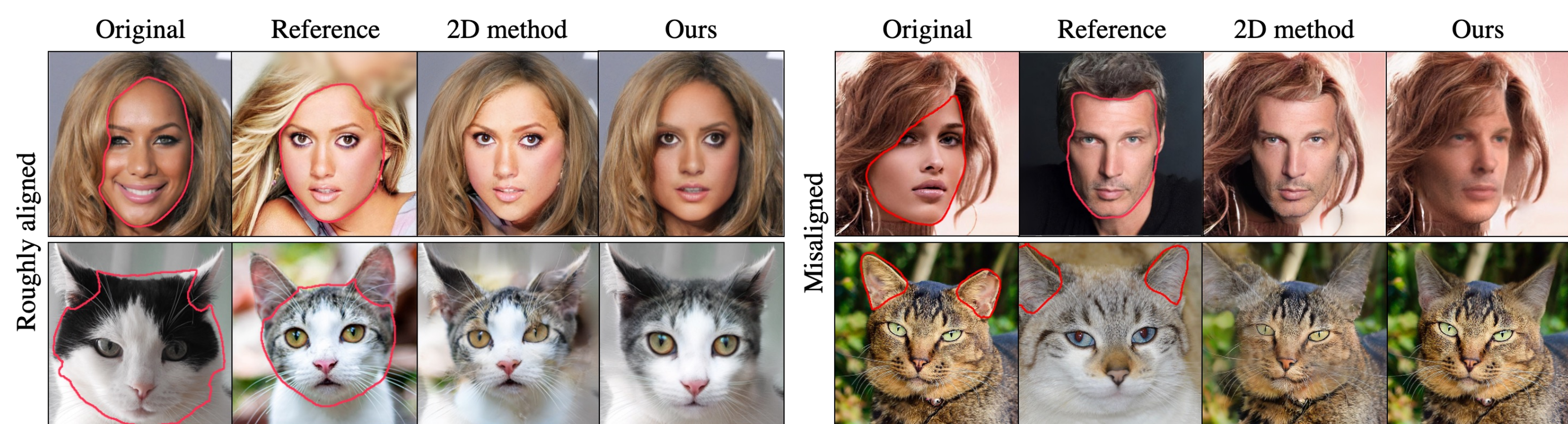


Abstract



Challenge: Seamless blending of images, especially with misalignments from camera poses and object shapes

Solution: 3D-aware blending using generative Neural Radiance Fields (NeRF)

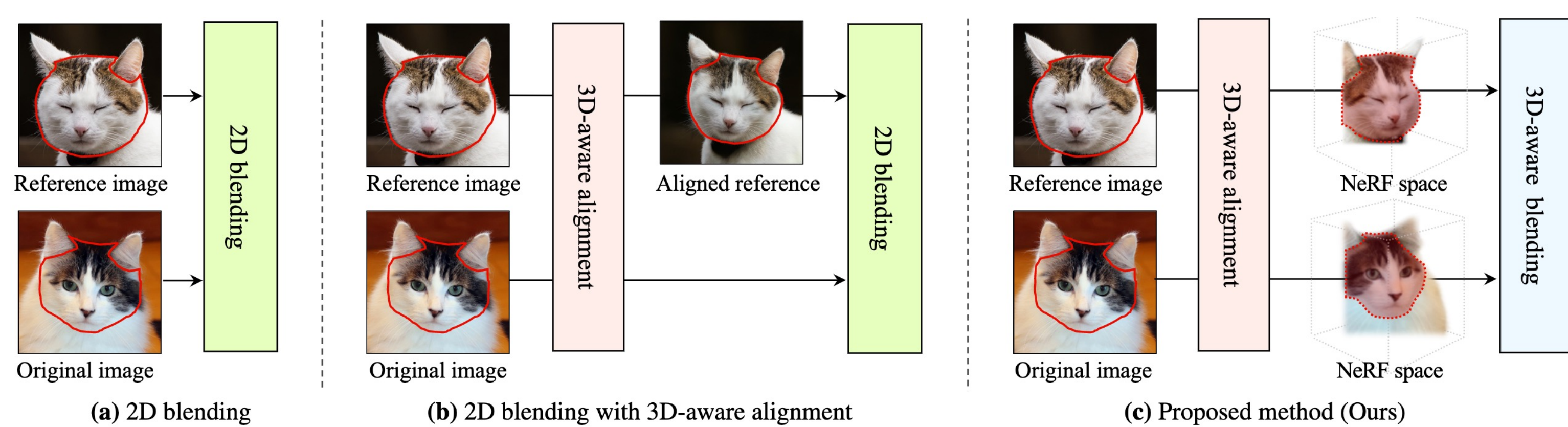
• 3D-aware Alignment:

- Estimate camera poses of the input images
- Perform pose alignment for objects

• 3D-aware Blending:

- Utilizes volume density rather than raw pixel space only
- Blend images in NeRF's latent representation space

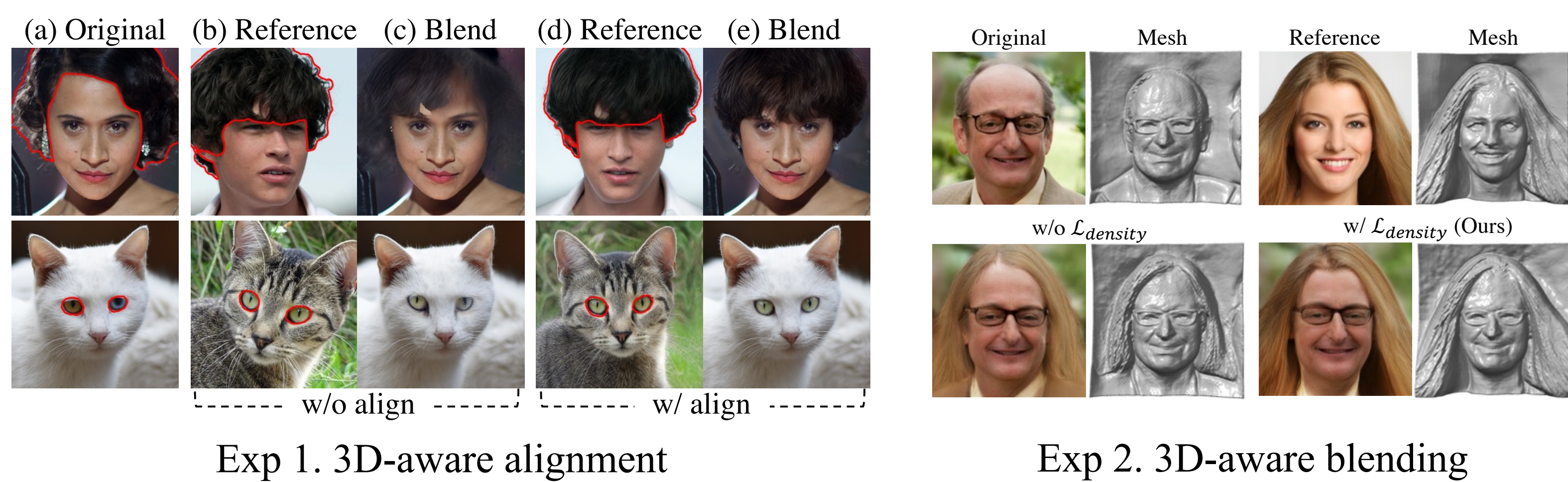
Comparison with Baselines



Red lines denote target blending parts.

- (a) **2D blending:** 2D blending methods compose two images without any 3D-aware alignment.
- (b) **2D blending with 3D-aware alignment:** To address misalignment, we apply our 3D-aware alignment method to existing 2D blending methods.
- (c) **Proposed method:** We propose 3D-aware blending after applying our 3D-aware alignment. Note that all methods do not use 3D labels or 3D morphable models.

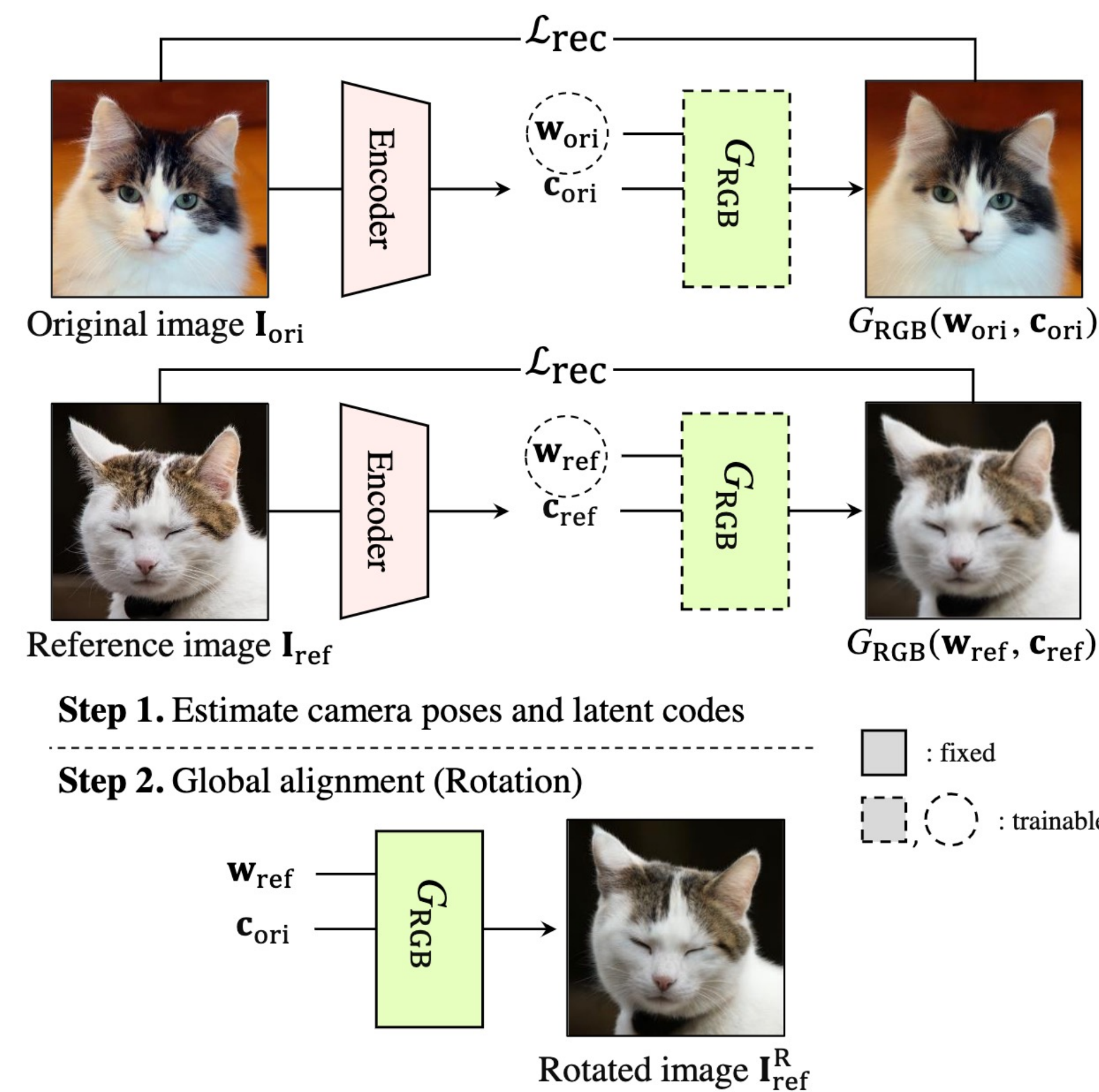
Ablation studies



Quantitative results

Method	w/o align (baseline only)			w/ 3D-aware align		
	KID ↓	LPIPS _m ↓	mL ₂ ↓	KID ↓	LPIPS _m ↓	mL ₂ ↓
Poisson Blending [72]	0.006	0.4203	0.0069	0.005	0.2355	0.0051
Latent Composition [11]	0.012	0.4735	0.0388	0.012	0.4487	0.0321
StyleGAN3 W [45]	0.016	0.4379	0.0353	0.017	0.3921	0.0307
StyleGAN3 W+ [45]	0.025	0.4634	0.0462	0.023	0.4086	0.0391
StyleMapGAN (32 × 32) [50]	0.007	0.3792	0.0118	0.006	0.1989	0.0045
SDEdit [64]	0.011	0.3857	0.0076	0.008	0.3427	0.0003
Ours	0.013	0.2046	0.0050	0.013	0.2046	0.0050
Ours + Poisson Blending	0.002	0.1883	0.0007	0.002	0.1883	0.0007

3D-aware alignment



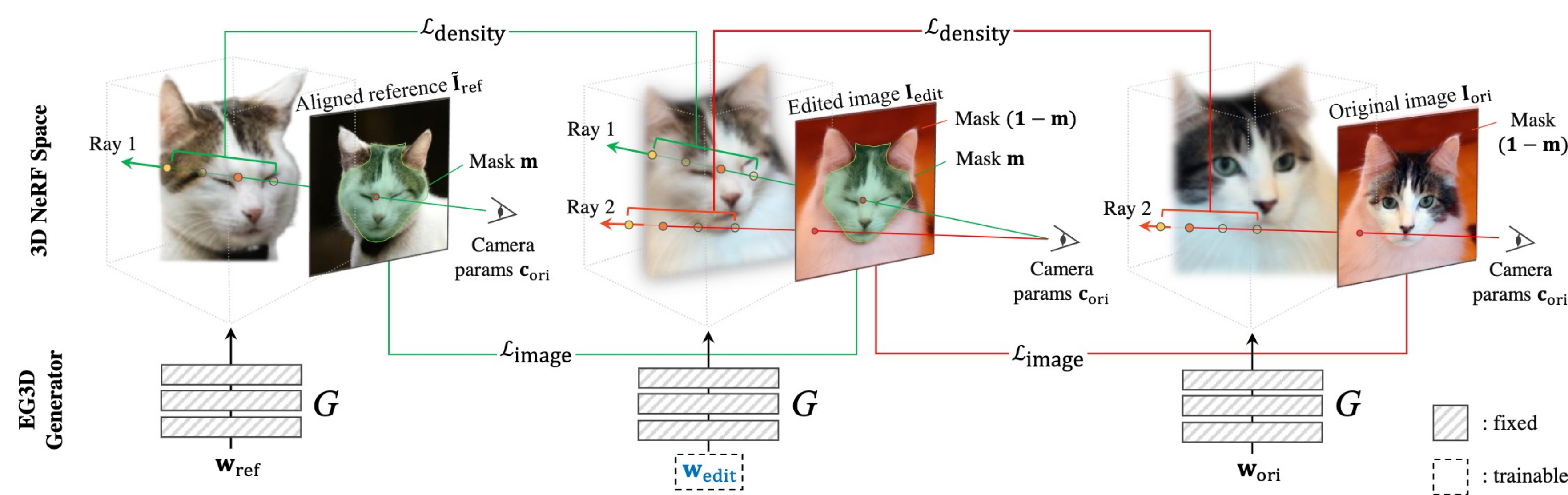
We first use a CNN encoder to infer the camera pose of each input image.

Step 1. Given the camera pose \mathbf{c} , we estimate the latent code \mathbf{w} for each input using a reconstruction loss \mathcal{L}_{rec} .

Step 2. Given the estimated camera pose \mathbf{c}_{ori} and latent code \mathbf{w}_{ref} , we align the reference image to match the pose of the original image.

$$I_{ref}^R = G_{RGB}(\mathbf{w}_{ref}, \mathbf{c}_{ori})$$

3D-aware blending



We aim to find the best latent code \mathbf{w}_{edit} to synthesize a seamless and natural output. To achieve this goal, we exploit both 2D pixel constraints (RGB value) and 3D geometric constraints (volume density). With the proposed image-blending and density-blending losses, we optimize the latent code \mathbf{w}_{edit} .

Image blending loss

$$\mathcal{L}_{image} = \|(1 - \mathbf{m}) \circ I_{edit} - (1 - \mathbf{m}) \circ I_{ori}\|_1 + \lambda_1 \mathcal{L}_{LPIPS}((1 - \mathbf{m}) \circ I_{edit}, (1 - \mathbf{m}) \circ I_{ori}) + \lambda_2 \mathcal{L}_{LPIPS}(\mathbf{m} \circ I_{edit}, \mathbf{m} \circ I_{ref}),$$

Density blending loss

$$\mathcal{L}_{density} = \sum_{\mathbf{r} \in \mathcal{R}_{ref}} \sum_{\mathbf{x} \in \mathcal{R}} \|G_{\sigma}(\mathbf{w}_{edit}; \mathbf{x}) - G_{\sigma}(\mathbf{w}_{ref}; \mathbf{x})\|_1 + \sum_{\mathbf{r} \in \mathcal{R}_{ori}} \sum_{\mathbf{x} \in \mathcal{R}} \|G_{\sigma}(\mathbf{w}_{edit}; \mathbf{x}) - G_{\sigma}(\mathbf{w}_{ori}; \mathbf{x})\|_1.$$

Blending results in various datasets

EG3D backbone : CelebA-HQ, AFHQ-Cat, ShapeNet-Car

StyleSDF backbone: FFHQ

