# Dynamic-eDiTor: Training-Free Text-Driven 4D Scene Editing with Multimodal Diffusion Transformer

Dong In Lee[1,2*‡]    Hyungjun Doh[1*]    Seunggeun Chi[1]    Runlin Duan[1]

Sangpil Kim[2†]    Karthik Ramani[1†]
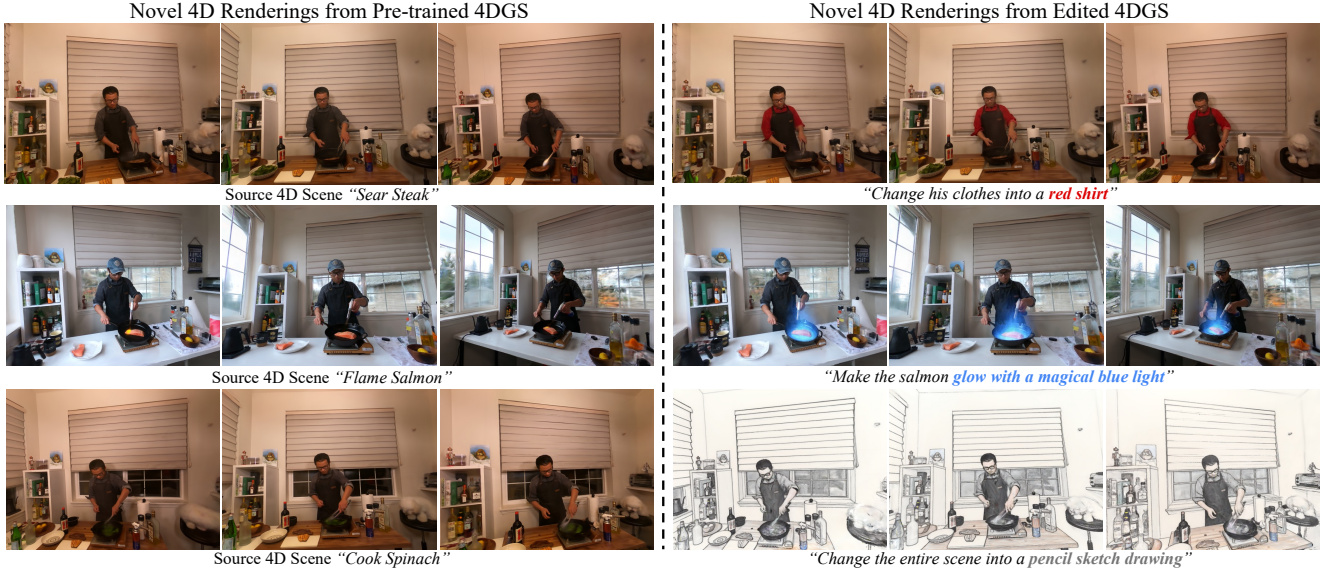
[1]Purdue University    [2]Korea University

Figure 1. We propose **Dynamic-eDiTor** enables flexible and high-quality editing of pre-trained 4D Gaussian Splatting [55] models leveraging Multimodal Diffusion Transformer [11, 54] guided solely by text instructions. Through its design focused on both multi-view and temporal consistency, our approach demonstrates robust performance, producing realistic and fine-grained 4D scene manipulation.

## Abstract

*Recent progress in 4D representations, such as Dynamic NeRF and 4D Gaussian Splatting (4DGS), has enabled dynamic 4D scene reconstruction. However, text-driven 4D scene editing remains under-explored due to the challenge of ensuring both multi-view and temporal consistency across space and time during editing. Existing studies rely on 2D diffusion models that edit frames independently, often causing motion distortion, geometric drift, and incomplete editing. We introduce Dynamic-eDiTor, a training-free text-driven 4D editing framework leveraging Multimodal Diffusion Transformer (MM-DiT) and 4DGS. This mechanism consists of Spatio-Temporal Sub-Grid Attention (STGA) for locally consistent cross-view and temporal fusion, and Context Token Propagation (CTP) for global propagation via token inheritance and optical-flow-guided token replacement. Together, these components allow Dynamic-eDiTor to perform seamless, globally con-sistent multi-view video without additional training and directly optimize pre-trained source 4DGS. Extensive experiments on multi-view video dataset DyNeRF demonstrate that our method achieves superior editing fidelity and both multi-view and temporal consistency prior approaches. Project page for results and code:: https://di-lee.github.io/dynamic-eDiTor/*

## 1. Introduction

Recent advances in 3D representations, such as Neural Radiance Field (NeRF) [36] and 3D Gaussian Splatting (3DGS) [23], have achieved significant progress in photo-realistic 3D reconstruction of real-world scenes. More recently, 4D representations such as Dynamic NeRF [41] and 4D Gaussian Splatting (4DGS) [55] extend 3D representations into the time domain, enabling spatio-temporally coherent reconstruction. However, text-driven 4D scene editing remains under-explored, primarily due to the difficulty of maintaining both multi-view and temporal consistency across space and time during editing.

In this work, we focus on the multi-view video setting of 4DGS, which provides richer viewpoint coverage but fur-

---

*Co-first authors.

†Co-corresponding authors.

‡Work done at Purdue University as a visiting scholar.

ther amplifies the difficulty of achieving both multi-view and temporal consistency during editing. While 3D editing primarily focuses on multi-view consistency, 4D editing introduces the further challenge of ensuring both multi-view and temporal consistency across viewpoints and time.

Recent studies [18, 27, 37] have attempted 4D scene editing by combining 2D diffusion models with 4D representations [46, 55]. However, these methods typically perform frame-wise editing or require per-scene finetuning of the 2D diffusion model [4], lacking a unified mechanism to jointly process information across views and time. Consequently, they struggle with non-rigid content manipulation and are often limited to style-oriented edits, leading to motion distortions, geometric drift, and incomplete editing results.

To address these limitations, we propose **Dynamic-eDiTor**, a novel training-free, text-driven 4D editing framework that leverages Multimodal Diffusion Transformer (MM-DiT) [11, 54] and 4DGS. Our goal is to maintain globally coherent motion and geometry while enabling flexible, semantically grounded edits. To this end, we propose Grid-based Spatio-Temporal Propagation, which represents the entire multi-view video as a unified camera-time grid and jointly aggregates spatial and temporal information and propagates the fused features across throughout the grid.

As its foundation, we introduce Spatio-Temporal Sub-Grid Attention (STGA), which extends MM-DiT's dual-stream self-attention to operate on localized spatio-temporal sub-grids. By jointly attending to adjacent viewpoints and neighboring time steps, STGA enables coherent local feature fusion without additional training. We additionally identify a vital layer range in MM-DiT where incorporating STGA yields the strongest improvements in multi-view and temporal consistency.

To ensure that the fused information is globally propagated throughout the multi-view video, we further introduce Context Token Propagation (CTP), an explicit propagation mechanism that transfers fused tokens along a structured traversal path over the entire multi-view video. Tokens in overlapping regions are fully inherited, while non-overlapping temporal regions are updated through flow-guided token warping using optical flow [47]. This unified propagation strategy ensures coherent feature flow across views and time, reinforcing multi-view and temporal consistency and enabling stable, high-fidelity 4D optimization.

Finally, the edited frames are directly used to optimize the pre-trained 4DGS without the Iterative Dataset Update (IDU) [18, 37], resulting in globally consistent 4D content that faithfully reflects the desired edits.

We validate Dynamic-eDiTor on the multi-view video dataset DyNeRF [33], achieving superior editing fidelity, temporal smoothness, and robustness compared to state-of-the-art methods. Our key contributions are as follows:

- We present Dynamic-eDiTor, a novel training-free, text-

driven 4D editing framework that leverages MM-DiT [11, 54] and 4DGS [55] to enable spatially and temporally consistent dynamic 4D scene editing.
- We propose Spatio-Temporal Sub-Grid Attention (STGA), which jointly attends across adjacent viewpoints and neighboring time steps to integrate spatial and temporal features on a vital layer range in MM-DiT.
- We introduce Context Token Propagation (CTP), an explicit propagation mechanism that distributes fused spatio-temporal information across the entire multi-view video by inheriting tokens in overlapping regions and replacing non-overlapping regions via flow-based warping.
- Through extensive qualitative and quantitative experiments, we demonstrate that Dynamic-eDiTor significantly outperforms existing methods in 4D editing fidelity, spatio-temporal stability, and robustness.

## 2. Related Work

### 2.1. 2D Editing

2D diffusion models have demonstrated remarkable generalization and controllability for text-guided image editing. U-Net-based [43] models such as Prompt-to-Prompt [19], SDEdit [35], and InstructPix2Pix [4] enable text-guided image manipulation via controlled denoising. More recently, Diffusion Transformers (DiT) [39] replace the U-Net backbone with a Transformer architecture [10, 48], offering improved scalability and visual coherence. This approach has evolved into multimodal variants such as Multimodal Diffusion Transformer (MM-DiT) [11]. MM-DiT employs a dual-stream architecture, processing text and image tokens in parallel streams fused via joint attention. Building on this trend, MM-DiT-based editors such as FLUX [28, 29], HiDream [5], and Qwen-Image-Edit [54] achieve precise instruction-driven image editing. Leveraging the MM-DiT, our Dynamic-eDiTor extends a MM-DiT-based image editor to maintain consistency across time and viewpoints within a unified 4D framework.

### 2.2. 3D Scene Editing

Neural Radiance Fields (NeRF) [36] and 3D Gaussian Splatting (3DGS) [23] have enabled high-fidelity 3D reconstruction and inspired extensive research on 3D scene editing. Instruct-NeRF2NeRF [17] introduces the Iterative Dataset Update (IDU) that edits the rendered image using a 2D diffusion model [4] while optimizing the underlying 3D model, NeRF. GaussianEditor [50] adopts IDU on 3DGS and explicitly controls 3D Gaussians. Recently, EditSplat [31] achieves high-fidelity 3D edits by integrating multi-view information into the diffusion process and explicitly pruning 3DGS. Despite these advances, existing 3D editing methods [6, 9, 22, 51, 56, 61, 62] primarily focus on enforcing multi-view consistency within static scenes, and cannot handle temporal dynamics across frames.
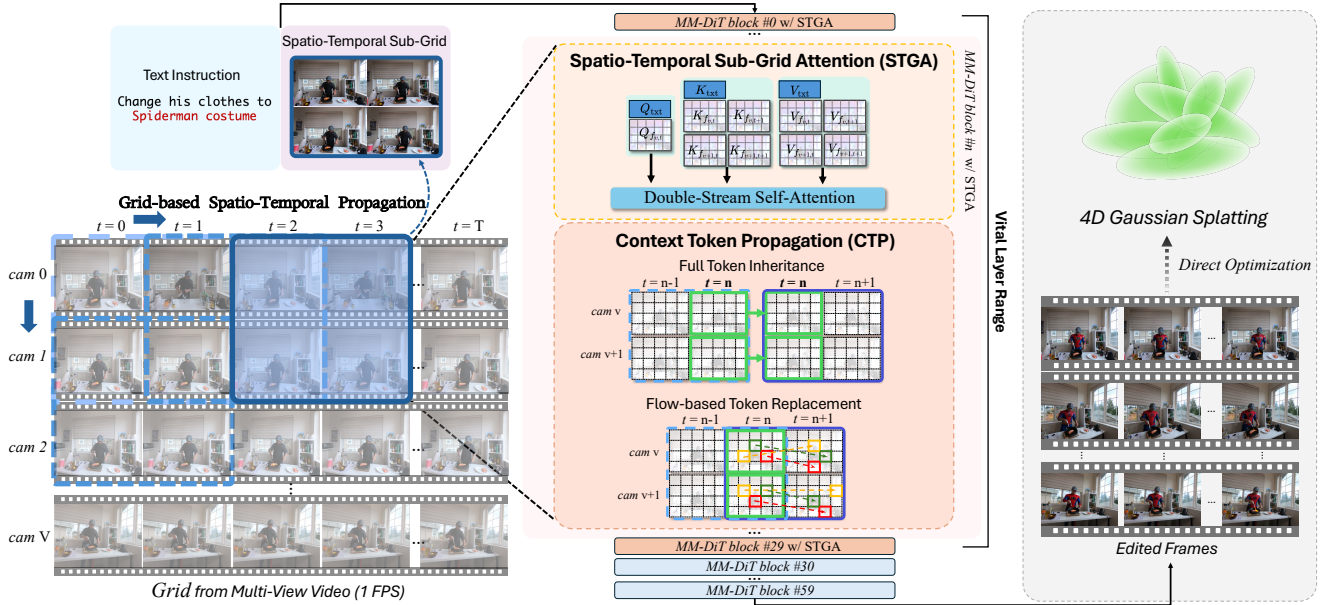
Figure 2. **Dynamic-eDiTor Overview.** We represent the multi-view video as a unified camera–time grid. Dynamic-eDiTor combines Spatio-Temporal Sub-Grid Attention (STGA), which performs locally coherent cross-view and temporal fusion within each sub-grid, with Context Token Propagation (CTP), which globally propagates the aggregated features across the grid via Full Token Inheritance and Flow-guided Token Replacement for robust spatio-temporal consistency enforcement. Together, these modules enable seamless, globally consistent multi-view video editing without additional training, while directly optimizing the pre-trained 4DGS.

## 2.3. 4D Scene Editing

Extending 3D scene editing to 4D representations introduces the additional challenge of maintaining temporal consistency while preserving multi-view coherence. Instruct 4D-to-4D [37] first applied diffusion-based 2D editing to sequentially rendered frames while optimizing the underlying NeRF-based 4D model, while 4D-Editor [21] incorporates spatial segmentation and motion-aware propagation for object-level editing. Control4D [45] enables 4D portrait editing by distilling the knowledge from a 2D diffusion into a 4D GAN [16] generator. CTRL-D [18] finetunes InstructPix2Pix [4] with prior-preservation loss [44] per scene for consistent 2D edits and iteratively optimizes 4DGS. Instruct4DGS [27] edits canonical Gaussians first and employs score-distillation-based [40] refinement for temporal smoothness. While these methods demonstrate notable progress, they still edit frames independently without simultaneously processing information across views and time, often causing motion or geometric distortion. In contrast, our Dynamic-eDiTor achieves consistent 4D editing by jointly editing cross-view and temporal frames, and by propagating these context tokens to the entire multi-view video.

## 3. Preliminary

### 3.1. 4D Scene Representation

3D Gaussian Splatting (3DGS) [23] represents a scene as a set of anisotropic Gaussian primitives $\mathcal{G} = \{(\mu_i, \Sigma_i, c_i, \alpha_i)\}_{i=1}^{N}$, each defined by its position, covariance, color, and opacity. Rendering is performed via differentiable alpha compositing. To model dynamic scenes,

4D Gaussian Splatting (4DGS) [55] extends this representation with a deformation field that maps canonical Gaussians to their deformed states over time. The field predicts offsets for position, rotation, and scale using MLPs $\phi_x$, $\phi_r$, and $\phi_s$: $\Delta x = \phi_x(z)$, $\Delta r = \phi_r(z)$, and $\Delta s = \phi_s(z)$, where $z$ is a temporal feature encoding the dynamic state of the scene. The final deformed Gaussian parameters are obtained as:

$$(x', r', s') = (x + \Delta x, \, r + \Delta r, \, s + \Delta s), \qquad (1)$$

yielding the time-varying Gaussian set $\mathcal{G}'$.

## 4. Method

We propose a novel training-free 4D editing framework, **Dynamic-eDiTor**, carefully designed to achieve spatially and temporally consistent 4D scene editing leveraging Multimodal Diffusion Transformer (MM-DiT) [11, 54]. Initially, our approach edits source multi-view video frames at 1 FPS corresponding to a pre-trained 4D Gaussian Splatting (4DGS) [55], ensuring both multi-view and temporal consistency. The edited frames are then used to directly optimize the underlying pre-trained 4DGS representation.

### 4.1. Grid-based Spatio-Temporal Propagation

To ensure both multi-view and temporal consistency in multi-view video editing, we introduce Grid-based Spatio-Temporal Propagation, which enables coherent feature flow across the entire scene through two components: (1) Spatio-Temporal Sub-Grid Attention (STGA) for local fusion across adjacent views and neighboring time steps and (2) Context Token Propagation (CTP) for globally propagating

the fused information through the entire multi-view video.

We begin by representing all multi-view video frames as a unified camera–time grid:

$$Grid = \{f_{v,t} \mid v \in [0, \ldots, V],\ t \in [0, \ldots, T]\}, \quad (2)$$

where $f_{v,t}$ is the frame captured by viewpoint $v$ at time index $t$. The Grid's rows correspond to different viewpoints and columns represent temporally aligned frames.

To enable localized spatio-temporal fusion, we partition the $Grid$ into overlapping $2 \times 2$ sub-grid $\mathcal{S}_{v,t}$ at position $(v, t)$ on $Grid$ defined as:

$$\mathcal{S}_{v,t} = \{f_{v,t},\ f_{v+1,t},\ f_{v,t+1},\ f_{v+1,t+1}\}, \quad (3)$$

each covering adjacent views and neighboring time steps. These sub-grids serve as the atomic processing units for STGA, allowing each local region to aggregate and share information across both the view and temporal axes.

To propagate information across the entire $V \times T$ $Grid$, we process the sub-grids sequentially using an asymmetric sliding pattern. We first sweep vertically along the spatial axis at $t=0$ to establish multi-view alignment, and then slide horizontally along the temporal axis to propagate consistency over time. The induced overlaps between neighboring sub-grids provide the structural linkage for STGA and CTP to effectively distribute information, enforcing globally coherent spatio-temporal editing.

#### 4.1.1. Spatio-Temporal Sub-Grid Attention (STGA)

Grid-based Spatio-Temporal Propagation's foundation is the fusion of information within local neighborhoods. We propose Spatio-Temporal Sub-Grid Attention (STGA), which extends the dual-stream self-attention mechanism in MM-DiT architecture [11] to jointly attend adjacent views and temporally neighboring frames.

Instead of processing each frame's feature independently as in standard MM-DiT, STGA operates on a local sub-grid $\mathcal{S}_{v,t}$, enabling each frame to aggregate features from its cross-view and temporal neighbors. Each sub-grid contains four frames, and every frame $f_i \in \mathcal{S}_{v,t}$ is sequentially processed as a query in turn. Following MM-DiT's dual-stream attention design, the attention calculation involves two parts–the text stream and image-feature stream. For given frame $f_i$, we use its image query $Q_{f_i}$. We then extend the image-feature stream by concatenating all frame-level features within the sub-grid $\mathcal{S}_{v,t}$ to construct joint key and value sets $K_{\mathcal{S}_{v,t}}$ and $V_{\mathcal{S}_{v,t}}$:

$$K_{\mathcal{S}_{v,t}} = [K_{f_{v,t}}, K_{f_{v+1,t}}, K_{f_{v,t+1}}, K_{f_{v+1,t+1}}],$$
$$V_{\mathcal{S}_{v,t}} = [V_{f_{v,t}}, V_{f_{v+1,t}}, V_{f_{v,t+1}}, V_{f_{v+1,t+1}}]. \quad (4)$$

The STGA for each frame $f_i$ integrates the text stream $Q_{txt}, K_{txt}, V_{txt}$ with our modified, spatio-temporal image stream. We then apply Rotary Position Embeddings (RoPE) to the image queries $Q_{f_i}, K_{\mathcal{S}_{v,t}}, V_{\mathcal{S}_{v,t}}$ to provide positional
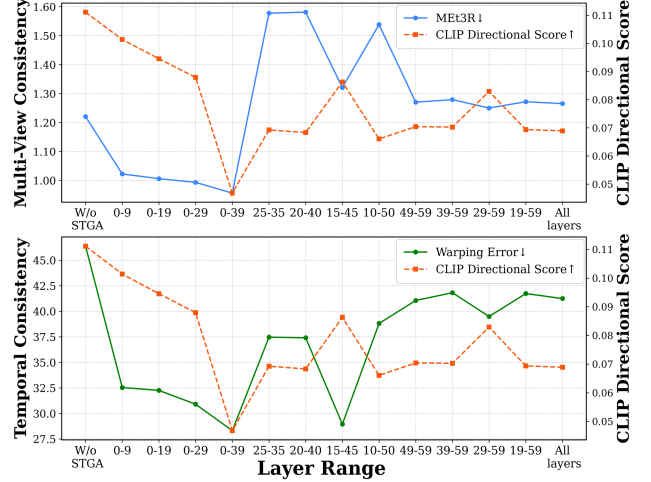


Figure 3. **Vital Layer Range Analysis.** We analyze the impact of applying Spatio-Temporal Sub-Grid Attention (STGA) across different layer ranges in MM-DiT [11, 54] during the multi-view video editing process. Performance is evaluated by temporal consistency (*Warping Error* [30]), multi-view consistency (*MEt3R* [2]), and editing fidelity (*CLIP Text-Image Directional Similarity [42]*). Applying STGA to the early ∼30 layers provides the best trade-off between consistency and editing fidelity.

information before the softmax operation:

$$\mathrm{STGA}(\mathcal{S}_{v,t}) = \mathrm{softmax}\Big( [Q_{txt}, \mathrm{RoPE}(Q_{f_{v,t}})] \cdot$$
$$[K_{txt}, \mathrm{RoPE}(K_{\mathcal{S}_{v,t}})]^\top / \sqrt{d_k} \Big) \cdot [V_{txt}, V_{\mathcal{S}_{v,t}}], \quad (5)$$

where $d_k$ denotes the dimensionality of the key vectors.

STGA encourages cross-view and temporal coherence, forming the foundation for globally consistent multi-view video editing. Unlike previous temporal-only extensions of self-attention [15, 57], our STGA enables each frame query to attend both spatially adjacent views and temporally neighboring frames within its saptio-temporal patch. While STGA operates locally within each sub-grid, the overlapping sliding pattern naturally leads to implicit propagation of fused information across adjacent sub-grids.

**Vital Layer Range Selection.** As illustrated in Fig. 7, applying STGA to all self-attention layers of MM-DiT [11, 54] leads to visual artifacts, as the STGA tends to over-attend within the local spatio-temporal patch, resulting in texture repetition and view-dependent inconsistencies [12]. While prior studies [3, 14, 24, 53] analyze layer importance in DiT-based models [39] by ranking individual layers. Instead, we investigate applying STGA across continuous layer ranges to capture this cumulative effect. We empirically observe that there is a vital layer range to apply STGA for effectively enforcing multi-view and temporal consistency while alleviating editing-quality degradation. As shown in Fig. 3, applying STGA to the first 30 layers achieves the best trade-off between consistency and fidelity, providing significant improvements in both multi-view and

4

Figure 4. **Qualitative Comparison.** Dynamic-eDiTor enables more robust non-rigid content manipulation and achieves more complete edits of the 4D scene. The top-row displays the original rendered frames, while the following rows show the edited 4DGS renderings produced by each baseline. Our method (bottom-row) outperforms all baselines in both text alignment and overall editing fidelity, while maintaining strong temporal and spatial consistency.

temporal coherence while maintaining editing quality.

### 4.1.2. Context Token Propagation (CTP)

While STGA achieves local cross-view and temporal coherence by jointly attending adjacent views and neighboring frames within the sub-grid, Context-Aware Propagation ensures this coherence is globally distributed across entire $V \times T$ $Grid$. As the sub-grid slides along the defined traversal path, we introduce Context Token Propagation (CTP), which explicitly propagates context information. This ensures that the coherent feature representations that contain spatial and temporal information computed by STGA in the previous sub-grid $S_{prev}$ are injected into the current sub-grid $S_{curr}$, ensuring coherent feature flow, enforcing global consistency and preventing information loss.

In this process, the token representation is defined as $\phi(\mathcal{S}_{v,t}) = \text{STGA}(\mathcal{S}_{v,t})$. We employ two Context Token Propagation strategies: Full Token Inheritance and Flow-guided Token Replacement. Full Token Inheritance is applied when the current sub-grid $S_{curr}$ shares frames with the previous sub-grid $S_{prev}$ along the temporal axis ($t = 1 \rightarrow T - 1$) or the spatial axis ($v = 1 \rightarrow V - 1$). We directly replace the entire current token $\phi(\mathcal{S}_{curr})$ in these overlapped frames with previous token $\phi(\mathcal{S}_{prev})$.

For a sub-grid along the temporal axis, a defined traversal path yields non-overlapped regions in the rightmost col-umn of the sub-grid. Thus, we apply Flow-guided Token Replacement to these regions, in which the tokens are replaced with tokens warped from the corresponding right-most column regions of the previous sub-grid. To ensure temporal alignment during warping, we estimate forward and backward optical flow between frames $f_t$ and $f_{t-1}$ using RAFT [47], and downsample the flow fields to match the token resolution. For each spatial location $(x, y)$, we use the downsampled forward flow $\mathbf{F}_{t \rightarrow t-1}(x, y)$ to backward-warp the tokens from the previous patch:

$$\hat{\phi}_{\text{r}}(\mathcal{S}_{v,t}) = \text{Warp}\big(\mathbf{F}_{t \rightarrow t-1}(x, y), \phi_{\text{r}}(\mathcal{S}_{v,t-1})\big). \quad (6)$$

where $\hat{\phi}_{\text{r}}(\mathcal{S}_{v,t})$ denotes the warped tokens in the rightmost column of the patch. During the replacement, we compute a validity mask $\text{M}(\text{x}, \text{y})$ via a forward–backward consistency check, inspired by [34, 58], to ensure precise replacement. With the mask $\text{M}$, tokens in valid regions are replaced by the warped tokens, while those in invalid regions retain the current frame's tokens:

$$\phi_{\text{r}}(\mathcal{S}_{v,t}) = \text{M} \odot \hat{\phi}_{\text{r}}(\mathcal{S}_{v,t}) + (1 - \text{M}) \odot \phi_{\text{r}}(\mathcal{S}_{v,t}), \quad (7)$$

where $\odot$ denotes element-wise multiplication.

This unified propagation mechanism enables efficient and robust propagation across both spatial and temporal dimensions. STGA provides locally coherent spatial-

| | Editing Fidelity | | User Study | | | | | | Reconstruction Fidelity | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Method | CLIP$_{dir}$ ↑ | CLIP$_{sim}$ ↑ | Overall Quality (%) | Motion Consist. (%) | Temporal Consist. (%) | Multi-view Consist. (%) | Prompt Align. (%) | Identity Preserv. (%) | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
| Instruct4D-to-4D [37] | 0.1077 | 0.6308 | 27.57 | 27.72 | 28.00 | 27.19 | 22.14 | 27.48 | 21.86 | 0.6978 | 0.2145 |
| Instruct-4DGS [27] | 0.1501 | 0.6342 | 10.48 | 10.52 | 11.05 | 10.48 | 9.29 | 11.05 | 20.62 | 0.6252 | 0.2869 |
| CTRL-D [18] | 0.1498 | 0.6141 | 13.00 | 14.62 | 15.57 | 14.14 | 11.71 | 13.95 | **31.06** | **0.8498** | **0.0970** |
| **Ours** | **0.1849** | **0.6397** | **48.95** | **47.14** | **45.38** | **48.19** | **56.86** | **47.52** | 29.25 | 0.8064 | 0.1006 |

Table 1. **Quantitative Comparison.** The evaluation spans three aspects: editing fidelity, user preference, and reconstruction fidelity. CLIP-based metrics [42] show that Dynamic-eDiTor achieves strong alignment with the editing prompts across 4D scenes, and user studies indicate a clear preference for our results over the baselines in terms of semantic alignment, perceptual realism, and coherent motion. Although reconstruction metrics (PSNR, SSIM [52], LPIPS [60]) are slightly lower, they remain competitive and do not detract from the method's overall superiority in semantic accuracy and perceptual edit quality.

temporal feature aggregation, while CTP propagates this coherence globally across the entire $Grid$. As a whole, they ensure consistent motion and geometry in multi-view videos, significantly improving stability in 4D editing.

## 4.2. Direct 4D Scene Optimization

Our approach produces multi-view and temporally consistent edited video frames, which can be directly used to optimize the pre-trained 4D representation. In contrast to prior works [18, 37] that rely on the Iterative Dataset Update (IDU), we directly optimize the 4D Gaussian model $\mathcal{G}'_{edit}$ using all edited frames $f^{edit}_{v,t}$ across the entire $Grid$. The optimization objective is defined as:

$$\mathcal{G}'_{edit} = \arg\min_{\mathcal{G}} \sum_{v,t \in V,T} \left\| \hat{f}_{v,t} - f^{edit}_{v,t} \right\| + \mathcal{L}_{tv}, \quad (8)$$

where both loss terms follow the 4DGS [55].

## 5. Experiment

### 5.1. Experimental Setup

**Dataset.** We evaluate our method on the real-world multi-view video dataset DyNeRF [33], which contains six dynamic scenes with 16-21 camera views capturing 10-second videos at 30 FPS. To evaluate editing consistency under sparse temporal sampling, we uniformly sample frames at 1FPS (160-210 frames per scene, compared to 4,800-6,300 frames at 30FPS). All experiments are conducted using 14 prompts covering all six scenes in the dataset.

**Baselines.** We compare our method against state-of-the-art 4D scene editing approaches, including Instruct 4D-to-4D [37], CTRL-D [18], and Instruct 4DGS [27]. Since our task focuses on text-based scene editing, we reproduce all baseline results using text prompt input only.

### 5.2. Implementation Details

Our method leverages the MM-DiT-based [11] image editing model, Qwen-Image-Edit [54] and 4D Gaussian Splatting [55]. For evaluation, we use all camera views in the dataset and hold out the final frame of each view as the test set. We evaluate our rendered results on this test set. The
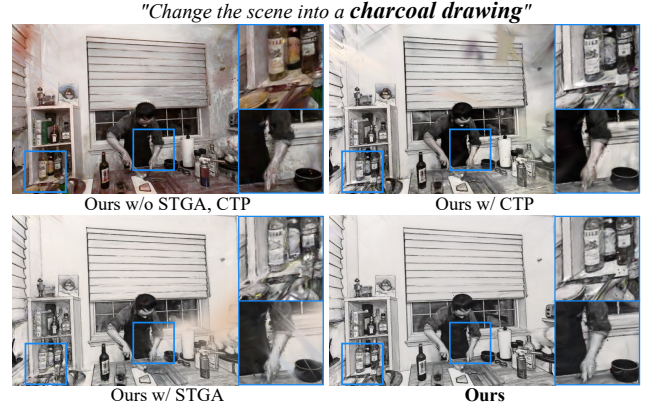
*"Change the scene into a **charcoal drawing**"*



Figure 5. **Qualitative Ablation Results.** The model lacking both components (top-left) suffers from severe artifacts and geometric drift. Adding only STGA or only CTP progressively improves the result, but still leaves residual motion blur and geometric drift. Our full method (bottom-right) successfully ensuring the spatio-temporal consistency to produce a stable and complete edit.

full 4D editing process takes approximately 51 minutes for the "coffee martini" scene on a single NVIDIA H100 GPU.

### 5.3. Qualitative Results

Fig. 1 illustrates Dynamic-eDiTor's ability to perform multi-view temporal scene editing. Our model effectively edits diverse scenes and local objects while maintaining strong temporal and spatial consistency. Dynamic-eDiTor is able to perform non-rigid appearance editing, semantic local editing, and consistent stylization while still preserving motion consistency across viewpoints and over time.

In Fig. 4, we compare Dynamic-eDiTor with recent 4D scene editing baselines [18, 27, 37]. We observe that most baseline methods fail to handle non-rigid content manipulation and are often limited to style-oriented edits. This core limitation leads to significant artifacts, such as motion distortions, geometric drift, and incomplete editing results. These failures are evident across the examples. When editing the scene into a "fire emergency", all baselines fail to generate plausible emergency-related elements, revealing weak text–scene alignment and incomplete editing. In the second column, Instruct-4DGS struggles with non-rigid

| STGA | CTP | 2D Consistency | | Reconstruction Fidelity | | | Editing Fidelity | |
|---|---|---|---|---|---|---|---|---|
| | | Warp-Err $_{10^{-3}}$ ↓ | MEt3R $_{10^{-1}}$ ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ | CLIP$_{dir}$ ↑ | CLIP$_{sim}$ ↑ |
| - | - | 56.98 | 1.0721 | 26.14 | 0.7445 | 0.1408 | <u>0.1930</u> | 0.5414 |
| ✓ | - | 38.64 | <u>0.9277</u> | 28.08 | 0.7875 | <u>0.1122</u> | 0.1872 | <u>0.6407</u> |
| - | ✓ | <u>29.44</u> | 1.0695 | <u>28.74</u> | <u>0.8013</u> | 0.1165 | **0.1944** | **0.6418** |
| ✓ | ✓ | **28.94** | **0.9074** | **29.25** | **0.8064** | **0.1006** | 0.1849 | 0.6397 |

Table 2. **Ablation Study.** Each component, Spatio-Temporal Sub-Grid Attention (STGA) and Context Token Propagation (CTP), helps preserve temporal and multi-view consistency, improving the 4D reconstruction quality. Our method prioritizes a globally stable 4D structure, yielding consistent temporal and spatial behavior and thus more robust reconstruction fidelity. Although CLIP-based metrics [42] show a slight drop due to the trade-off between semantic alignment and spatio-temporal coherence, our method still produces more stable and reliable 4D edits, avoiding the geometric and temporal artifacts seen in the ablated variants.

| CTP-Full | CTP-Flow | 2D Consistency | | Reconstruction Fidelity | | | Editing Fidelity | |
|---|---|---|---|---|---|---|---|---|
| | | Warp-Err $_{10^{-3}}$ ↓ | MEt3R $_{10^{-1}}$ ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ | CLIP$_{dir}$ ↑ | CLIP$_{sim}$ ↑ |
| - | - | 38.64 | 0.9277 | 28.08 | 0.7875 | 0.1122 | **0.1872** | **0.6407** |
| - | ✓ | <u>29.79</u> | 0.9205 | <u>28.97</u> | <u>0.7990</u> | <u>0.1034</u> | 0.1852 | <u>0.6402</u> |
| ✓ | - | 33.22 | <u>0.9094</u> | 28.19 | 0.7906 | 0.1089 | <u>0.1865</u> | 0.6400 |
| ✓ | ✓ | **28.94** | **0.9074** | **29.25** | **0.8064** | **0.1006** | 0.1849 | 0.6397 |

Table 3. **Ablation Study: Context Token Propagation (CTP).** This ablation study is conducted with STGA included to isolate the impact of CTP. Full Token Inheritance (CTP-Full) and Flow-Guided Token Replacement (CTP-Flow) play a critical role in reinforcing temporal and multi-view consistency, enabling more accurate reconstruction of the edited dynamic scene. Despite a slight trade-off in CLIP-based metrics [42], CTP substantially improves spatio-temporal coherence and overall 4D editing fidelity.

editing, causing clear motion distortions around the hand. Meanwhile, Instruct 4D-to-4D and CTRL-D introduce noticeable artifacts such as facial color shifts and blurring. In the third column, CTRL-D further demonstrates viewpoint inconsistencies and produces blurred edited regions, while other baselines result in incomplete edits. Instruct 4D-to-4D fails to modify the target scene, incorrectly altering surrounding objects. This indicates weak text alignment despite sharing the same diffusion backbone such as InstructPix2Pix [4] as other baselines. Overall, Dynamic-eDiTor outperforms all previous 4D scene editing methods by achieving superior editing completeness and effectively preserving temporal coherence and multi-view consistency, resulting in high-quality dynamic scene edits.

## 5.4. Quantitative Results

Tab. 1 presents a quantitative comparison with prior 4D scene editing methods, focusing on the 4D rendered image quality. Our evaluation is structured into two categories: text-prompt alignment and reconstruction fidelity.

To evaluate text-prompt alignment, we use CLIP-based [42] metrics. Specifically, the CLIP text-image directional similarity captures how changes in text captions correspond to changes between the source and rendered images in CLIP embedding space, while the CLIP text-image similarity directly measures alignment between the target text and rendered frames. Our method surpasses all baselines in these CLIP metrics, demonstrating that our rendered results are significantly better aligned with the user's intended edit.

To assess reconstruction fidelity, we report PSNR,

SSIM [52], and LPIPS [60]. These metrics are computed between the final rendered test frames and the corresponding 2D edited target frames, measuring how faithfully the 4D model reconstructs the target edits. Although Dynamic-eDiTor obtains slightly lower values than CTRL-D on these reconstruction metrics, this highlights that our method achieves a better balance by prioritizing faithful text alignment and reliable 4D scene editing. This trade-off is further supported by our vital layer analysis in Fig. 3.

Beyond reconstruction metrics, we evaluate perceptual quality through a user study with 150 participants on Amazon Mechanical Turk [1]. Participants were asked to compare our final 4D rendered videos against baseline methods [18, 27, 37]. As shown in Tab. 1, our method consistently outperforms the baselines in terms of overall visual quality, motion consistency, temporal and multi-view consistency, text-prompt alignment, and identity preservation.

## 5.5. Ablation Study

We conduct an ablation study on our Grid-based Spatio-Temporal Propagation and Context Token Propagation.
**2D Consistency.** We first assess each component's impact on 2D temporal and multi-view consistency, a key factor for high-quality 4D reconstruction. As shown in Fig. 6, STGA strengthens semantic alignment and view consistency, whereas CTP improves temporal coherence through information propagation. Collectively, they yield notable improvements in 2D spatial and temporal stability. The quantitative results in Tab. 6 support these findings. Our full method achieves superior spatio-temporal consistency,

*"Make the scene look like a volcano with glowing lava reflection"*

Cam 16 / t = 6   Cam 16 / t = 7   Cam 12 / t = 4   Cam 12 / t = 5

Figure 6. **Ablation Study: 2D Consistency.** Each component in our method strengthens temporal and multi-view consistency in 2D editing. STGA improves semantic alignment and preserves fine details across views, while CTP enhances coherence by propagating information across neighboring frames.

evidenced by the lowest warping error [30] and MEt3R [2], which further strengthens overall 4D fidelity.

**4D Fidelity.** Fig. 5 demonstrates how improved 2D consistency translates into higher-quality 4D scene edits. Without our components, the edited scene exhibits severe motion artifacts, especially around the man's hand, along with background degradation and incorrect text alignment, such as failing to produce "charcoal drawing" colors. Adding STGA reduces large-scale artifacts and stabilizes dynamic motion, while incorporating CTP further refines fine-grained details by leveraging temporal and multi-view cues from previous grids. With all components combined, Dynamic-eDiTor achieves consistent motion reconstruction and editing, effectively eliminating geometric drift in the 4D rendered output. Tab. 6 reveals that each component reinforces 4D fidelity. The PSNR, SSIM, and LPIPS scores validate this, demonstrating that Dynamic-eDiTor delivers highly coherent edits. We also note that CLIP-based metrics are slightly higher when STGA is removed. As mentioned in Sec. 4.1.1, this reflects the trade-off between semantic alignment and spatio-temporal coherence. balanced, spatio–temporally coherent 4D reconstruction, rather than optimizing semantic alignment alone. Our method prioritizes a balanced, stable, and coherent 4D reconstruction, whereas the ablated variant attains higher CLIP-based metrics by sacrificing this stability, resulting in geometrically and temporally unstable reconstruction.

**Context Token Propagation.** We further ablate our Context Token Propagation (CTP) components in Tab. 7. First, removing only the Full Token Inheritance leads to reduced 2D multi-view consistency and lower fidelity in the 4D rendered images. Next, removing only the Flow-guided Token Replacement results in a significant drop in temporal



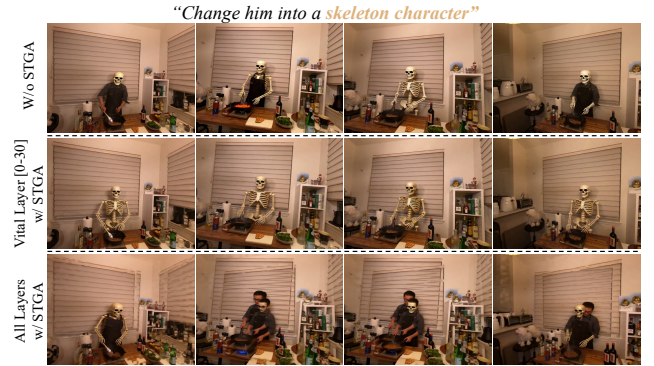*"Change him into a skeleton character"*

Figure 7. **Qualitative Analysis of Vital Layer Range.** Applying STGA to all layers introduces visual artifacts across views and time, while omitting STGA produces inconsistent multi-view and temporal edits. Restricting STGA to the vital range yields the most coherent and stable multi-view–time editing results.

consistency. Finally, removing the entire Context Token Propagation mechanism (both components) causes a severe degradation in performance, dramatically worsening both 2D consistency and 4D reconstruction fidelity. Similar to the previous ablation, these results reflect the inherent trade-off between semantic alignment and spatio-temporal coherence, confirming that both the Full Token Inheritance and Flow-guided Token Replacement are essential for achieving high-quality and consistent 4D editing.

## 6. Conclusion

We presented Dynamic-eDiTor, a training-free framework for text-driven 4D scene editing that achieves spatially and temporally consistent results across multi-view videos, enabling stable optimization of 4D representations with MM-DiT [11, 54] and 4DGS [55]. The core of our approach is Grid-based Spatio-Temporal Propagation, combining Spatio-Temporal Sub-Grid Attention (STGA) for localized view-time fusion and Context Token Propagation (CTP) for explicit global consistency. Together, these components ensure coherent geometry and motion, and high-fidelity dynamic edits. Extensive experiments on DyNeRF [33] demonstrate our method significantly outperforms prior work in editing fidelity, temporal smoothness, and robustness. We believe Dynamic-eDiTor represents a notable progression toward text-driven dynamic scene editing.

**Limitation.** While Dynamic-eDiTor effectively enforces multi-view and temporal consistency during text-driven edits, it is less suitable for large-scale geometric alterations such as substantial motion reconfiguration or topology-changing edits. Since our framework operates by propagating spatio-temporal features without modeling geometric deformation, edits that require significant structural changes remain challenging. Extending our approach to handle more drastic motion editing or geometry-changing transformations represents an important direction for future work.

## 7. Acknowledgement

## References

[1] Amazon mechanical turk. https://www.mturk.com/, 2005. 7, 2

[2] Mohammad Asim, Christopher Wewer, Thomas Wimmer, Bernt Schiele, and Jan Eric Lenssen. Met3r: Measuring multi-view consistency in generated images. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 6034–6044, 2025. 4, 8, 1

[3] Omri Avrahami, Or Patashnik, Ohad Fried, Egor Nemchinov, Kfir Aberman, Dani Lischinski, and Daniel Cohen-Or. Stable flow: Vital layers for training-free image editing. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 7877–7888, 2025. 4

[4] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18392–18402, 2023. 2, 3, 7, 1

[5] Qi Cai, Jingwen Chen, Yang Chen, Yehao Li, Fuchen Long, Yingwei Pan, Zhaofan Qiu, Yiheng Zhang, Fengbin Gao, Peihan Xu, et al. Hidream-i1: A high-efficient image generative foundation model with sparse diffusion transformer. *arXiv preprint arXiv:2505.22705*, 2025. 2

[6] Minghao Chen, Iro Laina, and Andrea Vedaldi. Dge: Direct gaussian 3d editing by consistent multi-view editing. In *European Conference on Computer Vision*, pages 74–92. Springer, 2024. 2

[7] Wilfrid J Dixon and Alexander M Mood. The statistical sign test. *Journal of the American Statistical Association*, 41(236):557–566, 1946. 3

[8] Hyungjun Doh, Dong In Lee, Seunggeun Chi, Pin-Hao Huang, Kwonjoon Lee, Sangpil Kim, and Karthik Ramani. Occlusion-aware temporally consistent amodal completion for 3d human-object interaction reconstruction. In *Proceedings of the 33rd ACM International Conference on Multimedia*, pages 52–61, 2025. 2

[9] Jiahua Dong and Yu-Xiong Wang. Vica-nerf: View-consistency-aware 3d editing of neural radiance fields. *Advances in Neural Information Processing Systems*, 36: 61466–61477, 2023. 2

[10] Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2

[11] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024. 1, 2, 3, 4, 6, 8

[12] Haoran Feng, Zehuan Huang, Lin Li, Hairong Lv, and Lu Sheng. Personalize anything for free with diffusion transformer. *arXiv preprint arXiv:2503.12590*, 2025. 4

[13] Hang Gao, Ruilong Li, Shubham Tulsiani, Bryan Russell, and Angjoo Kanazawa. Dynamic novel-view synthesis: A reality check. In *NeurIPS*, 2022. 5

[14] Sicheng Gao, Nancy Mehta, Zongwei Wu, and Radu Timofte. Ditvr: Zero-shot diffusion transformer for video restoration. *arXiv preprint arXiv:2508.07811*, 2025. 4

[15] Michal Geyer, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Tokenflow: Consistent diffusion features for consistent video editing. *arXiv preprint arXiv:2307.10373*, 2023. 4

[16] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. 3

[17] Ayaan Haque, Matthew Tancik, Alexei A Efros, Aleksander Holynski, and Angjoo Kanazawa. Instruct-nerf2nerf: Editing 3d scenes with instructions. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 19740–19750, 2023. 2, 1

[18] Kai He, Chin-Hsuan Wu, and Igor Gilitschenski. Ctrl-d: Controllable dynamic 3d scene editing with personalized 2d diffusion. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 26630–26640, 2025. 2, 3, 6, 7, 1

[19] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022. 2

[20] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022. 1

[21] Dadong Jiang, Zhihui Ke, Xiaobo Zhou, Tie Qiu, Xidong Shi, and Hao Yan. 4d-editor: Interactive object-level editing in dynamic neural radiance fields via semantic distillation. In *2025 International Conference on 3D Vision (3DV)*, pages 702–712. IEEE, 2025. 3

[22] Nazmul Karim, Hasan Iqbal, Umar Khalid, Chen Chen, and Jing Hua. Free-editor: zero-shot text-driven 3d scene editing. In *European Conference on Computer Vision*, pages 436–453. Springer, 2024. 2

[23] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 1, 2, 3

[24] Min-Jung Kim, Dongjin Kim, Seokju Yun, and Jaegul Choo. Tv-live: Training-free, text-guided video editing via layer informed vitality exploitation. *arXiv preprint arXiv:2506.07205*, 2025. 4

[25] Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 1

[26] Tobias Kirschstein, Javier Romero, Artem Sevastopolsky, Matthias Nießner, and Shunsuke Saito. Avat3r: Large animatable gaussian reconstruction model for high-fidelity 3d head avatars. *arXiv preprint arXiv:2502.20220*, 2025. 2

[27] Joohyun Kwon, Hanbyel Cho, and Junmo Kim. Efficient dynamic scene editing via 4d gaussian-based static-dynamic separation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 26855–26865, 2025. 2, 3, 6, 7, 1

[28] Black Forest Labs. Flux. `https://github.com/black-forest-labs/flux`, 2024. 2

[29] Black Forest Labs, Stephen Batifol, Andreas Blattmann, Frederic Boesel, Saksham Consul, Cyril Diagne, Tim Dockhorn, Jack English, Zion English, Patrick Esser, Sumith Kulal, Kyle Lacey, Yam Levi, Cheng Li, Dominik Lorenz, Jonas Müller, Dustin Podell, Robin Rombach, Harry Saini, Axel Sauer, and Luke Smith. Flux.1 kontext: Flow matching for in-context image generation and editing in latent space, 2025. 2

[30] Wei-Sheng Lai, Jia-Bin Huang, Oliver Wang, Eli Shechtman, Ersin Yumer, and Ming-Hsuan Yang. Learning blind video temporal consistency. In *Proceedings of the European conference on computer vision (ECCV)*, pages 170–185, 2018. 4, 8, 2

[31] Dong In Lee, Hyeongcheol Park, Jiyoung Seo, Eunbyung Park, Hyunje Park, Ha Dam Baek, Sangheon Shin, Sangmin Kim, and Sangpil Kim. Editsplat: Multi-view fusion and attention-guided optimization for view-consistent 3d scene editing with 3d gaussian splatting. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 11135–11145, 2025. 2

[32] Vincent Leroy, Yohann Cabon, and Jérôme Revaud. Grounding image matching in 3d with mast3r. In *European Conference on Computer Vision*, pages 71–91. Springer, 2024. 1

[33] 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*, pages 5521–5531, 2022. 2, 6, 8

[34] Simon Meister, Junhwa Hur, and Stefan Roth. Unflow: Unsupervised learning of optical flow with a bidirectional cen-

sus loss. In *Proceedings of the AAAI conference on artificial intelligence*, 2018. 5

[35] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021. 2

[36] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 1, 2

[37] Linzhan Mou, Jun-Kun Chen, and Yu-Xiong Wang. Instruct 4d-to-4d: Editing 4d scenes as pseudo-3d scenes using 2d diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20176–20185, 2024. 2, 3, 6, 7, 1

[38] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 1

[39] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023. 2, 4

[40] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 3, 1

[41] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10318–10327, 2021. 1

[42] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR, 2021. 4, 6, 7, 2, 5

[43] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 2

[44] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22500–22510, 2023. 3

[45] Ruizhi Shao, Jingxiang Sun, Cheng Peng, Zerong Zheng, Boyao Zhou, Hongwen Zhang, and Yebin Liu. Control4d: Efficient 4d portrait editing with text. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4556–4567, 2024. 3

[46] Liangchen Song, Anpei Chen, Zhong Li, Zhang Chen, Lele Chen, Junsong Yuan, Yi Xu, and Andreas Geiger. Nerf-

player: A streamable dynamic scene representation with decomposed neural radiance fields. *IEEE Transactions on Visualization and Computer Graphics*, 29(5):2732–2742, 2023. 2

[47] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *European conference on computer vision*, pages 402–419. Springer, 2020. 2, 5

[48] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 2

[49] Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul, Mishig Davaadorj, Dhruv Nair, Sayak Paul, William Berman, Yiyi Xu, Steven Liu, and Thomas Wolf. Diffusers: State-of-the-art diffusion models. https://github.com/huggingface/diffusers, 2022. 1

[50] Junjie Wang, Jiemin Fang, Xiaopeng Zhang, Lingxi Xie, and Qi Tian. Gaussianeditor: Editing 3d gaussians delicately with text instructions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 20902–20911, 2024. 2

[51] Yuxuan Wang, Xuanyu Yi, Zike Wu, Na Zhao, Long Chen, and Hanwang Zhang. View-consistent 3d editing with gaussian splatting. In *European conference on computer vision*, pages 404–420. Springer, 2024. 2

[52] 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. 6, 7, 2

[53] Tianyi Wei, Yifan Zhou, Dongdong Chen, and Xingang Pan. Freeflux: Understanding and exploiting layer-specific roles in rope-based mmdit for versatile image editing. *arXiv preprint arXiv:2503.16153*, 2025. 4

[54] Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng-ming Yin, Shuai Bai, Xiao Xu, Yilei Chen, et al. Qwen-image technical report. *arXiv preprint arXiv:2508.02324*, 2025. 1, 2, 3, 4, 6, 8

[55] 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*, pages 20310–20320, 2024. 1, 2, 3, 6, 8

[56] Jing Wu, Jia-Wang Bian, Xinghui Li, Guangrun Wang, Ian Reid, Philip Torr, and Victor Adrian Prisacariu. Gaussctrl: Multi-view consistent text-driven 3d gaussian splatting editing. In *European Conference on Computer Vision*, pages 55–71. Springer, 2024. 2

[57] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7623–7633, 2023. 4

[58] Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezatofighi, and Dacheng Tao. Gmflow: Learning optical flow via global matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8121–8130, 2022. 5

[59] Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin. Deformable 3d gaussians for high-fidelity monocular dynamic scene reconstruction. *arXiv preprint arXiv:2309.13101*, 2023. 2, 5

[60] 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*, pages 586–595, 2018. 6, 7, 2

[61] Canyu Zhao, Xiaoman Li, Tianjian Feng, Zhiyue Zhao, Hao Chen, and Chunhua Shen. Tinker: Diffusion's gift to 3d–multi-view consistent editing from sparse inputs without per-scene optimization. *arXiv preprint arXiv:2508.14811*, 2025. 2

[62] Zhe Zhu, Honghua Chen, Peng Li, and Mingqiang Wei. Coreeditor: Consistent 3d editing via correspondence-constrained diffusion. *arXiv preprint arXiv:2508.11603*, 2025. 2

[63] Qi Zuo, Xiaodong Gu, Yuan Dong, Zhengyi Zhao, Weihao Yuan, Lingteng Qiu, Liefeng Bo, and Zilong Dong. High-fidelity 3d textured shapes generation by sparse encoding and adversarial decoding. In *European Conference on Computer Vision*, pages 52–69. Springer, 2024. 2

# Dynamic-eDiTor: Training-Free Text-Driven 4D Scene Editing with Multimodal Diffusion Transformer

## Supplementary Material

## Overview

This supplementary material provides additional details, analyses, and experimental results for our proposed method, Dynamic-eDiTor.

- Fig. 11 and Fig. 12 present additional qualitative results, including extended comparisons with baseline methods.

- Sec. A and Sec. B provide further implementation details and descriptions of all evaluation metrics used in our experiments.

- Sec. C summarizes the user study protocol and provides a detailed analysis of participant responses.

- Sec. D details our Grid-based Spatio-Temporal Propagation mechanism, including Asymmetric Traversal Strategy and the accompanying algorithm in Algorithm 1.

- Sec. E offers an extended analysis of the *vital layer range* for Spatio-Temporal Sub-Grid Attention (STGA).

- Sec. F contains additional ablation studies analyzing Asymmetric Traversal Strategy of Dynamic-eDiTor.

- Sec. G presents additional qualitative results in the monocular video setting of 4D Gaussian Splatting (4DGS) [55].

## A. Implementation Details

For each scene, we first train the source 4D Gaussian Splatting [55] representation for 30,000 iterations using the Adam optimizer [25] with the same learning rate schedule as 4DGS. During the editing stage, we optimize the model for 20,000 iterations using the edited frames, following the original 4DGS hyperparameter configuration. All experiments are conducted on an NVIDIA H100 GPU; however, by employing local caching for Temporal Context Token Replacement, our method also runs efficiently on an NVIDIA A6000 GPU.

For the 2D MM-DiT [11, 54] image editor, we utilize Qwen-Image-Edit [54] from the Diffusers library [49]. To enhance computational efficiency, we incorporate the LoRA [20] weight Qwen-Image-Lightning-8steps-V1.1. All input images are resized to $768 \times 768$ before processing.

For baseline comparisons, we follow the official implementations of Instruct4D-to-4D [37] and Instruct-4DGS [27], since both are text-driven 4D editing methods comparable to ours. CTRL-D [18] requires an additional

---

**Algorithm 1** Asymmetric Sub-Grid Traversal

**Require:** Camera–time grid $\text{Grid} = \{f_{v,t}\}$ of size $V \times T$
**Ensure:** Ordered list of sub-grids $\Omega = \{\mathcal{S}^{(k)}\}$
1: $\Omega \leftarrow [\,]$        ▷ Initialize empty sub-grid sequence
2: **for** $v = 0$ **to** $V - 2$ **do**
3:     **if** $v$ is even **or** $v = V - 2$ **then**    ▷ Temporal sweep
4:         **for** $t = 0$ **to** $T - 2$ **do**
5:             $\mathcal{S}_{v,t} \leftarrow \{f_{v,t}, f_{v+1,t}, f_{v,t+1}, f_{v+1,t+1}\}$
6:             Append $\mathcal{S}_{v,t}$ to $\Omega$
7:         **end for**
8:     **else**               ▷ Cross-view alignment at $t = 0$
9:         $\mathcal{S}_{v,0} \leftarrow \{f_{v,0}, f_{v+1,0}, f_{v,1}, f_{v+1,1}\}$
10:       Append $\mathcal{S}_{v,0}$ to $\Omega$
11:     **end if**
12: **end for**
13: **return** $\Omega$

---

edited reference image that must be produced by choosing one of several diffusion-based editing modes, such as image-prompt editing, text-prompt editing, or mask-based editing, before fine-tuning its InstructPix2Pix [4] backbone. To ensure a fair and consistent comparison, we fix this pre-editing stage to use the standard InstructPix2Pix image editor, which is the backbone originally used in CTRL-D, when generating the reference edited image.

## B. Metric

To evaluate Dynamic-eDiTor, we use a combination of 2D consistency, 4D editing fidelity, and 4D reconstruction fidelity metrics.

For **2D consistency**, we evaluate both temporal and multi-view stability. Our 4D editing baselines such as Instruct4D-to-4D [37] and CTRL-D [18] rely on Iterative Dataset Update (IDU)[17], and Instruct-4DGS[27] uses an SDS-based [40] optimization strategy. However, these approaches do not generate temporally aligned or viewpoint-consistent 2D edited frames. Their updates are stochastic and occur directly in 3D or 4D space, which makes extracting coherent multi-view video sequences infeasible. Therefore, 2D consistency metrics cannot be fairly compared with these baselines and are used only within our ablation studies.

- **MEt3R** [2]: Evaluates multi-view consistency by comparing feature similarity between view-warped images. We employ the official MEt3R metric with MASt3R [32], DINOv2 [38] (FeatUp) features, 448 image resolution,
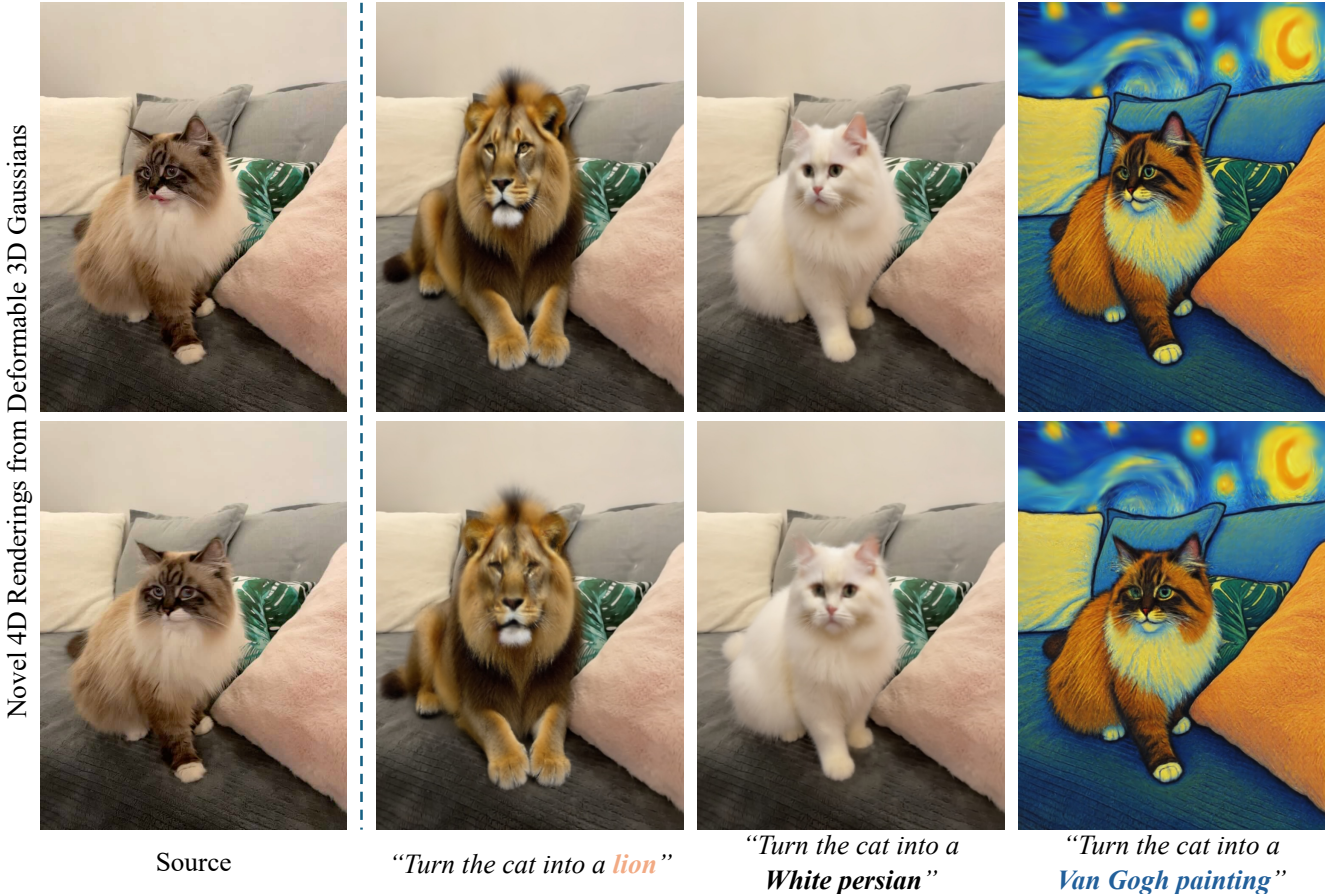
Figure 8. **Qualitative Results in Monocular video setting.** We evaluate the applicability of Dynamic-eDiTor in the challenging monocular video setting, where only a single moving-camera sequence is available and no multi-view redundancy exists. Using Deformable 3D Gaussian Splatting [59] as the underlying 4D representation, our method successfully performs text-driven appearance and object manipulations while maintaining stable geometry and consistent motion over time.

and cosine similarity. Lower values indicate more coherent appearance across viewpoints.

- **Warping Error** [30]: Measures temporal consistency by computing the discrepancy between frame $f_t$ and the optical-flow–warped version of frame $f_{t-1}$ using RAFT [47]. Lower scores indicate smoother temporal alignment and fewer motion artifacts.

For **4D editing fidelity**, we adopt CLIP-based metrics [42]. We compute both CLIP text-image directional similarity and CLIP text-image similarity using the *rendered images* produced by the edited 4D scene. The directional similarity evaluates whether the change described by the text prompt corresponds to the transformation from the source image to the edited rendering in CLIP embedding space. The CLIP text-image similarity, on the other hand, directly measures how well the rendered frames semantically align with the target text prompt.

For **4D reconstruction fidelity**, we report PSNR, SSIM [52], and LPIPS [60], following prior works [8, 26, 27, 63]. All three metrics are computed between the edited test-view image and the rendered test-view image from the same camera viewpoint, enabling a direct comparison of reconstruction quality.

## C. User Study Detail

To compare the editing performance of Dynamic-eDiTor against baseline methods, we conducted a user study with 150 participants on Amazon Mechanical Turk [1]. Each participant evaluated 14 scenarios, and for each scenario, they compared the 4D rendered video results produced by four systems across six subjective dimensions, as illustrated in Fig. 13. We designed six evaluation questions covering prompt alignment (Q1), temporal consistency (Q2), viewpoint consistency (Q3), motion consistency (Q4), identity preservation (Q5), and overall visual quality (Q6). For each dimension, participants selected the system they judged to perform best. To reduce human bias, the presentation order of the four systems was randomized for every question.

For analysis, we first counted how many times Dynamic-eDiTor was selected across the 14 scenarios and compared it with the best-performing baseline (best baseline) on each dimension. The results show that Dynamic-eDiTor consistently outperformed the best baseline across all six evaluation dimensions. For example, on overall quality (Q6),

| AGT | STGA | 2D Consistency | | | | Reconstruction Fidelity | | | Editing Fidelity | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Local Warp-Err $_{10^{-3}}$ ↓ | Global Warp-Err $_{10^{-3}}$ ↓ | Local MEt3R $_{10^{-1}}$ ↓ | Global MEt3R $_{10^{-1}}$ ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ | CLIP$_{dir}$ ↑ | CLIP$_{sim}$ ↑ |
| - | - | 56.98 | 58.47 | 1.0721 | 1.4312 | 26.14 | 0.7445 | 0.1408 | **0.1930** | 0.6414 |
| - | ✓ | **34.56** | 42.86 | **0.9266** | 1.2984 | 27.84 | 0.7793 | 0.1160 | 0.1876 | **0.6499** |
| ✓ | ✓ | 38.64 | **42.33** | 0.9277 | **1.2953** | **28.08** | **0.7875** | **0.1122** | 0.1872 | 0.6407 |

Table 4. **Ablation Study: Asymmetric Sub-Grid Traversal (AGT).** This evaluation is conducted without CTP to isolate the impact of Asymmetric Sub-Grid Traversal (AGT). The results show that sub-grids without AGT achieve slightly better local consistency metrics because all frames within each sub-grid are updated independently. However, the lack of linkage between sub-grids introduces discontinuities, weakening overall 4D reconstruction fidelity. In contrast, applying AGT improves global consistency by overlapping frames across sub-grids, even at the cost of some local editing precision, as it enables effective information propagation. This leads to more stable and reliable 4D edits, demonstrating that global consistency is ultimately more critical for 4D reconstruction fidelity.

Dynamic-eDiTor achieved an average selection rate of 0.49, compared to 0.28 for the best baseline (Instruct4d [37]). The advantage is even more pronounced for prompt alignment (Q1), with selection rates of 0.57 vs. 0.22 (Instruct4d). For other questions, Dynamic-eDiTor's average selection rate exceeded the best baseline by approximately 0.17–0.21, demonstrating stable and comprehensive improvements.

To examine whether these differences were statistically significant, we performed a two-sided signed [7] test for each question, pairing each participant's selection ratio for Dynamic-eDiTor with that of the best baseline. All six dimensions yielded p-values far below 0.01, specifically: Q1 ($p = 8.88 \times 10^{-16}$), Q2 ($p = 5.43 \times 10^{-3}$), Q3 ($p = 4.49 \times 10^{-5}$), Q4 ($p = 1.41 \times 10^{-4}$), Q5 ($p = 9.33 \times 10^{-5}$), and Q6 ($p = 2.10 \times 10^{-5}$).

These results confirm that human evaluators consistently prefer Dynamic-eDiTor over the best baseline. The static significance also reveals that the gains are robust rather than due to random variation.

# D. Asymmetric Sub-Grid Traversal with Overlapping Structure

In the main paper, Grid-based Spatio-Temporal Propagation is introduced as a mechanism that performs local fusion via Spatio-Temporal Sub-Grid Attention (STGA) and global propagation via Context Token Propagation (CTP). In this section, we provide additional details on how the camera–time grid $Grid$ is traversed and how overlapping sub-grids are constructed to enable stable spatio-temporal propagation. Algorithm 1 formalizes the sub-grid generation process used in our implementation.

## D.1. Overlapping Sub-Grids

Since neighboring sub-grids share frames on their boundaries, they form an overlapping tiling of the camera–time grid. This overlap is crucial for CTP: the shared regions act as "anchors" through which coherent token representations can be propagated from one sub-grid to the next.
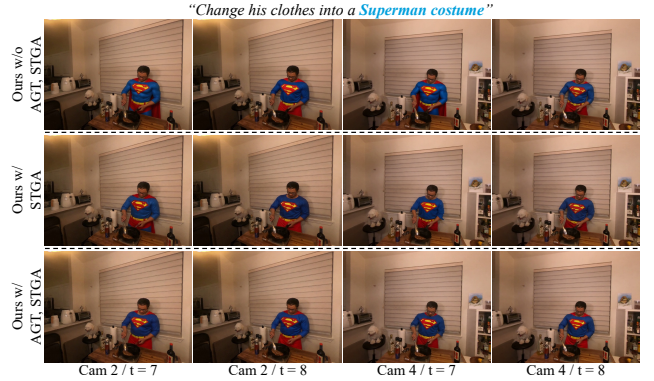


*"Change his clothes into a Superman costume"*

Figure 9. **Ablation Study: Asymmetric Sub-Grid Traversal (AGT).** This qualitative result clearly demonstrates that AGT preserves global multi-view and temporal consistency. Without AGT, noticeable discontinuities appear between sub-grids.

## D.2. Asymmetric Traversal Strategy

Rather than processing all $\mathcal{S}_{v,t}$ in a simple raster-scan order, we adopt an *asymmetric traversal* strategy that balances temporal propagation and cross-view alignment.

Concretely, we iterate over camera indices $v = 0, \ldots, V-2$, and for each $v$ we choose a different temporal traversal pattern:

- For **even** camera indices (and the last camera pair $v = V-2$), we perform a *full temporal sweep*, generating sub-grids $\mathcal{S}_{v,t}$ for all $t \in [0, T-2]$. This encourages strong temporal propagation along the time axis.
- For **odd** camera indices, we generate only $\mathcal{S}_{v,0}$ (i.e., using the first two time steps $t = 0, 1$). This enforces cross-view alignment between neighboring cameras while avoiding redundant temporal passes.

The resulting traversal order can be summarized as follows: even-indexed camera rows perform dense temporal coverage, while odd-indexed rows act as cross-view bridges at the initial time step. This pattern yields an overlapping chain of sub-grids that spans the entire $V \times T$ grid, ensuring that information fused by STGA in one region can be propagated to distant regions through CTP.

## D.3. Effect on STGA and CTP

This asymmetric, overlapping traversal has two key effects:

- **Local fusion via STGA.** Within each $\mathcal{S}_{v,t}$, STGA jointly attends over adjacent views and neighboring time steps, producing locally coherent spatio-temporal features. Due to the overlapping structure, boundary frames participate in multiple sub-grids, implicitly coupling neighboring regions.
- **Global propagation via CTP.** CTP operates along the traversal order $\Omega$, propagating tokens from $\mathcal{S}_{prev}$ to $\mathcal{S}_{curr}$ through inherited and flow-guided token replacement. Since the sub-grids overlap in both view and time, this propagation forms a connected path over the entire grid, enabling the fused information to spread globally while respecting camera–time structure.

Together, the asymmetric traversal and overlapping sub-grids provide a principled backbone for Grid-based Spatio-Temporal Propagation, ensuring that local STGA fusion and global CTP propagation jointly enforce consistent editing across all views and time steps.

## E. Vital Layer Range Analysis for STGA

In the main paper, we apply Spatio-Temporal Sub-Grid Attention (STGA) only to a vital layer range of MM-DiT in order to enhance spatio-temporal consistency without overly harming editing fidelity. Here, we provide a more detailed quantitative analysis of this design choice.

**Experimental setup.** We conduct a systematic study on the DyNeRF dataset using 3 scenes sampled at 1 FPS and 5 editing prompts per scene (15 sequences in total). For each configuration, we enable STGA on a different continuous range of MM-DiT layers and keep all other components fixed. We report (i) *Warping Error↓* [30] for temporal consistency, (ii) *MEt3R↓* [2] for multi-view consistency, and (iii) *CLIP Text-Image Directional Similarity↑* (CLIP$_{dir}$) [42] for editing fidelity.

**Effect of early-layer STGA.** Without STGA, both Warping Error and MEt3R are the worst (46.37 and 1.22), indicating pronounced temporal flicker and cross-view inconsistency. As we progressively introduce STGA from shallow layers (0–9, 0–19, 0–29), both consistency metrics steadily improve. In particular, enabling STGA on the first 30 layers (0–29) yields a strong reduction in temporal and multi-view error (Warping Error 30.90, MEt3R 0.99), while preserving a relatively high CLIP$_{dir}$ score (0.088). This configuration achieves the best overall trade-off: it significantly enhances spatio-temporal coherence compared to the baseline, yet maintains competitive editing fidelity.

| Layer Range | Warp-Err $_{10^{-3}}$ ↓ | MEt3R $_{10^{-1}}$ ↓ | CLIP$_{dir}$ ↑ |
|---|---|---|---|
| W/o STGA | 46.37 | 1.221 | 0.1111 |
| 0–9 | 32.54 | 1.022 | 0.1014 |
| 0–19 | 32.25 | 1.006 | 0.0946 |
| 0–29 | 30.90 | 0.993 | 0.0879 |
| 0–39 | 28.32 | 0.956 | 0.0468 |
| 25–35 | 37.47 | 1.578 | 0.0693 |
| 20–40 | 37.41 | 1.581 | 0.0683 |
| 15–45 | 28.96 | 1.321 | 0.0863 |
| 10–50 | 38.81 | 1.538 | 0.0661 |
| 49–59 | 41.06 | 1.270 | 0.0704 |
| 39–59 | 41.82 | 1.279 | 0.0703 |
| 29–59 | 39.49 | 1.250 | 0.0829 |
| 19–59 | 41.74 | 1.272 | 0.0694 |
| All layers | 41.25 | 1.265 | 0.0689 |

Table 5. **Detailed vital layer range analysis for STGA.** This table reports the exact numerical values corresponding to the trend shown in Figure 3 of the main paper.

**Applying STGA too deep.** Extending STGA too aggressively into deeper layers (e.g., 0–39 or mid-to-deep ranges such as 25–35, 20–40, 10–50, or 19–59) further reduces or even oscillates consistency metrics, but at the cost of a substantial drop in CLIP$_{dir}$. For example, 0–39 attains the lowest Warping Error and MEt3R among all settings, but its CLIP$_{dir}$ score collapses to 0.047, indicating that over-attending within local spatio-temporal neighborhoods can oversmooth edits, weaken text alignment, and lead to texture repetition or view-dependent artifacts, as shown in Fig. 7 of the main paper. Similarly, configurations that only activate STGA in deeper blocks (e.g., 39–59, All layers) neither recover the consistency of early-layer STGA nor preserve high editing fidelity, suggesting that late-stage modifications are less effective for enforcing stable geometry and motion.

**Chosen configuration.** Based on these observations, we adopt the 0–29 configuration as our default choice in the main paper. This vital layer range provides a balanced compromise: it substantially improves temporal and multi-view consistency over the baseline and deep-only variants, while incurring only a modest decrease in CLIP$_{dir}$ relative to the no-STGA setting. In practice, we find that this trade-off yields visually smoother 4D reconstructions and more reliable 4D scene editing, whereas configurations with either no STGA or overly deep STGA tend to produce flickering, geometric drift, or oversmoothed, weakly edited results.

## F. Additional Ablation Study

We conduct additional ablation study to evaluate the impact of Asymmetric Sub-Grid Traversal (AGT) in Dynamic-eDiTor.

4

| STGA | CTP | 2D Consistency | | | | Reconstruction Fidelity | | | Editing Fidelity | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Local Warp-Err $_{10^{-3}}$ ↓ | Global Warp-Err $_{10^{-3}}$ ↓ | Local MEt3R $_{10^{-1}}$ ↓ | Global MEt3R $_{10^{-1}}$ ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ | CLIP$_{dir}$ ↑ | CLIP$_{sim}$ ↑ |
| - | - | 56.98 | 58.37 | 1.0721 | 1.4312 | 26.14 | 0.7445 | 0.1408 | <u>0.1930</u> | 0.5414 |
| ✓ | - | 38.64 | 42.33 | <u>0.9277</u> | <u>1.2953</u> | 28.08 | 0.7875 | <u>0.1122</u> | 0.1872 | <u>0.6407</u> |
| - | ✓ | <u>29.44</u> | <u>31.63</u> | 1.0695 | 1.4364 | <u>28.74</u> | <u>0.8013</u> | 0.1165 | **0.1944** | **0.6418** |
| ✓ | ✓ | **28.94** | **30.69** | **0.9074** | **1.2657** | **29.25** | **0.8064** | **0.1006** | 0.1849 | 0.6397 |

Table 6. **Ablation Study: Local and Global Consistency.** Our method improves both local and global 2D consistency, ensuring that each sub-grid remains coherent both internally (local) and with its neighbors (global). Each component, STGA and CTP, helps maintain temporal and multi-view consistency, improving overall 4D reconstruction. By enforcing a globally stable 4D structure, our method achieves more consistent spatio-temporal behavior and higher reconstruction fidelity. Although CLIP-based metrics [42] show a slight drop due to the trade-off between semantic alignment and spatio-temporal coherence, our approach still delivers more stable and reliable 4D edits, avoiding the geometric and temporal artifacts seen in the ablated variants.

| CTP-Full | CTP-Flow | 2D Consistency | | | | Reconstruction Fidelity | | | Editing Fidelity | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Local Warp-Err $_{10^{-3}}$ ↓ | Global Warp-Err $_{10^{-3}}$ ↓ | Local MEt3R $_{10^{-1}}$ ↓ | Global MEt3R $_{10^{-1}}$ ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ | CLIP$_{dir}$ ↑ | CLIP$_{sim}$ ↑ |
| - | - | 56.98 | 58.37 | 1.0721 | 1.4312 | 26.14 | 0.7445 | 0.1408 | **0.1930** | **0.6407** |
| - | ✓ | <u>29.79</u> | <u>32.66</u> | 0.9205 | 1.2813 | <u>28.97</u> | <u>0.7990</u> | <u>0.1034</u> | 0.1852 | <u>0.6402</u> |
| ✓ | - | 33.22 | 37.36 | <u>0.9094</u> | <u>1.2736</u> | 28.19 | 0.7906 | 0.1089 | <u>0.1865</u> | 0.6400 |
| ✓ | ✓ | **28.94** | **30.69** | **0.9074** | **1.2657** | **29.25** | **0.8064** | **0.1006** | 0.1849 | 0.6397 |

Table 7. **Ablation Study: Context Token Propagation (CTP).** This ablation study is conducted with STGA included to isolate the effect of CTP. The results show that our method maintains both local (within each sub-grid) and global consistency. Full Token Inheritance (CTP-Full) and Flow-Guided Token Replacement (CTP-Flow) play a crucial role in reinforcing temporal and multi-view coherence, enabling more accurate reconstruction of edited dynamic scenes. Although CLIP-based metrics [42] show a slight trade-off, CTP significantly enhances spatio-temporal consistency and overall 4D editing fidelity.

**Asymmetric Sub-Grid Traversal (AGT)** For the ablation of AGT, we perform experiments without CTP, as it relies on the sliding mechanism introduced by AGT. AGT is designed to create overlapping regions between adjacent sub-grids, promoting smoother transitions and improving global consistency. We hypothesize that removing this overlap will yield results that remain locally consistent within each sub-grid but exhibit severe temporal and multi-view discontinuities across sub-grid boundaries.

To fairly assess global consistency, we introduce two new metrics: **Global Warping Error** and **Global MEt3R**. Unlike their standard versions, these metrics are computed only between the left boundary frames of adjacent sub-grids, directly quantifying the discontinuities that AGT aims to mitigate. Additionally, we define per-frame consistency metrics as **Local Warping Error** and **Local MEt3R**, which are the original 2D consistency metric used in the main paper.

As shown in Tab. 4, removing AGT yields slightly higher local consistency because each sub-grid updates all its frames during every iteration. However, without any connection between sub-grids, clear discontinuities emerge between them, degrading overall 4D reconstruction quality. In contrast, AGT introduces overlap across sub-grids,

enabling information to propagate between them and producing more coherent results, despite a slight reduction in local consistency, as only the non-overlapping frames are updated. Overall, these findings confirm that AGT plays a crucial role in preserving global coherence, enabling higher-quality and more faithful 4D reconstruction than using STGA alone.

We also observe that the non-sliding variant yields slightly higher CLIP scores, reflecting the inherent trade-off between consistency and text alignment. Without sliding, the model can more aggressively edit each isolated sub-grid to match the prompt, but this comes at the cost of failing to produce a globally coherent 4D scene, which is the primary objective of our method.

## G. Monocular Video Setting

4D dynamic scene editing typically refers to a multi-view video setting where sufficient spatio-temporal information is captured. However, to explore the applicability of Dynamic-eDiTor in a monocular setting, we evaluate our method on the DyCheck [13] dataset using Deformable 3D Gaussian Splatting model [59]. Since a monocular dataset contains only a single camera, we modify the camera–time grid to a purely temporal grid:

$$Grid_{temp} = \{f_t \mid t \in [0, \ldots, T]\}, \tag{9}$$

Accordingly, each sub-grid $\mathcal{S}_t$ consists of consecutive frames along the temporal axis:

$$\mathcal{S}_t = \{f_t, \ f_{t+1}, \ f_{t+2}, \ f_{t+3}\}. \tag{10}$$

Based on this modified sub-grid, we apply same Spatio-Temporal Sub-Grid Attention (STGA) mechanism. However, because only temporally adjacent frames are available, the key and value sets $K_{\mathcal{S}_t}$ and $V_{\mathcal{S}_t}$ become:

$$K_{\mathcal{S}_t} = [K_{f_t}, K_{f_{t+1}}, K_{f_{t+2}}, K_{f_{t+3}}],$$
$$V_{\mathcal{S}_t} = [V_{f_t}, V_{f_{t+1}}, V_{f_{t+2}}, V_{f_{t+3}}]. \tag{11}$$

Thus, STGA in the monocular setting becomes:

$$\mathrm{STGA}(\mathcal{S}_t) = \mathrm{softmax}\Big([Q_{\mathrm{txt}}, \ \mathrm{RoPE}(Q_{f_t})] \cdot$$
$$[K_{\mathrm{txt}}, \ \mathrm{RoPE}(K_{\mathcal{S}_t})]^\top / \sqrt{d_k}\Big) \cdot [V_{\mathrm{txt}}, \ V_{\mathcal{S}_t}], \tag{12}$$

where $d_k$ denotes the dimensionality of the key vectors.

For Context Token Propagation (CTP), where the token representation is defined as $\phi(\mathcal{S}_t) = \mathrm{STGA}(\mathcal{S}_t)$, we employ two Context Token Propagation strategies: Full Token Inheritance and Flow-guided Token Replacement. Since the sub-grids overlap by two temporal frames, we directly replace the entire current token $\phi(\mathcal{S}_{curr})$ in these overlapped frames with the previous token $\phi(\mathcal{S}_{prev})$. For the non-overlapped region, which corresponds to the two rightmost frames of the sub-grid, we apply flow-guided token replacement. To propagate the most recent temporal information, we warp tokens from the rightmost frame of the overlapped region and replace the tokens of the two non-overlapping frames:

$$\hat{\phi}_{\mathrm{r}}(\mathcal{S}_t) = \mathrm{Warp}\big(\mathbf{F}_{t \to t-1}(x, y), \ \phi_{\mathrm{r}}(\mathcal{S}_{t-1})\big), \tag{13}$$

where $\hat{\phi}_{\mathrm{r}}(\mathcal{S}_t)$ denotes the warped tokens in the rightmost column of the patch and $\mathbf{F}_{t \to t-1}(x, y)$ represents the down-sampled forward flow. To ensure precise replacement, we compute a validity mask M(x,y) and replace only tokens in valid regions with warped tokens.

As illustrated in Fig. 8, Dynamic-eDiTor achieves stable and reliable monocular scene editing. Our model effectively maintains the temporal consistency, and both STGA and CTP contribute significantly to producing temporally coherent non-rigid appearance edits and semantic local editing.
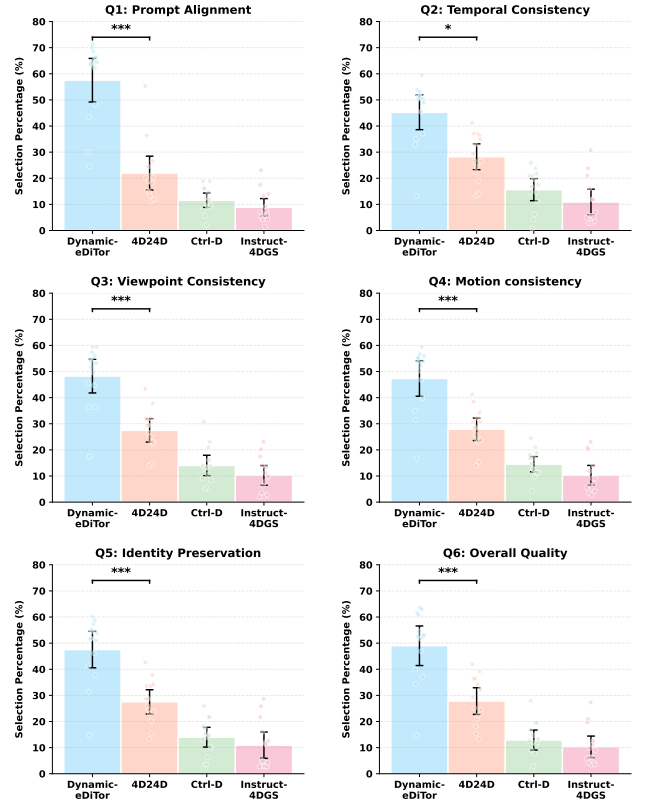


Figure 10. **User Study.** The user study results indicate a human preference for Dynamic-eDiTor, with superior ratings in both consistency and edited-quality categories compared to all baselines
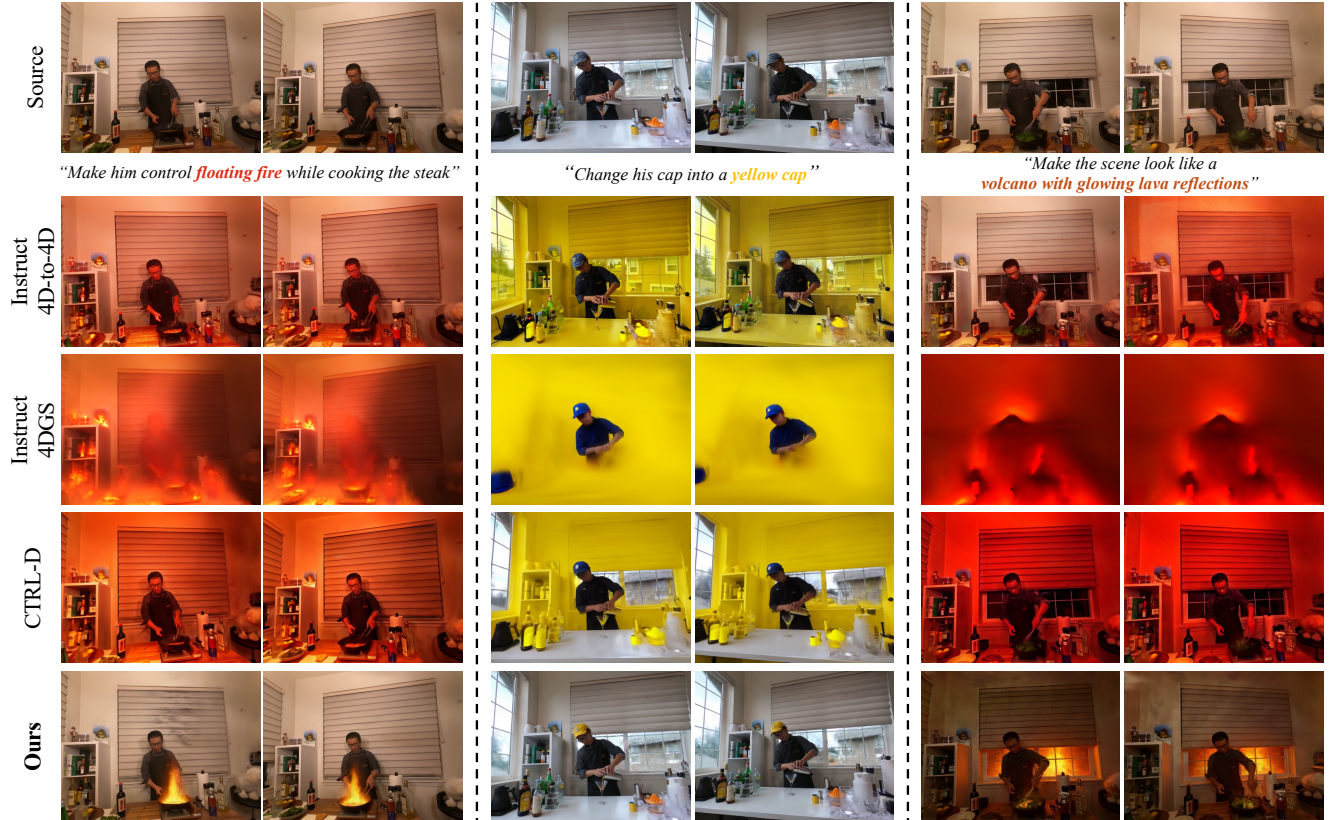
Figure 11. **Qualitative Results.** Dynamic-eDiTor enables higher-quality manipulation of non-rigid content and delivers more complete edits across the 4D scene. The upper row contains the original rendered frames, while the rows beneath show the edited 4DGS results from each baseline approach. Our method (bottom row) demonstrates superior correspondence to the text prompt and achieves strong edit fidelity while maintaining temporal and spatial consistency.

Novel 4D Renderings from Pre-trained 4DGS  ·  Novel 4D Renderings from Edited 4DGS

Source 4D Scene *"Coffee Martini"*  ·  *"Transform the kitchen into vaporwave aesthetic with pink color"*

Source 4D Scene *"Flame Steak"*  ·  *"Change his apron to a yellow apron"*

Source 4D Scene *"Cut Roasted Beef"*  ·  *"Make the kitchen scene appear as though it's underwater"*

Source 4D Scene *"Sear Steak"*  ·  *"Make the steak glow with a magical green light"*

Figure 12. **Qualitative Results.** Dynamic-eDiTor preserves both multi-view and temporal consistency, enabling high-quality text-driven editing of pre-trained 4D Gaussian Splatting. It is capable of performing effective edits across diverse scenes as well as on local objects.



Figure 13. **User study interface and questionnaire.** We illustrate the interface used in our user study. Each participant is first shown the original scene and text prompt, then presented with four edited 4D-rendered video results (A–D) generated by different methods. Participants watch the videos and select the best method for each of the six evaluation criteria: prompt alignment (Q1), temporal consistency (Q2), viewpoint consistency (Q3), motion consistency (Q4), identity preservation (Q5), and overall quality (Q6).