

ASE, EFM3D & EVL: Datasets, Models & Tools for NBV

Towards Relative Reconstruction Metrics for Next-Best-View

Jan Duchscherer

VCML Seminar WS24/25

Aria Synthetic Environments

Dataset for Egocentric 3D Scene
Understanding



Figure 1: [Ave+24]

ASE Dataset Overview

Dataset Content

- 100,000 unique multi-room interior scenes
- ~2-min egocentric trajectories per scene
- Populated with 8,000 3D objects
- Aria camera & lens characteristics

Ground Truth Annotations

- 6DoF trajectories
- RGB-D frames
- 2D panoptic segmentation
- Semi-dense SLAM PC w/ visibility info
- 3D floor plan (SSL format)
- *GT meshes* as .ply files

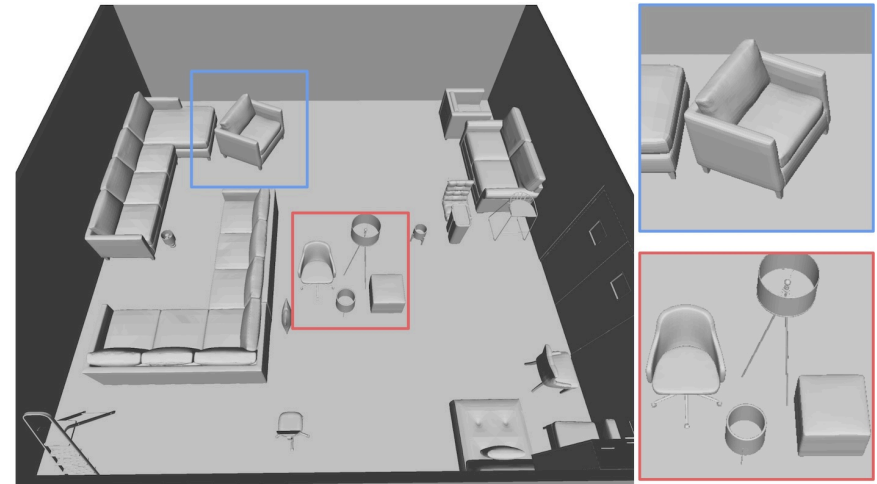
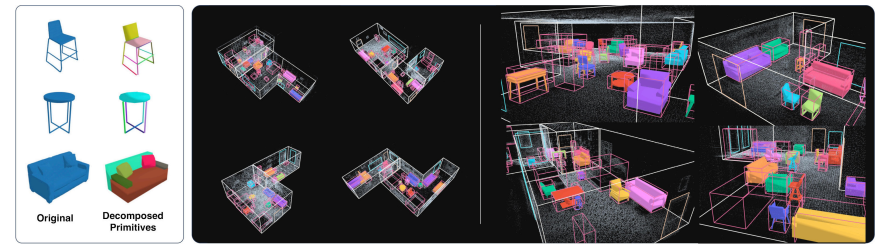


Figure 2: [Ave+24]

ASE Dataset Overview

Key Resources

- Project Aria Tools for data access
- [ASE documentation](#) [Ave+24, Met25a]

ASE Dataset Structure

```
1 scene_id/
2 |-- ase_scene_language.txt          # Ground truth scene layout in SSL format
3 |-- object_instances_to_classes.json # Mapping from instance IDs to semantic classes
4 |-- trajectory.csv                  # 6DoF camera poses along the egocentric path
5 |-- semidense_points.csv.gz         # Semi-dense 3D point cloud from MPS SLAM
6 |-- semidense_observations.csv.gz   # Point observations (which images see which points)
7 |-- rgb/                            # RGB image frames
8 |   |-- 000000.png
9 |   |-- ...
10 |-- depth/                          # Ground truth depth maps
11 |   |-- 000000.png
12 |   |-- ...
13 |-- instances/                       # Instance segmentation masks
14 |   |-- 000000.png
15 |   |-- ...
```

EFM3D Benchmark

3D Egocentric Foundation Model:
Egocentric Voxel Lifting (EVL)

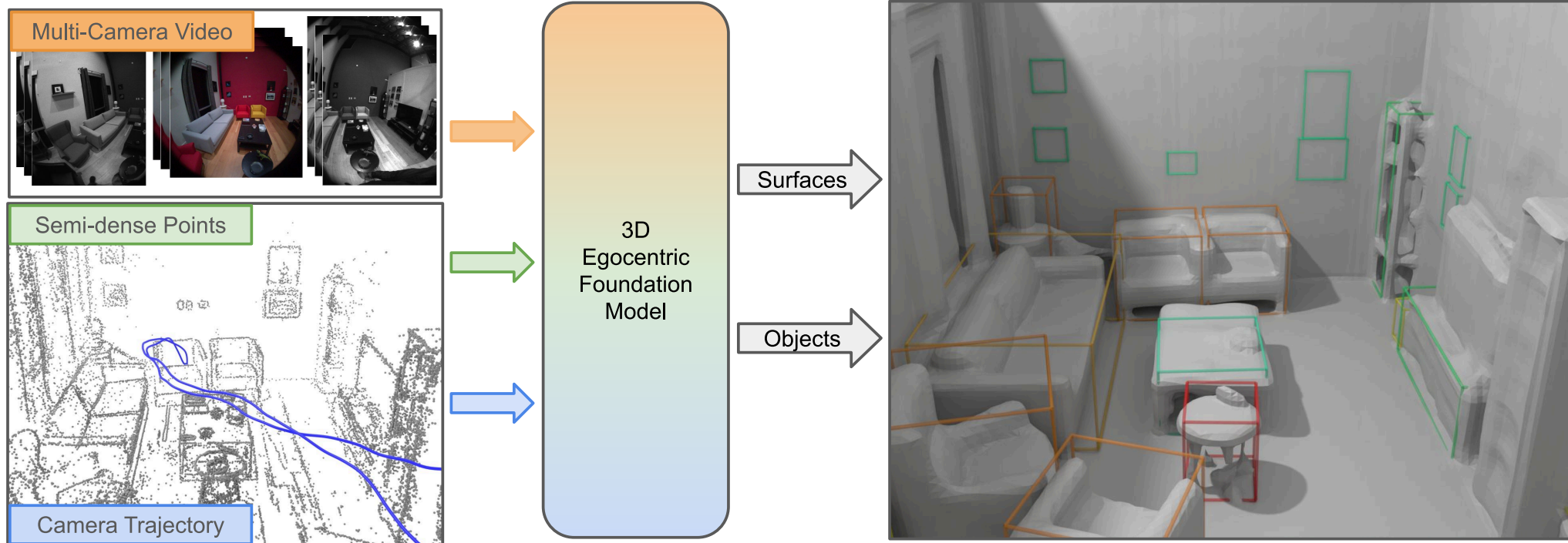


Figure 3: [Str+24]

EFM3D Tasks

- 3D object detection
- 3D surface regression (occupancy volumes)
 - on ASE, ADT¹, AEO² datasets

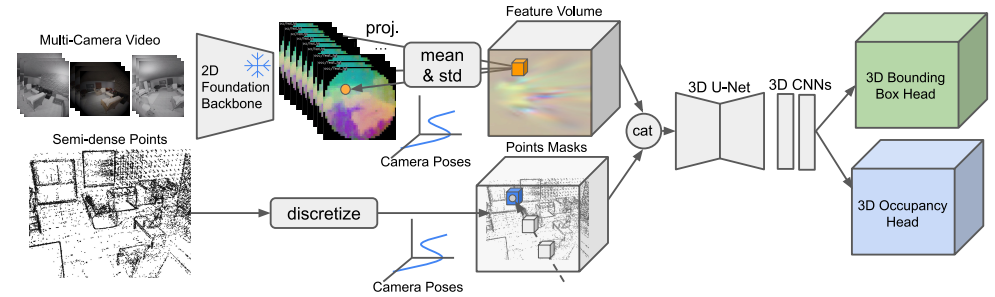


Figure 4: [Str+24]

¹ Aria Digital Twin

² Aria Everyday Objects: small-scale, real-world w/ 3D OBBS

EVL Architecture

- Utilizes **all** available egocentric modalities:
 - 1 multiple (rectified) RGB, grayscale, and semi-dense points inputs
 - 2 camera intrinsics and extrinsics
- **16.7M trainable** + 86.6M frozen params
- Inherits foundational capabilities from frozen 2D model (DinoV2.5) by lifting 2D features to 3D [Str+24]

EVL: Egocentric Voxel Lifting Architecture

HM ■

Model Overview

Egocentric Voxel Lifting (EVL): Multi-task 3D perception from egocentric video

Key Principle: Lift 2D image features to 3D voxel space using camera geometry

Input Formulation

$$\mathbf{X}_{\text{in}} = \{I_1, I_2, \dots, I_F, D_{\text{semi}}, \mathbf{K}, \mathbf{T}\}$$

Where:

- $I_f \in \mathbb{R}^{H \times W \times 3}$: RGB frames (F frames)
- $D_{\text{semi}} \in \mathbb{R}^{N \times 3}$: Semi-dense 3D points
- $\mathbf{K} \in \mathbb{R}^{3 \times 3}$: Camera intrinsics matrix
- $\mathbf{T}_f \in \text{SE}(3)$: Camera pose for frame f

Multiple camera streams supported:

- RGB (high-res)
- SLAM cameras (grayscale, rectified)

Output Formulation

3D Occupancy Volume:

$$\mathbf{V}_{\text{out}} \in \mathbb{R}^{D_x \times D_y \times D_z \times C}$$

- Voxel grid dimensions: $D_x \times D_y \times D_z$
- C channels for:
 - 1 Occupancy probability
 - 2 Object class scores
 - 3 Surface normals

Detected Objects:

$$\mathcal{O} = \{(\mathbf{b}_i^{3D}, c_i, s_i)\}_{i=1}^N$$

- $\mathbf{b}_i^{3D} \in \mathbb{R}^9$: Oriented bounding box
- c_i : Object class
- s_i : Confidence score

Feature Lifting Process

- 1 **2D Feature Extraction:** Frozen DinoV2.5 backbone

$$\mathbf{F}_{2D} = \varphi_{\text{DINOv2.5}}(\mathbf{I}_f) \in \mathbb{R}^{H' \times W' \times D_{\text{feat}}}$$

- 2 **3D Projection:** For each voxel $\mathbf{v} \in \mathbb{R}^3$, aggregate features from all frames

$$\mathbf{F}_{3D(\mathbf{v})} = \text{Aggregate} \left(\left\{ \pi \left(\mathbf{T}_f^{-1} \mathbf{v}, \mathbf{K}, \mathbf{F}_{2D}^f \right) \right\}_{f=1}^F \right)$$

where $\pi(\cdot)$ is the camera projection function

- 3 **3D Convolution:** Process lifted features

$$\mathbf{V}_{\text{out}} = \psi_{\text{3D-CNN}}(\mathbf{F}_{3D}, \mathbf{D}_{\text{semi}})$$

ATEK Toolkit

Streamlined ML Workflows for Aria Datasets

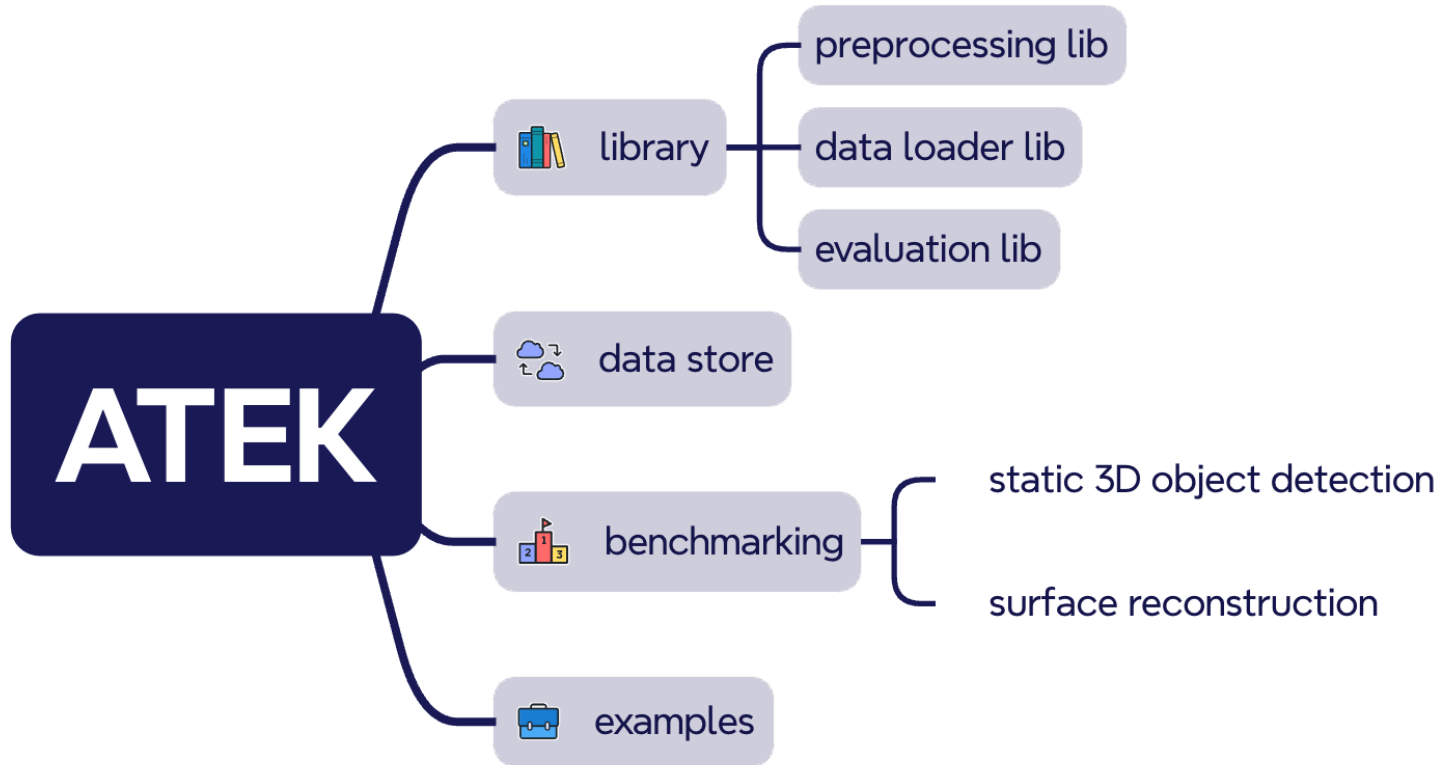


Figure 5: [Met25b]

ATEK Data Store

- Pre-processed for various tasks → ready for PyTorch training
- Local download or cloud streaming
- Eval metrics (accuracy, completeness, F-score) → adaptation for RRI
- Integration w/ Meta's MPS
- Various example notebooks

Provided Models

- *Cube R-CNN* [Bra+23] for OBBs
- *EFM* [Str+24] for OBBs & surface reconstruction

Resources

- [ATEK GitHub](#) [Met25c]
- [ECCV 2024 Tutorial: Ego-centric Research with Project Aria](#)
- Atek Context7 ID: /facebookresearch/atek

ATEK provides **streamlined ML workflows** for rapid prototyping and benchmarking on Aria datasets.

Literature Review

- Read Project Aria paper [Met25a]
- Study EFM3D & EVL architecture in depth [Str+24]
- Deep dive into GenNBV's multi-source embeddings [Che+24]
- Compare VIN-NBV vs. GenNBV: RRI prediction vs. coverage-based rewards

Technical Exploration

- Explore GT meshes (.ply files) in ASE dataset
- Get familiar with [ATEK](#) and [ATEK Data Store](#)
- Test mesh-based evaluation metrics (accuracy, completeness, F-score)
- Experiment with probabilistic 3D occupancy grids

Implementation Goals

- Implement ray-casting for mesh-based visibility computation
- Develop entity-wise RRI computation pipeline using GT meshes
- Design 5DoF action space for scene exploration
- Build multi-source state embedding (geometric + semantic + action)
- Prototype RRI prediction network architecture

Key Innovation

- First NBV method to directly optimize **reconstruction quality** (not coverage)
- Predicts **Relative Reconstruction Improvement (RRI)** without capturing new images
- 30% improvement over coverage-based baselines
- Trained 24h on 4 A6000 GPUs [Fra+25]

VIN Architecture

Predicts RRI from current reconstruction state:

$$\widehat{\text{RRI}}(q) = \text{VIN}_\theta(\mathcal{R}_{\text{base}}, \mathcal{C}_{\text{base}}, \mathcal{C}_q)$$

- **Input:** Partial point cloud + camera poses
- **Features:** Surface normals, visibility counts, depth, coverage
- **Output:** Predicted RRI via ordinal classification (15 bins)

Relative Reconstruction Improvement

For a candidate view q , RRI quantifies expected improvement:

$$\text{RRI}(q) = \frac{\text{CD}(\mathcal{R}_{\text{base}}, \mathcal{R}_{\text{GT}}) - \text{CD}(\mathcal{R}_{\text{base} \cup q}, \mathcal{R}_{\text{GT}})}{\text{CD}(\mathcal{R}_{\text{base}}, \mathcal{R}_{\text{GT}})}$$

- Range: $[0, 1]$ where higher = better view
- Normalized by current error \rightarrow scale-independent
- CD measures reconstruction quality

VIN-NBV demonstrates that *learning reconstruction-aware NBV policies* significantly outperforms traditional coverage-based approaches.

Key Innovations

- **5DoF free-space action space**: 3D position + 2D rotation (yaw, pitch)
- **Multi-source state embedding**: geometric, semantic, action representations
- **Probabilistic 3D occupancy grid** vs. binary (distinguishes unscanned from empty)
- Cross-dataset generalization: 98.26% coverage on Houses3K, 97.12% on OmniObject3D

State Representation

Geometric Embedding s_t^G :

- Probabilistic 3D occupancy grid from depth maps
- Bresenham ray-casting with log-odds update
- Three states: **occupied**, **free**, **unknown**

Semantic Embedding s_t^S :

- RGB images → grayscale → 2D CNN
- Helps distinguish holes from incomplete scans

Action Embedding s_t^A :

- Historical viewpoint sequence encoding

Combined: $s_t = \text{Linear}(s_t^G; s_t^S; s_t^A)$

GenNBV: Generalizable Next-Best-View Policy

HM ■

Action Space Design

$$\mathcal{A} = \underbrace{\mathbb{R}^3}_{\text{position}} \times \underbrace{SO(2)}_{\text{heading}}$$

- Approximately 20m x 20m x 10m position space
- Omnidirectional heading subspace
- *No hand-crafted constraints* (e.g., hemisphere)

RL-based framework with PPO. Reward: Δ CR between steps. [Che+24]

Surface-to-Surface Distance Metrics

Accuracy (Prediction \rightarrow GT):

$$\text{Acc} = \frac{1}{|\mathcal{P}|} \sum_{\mathbf{p} \in \mathcal{P}} \min_{\mathbf{q} \in \mathcal{M}_{\text{GT}}} \|\mathbf{p} - \mathbf{q}\|_2$$

Completeness (GT \rightarrow Prediction):

$$\text{Comp} = \frac{1}{|\mathcal{M}_{\text{GT}}|} \sum_{\mathbf{q} \in \mathcal{M}_{\text{GT}}} \min_{\mathbf{p} \in \mathcal{P}} \|\mathbf{p} - \mathbf{q}\|_2$$

Where:

- \mathcal{P} : Predicted PC from dense or semi-dense reconstruction or sampled from pred mesh
- \mathcal{M}_{GT} : Sampled points from GT mesh

Precision, Recall & F-score

At threshold τ (typically 5cm):

$$\text{Pr}_{@ \tau} = \frac{|\{\mathbf{p} \in \mathcal{P} : \min_{\mathbf{q} \in \mathcal{M}_{\text{GT}}} \|\mathbf{p} - \mathbf{q}\| < \tau\}|}{|\mathcal{P}|}$$

$$\text{Re}_{@ \tau} = \frac{|\{\mathbf{q} \in \mathcal{M}_{\text{GT}} : \min_{\mathbf{p} \in \mathcal{P}} \|\mathbf{p} - \mathbf{q}\| < \tau\}|}{|\mathcal{M}_{\text{GT}}|}$$

$$\text{F-score}_{@ \tau} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

Chamfer Distance (Bidirectional)

$$\text{CD}(\mathcal{P}, \mathcal{M}_{\text{GT}}) = \text{Acc} + \text{Comp}$$

Combines both directions of surface error

ATEK Implementation:

- `evaluate_single_mesh_pair()`¹ computes all metrics using:
 - [trimesh.Trimesh](#): Load meshes + sample surfaces uniformly
 - `compute_pts_to_mesh_dist()`: Point-to-mesh distance via batched triangle projection
 - `point_to_closest_tri_dist()`: Barycentric coordinate projection test + plane distance
 - Fallback: `point_to_closest_vertex_dist()` when projection fails

See [metrics.qmd](#) for detailed formulas and algorithm explanations.

¹[src](#)

RRI from GT Mesh

Given:

- \mathcal{M}_{GT} : GT mesh (from ASE .ply files)
- \mathcal{P}_t : Current reconstruction from first t views
- $\mathbf{q} \in \text{SO}(2) \times \mathbb{R}^3$: Candidate viewpoint, 5DoF (position + yaw, pitch)
- $\mathcal{P}_{t \cup \mathbf{q}}$: Updated reconstruction after capturing from \mathbf{q}

Mesh-based RRI (oracle):

$$\text{RRI}(\mathbf{q}) = \frac{\text{CD}(\mathcal{P}_t, \mathcal{M}_{\text{GT}}) - \text{CD}(\mathcal{P}_{t \cup \mathbf{q}}, \mathcal{M}_{\text{GT}})}{\text{CD}(\mathcal{P}_t, \mathcal{M}_{\text{GT}})}$$

RRI Oracle Pipeline

- 1 Load GT mesh from ASE
- 2 Build \mathcal{P}_t from captured views
 - dense PC from depth maps
 - or semi-dense SLAM PC¹
- 3 Simulate view from \mathbf{q}
 - Ray-cast to $\mathcal{M}_{\text{GT}} \rightarrow \mathcal{P}_{\mathbf{q}}$
- 4 Merge: $\mathcal{P}_{t \cup \mathbf{q}} = \mathcal{P}_t \cup \mathcal{P}_{\mathbf{q}}$
 - Voxel downsample for consistency (e.g., 1cm)
- 5 Compute RRI using CD metric

¹semidense_points.csv

Key Functions from EFM3D & ATEK

Point Cloud Generation:

- `dist_im_to_point_cloud_im()`: Depth \rightarrow 3D points
- `collapse_pointcloud_time()`: Merge temporal PCs
- `pointcloud_to_voxel_counts()`: PC \rightarrow density grid

Ray-Mesh Operations:

- `ray_obb_intersection()`: Ray-box intersection
- `sample_depths_in_grid()`: Sample depths along rays

Distance Computation:

- `compute_pts_to_mesh_dist()`: Min distance to triangles
- `eval_mesh_to_mesh()`: Full evaluation pipeline

RRI-based NBV for Scene-Level Reconstruction

HM ■

VIN with EVL Backbone

Our Approach: Adapt RRI prediction to **scene-level** environments with 5DoF action space

RRI with GT Meshes

Use ASE visibility data + GT meshes for *oracle RRI*:

$$\text{RRI}(\mathbf{q}) = \frac{d(\mathcal{P}_{\text{partial}}, \mathcal{M}_{\text{GT}}) - d(\mathcal{P}_{\text{partial}} \cup \mathbf{q}, \mathcal{M}_{\text{GT}})}{d(\mathcal{P}_{\text{partial}}, \mathcal{M}_{\text{GT}})}$$

where \mathcal{M} represents meshes, $d(\cdot, \cdot)$ is mesh distance

Proposed Pipeline

- 1 Reconstruct:** Build $\mathcal{P}_{\text{partial}}$ from historical trajectory
- 2 Sample:** Generate candidate viewpoints in free space around latest pose
- 3 Compute Features:** Extract geometric + semantic embeddings from **EVL**
- 4 Predict:** Use **VIN** to predict RRI per candidate
- 5 Select:** Choose NBV based on RRI

Key Challenge: Ray-casting from candidate views to compute visibility on GT meshes for entity-wise RRI computation

Extension Entity-wise RRI:

$$\text{RRI}_{\text{total}} = \sum_{e \in \mathcal{E}} w_e \cdot \text{RRI}_e \text{ where } \mathcal{E} = \{\text{walls, doors, objects, ...}\}$$

- This could be done by segmenting the GT meshes and PCs per entity type and computing the RRI separately.

Bibliography

- [Ave+24] A. Avetisyan *et al.*, “SceneScript: Reconstructing Scenes With An Autoregressive Structured Language Model.” [Online]. Available: <https://arxiv.org/abs/2403.13064>
- [Met25] Meta Platforms Inc., “Aria Synthetic Environments Dataset.” [Online]. Available: https://facebookresearch.github.io/projectaria_tools/docs/open_datasets/aria_synthetic_environments_dataset
- [Str+24] J. Straub, D. DeTone, T. Shen, N. Yang, C. Sweeney, and R. Newcombe, “EFM3D: A Benchmark for Measuring Progress Towards 3D Egocentric Foundation Models.” [Online]. Available: <https://arxiv.org/abs/2406.10224>
- [Met25] Meta Platforms Inc., “Aria Training and Evaluation toolkit (ATEK).” [Online]. Available: <https://github.com/facebookresearch/ATEK>
- [Bra+23] G. Brazil, A. Kumar, J. Straub, N. Ravi, J. Johnson, and G. Gkioxari, “Omni3D: A Large Benchmark and Model for 3D Object Detection in the Wild.” [Online]. Available: <https://arxiv.org/abs/2207.10660>

Bibliography

- [Met25] Meta Platforms Inc., “Aria Training and Evaluation toolkit (ATEK) documentation.” [Online]. Available: https://facebookresearch.github.io/projectaria_tools/docs/ATEK/about_ATEK
- [Che+24] X. Chen, Q. Li, T. Wang, T. Xue, and J. Pang, “GenNBV: Generalizable Next-Best-View Policy for Active 3D Reconstruction.” [Online]. Available: <https://arxiv.org/abs/2402.16174>
- [Fra+25] N. Frahm *et al.*, “VIN-NBV: A View Introspection Network for Next-Best-View Selection.” [Online]. Available: <https://arxiv.org/abs/2505.06219>