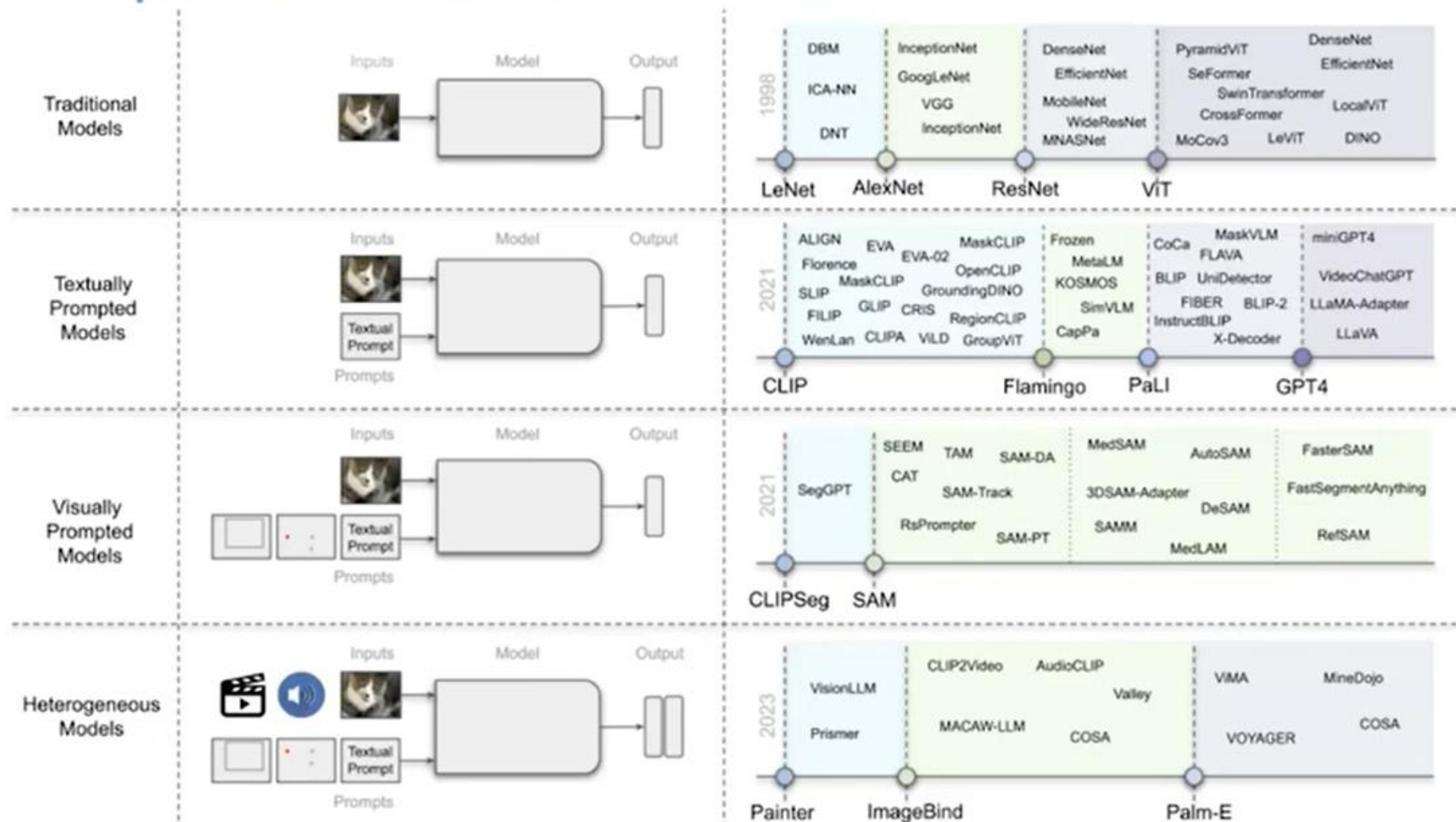


Computer Vision

Evolution of foundational model in computer vision



Hallucination



Generated Image

Shadow Errors

Detected Shadow Errors

Vanishing Point Errors

Detected Perspective Errors

Sarkar et al., Shadows Don't Lie and Lines Can't Bend! Generative Models don't know Projective Geometry...for now, CVPR 2024

Bias and fairness

- Biases, stereotypes, and prejudices related to race, underrepresented groups, minority cultures, and gender can make the models output biased predictions or exhibit skewed behavior



Figure 8: **Bias of CLIP.** We zero-shot classify a random COCO image with a male-looking and a female-looking individual, also adding a circle over each of them. We score the following sentences to the images: *This is an image of a {woman, man, missing person, suspected murderer}*. The apparent gender resolution is correct, but the circled images tend to be scored higher as missing or murderers. Blur added for privacy.

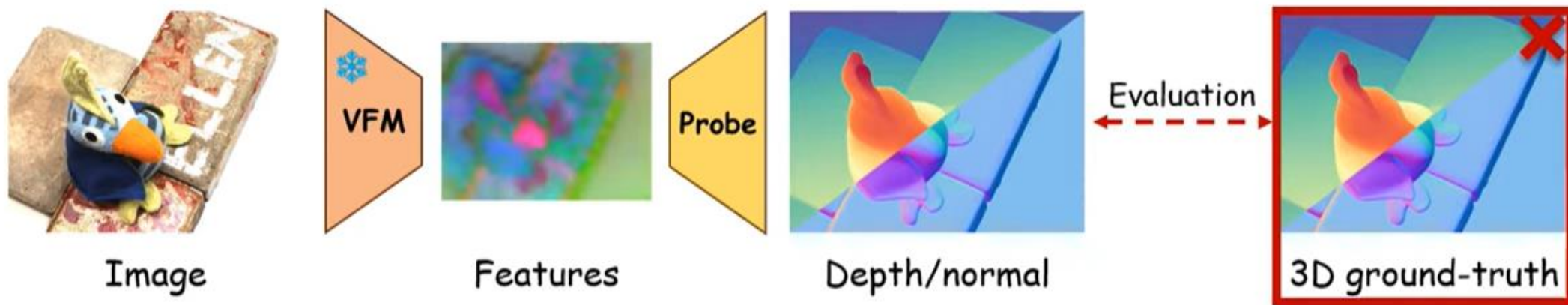


Visual Foundation Models

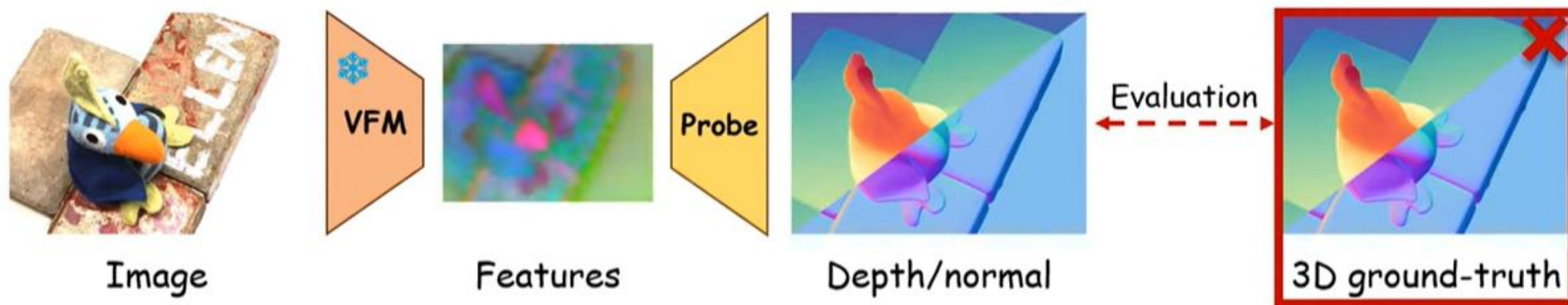
How well do they understand the 3D world?

3D Probing!

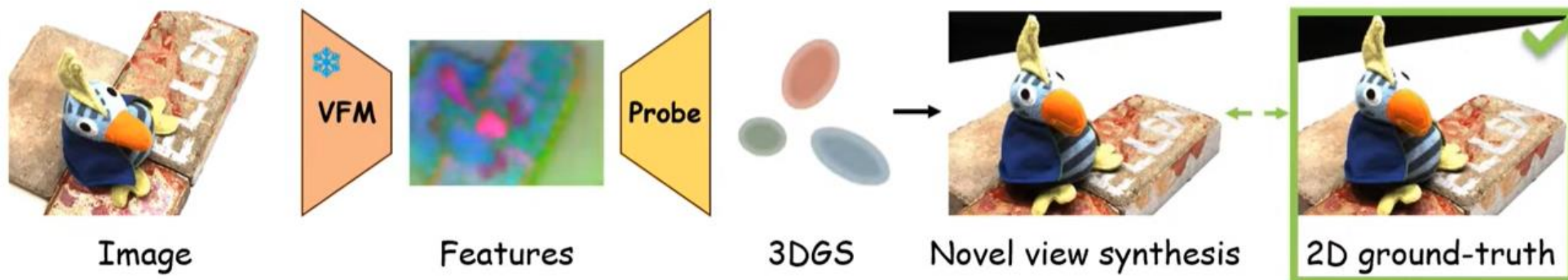
3D Probing



3D Probing



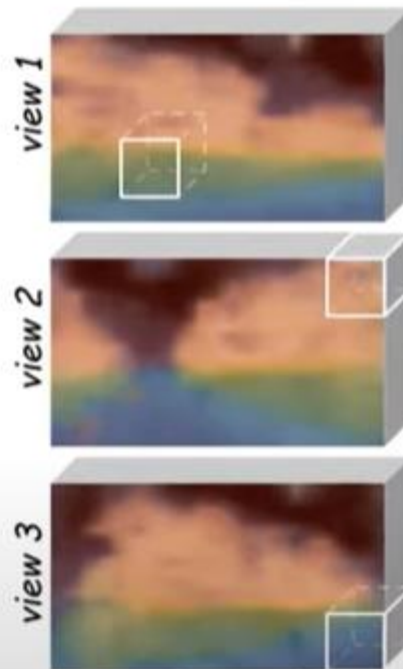
🔥 New Benchmark



Feat2GS as VFM Probe

Novel View Synthesis as a proxy task

Training Views



Readout



Readout

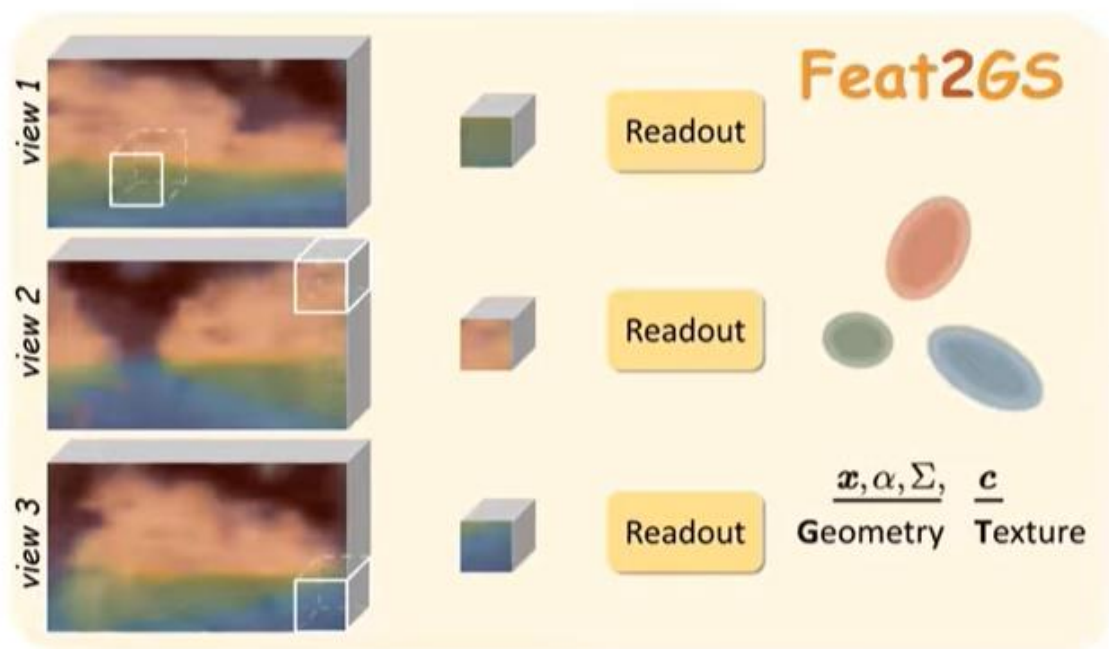


Readout



$\mathbf{x}, \alpha, \Sigma, \mathbf{c}$

3DGS parameters

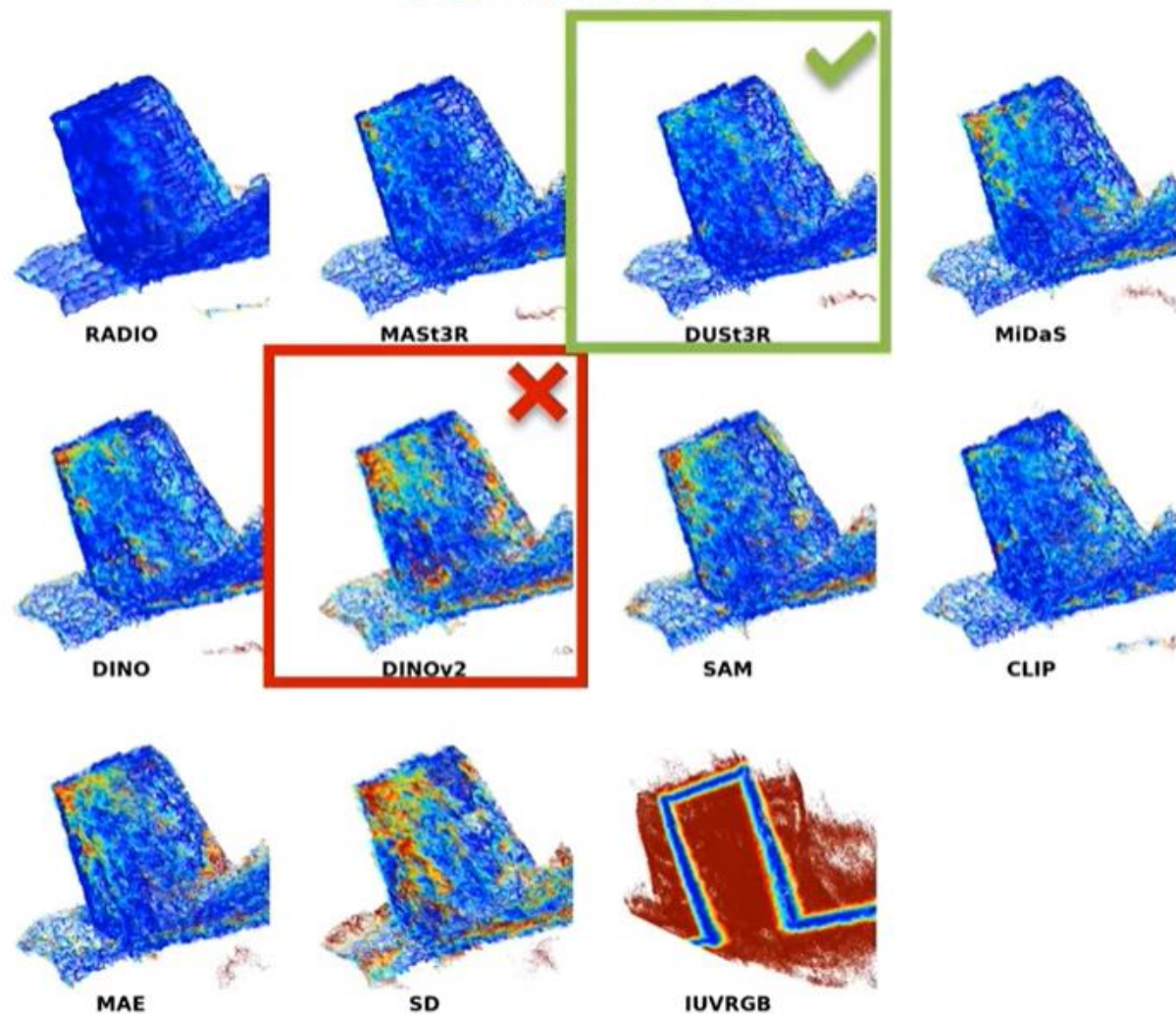


Splatting



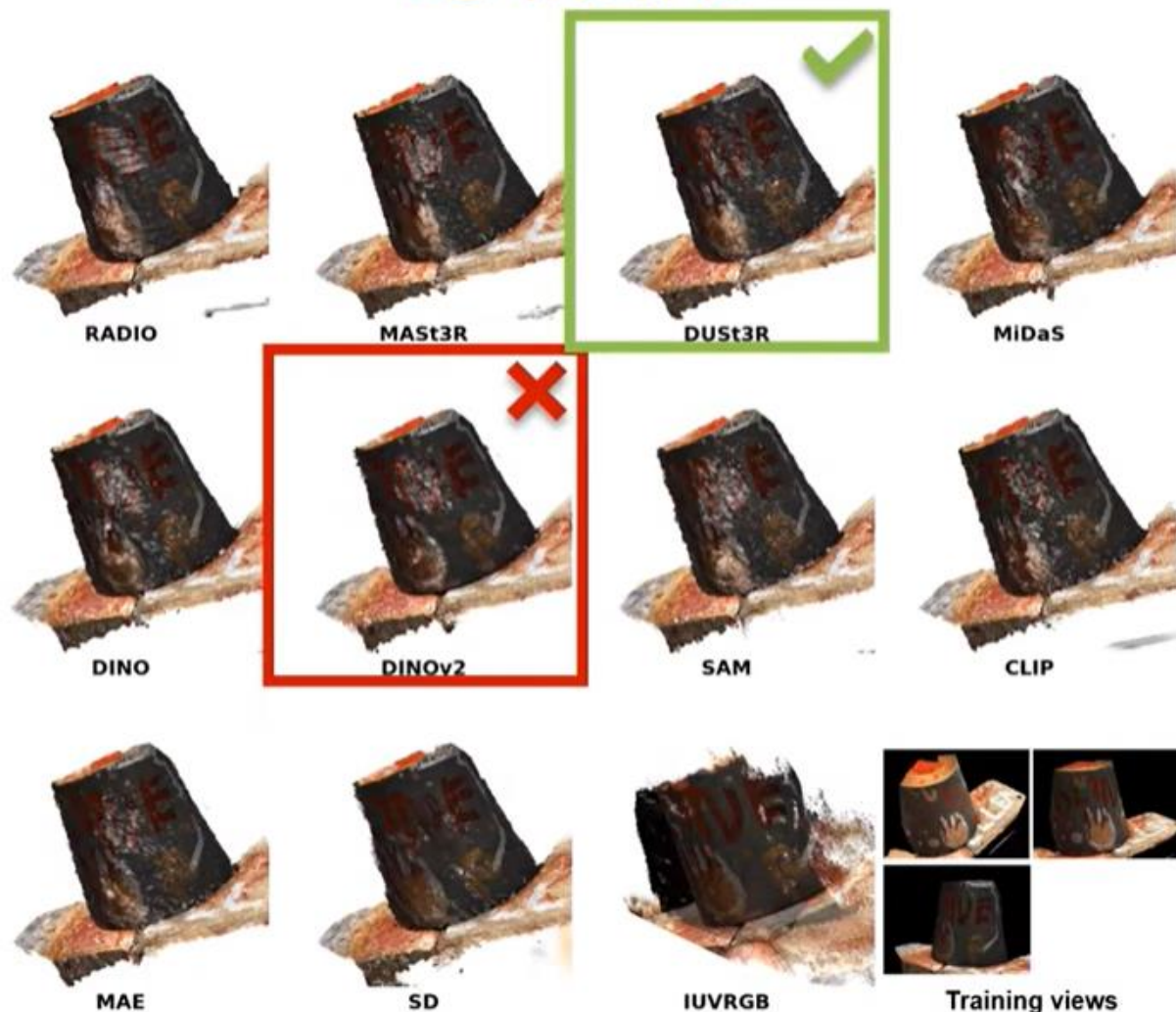
Novel View Synthesis

3D Metric



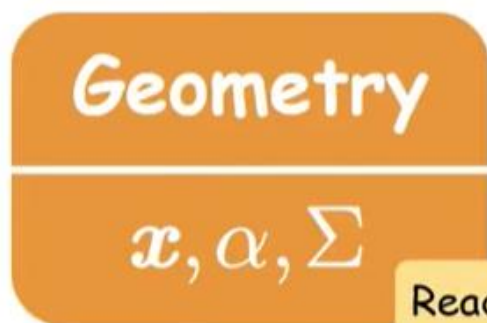
Pointcloud Error Map

2D Metric



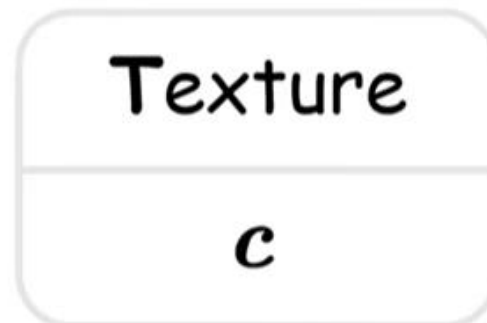
Novel View Synthesis

Geometry Probing



Readout

A small yellow rectangular label with the word 'Readout' in black text, positioned to the right of the bottom section of the 'Geometry' block.



Geometry Probing

Geometry

x, α, Σ

Readout



RADIO



MASt3R



DUST3R



MiDaS



DINO



DINOv2



SAM



CLIP



MAE



SD



IUVRGB



Training views



RADIO-Geometry



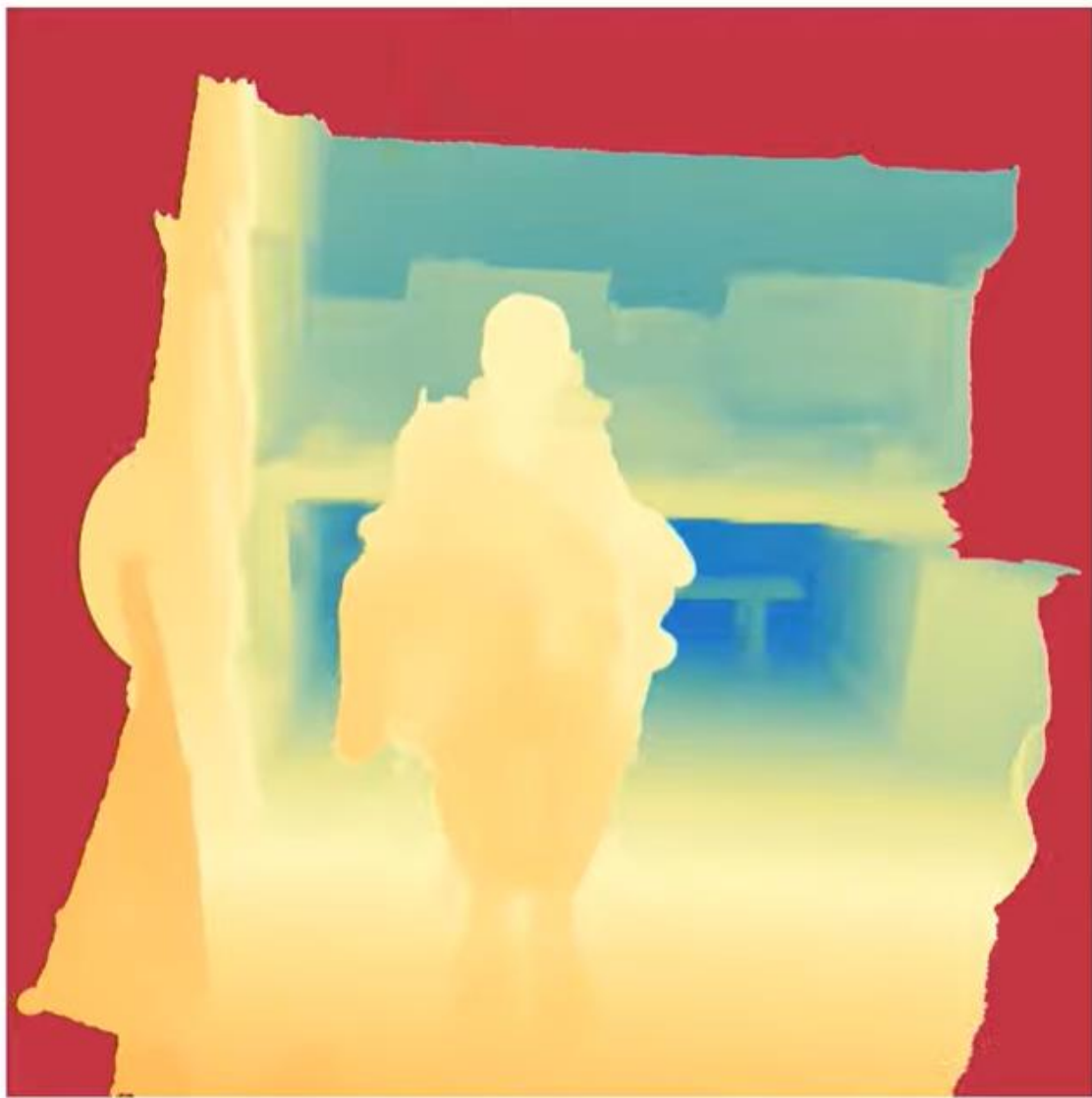
RADIO-All



Feat2GS-Geometry



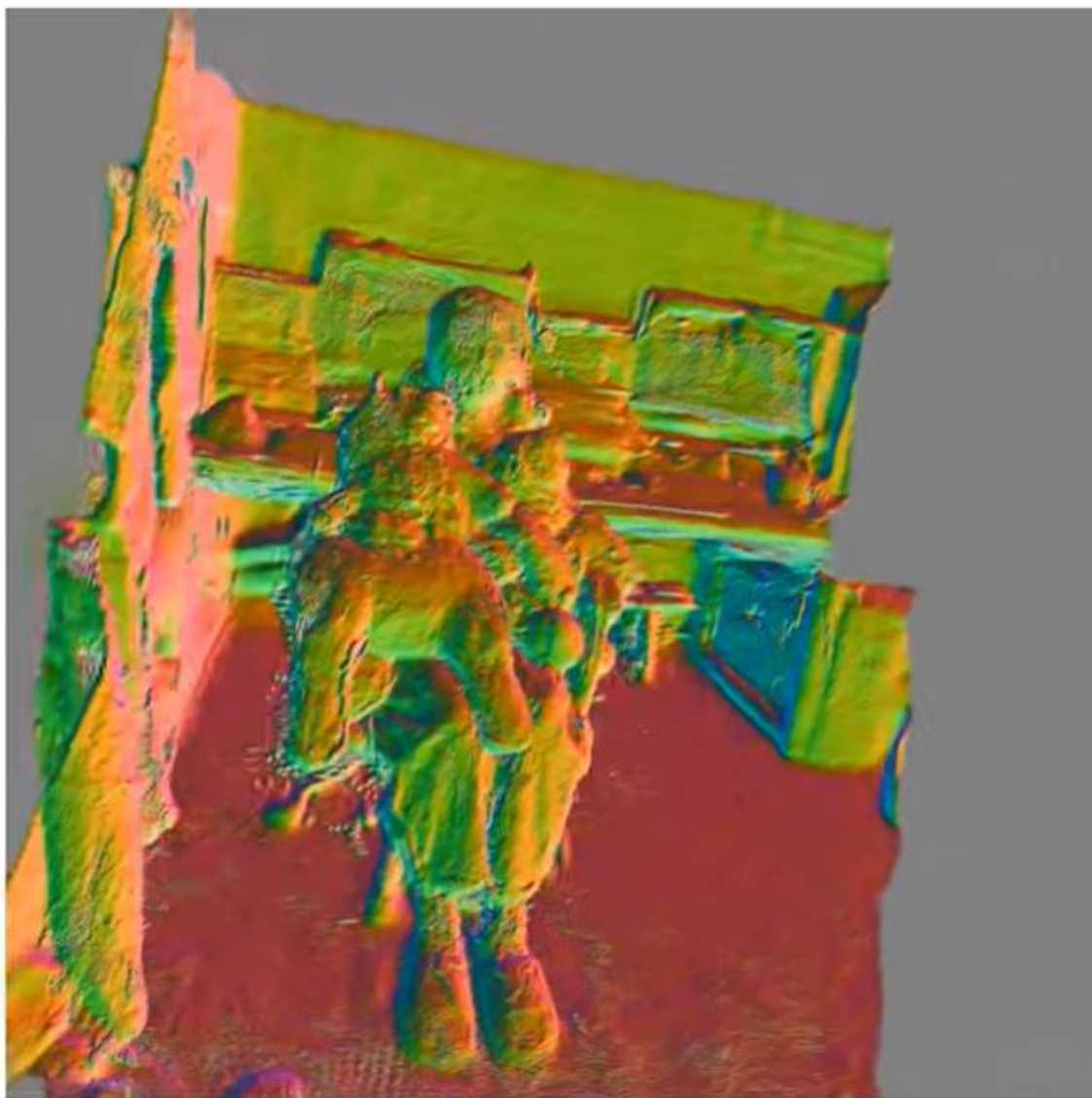
InstantSplat



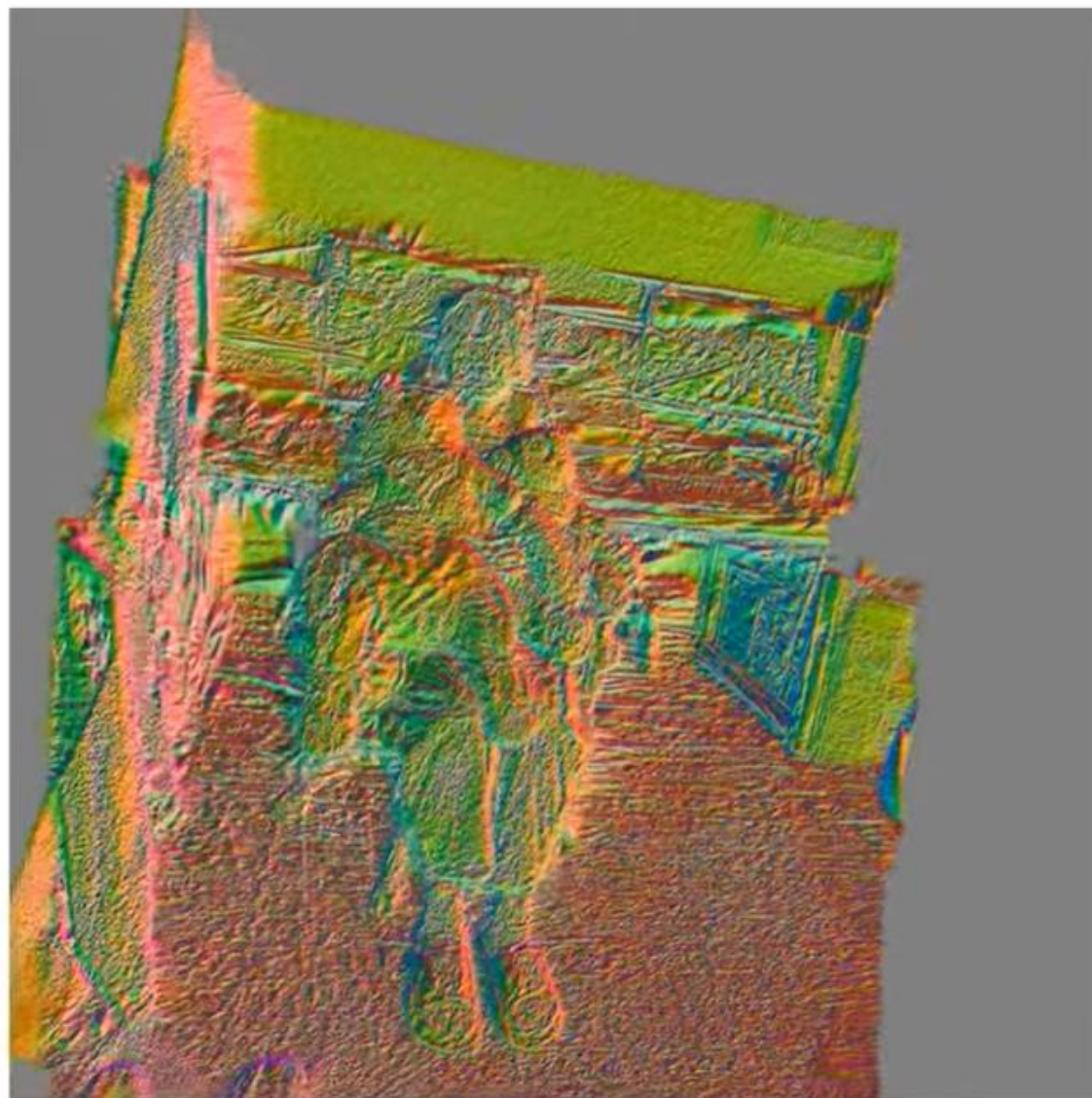
Feat2GS-Geometry



InstantSplat



Feat2GS-Geometry



InstantSplat

Feat2GS

Probing Visual Foundation Models with Gaussian Splatting



Yue Chen¹



Xingyu Chen¹



Anpei Chen^{1,3}



Gerard Pons-Moll^{3,4}



Yuliang Xiu^{1,2}

¹Westlake University

²Max Planck Institute for Intelligent Systems

³University of Tübingen, Tübingen AI Center

⁴Max Planck Institute for Informatics



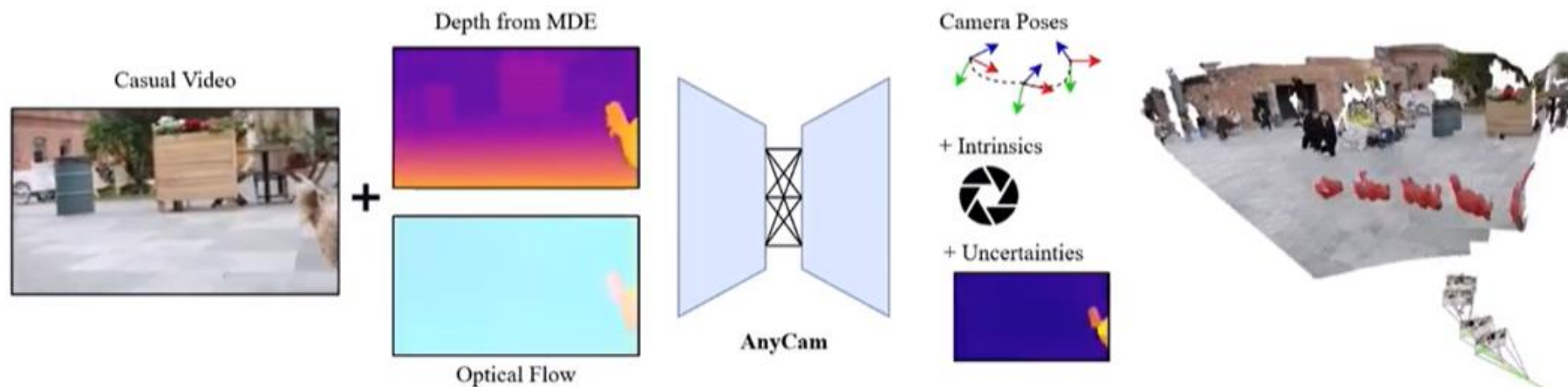
Tübingen AI Center
tuebingen.ai

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN



AnyCam: Learning to Recover Camera Poses and Intrinsic from Casual Videos

fwmb.github.io/anycam



Felix Wimbauer^{1,2,3} Weirong Chen^{1,2,3} Dominik Muhle^{1,2} Christian Rupprecht³ Daniel Cremers^{1,2}

¹Technical University of Munich ²MCML ³University of Oxford

AnyCam is a learning-based method that processes a **casual input video** and computes **camera poses** and **intrinsic**s.



Unlike other methods, ***AnyCam*** is **not** trained in a **supervised way**, but through **self-supervision on casual videos**.

AnyCam Training Datasets

RealEstate10K



WalkingTours



OpenDV



YouTube VOS



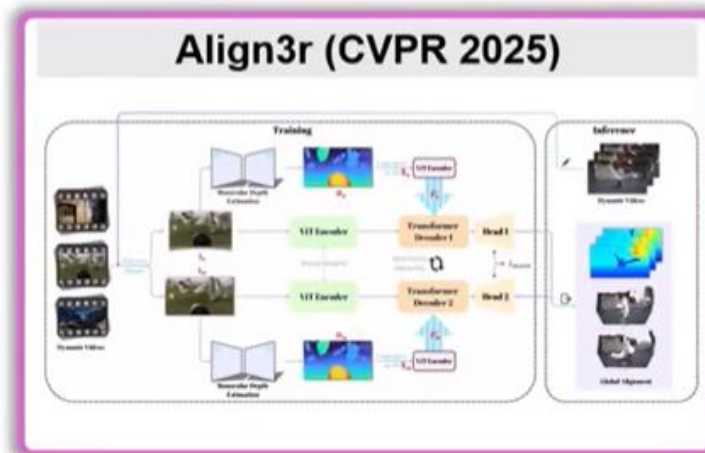
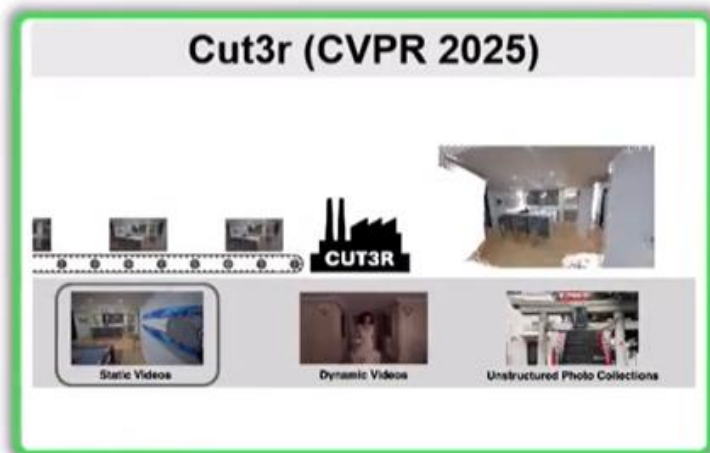
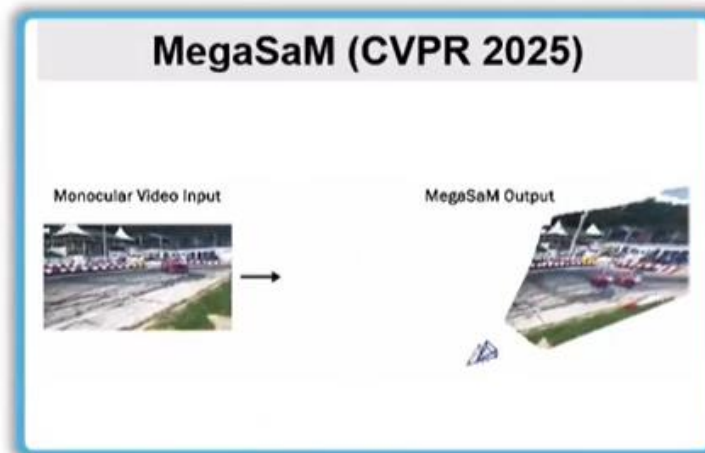
EpicKitchens



From YouTube

No ground truth data

Comparing to Existing Works



And many more

Comparing to Existing Works

Impressive results

but

Supervised training

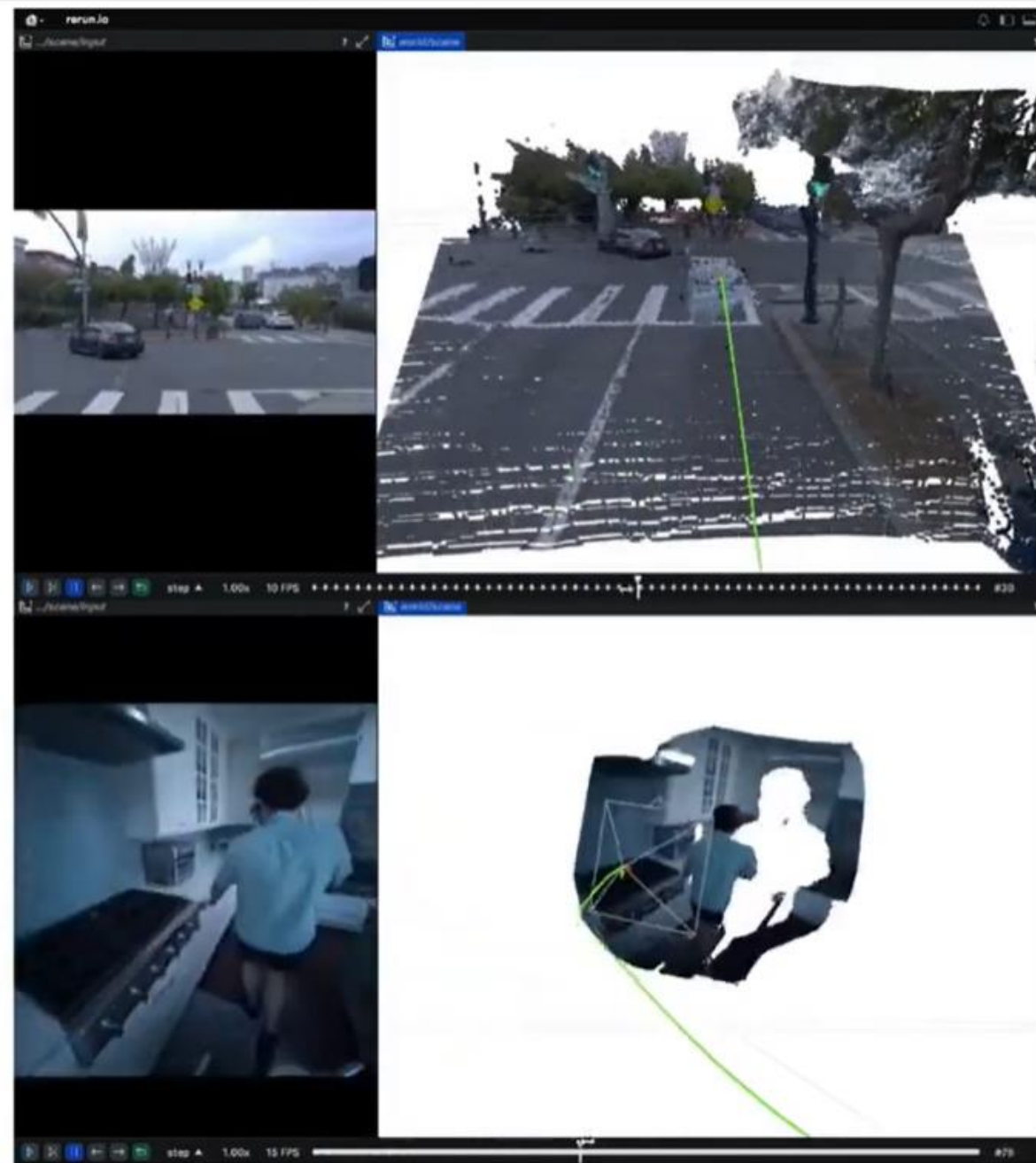


- ⚠ Expensive data collection
- ⚠ Dataset biases
- ⚠ Limited datasets
- ⚠ Synthetic-to-real gap



AnyCam solves these issues

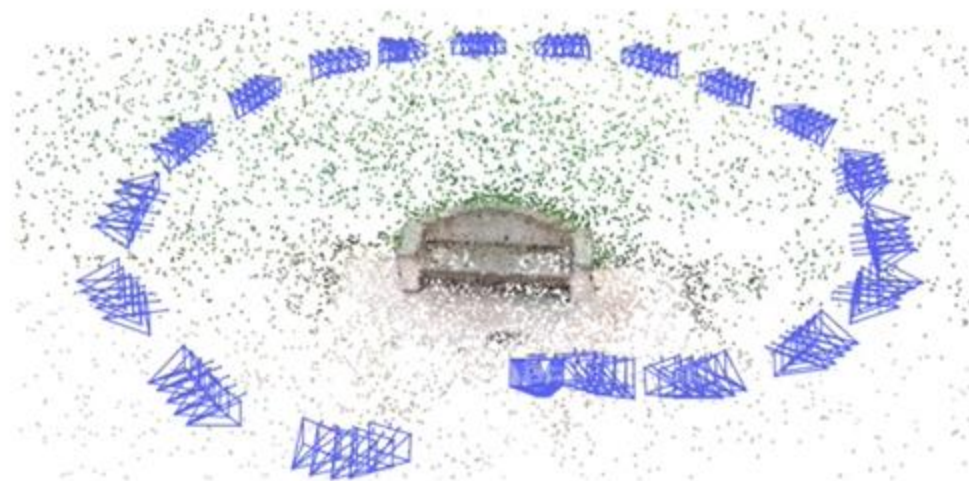




Joint optimization of NeRF and camera poses from a monocular video

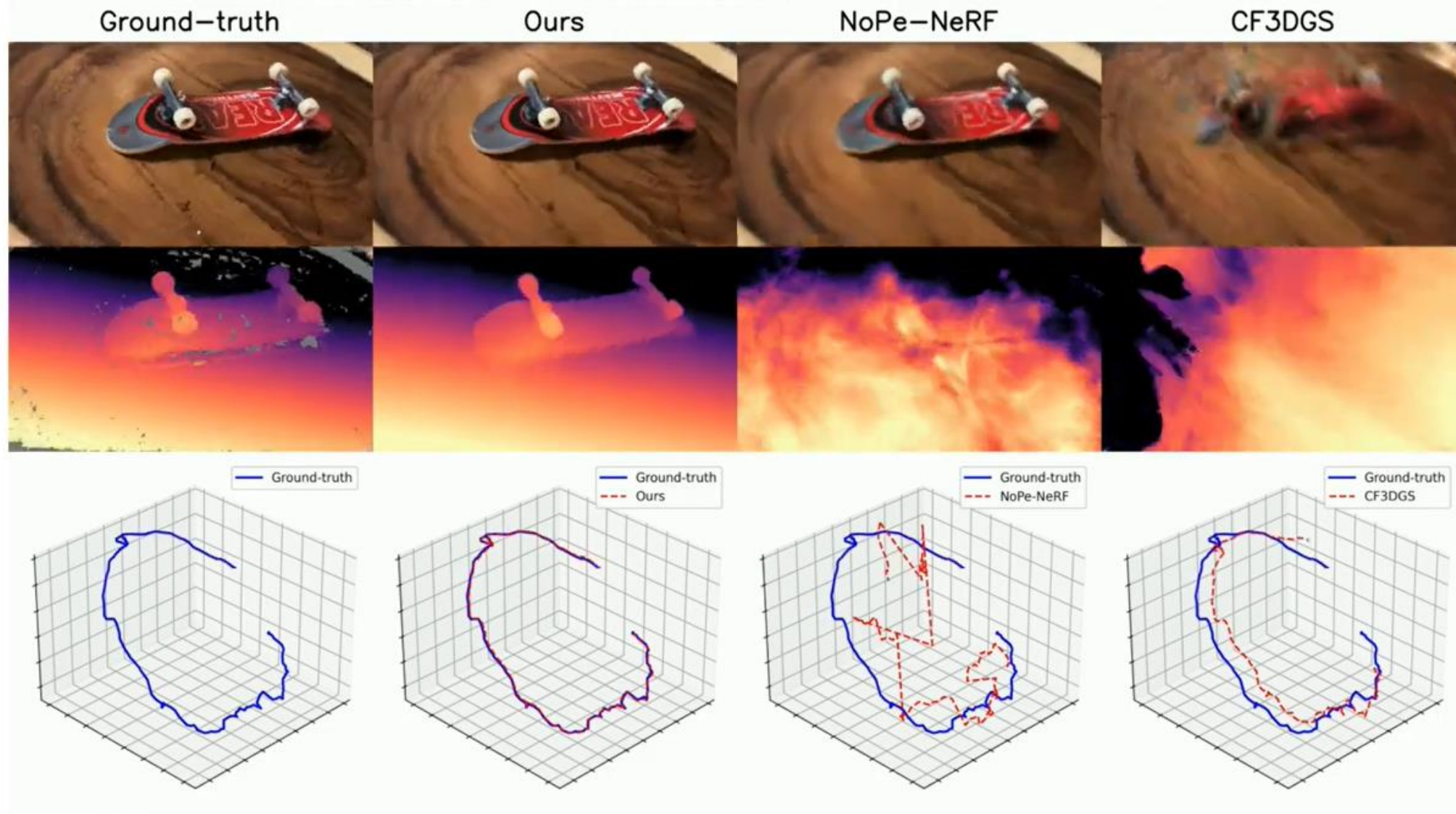


Input: Monocular video



Output: Scene geometry and camera poses.
The scene geometry is represented via NeRF.

Learned poses, rendered images and rendered depth maps for the train set



Arena Style



- Voters perform side-by-side comparisons of (typically) two anonymized models

<https://lmarena.ai/leaderboard>

- Challenges
 - Zero-sum, dynamic
 - Difficult to control for prompts, distributions, difficulty
 - Live traffic can be noisy over time
 - Relies on expensive human resources

Leaderboard Overview

See how leading models stack up across text, image, vision, and beyond. This page gives you a snapshot of each Arena, you can explore deeper insights in their dedicated tabs. [Learn more about it here.](#)

[View Blog](#) ↗

Text				WebDev			
Rank (UB) ↑	Model ↓	Score ↑	Votes ↑	Rank (UB) ↑	Model ↓	Score ↑	Votes ↑
1	gemini-2.5-pro-preview-06-05	1470	4,701	1	Gemini-2.5-Pro-Preview-06-05	1443	1,872
2	gemini-2.5-pro-preview-05-06	1446	10,386	1	Claude Opus 4 (20250514)	1412	2,466
2	o3-2025-04-16	1443	13,808	2	Gemini-2.5-Pro-Preview-05-06	1408	3,858
4	chatgpt-4o-latest-20250326	1431	18,302	2	Claude Sonnet 4 (20250514)	1389	2,078
4	gpt-4.5-preview-2025-02-27	1425	15,271	5	Claude 3.7 Sonnet (20250219)	1357	7,481
5	gemini-2.5-flash-preview-05-...	1419	9,970	6	Gemini-2.5-Flash-Preview-05-...	1312	2,626

Fixed Test Sets



- Annotate and collect examples with labelled ground truth

- Challenges

- Poor Annotation quality
- Training Data leakage
- Overfitting to the test set
- Biases
- Creating difficult examples
- Metrics for open-ended evaluation

“When a measure becomes a target, it ceases to be a good measure” – Goodhart’s Law.





Unbiasing through Textual Descriptions: Mitigating Representation Bias in Video Benchmarks

CVPR 2025



Nina Shvetsova
University of Tuebingen /
MPII Saarbrücken /
Goethe University Frankfurt



Arsha Nagrani
University of Oxford



Bernt Schiele
MPII Saarbrücken



Hilde Kuehne
University of Tuebingen /
MIT-IBM Watson AI Lab /
Goethe University Frankfurt



Christian Rupprecht
University of Oxford

References

- [\[CVPR 2025\] Feat2GS: Probing Visual Foundation Models with Gaussian Splatting \(https://www.youtube.com/watch?v=4fT5lzcAJqo\)](https://www.youtube.com/watch?v=4fT5lzcAJqo)
- [\[CVPR 2025\] AnyCam: Learning to Recover Camera Poses and Intrinsic from Casual Videos](#)
- [\[CVPR 2025\] Joint Opt. of NeRFs and Continuous Camera Motion from a Monocular Video \(Nguyen et al.\)](#)