

An aerial photograph of a city, likely Zurich, showing a river (Limmat) flowing through the center. The river is surrounded by dense urban development, including various residential and commercial buildings. A bridge crosses the river in the upper left. The foreground shows a mix of older stone buildings and more modern structures. The background features more greenery and distant cityscapes.

Stereo Video Depth Estimation in 2025 Research Proposal

3D Vision Group 18

Hepeng Fan, Qingrui Deng, Tong Su, Yixin Zhou

Supervisor: Haofei Xu

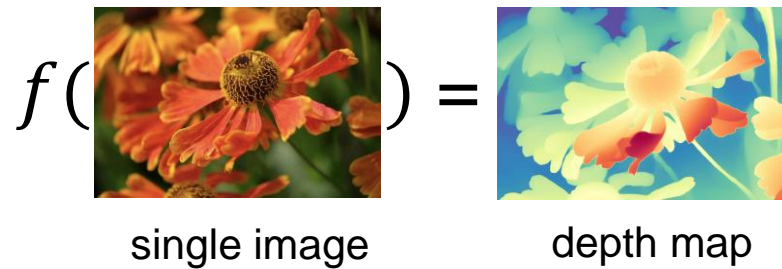
10.03.2025

Monocular Depth Estimation

- **Depth** refers to the distance from the **object point to the camera center**

- Two ways of describing depth: **Relative Depth** and **Metric Depth**

$$\text{relative depth} + \text{scale} = \text{metric depth}$$

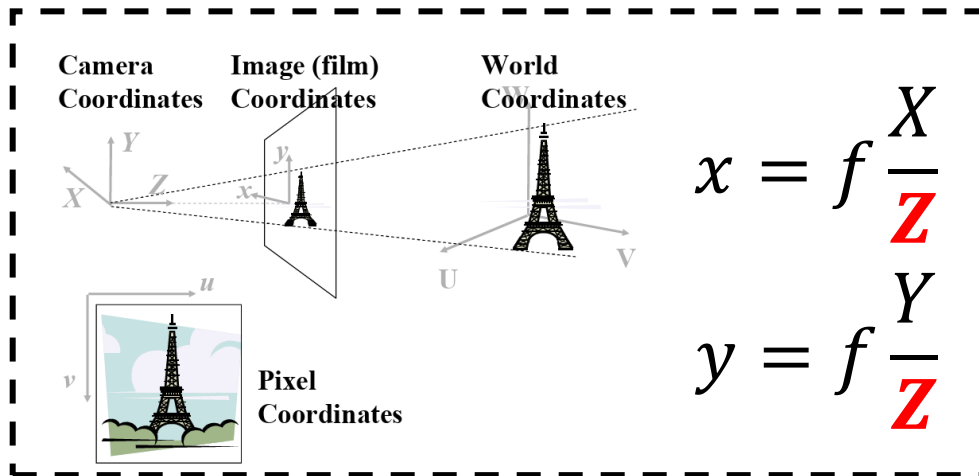


1.5	1.8	2.5	2.6
1.1	0.2	0.5	0.3
3.8	5.5	2.1	1.5

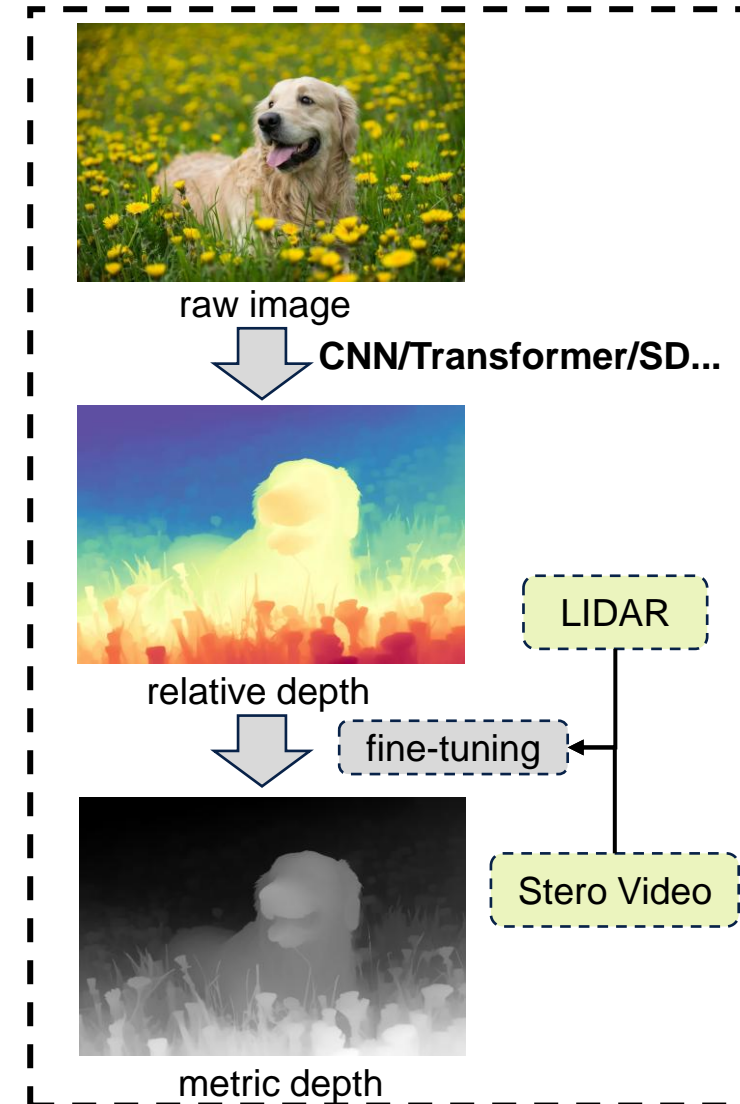
relative depth map

1.5m	1.8m	2.5m	2.6m
1.1m	0.2m	0.5m	0.3m
3.8m	5.5m	2.1m	1.5m

metric depth map



Pipeline for metric depth estimation





Raw Video



Image Depth Estimation
(Video Stitching)

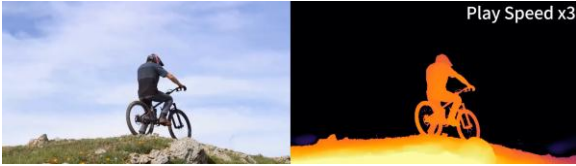


Consistent Video Depth Estimation

Our Project: Stereo **Video Metric** Depth Estimation

Baseline Model

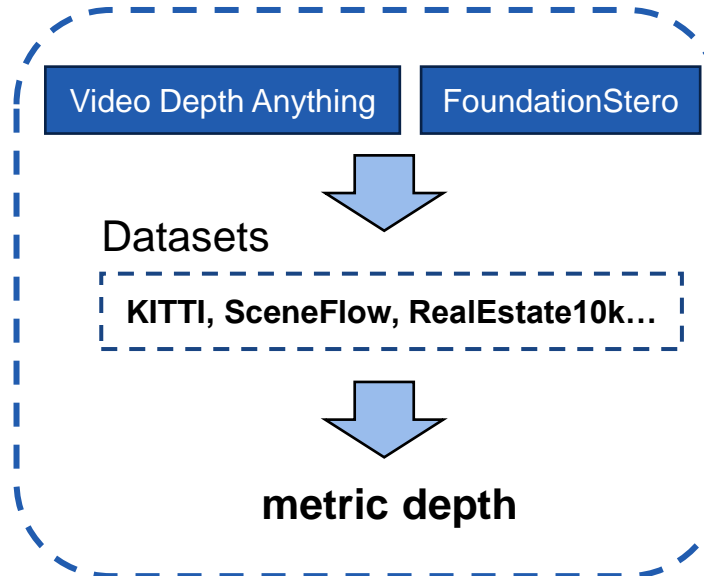
- Video Depth Anything



- FoundationStero



Model Design



Model Evaluation

Evaluation Metrics

- Abs.Rel
- End-Point Error
- Temporal Consistency Score

Comparion

SOTA Models: BIDA Stero, Dynamic Stero

Challenges and Potential Solutions

How to improve the temporal consistency of FoundationStero with the help of Video Depth Anything



Solution: Add one temporal-spatial layer to FoundationStero

Questions & Feedback