IMPLICIT 4D GAUSSIAN SPLATTING FOR FAST MOTION WITH LARGE INTER-FRAME DISPLACEMENTS

Anonymous authors

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027 028 029

031

033

034

037 038

040 041

042

043

044

046 047

048

051

052

Paper under double-blind review

ABSTRACT

Recent 4D Gaussian Splatting (4DGS) methods often fail under fast motion with large inter-frame displacements, where Gaussian attributes are poorly learned during training, and fast-moving objects are often lost from the reconstruction. In this work, we introduce Spatiotemporal Position Implicit Network for 4DGS, coined SPIN-4DGS, which learns Gaussian attributes from explicitly collected spatiotemporal positions rather than modeling temporal displacements, thereby enabling more faithful splatting under fast motions with large inter-frame displacements. To avoid the heavy memory overhead of explicitly optimizing attributes across all spatiotemporal positions, we instead predict them with a lightweight feed-forward network trained under a rasterization-based reconstruction loss. Consequently, SPIN-4DGS learns shared representations across Gaussians, effectively capturing spatiotemporal consistency and enabling stable high-quality Gaussian splatting even under challenging motions, while also reducing storage overhead by avoiding the need for explicit parameter storage. Across extensive experiments, SPIN-4DGS consistently achieves higher fidelity under large displacements, with clear improvements in PSNR and SSIM on challenging sports scenes from the CMU Panoptic dataset. For example, SPIN-4DGS notably outperforms the strongest baseline, D3DGS, by achieving +1.83 higher PSNR on the Basketball scene.



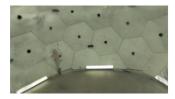
Figure 1: Faithful reconstruction of fast motion with large inter-frame displacements. Existing 4DGS approaches often produce blurred or incomplete reconstructions of fast-moving objects. In contrast, ours successfully reconstructs clear and accurate details, such as the basketball in the scene.

1 Introduction

Rendering fast motions with large inter-frame displacements remains challenging for dynamic scene reconstruction, despite its importance for a wide range of real-world applications. Recent advances in 4D Gaussian Splatting (4DGS) have shown remarkable efficiency and visual quality, making it a promising framework for dynamic scene reconstruction. In particular, existing 4DGS methods (Yang et al., 2023; Duan et al., 2024; Wu et al., 2024; Xu et al., 2024) achieves strong results on dynamic scenes with small displacements (*e.g.*, Neu3D (Li et al., 2022)) across video frames.

However, as motions become faster and inter-frame shifts grow larger (e.g., Panoptic Sports (Joo et al., 2015)), existing 4DGS methods, including explicit 4D parametrization (Yang et al., 2023; Duan et al., 2024) and deformable approaches (Wu et al., 2024; Xu et al., 2024; Kwak et al., 2025), often fail to capture the rapid dynamics, producing blurred or even vanished objects. To be specific, in deformable approaches, Gaussians are defined in a static canonical space and transformed over time through learned deformations. However, they often fail to assign initial Gaussians for fast-moving objects in the canonical space, causing those objects to remain unseen during deformation training. On the other hand, although explicit parameterization roughly tracks Gaussian positions





(a) Left: 15K training iteration, Right: 30K training iteration.

(b) Canonical space

Figure 2: Failure modes on fast motions with large inter-frame displacements. We visualize failure modes of existing frameworks; (2a) explicit parameterization and (2b) deformable methods. Figure (2a) shows drastic degradation on training iterations (i.e., $15K \rightarrow 30K$), and (2b) shows the canonical space of deformable initialization fails to assign Gaussians for fast motions.

that correspond to fast motions at the early stage of training, their attributes, including color, opacity, scale, and rotation, rapidly collapse in later training, leading to drastic degradation. As a result, both approaches show failure modes with blurred or even vanished motions, as shown in Figure 2.

To this end, we focus on addressing the failure in learning Gaussian attributes for fast-moving objects with large displacements. Although positions remain sufficiently accurate to capture motion, other attributes collapse more easily. This degradation arises because reconstruction loss is dominated by background Gaussians; fast-moving Gaussians at new positions incur higher reconstruction errors, while static backgrounds are easier to fit. As a result, both deformable and explicit 4DGS approaches tend to bias learning toward background fitting, eventually leading dynamic objects to disappear. In addition, such large displacements can cause cross-frame interference during frame-by-frame rasterization. Although a 4D Gaussian (Wu et al., 2024) can be sliced differently at each time, its parameters are shared, so optimizing for one frame makes other slices suboptimal unless we separate Gaussians by explicit spatiotemporal positions (x, y, z, t). This observation motivates us to leverage the explicit spatiotemporal positions of fast-moving Gaussians as inputs for generating their attributes, thereby achieving more faithful splatting under large displacements.

In this paper, we introduce SPIN-4DGS (Spatiotemporal Position Implicit Network for 4DGS), a lightweight yet effective framework designed to handle fast motions with large inter-frame displacements. To be specific, we first estimate high-quality spatiotemporal positions (x,y,z,t) of Gaussians that can serve as the inputs for later attribute prediction, initially gathering them across the entire scene before refining them in a frame-wise manner. Then, we leverage these positions to predict Gaussian attributes through a lightweight feed-forward network, avoiding the heavy memory overhead of explicitly optimizing attributes over all spatiotemporal positions. This design learns a shared implicit representation across all Gaussians and decodes attributes directly from positions under rasterization loss. As a result, learned attributes remain consistent across positions, and dynamic objects can stay stable even under large displacements. Meanwhile, attributes are stored implicitly in network parameters, rather than explicitly for each Gaussian, which significantly improves memory efficiency on a large number of spatiotemporal positions.

To validate the effectiveness of the proposed SPIN-4DGS, we perform experiments on various sports scenes from the CMU Panoptic Sports dataset, where human motions are rapid and small objects move across large inter-frame displacements. Across six sports scenes, SPIN-4DGS achieves the best performances, significantly outperforming all baselines. For example, SPIN-4DGS achieves a +1.83 PSNR dB improvement compared to the strongest baseline, D3DGS (Luiten et al., 2024), with a higher SSIM of 0.92 on the Basketball scene. Specifically, qualitative results in Figure 4 and Table 1 demonstrate that prior methods often blur Gaussians or fail to capture fast-moving objects (e.g., Basketball), whereas SPIN-4DGS preserves them with sharp and stable reconstructions. Interestingly, we further observe that SPIN-4DGS significantly improves performance when reusing pre-trained Gaussian positions from strong baselines, such as 4DGS (Yang et al., 2023) and D3DGS.

Overall, our work introduces SPIN-4DGS, the first framework that learns Gaussian attributes directly from spatiotemporal positions, enabling stable and high-quality 4DGS under large inter-frame displacements. In addition, our design preserves the rendering efficiency of prior 4DGS methods while improving storage and training stability, thereby enhancing the efficiency—quality balance required for practical deployment. We believe SPIN-4DGS mitigates a fundamental challenge of 4DGS by enabling stable learning in various real-world scenarios and opens a new direction for advancing dynamic scene representation.

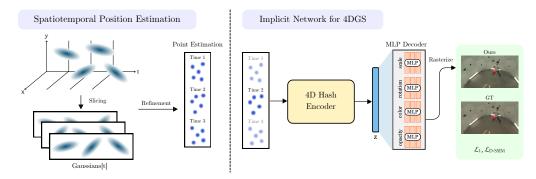


Figure 3: **Illustration of the overall framework of SPIN-4DGS.** SPIN-4DGS consists of two stages of (a) Spatiotemporal Position Estimation and (b) Implicit Network for 4DGS. Specifically, (a) we slice Gaussians along the temporal axis to obtain spatiotemporal position sets and refine them with rasterization loss. Then, (b) the refined positions are normalized and passed through a 4D Hash encoder and multi-branch decoders to predict Gaussian attributes (scale, rotation, color, and opacity).

2 Method

In this section, we present the Spatiotemporal Position Implicit Network for 4DGS (SPIN-4DGS), a framework for reconstructing dynamic scenes under motions with large inter-frame displacements. We first review 3D Gaussian Splatting as preliminaries in section 2.1. Then, section 2.2 describes how we obtain explicit spatiotemporal Gaussian positions, and section 2.3 details how their attributes are predicted via a feed-forward implicit network. The overall framework is illustrated in Figure 3.

2.1 PRELIMINARY: 3D GAUSSIAN SPLATTING

3D Gaussian Splatting (3DGS; Kerbl et al. (2023)) provides a differentiable volume rendering representation. It introduces anisotropic Gaussians where the parameters (*i.e.*, position, scale, rotation, color, and opacity) of each Gaussian are defined and optimized for tile-based rendering. Structure-from-Motion (SfM; Schonberger & Frahm (2016)) techniques estimate the initial positions and colors of Gaussians from input images. For a given arbitrary point $\mathbf{x} \in \mathbb{R}^3$ in the 3D scene, Gaussian is defined as follows:

$$G^{3D}(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\top} \boldsymbol{\Sigma}_{3D}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right), \tag{1}$$

where $G^{3D}(\mathbf{x})$ denotes the value of the Gaussian at arbitrary point \mathbf{x} . The parameters of each Gaussian include the position $\boldsymbol{\mu} \in \mathbb{R}^3$, opacity $o \in \mathbb{R}$, color $\mathbf{c} \in \mathbb{R}^3$ represented using spherical harmonics coefficients, and the covariance matrix $\mathbf{\Sigma} \in \mathbb{R}^{3\times 3}$, which is defined in terms of a diagonal scale matrix $\mathbf{S} = \operatorname{diag}(s_1, s_2, s_3)$ and a rotation matrix $\mathbf{R} \in \mathrm{SO}(3)$, obtained from quaternions. To ensure that the covariance matrix remains $\mathbf{\Sigma}$ positive semi-definite, it is computed as follows:

$$\Sigma_{3D} = \mathbf{R} \mathbf{S} \mathbf{S}^{\mathsf{T}} \mathbf{R}^{\mathsf{T}}.$$
 (2)

To render via rasterization, the 3D Gaussians are first projected onto the 2D image plane. This is done by applying the viewing transformation **W** and the Jacobian matrix (Zwicker et al., 2001) **J** to compute the 2D covariance matrix Σ_{2D} as follows:

$$\Sigma_{2D} = \mathbf{J} \, \mathbf{W} \, \Sigma_{3D} \, \mathbf{W}^{\mathsf{T}} \mathbf{J}^{\mathsf{T}} \tag{3}$$

During the rendering process, pixel values are computed via alpha blending. The alpha value (i.e., opacity) of each Gaussian is obtained by projecting the 3D Gaussian G_i into 2D as G_i^{2D} . Specifically, for each pixel, the alpha value α_i' and the resulting color C are computed as

$$\alpha_i' = o_i G_i^{2D}(\mathbf{x}), \quad \mathbf{C}(\mathbf{x}) = \sum_{i=1}^N c_i \alpha_i' \prod_{j=1}^{i-1} (1 - \alpha_j'), \tag{4}$$

where o_i is the intrinsic opacity of the *i*-th Gaussian, c_i its color, and N the total number of Gaussians contributing to the pixel.

2.2 SPATIOTEMPORAL GAUSSIAN POSITIONS

In this section, our goal is to construct a Gaussian point set that fully spans the trajectory of dynamic objects, ensuring sufficient coverage for faithful reconstruction. Prior methods perform reasonably well in regions where Gaussians are densely clustered, but under large inter-frame displacements, they often fail to maintain enough points around fast-moving objects. Even when positions are sufficiently available, attributes remain challenging. To be specific, when Gaussians spanning the entire trajectory are optimized jointly, frame-specific supervision signals interfere with one another. Because rasterization is performed frame by frame, updates that make a Gaussian optimal for one timestamp can render it suboptimal for others. This cross-frame interference can potentially weaken the learning signal and degrade the quality of the reconstruction.

To overcome these issues, we construct Gaussian sets independently at each time step, explicitly separating them by spatiotemporal positions. This formulation avoids interference across frames and enables more reliable attribute learning for fast-moving objects under large displacements.

$$u_t = f_\theta(\mathbf{x}, \mathbf{y}, \mathbf{z}, \mathbf{t}), \qquad \mathbf{c}_t = g_\phi(\mathbf{x}, \mathbf{y}, \mathbf{z}, \mathbf{t}) \in \{0, \dots, 255\}^3, \qquad t \in \{1, \dots, T\}.$$
 (5)

Here, F_{θ} and g_{ϕ} are explicit functions of the temporal axis predicting the position and color corresponding to each frame t. For example, the explicit approaches (Duan et al., 2024; Yang et al., 2023) construct points by performing time slicing along the temporal axis, whereas a deformable approach (Wu et al., 2024; Xu et al., 2024; Bae et al., 2024) network can learn to construct them as time-varying structures. For SPIN-4DGS, we estimate spatiotemporal positions by an explicit method (Yang et al., 2023) as a default. Lastly, we further refine the estimated position by utilizing rasterization loss with corresponding colors at every frame to densify salient points and prune unnecessary ones.

$$u_t \leftarrow \text{Refine}(u_t, \mathbf{c}_t; t), \qquad t = 0, \dots, T.$$
 (6)

After refinement, the Gaussian points at each timestamp t are fixed as explicit spatiotemporal positions and directly fed into our implicit network to learn the corresponding attributes.

2.3 IMPLICIT NETWORK FOR 4D GAUSSIAN SPLATTING

In this section, we describe the proposed implicit network for predicting 4D Gaussian parameters, including opacity $o \in \mathbb{R}^1$, spherical harmonics coefficients $sh \in \mathbb{R}^{48}$, scale $\mathbf{s} \in \mathbb{R}^3$, and rotation (unit quaternion) $\mathbf{r} \in \mathbb{R}^4$, from inputs of spatiotemporal Gaussian positions $(\mu, t) \in \mathbb{R}^4$.

Input position normalization. At each frame t, the spatiotemporal Gaussian positions serve as inputs to a feed-forward network for learning implicit representations of other Gaussian parameters. Since the raw 3D Gaussian positions $\mathbf{u} \in \mathbb{R}^3$ on the scene are unbounded, we apply a normalization step to stabilize learning and preserve representational capacity. Specifically, following the scene contraction strategy of Mip-NeRF (Barron et al., 2021), we first compress the coordinates into a finite ball and then map them to the normalized range $[0,1]^3$, as defined in equation 7.

$$\operatorname{contract}(\boldsymbol{\mu}) = \begin{cases} \boldsymbol{\mu}, & \|\boldsymbol{\mu}\| \le 1, \\ \left(2 - \frac{1}{\|\boldsymbol{\mu}\|}\right) \frac{\boldsymbol{\mu}}{\|\boldsymbol{\mu}\|}, & \|\boldsymbol{\mu}\| > 1, \end{cases} \quad \bar{\boldsymbol{\mu}} = \frac{1}{4} \operatorname{contract}(\boldsymbol{\mu}) + \frac{1}{2} \in [0, 1]^3. \quad (7)$$

To keep spatial and temporal scales comparable, we normalize time to [0,1] to match the scale of the spatial embedding. Concretely, we use the current timestamp divided by the total duration:

$$t_{\text{norm}} = \frac{t - t_{\min}}{t_{\max} - t_{\min}} \in [0, 1] \text{ and } \tilde{\mathbf{x}} = \left[\bar{\boldsymbol{\mu}}^{\top}, t_{\text{norm}}\right]^{\top} \in [0, 1]^{4}.$$
 (8)

Our network adopts an encoder-decoder architecture, where the encoder f_{enc} maps the normalized input $\tilde{\mathbf{x}}$ to a latent embedding $f_{enc}(\tilde{\mathbf{x}})$ which the decoder takes as input.

Encoder architecture for shared latent representation. We directly extend the widely used 3D Instant-NGP (Müller et al., 2022) multi-hash grid to 4D by appending the temporal axis, producing a compact latent vector for each input. Given the normalized input $\tilde{\mathbf{x}} = [\bar{\boldsymbol{\mu}}^{\mathsf{T}}, t_{\text{norm}}]^{\mathsf{T}}$, we map it into a unified 4D embedding via a single spatio-temporal encoder. In contrast to low-rank decompositions (*e.g.*, planar), which may reduce computation but degrade expressiveness as the number of

Gaussians grows and incur extra cost from per-level hash management, we avoid such factorizations and employ the 4D hash encoder (Chen et al., 2025) directly as follows:

$$z = f_{enc}(\tilde{\mathbf{x}}) \in \mathbb{R}^{LF} \tag{9}$$

where L denotes the number of hash levels and F the number of channels per level. We concatenate features across all levels to form $z \in \mathbb{R}^{LF}$, which we use as the latent representation.

Decoder architecture for attribute prediction. Given a latent vector \mathbf{z} , we use a multi-branch decoder to produce Gaussian parameters, scale, rotation, spherical harmonics coefficients, and opacity. Each attribute is predicted by a separate head (three-layer MLP) taking the shared encoder output $\mathbf{z} \in \mathbb{R}^{LF}$ as input. All decoder heads use GELU (Hendrycks & Gimpel, 2016) activations instead of ReLU (Agarap, 2018).

$$(\hat{\mathbf{s}}, \hat{\mathbf{r}}, \hat{\mathbf{sh}}, \hat{\mathbf{o}}) = (f_{\text{scale}}(\mathbf{z}), f_{\text{rot}}(\mathbf{z}), f_{\text{sh}}(\mathbf{z}), f_{\text{opacity}}(\mathbf{z})). \tag{10}$$

We convert the raw outputs into valid parameter ranges using attribute-specific activations and post-processing:

$$(\mathbf{s}, \mathbf{r}, \mathbf{sh}, \mathbf{o}) = \left(\exp(\hat{\mathbf{s}}), \frac{\hat{\mathbf{r}}}{\|\hat{\mathbf{r}}\|_2}, \hat{\mathbf{sh}}, \sigma(\hat{\mathbf{o}})\right). \tag{11}$$

We also follow the Gaussian post-processing (Kerbl et al., 2023) pipeline, with a few additional steps for stable training.

Scale post-processing. To prevent gradient explosion from exponential growth, we clip the prescale in the backward pass: $\hat{s} \leftarrow \text{clip}(\hat{s}, \text{max} = 20)$. We also initialize the final-layer bias to -5 to start from a small scale, *i.e.*, $\exp(-5) \approx 0.0067$, and keep it trainable.

Rotation initialization. To bias the predicted quaternion toward the identity at the start of training, we set the final-layer bias to (1,0,0,0), *i.e.*, the first element to 1.0 and the remaining elements to 0, so that the initial output satisfies $\hat{\mathbf{r}} \approx (1,0,0,0)$. This initialization lets the network learn rotations progressively while preserving a stable initial structure.

Opacity post-processing. To stabilize early training and encourage a near-transparent start, we set the final-layer bias to $logit(0.1) \approx -2.197$ (trainable), initializing \hat{o} at ≈ 0.1 . The network then learns to increase opacity only where needed, focusing on informative points.

Color post-processing. The decoder directly regresses color coefficients from encoder embeddings, with no special initialization, to capture the high-dimensional color required for SH-based rendering.

Loss objectives. Finally, frame images are rendered via rasterization, and we optimize the model using the standard 3DGS reconstruction loss as follows:

$$\mathcal{L} = (1 - \lambda) \mathcal{L}_1 + \lambda \mathcal{L}_{\text{D-SSIM}}. \tag{12}$$

where λ is a hyperparameter; we set $\lambda = 0.2$, as following prior works (Kerbl et al., 2023).

3 EXPERIMENTS

In this section, we demonstrate the effectiveness of the proposed method, SPIN-4DGS. Specifically, we choose to employ the recent explicit parameterization method for 4DGS (Yang et al., 2023), which is publicly available, only in the early training stage to estimate spatiotemporal Gaussian positions. These positions are then used in our framework to learn new Gaussian attributes. We then evaluate its ability to capture fast motion on various sports scene benchmarks from the CMU Panoptic Sports dataset (Joo et al., 2015), comparing it with existing 4DGS baselines.

Implementation details. The network uses a hidden dimension of 64, and the encoder is configured with L=16 levels and F=4 features per level; the hash map size 21. Parameter groups in both the encoder and the decoder utilize separate learning rates, as Gaussian attributes (e.g., scale and rotation) are sensitive, and a uniform learning rate often leads to unstable optimization. The encoder learning rate is initialized to 8×10^{-3} ; decoder learning rates are 1×10^{-3} for color, 3×10^{-4} for scale, 3×10^{-5} for rotation, and 8×10^{-4} for opacity. position parameters use the same learning rate as 3DGS. We use Adam (Kinga et al., 2015) with a linear warm-up and cosine

decay schedule, under which every parameter's learning rate decays to 1×10^{-7} . Experiments were conducted on an RTX 4090 GPU (24 GB), and the implementation utilized PyTorch (Paszke et al., 2019) 2.1 with CUDA 11.8. By default, we train for 40K iterations with a batch size of 3. We perform quantitative evaluations using PSNR (Peak Signal-to-Noise Ratio), SSIM (Wang et al. (2004); Structural Similarity Index), and LPIPS (Zhang et al. (2018); Learned Perceptual Image Patch Similarity), and additionally report frames per second (FPS) to assess rendering speed.

Datasets. We employ the CMU Panoptic Sports dataset (Joo et al., 2015) to validate scenarios of fast motions with large inter-frame displacements in our experiments. Specifically, the Panoptic Sports dataset is a challenging benchmark containing six sports scenes: juggle, basketball, boxing, football, softball, and tennis. Each scene is recorded at 30 FPS for 5 seconds (150 frames per scene). A total of 31 cameras were used, and the native resolution of 640×360 was retained. Evaluation was referenced to four fixed test cameras (IDs of 0, 10, 15, and 30).

4DGS baselines. We compare our method with recent 4DGS baselines across various strategies, including (a) explicit parameterizations for 4DGS: 4DGS (Yang et al., 2023), (b) deformable 4DGS: Grid4D (Xu et al., 2024), 4D Gaussian (Wu et al., 2024), and (c) deformable with external supervision: D3DGS (Luiten et al., 2024), TC3DGS (Javed et al., 2024). We note that D3DGS and TC3DGS are designed to handle large inter-frame displacements, such as the Panoptic Sports dataset. All baselines are evaluated following the dynamic scene setups specified in their papers.

Table 1: Comparisons on dynamic sports scenes in the CMU Panoptic Sports dataset. We evaluate ours with existing 4DGS baselines on benchmarks containing fast motions with large interframe displacements. We report PSNR and SSIM for six sports scene sequences across all baselines.

Method	4DGS Category	Bask	etball	Bo	xes	Foo	tball	Jug	gle	Soft	ball	Ten	nis		Av	g.	
		PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	FPS	Storage
Grid4D	Deformable	25.82	0.89	26.81	0.91	27.61	0.91	28.09	0.92	26.91	0.91	27.10	0.91	27.06	0.91	146	333
4DGaussian	Deformable	27.28	0.90	27.32	0.91	28.71	0.91	26.94	0.91	27.24	0.91	27.66	0.91	27.53	0.91	40	62
TC3DGS	External supervision	27.92	0.89	28.28	0.89	28.00	0.89	29.15	0.90	27.96	0.89	25.97	0.89	27.88	0.89	890	49
D3DGS	External supervision	28.22	0.91	29.46	0.91	28.49	0.91	29.48	0.91	28.43	0.91	28.11	0.91	28.70	0.91	760	1994
4DGS	Explicit	27.89	0.92	28.17	0.93	28.35	0.93	28.68	0.93	28.67	0.93	28.50	0.93	28.38	0.93	197	1293
Ours		30.05	0.92	29.91	0.93	29.99	0.93	30.31	0.93	30.24	0.93	30.14	0.93	30.11	0.93	104	1261

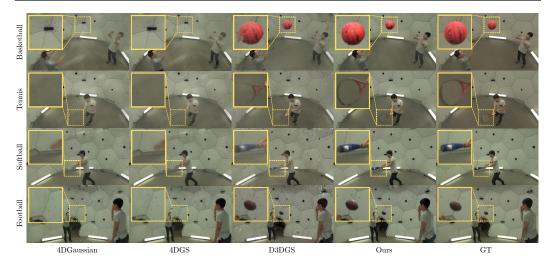


Figure 4: **Visualizations on dynamic sports scenes in the CMU Panoptic Sports dataset.** Compared to prior 4D Gaussian baselines, where fast-moving objects (*e.g.*, ball, bat, racket) often disappear or become corrupted, our method successfully preserves these objects throughout the sequence, producing more faithful and consistent reconstructions.

3.1 EXPERIMENTAL RESULTS ON PANOPTIC SPORTS

In this section, we evaluate our method with various 4DGS baselines on six dynamic sports scenes (Basketball, Boxes, Football, Juggle, Softball, and Tennis) from the CMU Panoptic Sports dataset, which involve rapid human motions and small objects undergoing large inter-frame displacements.

Table 2: **Ablation study on spatiotemporal positions.** We perform ablation studies on (a) varying the early training duration used to estimate spatiotemporal Gaussian positions from the explicit baseline, 4DGS (Yang et al., 2023), evaluated without subsequent refinement, and (b) varying the number of frame-wise refinement iterations applied after position estimation.

(a) Effect of early training duration

(b) Effect of refinement iterations

Iteration	PSNR	SSIM	LPIPS	Time
15K	29.57	0.92	0.15	10m
30K	29.39	0.91	0.15	30m

Iteration	PSNR	SSIM	LPIPS	Time
0.5K	29.86	0.92	0.14	33m
1K	29.89	0.92	0.14	55m
2K	30.05	0.92	0.14	1h 40m

As reported in Table 1, SPIN-4DGS achieves the best PSNR on all six scenes with an average of 30.11 dB, outperforming the strongest explicit baseline 4DGS (28.38 dB) by +1.73 dB and the deformable baseline D3DGS (28.70 dB) by +1.41 dB. SSIM also remains consistently high at 0.93, confirming stable reconstruction quality. In particular, our method even surpasses models trained with external supervision such as segmentation maps (*e.g.*, D3DGS, TC3DGS), showing that high fidelity can be achieved without costly priors. Beyond averages, the improvements are most significant on challenging scenes with extreme motion. On Basketball, SPIN-4DGS improves PSNR by over +1.83 dB against the best baseline, D3DGS (Luiten et al., 2024), while on Tennis, where rackets move rapidly and cover large displacements, the margin also surpasses +1.64 dB against the best baseline, 4DGS (Yang et al., 2023). Overall, those results highlight the robustness of our approach to challenging fast-motion scenarios where prior methods often fail.

Qualitative comparisons in Figure 4 further illustrate these differences. Deformable baseline (*i.e.*, 4DGaussian, first column) results blur fine-grained structures, explicit 4DGS (second) results frequently lose moving objects entirely, and even externally supervised (*i.e.*, D3DGS, third) results suffer degradation on small, fast objects like a tennis racket. In contrast, SPIN-4DGS preserves sharp and stable reconstructions across all frames, closely matching the ground truth.

While lightweight deformable baselines exhibit higher FPS, their reconstructions frequently collapse on fast-moving objects, as shown in Figure 4, making them less competitive for challenging scenarios like sports scenes with rapid motions and small objects. We therefore emphasize comparisons with the stronger baselines, 4DGS and D3DGS. Both require substantial storage (1293MB and 1994MB, respectively), whereas SPIN-4DGS achieves higher fidelity with a smaller footprint of 1261MB. Moreover, while our FPS is relatively lower than 4DGS, it remains within the practical real-time range and is achieved with remarkably better reconstruction quality. We argue that these results demonstrate a more balanced trade-off between efficiency and fidelity; SPIN-4DGS provides more accurate and temporally consistent reconstructions under large inter-frame displacements, without incurring excessive storage or rendering cost.

3.2 ABLATION STUDY AND ANALYSIS

We conduct a series of ablation experiments to demonstrate the proposed method further. First, we analyze the impacts of estimated spatiotemporal position quality and refinement in Table 2 and 3. We also examine a position-reuse setting in Table 4, where positions from existing 4DGS models are provided, to validate the effectiveness of our attributes learning scheme. Finally, we validate the effect of each component of our implicit network in Table 5.

Ablation on spatiotemporal positions. We investigate how the quality of estimated spatiotemporal positions and the amount of refinement affect the final reconstruction. Results are summarized in Table 2, evaluated on the Basketball sequence.

- (a) Effect of early training duration. Using positions extracted after 15K iterations of 4DGS (Yang et al., 2023) already achieves strong performance (29.57 PSNR, 0.92 SSIM) with only 10 minutes of cost. Extending training to 30K iterations not only triples the runtime but also slightly degrades the quality, indicating that long optimization is unnecessary once positions are sufficiently stable.
- (b) Effect of refinement iterations. We then vary the number of frame-wise refinement steps after position estimation. Increasing refinement from 0.5K to 2K iterations improves PSNR from 29.86 to 30.05, confirming that moderate refinement shows consistent benefits. In practice, we also prune redundant Gaussians during refinement, which reduces background clutter and focuses updates on salient regions, further enhancing efficiency.

Figure 5: **Qualitative comparison of spatiotemporal slicing.** Without slicing, Gaussians simultaneously represent multiple time steps, causing their contributions to overlap and conflict during rasterization. This results in blurred faces and distorted fast-moving objects. Our slicing explicitly separates Gaussians over time, enhancing qualities with temporally consistent reconstructions.

Effects of input position formulation. We further analyze the effect of spatiotemporal slicing in Table 3 and Figure 5. For a fair comparison, we fix the batch size to 1 and run all experiments on the football sequence, varying only the position design. We compare two settings: (a) w/o spatiotemporal slicing, where all Gaussians are optimized jointly without slicing, and (b) spatiotemporal slicing (ours), where

Table 3: **Ablation on spatiotemporal slicing.** We compare a unified 4D formulation (*i.e.*, w/o slicing), where Gaussian positions are optimized jointly across space–time, against our spatiotemporal slicing strategy that assigns positions per frame.

Spatiotemporal Slicing	$\mathbf{PSNR} \!\!\uparrow$	$\textbf{SSIM} \!\!\uparrow$	Time↓	Train Cost↓
Х	27.48	0.89	1h 20m	18GB
✓	28.96	0.92	25m	9GB

positions are sliced frame by frame and aligned along the time axis.

Without slicing, as shown in Table 3, Gaussians are optimized jointly in a unified 4D space. In this formulation, a single Gaussian must simultaneously explain multiple timestamps. However, rasterization is performed frame by frame, so optimization that reduces the loss for one frame inevitably makes the Gaussian suboptimal for others. This cross-frame interference weakens the supervision signal, slows convergence, and increases both training time and memory usage. In contrast, spatiotemporal slicing explicitly separates Gaussians by (x,y,z,t) and filters out irrelevant points at each frame. This avoids interference across frames and ensures that optimization is focused on the relevant Gaussians for each time step. As a result, slicing achieves both higher reconstruction fidelity and significantly lower training cost.

Qualitative comparisons in Figure 5 confirm these findings. Under the unified 4D formulation, attributes collapse as Gaussians attempt to describe other timestamps, producing blurred faces and distorted fast-moving objects (*e.g.*, football in the scene). Our sliced formulation decouples Gaussians over time, showing sharper details and temporally consistent reconstructions.

Table 4: **Compatibility with existing 4DGS baselines.** We reuse pre-trained positions from D3DGS (Luiten et al., 2024) and 4DGS (Yang et al., 2023), replacing their attribute optimization with our proposed implicit network training. SPIN-4DGS consistently improves PSNR/SSIM across all scenes, highlighting the compatibility and effectiveness of the proposed implicit 4DGS scheme.

Method	Basketball		Boxes		Football		Juggle		Softball		Tennis	
	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
4DGS	27.89	0.92	28.17	0.93	28.35	0.93	28.68	0.93	28.67	0.93	28.50	0.93
4DGS + Ours	29.57	0.92	29.63	0.92	29.42	0.92	29.98	0.92	29.94	0.93	29.83	0.93
D3DGS	28.22	0.91	29.46	0.91	28.49	0.91	29.48	0.92	28.43	0.91	28.11	0.91
D3DGS + Ours	30.50	0.94	30.26	0.94	29.21	0.94	31.04	0.95	30.92	0.95	30.51	0.95

Impacts of implicit 4DGS scheme. We further validate the impacts of the proposed implicit 4DGS scheme by reusing positions from strong baselines and retraining only attributes with our framework. In this setting, we extract pre-trained spatiotemporal positions from D3DGS (Luiten

et al., 2024) and 4DGS (Yang et al., 2023), but discard all their learned attributes. SPIN-4DGS then learns new attributes through its implicit network, while keeping the imported positions learnable to allow for refinement during training. This setup highlights not only the effectiveness of our attribute learning design but also its plug-and-play compatibility with existing 4DGS pipelines, where our scheme consistently enhances reconstruction quality regardless of the underlying position estimator.

As shown in Table 4, SPIN-4DGS consistently improves performance even when initialized with external positions. With D3DGS positions, our method improves PSNR by +2.49 dB and SSIM to 0.95 on the *Softball* scene. Notably, while D3DGS collapses on *Tennis* (28.11/0.91), SPIN-4DGS still achieves 30.51/0.95 with sharp reconstructions. Even with 4DGS positions, our framework still achieves consistent gains overall, demonstrating the effectiveness of the proposed scheme. While SSIM occasionally drops slightly in some cases, we find that this minor loss is easily resolved by our refinement step, showing that performance remains robust and stable.

We find that high-quality positions alone are insufficient for stable reconstructions, as attributes play a critical role in maintaining fidelity under large inter-frame displacements. We also observe that our implicit 4DGS formulation generalizes seamlessly across different baselines, showing strong compatibility and consistent improvements regardless of the source of positions. Taken together, these results highlight SPIN-4DGS as an effective and extensible 4DGS scheme that can enhance a wide range of dynamic scene reconstruction tasks.

Table 5: **Ablation on our implicit network design components.** Starting from the original 4D Hash encoder (Chen et al., 2025) encoder, we analyze the effects of making positions trainable, applying input position normalization, and changing activation functions. Each modification progressively enhances reconstruction quality. All experiments are performed on the Basketball scene.

	Results					
Trainable position	Normalization	Activation	Hash Map Size	PSNR↑	SSIM↑	LPIPS↓
		ReLU	21	18.64	0.78	0.45
✓		ReLU	21	19.17	0.78	0.46
✓	✓	ReLU	21	29.89	0.92	0.16
✓	✓	GELU	21	30.05	0.92	0.14
✓	1	GELU	23	30.25	0.93	0.13

Component analysis. We conduct an ablation study on the design components of our implicit network, summarized in Table 5. Starting from the original 4D Instant-NGP (Chen et al., 2025) encoder with fixed positions and ReLU activations, reconstruction quality is notably poor (18.64 PSNR, 0.78 SSIM, 0.45 LPIPS). Making positions trainable alone achieves only marginal improvement, indicating that position refinement is insufficient without additional modifications. Applying input position normalization shows a significant gain, boosting PSNR by over +10 dB and reducing LPIPS from 0.45 to 0.16, highlighting its importance for stabilized training. Replacing ReLU with GELU further improves fidelity, enlarging the hash map size further achieves the best overall performances (30.25 PSNR, 0.93 SSIM, 0.13 LPIPS). These results confirm that each component contributes to stability and accuracy, with their combination being crucial for achieving high-quality reconstructions.

4 Conclusion

In this work, we addressed the challenge of dynamic scene reconstruction under fast motions with large inter-frame displacements. Existing 4DGS methods, including deformable and explicit approaches, often fail to maintain Gaussian attributes in such regimes, resulting in blurred or vanished objects. We introduced SPIN-4DGS, a new framework that learns Gaussian attributes directly from explicit spatiotemporal positions via a feed-forward network, rather than relying on temporal displacement modeling. By decoupling Gaussian position estimation from Gaussian attribute learning, SPIN-4DGS offers a simple yet effective design principle for representing dynamic scenes under fast motions. SPIN-4DGS achieved high-quality reconstructions of fast-moving objects under large inter-frame displacements, as demonstrated through extensive experiments across dynamic sports scenes in the CMU Panoptic Sports dataset. Overall, we believe SPIN-4DGS advances 4D Gaussian Splatting toward practical deployment in challenging real-world scenarios.

REFERENCES

- Abien Fred Agarap. Deep learning using rectified linear units (relu). *arXiv preprint* arXiv:1803.08375, 2018.
- Jeongmin Bae, Seoha Kim, Youngsik Yun, Hahyun Lee, Gun Bang, and Youngjung Uh. Pergaussian embedding-based deformation for deformable 3d gaussian splatting. In *European Conference on Computer Vision*, pp. 321–335. Springer, 2024.
- Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields, 2021.
- Ang Cao and Justin Johnson. Hexplane: A fast representation for dynamic scenes. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 130–141, 2023.
- Jie Chen, Zhangchi Hu, Peixi Wu, Huyue Zhu, Hebei Li, and Xiaoyan Sun. Dash: 4d hash encoding with self-supervised decomposition for real-time dynamic scene rendering. arXiv preprint arXiv:2507.19141, 2025.
- Yuanxing Duan, Fangyin Wei, Qiyu Dai, Yuhang He, Wenzheng Chen, and Baoquan Chen. 4d-rotor gaussian splatting: towards efficient novel view synthesis for dynamic scenes. In *ACM SIGGRAPH 2024 Conference Papers*, pp. 1–11, 2024.
- Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12479–12488, 2023.
- Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint* arXiv:1606.08415, 2016.
- Saqib Javed, Ahmad Jarrar Khan, Corentin Dumery, Chen Zhao, and Mathieu Salzmann. Temporally compressed 3d gaussian splatting for dynamic scenes. *arXiv preprint arXiv:2412.05700*, 2024.
- Hanbyul Joo, Hao Liu, Lei Tan, Lin Gui, Bart Nabbe, Iain Matthews, Takeo Kanade, Shohei Nobuhara, and Yaser Sheikh. Panoptic studio: A massively multiview system for social motion capture. In *Proceedings of the IEEE international conference on computer vision*, pp. 3334–3342, 2015.
- Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. ACM Transactions on Graphics, 42(4), July 2023.
- Diederik Kinga, Jimmy Ba Adam, et al. A method for stochastic optimization. In *International conference on learning representations (ICLR)*, volume 5. California;, 2015.
- Sangwoon Kwak, Joonsoo Kim, Jun Young Jeong, Won-Sik Cheong, Jihyong Oh, and Munchurl Kim. Modec-gs: Global-to-local motion decomposition and temporal interval adjustment for compact dynamic 3d gaussian splatting. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 11338–11348, 2025.
- Junoh Lee, Chang Yeon Won, Hyunjun Jung, Inhwan Bae, and Hae-Gon Jeon. Fully explicit dynamic gaussian splatting. *Advances in Neural Information Processing Systems*, 37:5384–5409, 2024.
- Tianye Li, Mira Slavcheva, Michael Zollhoefer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, Richard Newcombe, et al. Neural 3d video synthesis from multi-view video. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 5521–5531, 2022.
- Youtian Lin, Zuozhuo Dai, Siyu Zhu, and Yao Yao. Gaussian-flow: 4d reconstruction with dynamic 3d gaussian particle. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21136–21145, 2024.
- Zhicheng Lu, Xiang Guo, Le Hui, Tianrui Chen, Min Yang, Xiao Tang, Feng Zhu, and Yuchao Dai. 3d geometry-aware deformable gaussian splatting for dynamic view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8900–8910, 2024.

- Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. In 2024 International Conference on 3D Vision (3DV), pp. 800–809. IEEE, 2024.
- Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics*, 41(4):1–15, July 2022. ISSN 1557-7368. doi: 10.1145/3528223.3530127.
- Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *ICCV*, 2021.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019.
- Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings* of the IEEE conference on computer vision and pattern recognition, pp. 4104–4113, 2016.
- Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Xinggang Wang. 4d gaussian splatting for real-time dynamic scene rendering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 20310–20320, 2024.
- Jiawei Xu, Zexin Fan, Jian Yang, and Jin Xie. Grid4d: 4d decomposed hash encoding for high-fidelity dynamic gaussian splatting. *Advances in Neural Information Processing Systems*, 37: 123787–123811, 2024.
- Zeyu Yang, Hongye Yang, Zijie Pan, and Li Zhang. Real-time photorealistic dynamic scene representation and rendering with 4d gaussian splatting. *arXiv preprint arXiv:2310.10642*, 2023.
- Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin. Deformable 3d gaussians for high-fidelity monocular dynamic scene reconstruction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 20331–20341, 2024.
- Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 586–595, 2018.
- Ruijie Zhu, Yanzhe Liang, Hanzhi Chang, Jiacheng Deng, Jiahao Lu, Wenfei Yang, Tianzhu Zhang, and Yongdong Zhang. Motiongs: Exploring explicit motion guidance for deformable 3d gaussian splatting. *Advances in Neural Information Processing Systems*, 37:101790–101817, 2024.
- Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross. Ewa volume splatting. In *Proceedings Visualization*, 2001. VIS'01., pp. 29–538. IEEE, 2001.

A RELATED WORK

Learning temporal dynamics for 3D Gaussian. Recent advances have shown significant progress in extending 3D Gaussian representations into the 4D domain by learning temporal dynamics. Early attempts (Luiten et al., 2024; Javed et al., 2024) incrementally propagated 3D Gaussians from the first frame, demonstrating the potential of 4D Gaussian splatting, but suffering from sequential inefficiency and overfitting due to numerous per-frame iterations. However, their reliance on external supervision, such as segmentation masks, substantially increases training costs and computational overhead. Meanwhile, inspired by NeRF approaches (Cao & Johnson, 2023; Fridovich-Keil et al., 2023; Park et al., 2021), deformable frameworks (Wu et al., 2024; Yang et al., 2024; Lu et al., 2024; Zhu et al., 2024; Lu et

Explicit 4D Gaussian Parameterizations. Another promising direction is to focus on directly parameterizing Gaussians in 4D from scratch (Yang et al., 2023; Lee et al., 2024; Duan et al., 2024), rather than optimizing each frame separately in the 3D domain. Such explicit 4D modeling unifies space and time into a continuous field and encodes dynamics as 4D Gaussian splats, which can be temporally sliced for rendering, in contrast to deformable-based approaches. For instance, RotorGS (Duan et al., 2024) performs temporal slicing of 4D Gaussians at each timestamp to obtain dynamic 3D Gaussians, which are then projected to the image plane. While explicit 4D parameterizations achieve higher frame rates (FPS) and improved rendering quality, they require longer training times and larger storage requirements. Moreover, under fast motion with large displacements, we observed that corresponding Gaussian attributes gradually blur due to cross-frame interference. A single 4D Gaussian must simultaneously explain all timestamps; however, rasterization is performed frame by frame, so optimization for one frame inevitably affects the others. This inherent conflict in the rasterization process makes it challenging to maintain temporal consistency in dynamic regions. In contrast, our method avoids such cross-frame degradation by directly decoding Gaussians at explicit spatiotemporal positions in a feed-forward manner. This design explicitly separates largedisplacement Gaussians at each timestamp, which prevents cross-frame degradation and leads to both stable training and consistent rendering quality.

B More Qualitative Results

Qualitative results on spatiotemporal position refinement. Figure 6 visualizes the ablation study on the position refinement in Table 2. With only 0.5K points, SPIN-4DGS already captures stable object structures without collapse, while increasing the refinement iterations to 1K and 2K further improves fine details such as ball edges and racket nets. These results show that SPIN-4DGS not only enhances reconstruction quality under minimal refinement but also scales effectively with additional iterations to capture finer details.

Qualitative results on compatibility with existing 4DGS baselines. Figure 7 presents qualitative results for the compatibility study in Table 4. Across scenes, our method reconstructs stable backgrounds and more temporally consistent details, while the original baselines often blur or lose fast-moving objects. These results confirm that the proposed scheme integrates effectively with existing 4DGS approaches and potentially improves their reconstruction fidelity under large interframe displacements.

C LLM POLICY

LLMs were used only for editing (grammar and typographical errors) and did not contribute to idea generation, analysis, experimental design, or interpretation.

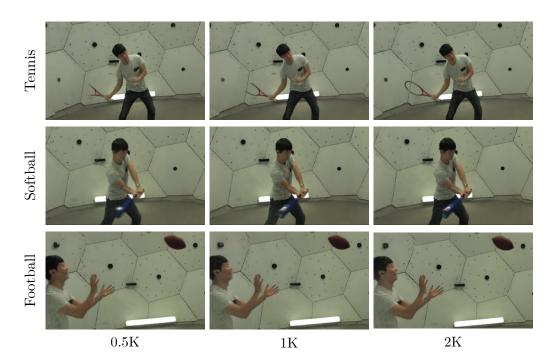


Figure 6: **Visualizations on spatiotemporal position refinement.** We visualize the effect of varying refinement iterations (0.5K, 1K, 2K) corresponding to Table 2b. Even with only 0.5K iterations, SPIN-4DGS retains consistent object structures without collapse, while progressively increasing the number of iterations recovers fine details such as ball edges and racket nets.

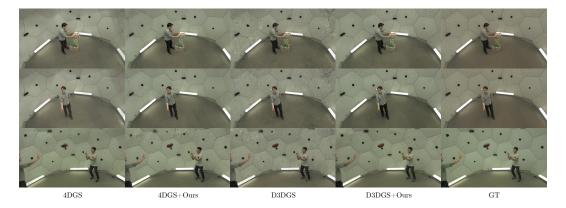


Figure 7: **Visualizations on compatibility with existing 4DGS baselines.** We visualize the results of compatibility ablation study in Table 4, where pre-trained positions from 4DGS and D3DGS are reused while attributes are retrained with SPIN-4DGS. Our method consistently produces clear and more stable reconstructions than the original baselines.