

Rethinking Camera Parametrization for (Sparse-view) Pose Prediction



Shubham Tulsiani Assistant Professor, CMU





Goal: Virtualize a real object (for AR/VR, or posterity etc.)





Dense Multi-view Images





Single image 'in-the-wild'



'Ease' of Capture







Faithfullness

Precisely capture the underlying instance

Hallucination/ Approximation



NeRF, IDR, NeUS, Gaussian Splatting etc.

Faithfullness

















De .





CMR, Zero-1-to-3, TripoSR etc.





















(circa 2021)

I want to reconstruct Craigslist objects Jason:

> Ok, let's first setup a NeRF baseline. Run SfM to get poses, and ...

Well, I tried really hard, but COLMAP keeps crashing...

me:













COLMAP

Images

3D from images: the typical recipe





3D Representation

Cameras

Challenge: Classical SfM is not robust for Sparse-view Pose

202 Input Images











Cameras

Images

10 Input Images











Let's learn sparse-view pose prediction!





 $\{\mathbf{I}_n\}$

*essentially, learning-based localization (without mapping)

Distribution of $R_{\rm rel}$ 60° 30°, ∕-90° 90° **0**° 0° -30^o



RelPose, RelPose++ (Lin*, Zhang* et. al.)





Let's learn sparse-view pose prediction!



SparsePose (Sinha et. al.)

Regression + 3D-based refinement



RelPose, RelPose++ (Lin*, Zhang* et. al.) Energy-based modeling

PoseDiffusion (Wang et. al.) Denoising diffusion



Innovations in uncertainty modeling, joint 3D inference, etc..



SparsePose (Sinha et al.) Can we re-parametrize camera estimation as a local prediction task? (Wang et Regression + 3D-based refinement Denoising diffusion



Parametrization for Prediction: global extrinsics (+ intrinsics)

But is this a good representation for learning-based prediction?

(difficult for global models to leverage local cues e.g. correspondences)



Representing Cameras via Ray Bundles

Key Idea: Represent cameras via per-pixel rays in a common coordinate frame



Plucker Ray Parametrization

A distributed and generic representation

(inspired by work in camera calibration e.g. Grossberg & Nayar, Schops et. al.)

Can analytically recover pinhole camera parameters given predicted rays

Cameras as Rays: Pose Estimation via Ray Diffusion. Zhang*, Lin*, Ramanan, Tulsiani. In ICLR 2024

Camera Prediction via Ray Regression



Cameras as Rays: Pose Estimation via Ray Diffusion. Zhang*, Lin*, Ramanan, Tulsiani. In ICLR 2024

Camera Prediction via Ray Diffusion



Cameras as Rays: Pose Estimation via Ray Diffusion. Zhang*, Lin*, Ramanan, Tulsiani. In ICLR 2024

Visualizing Reverse Diffusion

Input Images



Directions



Moments



3D Rays





Visualizing Reverse Diffusion

Input Images



Directions

Moments



3D Rays



Visualizing Reverse Diffusion

Input Images



Directions Moments



3D Rays



Images









Pose Diffusion













RelPose++

Ray Regression (Ours)

Ray Diffusion (Ours)

Quantitative: Rotation Accuracy ($\% < 15^{\circ}$)



8

- COLMAP (SP+SG)
- RelPose
- PoseDiffusion **—**—
- Pose Regression **—**—
- RelPose++ **—**——
- Ray Regression (Ours)
- Ray Diffusion (Ours) **—**—

Quantitative: Cam Center Accuracy (% < 0.1)



- COLMAP (SP+SG)
- PoseDiffusion **___**
- Pose Regression
- RelPose++ **___**
- Ray Regression (Ours) **___**
- Ray Diffusion (Ours) **—0**—

Modeling Uncertainty

Input Images





Ray Regression

Ray Diffusion (100 Samples)



Early Stopping in Reverse Diffusion



At each diffusion timestep, we obtain a denoised ray prediction

Pose accuracy increases initially, but drastically drops later

Hypothesis: For accuracy metrics, we want modes, not samples





In-the-wild Generalization













Dense Prediction for Pose (and 3D) Estimation



Dust3r (CVPR 24), Mast3r (ECCV 24)

Pointmap prediction from image pairs + global alignment



Ace-Zero (ECCV 24)

Rethinking SfM pipeline via Scene coordinate regression

Were cameras the only missing piece for Sparse-view 3D?





RayDiffusion, Dust3r



Images







Cameras



Challenges due to unobserved aspects, pose outliers



Analysis by Generative Synthesis

Input Images $\{I_i\}$





Off-the-shelf Pose Estimator Camera Poses $\{\pi_i\}$









Analysis by Generative Synthesis

Input Images $\{I_i\}$





Off-the-shelf Pose Estimator Camera Poses $\{\pi_i\}$



















































Input

Initial Pose















Thank you!



Jason Amy Qitao

