# E-M3RF: An Equivariant Multimodal 3D Re-assembly Framework

Adeela Islam[1,2]   Stefano Fiorini[1]   Manuel Lecha[1]   Theodore Tsesmelis[1]   Stuart James[3]
Pietro Morerio[1]   Alessio Del Bue[1]

[1]Fondazione Istituto Italiano di Tecnologia   [2]University of Genova   [3]Durham University

{adeela.islam, pietro.morerio}@iit.it

## Abstract

*3D reassembly is a fundamental geometric problem, and in recent years it has increasingly been challenged by deep learning methods rather than classical optimization. While learning approaches have shown promising results, most still rely primarily on geometric features to assemble a whole from its parts. As a result, methods struggle when geometry alone is insufficient or ambiguous, for example, for small, eroded, or symmetric fragments. Additionally, solutions do not impose physical constraints that explicitly prevent overlapping assemblies. To address these limitations, we introduce E-M3RF, an equivariant multimodal 3D reassembly framework that takes as input the point clouds, containing both point positions and colors of fractured fragments, and predicts the transformations required to reassemble them using $SE(3)$ flow matching. Each fragment is represented by both geometric and color features: i) 3D point positions are encoded as rotation-consistent geometric features using a rotation-equivariant encoder, ii) the colors at each 3D point are encoded with a transformer. The two feature sets are then combined to form a multimodal representation. We experimented on four datasets: two synthetic datasets, Breaking Bad and Fantastic Breaks, and two real-world cultural heritage datasets, RePAIR and Presious, demonstrating that E-M3RF on the RePAIR dataset reduces rotation error by **23.1%** and translation error by **13.2%**, while Chamfer Distance decreases by **18.4%** compared to competing methods.*

## 1. Introduction

Reassembly tasks play a fundamental role in many domains, including reconstructing archaeological 3D artifacts [26], piecing together shredded documents [18], molecular docking [3], and solving jigsaw or 3D puzzles [7, 22]. Accurate reassembly of fractured fragments requires precise estimation of each piece's $SE(3)$ transformation that rebuilds the whole 3D object shape. This problem becomes particularly challenging for symmetric or eroded fragments
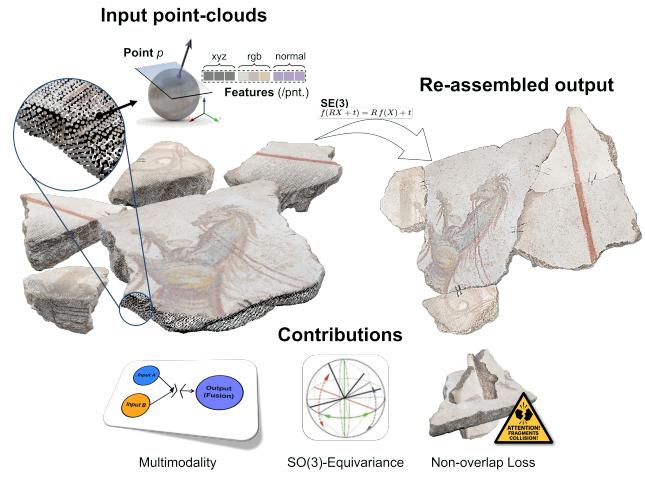


Figure 1. **E-M3RF overview.** Fragments are represented as colored point clouds with per-point features $(x, y, z, \text{rgb}, \mathbf{n})$. A multimodal SE(3)-equivariant backbone fuses geometry and color and predicts per-fragment rigid transforms $(R, t)$ via flow-based estimation, while a differentiable non-overlap loss enforces physical plausibility. The predicted poses reassemble the object, aligning fracture boundaries and color patterns into a coherent result.

due to the scarcity of distinctive geometric features. In recent years, 3D fracture reassembly has remained a fundamentally geometric optimization task but has increasingly been approached using deep learning [21, 22] rather than classical optimization [14, 33]. Learning-based methods have progressed along two complementary trends: improving fragment representations via powerful encoders and pretraining schemes that emphasize fine-scale geometry and surface appearance [9, 12], and reframing pose estimation as a generative task, where flow- or diffusion-based pipelines iteratively refine fragment alignments [21]. Despite these recent advances, geometry-centric designs have practical limitations: geometric information alone can be ambiguous [26], while real fragments often exhibit appearance details, such as paint, grain, or tool marks. These discriminative cues are ignored by geometry-only pipelines, overlooking an important source of information. Moreover, geometry-only methods often fail to account for the for-
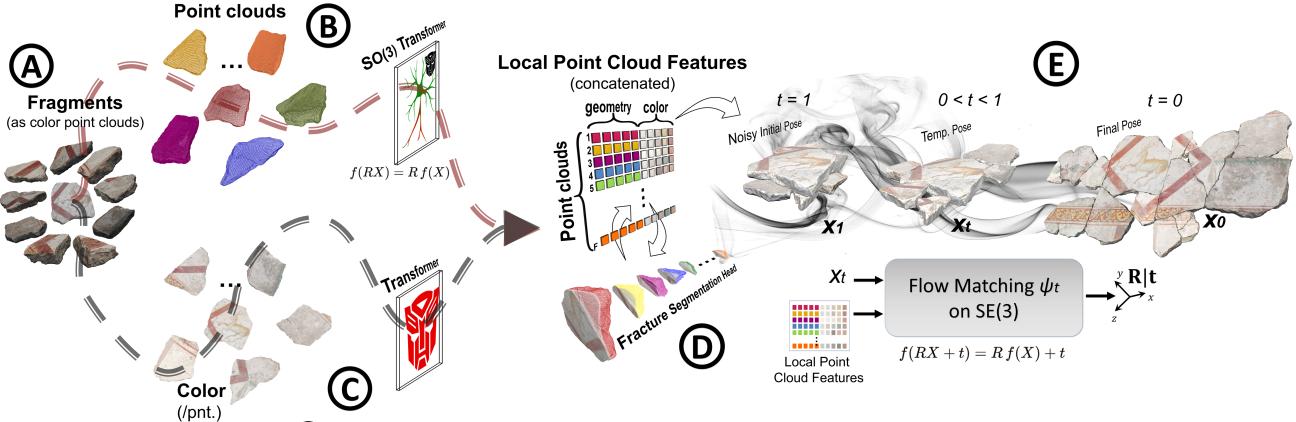
Figure 2. **E-M3RF's pipeline.** (A) A set of textured fragments (left) loaded as colored point clouds (w/ per-point RGB). Two transformer encoders extract (B) **fracture-aware geometric features**, which are geometrically equivariant, since they are processed by an SO(3)-**equivariant Transformer** so representations transform consistently under motions (rotation/ translation) and (C) **color features** through a Transformer over per-point colors which extracts dense color descriptors. The streams are concatenated into (D) **local point-cloud features** (xyz/rgb/normals tokens); the geometry stream is guided by a fracture-segmentation head during pretraining. Using these fused features, we perform (E) **flow matching on** SE(3): a time-dependent vector field $\psi_t$ transports fragment poses from a noisy initialization $x_1$ through intermediate states to the assembled configuration $x_0$, yielding per-fragment transforms $(R, t)$. During training, a **non-overlap loss** penalizes interpenetration to enforce physically plausible assemblies (omitted from the illustration for clarity).

mal group structure of transformations acting on the pieces. This is compounded by state-of-the-art models' struggle with physical constraints, which often results in intersecting fragments [9, 21].

To address the limitations of geometry-only approaches and handle complex fractures, we introduce E-M3RF, a multimodal 3D reassembly framework that takes as input the colored point clouds of fractured fragments and predicts the corresponding $SE(3)$ transformations via $SE(3)$ flow matching [1, 2, 15], see Fig. 1. The model jointly leverages fragment color and symmetry-aware geometry: an $SO(3)$-equivariant encoder extracts rotation-consistent geometric features from origin-centered point clouds (*i.e.*, no translation involved), while a parallel color branch encodes per-fragment appearance. The resulting features are concatenated into a joint multimodal embedding per fragment point. Both rotations and translations are predicted with a transformer-based architecture trained via a flow-matching objective on the $SE(3)$ manifold. A differentiable no-overlap loss further enforces physical plausibility by penalizing fragment intersections, as illustrated in Figure 2.

We evaluate E-M3RF on a combination of synthetic and real fracture datasets, including the large-scale synthetic Breaking Bad [22] and Fantastic Breaks [10], as well as the real-world RePAIR [26] and Presious [23]. In all four datasets, we consistently outperform geometry-only baselines, achieving higher part accuracy and lower rotation, translation, and Chamfer Distance errors. To quantify the impact of our design choices, we perform ablations over each component.

**The main contributions of this work are**:

- An end-to-end multimodal fragment representation that using both geometric and color features, producing a joint appearance–geometry embedding that resolves ambiguities and leverages surface patterns when available.
- We leverage an inductive bias in the form of a rotation-equivariant geometric backbone. Specifically, $SO(3)$-equivariant layers ensure that features transform consistently under rotations, improving robustness to arbitrary fragment orientations and enhancing generalization.
- We enforce physical plausibility via a differentiable no-overlap loss that penalizes fragment intersections.

## 2. Related Work

We review the state of the art in 3D reassembly and puzzle-solving problems, together with a deep dive into representation encoders commonly used in the solution of 3D tasks.

### 2.1. 3D Fracture Assembly

3D puzzle solving and fractured-object reassembly have been approached through a broad range of techniques, spanning classical geometric registration, learning-based pose estimation, and more recent generative formulations. Classical pipelines rely on hand-crafted local descriptors, correspondence hypothesis testing, and ICP-style refinement [17, 27]. These methods provide interpretable matching criteria grounded in geometry and remain effective for clean, low-ambiguity fragments. Learning-based methods shift to data-driven feature extraction and direct pose prediction. Recent works learn per-point or per-fragment embeddings and estimate 6-DoF alignments using regression or transformer-based architectures [12, 16, 21, 29], often

supported by large-scale synthetic fracture corpora such as Breaking Bad [22]. Generative approaches, including diffusion or flow-based estimators [12, 21, 29], treat assembly as sampling or transporting distributions over $SE(3)$, enabling iterative refinement and uncertainty-aware reasoning. These methods represent a growing direction for globally coherent reconstruction.

Although real-world fragments contain useful color cues, existing learning-based reassembly methods rely only on geometric information. E-M3RF addresses this limitation by fusing geometric and color features into a unified multimodal representation, enabling more discriminative fragment descriptors and more reliable alignment.

## 2.2. Representation (Encoder)

PointNet [19] introduced a shared-MLP encoder with global feature aggregation, proving that deep learning on point sets is effective for 3D tasks. PointNet++ [20] extended PointNet with hierarchical neighborhood grouping, allowing the network to capture local geometric features across multiple scales. Subsequent methods strengthened local feature modeling through convolutional and graph-based operators. EdgeConv [28] captures point-to-point relationships via dynamic graphs, while KPConv [25] employs kernel-point convolutions as a continuous 3D filtering mechanism, providing strong geometric inductive bias. Transformer-based encoders have recently gained prominence in 3D vision. Point Transformer [32] and PCT [8] use self-attention to capture long-range point dependencies, while hierarchical variants like Swin3D [6] apply windowed attention to 3D grids. However, these models are not inherently rotation-aware and do not ensure consistent features under arbitrary $SO(3)$ transformations. To address this issue, Vector Neurons [4] lift point features from scalars to 3D vectors, allowing rotations to act via matrix multiplication. Building on this representation, the VN-Transformer [5] extends self-attention to vector features, yielding a rotation-equivariant attention mechanism that preserves $SO(3)$ consistency throughout the network and enables it to reason about 3D geometry while respecting rotational symmetries.

To effectively address reassembly scenarios where fragments undergo arbitrary rotations, we leverage this rotation-equivariant inductive bias so that the backbone inherits $SO(3)$ transformations directly into the learned features, providing E-M3RF with representations that remain geometrically consistent under all input orientations.

## 3. Problem Statement

We formalize the geometric object reassembly problem given a finite collection of disjoint fragments.

**Fragments.** Let $P = \{p_i\}_{i=1}^N$ be a finite set of points. A fragment is defined as the tuple

$$\mathcal{F} = \left( P, \ \mathcal{X}, \ \{\mathcal{G}_k\}_{k=1}^{K_{\text{geo}}}, \ \{\mathcal{A}_k\}_{k=1}^{K_{\text{app}}} \right). \tag{1}$$

The map $\mathcal{X} : P \to \mathbb{R}_{\text{xyz}}^3$ assigns spatial coordinates, *i.e.* the 3D coordinates for each point. Each geometric attribute $\mathcal{G}_k : P \to \mathbb{R}^3$ transforms under the standard vector representation of $SE(3)$ (*e.g.*, normals), while each appearance attribute $\mathcal{A}_k : P \to \mathbb{R}^{d_k}$ is invariant under rotations (*e.g.*, color). In general, there can be $K_{\text{geo}} \in \mathbb{N}_0$ geometric attributes and $K_{\text{app}} \in \mathbb{N}_0$ appearance attributes. Fixing an ordering of the points induces two structured feature spaces:

$$V_{\text{geo}} = \mathbb{R}^N \otimes \mathbb{R}^{K_{\text{geo}}+1} \otimes \mathbb{R}^3, \quad V_{\text{app}} = \mathbb{R}^N \otimes \mathbb{R}^{d_{\text{app}}}, \tag{2}$$

where the $K_{\text{geo}} + 1$ geometric channels include $\mathcal{X}$ as a dedicated vector channel, and $d_{\text{app}} = \sum_{k=1}^{K_{\text{app}}} d_k$ is the total dimensionality of all appearance channels. The full fragment representation lies in

$$V_F = V_{\text{geo}} \oplus V_{\text{app}}, \tag{3}$$

and we write $F \in V_F$ for the feature tensor associated with a fragment. We denote its positional component by $F_{\text{xyz}} \in \mathbb{R}^N \otimes \mathbb{R}^3$.

**Unassembled Object.** An unassembled object is represented as a collection of $M$ fragments $\mathcal{F}$, each recentered so that its centroid lies at the origin:

$$\hat{\mathcal{O}} = \left\{ \hat{F}_i \in V_F \ \Big| \ \frac{1}{N_i} \mathbf{1}_{N_i}^\top \hat{F}_{\text{xyz},i} = \mathbf{0} \in \mathbb{R}^3 \right\}_{i=1}^M,$$

where $\mathbf{1}_{N_i}$ is the all-ones vector of length $N_i$.

**Assembled Object.** We assume the existence of a canonical assembled configuration of the object. Each centered fragment $\hat{F}_i$ is placed in this configuration by a rigid transformation $g_i \in SE(3)$, acting exclusively on its positional channel. The assembled object is then denoted by

$$\mathcal{O} = \left\{ F_i = g_i \cdot \hat{F}_i \right\}_{i=1}^M.$$

**Object Reassembly.** Given any unassembled object $\hat{\mathcal{O}} = \{\hat{F}_i\}_{i=1}^M$, the goal is to recover rigid transformations $\tilde{g}_i \in SE(3)$ that approximate the unknown ground-truth $g_i$. We therefore seek a single model

$$\Phi : \hat{\mathcal{O}} \to SE(3)^M, \qquad \Phi(\hat{\mathcal{O}}) \approx \{ g_i \}_{i=1}^M,$$

that operates on arbitrary unassembled objects and whose predictions transform and assemble the fragments:

$$\Phi(\hat{\mathcal{O}}) \cdot \hat{\mathcal{O}} \approx \mathcal{O}.$$

## 4. Methodology

We identify two capabilities required for object reassembly: (1) a representation expressive enough to capture

the geometric and appearance cues that govern how fragments fit together, and (2) a mechanism that uses this representation to predict the rigid transformations realizing the correct assembly. Having formalized the task in Section 3, our methodology is built to satisfy these requirements in a principled, symmetry–aware way, see Fig. 3. We first learn fragment–level descriptors using an $S_N \times SO(3)$–equivariant encoder, ensuring that the learned representations respect the intrinsic permutation and rotational symmetries of the fragment space (details in the suppl. material). These descriptors condition a geometric generative model—a conditional $SE(3)$ Riemannian flow-matching network—which predicts a coherent set of rigid transformations aligning fragments into their assembled configuration. The full pipeline couples equivariant representation learning with flow-based geometric reasoning, yielding a unified approach that maps any unassembled object $\hat{\mathcal{O}}$ to an assembled one.

## 4.1. Encoding Fragments

**Colored Fragments.** We consider fragments whose per-point attributes consist of positions, surface normals, and RGB color. Following the formulation in Eq. 1 of Sec. 3 , each fragment is represented as

$$\mathcal{F} = (P, \mathcal{X}, \{\mathcal{N}\}, \{\mathcal{C}\}), \quad \mathcal{N} : P \to \mathbb{R}^3, \quad \mathcal{C} : P \to \mathbb{R}^3,$$

where $\mathcal{X} : P \to \mathbb{R}^3$ gives point coordinates, $\mathcal{N}$ gives surface normals, and $\mathcal{C}$ gives per-point color. Hence, the geometric and appearance feature spaces in Eq. 2 specialize to

$$V_{\text{geo}} = \mathbb{R}^N \otimes \mathbb{R}^2 \otimes \mathbb{R}^3, \qquad V_{\text{rgb}} = \mathbb{R}^N \otimes \mathbb{R} \otimes \mathbb{R}^3,$$

where the two geometric vector channels in $V_{\text{geo}}$ correspond to the coordinates $\mathcal{X}$ and the normals $\mathcal{N}$. Each fragment is represented by an element $F \in V_F$ in the joint space of Eq. 3, with canonical projections: $F_{\text{geo}} \in V_{\text{geo}}$ collecting the geometric vector channels; $F_{\text{xyz}} \in \mathbb{R}^N \otimes \mathbb{R}^3$ giving the point coordinates; $F_{\text{n}} \in \mathbb{R}^N \otimes \mathbb{R}^3$ giving the normals; and $F_{\text{rgb}} \in V_{\text{rgb}}$ giving the per-point color vectors.

**Fragment Symmetries.** Since each fragment in $\hat{\mathcal{O}}$ is centered at the origin, its continuous geometric symmetry reduces to global rotations $SO(3) \subset SE(3)$. In addition, the indexing of points induces a discrete permutation symmetry described by the symmetric group $S_N$. The resulting symmetry group acting on geometric features is therefore

$$\mathbf{G} := S_N \times SO(3),$$

where $S_N$ permutes point indices and $SO(3)$ acts on every geometric $\mathbb{R}^3$ channel. For $\sigma \in S_N$, let $P_\sigma$ be the associated permutation matrix, and let $I_2$ denote the identity on the two geometric channels. The induced action on $F_{\text{geo}} \in V_{\text{geo}}$ is

$$(\sigma, R) \cdot F_{\text{geo}} = (P_\sigma \otimes I_2 \otimes R) F_{\text{geo}}, \quad (\sigma, R) \in \mathbf{G}, \quad (4)$$
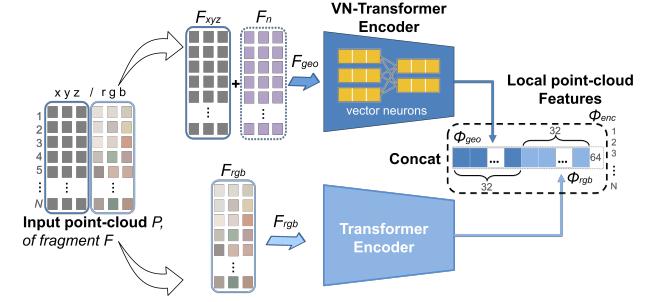


Figure 3. **Geometry/color feature extraction and fusion.** The input point cloud $P$ (per-point $xyz$, $rgb$ and normals $F_{\text{n}}$), of fragment $\mathcal{F}$ is split into $F_{\text{geo}}$ and $F_{\text{rgb}}$. A **VN-Transformer encoder** ($SO(3)$-equivariant via vector neurons) produces geometric features $\Phi_{\text{geo}}$, while a standard **Transformer encoder** produces color features $\Phi_{\text{rgb}}$. The features are concatenated per point to yield $\Phi_{\text{enc}} = [\Phi_{\text{geo}} \| \Phi_{\text{rgb}}]$, which serve as the local point-cloud representation for downstream pose estimation.

where $\otimes$ is the Kronecker product implementing the tensor-product representation. In pointwise notation, for $i \in \{1, \dots, N\}$,

$$((\sigma, R) \cdot F_{\text{geo}})_i = \begin{bmatrix} R\,(F_{\text{xyz}})_{\sigma^{-1}(i)} \\ R\,(F_{\text{n}})_{\sigma^{-1}(i)} \end{bmatrix} \in \mathbb{R}^2 \otimes \mathbb{R}^3,$$

where the first component corresponds to permuted and rotated coordinates and the second to permuted and rotated normals. For color features $F_{\text{rgb}} \in V_{\text{rgb}}$, the $SO(3)$ factor acts trivially. Writing the induced action using the same tensor-product structure,

$$(\sigma, R) \cdot F_{\text{rgb}} = (P_\sigma \otimes I_1 \otimes I_3) F_{\text{rgb}}, \quad (\sigma, R) \in \mathbf{G}, \quad (5)$$

reflecting that RGB vectors are rotation-invariant and only their pointwise ordering is affected.

**Geometric Encoder** $\Phi_{\text{geo}}$**.** The geometric branch implements a $\mathbf{G}$-equivariant map using a VN-Transformer [5],

$$\Phi_{\text{geo}} : V_{\text{geo}} \to \mathbb{R}^N \otimes \mathbb{R}^{C_{\text{geo}}} \otimes \mathbb{R}^3,$$

which outputs $C_{\text{geo}}$ equivariant vector channels, each with values in $\mathbb{R}^3$. The $\mathbf{G}$-action on this output space is obtained by lifting the action in Eq. 4 from the two input geometric channels to the $C_{\text{geo}}$ learned channels. As a result, $\Phi_{\text{geo}}$ satisfies

$$\Phi_{\text{geo}}\big((P_\sigma \otimes I_2 \otimes R)\,F_{\text{geo}}\big) = (P_\sigma \otimes I_{C_{\text{geo}}} \otimes R)\,\Phi_{\text{geo}}(F_{\text{geo}}),$$

for all $(\sigma, R) \in \mathbf{G}$. Here $I_{C_{\text{geo}}}$ denotes the identity in $\mathbb{R}^{C_{\text{geo}}}$.

**Color Encoder** $\Phi_{\text{rgb}}$**.** Color features are encoded by a Transformer

$$\Phi_{\text{rgb}} : V_{\text{rgb}} \to \mathbb{R}^N \otimes \mathbb{R}^{C_{\text{rgb}}} \otimes \mathbb{R}^3,$$

which operates on rotation–invariant channels. Consequently, the only nontrivial component of the $\mathbf{G}$-action is

the permutation action of $S_N$. Thus, the $SO(3)$ factor of $\mathbf{G}$ acts trivially on $V_{\mathrm{rgb}}$, and the only nontrivial group action arises from permutations of point indices. Lifting the permutation action of Eq. 5 to the output space, $\Phi_{\mathrm{rgb}}$ satisfies

$$\Phi_{\mathrm{rgb}}((P_\sigma \otimes I_1 \otimes I_3)F_{\mathrm{rgb}}) = (P_\sigma \otimes I_{C_{\mathrm{rgb}}} \otimes I_3)\, \Phi_{\mathrm{rgb}}(F_{\mathrm{rgb}}),$$

for all $\sigma \in S_N$.

**Fused Encoder $\Phi_{\mathrm{enc}}$.** Let $C = C_{\mathrm{geo}} + C_{\mathrm{rgb}}$. We define the fused encoder

$$\Phi_{\mathrm{enc}} : V_F \longrightarrow \mathbb{R}^N \otimes \mathbb{R}^C \otimes \mathbb{R}^3,$$

by

$$\Phi_{\mathrm{enc}}(F) = \big[\Phi_{\mathrm{geo}}(F_{\mathrm{geo}}) \,\|\, \Phi_{\mathrm{rgb}}(F_{\mathrm{rgb}})\big],$$

where $\|$ denotes concatenation along the feature–channel axis. For each point $i$, the geometric block $\Phi_{\mathrm{geo}}(F_{\mathrm{geo}})_i \in \mathbb{R}^{C_{\mathrm{geo}}} \otimes \mathbb{R}^3$ and the color block $\Phi_{\mathrm{rgb}}(F_{\mathrm{rgb}})_i \in \mathbb{R}^{C_{\mathrm{rgb}}}$ are merged into a single $C$-channel tensor. Since both branches respect the $\mathbf{G}$-action on their respective representation spaces, the fused encoder $\Phi_{\mathrm{enc}}$ is $\mathbf{G}$-equivariant by construction.

### 4.2. Fracture-Boundary Segmentation Pretraining

The encoder is pretrained via an auxiliary fracture-boundary segmentation task, which encourages the learned fragment representations to capture the geometric signatures of fracture interfaces. For each fragment $\hat{F}_i$, the encoder produces the pointwise embedding

$$H_i = \Phi_{\mathrm{enc}}(\hat{F}_i) \in \mathbb{R}^{N_i} \otimes \mathbb{R}^C \otimes \mathbb{R}^3. \qquad (6)$$

A lightweight segmentation head

$$\Psi_{\mathrm{seg}} : H_i \longrightarrow \hat{y}_i, \qquad \hat{y}_i \in [0,1]^{N_i},$$

predicts per-point boundary probabilities, supervised by binary ground-truth masks. The pretraining objective optimizes the composition $\Psi_{\mathrm{seg}} \circ \Phi_{\mathrm{enc}}$ to correctly identify fracture boundaries, forcing $\Phi_{\mathrm{enc}}$ to learn geometry–appearance interactions characteristic of fracture morphology. This procedure yields symmetry-preserving embeddings that serve as inputs to the subsequent flow-matching model on $SE(3)$.

### 4.3. Fragment Assembly

Given the encoded fragment features (Eq.6), $H_i = \Phi_{\mathrm{enc}}(\hat{F}_i)$, we enrich each fragment with geometric priors using a NeRF-style positional encoding. Specifically, we concatenate multi-frequency sinusoidal embeddings $\mathrm{PE}(\cdot)$ [24] of point cloud coordinates $\hat{F}_{\mathrm{xyz}}$, normals $\hat{F}_{\mathrm{n}}$, and scale information $s_i \in \mathbb{R}$ along the channel dimension. Let $\mathcal{H} := \mathbb{R}^{N_\star} \otimes \mathbb{R}^{C_{\mathrm{enc}}} \otimes \mathbb{R}^3$, where $N_\star$ is a fixed maximum number of points. The enriched fragment representation becomes:

$$H_i := f_{shape}\left(\big[H_i \,\|\, \mathrm{PE}((\hat{F}_{\mathrm{xyz}})_i) \,\|\, \mathrm{PE}((\hat{F}_{\mathrm{n}})_i) \,\|\, \mathrm{PE}(s_i)\big]\right),$$

where $f_{shape}(\cdot)$ is the shape embedding function. This provide high-frequency geometric cues that improve pose estimation. We additionally introduce a learned *anchor token* that designates the reference fragment around which all poses are predicted. Conditioned on $H_i$, we define a conditional Riemannian flow-matching model on the Lie group $SE(3)$ which learns a continuous velocity field whose integration yields the rigid transformation $g_i = (R_i, \beta_i) \in SE(3)$ required to place the fragment $i$ in its correct location within the assembled object. Formally, the model parametrizes a time-dependent velocity field

$$v : SE(3) \times [0,1] \times \mathcal{H} \longrightarrow \mathfrak{se}(3) = \mathfrak{so}(3) \times \mathbb{R}^3,$$

which associates to each element $g \in SE(3)$, time $t \in [0,1]$, and conditioning features $H_i \in \mathcal{H}$ a tangent vector $v(g,t;H_i) \in \mathfrak{se}(3)$. During training, the flow matching objective aligns this vector field with the target probability path interpolating between an initial distribution

$$p_0(g) = \mathcal{U}(SO(3)) \otimes \mathcal{N}(\mathbf{0}, I_3)$$

and the empirical distribution of assembled configurations. At inference time, the learned vector field is integrated over the unit interval as

$$\dot{g}_i^t = v\big(g_i^t, t; H_i\big), \qquad g_i^0 = (I, \mathbf{0}),$$

yielding the final transformations $g_i^1$ reassembling the object:

$$\mathcal{O} = \big\{\, g_i^1 \cdot \hat{F}_i \,\big\}_{i=1}^M.$$

This formulation defines a continuous, symmetry-preserving flow on $SE(3)$ that transports fragments from their initial centered configuration to globally consistent positions in the assembled object. In detail, the conditional trajectory between an initial pose $g_i^0$ and the target pose $g_i^1$ follows a geodesic interpolation in rotation and a linear interpolation in translation:

$$R_i^t = \exp_{R_i^0}\big(t \log_{R_i^0}(R_i^1)\big), \quad \beta_i^t = (1-t)\,\beta_i^0 + t\,\beta_i^1.$$

Accordingly, let us denote the rotational and translational velocity fields as

$$v_R^{(i)}(g_i^t, t) \in \mathfrak{so}(3), \qquad v_\beta^{(i)}(g_i^t, t) \in \mathbb{R}^3.$$

Following the flow-matching formulation, these predicted velocities are trained to match the target geodesic residuals scaled by $(1-t)^{-1}$. Hence, the objective is expressed as:

$$\mathcal{L}_{\mathrm{flow}} = \mathbb{E}_{t,\, p_1(g_i^1),\, p_t(g_i^t|g^1)}\left[\frac{1}{N}\sum_{i=1}^N \left(\Big\|v_\beta^{(i)}(g_i^t, t) - \tfrac{\beta_i^1 - \beta_i^t}{1-t}\Big\|_2^2\right.\right.$$

$$\left.\left. + \lambda \Big\|v_R^{(i)}(g^t, t) - \tfrac{1}{1-t}\log_{R_i^t}\big(R_i^1\big)\Big\|_2^2\right)\right], \qquad (7)$$

where $\log_{R_i^t}\left(R_i^1\right) \in \mathfrak{so}(3) \cong \mathbb{R}^3$ denotes the Lie-algebra residual from $R_i^t$ to $R_i^1$, and $\lambda$ balances rotational and translational terms. This formulation ensures that predicted fragment transformations follow the correct geodesic trajectories on $SE(3)$, maintaining smooth and consistent rotational and translational flows.

## 4.4. No-Overlap Loss

To prevent implausible fragment intersections, we introduce a differentiable no-overlap loss. For each pair of fragments $(i, j)$, we compute soft occupancy masks $M_i(x), M_j(x) \in [0, 1]$ denote the soft occupancy values of fragments $i$ and $j$ at volume cell (grid location) $x$ in a regular volumetric grid. We define the pairwise IoU:

$$\text{IoU}_{ij} = \frac{\sum_x \min\left(M_i(x), M_j(x)\right)}{\sum_x \max\left(M_i(x), M_j(x)\right) + \varepsilon}, \qquad (8)$$

where $\varepsilon > 0$ ensures numerical stability. The no-overlap loss aggregates these pairwise terms:

$$\mathcal{L}_{\text{no-overlap}} = \frac{1}{|P|} \sum_{(i,j) \in P} \text{IoU}_{ij}. \qquad (9)$$

Minimizing $\mathcal{L}_{\text{no-overlap}}$ enforces non-overlap between fragments, improving visual quality and physical plausibility.

The full loss combines the flow-matching term and the no-overlap penalty:

$$\mathcal{L} = \mathcal{L}_{\text{flow}} + \alpha\, \mathcal{L}_{\text{no-overlap}}, \qquad (10)$$

where $\alpha$ balances pose accuracy with overlap prevention. Together, the two-stage training and differentiable no-overlap loss enable the network to produce $SE(3)$ predictions that are both accurate and physically consistent.

## 5. Experimental Evaluation

We evaluate our method on four different dataset. The datasets, evaluation metrics, and competing methods are described in Section 5.1, while Section 5.2 reports the performance of all methods along with some qualitative results. In Section 5.3, we present an ablation study examining how E-M3RF generalizes across different objects and datasets[1].

### 5.1. Dataset and Evaluation Metrics

**Datasets.** Similarly to the evaluation protocol described in [12], we use a combination of large-scale synthetic and real scanned datasets to probe complementary aspects of 3D assembly. For the colorless (geometry-only) setting, we us (i) Breaking Bad [22], a large-scale synthetic dataset of over a million fractured instances derived from around 10k base shapes (PartNet/Thingi10k), which provides clean ground-truth poses and diverse geometry for controlled ablations and large-scale pretraining; and (ii) Fantastic Breaks [10], offering paired scans of intact and broken real objects that capture realistic micro-fracture morphology and scanning artifacts for high-fidelity geometric evaluation. For the colored fragments (multimodal setting), we evaluate on RePAIR [26] and the Presious [23] datasets. The former is a collection of archaeological fresco fragments with high-resolution RGB imagery and 3D geometry, and the latter is a set of several 3D cultural heritage fragment groups. Both are well suited to assessing color-based cues and multimodal fusion challenges, like erosion and missing-parts. Due to the size of the Presious dataset, *i.e.* only six sets, we use it only for testing the generalization accuracy of E-M3RF. We also reference FRACTURA[2] [12], curated with archaeologists includes different types of fragments for evaluating generalization to realistic breakage. Accordingly, we report the all "with color" results only on RePAIR and Presious datasets, while geometry-only ablations and generalization studies are conducted on all aforementioned datasets. Following [12, 29], we sample 5k points per object, ensuring uniform per-fragment density.

**Metrics.** We report four measures: **RMSE** ($\mathcal{R}^\circ$) is the root-mean-square rotation error (in degrees); **RMSE** ($\mathcal{T}_{\text{mm}}$) is the root-mean-square translation error (in millimeters); **PA** (Part Accuracy) is the percentage of fragments whose per-fragment Chamfer Distance to the ground-truth placement falls below 0.01%; and **CD** is the Chamfer Distance between the assembled object point cloud and the ground-truth assembly. We keep the same definitions and threshold as in [12] to ensure a fair comparison across datasets and baselines. The best performance metric is highlighted in **bold**, while the second-best competitor is <u>underlined</u>.

### 5.2. Results

Table 1 reports a head-to-head comparison on the color-enabled RePAIR dataset. The full E-M3RF model (last row) achieves the best score on all four metrics. It reduces rotation and translation errors by 23.1% and 13.2% respectively, and Chamfer Distance by 18.4%, while increasing Part Accuracy by 15.3% relative to the second best competitor model, GARF [12]. The margins over the rest geometric based methods are substantially larger across all metrics, underscoring that a geometry-only pipeline leaves disambiguating information untapped on colored fragments.

E-M3RF exploits its superiority mainly via the dense per-point color features fused with fracture-aware geometry, while its $SO(3)$-equivariant backbone makes the learned features transform consistently under pose, improving rotation estimates in particular. The non-overlap loss further suppresses physically implausible interpenetrations,

---

[1]Further ablations—including point sampling and runtime—appear in the suppl. material.

[2]Not publicly available at submission; excluded from our experiments.

Table 1. Quantitative results on RePAIR dataset. See Figure 4 for corresponding qualitative results.

| Method | RMSE ($\mathcal{R}°$)↓ | RMSE ($\mathcal{T}_{mm}$)↓ | PA (%)↑ | CD↓ |
|---|---|---|---|---|
| DiffAssemble [21] | 69.54 | 67.96 | 17.92 | 4.18 |
| PMTR [11] | 76.72 | 71.40 | 12.24 | 5.23 |
| PF++ [29] | 67.75 | 70.49 | 20.82 | 18.97 |
| GARF [12] | 49.31 | 32.19 | 31.14 | 2.66 |
| **E-M3RF** (w/o color) | 40.17 | 29.01 | 33.71 | 2.98 |
| **E-M3RF** (w/o non-overl.) loss | 45.81 | 25.01 | 32.91 | 3.83 |
| **E-M3RF** (w/ rotation-inv.) | 43.27 | 41.60 | 32.50 | 3.71 |
| **E-M3RF** (w/o rotation-equiv.) | 43.22 | 29.28 | 32.26 | 3.07 |
| **E-M3RF** | **37.91** | **27.93** | **35.91** | **2.17** |

reflected in the lower CD despite tighter placement (higher PA). This is corroborated by the ablations that we have conducted, where we examine the affect of each of our contributions to the overall pipeline:

- **w/o color:** Removing the color branch causes E-M3RF to rely solely on geometric features. PA decreases by 2.2 percentage points (pp), CD increases by 37%, rotation error increases by 6%, and translation error increases by 3.9% in comparison to the full model. These results indicate that color provides complementary cues to geometry, helping the model distinguish similar surfaces and resolve ambiguities geometry alone cannot.
- **w/o non-overlap loss:** Removing the non-overlap loss causes the largest increase in CD (76%) and a decrease in PA by 3.0 pp. Interestingly, translation RMSE decreases (27.93→25.01 mm). However, this improvement is misleading: without the non-overlap loss, fragments interpenetrate, yielding closer centroids but worse CD and PA. These results confirm that the non-overlap loss is essential for plausible placements and accurate surface alignment, as it can be clearly seen in Fig. 4 where in none of our solutions there is an overlap.
- **w/ rotation-inv.:** We replace the rotation-equivariant backbone with a rotation-invariant one, removing the inductive bias toward learning orientation-consistent features. Rotation RMSE increases by 14.1%, translation RMSE rises by 48.9%, PA decreases by 9.5 pp, and CD worsens by 71.0%. These results show that rotation-invariance degrades performance, as it removes orientation cues essential for alignment and thus the model cannot distinguish rotated inputs, resulting in larger errors.
- **w/o rotation-equiv.:** We replace the rotation-equivariant backbone with a standard architecture lacking any inductive bias toward orientation consistency. CD increases most (29%), while rotation error increases by 12%, PA decreases by 3.6 pp, and translation error increases by 5%. These results highlight that rotation-equivariant features act as a crucial inductive bias for this task, enabling consistent reasoning under arbitrary orientations and improving alignment stability and generalization.

Table 2 summarizes the results on the Breaking Bad

Table 2. Quantitative results on Breaking Bad dataset.

| Methods | RMSE ($\mathcal{R}°$)↓ | RMSE ($\mathcal{T}_{mm}$)↓ | PA (%)↑ | CD↓ |
|---|---|---|---|---|
| *Tested on the Everyday Subset* | | | | |
| Global [13] | 80.50 | 14.60 | 28.70 | 13.00 |
| LSTM [30] | 82.70 | 15.10 | 27.50 | 13.30 |
| DGL [31] | 80.30 | 13.90 | 31.60 | 11.80 |
| Jigsaw [16] | 42.19 | 6.85 | 68.89 | 8.22 |
| PMTR [11] | 31.57 | 9.95 | 70.60 | 5.56 |
| PF++ [29] | 35.61 | 6.05 | 76.17 | 2.78 |
| GARF-mini [12] | 6.68 | 1.34 | 94.77 | 0.25 |
| GARF [12] | 6.10 | 1.22 | 95.33 | 0.22 |
| **E-M3RF** (w/o color) | **5.31** | **1.14** | **96.20** | **0.18** |
| *Tested on the Artifact Subset* | | | | |
| Jigsaw [16] | 43.75 | 7.91 | 65.12 | 8.50 |
| PF++ [29] | 47.03 | 10.63 | 57.97 | 8.24 |
| GARF-mini [12] | 7.67 | 1.77 | 93.34 | 0.81 |
| GARF [12] | 5.82 | 1.27 | 95.04 | 0.42 |
| **E-M3RF** (w/o color) | **4.63** | **1.07** | **96.80** | **0.20** |

dataset. Because this dataset lacks color information, we evaluate E-M3RF with the geometry stream only of the architecture. On the Everyday subset, our method achieves the best overall performance, reducing the rotational error by 13.0%, the translational error by 6.6% and the CD by 18.2% compared to the strongest baseline. The PA improves by 0.87 pp. On the Artifact subset, our method again leads on all metrics, yielding up to 16.3% and 15.8% lower rotational and translation error respectively and a massive 35.5% reduction in CD while PA improves by 1.56 pp. On average across the two subsets, the improvements correspond to roughly 15% lower rotation error, 11% lower translation error, 27% lower CD, and 1.2 pp PA compared to the second best competitor model, GARF.

## 5.3. Ablation Studies

We assess E-M3RF 's zero-shot generalization on Presious (color+geometry) and Fantastic Breaks (geometry-only), and fine-tuning from RePAIR and Breaking Bad for the respective scenarios.

Table 3 shows the results on Fantastic Breaks a realscan, geometry-only benchmark, indicating that E-M3RF's equivariant geometry backbone successfully generalizes beyond the training distribution. Despite the lack of color, our method outperforms [12] in rotation and translation by 25.0%, while being on par for the other metrics.

Table 3. Quantitative results on Fantastic Breaks dataset. This includes manually collected real-world objects.

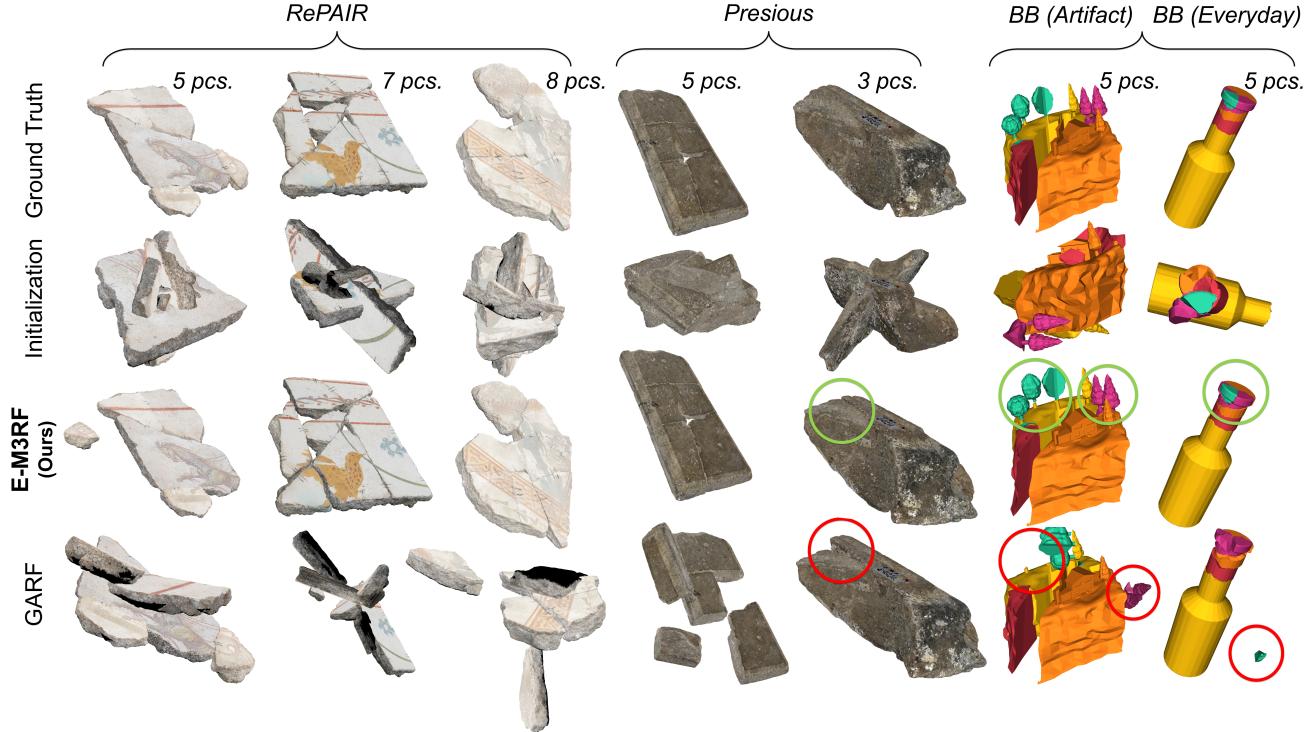| Methods | RMSE ($\mathcal{R}°$)↓ | RMSE ($\mathcal{T}_{mm}$)↓ | PA (%)↑ | CD↓ |
|---|---|---|---|---|
| Jigsaw [16] | 26.30 | 6.43 | 73.64 | 10.47 |
| PF++ [29] | 20.68 | 4.37 | 83.33 | 6.68 |
| GARF [12] | 10.62 | 2.10 | 91.00 | **2.12** |
| **E-M3RF** (w/o color) | **7.96** | **2.02** | **92.01** | 2.19 |

Figure 4. **Qualitative Comparisons on the RePAIR, Presious and Breaking Bad (BB).** E-M3RF consistently produces more accurate re-assemblies. Especially on the Presious scenes, is demonstrating strong generalization to unseen object. Green circles denote fine, ambiguous contact regions correctly recovered by our method. Additional results are available in the supplementary material.

Without any training or fine-tuning on Presious, Table 4 shows that E-M3RF preserves the same trends observed in Table 1. The full model consistently outperforms across all metrics except the Chamfer Distance where it follows by close, while again each ablation results in the expected performance degradation. Removing color (w/o color) lowers PA and raises CD (appearance cues no longer help disambiguate matches); removing the non-overlap loss (w/o non-overlapp. loss) increases CD by allowing overlapping even when poses look plausible; replacing the rotation-equivariant backbone with a rotation-invariant one (w/ rotation-inv) negatively impacts both rotation and translation performance; removing rotation equivariance (w/ rotation-equiv) primarily hurts rotation accuracy and slightly destabilizes translation. Geometry-only baselines such as GARF trail the full E-M3RF, that our color with an equivariant backbone fusion transfers reliably to previously unseen data without any dataset-specific training, evidencing strong out-of-distribution generalization, see Fig. 4.

## 6. Conclusion

We introduce E-M3RF, a multimodal 3D reassembly framework that fuses fracture-aware geometry via an $SO(3)$-equivariant backbone with color cues, and enforces phys-

Table 4. Quantitative results on 3D puzzle solving generalization on Presious dataset (model trained on RePAIR).

| Method | RMSE ($\mathcal{R}°$)↓ | RMSE ($\mathcal{T}_{mm}$)↓ | PA (%)↑ | CD↓ |
|---|---|---|---|---|
| DiffAssemble [21] | 56.65 | 176.41 | 24.72 | 37.73 |
| PMTR [11] | 80.01 | 68.33 | 14.27 | 18.21 |
| PF++ [29] | 70.61 | 146.84 | 24.79 | 41.58 |
| GARF [12] | 43.33 | 14.24 | 50.27 | **18.01** |
| **E-M3RF** (w/o color) | 35.12 | 13.85 | 52.91 | 21.18 |
| **E-M3RF** (w/o non-overlapp.) loss | 36.12 | 11.91 | 50.05 | 21.43 |
| **E-M3RF** (w/ rotation-inv.) | 41.26 | 45.01 | 47.12 | 24.62 |
| **E-M3RF** (w/o rotation-equiv.) | 39.67 | 15.72 | 46.11 | 23.71 |
| **E-M3RF** | **30.26** | 12.85 | **57.49** | 20.95 |

ical plausibility through a non-overlap loss. Evaluated on two real-world fresco collections with colored fractures and two synthetic, without color benchmarks, E-M3RF consistently outperforms all competing methods. These results highlight the potential of E-M3RF for robust, generalizable 3D reconstruction in both synthetic and real-world settings.
**Limitations & Future Work.** Point count per object affects accuracy (see suppl. material), but higher counts raise compute and memory cost. We will prioritize scalability via linearized attention and other efficiency-oriented architectures. A further limitation is our point-cloud input, which can miss high-resolution geometric and appearance detail; to address this, we plan a mesh-first variant using triangle meshes and texture maps to exploit connectivity and perface cues.

# Acknowledgements

# References

[1] Joey Bose, Tara Akhound-Sadegh, Guillaume Huguet, Kilian FATRAS, Jarrid Rector-Brooks, Chenghao Liu, Andrei Nica, Maksym Korablyov, Michael Bronstein, and Alexander Tong. Se(3)-stochastic flow matching for protein backbone generation. In *International Conference on Representation Learning*, pages 22590–22621, 2024. 2

[2] Ricky T. Q. Chen and Yaron Lipman. Flow matching on general geometries. In *The Twelfth International Conference on Learning Representations*, 2024. 2

[3] Gabriele Corso, Hannes Stärk, Bowen Jing, Regina Barzilay, and Tommi Jaakkola. Diffdock: Diffusion steps, twists, and turns for molecular docking. *International Conference on Learning Representations (ICLR)*, 2023. 1

[4] Congyue Deng, Or Litany, Yueqi Duan, Adrien Poulenard, Andrea Tagliasacchi, and Leonidas J Guibas. Vector neurons: A general framework for so (3)-equivariant networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12200–12209, 2021. 3

[5] Yu Deng, Or Litany, Yueqi Duan, Yu Qiao, and Hao Su. Vntransformer: Rotation-equivariant attention for point cloud processing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 3, 4

[6] Qi Fan, Shaoshuai Sun, Ping Chen, Xiaolong Zhang, Jiajun Wang, Jianbo Wang, Gao Huang, et al. Swin3d: A pretrained transformer backbone for 3d indoor scene understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3

[7] Herbert Freeman and L Garder. Apictorial jigsaw puzzles: The computer solution of a problem in pattern recognition. *IEEE Transactions on Electronic Computers*, (2):118–127, 2006. 1

[8] Meng-Hao Guo, Jun-Xiong Cai, Zhong-Qi Liu, Tong Mu, Ralph R Martin, and Shi-Min Hu. Pct: Point cloud transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPR-W)*, 2021. 3

[9] Adeela Islam, Stefano Fiorini, Stuart James, Pietro Morerio, and Alessio Del Bue. Reassemblenet: Learnable keypoints and diffusion for 2d fresco reconstruction. *ArXiv*, abs/2505.21117, 2025. 1, 2

[10] Nikolas Lamb, Cameron Lowell Palmer, Benjamin Molloy, Sean Banerjee, and Natasha Kholgade Banerjee. Fantastic breaks: A dataset of paired 3d scans of real-world broken objects and their complete counterparts. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4681–4691, 2023. 2, 6

[11] Nahyuk Lee, Juhong Min, Junha Lee, Seungwook Kim, Kanghee Lee, Jaesik Park, and Minsu Cho. 3d geometric shape assembly via efficient point cloud matching. *ArXiv*, abs/2407.10542, 2024. 7, 8, 3, 4

[12] Sihang Li, Zeyu Jiang, Grace Chen, Chenyang Xu, Siqi Tan, Xue Wang, Irving Fang, Kristof Zyskowski, Shannon J. P. McPherron, Radu Iovita, Chen Feng, and Jing Zhang. Garf: Learning generalizable 3d reassembly for real-world fractures. *ArXiv*, abs/2504.05400, 2025. 1, 2, 3, 6, 7, 8, 4

[13] Yichen Li, Kaichun Mo, Lin Shao, Minhyuk Sung, and Leonidas Guibas. Learning 3d part assembly from a single image. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VI 16*, pages 664–682. Springer, 2020. 7

[14] Sheng-hui Liao, Chao Xiong, Shu Liu, Ying-qi Zhang, and Chun-lin Peng. 3d object reassembly using region-pair-relation and balanced cluster tree. *Computer Methods and Programs in Biomedicine*, 197:105756, 2020. 1

[15] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*, 2023. 2

[16] Jiaxin Lu, Yifan Sun, and Qi-Xing Huang. Jigsaw: Learning to assemble multiple fractured objects. *ArXiv*, abs/2305.17975, 2023. 2, 7, 3, 4

[17] Xiaolei Niu, Qifeng Wang, B. Liu, and Jianxin Zhang. An automatic chinaware fragments reassembly method framework based on linear feature of fracture surface contour. *ACM Journal on Computing and Cultural Heritage*, 16:1 – 22, 2022. 2

[18] Thiago M Paixao, Rodrigo F Berriel, Maria Boeres, Alessandro L Koerich, Claudine Badue, Alberto F De Souza, and Thiago Oliveira-Santos. Fast (er) reconstruction of shredded text documents via self-supervised deep asymmetric metric learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14343–14351, 2020. 1

[19] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 3

[20] Charles R Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. 3

[21] Gianluca Scarpellini, Stefano Fiorini, Francesco Giuliari, Pietro Moreiro, and Alessio Del Bue. Diffassemble: A unified graph-diffusion model for 2d and 3d reassembly. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 28098–28108, 2024. 1, 2, 3, 7, 8, 4

[22] Silvia Sell'an, Yun-Chun Chen, Ziyi Wu, Animesh Garg, and Alec Jacobson. Breaking bad: A dataset for geometric fracture and reassembly. *ArXiv*, abs/2210.11463, 2022. 1, 2, 3, 6

[23] Theoharis T. and G. Papaioannou. PRESIOUS 3d cultural heritage fragments, 2013. 2, 6

[24] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *Advances in neural information processing systems*, 33:7537–7547, 2020. 5

[25] Hugues Thomas, Charles R Qi, Jean-Emmanuel Deschaud, Beatriz Marcotegui, François Goulette, and Leonidas J Guibas. Kpconv: Flexible and deformable convolution for point clouds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. 3

[26] Theodore Tsesmelis, Luca Palmieri, Marina Khoroshiltseva, Adeela Islam, Gur Elkin, Ofir Itzhak Shahar, Gianluca Scarpellini, Stefano Fiorini, Yaniv Ohayon, Nadav Alali, Sinem Aslan, Pietro Morerio, Sebastiano Vascon, Elena Gravina, Maria Cristina Napolitano, Giuseppe Scarpati, Gabriel Zuchtriegel, Alexandra Spuhler, Michel E. Fuchs, Stuart James, Ohad Ben-Shahar, Marcello Pelillo, and Alessio Del Bue. Re-assembling the past: The repair dataset and benchmark for real world 2d and 3d puzzle solving. *ArXiv*, abs/2410.24010, 2024. 1, 2, 6, 4

[27] Jie Wang, Congyi Zhang, Peng Wang, X. Li, Peter J. Cobb, Christian Theobalt, and Wenping Wang. Batch-based model registration for fast 3d sherd reconstruction. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 14473–14483, 2022. 2

[28] Yue Wang and Justin Solomon. Dynamic graph cnn for learning on point clouds. *ACM Transactions on Graphics (TOG)*, 38(5):1–12, 2019. 3

[29] Zhengqing Wang, Jiacheng Chen, and Yasutaka Furukawa. Puzzlefusion++: Auto-agglomerative 3d fracture assembly by denoise and verify. *ArXiv*, abs/2406.00259, 2024. 2, 3, 6, 7, 8, 4

[30] Rundi Wu, Yixin Zhuang, Kai Xu, Hao Zhang, and Baoquan Chen. Pq-net: A generative part seq2seq network for 3d shapes. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 829–838, 2020. 7

[31] Guanqi Zhan, Qingnan Fan, Kaichun Mo, Lin Shao, Baoquan Chen, Leonidas J Guibas, Hao Dong, et al. Generative 3d part assembly via dynamic graph learning. *Advances in Neural Information Processing Systems*, 33:6315–6326, 2020. 7

[32] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 3

[33] SY Zheng, RY Huang, J Li, and Z Wang. Reassembling 3d thin fragments of unknown geometry in cultural heritage. *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, 2:393–399, 2014. 1

# E-M3RF: An Equivariant Multimodal 3D Re-assembly Framework

## Supplementary Material

This document expands the main paper along three areas.

**Section A (Groups, Symmetries and Equivariance).** We formalize the symmetry principles behind our model where a fragment collection is treated as an *unordered set* acted on by the symmetric group $S_N$, meanwhile geometry lives in $\mathbb{R}^3$ under rigid motions $SE(3)$. We justify that the learned tokenization and message passing are *permutation–equivariant* ($f(\pi X) = \pi f(X)$ for any $\pi \in S_N$) and that feature fields are *rotation–equivariant* ($f(RX) = R f(X)$ for $R \in SO(3)$), while scalar heads remain invariant. Short derivations, notation, and sanity checks are provided.

**Section B (Implementation & ablations).** We add training and evaluation details and report extended ablations, including points-per-fragment, memory/compute trade-offs, fracture-surface thresholds, and the effect of the non-overlap weight.

**Section C (Additional qualitative results).** We present more visual assemblies, side-by-side comparisons to the other baselines, and representative failure cases, complementing the quantitative and qualitative results presented in the main paper.

## A. Groups, Symmetries and Equivariance

In this section we recall the basic notions linking group theory, symmetries, and equivariance. Abstract group theory (Sec. A.1) provides an abstract framework for describing symmetries as transformations that leave certain structures invariant. Group *representations* (Sec. A.2) realize these abstract symmetries as linear transformations on vector spaces, enabling their analysis via linear algebra. Building on these concepts, equivariance (Sec. A.3) characterizes maps that are compatible with a given group action, in the sense that applying a symmetry before or after the map yields consistent results under some group action.

### A.1. Groups

**Definition A.1 (Group)** *A group is a pair $(G, \cdot)$ consisting of a set $G$ together with a binary operation $\cdot : G \times G \to G$ satisfying the following axioms:*
*(A1) Associativity: $(a \cdot b) \cdot c = a \cdot (b \cdot c)$ for all $a, b, c \in G$;*
*(A2) Identity: there exists an element $e \in G$ such that $e \cdot a = a \cdot e = a$ for all $a \in G$;*
*(A3) Inverses: for each $a \in G$, there exists an element $a^{-1} \in G$ satisfying $a^{-1} \cdot a = a \cdot a^{-1} = e$.*
*We typically write the operation and omit the symbol "$\cdot$" when unambiguous. If, in addition, the binary operation $\cdot$ is commutative, that is, if $a \cdot b = b \cdot a$ for all $a, b \in G$, we say*

*the group is abelian.*

We present the following groups, which formalize discrete and continuous symmetries of interest for the fragments introduced in the main paper (Sec. 4.1):

**Definition A.2 (Symmetric group $S_N$)** *For a finite set $[N] = 1, \dots, N$, the set of all bijections $\sigma : [N] \to [N]$ forms the* symmetric group *on $N$ elements,*

$$S_N = \{ \, \sigma : [N] \to [N] \mid \sigma \text{ is bijective } \, \},$$

*under composition of maps. The identity element is the identity map* id.

This notion extends naturally from discrete sets to linear spaces.

**Definition A.3 (General Linear Group GL(V))** *Let $V$ be a finite-dimensional real vector space. The* general linear group *of $V$ is*

$$GL(V) = \{ \, T : V \to V \mid T \text{ is linear and invertible} \, \},$$

*with composition as the group operation. Choosing a basis of $V$ identifies each $T \in GL(V)$ with an invertible matrix, yielding the matrix group*

$$GL_N(\mathbb{R}) = \{ \, A \in \mathbb{R}^{N \times N} \mid \det(A) \neq 0 \, \},$$

*where $N = \dim(V)$.*

**Definition A.4 (Special Orthogonal Group $SO(V)$)** *Let $V$ be an $N$-dimensional real inner-product space with inner product $\langle \cdot, \cdot \rangle$. The* special orthogonal group *of $V$ is*

$$SO(V) = \{ \, T \in GL(V) \mid \langle Tv, Tw \rangle = \langle v, w \rangle \\ \text{for all } v, w \in V, \det(T) = 1 \, \}.$$

*Its elements are the orientation-preserving linear isometries of $V$.*

This definition has no coordinates in it; the structure depends only on the inner product and the notion of orientation. A basis simply provides a concrete representation.

**Observation A.5** *Choosing an orthonormal basis of $V$ identifies $SO(V)$ with the matrix group*

$$SO(N) = \{ \, R \in \mathbb{R}^{N \times N} \mid R^\top R = I, \ \det(R) = 1 \, \},$$

*since linear isometries correspond exactly to matrices preserving the standard inner product.*

Specializing to three dimensions gives the familiar rotation group in Euclidean space.

**Definition A.6 (Special Orthogonal Group SO(3))** *The* special orthogonal group *in three dimensions is*

$$SO(3) = \{\, R \in \mathbb{R}^{3\times 3} \mid R^\top R = I,\ \det(R) = 1\,\},$$

*the group of proper rotations of $\mathbb{R}^3$.*

**Definition A.7 (Special Euclidean Group SE(V))** *Let $V$ be an $N$-dimensional real inner-product space. The* special Euclidean group *of $V$ is the group of orientation-preserving isometries $f : V \to V$. Each such isometry has a unique decomposition*

$$f(x) = Tx + t, \qquad T \in SO(V),\ t \in V,$$

*and composition is given by*

$$(T_1, t_1)(T_2, t_2) = (T_1 T_2,\ T_1 t_2 + t_1).$$

**Definition A.8 (Special Euclidean Group SE(3))**
*Specializing to $V = \mathbb{R}^3$ under a basis, one obtains*

$$SE(3) = \{(R, t) \mid R \in SO(3),\ t \in \mathbb{R}^3\},$$

*acting on $x \in \mathbb{R}^3$ by $x \mapsto Rx + t$, with the composition law*

$$(R_1, t_1)(R_2, t_2) = (R_1 R_2,\ R_1 t_2 + t_1).$$

**Observation A.9 (Semidirect-product form)** *The map*

$$SO(3) \ltimes \mathbb{R}^3 \ \longrightarrow\ SE(3), \qquad (R, t) \mapsto (x \mapsto Rx + t),$$

*is a group isomorphism. Thus $SE(3)$ is naturally identified with the semidirect product $SO(3) \ltimes \mathbb{R}^3$, where $SO(3)$ acts on $\mathbb{R}^3$ by its standard linear action.*

## A.2. Representations

Representations make concrete the notion of symmetry by describing how elements of a group act linearly on a vector space. Each group element is associated with an invertible linear transformation, transferring the group's algebraic structure into linear operators.

**Definition A.10 (Linear Group Action)** *Let $G$ be a group and $V$ a finite-dimensional real vector space. A* representation *of $G$ on $V$ is a group homomorphism*

$$\rho_V : G \to GL(V),$$

*which defines the action*

$$g \cdot v := \rho_V(g)\, v, \qquad g \in G,\ v \in V.$$

**Definition A.11 (Linear Group Action)** *Let $V$ be a finite-dimensional vector space, and let $G$ be a finite group. A* representation $V$ *or linear group action is a group homomorphism*

$$\rho^V : G \to GL(V),$$

*which defines the action*

$$\rho^V(g) \cdot v \quad \text{for all} \quad g \in G,\, v \in V.$$

The symmetric group (Definition A.2) admits a representation permuting coordinate entries.

**Definition A.12 (Standard Representation of $S_N$)** *Let $V = \mathbb{R}^N$ with canonical basis $\{e_i\}_{i=1}^N$. The* standard representation *of the symmetric group $S_N$ acts by permuting coordinates:*

$$(\sigma \cdot x)_i = x_{\sigma^{-1}(i)}, \qquad \forall \sigma \in S_N,\, x = [x_1, \ldots, x_N] \in \mathbb{R}^N.$$

*Equivalently, this action can be expressed as left multiplication by the permutation matrix $P_\sigma \in \mathrm{GL}_N(\mathbb{R})$, defined by*

$$(P_\sigma)_{ij} = \begin{cases} 1, & \text{if } i = \sigma(j), \\ 0, & \text{otherwise.} \end{cases}$$

*This realizes the standard representation as a concrete matrix representation of $S_N$.*

Continuous groups admit analogous linear actions. For $V = \mathbb{R}^3$, the inclusion $SO(3) \subset GL_3(\mathbb{R})$ defines the standard representation

$$R \cdot x = Rx, \qquad R \in SO(3),\ x \in \mathbb{R}^3,$$

given by left multiplication by rotation matrices.

## A.3. Equivariance and Invariance

Equivariance and invariance describe how maps between representations interact with group actions. They formalize when a transformation "respects" symmetry. Such description in exploited in section 4.1 where the Geometric and Color encoders are introduced.

**Definition A.13 (Equivariant Map)** *Let $(V, \rho_V)$ and $(W, \rho_W)$ be representations of a group $G$. A (not necessarily linear) map $f : V \to W$ is* G-equivariant *if*

$$f(\rho_V(g)\, v) = \rho_W(g)\, f(v) \qquad \text{for all } g \in G,\ v \in V.$$

*Equivariance means that applying a group action before or after $f$ yields the same outcome.*

**Definition A.14 (Invariant Map)** *Let $(V, \rho_V)$ and $(W, \rho_W)$ be representations of a group $G$. A (not necessarily linear) map $f : V \to W$ is* G-invariant *if*

$$f(\rho_V(g)\, v) = f(v) \qquad \text{for all } g \in G,\ v \in V.$$

*Invariance means that $f$ is unaffected by the action of $G$ on its input.*

# B. Additional Info and Ablation Studies

Here we present additional details related to E-M3RF's training and evaluation pipelines (Section B.1) and extended ablations on the key factors that drive accuracy and cost (Section B.2). We first study the number of points sampling, showing how increasing point count improves performance while exhibiting diminishing returns. We then quantify memory usage as point density rises. Next, we analyze fracture-surface ground-truth sensitivity on RePAIR by varying the distance threshold used in pretraining, which is important given the dataset's random erosion, and report its effect on downstream accuracy. Finally, we sweep the no-overlap loss weight $\alpha$ to characterize the trade-off between collision suppression (CD/intersections) and part accuracy (PA), selecting a stable default.

## B.1. Additional Info

**Competing Methods.** We compare against GARF [12], DiffAssemble [21], PMTR [11], Puzzle-Fusion++ (PF++) [29] and Jigsaw [16] as the most recent and representative state-of-the-art for multi-fragment 3D assembly on point clouds, and they collectively span the main methodological families: *(i)* generative diffusion/flow pipelines for pose refinement (DiffAssemble), *(ii)* matching-based pose estimation with efficient point-cloud correspondences (PMTR), *(iii)* search/verify agglomerative assembly (PF++ and Jigsaw), and *(iv)* a generalization-oriented geometry SOTA with strong results on real fractures (GARF).

## B.2. Ablation Studies

### B.2.1. Number of Keypoints & Memory

We compare E-M3RF with GARF [12] across different numbers of points on the RePAIR dataset (Table 5). E-M3RF consistently outperforms GARF across all metrics, achieving lower rotation and translation errors, higher part accuracy, and reduced Chamfer Distance.

Notably, E-M3RF exhibits a positive correlation between the number of points and performance, highlighting our method's ability to leverage denser point clouds for more accurate fragment alignment. In contrast, this trend is not observed in GARF, where adding more points does not necessarily aid its reconstruction capabilities. We attribute this behavior to the fact that GARF does not utilize color information. By effectively leveraging the richer color signals present in denser point clouds, E-M3RF maximizes reconstruction accuracy.

**Computational Cost.** It is important to note that higher point densities naturally demand greater computational resources. As shown in Figure 5, memory usage for both models increases approximately linearly with the number of

points. Although E-M3RF exhibits slightly higher memory consumption than GARF, primarily due to the additional color features, the peak usage remains comfortably within our hardware limits (NVIDIA A100, 40GB).

Table 5. Quantitative results on the RePAIR dataset with varying number of points (5K–20K). **Bold** indicates the best result, while <u>underlined</u> values denote the competitor best result.

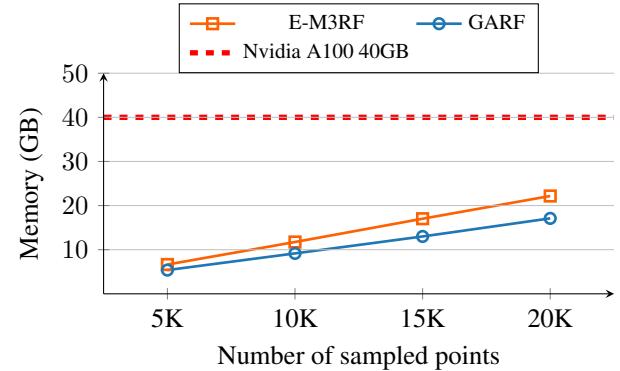| #Points | Method | RMSE ($\mathcal{R}°$)↓ | RMSE ($\mathcal{T}_{mm}$)↓ | PA (%)↑ | CD↓ |
|---|---|---|---|---|---|
| 5K | GARF [12] | 49.31 | 32.19 | 31.14 | 2.66 |
| | **E-M3RF** | 37.91 | 27.93 | 35.91 | 2.17 |
| 10K | GARF [12] | <u>43.92</u> | 30.12 | 32.55 | <u>2.52</u> |
| | **E-M3RF** | 35.82 | 26.01 | 37.15 | 2.12 |
| 15K | GARF [12] | 48.27 | 28.02 | 33.10 | 3.12 |
| | **E-M3RF** | 37.21 | 26.04 | 36.51 | **2.10** |
| 20K | GARF [12] | 45.56 | <u>27.94</u> | <u>34.21</u> | 3.01 |
| | **E-M3RF** | **34.92** | **25.03** | **37.21** | 2.11 |



Figure 5. GPU memory consumption on the RePAIR dataset as a function of the varying number of sampled points (5K–20K).

### B.2.2. Distance Threshold for Fracture-surface

The Fracture-Boundary Segmentation Pretraining task is introduced as a large-scale pretraining approach to learn fragment representations to effectively capture the geometric signatures of the fracture interfaces, crucial for accurate reassembly. The ground truth label for this task is derived by leveraging a distance threshold $\tau$ to identify contact points between two fragments, where $\tau$ controls how close points need to be to count as contact.

For the RePAIR dataset experiments reported in the main paper, we set the distance threshold to $0.4$. This value provides a balance between robustness to gaps caused by random erosion or scanning noise and precision in detecting true contacts. We note that smaller thresholds risk miss valid connections, while larger thresholds can falsely connect fragments that are not actually adjacent. In Table 6, we present an ablation study varying the distance threshold to analyze its effect on fragment connectivity and reconstruction. This experiment demonstrates how the choice of

threshold impacts the detection of shared surfaces and the resulting reassembly.

Table 6. Ablation study on fracture surface ground-truth generation for RePAIR [26]. Different distance thresholds for defining fracture-surface ground truth are evaluated.

| $\tau$ | RMSE ($\mathcal{R}°$)↓ | RMSE ($\mathcal{T}_{mm}$)↓ | PA (%)↑ | CD↓ |
|---|---|---|---|---|
| 0.2 | <u>40.15</u> | 57.22 | 27.02 | 4.14 |
| 0.4 | **37.91** | <u>27.93</u> | **35.91** | **2.17** |
| 0.6 | 41.02 | **27.42** | <u>33.67</u> | <u>2.18</u> |
| 0.8 | 42.71 | 34.31 | 31.29 | 3.01 |

### B.2.3. Alpha for No-Overlap Loss

In Table 7, we analyze the effect of the weighting factor $\alpha$, which balances assembly accuracy with overlap prevention. When $\alpha$ is too small, the network produces highly accurate assemblies but occasionally predicts overlapping parts. Conversely, a large $\alpha$ enforces physical consistency more strictly but can reduce assembly precision. Our experiments show that an intermediate value, specifically $\alpha = 0.3$, achieves the best trade-off, yielding $SE(3)$ predictions that are both accurate and physically feasible, as clearly demonstrated by the results in Table 7. This highlights the importance of the no-overlap loss term in achieving physically consistent training.

Table 7. Ablation study on the effect of the no-overlap loss weight ($\alpha$) on the RePAIR [26] dataset.

| $\alpha$ | RMSE ($\mathcal{R}°$)↓ | RMSE ($\mathcal{T}_{mm}$)↓ | PA (%)↑ | CD↓ |
|---|---|---|---|---|
| 0.1 | 43.26 | **25.92** | 33.07 | 3.42 |
| 0.3 | **37.91** | <u>27.93</u> | **35.91** | **2.17** |
| 0.5 | <u>37.98</u> | 28.22 | <u>35.02</u> | <u>3.01</u> |
| 0.7 | 40.12 | 33.75 | 31.05 | 4.21 |

## C. Additional Visualizations

This section presents qualitative results corresponding to the four evaluation datasets—RePAIR, Presious, Breaking Bad, and Fantastic Breaks—and illustrate how E-M3RF behaves under different appearance and fracture variations and across varying piece counts (note that Fantastic Breaks contains only 2-piece puzzles by design).

On the RePAIR dataset, Fig. 6, examples with several pieces show that our multimodal fusion aligns both fracture boundaries and colored motifs. On the other hand competing methods often leave slight pose drift or texture misalignment on thin or eroded edges.

On the Presious dataset, Fig. 7, despite the surface wear and uneven point density, our reconstructions minimize interpenetration and recover consistent contact along long,

low-curvature breaks; baselines tend to over- or under-insert pieces, leaving visible gaps or overlaps.

On the only geometry datasets and specifically on Breaking Bad, Fig. 8, With complex, multi-piece assemblies, the equivariant backbone maintains stable orientations and lowers rotational error, producing globally coherent placements where alternatives seem to get trapped in near-symmetric configurations.

Finally, on the Fantastic Breaks (2-piece only, and real only geometry related scans), Fig. 9, it seems that our no-overlap loss prevents subtle collisions while achieving tight, flush joins; while others show slight interpenetration or residual offsets.

Overall, the visuals echo the quantitative results: E-M3RF consistently achieves tighter pose alignment, fewer collisions, and—when the color cues are available—better cross-shard pattern continuity, while remaining robust on geometry-only data. Figures show qualitative comparisons of E-M3RF against GARF [12], DiffAssemble [21], PMTR [11], PF++ [29] and Jigsaw [16] on the four benchmarks.
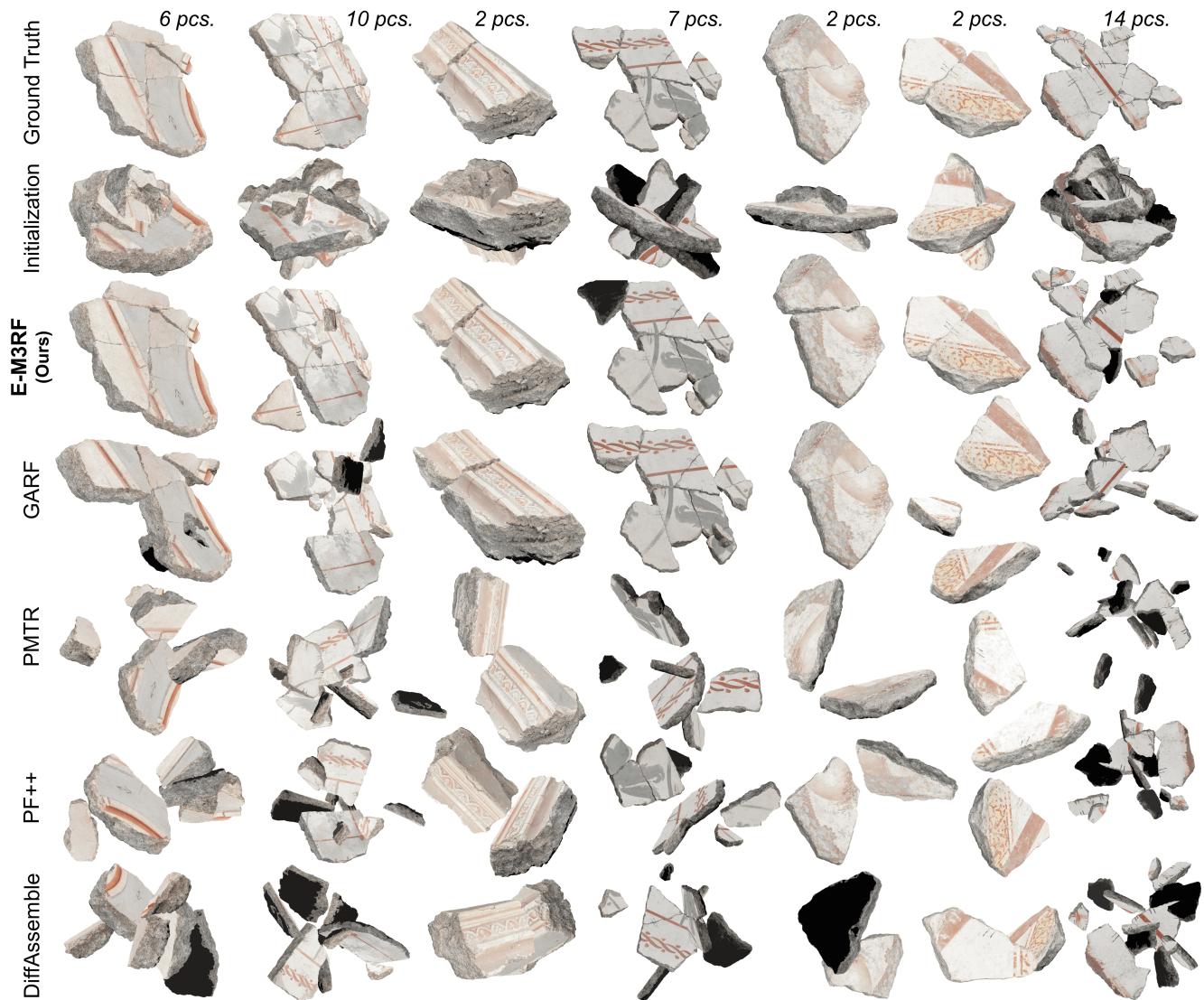
Figure 6. Qualitative Comparisons on the RePAIR. E-M3RF consistently produces more accurate re-assemblies.
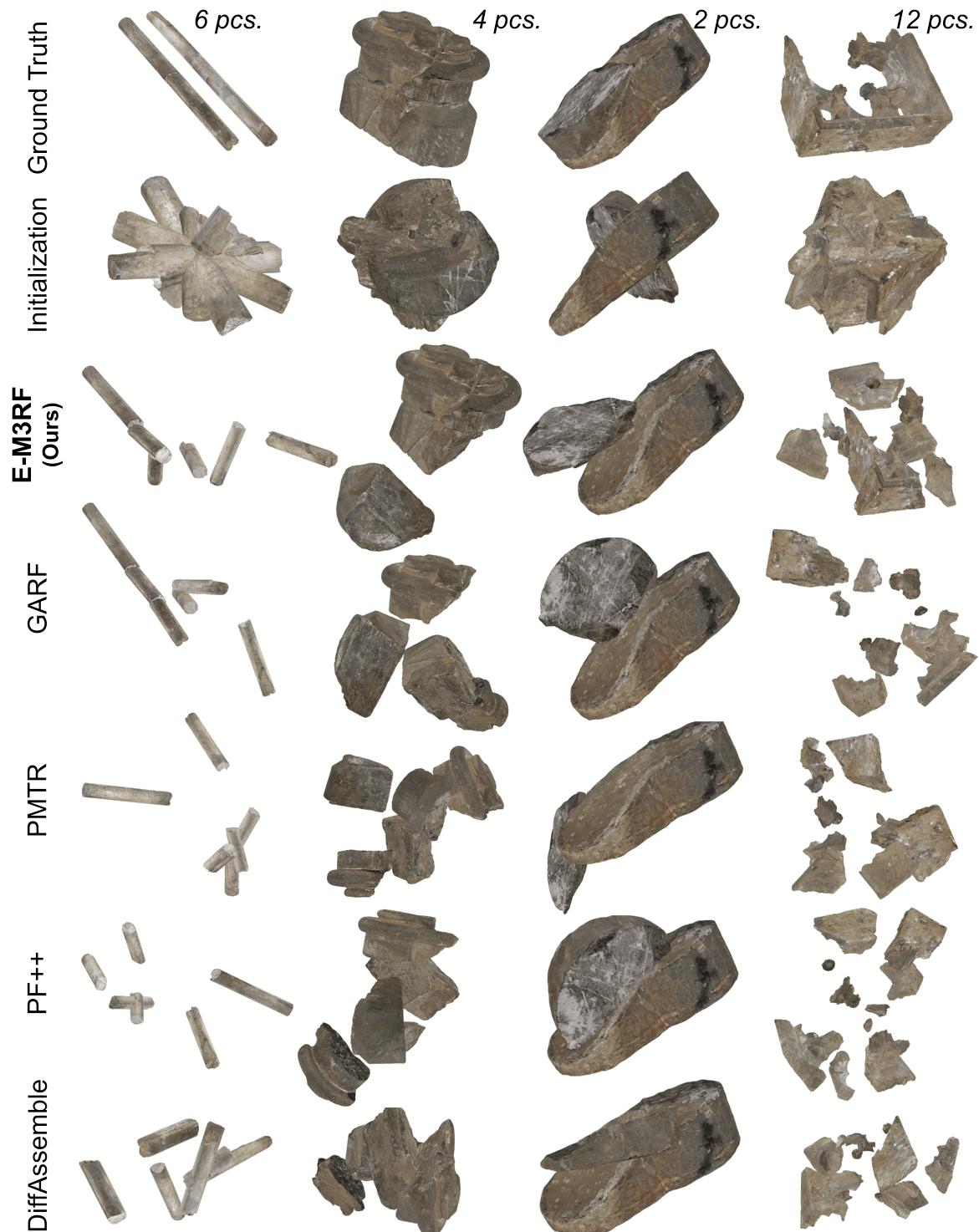
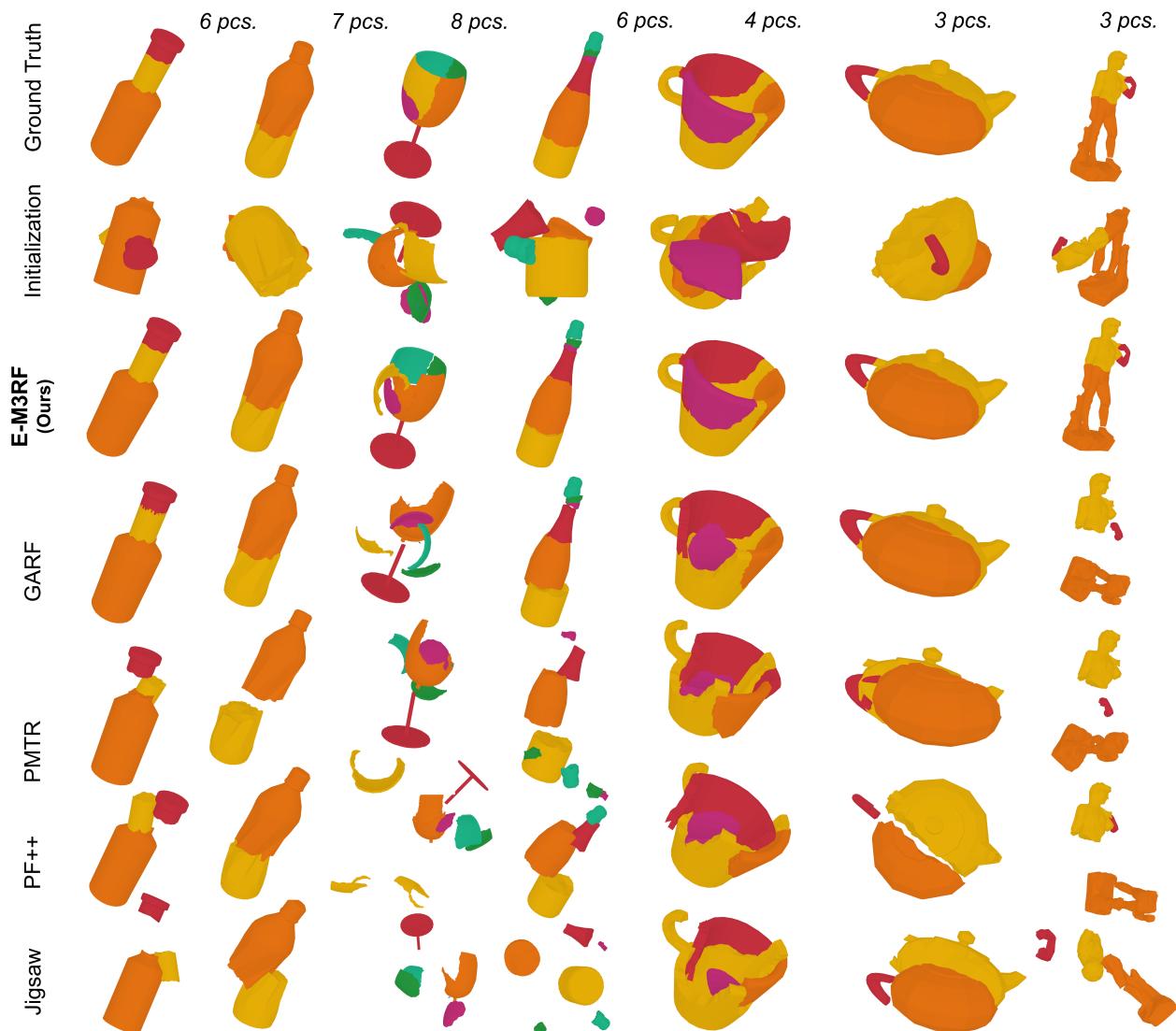Figure 7. Qualitative Comparisons on the remaining four puzzles of the Presious dataset.

Figure 8. Qualitative Comparisons on the Breaking Bad (BB). E-M3RF consistently produces more accurate re-assemblies.
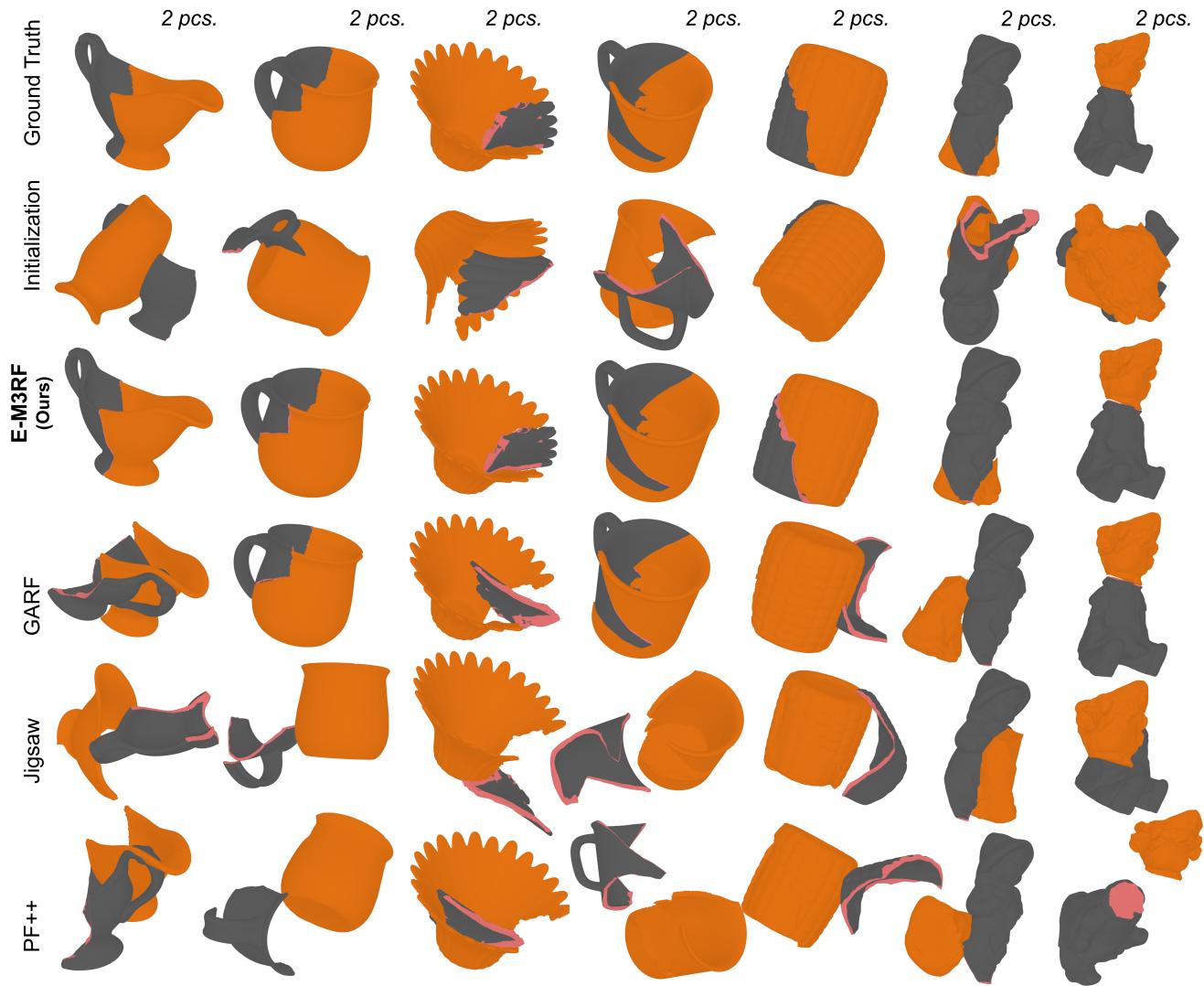
Figure 9. Qualitative Comparisons on the Fantastic Breaks. E-M3RF consistently produces more accurate re-assemblies. Note, the Fantastic Breaks dataset is composed of only 2-piece puzzles by design.