-tuning.

**SAB3R: Semantic-Augmented Backbone in 3D Reconstruction**
*UVA CV LAB supervised Prof. Zezhou Cheng*                                                                   *Nov. 2024*
- Developed SAB3R, a method to distill 2D semantic features into 3D vision foundation models, enhancing semantic understanding while retaining spatial reasoning.
- Introduced a novel task, *Map and Locate*, enabling multi-view 3D open vocabulary semantic segmentation.
- Validated on depth estimation and pose regression, achieving improved 2D semantics without compromising 3D performance.

**3D-GRAND A Million-Scale Dataset for 3D-LLMs with Better Grounding and Less Hallucination**
*SLED lab in the University of Michigan supervised Prof. Joyce Chai & Prof. David Fouhey*              *Aug. 2024*
- Introduced 3D-GRAND, a large-scale dataset with 40,087 household scenes and 6.2 million densely grounded scene-language instructions to improve 3D-Language models (3D-LLMs).
- Proposed 3D-POPE, a benchmark to evaluate hallucinations in 3D-LLMs, enabling fair comparisons across models.
- Demonstrated that instruction tuning with 3D-GRAND significantly enhances grounding capabilities, emphasizing the importance of large-scale 3D-text datasets for advancing embodied AI research.

**Multi-Object Hallucination in Vision-Language Models**                                       **NeurIPS 2024**
*SLED lab in the University of Michigan supervised Prof. Joyce Chai & Prof. David Fouhey*              *July 2024*
- Investigated multi-object hallucination in Large Vision Language Models (LVLMs) using Recognition-based Object Probing Evaluation (ROPE), focusing on the distribution of object classes within a single image and visual referring prompts.
- Found that LVLMs exhibit more hallucinations when tasked with recognizing multiple objects compared to a single object, influenced by object class distribution and model behaviors.
- Identified key factors such as salience, frequency, and model intrinsic behaviors that contribute to hallucination, aiming to improve LVLMs' recognition and reasoning capabilities in complex visual scenes.

**LLM-Grounder: Open-Vocabulary 3D Visual Grounding with Large Language Model as an Agent. ICRA 2024**
*SLED lab in the University of Michigan supervised Prof. Joyce Chai*                                         *Aug. 2023*
- Present the first method capable of localizing novel objects in 3D scenes using Neural Radiance Field (NeRF) and Large Language Models (LLMs) through iterative, natural language-based interactions.
- Enables a more human-like interaction with 3D objects in a learned 3D scene representation.
- Evaluated and shown that dynamic grounding outperforms static grounding in terms of accuracy, 3DIoU, and human ratings.