

Temporally Consistent Amodal Completion for 3D Human-Object Interaction Reconstruction

Hyungjun Doh*

Purdue University

West Lafayette, Indiana, USA

hdoh@purdue.edu

Dong In Lee*†

Purdue University

West Lafayette, Indiana, USA

Korea University

Seoul, Republic of Korea

dilee99@korea.ac.kr

Seunggeun Chi*

Purdue University

West Lafayette, Indiana, USA

chi65@purdue.edu

Pin-Hao Huang, Kwonjoon Lee

Honda Research Institute USA

San Jose, California, USA

pin-hao_huang@honda-ri.com

kwonjoon_lee@honda-ri.com

Sangpil Kim

Korea University

Seoul, Republic of Korea

spk7@korea.ac.kr

Karthik Ramani

Purdue University

West Lafayette, Indiana, USA

ramani@purdue.edu

Human-Object Interaction Monocular Video



Temporally Consistent Amodal Completion



Animatable 3D Object Reconstruction for Novel View Synthesis



Figure 1: (left) In human-object interaction (HOI) scenarios, occlusions frequently affect both the human and the object. (middle) We inpaint the occluded regions while preserving temporal consistency for both entities across frames. (right) Leveraging the temporally consistent image sequences, we reconstruct the human and object using a 3D Gaussian splatting representation, enabling animatable 3D HOI applications.

Abstract

We introduce a novel framework for reconstructing dynamic human-object interactions from monocular video that overcomes challenges associated with occlusions and temporal inconsistencies. Traditional 3D reconstruction methods typically assume static objects or full visibility of dynamic subjects, leading to degraded

*Three authors contributed equally to this research.

†Work done at Purdue University while a visiting scholar.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ACM MM '25, Dublin, Ireland

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM

<https://doi.org/XXXXXX.XXXXXXX>

performance when these assumptions are violated—particularly in scenarios where mutual occlusions occur. To address this, our framework leverages amodal completion to infer the complete structure of partially obscured regions. Unlike conventional approaches that operate on individual frames, our method integrates temporal context, enforcing coherence across video sequences to incrementally refine and stabilize reconstructions. This template-free strategy adapts to varying conditions without relying on predefined models, significantly enhancing the recovery of intricate details in dynamic scenes. We validate our approach using 3D Gaussian Splatting on challenging monocular videos, demonstrating superior precision in handling occlusions and maintaining temporal stability compared to existing techniques.

CCS Concepts

- Computing methodologies → Shape inference; Reconstruction.

Keywords

Human-Object Interaction, Amodal Completion, Temporal Consistency

ACM Reference Format:

Hyungjun Doh, Dong In Lee, Seunggeun Chi, Pin-Hao Huang, Kwonjoon Lee, Sangpil Kim, and Karthik Ramani. 2025. Temporally Consistent Amodal Completion for 3D Human-Object Interaction Reconstruction. In *Proceedings of Proceedings of the 33rd ACM International Conference on Multimedia (ACM MM '25)*. ACM, New York, NY, USA, 14 pages. <https://doi.org/XXXXXXX>

1 Introduction

Understanding the complex interplay between humans and objects is a fundamental challenge in computer vision and robotics, with broad implications for both theory and real-world applications. This understanding is critical for enabling technologies such as autonomous robots that navigate complex human environments and augmented or virtual reality systems that rely on precise 3D reconstructions. Despite significant progress, reliably interpreting human-object interactions (HOI) in dynamic, real-world settings remains an open problem. Visual complexities—especially mutual occlusions—often obscure critical details, underscoring the need for enhanced 3D reconstruction techniques that capture and animate these interactions effectively.

Multi-view 3D reconstruction [16, 27, 33] methods have proven effective in generating detailed models from multiple images. Object reconstruction techniques [13, 35, 42, 49] are generally designed for static scenes, where the object remains stationary. In contrast, human reconstruction methods [10, 15, 19] typically handle dynamic subjects assuming full visibility. However, in HOI scenarios, these assumptions break down. Both humans and objects are in motion, leading to frequent and mutual occlusions as illustrated in Figure 1, resulting in ambiguous visual cues that complicate the reconstruction process. The presence of dynamic and mutual occlusions introduces complexities beyond the scope of static object or dynamic full-visibility human reconstruction methods.

Amodal completion [1, 2, 5, 22, 29, 40, 48] has emerged as a promising strategy to mitigate these occlusion challenges by inferring the complete structure of partially hidden regions. However, most existing approaches are designed for static scenes. Although they perform well in controlled environments, they often falter in dynamic, real-world settings. A key limitation is their reliance on image-level completion, which overlooks the rich temporal context provided by adjacent sequences in HOI cases. Without leveraging temporal information, reconstructions derived from amodal-completed images can suffer from inconsistencies, ultimately compromising the overall quality and stability of the models.

Recent research [6, 14, 37, 38, 43, 45, 47] suggests that incorporating temporal consistency can effectively address these challenges. By enforcing coherence across consecutive video frames, models can integrate information from multiple viewpoints more effectively, leading to more accurate recovery of occluded regions. This strategy also enables the incremental refinement of hidden details over time, yielding reconstructions that more faithfully capture the dynamic nature of human-object interactions.

Building on these insights, we propose a novel, temporally consistent amodal completion framework tailored for monocular video

input. By integrating temporal context into the amodal completion process, our approach overcomes the limitations of traditional frame-by-frame methods, significantly enhancing both stability and realism. Specifically, our framework consists of three key components. First, we introduce *Bidirectional Temporal Feature (BTF) Warping*, which leverages optical flow to warp latent features from both past and future frames into the current frame's latent space, enabling effective propagation of temporal cues and aligns spatially meaningful information across time. Second, we present a *Temporal Fusion Attention* mechanism that adaptively aggregates these temporally aligned features, producing a coherent and semantically enriched latent representation. Third, we propose a *Template-free Occlusion Identification* strategy that combines 2D segmentation and 3D projections to localize occlusions without relying on predefined templates. Finally, we apply amodal completion, guided by the temporal-aware latent features and precise occlusion masks, to complete missing regions with high fidelity. This unified pipeline enables robust recovery of occluded structures across time, ensuring temporal consistency and accurate amodal completion. Unlike prior approaches, our method generalizes effectively to dynamic, real-world human-object interaction scenarios.

We validate the effectiveness of our method through extensive experiments on two public datasets, BEHAVE [4] and InterCap [11, 12], which present severe occlusions and diverse human-object interactions. Our framework demonstrates consistently superior performance in both quantitative metrics and qualitative comparisons. Moreover, we extend our evaluation to reconstructing 3D human-object interactions from monocular video sequences using 3D Gaussian Splatting [16] (3DGS), enabling various 3D application of animatable HOI including novel-view synthesis.

In summary, our contributions include:

- **Temporally Consistent Completion Method:** We carefully design a Bidirectional Temporal Feature (BTF) Warping and a Temporal Fusion Attention mechanism that incorporate temporal context to amodal completion by adaptive fusion of aligned features across past and future frames.
- **Template-free Occlusion Identification Strategy:** We propose a hybrid approach that integrates 2D segmentation with 3D point cloud projection to accurately localize occluded regions. This method introduces a novel masking technique that identifies occlusions efficiently, without the need for predefined templates.
- **Application to 3D Reconstruction:** To the best of our knowledge, our work is the first approach to reconstruct photo-realistic and animatable 3D human-object interactions from monocular videos. By leveraging 3D Gaussian Splatting [16], we demonstrate that our pipeline is applicable to various subtasks, including novel-view and novel-pose synthesis for HOI.

2 Related Work

2.1 3D Reconstruction

3D reconstruction [16, 27, 33] from multi-view images has demonstrated versatility across various fields. In particular, 3D reconstruction using explicit representations [16] has emerged as an effective method for capturing 3D information. This representation has been

successfully employed in 3D scene reconstruction [8, 23, 25], object reconstruction [13, 35, 42, 49], and human reconstruction [10, 15, 19]. However, these methods often face challenges when occlusions occur. To mitigate the occlusion problem in human reconstruction, few research proposed [21, 34] a diffusion-based method. However, these methods implicitly secure temporal consistency through a 2D diffusion prior, not guaranteeing temporal consistency. In contrast to existing approaches that primarily address either human or static objects, our method specifically targets cases involving human-object interactions via template-free amodal completion.

2.2 Amodal Completion

Amodal completion focuses on inferring the hidden shapes and appearances of occluded objects. Recently, various works leverage generative models to reconstruct missing regions. For instance, pix2gestalt [2] synthesizes “wholes” from partial observations, while Progressive Mixed Context Diffusion [1] expands occluded areas with context-awareness. Other studies propose transformer-based amodal segmentation [48], leverage 3D shape priors to guide amodal masks [5], or learn compositional scene representations to handle complex occluders [22, 29]. Self-supervised frameworks also emerged for amodal completion by aligning visible parts with plausible hidden geometry [17, 40]. However, these single-image methods often struggle with dynamic human-object interactions, where severe occlusions may shift from frame to frame—leading to temporally inconsistent reconstructions. As an alternatives, video inpainting methods [24] has also been introduced, however these methods are usually for removing subjects well fitted to background and environment, showing limitation on HOI cases. Our approach addresses these issues by integrating temporal cues, ensuring stable amodal completions across sequences.

2.3 Temporal Consistency

Ensuring temporal consistency across image sequences is essential for stable and coherent video processing. Traditional approaches often rely on optical flow to propagate features between frames [14, 37, 38]. Recent advancements in video diffusion have incorporated explicit temporal modeling to improve alignment and consistency across frames. Several state-of-the-art techniques [6, 43, 45, 47] adopt flow-guided recurrent architectures, enabling temporally aware representations for applications such as video super-resolution, inpainting, video-to-video translation, and text-driven video editing. Complementary strategies include transformer-based frame propagation [46] and cross-frame feature matching through diffusion-based tokens [9], both of which contribute to enhanced coherence in complex temporal dynamics. While these methods primarily focus on video generation and restoration tasks, our approach uniquely leverages optical flow and feature warping mechanisms to specifically address occlusions and ensure temporally consistent reconstructions in 3D HOI scenarios.

3 Method

Our goal is to achieve temporal consistency while completely recovering occluded regions for accurate 3D reconstruction of Human–Object Interactions (HOI) from monocular video. This is challenging due to frequent and severe occlusions in HOI scenarios. To

address this challenge, we propose a novel pipeline that integrates temporal context into the amodal completion process.

We begin by formulating the temporally consistent amodal completion problem, clearly defining the inputs, outputs, and objectives in Section 3.1. Based on this formulation, our method consists of the following four main components, as illustrated in Figure 2. First, in Section 3.2, we introduce Bidirectional Temporal Feature (BTF) warping, which aligns features across neighboring frames using optical flow and latent feature warping. Second, in Section 3.3, we present Temporal Fusion Attention, a mechanism that dynamically fuses the temporally aligned features into a coherent representation. Third, in Section 3.4, we describe a template-free occlusion identification strategy that accurately localizes occluded regions. The last component, explained in Section 3.5, is temporally-aware amodal completion mechanism. Consequently, our method enables reliable 3D reconstruction through temporally consistent amodal completion, as detailed in the supplementary materials.

3.1 Problem Formulation

We address the problem of filling in occluded regions in a sequence of video frames such that the completed regions are both semantically plausible and temporally coherent. For each frame, let $I_{in} \in \mathbb{R}^{H \times W \times 3}$ be the segmented image containing only the visible regions, and let $M_{in} \in \{0, 1\}^{H \times W}$ be a binary mask, where 1 indicates the occluded regions to be inpainted. In addition, we provide a temporally-aware latent feature z_{in} , aggregated from temporally aligned features across neighboring frames, and a text prompt P as textual guidance. We formalize the inpainting process with the diffusion model [32] guided by both latent and mask inputs as:

$$I_{out} = F_{s \rightarrow e}(I_{in}, M_{in}, z_{in}, P), \quad (1)$$

where $F_{s \rightarrow e}$ denotes the diffusion process from denoising time step s to e . The effectiveness of our framework critically depends on two factors: (1) the quality of the latent feature z_{in} , which provides temporally consistent and semantically rich context, and (2) the accuracy of the occlusion mask M_{in} , which precisely localizes the region to be completed. In the following sections, we describe how each component is constructed to ensure effective amodal completion under complex human-object interaction scenarios.

3.2 Bidirectional Temporal Feature Warping

To ensure temporal consistency across consecutive frames, we introduce a Bidirectional Temporal Feature (BTF) warping method that aligns latent features from both past and future frames to the current frame. As part of this approach, we define a temporal support set $\mathcal{S}(t)$ at each time step t , consisting of past and future neighboring frames:

$$\mathcal{S}(t) = \{t-n, \dots, t-1, t+1, \dots, t+n\}. \quad (2)$$

We empirically choose $n = 7$ aggregating total 14 neighboring features. For each element in the support set $\tau \in \mathcal{S}(t)$, we estimate optical flow from frame I^τ to the reference frame I^t to spatially align their latent representations. The optical flow is computed as:

$$o_\tau^t = \text{Flow}(I^\tau, I^t), \quad \forall \tau \in \mathcal{S}(t), \quad (3)$$

where $\text{Flow}(\cdot, \cdot)$ denotes an off-the-shelf optical flow estimation network [38].

Amodal Completion with Temporal Consistency

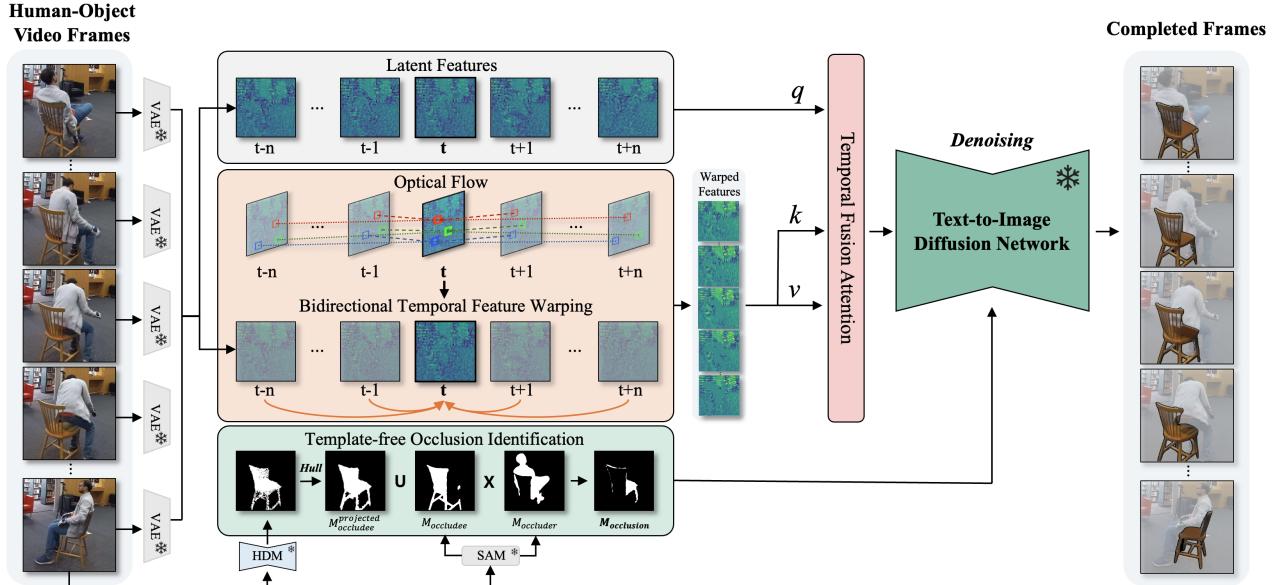


Figure 2: Overview of Our Framework: Given a human-object interaction (HOI) monocular video, our framework performs amodal completion through (1) Bidirectional Temporal Feature Warping via optical flow, (2) Temporal Fusion Attention for aggregating multi-frame context, (3) Template-free Occlusion Identification using 2D and 3D cues, and (4) temporally-aware amodal completion. This design enables temporally consistent and accurate amodal completion in complex HOI scenarios.

We then encode latent features using a pretrained Variational Autoencoder (VAE) [18] \mathcal{E} as utilized in diffusion models [32], defined as $z^\tau = \mathcal{E}(I^\tau)$. To spatially align the latent representations across frames, the encoded feature z^τ is subsequently warped into the coordinate space of frame t using the estimated optical flow o_t^τ which is bilinearly scaled to match the resolution of z^τ :

$$\hat{z}_\tau^t = \text{Warp}(z^\tau, o_t^\tau), \quad \forall \tau \in \mathcal{S}(t), \quad (4)$$

The warped features \hat{z}_τ^t serve as temporally aligned observations centered at frame I^t , enabling the aggregation of multi-frame context from both past and future frames. By referencing complementary information from adjacent frames, our framework can recover regions that are occluded or missing in the current view in complex human-object interaction scenarios.

3.3 Temporal Fusion Attention

While the warped features are temporally aligned, it is crucial to effectively incorporate temporal context from each warped feature into a unified latent representation. Therefore, we propose a Temporal Fusion Attention mechanism that selectively integrates complementary cues from neighboring frames. Specifically, we employ a cross-attention strategy to aggregate temporally aligned latent features within the support set $\mathcal{S}(t)$. At each time step t , the queries, keys, and values are defined as:

$$Q^t = z^t, \quad K^t = V^t = \text{concat}(\hat{z}_{t-n}^t, \dots, \hat{z}_{t+n}^t), \quad (5)$$

where $Q^t \in \mathbb{R}^{1 \times d}$ represents the latent feature from the current frame, and $K^t \in \mathbb{R}^{2n \times d}$, $V^t \in \mathbb{R}^{2n \times d}$ represent the warped latent

features aligned from adjacent frames. The attention-weighted representation is computed via scaled dot-product attention:

$$z_{fused}^t = \text{Attn}(Q^t, K^t, V^t) = \text{Softmax}\left(\frac{Q^t \cdot K^t}{\sqrt{d}}\right) V^t, \quad (6)$$

where d denotes the feature dimensionality. By leveraging temporally warped features from both past and future frames, this mechanism selectively integrates complementary information that is not visible in the current view, thereby improving occlusion completion. This yields z_{fused}^t , a semantically rich context that captures occluded structures and interaction context from neighboring frames. By emphasizing temporally relevant and contextually informative features, this attention mechanism significantly enhances coherence and robustness across frames, thereby facilitating high-quality and stable amodal completion.

3.4 Template-free Occlusion Identification

To accurately localize the occluded area, it is essential to estimate the complete shape of the target—referred to as the occludee—which is the object of interest partially hidden in the scene. Therefore, We employ the Hierarchical Diffusion Model (HDM) [41], a template-free approach specifically designed for reconstructing human-object interactions. HDM facilitates the inference of a comprehensive 3D point cloud representing complete geometry of human and object. We then project this point cloud onto the 2D image plane to obtain a set of projected points: $C = \{p_1, p_2, \dots, p_n\}$, where each point $p_i = (x_i, y_i) \in \mathbb{R}^2$ corresponds to a position in the 2D image plane. To extract a coherent shape from these points, we apply a concave hull algorithm [3] that constructs a tight boundary around the set. Unlike a convex hull, which encloses the outermost points with the

smallest convex polygon, the concave hull denoted as $\text{Hull}(\cdot)$ allows for inward indentations, enabling the capture of non-convexities and finer geometric details. This leads to a more accurate approximation of the object’s silhouette. The projected full shape mask of the target is then defined as:

$$M_{\text{occludee}}^{\text{projected}} = \text{Hull}(C). \quad (7)$$

To ensure a more robust estimation of the target’s shape, we construct a fused mask, M_{union} , by combining the segmented visible mask $M_{\text{occludee}}^{\text{visible}}$, obtained by SAM2 [31], with the projected full shape mask $M_{\text{occludee}}^{\text{projected}}$:

$$M_{\text{union}} = M_{\text{occludee}}^{\text{visible}} \cup M_{\text{occludee}}^{\text{projected}}, \quad (8)$$

which allows the two masks to effectively complement one another, integrating both observed and inferred regions of the target. The Figure 2 demonstrates how the visible and projected masks are combined to yield a more complete representation.

Finally, we introduce an occluder mask M_{occluder} , also segmented using SAM2, representing the region that occludes the target. The occluder can be either a human or an object, depending on the target of interest. The actual occlusion mask $M_{\text{occlusion}}$, which delineates the occluded area of the target, is determined by intersecting the fused mask M_{union} with the occluder mask M_{occluder} :

$$M_{\text{occlusion}} = M_{\text{union}} \cap M_{\text{occluder}}. \quad (9)$$

With this approach, we precisely identify the occluded region of the target object without relying on any predefined templates, significantly enhancing the generalizability and robustness of our approach.

3.5 Temporally-aware Amodal Completion

With temporally fused features z_{fused}^t and the occlusion regions $M_{\text{occlusion}}^t$ identified in previous stages, we inpaint the occluded region of the target using the diffusion inpainting pipeline. To achieve temporally consistent amodal completion, we first apply occludee mask M_{occludee}^t to the input image I thereby removing background distractions and encouraging the model to focus on the target:

$$I_{\text{occludee}}^t = M_{\text{occludee}}^t \odot I^t. \quad (10)$$

Next, we use the occlusion mask $M_{\text{occlusion}}^t$ to specify the region that should be completed. To incorporate temporal context, we leverage a temporally-aware latent feature z_{fused}^t , applying mask \hat{M}_{union}^t to ensure that latent representation only affects the target. Note that \hat{M}_{union}^t is bilinearly downsampled mask to match the dimension of z_{fused}^t . A corresponding text prompt P is also provided to guide the inpainting process. This process can be written as:

$$I_{\text{out}}^t = \mathbf{F}_{s \rightarrow e}(I_{\text{occludee}}^t, M_{\text{occlusion}}^t, z_{\text{fused}}^t \odot \hat{M}_{\text{union}}^t, P). \quad (11)$$

Through the integration of temporal cues, the diffusion inpainting pipeline reliably fills in the occluded regions while maintaining temporal consistency across frames. Consequently, our framework effectively resolves ambiguities due to occlusions and complex dynamics inherent in realistic human-object interactions, improving the quality and stability of amodal completion outcomes.

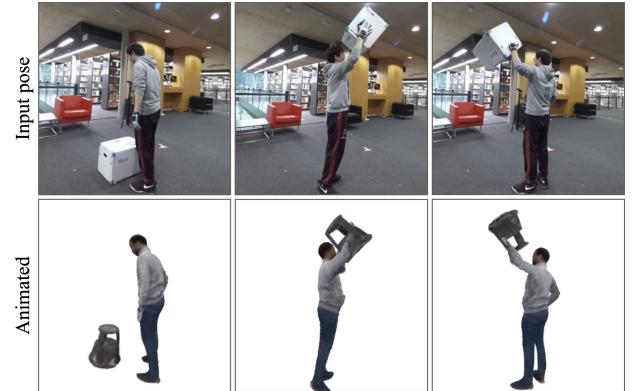


Figure 3: Animatable 3D Reconstruction of Human–Object Interactions: we train 3DGS [16] using temporally consistent amodal completions from our pipeline, sampled at 1 fps (15–47 frames per video). Conditioned on motion trajectories from the input video, our method enables realistic animation of novel human–object pairs while preserving geometry, appearance, and temporal coherence.

3.6 Application in 3D Reconstruction

To demonstrate the utility of our temporally consistent amodal completion in downstream tasks, we reconstruct photo-realistic and animatable 3D human–object interaction scenes from monocular video using 3D Gaussian Splatting (3DGS) [16]. Our pipeline produces temporally inpainted frames I_{out}^t , which serve as dense appearance supervision even under severe occlusions.

For object reconstruction, we follow a GS-Pose [26], using provided 6-DoF object and camera poses. To ensure consistent scale and visibility, each object in image is cropped and re-centered in every frame \tilde{I}_{out}^t . A 3D Gaussian representation G_{obj} is then optimized by minimizing a photometric loss between rendered images and inpainted frames:

$$\mathcal{L}_{\text{obj}} = \sum_t \mathcal{L}_{\text{photo}}(R(G_{\text{obj}}, H_t \circ X_t), \tilde{I}_{\text{out}}^t), \quad (12)$$

where $R(\cdot)$ denotes the Gaussian renderer, and H_t, X_t are the camera and object poses. The $\mathcal{L}_{\text{photo}}$ is a combination of L1 and SSIM terms:

$$\mathcal{L}_{\text{photo}}(I_r, I_{\text{gt}}) = \|I_r - I_{\text{gt}}\|_1 + \lambda \cdot (1 - \text{SSIM}(I_r, I_{\text{gt}})), \quad (13)$$

where I_r is the rendered image, I_{gt} is the ground truth image, and λ controls the balance between pixel-wise and perceptual similarity. In our experiments, we set $\lambda = 0.2$ to balance these terms effectively.

For human reconstruction, we adopt GaussianAvatar [10], which learns a canonical 3D Gaussian representation and a pose-conditioned deformation field based on SMPL parameters. The model is trained with the same inpainted sequences, allowing robust geometry and appearance recovery even under partial occlusions.

Together, as visualized in Figure 3, these reconstructions validate that our temporally consistent amodal completion provides a reliable and expressive foundation for photo-realistic, animatable 3D human-object modeling from monocular video.

Method	BEHAVE [4]				InterCap [11, 12]			
	Amodal Completion		Temporal Consistency		Amodal Completion		Temporal Consistency	
	IoU ↑	CLIP ↑	Warp-err ($\times 10^{-3}$) ↓	TC Score ↑	IoU ↑	CLIP ↑	Warp-err ($\times 10^{-3}$) ↓	TC Score ↑
Pix2gestalt [2]	61.29%	26.77	141.08	95.97	64.91%	25.32	112.75	95.83
SD Inpainting [32]	60.81%	27.63	9.87	96.78	45.16%	27.23	7.04	96.71
LaMa [36]	43.54%	26.08	6.34	96.96	56.42%	26.02	3.78	98.01
VDT [24]	55.69%	26.48	5.84	96.58	59.96%	26.11	3.42	96.39
Ours (HDM Mask)	61.75%	27.64	<u>6.26</u>	97.19	70.18%	27.65	3.22	<u>96.95</u>
Ours (Ground Truth Mask)	70.90%	28.01	6.09	98.15	74.79%	27.66	2.93	97.98

Table 1: Quantitative Comparison on BEHAVE [4] and InterCap [11, 12] for Amodal Completion and Temporal Consistency: Our method achieves consistently strong performance, highlighting its robustness to occlusion and temporal challenges. Bold and underline denote the best and second-best scores.

4 Experiment

4.1 Datasets

BEHAVE The BEHAVE dataset [4] consists of 321 RGB-D sequences capturing indoor human–object interactions, recorded using 4 Kinect cameras. The test set includes 3 subjects interacts with 20 objects. Among these, we select 18 videos for 18 objects excluding keyboard and basketball due to the lacking of ground truth pose annotation in 30 fps video sequence. For human, we selected 3 videos for 3 subjects in the test set. As a results, we applied our pipeline to approximately 27,000 frames and evaluate these results.

InterCap InterCap [11, 12] contains 223 RGB-D videos of human-object interactions, captured from 6 distinct viewpoints with 10 subjects and 10 objects. From this dataset, we extract 10 videos that collectively represent all 10 objects.

4.2 Evaluation Metrics

Amodal Completion We evaluate the performance of our amodal completion using two metrics: the CLIP score [30] and Intersection over Union (IoU). The CLIP score evaluates the alignment between the generated images and the corresponding category prompts, while IoU measures the overlap between the predicted and groundtruth amodal masks. Specifically, we calculate the CLIP score for each inpainted frame within tight bounding boxes to minimize the influence of background pixels. For the IoU evaluation, we segment the inpainted frame using SAM2 [31] to obtain masks and then calculate the overlap with the object’s groundtruth masks.

Temporal Consistency We assess temporal consistency using two metrics: the Temporal Consistency score (TC score) [7] and the Flow Warping Error [20]. The TC score is computed by extracting CLIP image embeddings for all output video frames and calculating the average cosine similarity between pairs of consecutive frames. Flow Warping Error, commonly used in optical flow estimation and video processing, measures the discrepancy between a warped image and its corresponding target frame. In our approach, we compute optical flows between consecutive ground-truth object masks using SEA-RAFT [38]. We then warp the previously inpainted frame to align with the current frame, and measure the discrepancy between this warped previous frame and the actual current frame.

A higher Temporal Consistency Score and a lower Flow Warping Error both indicate better temporal consistency.

3D Reconstruction Since our method reconstructs only the object or the human—excluding the background—we evaluate reconstruction quality using masked versions of standard metrics: Peak Signal-to-Noise Ratio (PSNR-M), Structural Similarity Index Measure (SSIM-M) [39], and Learned Perceptual Image Patch Similarity (LPIPS-M) [44]. Inspired by [28], PSNR-M, SSIM-M and LPIPS-M are calculated within tight bounding boxes around the reconstructed regions to minimize the influence of background pixels.

4.3 Temporally Consistent Amodal Completion

Quantitative Results We quantitatively evaluate our approach against several state-of-the-art baselines, including Pix2gestalt [2], Stable Diffusion Inpainting [32], LaMa [36], and VDT [24], initialized with random noise. The results are summarized in Table 1, covering both human and object categories, as the human is treated as one of 19 object classes in our experiments, and only frames with an object occlusion ratio between 15% and 70% are considered. Our method consistently achieves the best performance across amodal completion metrics, IoU and CLIP, on both BEHAVE and InterCap. This indicates a higher fidelity in reconstructing occluded regions, validating the effectiveness of our template-free regioning strategy in complex HOI cases. To assess ideal condition performance, we report results using ground-truth masks in the last row.

Regarding temporal consistency, our method achieves the lowest warping error among all methods, reflecting superior temporal alignment between consecutive frames. This highlights our model’s strength in preserving temporal coherence throughout the video sequence. Interestingly, although LaMa reports the highest TC score, it ranks lowest in IoU and performs poorly in terms of visual quality. This may be attributed to its tendency to overwrite masked regions with background-matching content, which is especially beneficial in static-camera datasets like BEHAVE and InterCap. In contrast, our method reconstructs occlusions by explicitly leveraging spatial-temporal information from adjacent frames, leading to more semantically faithful and temporally robust completions.

Qualitative Results Figure 4 presents a qualitative comparison across various methods on three representative HOI categories: *Square Table*, *Small Table*, and *Skateboard*. Compared to baseline

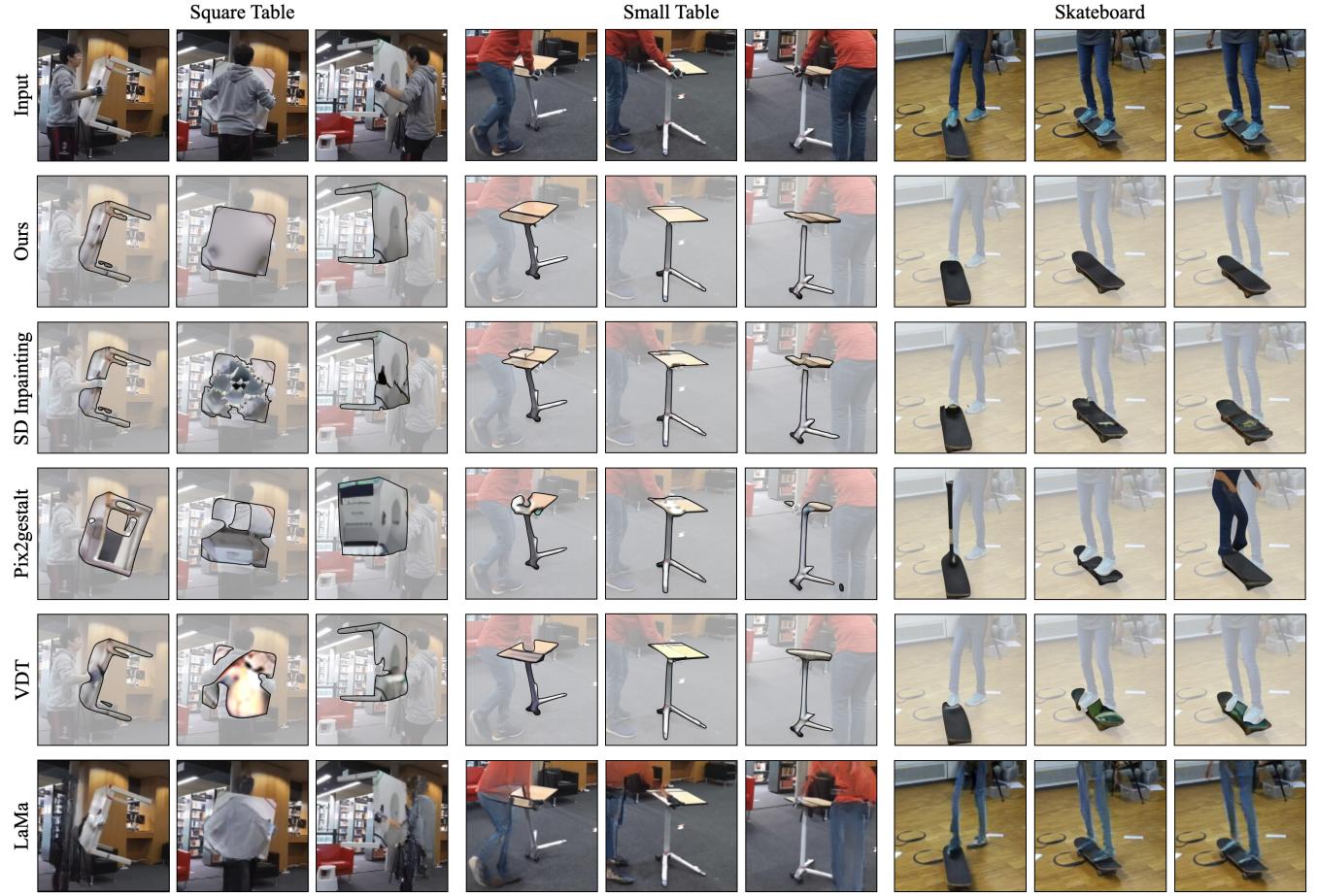


Figure 4: Qualitative comparison on BEHAVE [4] (Square Table, Small Table) and InterCap [11, 12] (Skateboard). Our method produces accurate and temporally consistent completions of occluded regions.

methods, our approach produces semantically faithful and visually coherent reconstructions, especially in the occluded regions. For *Square Table*, Pix2gestalt generates the implausible shape and VDT fails to complete the proper region. In contrast, our method preserves object integrity with precise geometry and sharper texture completion. In the *Skateboard* scenario, our method excels in preserving object shape and continuity under occlusion. Competing methods often hallucinate incorrect geometry or texture, especially under relatively simple occlusions caused by human limbs. Overall, the visual results confirm that our method outperforms prior approaches in generating amodally completed outputs that are both spatially and temporally coherent, reinforcing the strength of our template-free, temporally consistent completion framework.

4.4 3D Reconstruction

To the best of our knowledge, this is the first method to enable animatable and photo-realistic 3D reconstruction in human-object interaction (HOI) scenarios through temporally consistent amodal completion. We use the completed monocular video frames generated by our pipeline to train 3D Gaussian Splatting [16] (3DGS), enabling high-quality reconstruction from monocular inputs. For training, we extract 40 to 47 frames per sequence at 1 fps, a sparse set of inpainted frames.

We compare three settings: (1) our full pipeline with temporally-aware amodal completion, (2) Stable Diffusion inpainting without our temporal modules, and (3) the original input sequence without any completion. As shown in Table 2, our method achieves the highest PSNR-M, SSIM-M and LPIPS-M, outperforming all comparison methods. As illustrated in Figure 5, our approach improves both geometric accuracy and temporal consistency, which are crucial for smooth and realistic novel-view rendering. Furthermore, when conditioned on SMPL body poses or 6D object trajectories, the trained model produces photo-realistic animations of both humans and objects, highlighting the robustness and completeness of our pipeline. In contrast, without our temporally consistent amodal completion, occluded regions remain unresolved and frame-to-frame inconsistencies persist, leading to degraded 3D reconstruction quality.

4.5 Ablation Study

To assess the contributions of each component in our temporal consistency framework, we conduct ablation studies on the BEHAVE dataset, as shown in Table 3. We analyze the impact of two key components: feature warping and cross-frame attention. Without either feature warping or cross attention (first row), the model performs reasonably well, but shows higher warping error and lower temporal consistency. Introducing cross attention alone (second

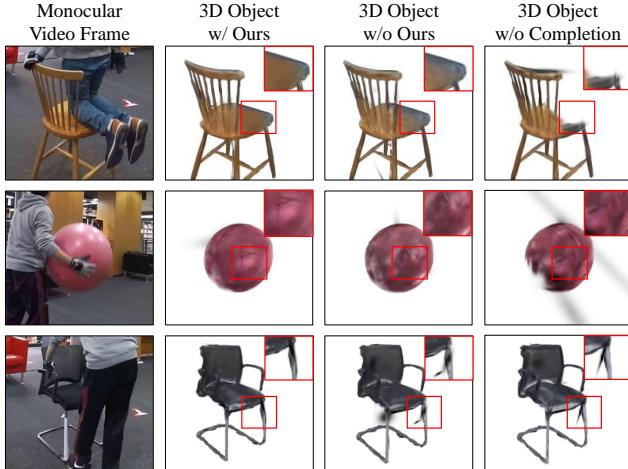


Figure 5: Qualitative results of 3D Reconstruction: Training on our temporally consistent completions (left) produces clearer geometry and fewer artifacts than single-frame (middle) or no completion (right).

Method	PSNR-M ↑	SSIM-M ↑	LPIPS-M ↓
w/o Completion	13.48	0.5679	0.3604
SD Inpainting [32]	15.38	0.6186	0.2455
Ours	15.43	0.6240	0.2383

Table 2: Quantitative results on 3D reconstruction from the BEHAVE dataset [4]: We report PSNR-M, SSIM-M, and LPIPS-M, computed within tight bounding boxes around humans or objects. Our method outperforms all baselines across metrics.

Feature	Cross Attn.	IoU ↑	CLIP ↑	Warp-err ($\times 10^{-3}$) ↓	TC ↑
✗	✗	60.81	27.63	9.87	96.78
✗	✓	60.82	27.64	6.56	97.06
✓	✗	58.35	27.60	6.23	97.16
✓	✓	61.75	27.64	<u>6.26</u>	97.19

Table 3: Ablation Studies on temporal consistency strategy for on BEHAVE dataset [4]. Bold and underline denote the best and second-best scores.

row) slightly improves all metrics, particularly reducing the warping error. Interestingly, feature warping alone (third row) leads to the lowest warping error (6.23) and improved TC score, demonstrating its effectiveness in aligning features across time. However, it degrades amodal completion accuracy (IoU). Combining both feature warping and cross attention (last row) yields the best overall performance, achieving the highest IoU (61.75), CLIP score (27.64), and TC score (97.19). This confirms that both components are complementary and essential for generating temporally consistent and semantically meaningful completions.

We conduct an additional ablation study to evaluate the effectiveness of our template-free occlusion identification strategy, as shown in Table 3. Specifically, we compare our method against a baseline that relies on predefined human masks to guide the inpainting process. Our approach derives the occlusion mask $M_{\text{occlusion}}$ by combining 2D segmentation with 3D point cloud projection,

Mask	IoU ↑	CLIP ↑	Warp-err ($\times 10^{-3}$) ↓	TC ↑
Human Mask	53.77	28.31	8.90	97.60
Ours ($M_{\text{occlusion}}$)	61.75	27.64	<u>6.26</u>	97.19

Table 4: Ablation Studies on Template-free Occlusion Identification strategy on BEHAVE dataset [4].

leading to significantly better performance in terms of IoU (61.75 vs. 53.77) and warping error (6.26 vs. 8.90). These results highlight our method's superior ability to accurately localize occluded regions and maintain temporal coherence. Although the human-mask baseline yields a slightly higher CLIP score and TC Score, this is likely due to its tendency to inpaint overly large regions, which degrades spatial accuracy and temporal consistency, as reflected in the lower IoU and higher warping error. Qualitative results of the ablation study are provided in the supplementary material.

5 Discussion and Limitations

Dynamic human–object interactions present significant challenges for optical flow-based shape estimation, particularly due to non-linear motion and frequent occlusions. Our approach leverages HDM for amodal completion, enabling the reconstruction of full object shapes—including occluded regions. However, the method may not always accurately infer the geometry of unseen parts. Future directions include integrating multi-view geometry techniques, such as Structure-from-Motion (SfM) and Multi-View Stereo (MVS), to densify and refine reconstructions, thereby enhancing the robustness of our pipeline.

Currently, our method assumes that each input video contains only a single human and a single object. While effective for controlled scenarios, this assumption limits applicability in more complex settings, such as multi-person interactions or interactions involving multiple objects. Extending the framework to handle such cases is a promising direction for future work.

Additionally, our method relies heavily on the inpainting capabilities of the Stable Diffusion model, which may struggle to accurately reconstruct objects or appearances that were not well represented during its training. We also observed that the quality of inpainting is highly sensitive to the accuracy of the occlusion masks. In particular, when the masks fail to fully exclude occluding regions—leaving behind small residual artifacts—the inpainting output is noticeably degraded. These limitations highlight the critical importance of precise occlusion localization, especially in scenes involving fine-grained interactions or partial visibility.

6 Conclusion

We proposed a novel framework for reconstructing dynamic human–object interactions from monocular video, with a focus on addressing occlusions and temporal inconsistency. Our method achieves state-of-the-art performance across multiple benchmarks, demonstrating both quantitative superiority and qualitative robustness. It also generalizes well to real-world scenarios. Notably, by integrating our approach with 3D Gaussian Splatting, we enable the generation of photorealistic and animatable 3D reconstructions from monocular input, supporting advanced applications such as novel-view and novel-pose synthesis in complex HOI scenes.

Acknowledgments

We wish to thank all the reviewers for their invaluable feedback. This work is partially supported by the NSF under the Future of Work at the Human-Technology Frontier (FW-HTF) 1839971 and Partnership for Innovation: Technology Transfer (PFI-TT) 2329804. Additional support for this work is provided by the Culture, Sports, and Tourism R&D Program through the Korea Creative Content Agency grant funded by the Ministry of Culture, Sports and Tourism in 2024 (International Collaborative Research and Global Talent Development for the Development of Copyright Management and Protection Technologies for Generative AI, RS-2024-00345025), Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (No. RS-2019-II190079, Artificial Intelligence Graduate School Program (Korea University). We also acknowledge the Feddersen Distinguished Professorship Funds and a gift from Thomas J. Malott. Any opinions, findings, and conclusions expressed in this material are those of the authors and do not necessarily reflect the views of the funding agency.

References

- [1] Firstname Author and Secondname Author. 2023. Amodal Completion via Progressive Mixed Context Diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 1–10.
- [2] Firstname Author, Secondname Author, and Thirdname Others. 2023. pix2gestalt: Amodal Segmentation by Synthesizing Wholes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 1–10.
- [3] K. Bellock. [n. d.]. alphashape. <https://github.com/bellockk/alphashape>. GitHub repository, accessed 2025-04-11.
- [4] Bharat Lal Bhatnagar, Xianghui Xie, Ilya A Petrov, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. 2022. Behave: Dataset and method for tracking human object interactions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 15935–15946.
- [5] A. Chen, B. Smith, and C. Lee. 2023. Amodal 3D Shape from Partial Views. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 4567–4576.
- [6] Yuren Cong, Mengmeng Xu, Christian Simon, Shoufa Chen, Jiawei Ren, Yanping Xie, Juan-Manuel Perez-Rua, Bodo Rosenhahn, Tao Xiang, and Sen He. 2023. Flatten: optical flow-guided attention for consistent text-to-video editing. *arXiv preprint arXiv:2310.05922* (2023).
- [7] Patrick Esser, Johnathan Chiu, Parmida Atighehchian, Jonathan Granskog, and Anastasis Germanidis. 2023. Structure and content-guided video synthesis with diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*. 7346–7356.
- [8] Yang Fu, Sifei Liu, Amey Kulkarni, Jan Kautz, Alexei A. Efros, and Xiaolong Wang. 2024. COLMAP-Free 3D Gaussian Splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 20796–20805.
- [9] Michal Geyer, Omer Bar-Tal, Shai Bagon, and Tali Dekel. 2023. Tokenflow: Consistent diffusion features for consistent video editing. *arXiv preprint arXiv:2307.10373* (2023).
- [10] Liangxiao Hu, Hongwen Zhang, Yuxiang Zhang, Boyao Zhou, Boning Liu, Shengping Zhang, and Liqiang Nie. 2024. Gaussianavatar: Towards realistic human avatar modeling from a single video via animatable 3d gaussians. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 634–644.
- [11] Yinghao Huang, Omid Taheri, Michael J Black, and Dimitrios Tzionas. 2022. InterCap: Joint markerless 3D tracking of humans and objects in interaction. In *DAGM German Conference on Pattern Recognition*. Springer, 281–299.
- [12] Yinghao Huang, Omid Taheri, Michael J Black, and Dimitrios Tzionas. 2024. InterCap: Joint Markerless 3D Tracking of Humans and Objects in Interaction from Multi-view RGB-D Images. *International Journal of Computer Vision* (2024), 1–16.
- [13] Ajay Jain, Matthew Tancik, and Pieter Abbeel. 2021. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5885–5894.
- [14] Jisoo Jeong, Jamie Menjay Lin, Fatih Porikli, and Nojun Kwak. 2022. Imposing consistency for optical flow estimation. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*. 3181–3191.
- [15] Wei Jiang, Kwang Moo Yi, Golnoosh Samei, Oncel Tuzel, and Anurag Ranjan. 2022. Neuman: Neural human radiance field from a single video. In *European Conference on Computer Vision*. Springer, 402–418.
- [16] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 2023. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM Transactions on Graphics* 42, 4 (July 2023). <https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/>
- [17] M. Kim, J. Park, and K. Lee. 2023. Monocular Differentiable Rendering for Self-Supervised 3D Amodal Masks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 789–798.
- [18] Diederik P Kingma, Max Welling, et al. 2013. Auto-encoding variational bayes.
- [19] Muhammed Kocabas, Jen-Hao Rick Chang, James Gabriel, Oncel Tuzel, and Anurag Ranjan. 2024. Hugs: Human gaussian splats. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 505–515.
- [20] Wei-Sheng Lai, Jia-Bin Huang, Oliver Wang, Eli Shechtman, Ersin Yumer, and Ming-Hsuan Yang. 2018. Learning blind video temporal consistency. In *Proceedings of the European conference on computer vision (ECCV)*. 170–185.
- [21] Inhee Lee, Byungjun Kim, and Hanbyul Joo. 2024. Guess the unseen: Dynamic 3d scene reconstruction from partial 2d glimpses. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 1062–1071.
- [22] P. Li, Q. Zhang, and R. Others. 2022. Compositional Models for Amodal Layout Completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2345–2354.
- [23] Jiaqi Lin, Zhihao Li, Xiao Tang, Jianzhuang Liu, Shiyong Liu, Jiayue Liu, Yangdi Lu, Xiaofei Wu, Songcen Xu, Youliang Yan, et al. 2024. Vastgaussian: Vast 3d gaussians for large scene reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5166–5175.
- [24] Haoyu Lu, Guoxing Yang, Nanyi Fei, Yuqi Huo, Zhiwu Lu, Ping Luo, and Mingyu Ding. 2023. Vdt: General-purpose video diffusion transformers via mask modeling. *arXiv preprint arXiv:2305.13311* (2023).
- [25] Tao Lu, Mulini Yu, Liming Xu, Yuanbo Xiangli, Limin Wang, Dahua Lin, and Bo Dai. 2024. Scaffold-gs: Structured 3d gaussians for view-adaptive rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 20654–20664.
- [26] Jilan Mei, Junbo Li, and Cai Meng. 2024. GS2Pose: Tow-stage 6D Object Pose Estimation Guided by Gaussian Splatting. *arXiv preprint arXiv:2411.03807* (2024).
- [27] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. *Commun. ACM* 65, 1 (2021), 99–106.
- [28] Michał Nazarczuk, Thomas Tanay, Sibi Catley-Chandar, Richard Shaw, Radu Timofte, and Eduardo Pérez-Pellitero. 2024. AIM 2024 sparse neural rendering challenge: Dataset and benchmark. *arXiv preprint arXiv:2409.15041* (2024).
- [29] H. Nguyen, T. Davis, and X. Xu. 2022. Learning Disentangled Shape-Texture for Amodal Completion. In *Advances in Neural Information Processing Systems (NeurIPS)*. 1–12.
- [30] Alex Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.
- [31] Nikhil Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädele, Chloe Rolland, Laura Gustafson, Eric Minturn, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollar, and Christoph Feichtenhofer. 2024. SAM 2: Segment Anything in Images and Videos. *arXiv preprint arXiv:2408.00714* (2024). <https://arxiv.org/abs/2408.00714>
- [32] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 10684–10695.
- [33] Johannes L Schonberger and Jan-Michael Frahm. 2016. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4104–4113.
- [34] Adam Sun, Tiange Xiang, Scott Delp, Fei-Fei Li, and Ehsan Adeli. 2024. Occfusion: Rendering occluded humans with generative diffusion priors. *Advances in Neural Information Processing Systems* 37 (2024), 92184–92209.
- [35] Cheng Sun, Min Sun, and Hwann-Tzong Chen. 2022. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 5459–5469.
- [36] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshit Goka, Kiwoong Park, and Victor Lempitsky. 2022. Resolution-robust large mask inpainting with fourier convolutions. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*. 2149–2159.
- [37] Z. Teed and J. Deng. 2020. RAFT: Recurrent All-Pairs Field Transforms for Optical Flow. In *European Conference on Computer Vision (ECCV)*. 402–419.
- [38] Yihan Wang, Lahav Lipson, and Jia Deng. 2024. Sea-raft: Simple, efficient, accurate raft for optical flow. In *European Conference on Computer Vision*. Springer, 36–54.
- [39] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* 13, 4 (2004), 600–612.

- [40] J. Wu, Z. Yang, and H. Kim. 2022. Self-Supervised Amodal Reconstruction from Single Images. In *European Conference on Computer Vision (ECCV)*. 341–356.
- [41] Xianghui Xie, Bharat Lal Bhatnagar, Jan Eric Lenssen, and Gerard Pons-Moll. 2024. Template Free Reconstruction of Human-object Interaction with Procedural Interaction Generation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [42] Chen Yang, Sikuang Li, Jiemin Fang, Ruofan Liang, Lingxi Xie, Xiaopeng Zhang, Wei Shen, and Qi Tian. 2024. GaussianObject: High-Quality 3D Object Reconstruction from Four Views with Gaussian Splatting. *ACM Transactions on Graphics (TOG)* 43, 6 (2024), 1–13.
- [43] Shuai Yang, Yifan Zhou, Ziwei Liu, and Chen Change Loy. 2023. Rerender a video: Zero-shot text-guided video-to-video translation. In *SIGGRAPH Asia 2023 Conference Papers*. 1–11.
- [44] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. 2018. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 586–595.
- [45] Zhixing Zhang, Bichen Wu, Xiaoyan Wang, Yaqiao Luo, Luxin Zhang, Yinan Zhao, Peter Vajda, Dimitris Metaxas, and Licheng Yu. 2024. Avid: Any-length video inpainting with diffusion model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7162–7172.
- [46] Shangchen Zhou, Chongyi Li, Kelvin CK Chan, and Chen Change Loy. 2023. Propainter: Improving propagation and transformer for video inpainting. In *Proceedings of the IEEE/CVF international conference on computer vision*. 10477–10486.
- [47] Shangchen Zhou, Peiqing Yang, Jianyi Wang, Yihang Luo, and Chen Change Loy. 2024. Upscale-a-video: Temporal-consistent diffusion model for real-world video super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2535–2545.
- [48] X. Zhou, Y. Li, Z. Wang, and T. Others. 2023. Amodal Instance Segmentation with Transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 1234–1243.
- [49] Zehao Zhu, Zhiwen Fan, Yifan Jiang, and Zhangyang Wang. 2024. Fsgs: Real-time few-shot view synthesis using gaussian splatting. In *European conference on computer vision*. Springer, 145–163.

Temporally Consistent Amodal Completion for 3D Human-Object Interaction Reconstruction

Supplementary Material

Overview

This supplementary material introduces further details and experimental results of our paper, Temporally Consistent Amodal Completion for 3D Human-Object Interaction Reconstruction.

- Section A elaborates on the preliminary experiments conducted prior to proposing our full pipeline.
- Section B describes the datasets used, including BEHAVE [4] and InterCap [11, 12].
- Section C details the implementation of amodal completion and 3D reconstruction.
- Section D provides a detailed explanation of data selection process and the evaluation metrics for amodal completion, temporal consistency, and 3D reconstruction.
- Section E presents additional experiments on temporal consistency and occlusion identification.

A Preliminary Experiments

Before proposing our full pipeline, we conducted preliminary experiments to explore temporal consistency and occlusion identification across frames. Specifically, we examined two strategies for achieving temporal consistency, fixing the random seed and applying a latent shift, and one strategy for occlusion identification, using a human mask and a background mask.

Temporal Consistency with Fixed Random Seed. Figure 6 illustrates the effect of using a fixed random seed on temporal consistency. We generate two different frames with the same random seed, meaning that Stable Diffusion performs amodal completion from identical initial noise. As shown in the figure, the completion results lack consistency, even though the object remains static during the human-object interaction in both frames. This experiment demonstrates that fixing the random seed alone is insufficient to ensure temporal consistency.

Temporal Consistency with Latent Shift. Since fixing the random seed was ineffective in maintaining temporal consistency, we hypothesized that leveraging latent features extracted from the original frame could help preserve consistency across frames. To isolate the effect of the latent representation, we applied a mask to retain only the target object’s region.

However, due to the dynamic nature of human-object interaction, the target object moves across frames. To address this, we hypothesized that spatially shifting the latent features could mitigate temporal inconsistency. As shown in Figure 7, we applied the latent from the first frame to subsequent frames by aligning it with the target object’s new position, without accounting for rotation. While the shifted latent features improved temporal consistency to some extent, the results suggest that spatial shifting alone—without

adjusting for the object’s rotation—is insufficient to fully preserve temporal consistency and accurately complete occluded regions.

Occlusion Identification with Human Mask and Background Mask. According to Figure 6, applying a human mask does not reliably capture occluded regions. Moreover, using the entire background as a mask, excluding only the unoccluded region of the target object, results in an entirely new object that aligns only with the text prompt, as illustrated in Figure 8. As a result, the completed areas often deviate from the true object shape, indicating that the outputs are not temporally consistent and geometrically accurate.

B Dataset

BEHAVE [4]. The BEHAVE dataset consists of 321 RGB-D video sequences of indoor human–object interactions, recorded using four Kinect cameras. The test set includes 3 human subjects interacting with 20 different objects. For evaluation, we select 18 object sequences—excluding “keyboard” and “basketball” due to missing ground-truth pose annotations in the 30 FPS version. Additionally, we select 3 sequences, one for each subject, to evaluate human performance. In total, our pipeline is applied to approximately 27,000 frames across all selected sequences.

InterCap [11, 12]. The InterCap dataset comprises 223 RGB-D videos of human–object interactions, captured from 6 camera viewpoints and involving 10 objects and 10 human subjects. We extract 10 representative videos that collectively cover all 10 object categories used in the dataset.

C Implementation Details

C.1 Baselines

We compare our approach against several state-of-the-art baselines for amodal completion and video inpainting. The selected methods represent diverse design philosophies and serve to evaluate different aspects of our method, including spatial plausibility and temporal consistency.

Pix2Gestalt [2]. This method performs amodal completion by hallucinating occluded object parts via semantic segmentation and image generation. Although designed for static images, it is widely adopted in amodal completion literature and serves as a strong baseline for appearance-level plausibility under occlusion.

Stable Diffusion Inpainting (SD Inpainting) [32]. A popular diffusion-based inpainting method that generates high-quality content guided by masked regions. We adopt this model as a zero-shot inpainting baseline to assess the generative quality of completed regions. However, it lacks temporal modeling, making it susceptible to frame-wise inconsistency in videos.

LaMa [36]. LaMa is a fast, high-resolution image inpainting model that leverages fast Fourier convolutions. Its strong performance

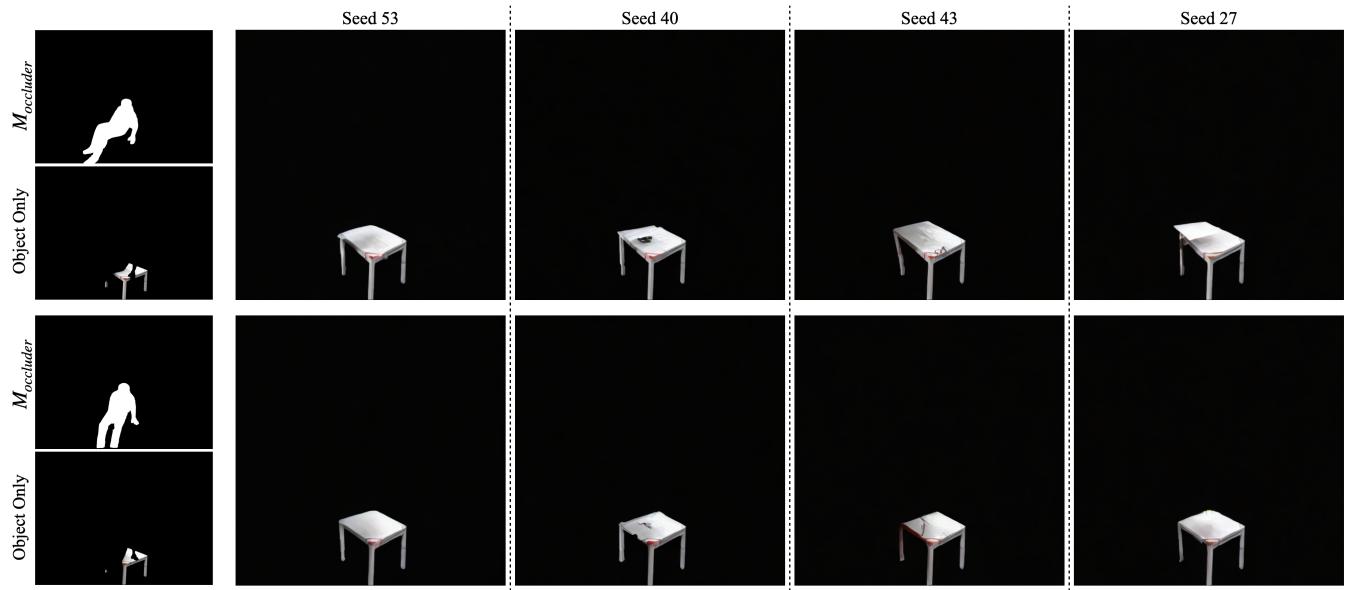


Figure 6: Fixed seed preliminary experiment for temporal consistency. Two frames are completed using four different random seeds to assess consistency.

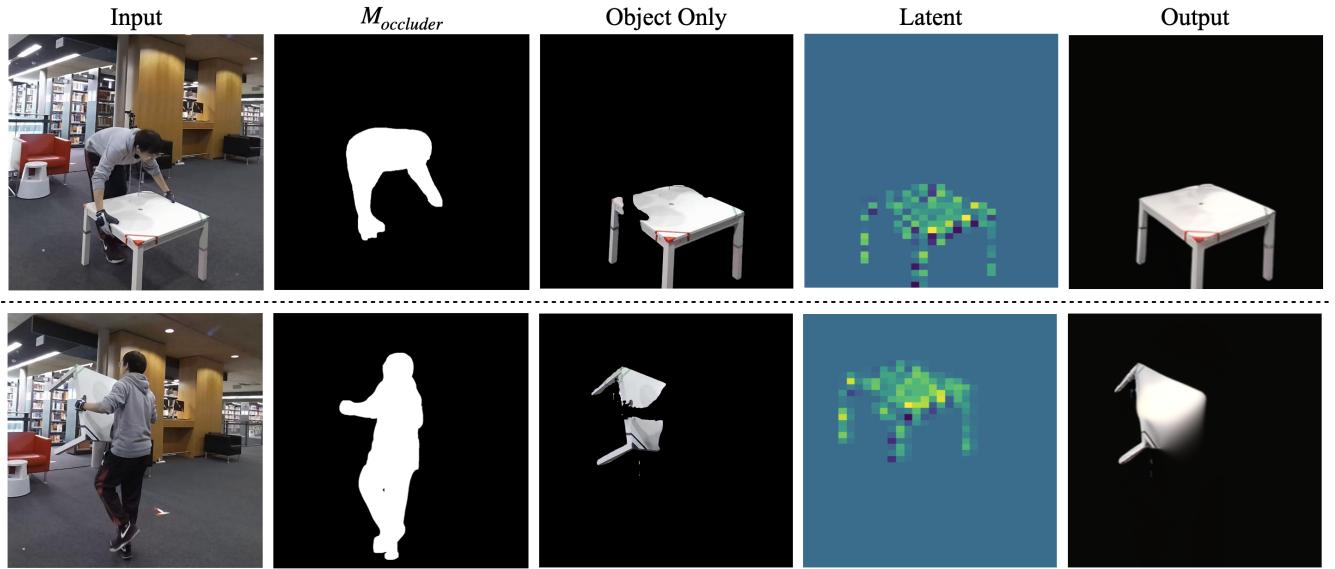


Figure 7: Preliminary experiment using latent shift for temporal consistency. The latent feature from the first frame is shifted to the second frame to evaluate whether temporal consistency can be maintained.

in image-level inpainting tasks makes it a competitive baseline, especially for spatial reconstruction fidelity. Like SD Inpainting, it does not consider temporal coherence.

VDT [24]. VDT is a video diffusion transformer trained for temporally consistent video inpainting. It leverages a space-time attention mechanism and is specifically designed for video scenarios. This method serves as a strong baseline for evaluating the temporal stability of inpainted sequences.

C.2 Amodal Completion.

We use Stable Diffusion Inpainting [32] with a guidance scale of 6.0, keeping all other hyperparameters consistent with the default configuration. For Bidirectional Temporal Feature (BTF) Warping, we adopt SeaRaft [38] for optical flow estimation and follow its official implementation settings. The window size n is empirically set to 7 based on performance comparison, as shown in Table 6; smaller window sizes ($n = 1, 3$) result in inferior performance across all evaluation metrics.

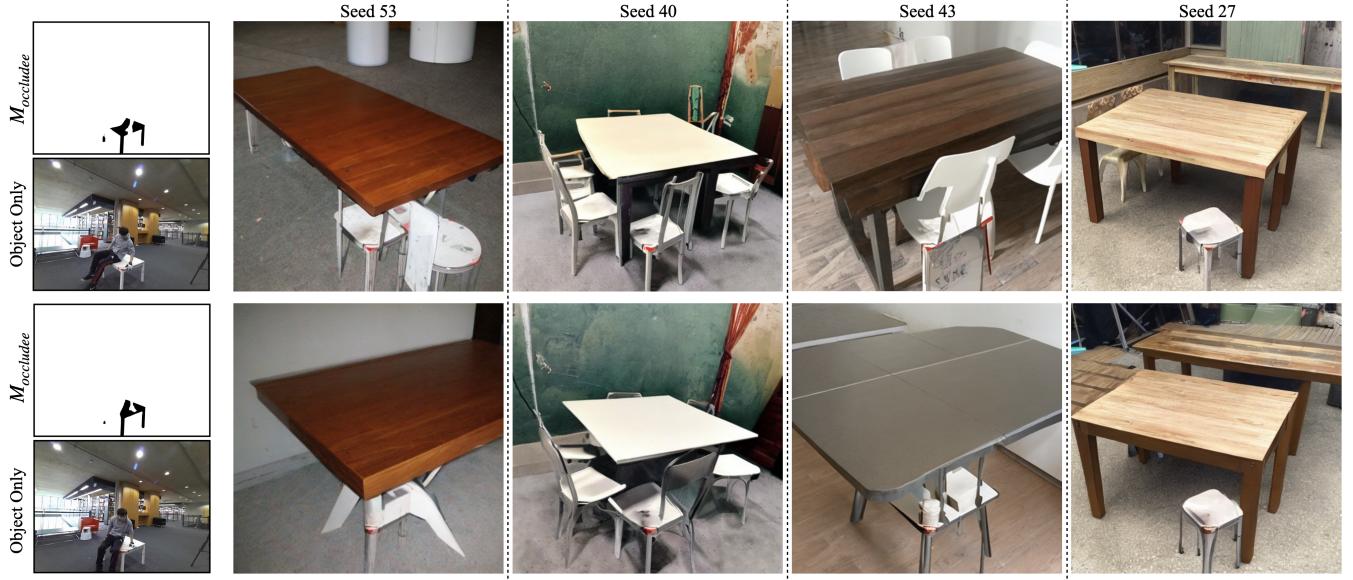


Figure 8: Preliminary experiment using a background mask for occlusion identification. The mask excludes only the unoccluded region of the target object.

C.3 3D reconstruction.

For object reconstruction, we use 30,000 training iterations and set the weighting parameter $\lambda = 0.2$ for the photometric loss \mathcal{L}_{photo} , which combines L1 and SSIM terms. All other hyperparameters follow the settings of GS-Pose [26] for objects and GaussianAvatar [10] for humans.

Since the 30 FPS object poses and SMPL parameters provided in the BEHAVE dataset [4] are not temporally aligned or spatially accurate, we observed noticeable misalignment between the rendered mesh and the inpainted object masks generated by SAM2 [31]. To mitigate this issue during 3D Gaussian Splatting [16] training, we use the provided object masks aligned with 6-DoF object poses instead of our predicted masks for object reconstruction. This assumes the inpainting process yields complete and reliable textures. Nevertheless, for evaluation purposes, we evaluate reconstruction performance using outputs generated solely by our full pipeline.

D Evaluation

Data Selection. To ensure rigorous and fair evaluation, we carefully curated samples from the BEHAVE [4] and InterCap [11, 12] datasets based on the level of object occlusion. Specifically, we include only frames in which the object occlusion ratio falls between 15% and 70%. Frames with minimal occlusion provide limited challenge for amodal completion, while those with severe occlusion lack sufficient visual cues for meaningful reconstruction. This strict filtering balances task difficulty, ensures evaluation stability, and prevents inflation of dataset size with low-signal or overly ambiguous cases.

Amodal Completion. To evaluate the quality of amodal completion, we use two complementary metrics: CLIP Score [30] and Intersection over Union (IoU). The CLIP score quantifies the semantic alignment between each inpainted image and a corresponding

Object	Size	Volume
tablesquare	557.56	1.68
chairwood	396.51	0.68
chairblack	294.38	1.56
tablesmall	233.93	0.42
yogaball	196.26	8.27
monitor	150.48	1.00
boxlarge	139.68	6.35
plasticcontainer	132.87	0.82
yogamat	125.55	3.24
boxlong	92.17	0.96
stool	87.29	2.64
backpack	83.14	2.98
suitcase	67.32	3.65
boxmedium	40.44	2.96
boxsmall	26.73	1.29
trashbin	30.45	1.75
toolbox	13.34	0.52
boxtiny	6.96	0.27

Table 5: 3D bounding box size and volume for each object in BEHAVE [4], sorted by descending size. *boxtiny* and *toolbox* have the smallest spatial extent and are excluded from 3D reconstruction experiments.

category-level textual prompt using CLIP’s vision–language embedding space. This allows us to assess whether the completed region semantically resembles the intended object class. To reduce the impact of unrelated background pixels, we compute CLIP scores within tight bounding boxes centered on the reconstructed objects.

For spatial accuracy, we report the IoU between the predicted amodal masks and the ground-truth masks. Since our model outputs



Figure 9: Qualitative Ablation of Our methods.

RGB images rather than masks, we employ SAM2 [31] to segment the inpainted regions and generate amodal predictions. IoU is then computed as the ratio of the intersection and union areas between the predicted and ground-truth object masks.

Temporal Consistency. We assess temporal consistency using two metrics: the Temporal Consistency score (TC score) [7] and Flow Warping Error [20].

The TC score is designed to measure perceptual smoothness across time. Specifically, we extract CLIP image embeddings from each frame of the output video and compute the average cosine similarity between embeddings of consecutive frames. A higher TC score implies better visual consistency over time, as perceptually similar frames yield higher similarity in CLIP space.

Flow Warping Error captures temporal misalignment at the pixel level. We first estimate optical flow between pairs of consecutive ground-truth object masks using SEA-RAFT [38]. Using the flow, the previously inpainted frame is warped to align with the current frame. The pixel-wise discrepancy between the warped frame and the actual current frame is then measured as the Flow Warping Error. A lower Flow Warping Error indicates higher temporal coherence and fewer frame-to-frame inconsistencies.

3D reconstruction. We evaluate 3D reconstruction on 16 out of the 18 objects used in our amodal completion and temporal consistency experiments. Following our setup, we sample video frames at 1 FPS, using only 15–47 frames per sequence. Two objects (“box tiny” and “tool box”) are excluded from evaluation due to their

Window Size n	IoU \uparrow	CLIP \uparrow	Warp-err ($\times 10^{-3}$) \downarrow	TC \uparrow
1	60.84	27.80	5.73	97.17
3	60.75	27.80	5.70	97.17
7 (Ours)	61.23	27.81	5.69	97.18

Table 6: Effect of window size n in Bidirectional Temporal Feature (BTF) Warping. Using a larger window ($n = 7$) provides improved alignment and better performance across IoU, CLIP score, warp error, and temporal consistency (TC), compared to smaller window sizes.

extremely small projected area and volume, which result in poor reconstruction quality (see Table 5 for object statistics).

Since our method reconstructs only the object or the human the background—we evaluate reconstruction quality using masked variants of standard metrics: Peak Signal-to-Noise Ratio (PSNR-M), Structural Similarity Index Measure (SSIM-M)[39], and Learned Perceptual Image Patch Similarity (LPIPS-M)[44]. Following the evaluation protocol in [28], we compute these metrics within tight axis-aligned bounding boxes surrounding the reconstructed human and object regions, rather than over the entire frame. To minimize background influence and better focus on the target regions, we crop around the reconstructed masks using a tight object mask with a small scale margin ($\times 1.2$), ensuring a highly localized evaluation.

This cropped evaluation is crucial because baseline methods often produce geometrically inconsistent reconstructions that fail to preserve the spatial extent of the original subject. As a result, evaluating over the full image would unfairly penalize such methods due to misalignment with the ground truth. By focusing the evaluation within localized bounding boxes, we reduce the influence of irrelevant background pixels and obtain a more faithful assessment of perceptual and geometric reconstruction quality.

E Additional Experiments

Ablation Study on Temporal Consistency. Figure 9 shows qualitative results comparing different ablations. Our method maintains consistent appearances across consecutive frames, while ablated variants suffer from noticeable temporal inconsistencies. This highlights the effectiveness of our temporal consistency strategy.

Effect of Temporal Window Size. To evaluate the impact of temporal context, we conduct an ablation study on the window size n used in the Bidirectional Temporal Feature (BTF) Warping module, as shown in Table 6. We compare three settings: $n = 1$, $n = 3$, and $n = 7$ (our default). Results show that increasing the window size leads to consistent improvements across all metrics. Specifically, our setting ($n = 7$) achieves the highest IoU (61.23), CLIP score (27.81), and temporal consistency (97.18), while also yielding the lowest warping error (5.69). This demonstrates that incorporating a broader temporal context enables more accurate feature alignment and robust reconstruction, particularly in the presence of complex motion and occlusion.