

# A “**Splatacular**” Year of 3D Reconstruction

Songyou Peng

**Google DeepMind**

Stanford University

May 1, 2025

The background image shows a modern living room with large windows overlooking a city skyline. The room is furnished with a white sofa, a brown armchair, and a wooden coffee table. A television is mounted on a wall covered in green plants. The overall aesthetic is clean and contemporary.

Intelligent systems interact with **3D environments**

## 3D Reconstruction

Create digital twins from real scenes

## 3D Scene Understanding

Analyze the scene digitally

# My Research During PhD

## Learn to Reconstruct and Understand 3D World

ConvOccNet  
ECCV 2020 (Spotlight)

MonoSDF  
NeurIPS 2022

Shape As Points  
NeurIPS 2021 (Oral)

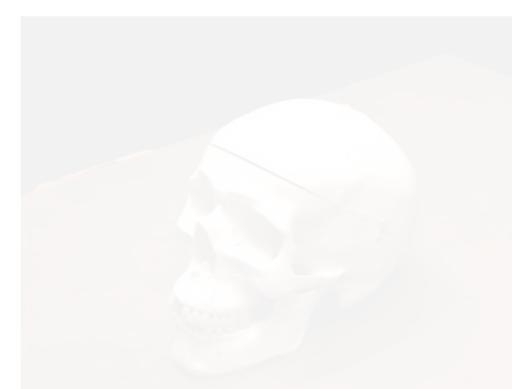
KiloNeRF  
ICCV 2021  
runs now at 50 ips on a GTX 1080 Ti



NICE-SLAM  
CVPR 2022



NICER-SLAM  
3DV 2024 (Oral)



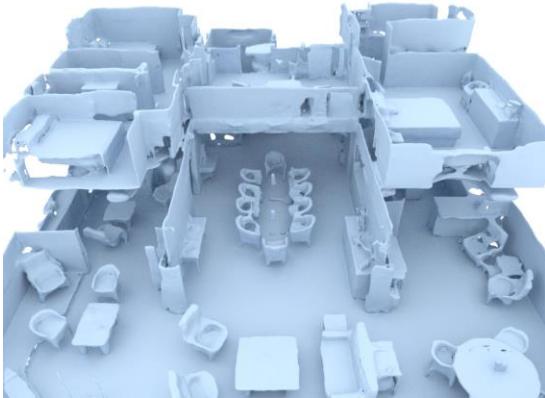
UNISURF  
ICCV 2021 (Oral)



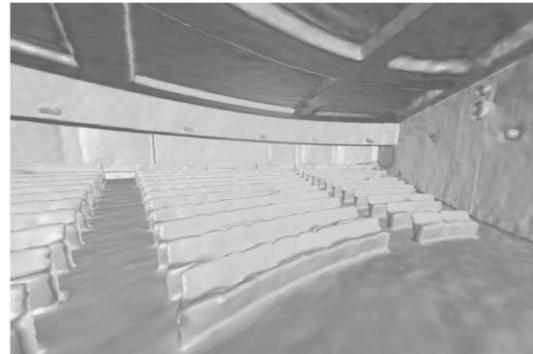
OpenScene  
CVPR 2023

# My Research During PhD

Learn to Reconstruct and Understand 3D World



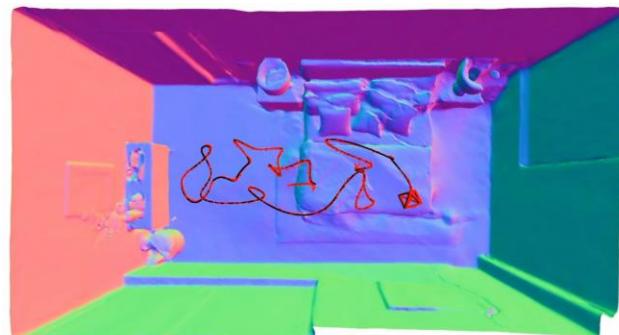
**ConvOccNet**  
ECCV 2020 (Spotlight)



**MonoSDF**  
NeurIPS 2022



**NICE-SLAM**  
CVPR 2022



**NICER-SLAM**  
3DV 2024 (Best Honor. Men.)



**UNISURF**  
ICCV 2021 (Oral)

**OpenScene**  
CVPR 2023

# My Research During PhD

Learn to Reconstruct and Understand 3D World

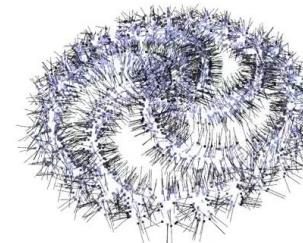
## Topic #2: Fast Inference

ConvOccNet  
ECCV 2020 (Spotlight)

MonoSDF  
NeurIPS 2022

**Shape As Points**  
NeurIPS 2021 (Oral)

KiloNeRF  
ICCV 2021



runs now at 50 fps on a GTX 1080 Ti

NICE-SLAM  
CVPR 2022

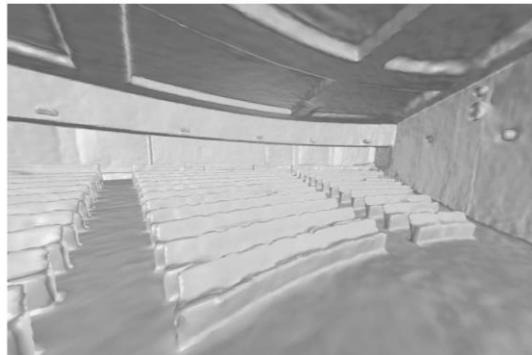
NICER-SLAM  
3DV 2024 (Oral)

UNISURF  
ICCV 2021 (Oral)

OpenScene  
CVPR 2023

# My Research During PhD

Learn to Reconstruct and Understand 3D World



ConvOccNet  
ECCV 2020 (Spotlight)

MonoSDF  
NeurIPS 2022

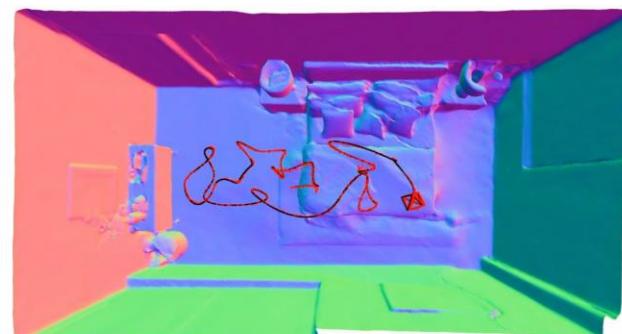
Shape As Points  
NeurIPS 2021 (Oral)

runs now at 50 ips on a GTX 1080 Ti

KiloNeRF  
ICCV 2021



NICE-SLAM  
CVPR 2022



NICER-SLAM  
3DV 2024 (Best Paper Honorable)



UNISURF  
ICCV 2021 (Oral)

Topic #3:  
Reconstruct from 2D Observations



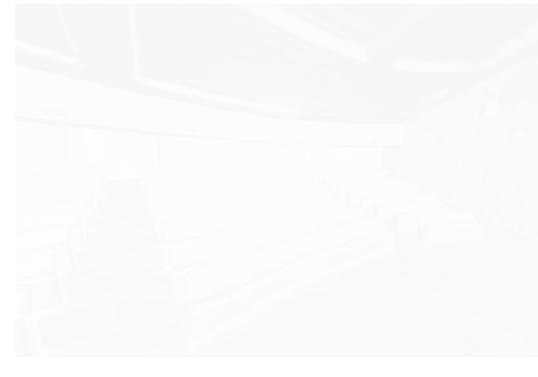
OpenScene  
CVPR 2023

# My Research During PhD

Learn to Reconstruct and Understand 3D World



**ConvOccNet**  
ECCV 2020 (Spotlight)



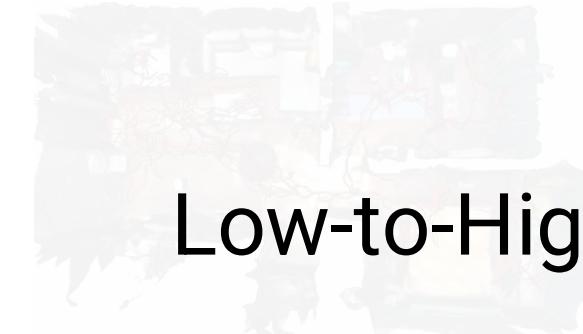
**MonoSDF**  
NeurIPS 2022



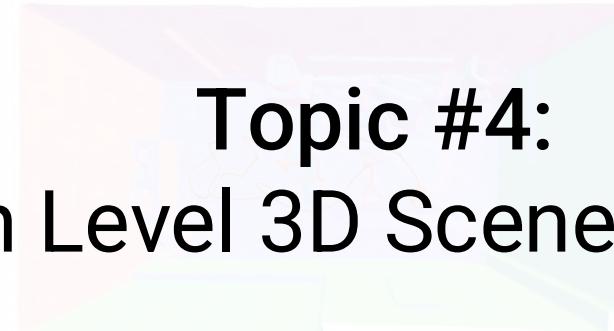
**Shape As Points**  
NeurIPS 2021 (Oral)



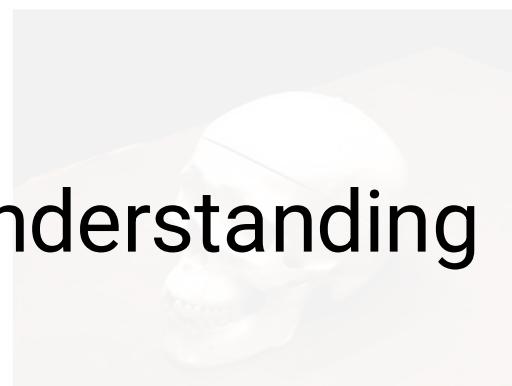
**KiloNeRF**  
ICCV 2021  
runs now at 50 ips on a GTX 1080 Ti



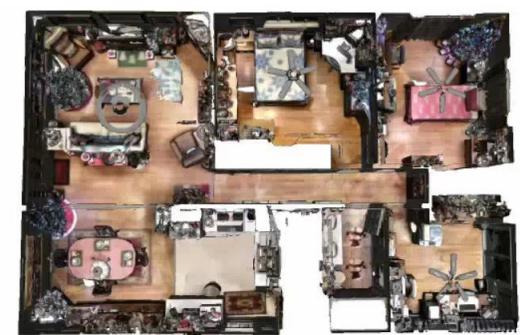
**Topic #4:**  
**Low-to-High Level 3D Scene Understanding**



**NICE-SLAM**  
CVPR 2022



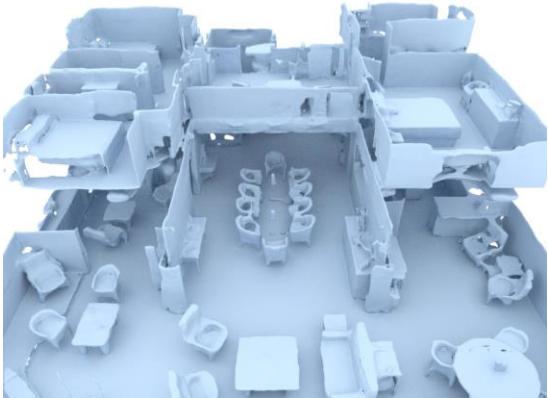
**UNISURF**  
ICCV 2021 (Oral)



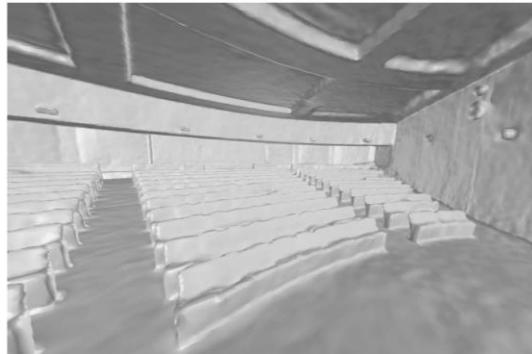
**OpenScene**  
CVPR 2023 6

# My Research During PhD

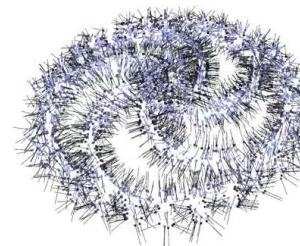
Learn to Reconstruct and Understand 3D World



**ConvOccNet**  
ECCV 2020 (Spotlight)



**MonoSDF**  
NeurIPS 2022

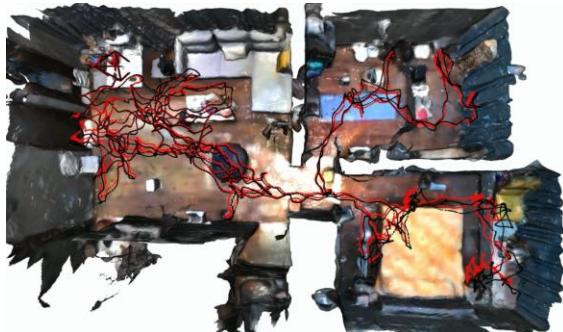


**Shape As Points**  
NeurIPS 2021 (Oral)

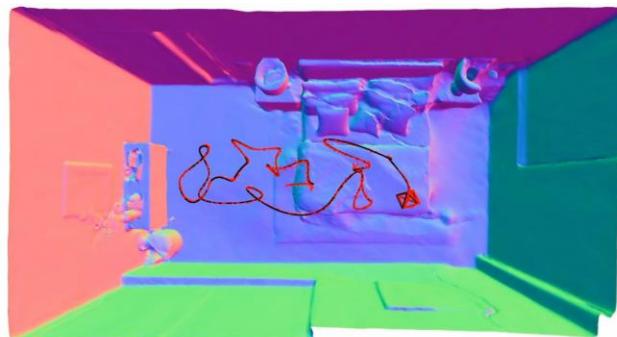


runs now at 50 fps on a GTX 1080 Ti

**KiloNeRF**  
ICCV 2021



**NICE-SLAM**  
CVPR 2022



**NICER-SLAM**  
3DV 2024 (Best Paper Honorable)



**UNISURF**  
ICCV 2021 (Oral)



**OpenScene**  
CVPR 2023

# My PhD Thesis

Already Tackling Some Challenges in 3D Reconstruction

- Reconstruct **at scale**
- Reconstruct **at speed**
- Reconstruct **from 2D observations**



# An Ideal 3D Reconstruction Pipeline

Instant, Pose-Free, Real-World 3D Everywhere



Feedforward



Pose-Agnostic



Dynamic



Fast Rendering



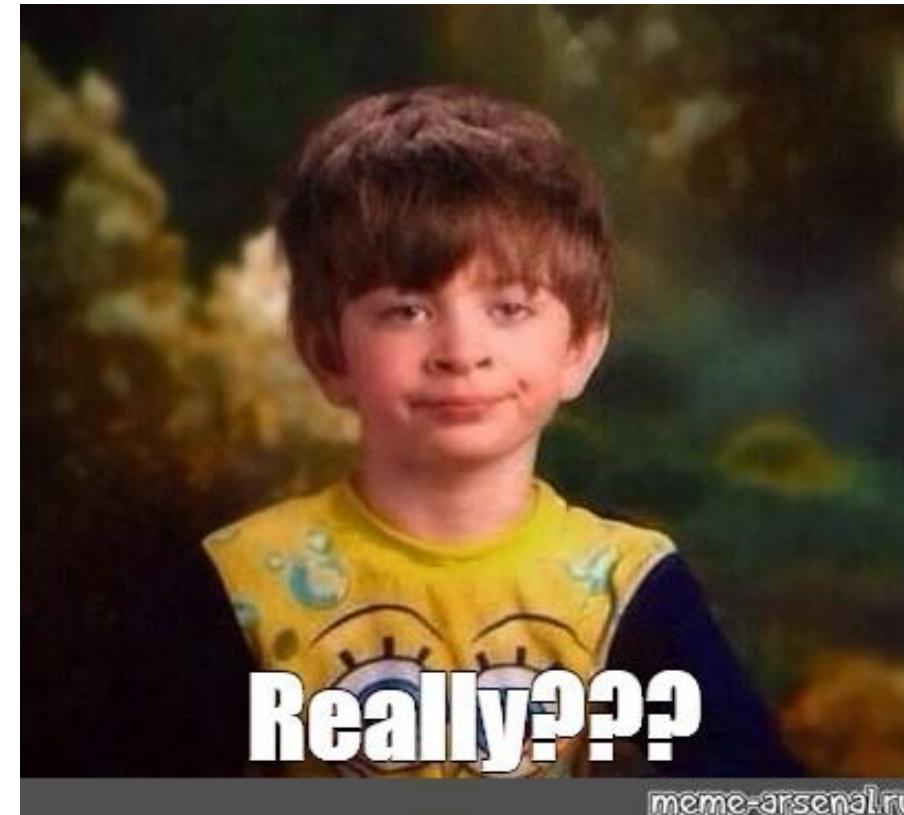
Arbitrary Lengths



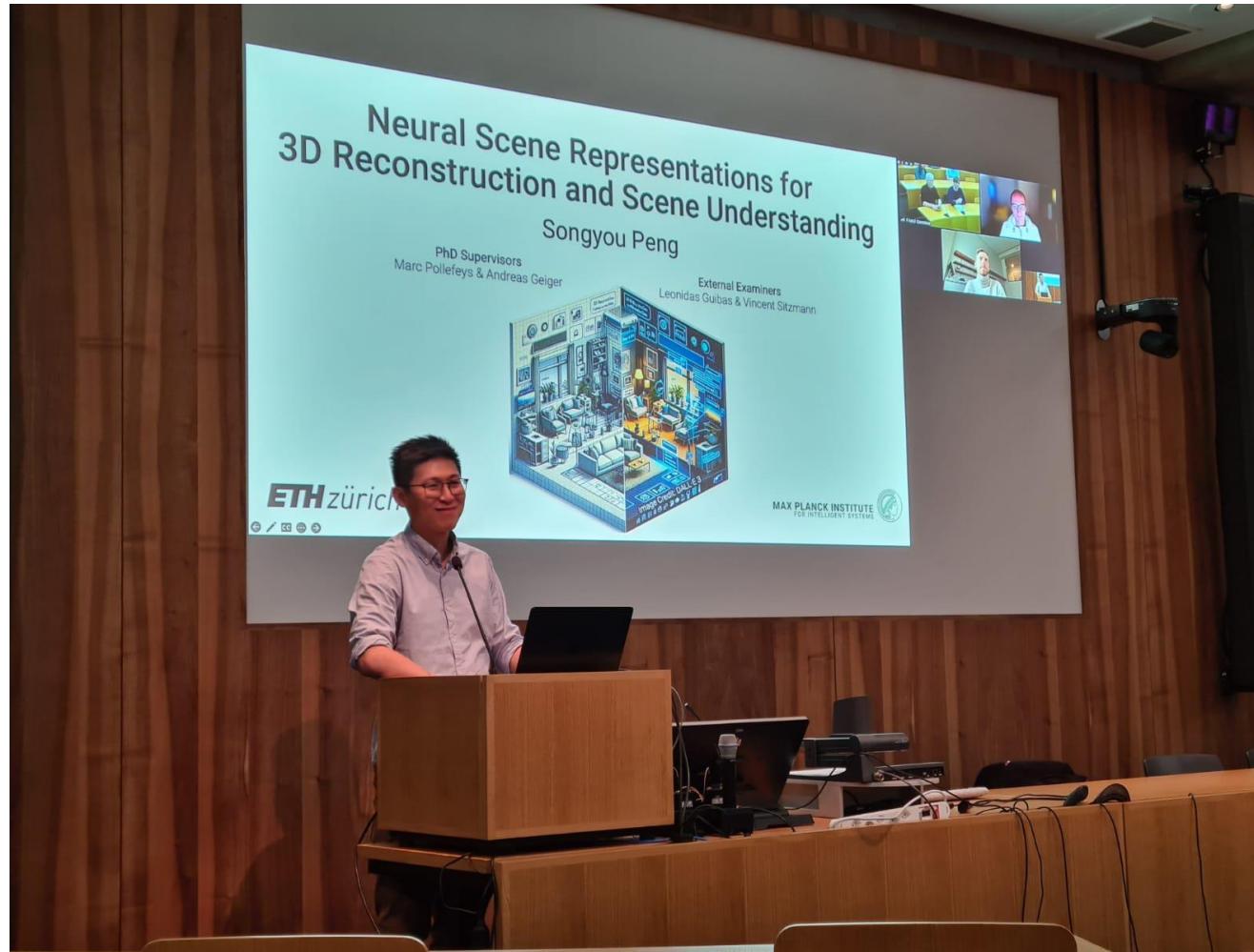
Lighting-Robust

“People overestimate what they can do in one year, and underestimate what they can do in ten years.”

--- Bill Gates



# I Defended on Nov 2023



How much can one push forward until **Nov 2024**?

# What Came Up in 2023?



3DGS  
→



Input Posed Images

# What Came Up in 2023?



# An Ideal 3D Reconstruction Pipeline

Instant, Pose-Free, Real-World 3D Everywhere



Feedforward



Pose-Agnostic



Dynamic



Fast Rendering



Arbitrary Lengths



Lighting-Robust

# An Ideal 3D Reconstruction Pipeline

Instant, Pose-Free, Real-World 3D Everywhere



Feedforward



Pose-Agnostic



Dynamic



Fast Rendering



Arbitrary Lengths



Lighting-Robust

# Feedforward 3D Gaussian Splatting





# DepthSplat

## Connecting Gaussian Splatting and Depth

[haofeixu.github.io/depthsplat/](http://haofeixu.github.io/depthsplat/)



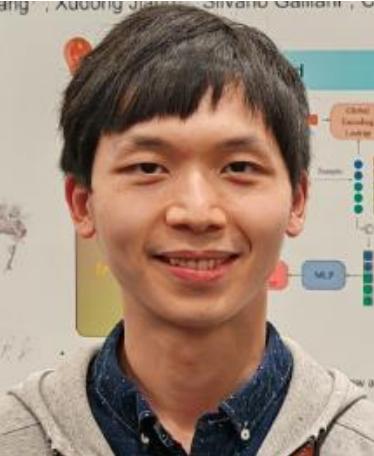
### CVPR 2025



Haofei Xu



Songyou Peng



Fangjinhua Wang



Hermann Blum



Daniel Barath



Andreas Geiger



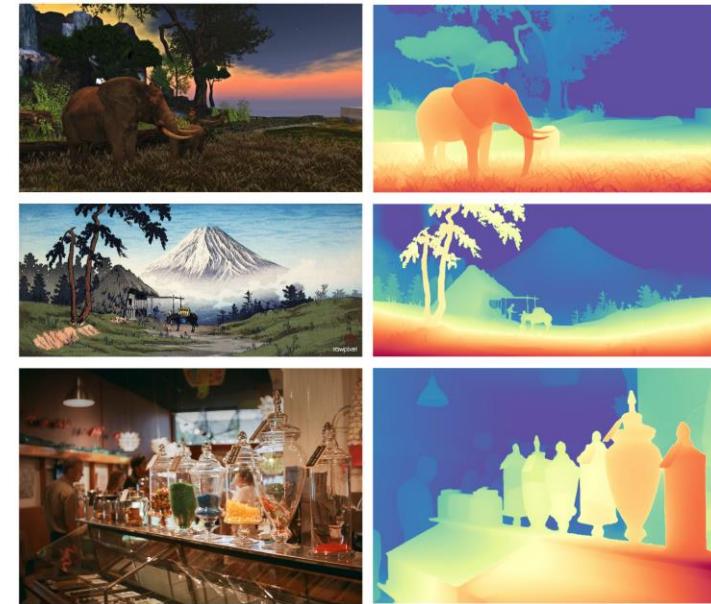
Marc Pollefeys

# Motivation



pixelSplat / MVSplat

- + Multi-view Consistent
- Lack robustness

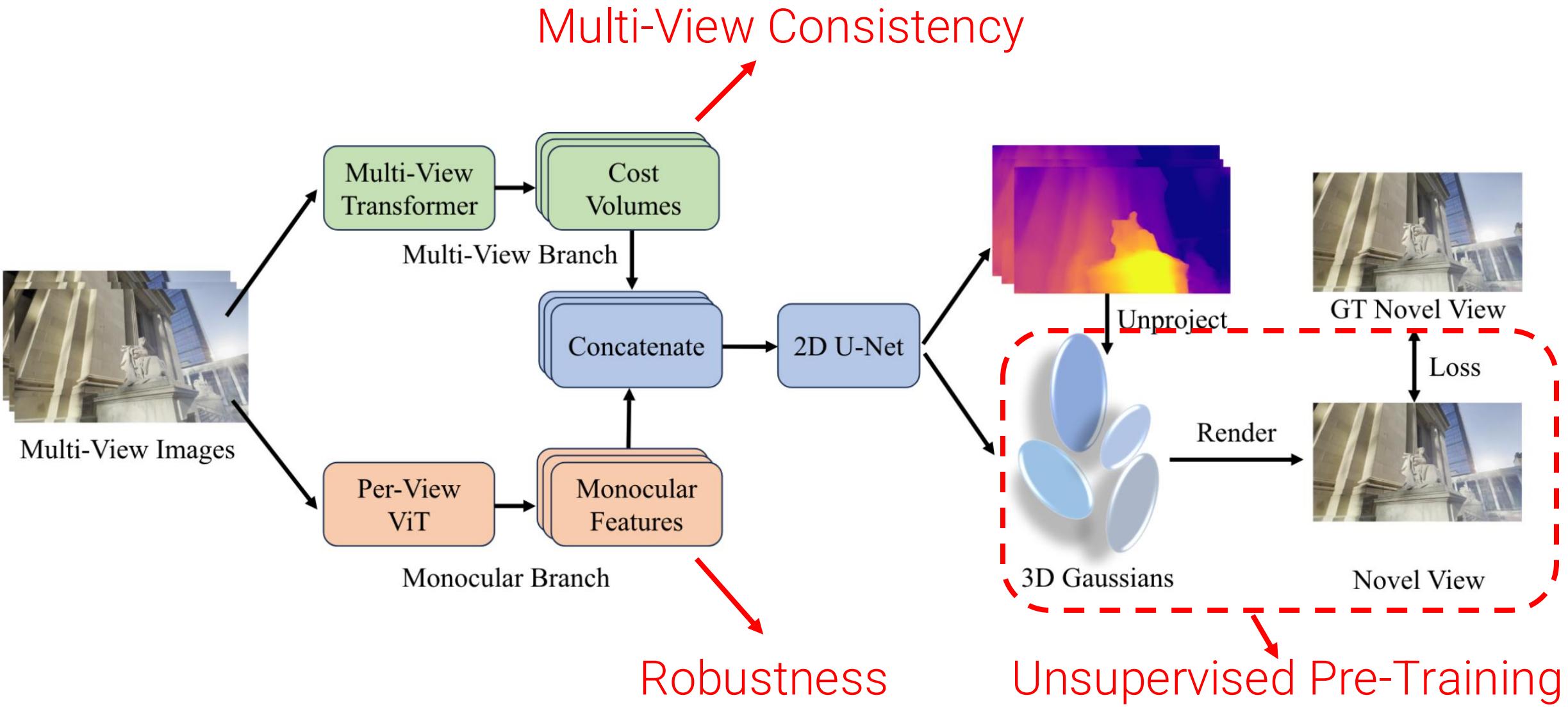


Depth Anything v2

- + Robust
- Unknown scale & shift

Can they benefit each other?

# Pipeline



# Feedforward View Synthesis



6 Input Views

DepthSplat



Rendering

# Feedforward View Synthesis



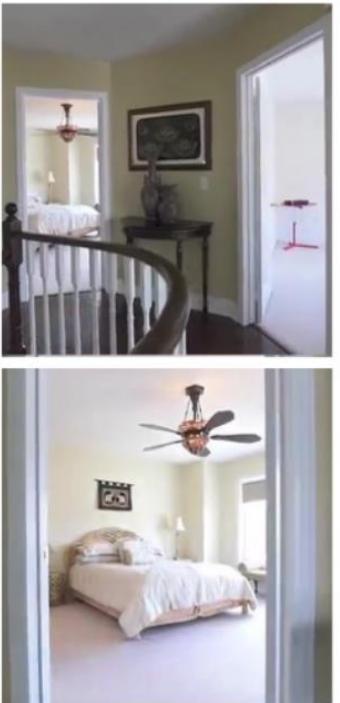
12 Input Views

DepthSplat



Rendering

# Comparison



Input



pixelSplat



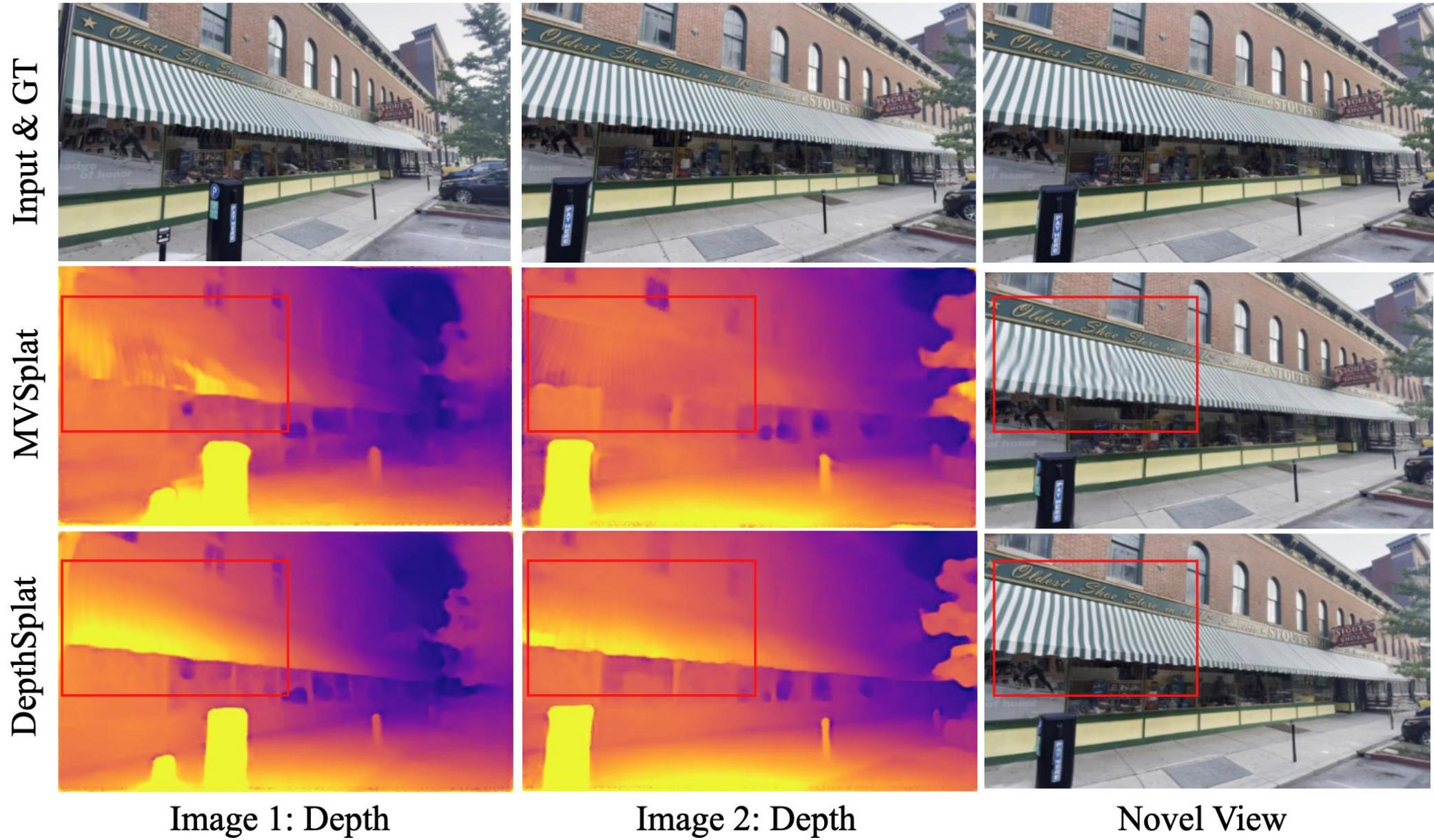
MVSplat



DepthSplat

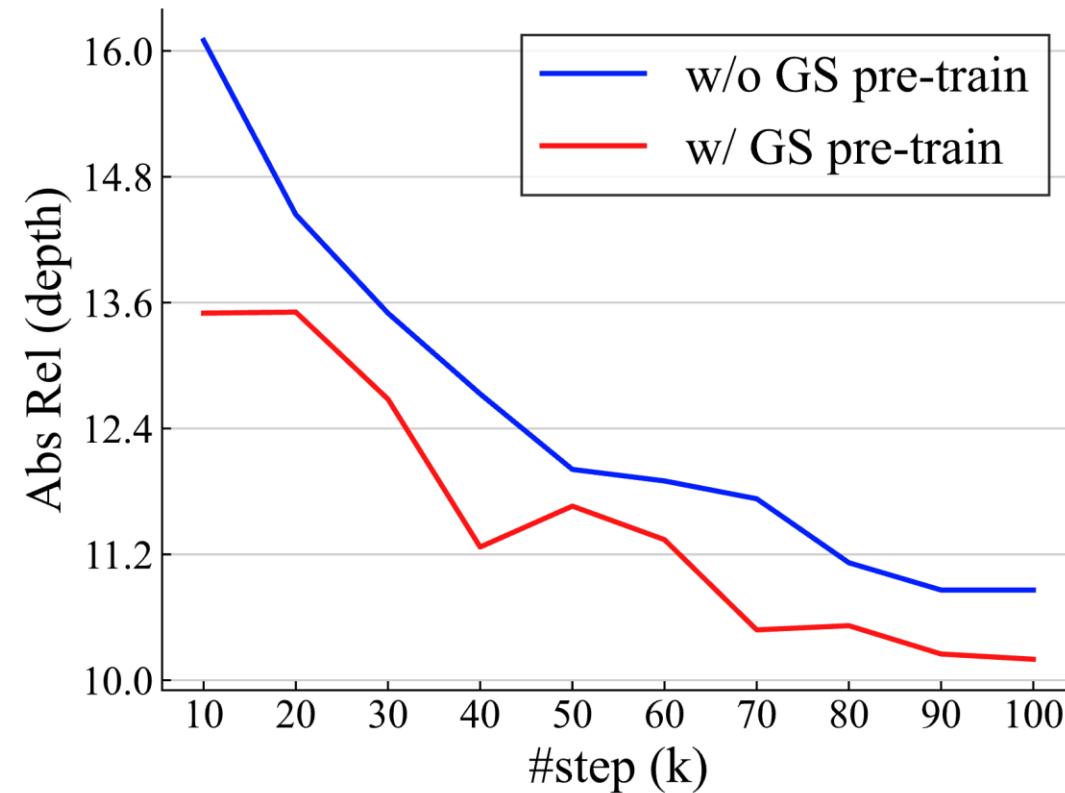
DepthSplat is significantly **more robust!**

# Depth → Gaussian Splatting



# Gaussian Splatting → Depth

- Unsupervised depth pre-training on RealEstate10K
- Supervised depth fine-tuning on TartanAir & VKITI2



Validation curves of depth prediction error

# Take-home Messages

- Depth and Gaussian splatting are helping each other!
- Feed-forward Gaussian splatting for large-scale scenes
- A few posed images as input... **This is not a practical setting!**

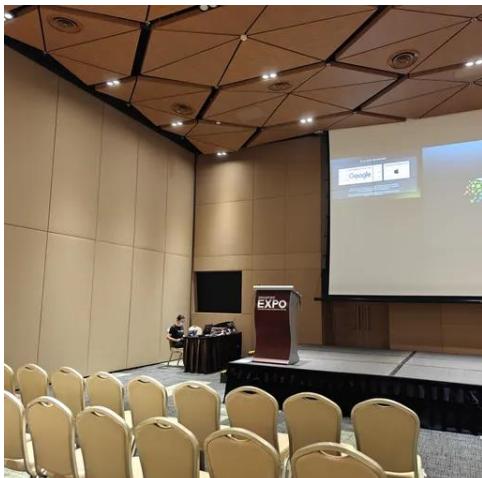


Image 1



Image 2

**Non-trivial to get camera poses!**

# An Ideal 3D Reconstruction Pipeline

Instant, Pose-Free, Real-World 3D Everywhere



Feedforward



Pose-Agnostic



Dynamic



Fast Rendering



Arbitrary Lengths



Lighting-Robust

# An Ideal 3D Reconstruction Pipeline

Instant, Pose-Free, Real-World 3D Everywhere



Feedforward



Pose-Agnostic



Dynamic



Fast Rendering

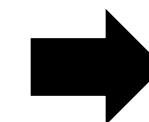
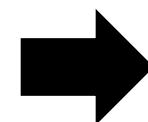


Arbitrary Lengths



Lighting-Robust

# Goal: Unposed Feedforward 3DGS



3D Gaussians

Novel Views

Input Images **w/o** poses

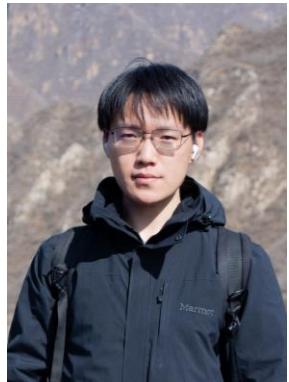
# No Pose, No Problem 🎭

## Surprisingly Simple 3D Gaussian Splats from Sparse Unposed Images

### (a.k.a NoPoSplat)



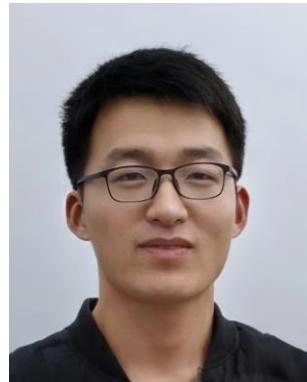
ICLR 2025 (**Oral**, top 1.8%)



Botao Ye



Sifei Liu



Haofei Xu



Xuetong Li



Marc Pollefeys

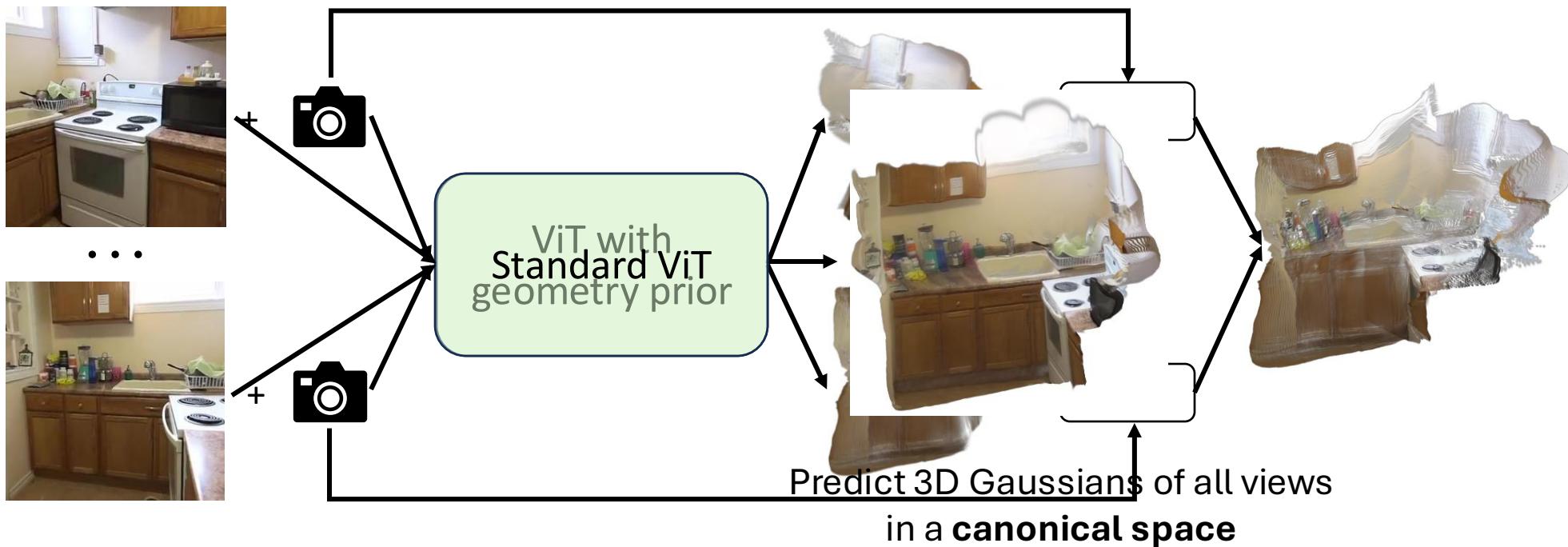


Ming-Hsuan Yang

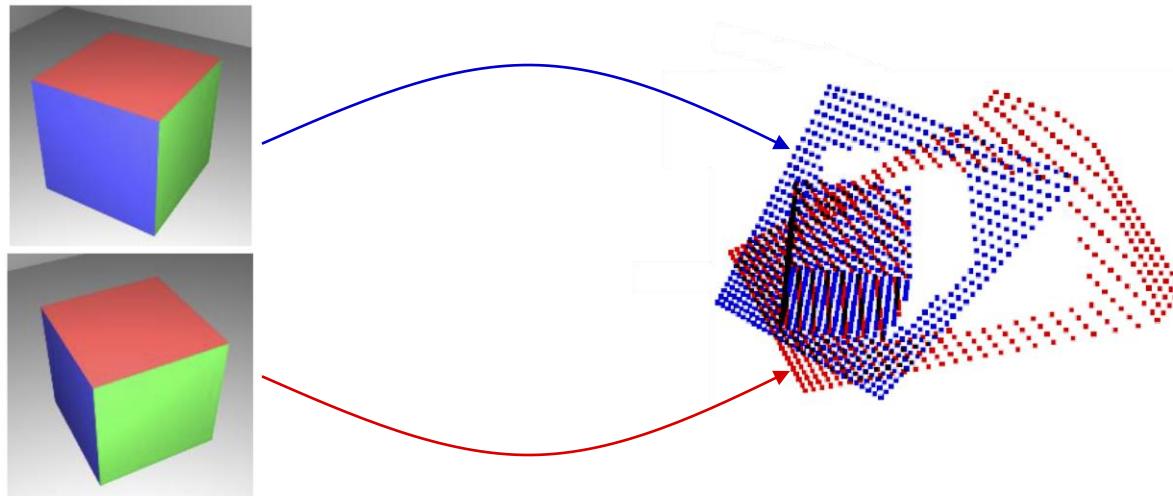


**Songyou Peng**

# Previous Feed-forward 3DGS



# Canonical Prediction



🤔 Does a similar philosophy apply to Gaussians?

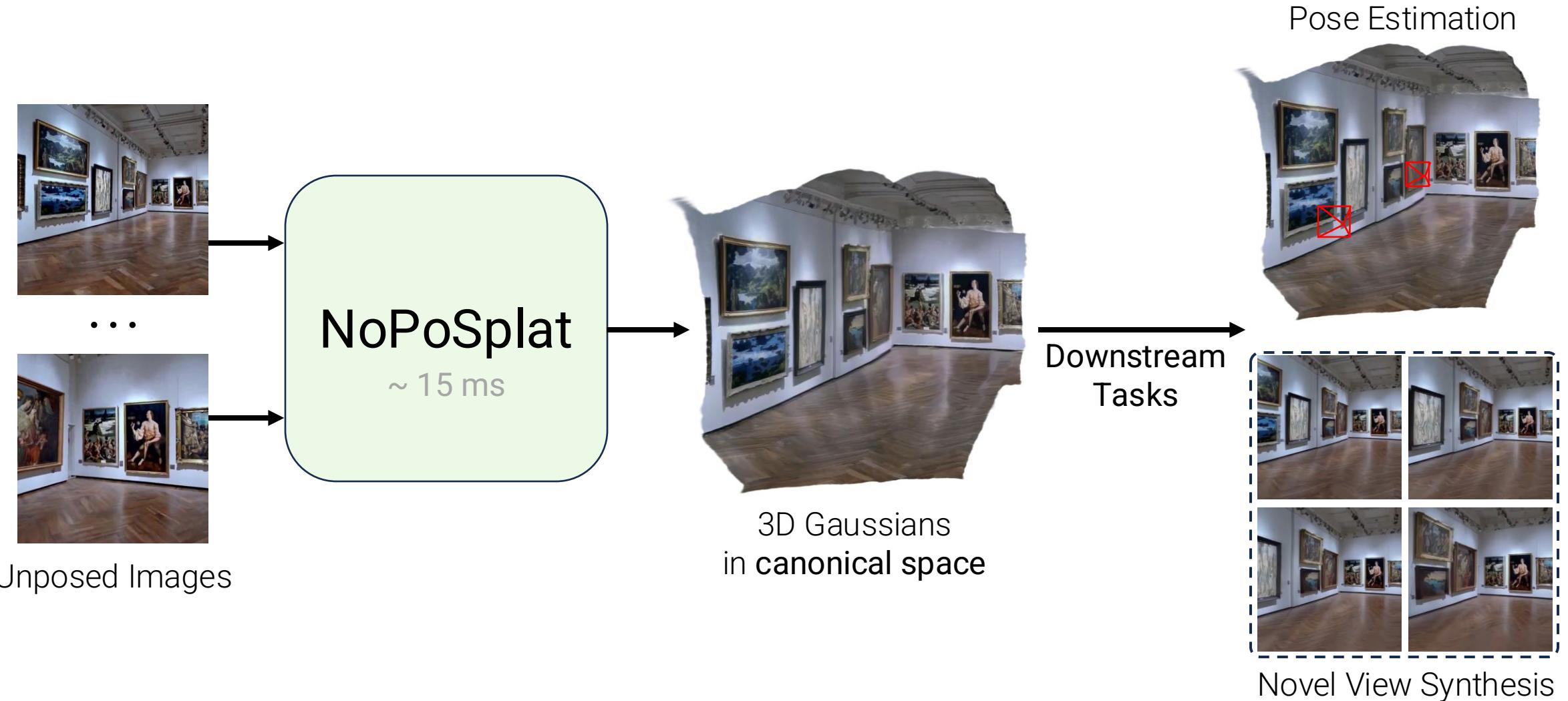
## Point Maps

- Discrete representation
- Ground truth depth needed

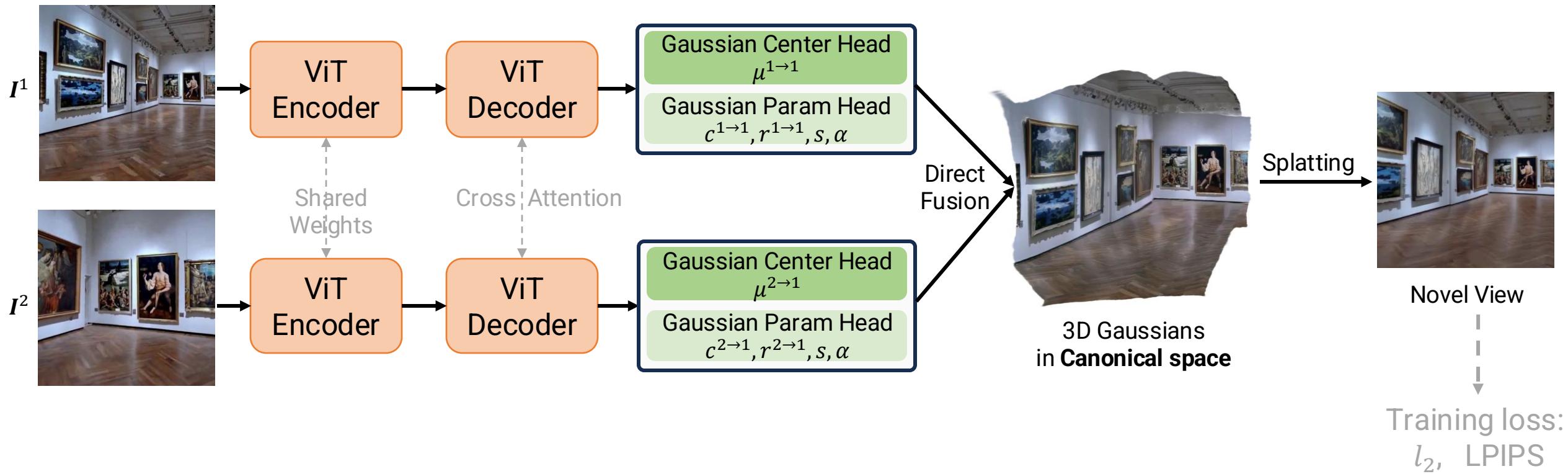
## 3DGS

- + Novel view synthesis
- + Video data for training

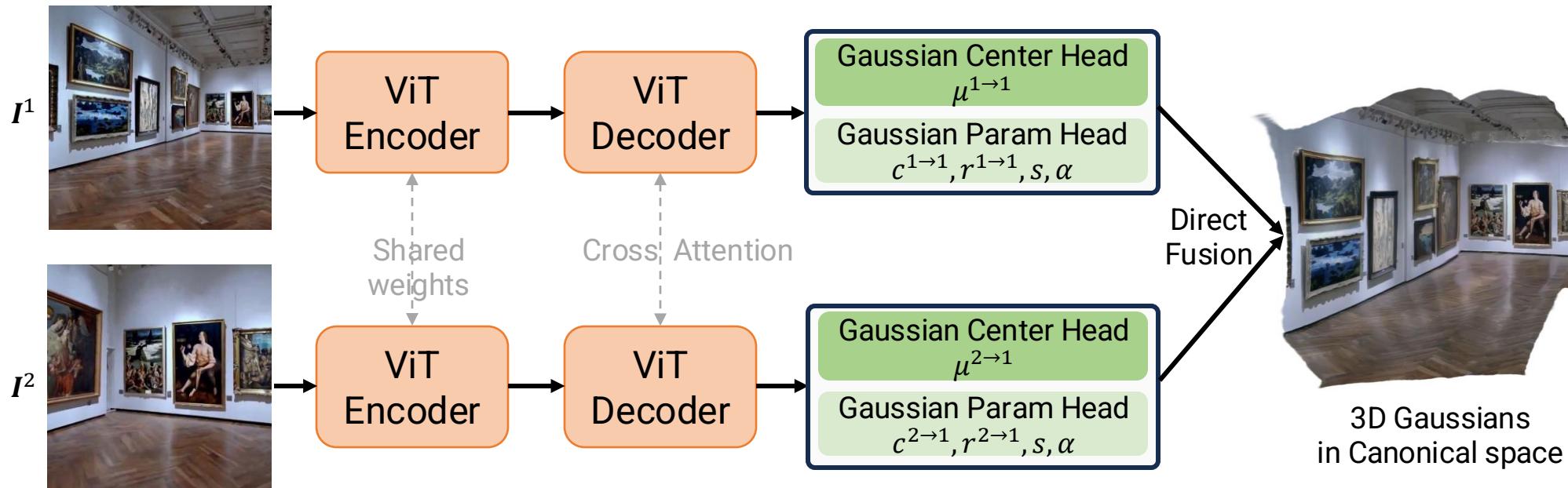
# NoPoSplat



# Architecture

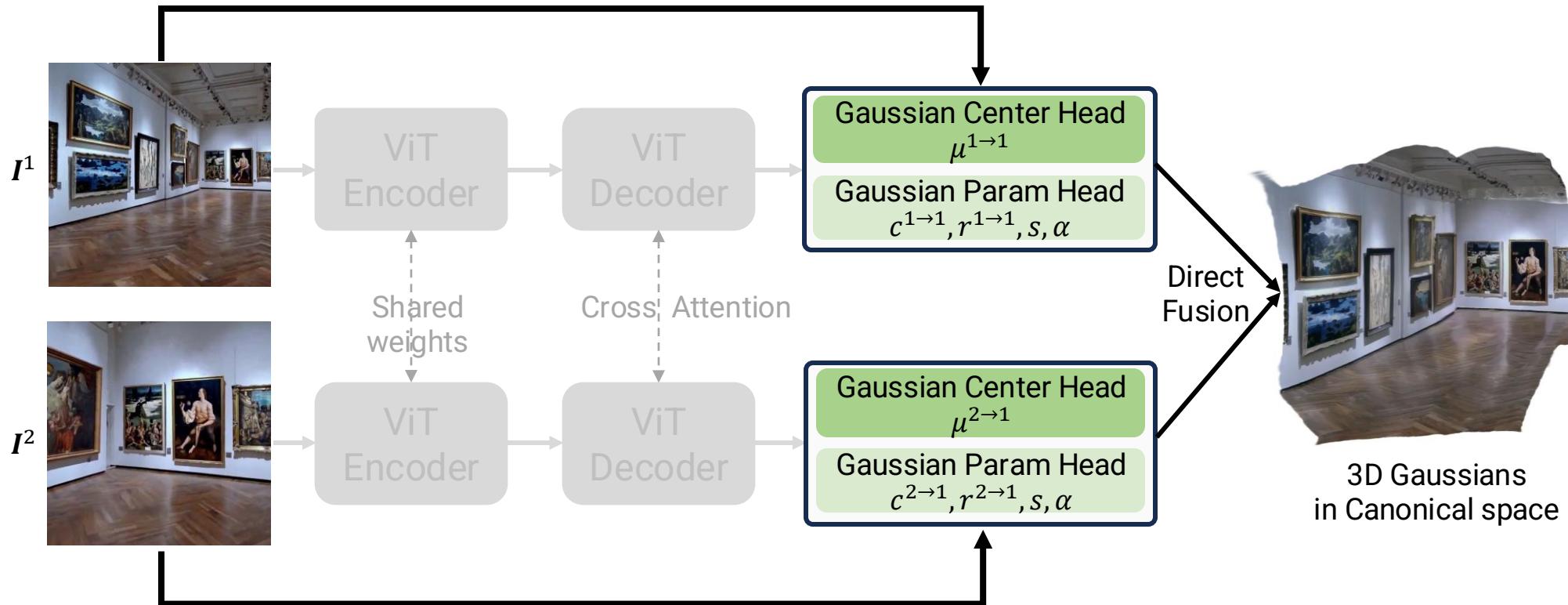


# Issue 1: Blurry Rendering



# Issue 1: Blurry Rendering

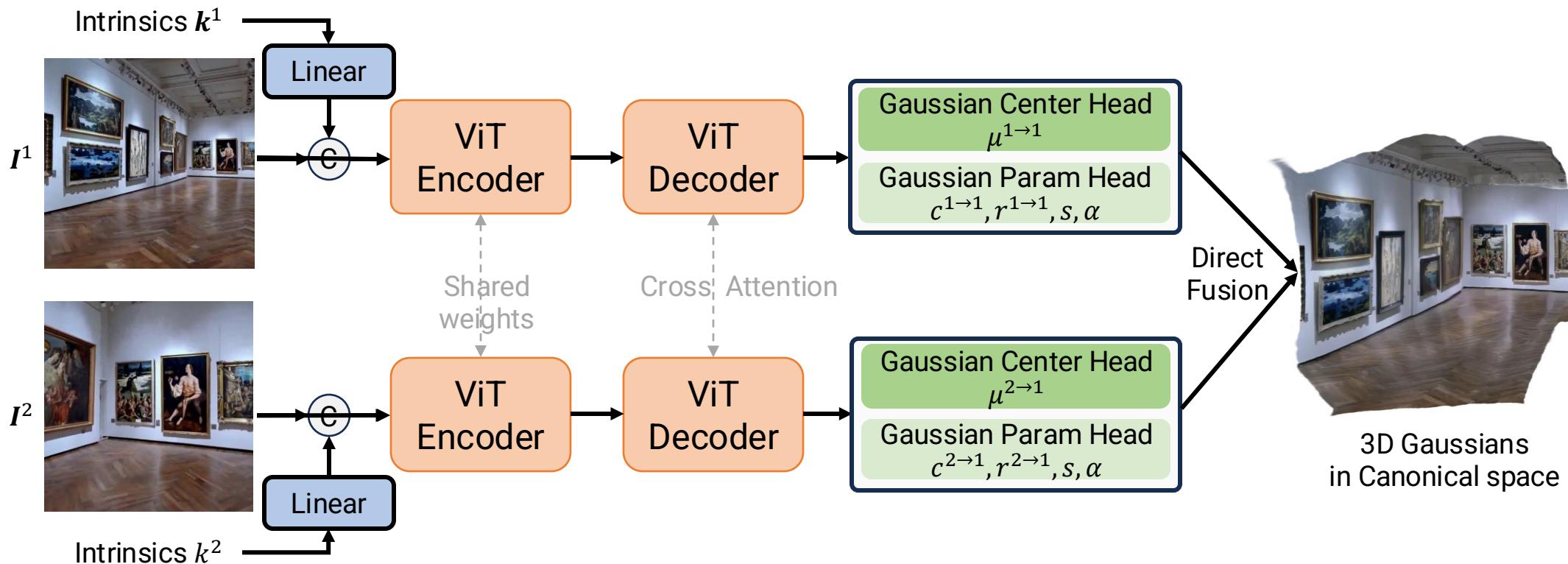
Solution: Add a shortcut!



# Issue 2: Scale Ambiguity

Solution: Add the intrinsic embeddings!

$$p = K(RP + t)$$



# Issue 3: Inaccurate Pose Estimation

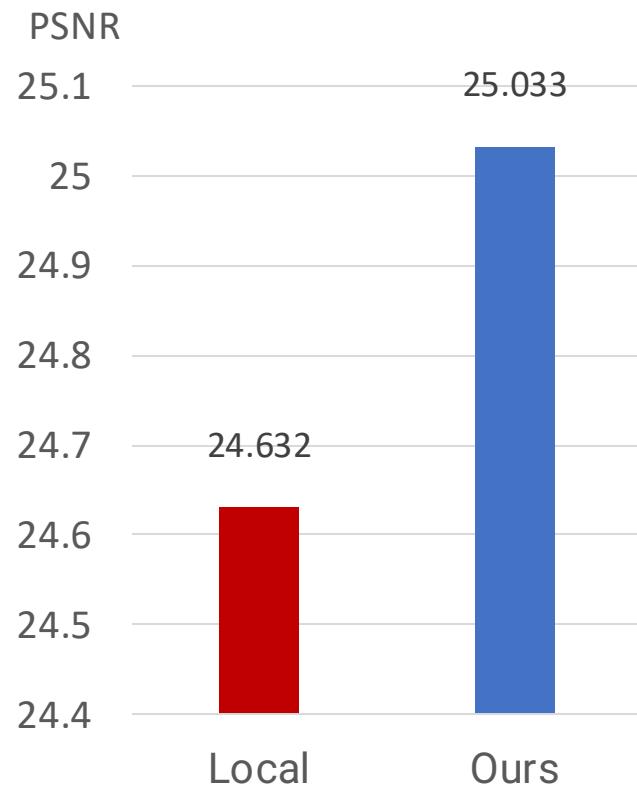
**Solution:** coarse-to-fine estimation

- Coarse stage: run RANSAC-PnP on Gaussian centers
- Refine stage: optimize with photometric loss

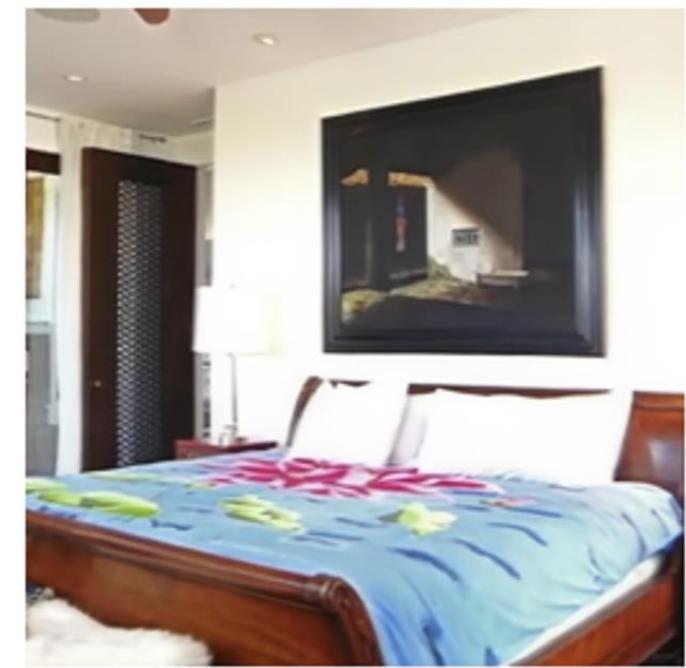
PnP	Photometric	5°	10°	20°
✓	✓	0.318	0.538	0.717
✓		0.287	0.506	0.692
	✓	0.017	0.027	0.051

# Ablation

## Canonical Gaussian prediction



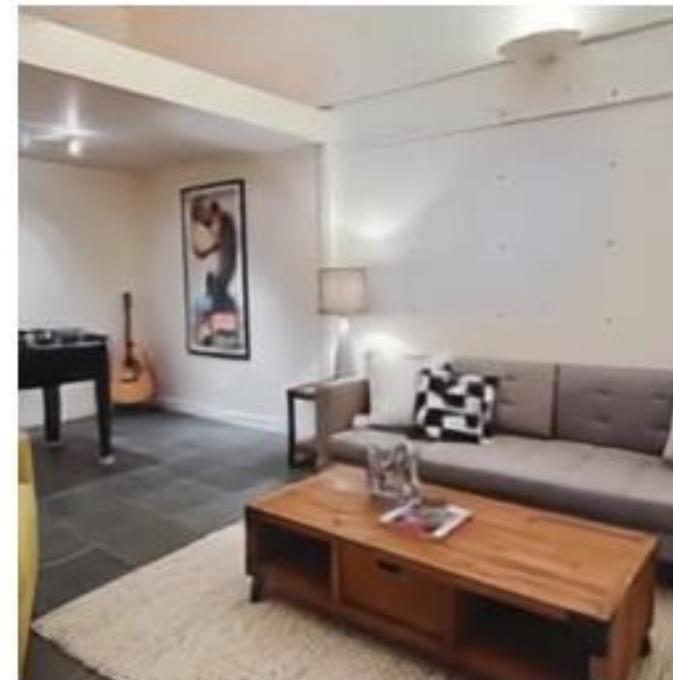
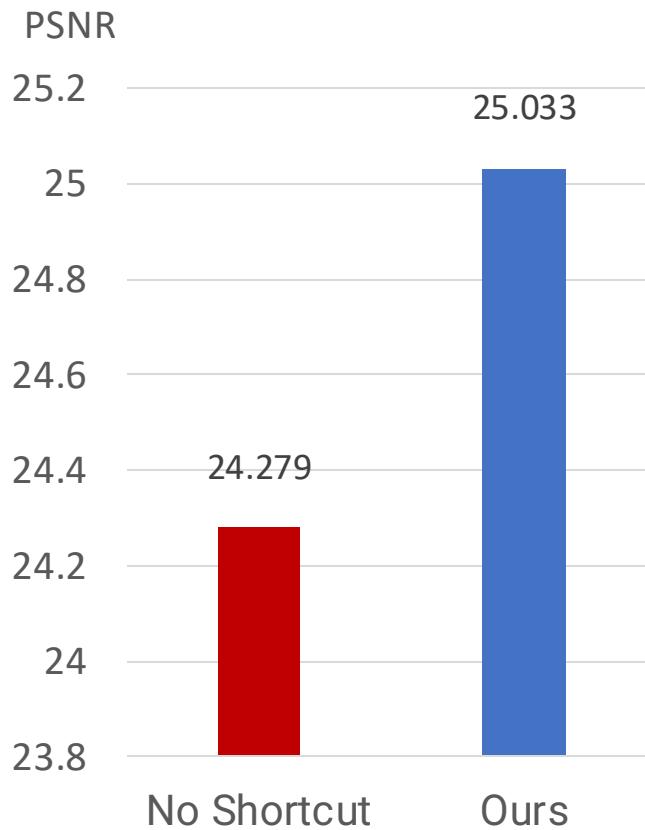
Local



Canonical

# Ablation

Image shortcut leads to sharper details



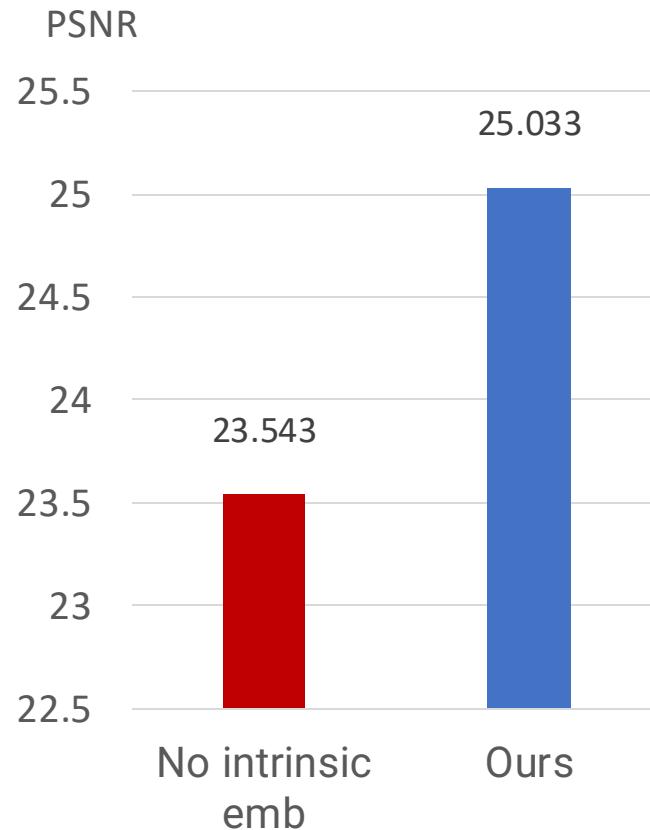
No Shortcut



Ours

# Ablation

## Intrinsic embedding



No Intrinsic  
Emb



Ours



GT

**What is More...**

# Accurate Pose Estimation

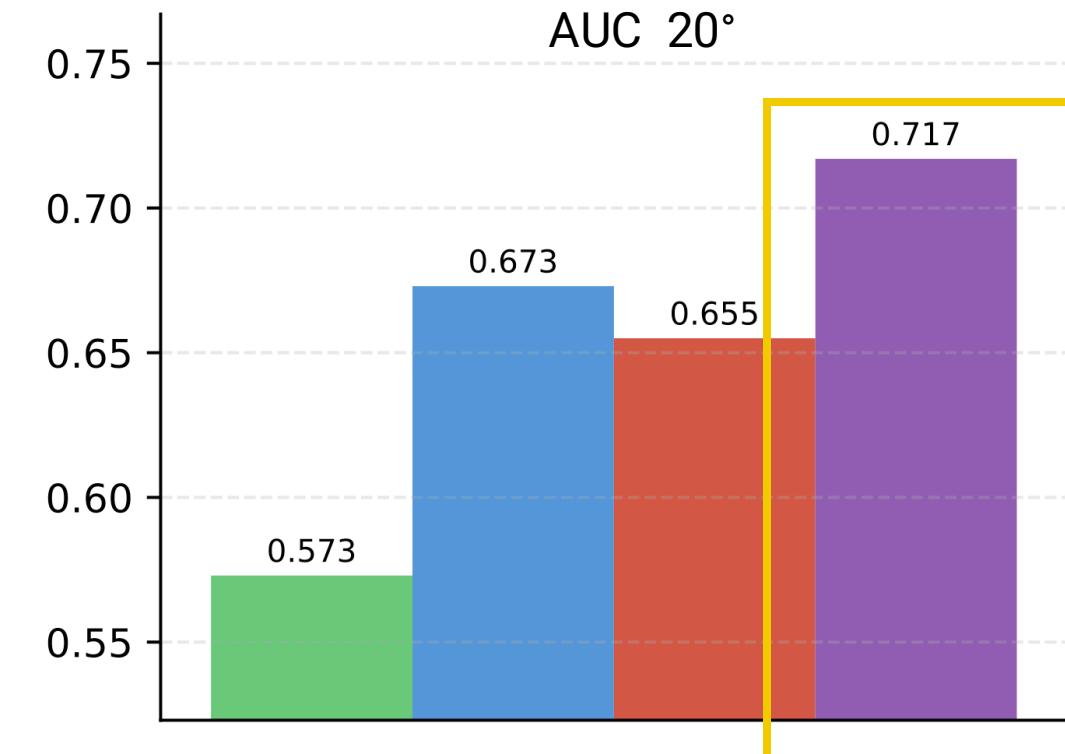
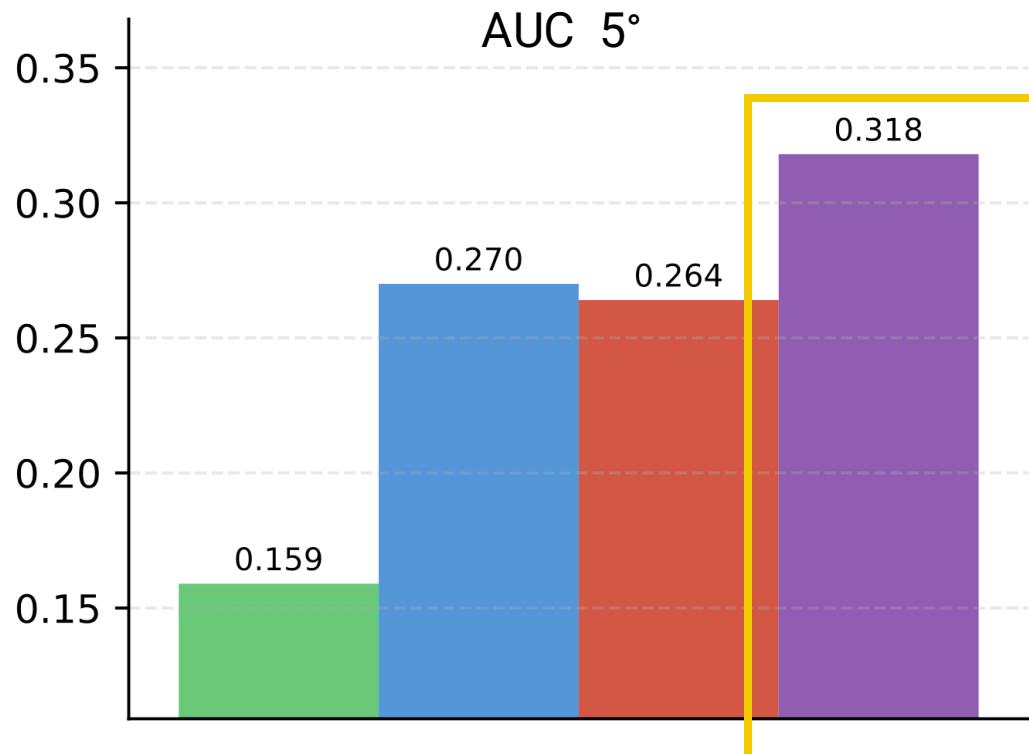
Evaluation on ScanNet

MASt3R

RoMa

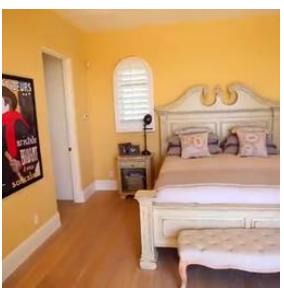
NoPoSplat (Trained on Re10k)

NoPoSplat (Trained on Re10k + DL3DV)



# High Quality Geometry

Input Images



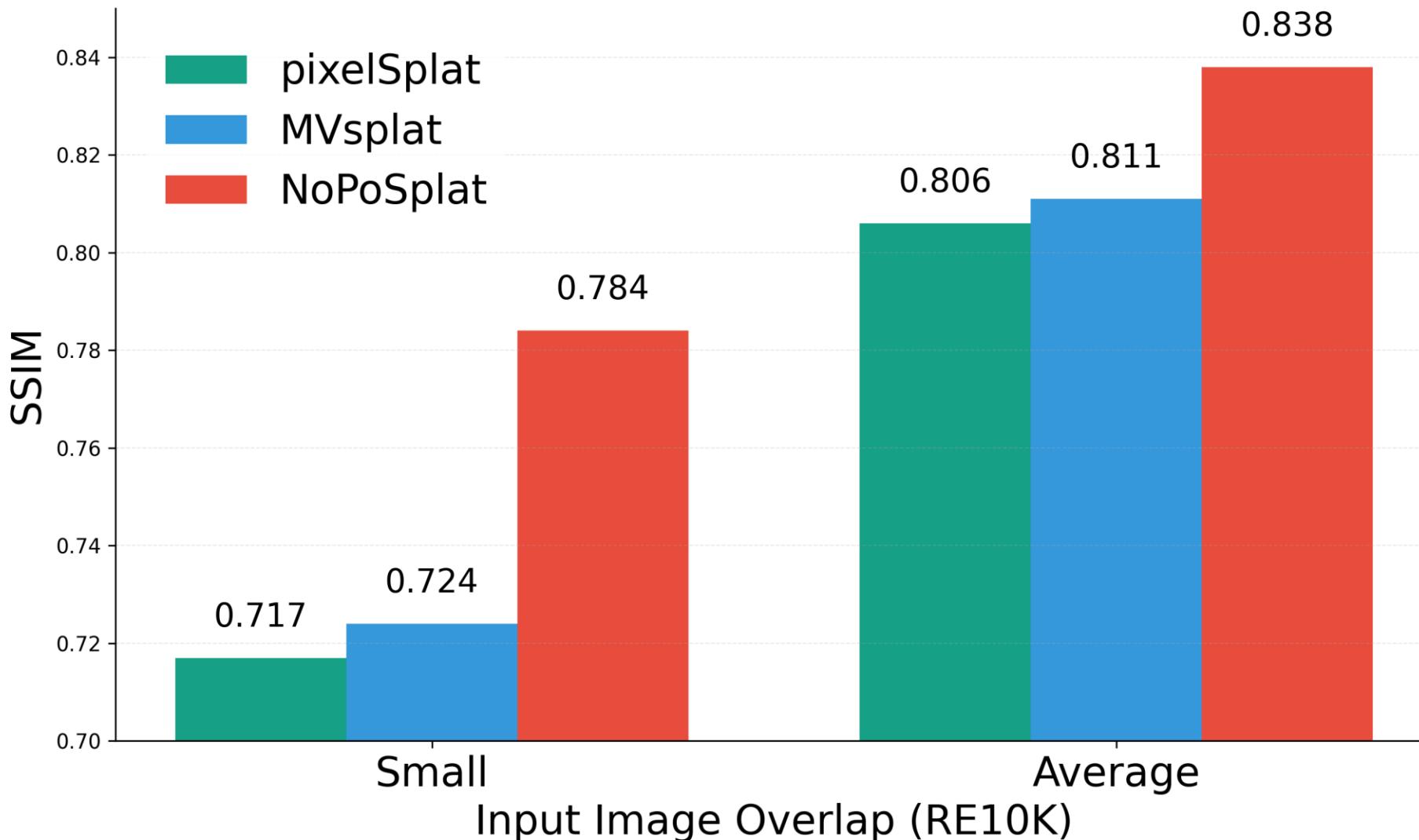
NoPoSplat (pose-free)



MVSplat (pose-required)

# Appearance Quality

Better even than pose-required methods!



# Appearance Quality

Input Views



MVSplat



NoPoSplat



# Cross-Dataset Generalization

RE10K → DTU

Input Views

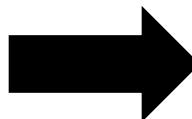


RE10K → ScanNet++



# In-the-Wild Data

Images extracted from OpenAI Sora

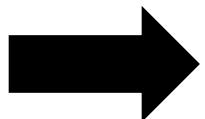
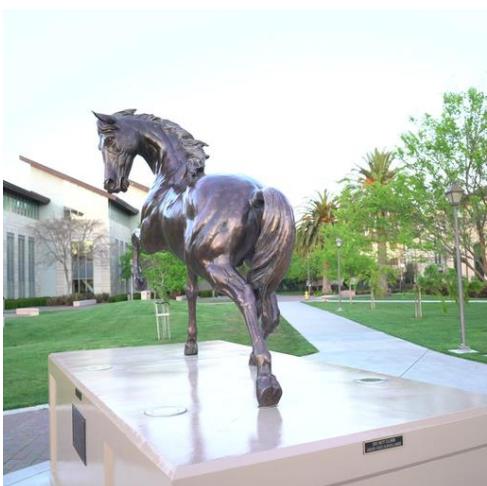


Input Images

Novel Views

# In-the-Wild Data

## Images from Tanks & Temples

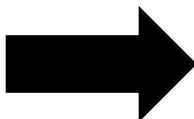


Input Images

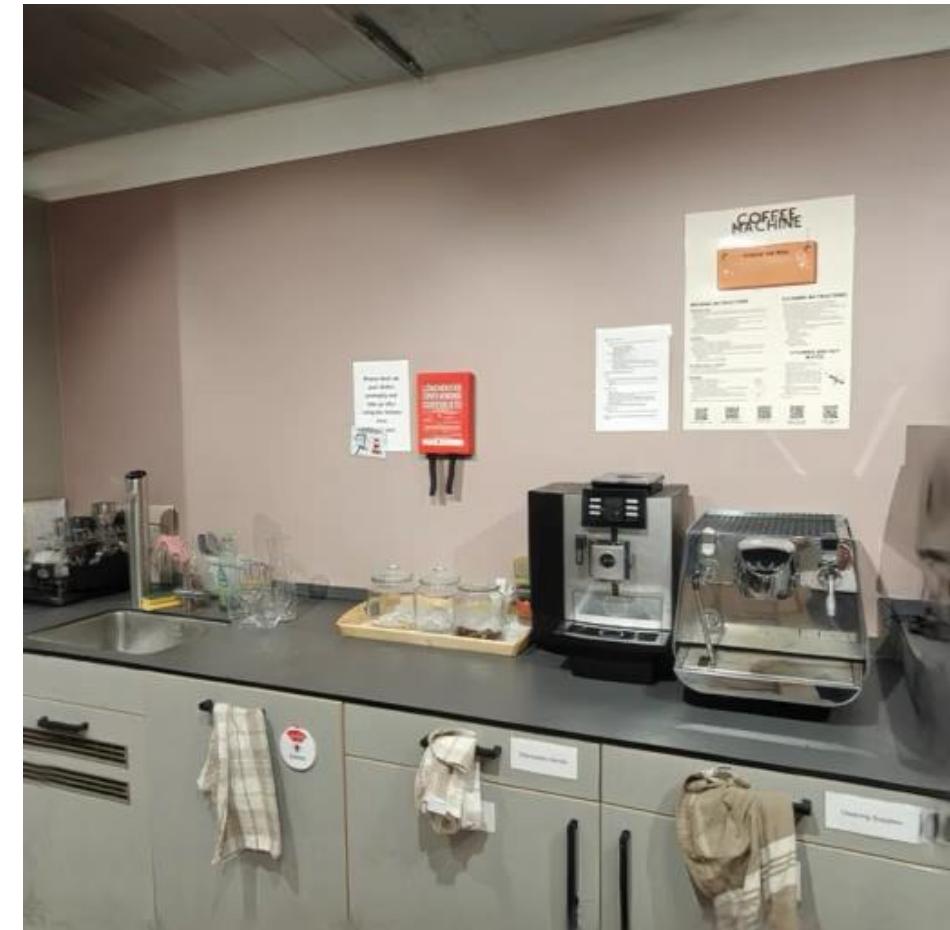
Novel Views

# In-the-Wild Data

## iPhone images



Input Images



Novel Views

# Take-home Messages

- Feedforward NVS can be surprisingly simple!
- Side product: SoTA relative pose estimation
- Works well with any static scenes **Not practical enough!**



THE  
PARISIAN  
DEPARTMENT  
STORE

TOUR EIFFEL · CHAMP DE MARS · MUSÉE DU LOUVRE · NOTRE-DAME · MUSÉE D'ORSAY · OPÉRA GARNIER · CHAMPS-ÉLYSÉES · GRAND PALAIS · TROCADÉRO  
**BIGBUS PARIS · LES CARS ROUGES**

**BIG  
BUS**

HOP-ON HOP-O

# Motivation



How to obtain **distractor-free 3D reconstruction**  
from **casually captured & long** image sequences **in the wild?**

# An Ideal 3D Reconstruction Pipeline

Instant, Pose-Free, Real-World 3D Everywhere



Feedforward



Pose-Agnostic



Dynamic



Fast Rendering



Arbitrary Lengths



Lighting-Robust

# An Ideal 3D Reconstruction Pipeline

Instant, Pose-Free, Real-World 3D Everywhere



Feedforward



Pose-Agnostic



Dynamic



Fast Rendering



Arbitrary Lengths

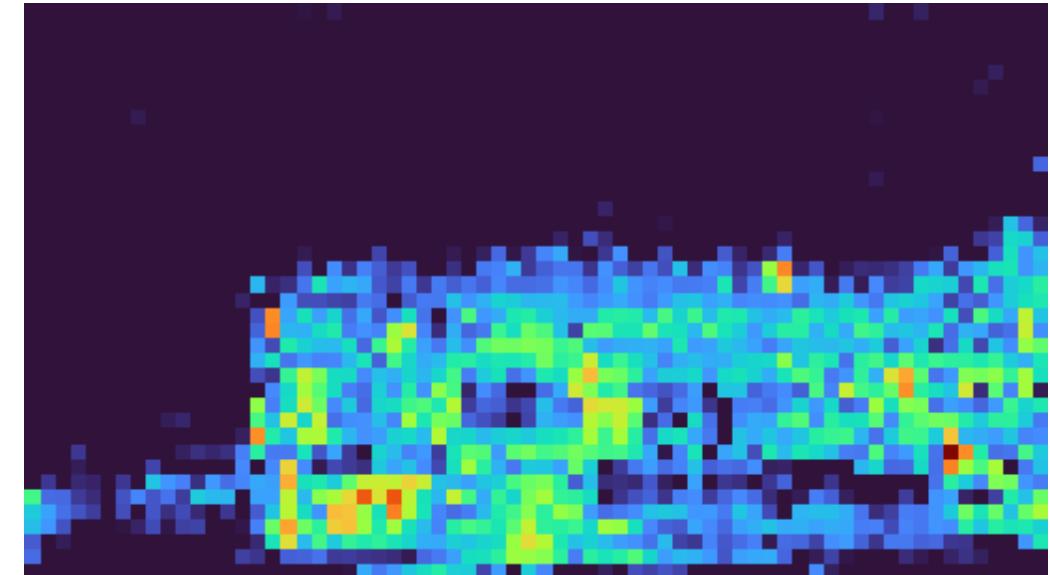


Lighting-Robust

# Uncertainty



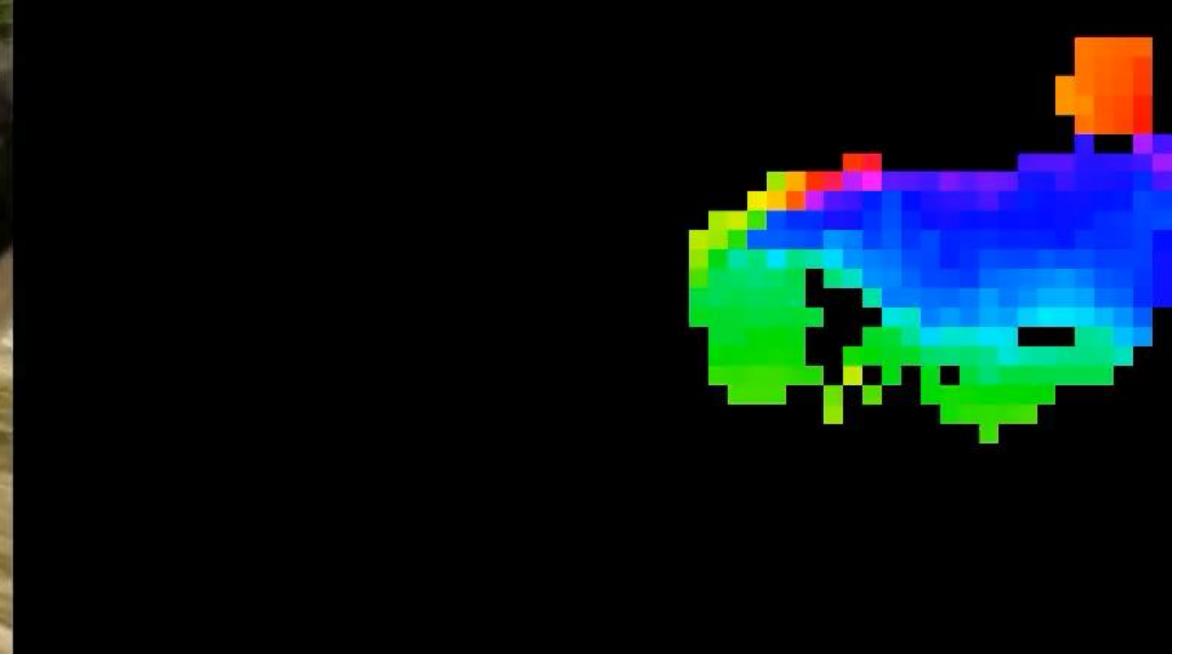
Input RGB



Uncertainty Map

How to learn a good uncertainty map?

# DINO v2



- A 2D foundation model producing **universal features**
- Preserve temporal-spatial consistency

How to leverage the DINO v2 to model uncertainty for  
3D reconstruction?



# NeRF *On-the-go*

Exploiting Uncertainty for Distractor-free NeRFs  
in the Wild



Weining  
Ren\*



Zihan  
Zhu\*

CVPR 2024



Boyang  
Sun



Julia  
Chen



Marc  
Pollefeys



Songyou  
Peng

# Pipeline

DINOv2 Feature Map



DINOv2



RGB Input

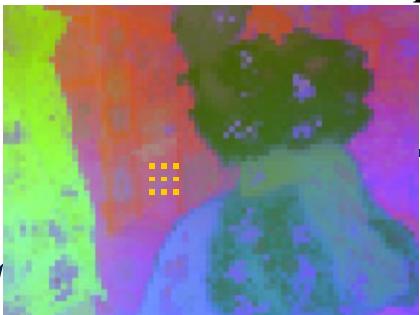
$G$   
Uncertainty  
MLP

$$\beta(r)$$

$C(r)$   
NeRF  
Representation

# Pipeline

DINOv2 Feature Map



DINOv2



RGB Input

$G$   
Uncertainty  
MLP

$\beta(r)$

$C(r)$

SSIM

$\mathcal{L}_{\text{uncer}}$

NeRF  
Representation

$$\mathcal{L}_{\text{uncer}}(r) = \frac{\mathcal{L}_{\text{SSIM}}}{2\beta(r)^2} + \lambda_1 \log \beta(r)$$

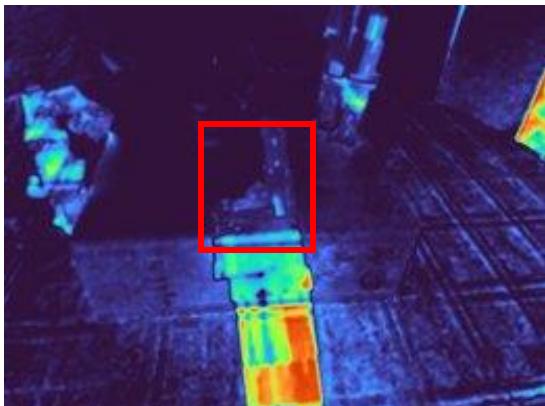
# To Learn the Uncertainty MLP...



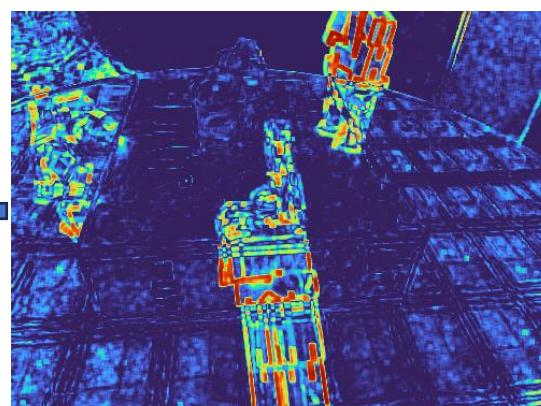
Rendered RGB



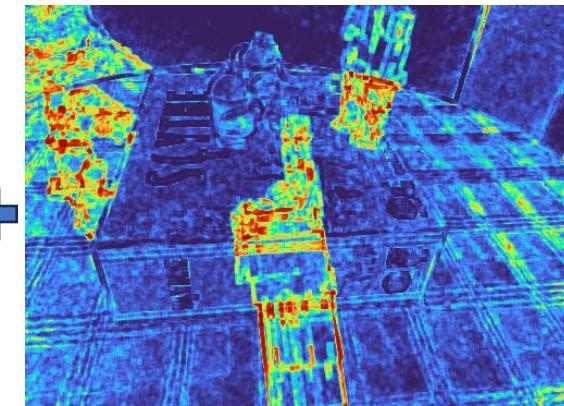
Train RGB



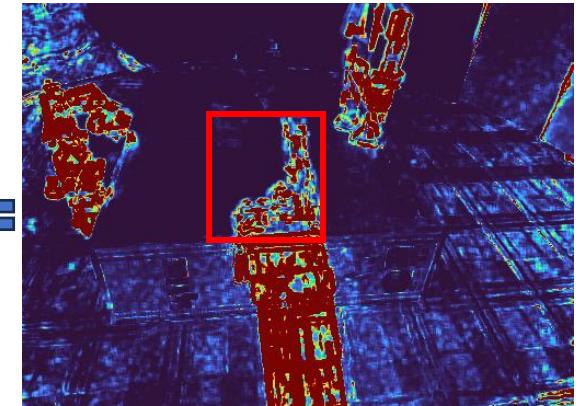
Luminance



Contrast



Structure



SSIM Error

## Why SSIM?

Leverage structure information when RGB is similar!

# Pipeline

DINOv2 Feature Map



DINOv2



RGB Input

$G$   
Uncertainty  
MLP

$\beta(r)$

$C(r)$

SSIM

NeRF  
Representation

$\mathcal{L}_{\text{uncer}}$

$\mathcal{L}_{\text{nerf}}$

$$\mathcal{L}_{\text{uncer}}(r) = \frac{\mathcal{L}_{\text{SSIM}}}{2\beta(r)^2} + \lambda_1 \log \beta(r)$$

$$\mathcal{L}_{\text{nerf}}(r) = \frac{\|C(r) - \hat{C}(r)\|^2}{2\beta^2(r)}$$

# **Results**



Train Station - Input Images



Train Station - Rendering Comparisons

Occlusion  
Ratio: **High**



Patio-High - Input



NeRF On-the-go  
(Ours)

Patio-High - Rendering Comparisons

# Analysis

# Analysis - Efficiency



25K

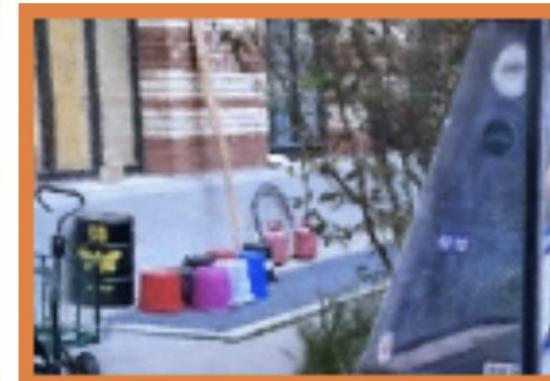
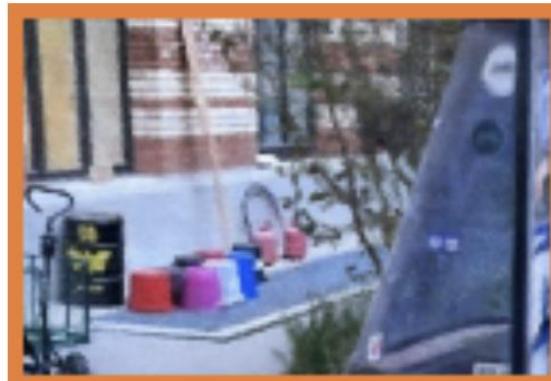
50K

100K

250K

**RobustNeRF**

# Analysis - Efficiency



25K

50K

100K

250K

**NeRF *On-the-go***  
**(Ours)**

# Analysis – Static Scene



0.447



0.376



0.374



RobustNeRF



Ours



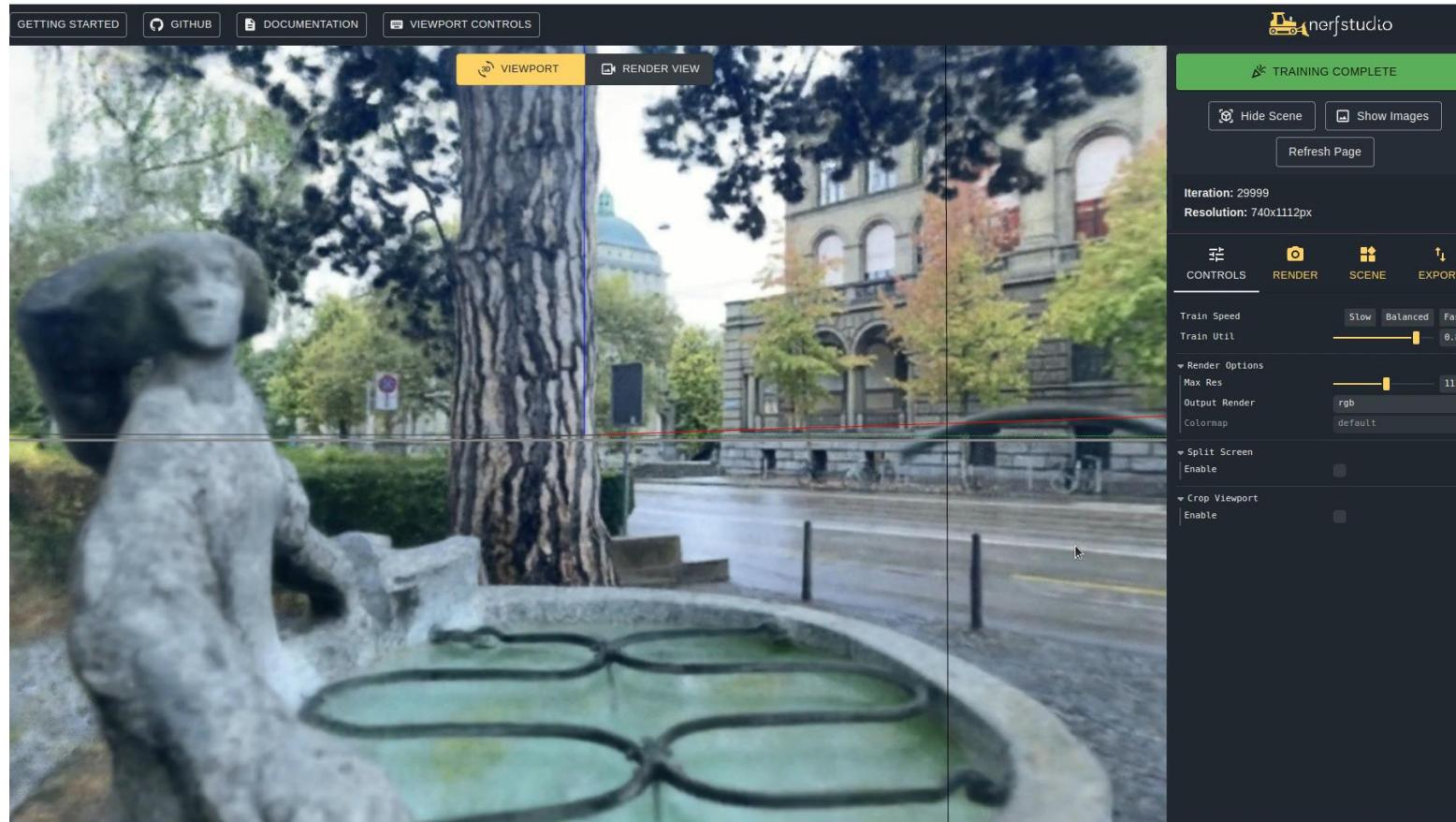
MipNeRF 360



GT

# Take-home Messages

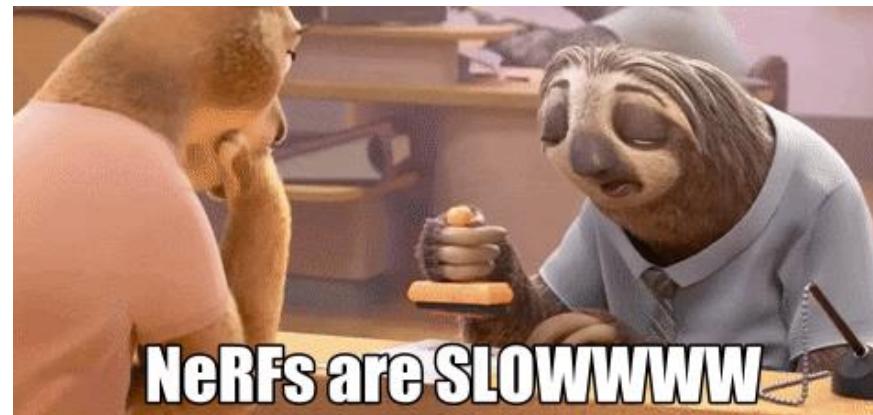
- ***On-the-go*** module is plug-and-play for all NeRF methods
  - Integrated into NeRFStudio



# Take-home Messages

- ***On-the-go*** module is plug-and-play for all NeRF methods
  - Integrated into NeRFStudio
- **2D foundation model** (DINOv2) rocks!

However, it is VERY SLOW



# In-the-Wild



# An Ideal 3D Reconstruction Pipeline

Instant, Pose-Free, Real-World 3D Everywhere



Feedforward



Pose-Agnostic



Dynamic



Fast Rendering



Arbitrary Lengths



Lighting-Robust

# An Ideal 3D Reconstruction Pipeline

Instant, Pose-Free, Real-World 3D Everywhere



Feedforward



Pose-Agnostic



Dynamic



Fast Rendering



Arbitrary Lengths



Lighting-Robust

# 3DGS

from **casually captured** images **in the wild**

- 😊 Robustly handle **arbitrary occlusions**
- 😊 Model any **illumination changes**
- 😊 **Real-time** rendering!



# WildGaussians

## 3D Gaussian Splatting In the Wild

NeurIPS 2024



Jonas  
Kulhanek



Songyou  
Peng



Zuzana  
Kukelova



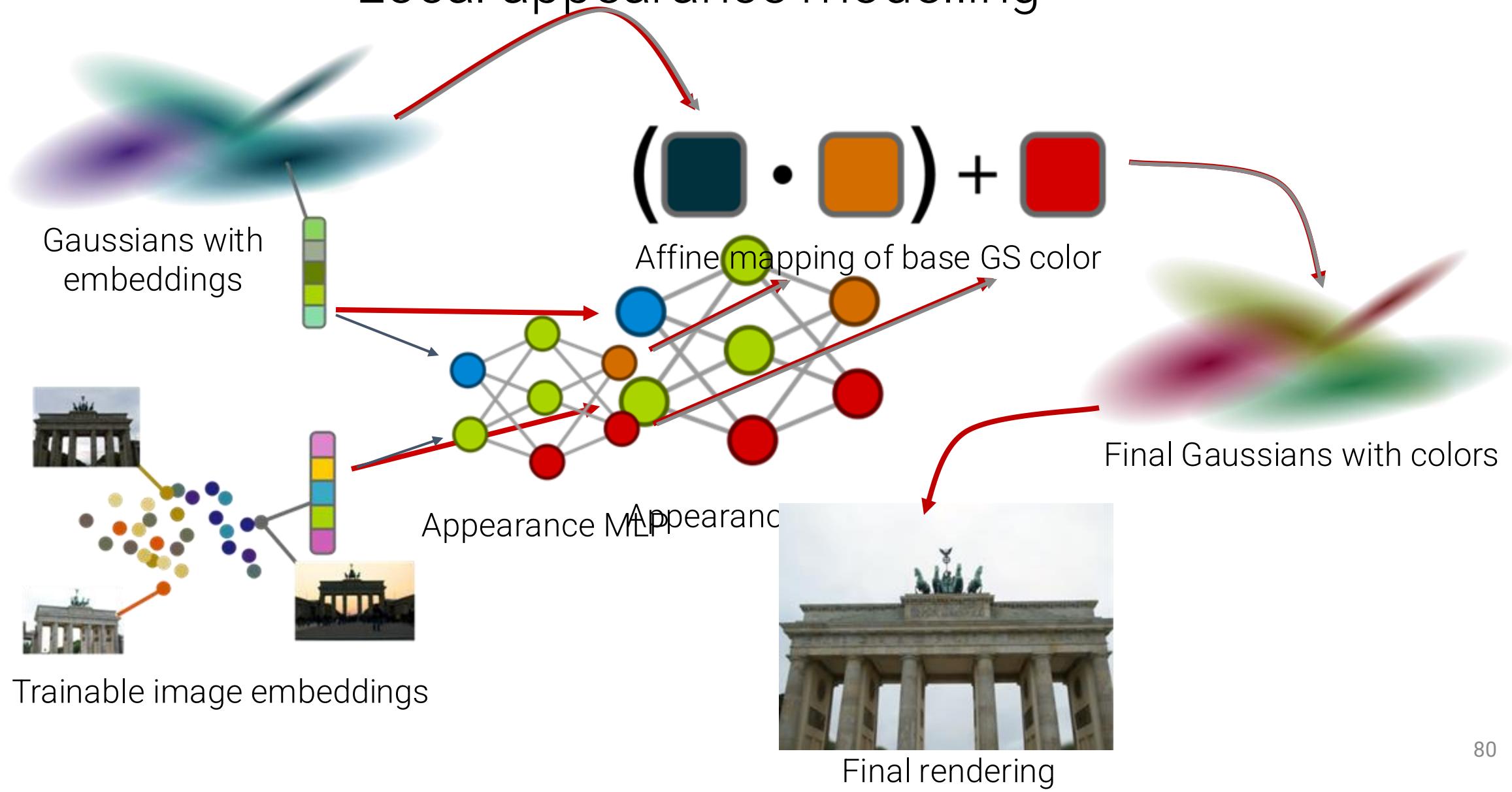
Marc  
Pollefeys



Torsten  
Sattler

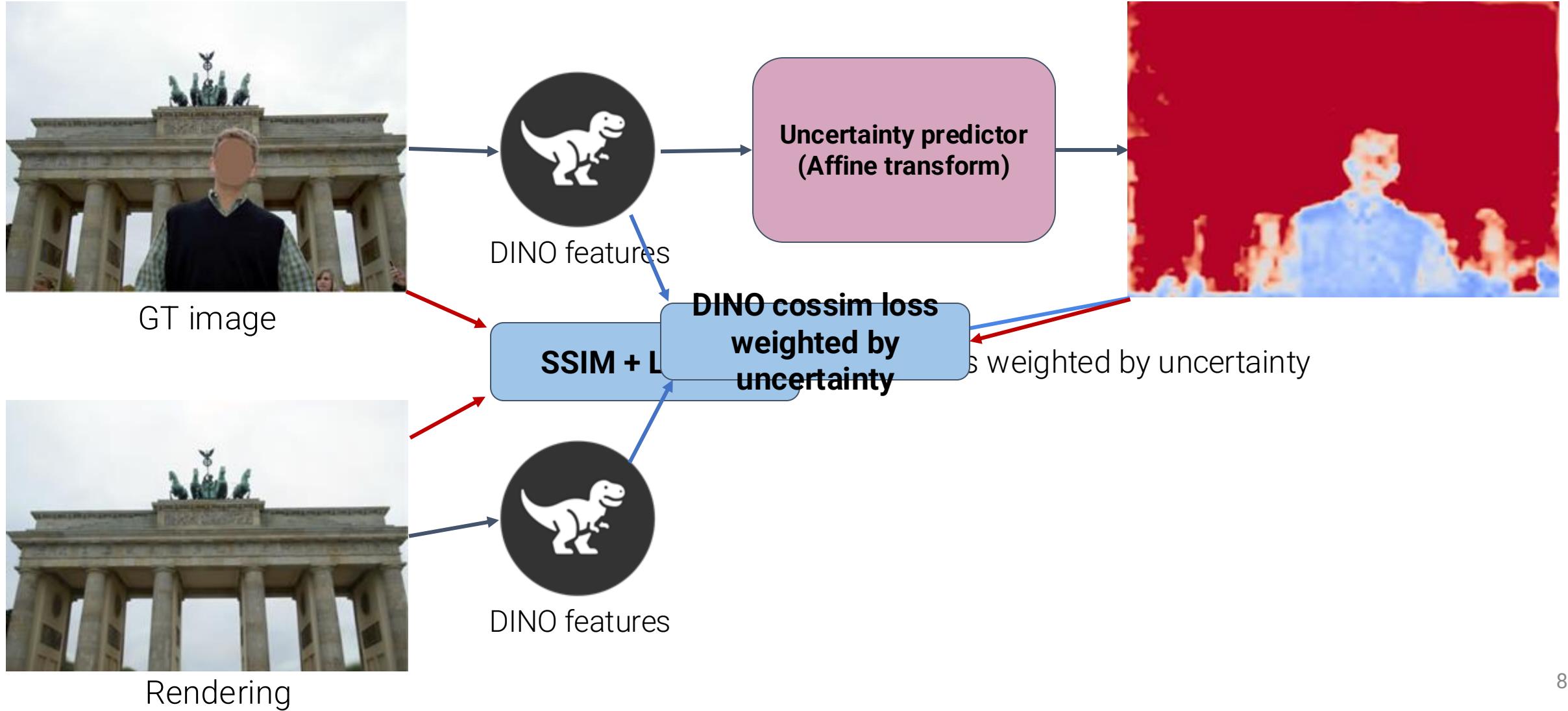
# Pipeline

## Local appearance modelling



# Pipeline

Optimization with certainty



# Phototourism Dataset

Wild Gaussians



# K-Planes

FPS 0.016

# Wild Gaussians

FPS 31

Rendering Speed



# An Ideal 3D Reconstruction Pipeline

Instant, Pose-Free, Real-World 3D Everywhere



Feedforward



Pose-Agnostic



Dynamic



Fast Rendering



Arbitrary Lengths



Lighting-Robust

# An Ideal 3D Reconstruction Pipeline

Instant, Pose-Free, Real-World 3D Everywhere



Feedforward



Pose-Agnostic



Dynamic



Fast Rendering



Arbitrary Lengths



Lighting-Robust



# WildGS-SLAM

## Monocular Gaussian Splatting SLAM in Dynamic Environments

CVPR 2025



Jianhao  
Zheng\*



Zihan  
Zhu\*



Valentin  
Bieri



Marc  
Pollefeys



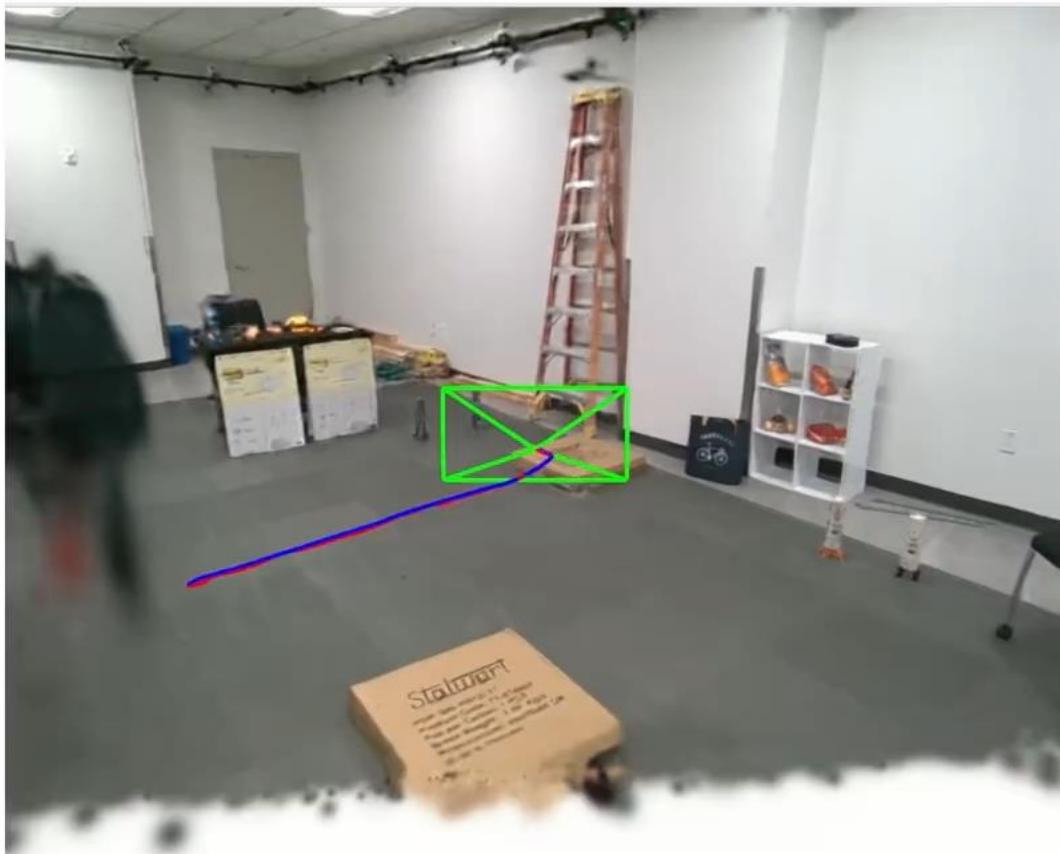
Songyou  
Peng



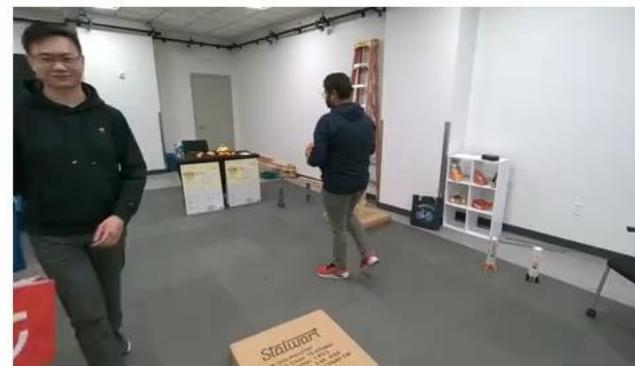
Iro  
Armeni

# Online mapping

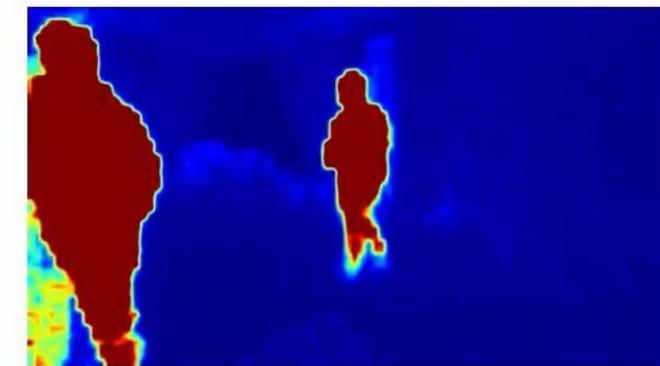
— GT traj — Est. traj



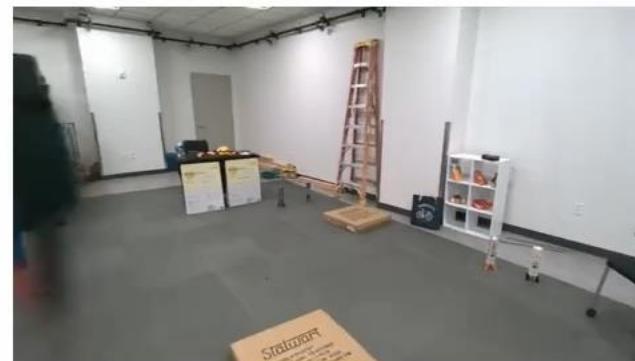
Input: RGB frames



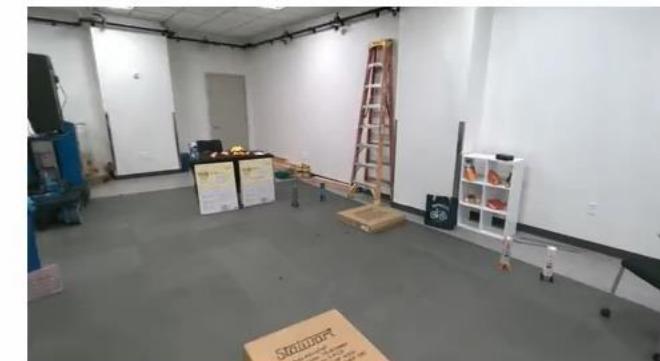
Uncertainty



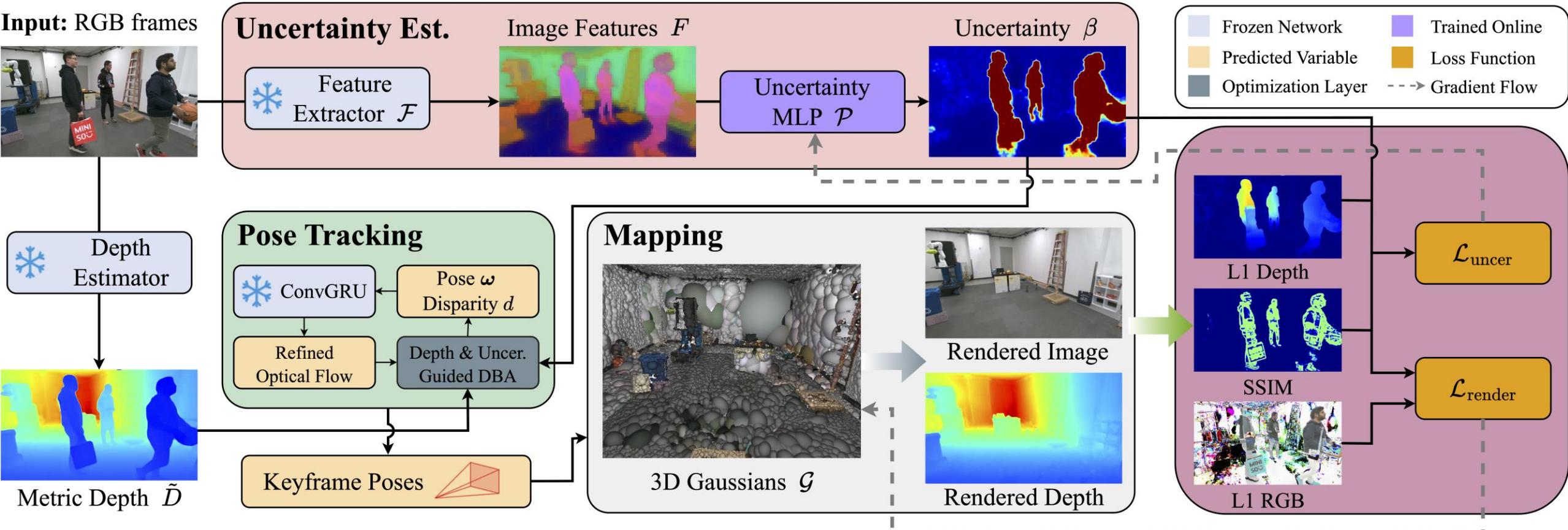
Rendered (online)



Rendered (final)



# WildGS-SLAM



# Results

## Umbrella



Input Frames



Traj. Error Colormap



MonoGS [CVPR' 24]



Splat-SLAM [arXiv' 24]



WildGS-SLAM (Ours)



# Results

Tower



Input Frames



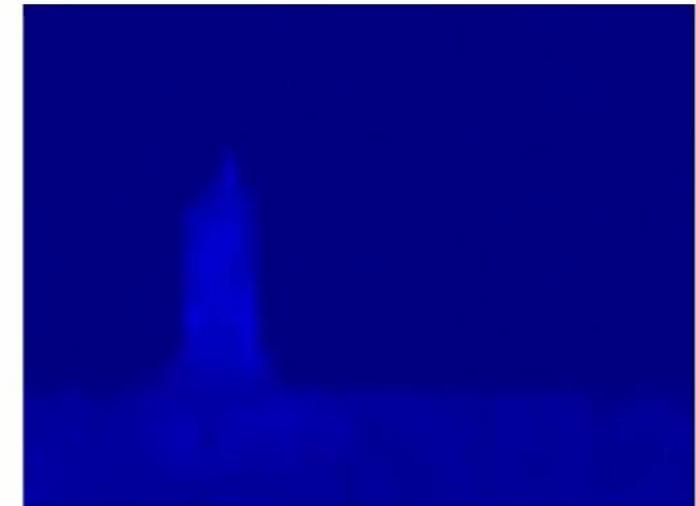
Rendered (MonoGS)



Rendered (Splat-SLAM)



Rendered (Ours)



Uncertainty (Ours)

# An Ideal 3D Reconstruction Pipeline

Instant, Pose-Free, Real-World 3D Everywhere



Feedforward



Pose-Agnostic



Dynamic



Fast Rendering



Arbitrary Lengths



Lighting-Robust

# Opinion

3D static reconstruction is almost at the last mile

128 Views



3.92 s



**VGGT**  
[CVPR'25 Oral]

# Opinion

3D static reconstruction is almost at the last mile

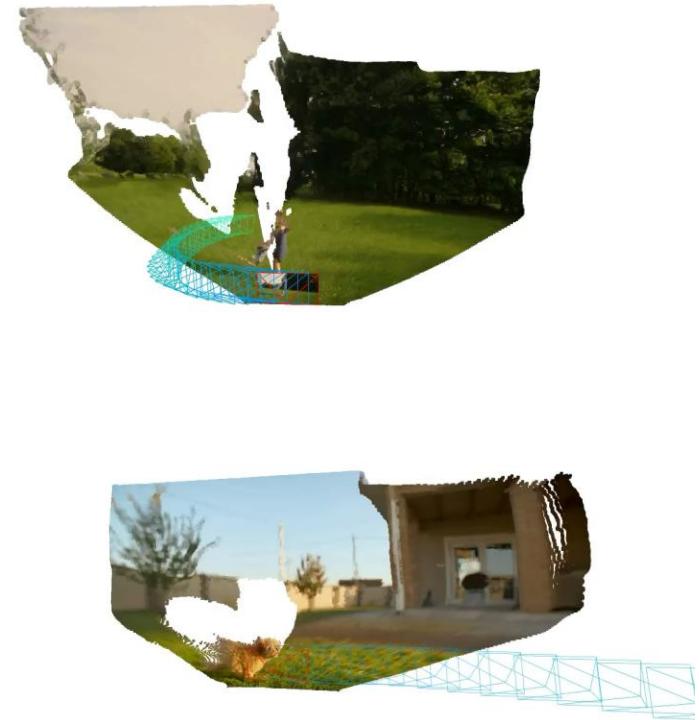
128 Views



3.92 s



**VGGT**  
[CVPR'25 Oral]



**MegaSaM**  
[CVPR'25 Oral]

# Challenges & Opportunities

How to robustly handle long sequences?



Static Videos



Dynamic Videos

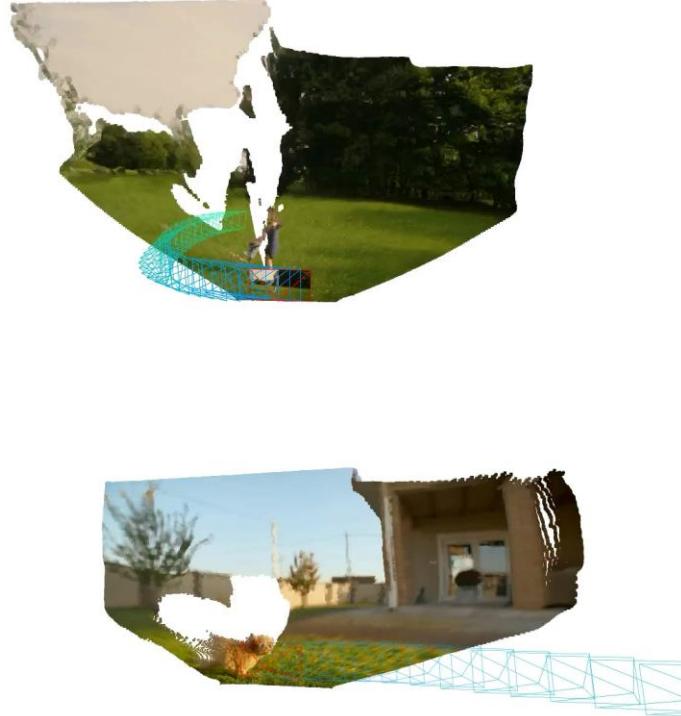


Unstructured Photo Collections

CUT3R  
[CVPR'25 Oral]

# Challenges & Opportunities

4D (3D dynamic) reconstruction is still very hard, especially for NVS



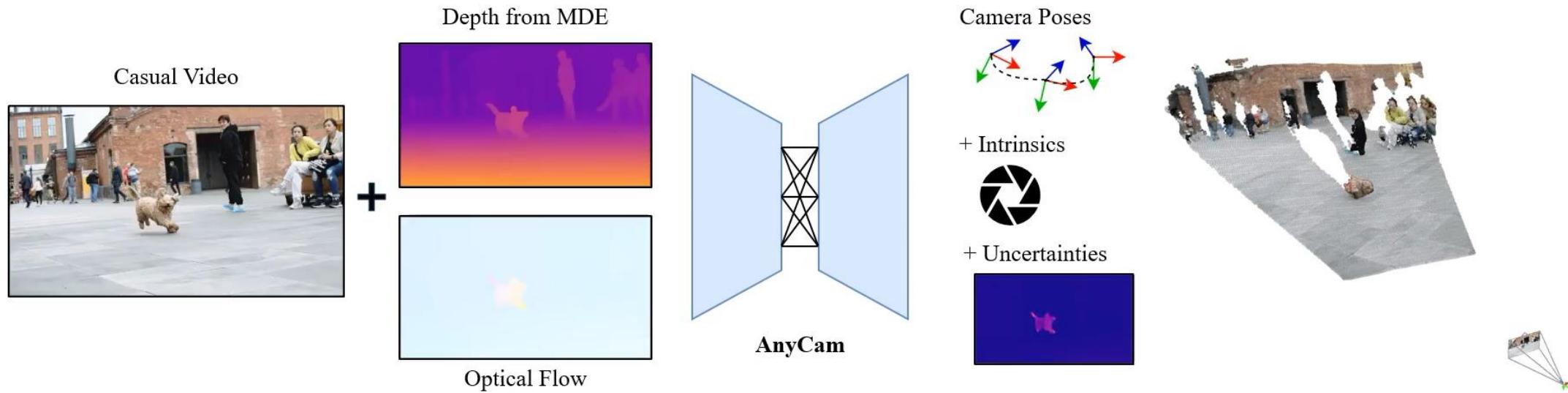
**MegaSaM**  
[CVPR'25 Oral]



**Gaussian-Flow**  
[CVPR'24 Highlight]

# Challenges & Opportunities

Feedforward pose estimation is not solved



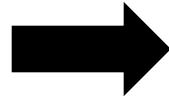
**AnyCam**  
[CVPR'25]

# Challenges & Opportunities

How to continuously update your 3D reconstruction?



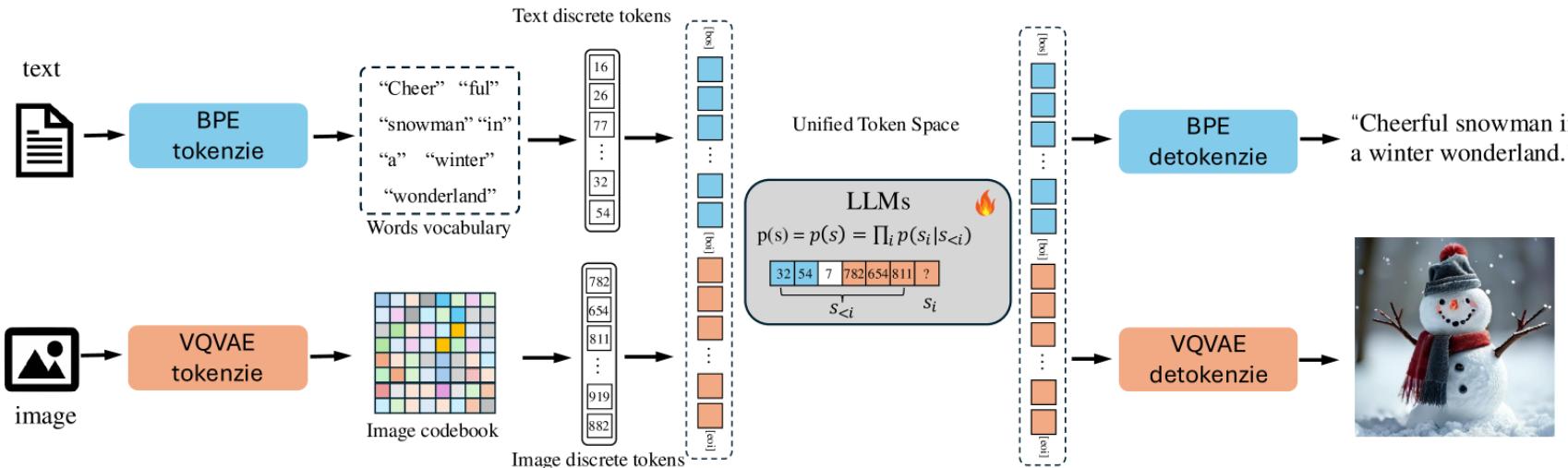
1 PM



7 PM

# Challenges & Opportunities

Can we unify 3D reconstruction/SLAM into LLM?



Liquid  
[arXiv'24]

# Challenges & Opportunities

Interaction with 3D scenes at speed!



World Labs Demo

“People overestimate what they can do in one year, and underestimate what they can do in ten years.”

--- *Bill Gates*

“People overestimate what they can do in one year, and underestimate what they can do in ten years.”

--- Bill Gates

“People overestimate what they can do in one week, and underestimate what they can do in one year.”

--- Lars Mescheder



