

AutoVFX: Physically Realistic Video Editing from Natural Language Instructions

Hao-Yu Hsu Chih-Hao Lin Albert J. Zhai Hongchi Xia Shenlong Wang
University of Illinois at Urbana-Champaign
<https://haoyuhsu.github.io/autovfx-website/>

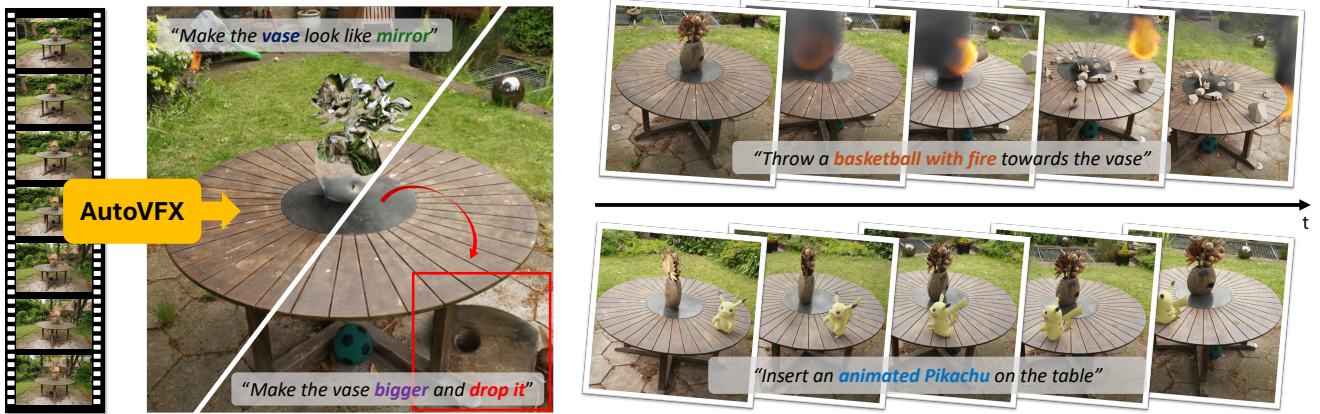


Figure 1. **AutoVFX** takes a video and language instructions as input, and automatically generates programs to produce visual effects and render a new video according to the instructions. It can modify appearance and geometry, enable dynamic interactions, apply particle effects, and even insert animated characters, producing results that are photorealistic, physically-plausible, and easily controllable.

Abstract

Modern visual effects (VFX) software has made it possible for skilled artists to create imagery of virtually anything. However, the creation process remains laborious, complex, and largely inaccessible to everyday users. In this work, we present AutoVFX, a framework that automatically creates realistic and dynamic VFX videos from a single video and natural language instructions. By carefully integrating neural scene modeling, LLM-based code generation, and physical simulation, AutoVFX is able to provide physically-grounded, photorealistic editing effects that can be controlled directly using natural language instructions. We conduct extensive experiments to validate AutoVFX’s efficacy across a diverse spectrum of videos and instructions. Quantitative and qualitative results suggest that AutoVFX outperforms all competing methods by a large margin in generative quality, instruction alignment, editing versatility, and physical plausibility.

1. Introduction

Visual effects (VFX) combine realistic video footage with computer-generated imagery to create novel, photorealistic visuals. Recent advances in graphics, vision, and physi-

cal simulation have made it possible to produce VFX that depict virtually anything—even those that are too costly, time-consuming, dangerous, or impossible to capture in real life. As a result, VFX have become essential in modern filmmaking, ads, simulation, AR/VR, etc. However, the process remains laborious, complex, and expensive, requiring expert skills and professional software [3, 15, 16, 38], making it largely inaccessible to everyday users.

A promising approach to democratizing VFX is to treat it as a generative video editing problem, where raw video and language prompts are used to generate new videos reflecting the original content and given instructions [4, 8, 17, 29, 42, 56, 57, 60, 81, 83, 91]. This method leverages advances in generative modeling, learning from large-scale internet data to produce controllable video. Successes have been seen in deepfake videos, fashion, driving, and robotics [13, 28, 44, 77, 90]. However, this purely data-driven generative editing approach hasn’t yet replaced traditional VFX pipelines due to challenges in achieving guaranteed physical plausibility, precise 3D-aware control, and various special effects.

Another appealing alternative is to build a 3D representation from video input, apply edits like object insertion or texture changes, and then render the final output [10, 11, 20, 22, 23, 32, 45, 47–49, 55, 62, 73, 88, 100, 101]. While this approach aligns well with the VFX pipeline, it is often limited in editing capabilities and still requires manual inter-

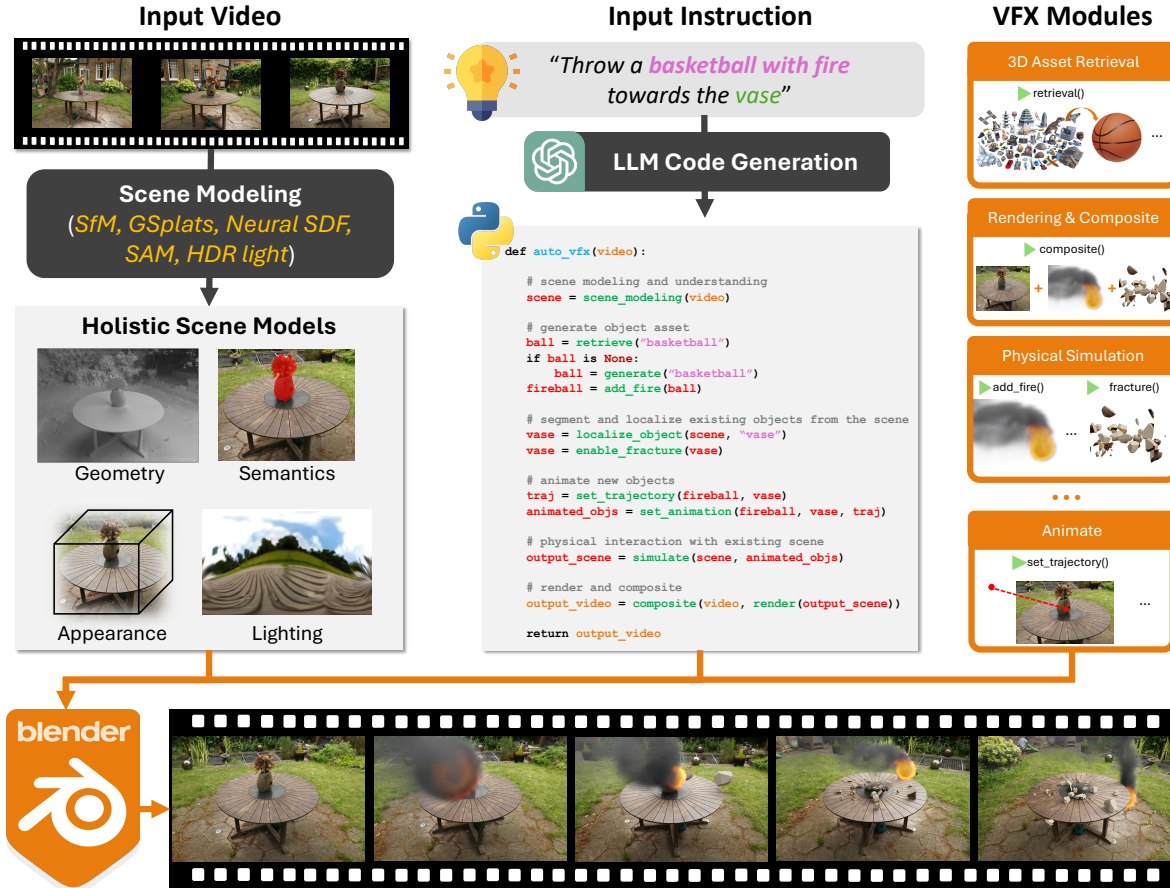


Figure 2. **AutoVFX framework.** Our instruction-guided video editing framework consists of three main modules: (1) **3D Scene Modeling** (left), which integrates 3D reconstruction and scene understanding models; (2) **Program Generation** (middle), where LLMs generate editing programs based on user instructions; and (3) **VFX Modules** (right), which include predefined functions specialized for various editing tasks. These components are integrated with a physically-based simulation and rendering engine (e.g., Blender) to generate the final video.

action with cumbersome interfaces, making it difficult for everyday users. Bridging this gap is essential to make 3D scene editing capable of handling most visual effects while remaining accessible to everyone.

In this work, we present AutoVFX, a framework that automatically creates realistic and dynamic VFX videos from a single video and natural language instructions. At the core of our method is a novel integration of neural scene modeling, LLM-based code generation, and physical simulation. First, we establish a holistic scene model that encodes rich geometry, appearance, and semantics from the input video. This model serves as the foundation for a variety of scene editing, simulation, and rendering capabilities, which we organize into a collection of executable functions. Next, AutoVFX takes simple language editing instructions and converts them into programs using large language models (LLMs). These programs consist of a sequence of calls to our predefined functions. Finally, the generated code is executed, producing a free-viewpoint video that reflects the instructed changes. Fig. 2 illustrates the overall framework.

AutoVFX combines the strengths of generative editing and physical simulation, yet is uniquely set apart from both. Like traditional VFX, AutoVFX produces videos with physics-grounded, controllable, and photorealistic effects. At the same time, similar to generative editing, we support open-world natural language instructions, allowing anyone to edit a video by simply describing the desired effects.

We conduct extensive experiments to validate AutoVFX’s efficacy across a diverse spectrum of videos and instructions. We also perform user studies and qualitative and quantitative comparisons with existing video and scene editing methods. Experimental results suggest AutoVFX outperforms all competing methods by a large margin in generative quality, instruction alignment, editing versatility, and physical plausibility. This demonstrates the effectiveness and convenience of our approach, highlighting its potential as a valuable framework to democratize VFX and pave the way for future integration of even more capabilities to further enhance realism in automatic VFX.

Table 1. **Comparison of existing and proposed methods for visual editing.** Generative editing models lack physical plausibility and precise controllability. Existing physics-based editing methods have complicated interfaces and are limited in their range of editing capacities. Our method, AutoVFX, enjoys a convenient natural language interface while providing the widest range of capabilities.

Method	Input & Output				Editing Capacities						
	Real World Video Editing	Free-Viewpoint Rendering	Editing Interface	Open-world Query	Object Insertion	Object Removal	Object Rearrange	Appearance Change	Animated Objects	Physics Simulation	Particle Effects
Visual Programming [31]	✓	✗	Natural Language	✓	✓	✓	✓	✓	✗	✗	✗
FRESCO [91]	✓	✗	Natural Language	✓	✗	✗	✗	✓	✗	✗	✗
ClimateNeRF [45]	✓	✓	Predefined Scripts	✗	✗	✗	✗	✓	✗	✗	✓
Feature Splatting [64]	✓	✓	Predefined Scripts	✓	✓	✓	✓	✓	✓	✓	✗
GaussianEditor [11]	✓	✓	Graphical	✓	✓	✓	✗	✓	✗	✗	✗
Gaussian Grouping [93]	✓	✓	Graphical	✓	✓	✓	✓	✓	✗	✗	✗
PhysGaussian [86]	✓	✓	Graphical	✗	✗	✗	✗	✗	✗	✓	✗
VR-GS [39]	✓	✓	Graphical	✗	✓	✓	✓	✗	✗	✓	✗
Gaussian Splashing [25]	✓	✓	Graphical	✗	✓	✓	✓	✗	✗	✓	✓
DMRF [62]	✓	✓	Graphical	✗	✓	✗	✗	✗	✓	✓	✓
Instruct-N2N [32]	✓	✓	Natural Language	✓	✗	✗	✗	✓	✗	✗	✗
DGE [10]	✓	✓	Natural Language	✓	✗	✗	✗	✓	✗	✗	✗
Chat-Edit-3D [23]	✓	✓	Natural Language	✓	✓	✓	✓	✓	✗	✗	✗
ChatSim [80]	✓	✓	Natural Language	✗	✓	✓	✓	✓	✓	✗	✗
AutoVFX(Ours)	✓	✓	Natural Language	✓	✓	✓	✓	✓	✓	✓	✓

2. Related Work

Our framework is closely related to several areas, including physical simulation on NeRFs, instruction-guided visual editing, and LLMs for code generation, integrating aspects of all three. Next, we will discuss these areas and highlight and contrast notable works in Tab. 1.

Physical Simulation on NeRFs and 3D Gaussians Integrating physics simulation into NeRFs and 3D Gaussians enables immersive and convincing dynamic effects within captured scenes. Several lines of work have explored various physical interactions, including rigid body object interaction [80, 85], particle physics effects such as flooding and fog [25, 45], elastic deformable objects [98], and plastic objects [39, 61, 86]. The key idea is to enable captured scenes to faithfully interact with new events or entities through physical simulation. However, this is challenging for vanilla neural implicit models, as conventional simulation often requires high-fidelity surface geometry, which is not explicit in these models. Therefore, various approaches seek to extract meshes from NeRFs [14, 39, 58, 61, 87, 95] to facilitate simulation, while others adapt implicit or particle-based simulation so that it can be directly applied to implicit models or Gaussians [24, 25, 43, 86]. AutoVFX explores a hybrid representation where meshes are used for physical interaction and Gaussians are stacked on the mesh surface for rendering, combining the best of both worlds. Another challenge is that some physical interactions require an understanding of physical properties. Various approaches address this through inverse physics [43, 98], common sense knowledge in large foundation models [26, 51, 96], or generative models [98]. Most works on physical simulation, however, are driven by domain-specific scripts rather than natural language instructions, often restricting them to specific physical effects and limiting their user base. AutoVFX seeks to bridge this gap

by using LLMs to convert language instructions into simulation programs and supporting numerous dynamical effects through off-the-shelf simulators.

Instruction-based Visual Editing Recent advancements in visual-language models have made visual editing more accessible by allowing users to edit a wide range of content, such as images, videos, and 3D scenes, using language instructions [7, 26, 32, 42, 91] instead of relying on GUI interactions or script programs [11, 39, 45, 64, 93]. Text- and image-conditioned generative models, particularly diffusion-based approaches [33, 70], have been explored for text-guided image [6, 7, 54, 70, 97] and video [4, 8, 17, 29, 42, 56, 57, 60, 81, 83, 91] editing. Language-embedded NeRFs extend this generative editing capability to 3D scenes [10, 11, 20, 22, 32, 55, 75, 88, 100, 101]. However, mapping text instructions to desired edits in videos and scenes purely through diffusion models can be challenging, particularly when tasks involve complex steps, dynamic interactions, or physics, which can affect instruction alignment, physical plausibility, or realism. To address this, some methods use large language models (LLMs) to break down tasks into subtasks [22, 23, 80] or generate executable programs based on instructions [31, 35, 52]. AutoVFX belongs to this latter category but demonstrates significantly richer capabilities, such as dynamic visual effects, animated objects and physical interaction, compared to these methods.

LLMs for Code Generation The powerful capabilities of large language models (LLMs) have revolutionized code generation based on natural language descriptions. By providing in-context examples, LLMs can generate code snippets in specific formats or syntaxes. Studies such as [2, 9, 21] have explored the effectiveness of LLMs in solving math and code problems. LLM-based code generation has recently been investigated in vision and robotics. Many works adopt LLMs

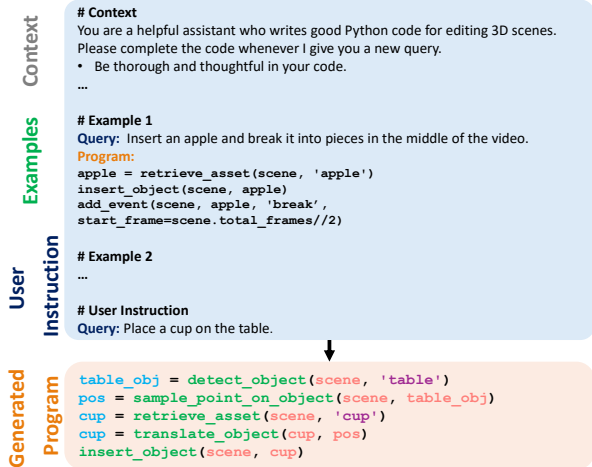


Figure 3. **Program generation.** The LLM generates the editing program through in-context learning. With provided context and examples, it learns to call VFX modules and, given unseen user instructions (blue block), generates the program (orange block).

for decision-making in embodied AI, including tasks like manipulation [36, 46, 74, 89] and navigation [53]. Recently, LLM-based programs for visual content creation have gained attention, with notable progress in 2D understanding and editing [31], driving simulation [23], video generation [52], and procedural 3D scene generation [27, 34, 66, 67, 99]. In our method, we harness GPT-4 to interpret natural language descriptions into executable programs for creating diverse visual effects for generic real-world videos.

3. Text-Driven VFX Creator

AutoVFX takes as input a video and a natural language editing prompt, and outputs an edited free-viewpoint video. The core idea is to uniquely combine the code generation capabilities of LLMs with 3D scene modeling and physics-based simulation techniques. Fig. 2 depicts the overall framework. First, we harness various 3D vision methods to estimate key scene properties from the input video (Sec. 3.1). This lays the foundation for a variety of scene editing, simulation, and rendering capabilities, which we organize into a collection of executable modules (Sec. 3.2, Sec. 3.3). An LLM is used to convert the natural language editing instructions into a program calling these functions (Sec. 3.4). Finally, the generated program is executed, producing a free-viewpoint video that reflects the instructed changes.

3.1. 3D Scene Modeling

Photo-realistic, physics-based VFX creation requires modeling several key properties of the captured scene, namely geometry, appearance, semantics, and lighting. We employ a variety of recent scene understanding models to estimate these properties, including separate models for geometry and appearance in order to achieve both high simulation fidelity

and photorealistic rendering.

Geometry Modeling the 3D geometry of the scene is essential for any sort of object insertion, removal, or simulation. We first run COLMAP [72] to infer the camera poses of each frame. We then capture the geometry in the form of a triangle mesh produced by BakedSDF [92], a multi-view reconstruction method that optimizes a hybrid implicit neural representation that encodes the signed distance field of the scene, and then bakes the representation onto a triangulated mesh. It achieves a desirable balance between surface accuracy, completeness, and efficiency. We choose to use a mesh representation because it can be directly loaded into a standard VFX pipeline, rendered efficiently, and further enables accurate physical simulation. Moreover, it serves as the geometry proxy for object instance extraction.

Appearance We capture the appearance of the scene in two ways. First, we use SuGaR [30], a Gaussian Splatting [40] based novel view synthesis method, to enable free-view rendering. Although SuGaR can provide realistic renderings, it cannot be directly incorporated into physical-based rendering, making it unsuitable for modeling reflective effects on inserted objects or material editing. Thus, we also represent the scene by texturing the BakedSDF mesh, which has lower visual fidelity but can be integrated with physically-based rendering. This textured mesh is used for shadow mapping and encoding multi-bounce effects.

Semantics In many cases, a user would like to perform an edit localized to a specific semantic region of the scene, for example “make the car on fire”. To enable such edits, we use Grounding SAM [50] to perform open-vocabulary instance segmentation and DEVA [12] to associate the instances across frames. To lift video segmentation to 3D, we first un-project each pixel from the 2D segmentation mask into the 3D scene geometry. A voting mechanism is used to determine the visibility of mesh vertices across multiple camera views. By setting a threshold for visibility, we select mesh faces that meet or exceed this threshold. Next, we find 3D Gaussians that are closest to these selected mesh faces and render them to produce an alpha image. We then calculate the average mean Intersection over Union (mIoU) between the rendered alpha image and the original segmentation masks. Finally, we select the mesh faces and 3D Gaussians with the highest mIoU, representing the most accurate 3D segmentation.

Lighting Accurate lighting estimation ensures that all elements within the scene are coherently illuminated. We estimate the environmental lighting of a scene in two ways. For fully captured indoor scenes, such as those in ScanNet++ [94], we unproject the over-saturated image pixels

into space and use majority voting to determine the estimated emitter meshes. These meshes with emissions lights are subsequently imported into the renderer to serve as light sources. For partially captured indoor scenes and outdoor scenes, such as those in MipNeRF360 [5], we use DiffusionLight [59] to inpaint chrome balls in the center of initial frame at multiple exposure levels. An high dynamic range (HDR) map is then generated from these inpainted frames and imported into the renderer as an environmental light.

3.2. Scene Editing and Simulation

The multimodal 3D scene modeling described above paves the way for a wide assortment of editing, simulation, and rendering operations to be performed. We design a suite of intuitive modules that can be seamlessly composed together to provide a rich set of VFX capabilities. We note that this modular framework also allows new capabilities to be easily added via registering new modules. We describe the specific techniques used within each module below.

3D Asset Creation To enable diverse object insertions, we use the Objaverse 1.0 dataset [19] and a high-quality subset from Richdramer [63] with 280k annotated 3D assets. We adopt a twofold approach for 3D asset creation. We first rank 3D assets using Sentence-BERT [68] to match query text and identify the top K candidates, then refine the selection with CLIP [65] based on multi-view renderings to select best aligned asset. For text queries beyond existing descriptions, we use Meshy AI to generate high-quality 3D assets with PBR materials, expanding the range of insertable objects.

Insertion To achieve realistic object insertion, two critical properties must be accurately determined: position and scale. To ensure plausible positioning of objects, we sample the centers of triangles from the supporting mesh that are sufficiently flat to provide accurate support for the objects. For scaling, we utilize GPT-4V [1] models to estimate the real-world dimensions of 3D assets. Detailed prompts for scale estimation are illustrated in the supplementary material.

Removal To effectively remove a specific instance from a scene, we begin by extracting the target objects using semantic modules. We remove these Gaussian points along with their associated mesh faces. For geometric restoration, we employ a planar mesh to cover the exposed area on the bottom. For appearance restoration, we first use LaMa [78] to inpaint the missing regions across all video frames. Then, we fine-tune the current 3D Gaussian Splatting model on the inpainted frames to ensure a 3D-consistent recovery.

Material Editing Accurate material editing ensures 3D assets responding correctly to lighting, shading, and environmental conditions, helping them blend seamlessly with

live-action footage for convincing visual effects. We provide multiple options for material editing, all of which result in modifications to the material nodes of a 3D asset in the renderer. Users can adjust parameters such as metallic, specular, and roughness values of an asset. Users can also modify the overall color of an asset by altering the color intensity in the texture image of an 3D asset. We also support queries by material name, enabling a search across a material database sourced from PolyHaven. The queried material can then be imported and applied to all relevant material nodes.

Physical Simulation Because our scene model is directly compatible with Blender, we can leverage its powerful simulation capabilities by simply calling functions from its simulation library. We use these functions to enable simulations of rigid-body physics and particle effects. Rigid-body physics in Blender is based on the Bullet physics engine [18]. To achieve both accurate and realistic object interactions, we pre-compute the center of mass and convex hull for collision checking of any interactive objects. Particle effects, such as smoke and fire, rely upon mantaflow [79]. We modify the default simulation settings to ensure convincing effects. Further details are provided in the supplementary.

3.3. Scene Rendering and Video Compositing

We use a careful rendering and compositing scheme in order to produce a photorealistic video result. First, we render inserted objects while keeping the background mesh invisible to first-bounces during raytracing but visible to higher-order bounces. This allows the rendered objects to be affected by lighting from the background. Next, we set the background mesh to be visible and render twice: with and without the inserted objects. We use the ratio of the pixel values between these two as an approximation of the objects' effects on the background surfaces. This ratio is multiplied to either a SuGaR rendering or the original video frames, depending on if novel views are desired. Finally, we alpha-blend the inserted objects into the video, using depth maps of the background mesh for occlusion reasoning. If the objects contain fire (emissive transparent elements), we use premultiplied alpha-blending; otherwise, we use straight alpha-blending.

3.4. LLM Integration

AutoVFX aims to enable the creation of VFX directly from natural language instructions, providing a user-friendly interface accessible to anyone. Towards this goal, we integrate our editing modules into an API within an LLM-agent framework, drawing inspiration from recent works such as Code-as-Policies [46] and Visual Programming [31]. Leveraging GPT-4 [1], we prompt the model with in-context examples that pair editing instructions with corresponding programs composed of our predefined editing functions. The LLM

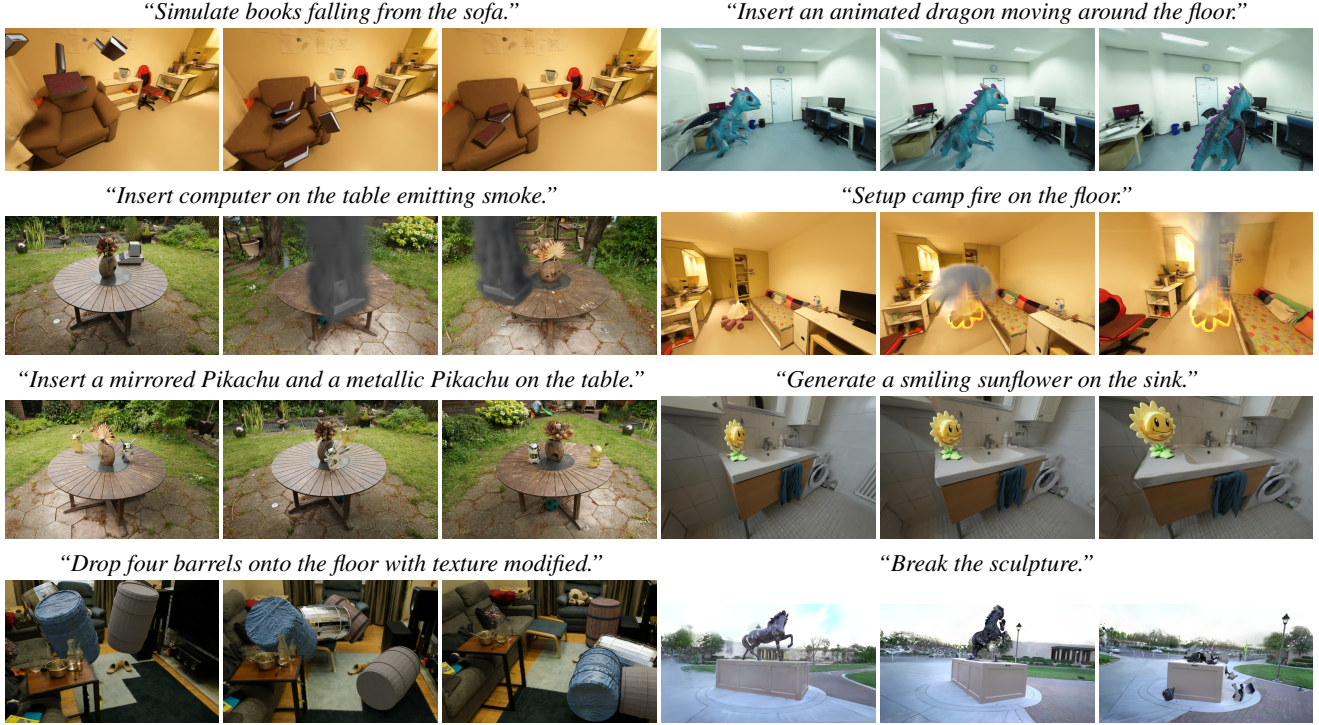


Figure 4. **Dynamic VFX video editing using AutoVFX.** Our approach enables physical interaction, articulated animation, particle effects, insertion of generated 3D assets, material editing, and geometry fracturing.

then generates a program that is directly executed to perform the specified scene edits.

Modular Function Encapsulation Our method encapsulates predefined editing modules into callable and executable functions, which can be combined to form comprehensive programs. Scene objects are represented as Python dictionaries, facilitating straightforward and interpretable edits. Each function’s parameters are fully transparent, allowing users with varying levels of programming expertise to manipulate the process. A full list of the predefined modules is available in the supplementary material.

LLM-driven Program Composition and Execution To ensure GPT-4 effectively composes these functions into executable programs, we design prompts that guide the model. These prompts include examples demonstrating the combination of predefined functions into Python-like scripts that represent the desired scene edits. Once generated, the program is executed within a Python interpreter, triggering the associated simulations or rendering processes to produce the final visual effect. As illustrated in Fig. 3, this approach simplifies the creation of complex visual effects, making it accessible to a broader audience through natural language instructions. The system’s modular design also provides flexibility and scalability, allowing users to customize and extend functionality as needed.

4. Experiments

We evaluate AutoVFX across a diverse set of scenes and editing prompts and provide both qualitative and quantitative comparisons with other related methods.

4.1. Experimental Details

Dataset & Preprocessing We adopt scenes from real-world datasets such as Mip-NeRF360 [5], Tanks & Temples [41], ScanNet++ [94], and Waymo dataset [76] to demonstrate our editing capabilities across diverse scenarios. We use COLMAP [71] to extract camera poses and sparse point clouds from images for GSplat initialization.

Baselines We compare our method with three text-based visual editing methods: Instruct-N2N [32], DGE [10] and FRESCO [91]. Instruct-N2N and DGE both perform edits on 3D scene representations based on text descriptions, with the former utilizing NeRF and the latter relying on 3D Gaussians. FRESCO, on the other hand, translates an input video to align with a target text prompt. For our experiments, we set the guidance scale to 12.5 for DGE to achieve more noticeable edits. Apart from this adjustment, all experimental setups adhere to the default settings of each method to ensure a fair comparison. The qualitative results for Instruct-N2N and DGE are obtained by rendering from the edited 3D representations.

Table 2. Quantitative comparison with other methods. We employ automatic metrics and human evaluation to evaluate the performance. AutoVFX consistently outperforms baseline methods across various metrics.

Method	Semantic Consistency Measures			Multimodal LLM Quality Evaluation				User Study	
	Object Detection	CLIP Similarity	CLIP Directional Similarity	Photo-realism	Text Alignment	Structure Preservation	Overall Quality	Text Alignment	Video Quality
Instruct-N2N [32]	0.343	0.209	0.019	0.402	0.329	0.440	0.043	0.07	0.04
DGE [10]	0.347	0.195	0.278	0.562	0.312	0.619	0.106	0.06	0.03
FRESCO [91]	0.373	0.214	0.009	0.622	0.427	0.632	0.204	0.04	0.02
AutoVFX (Ours)	0.537	0.206	0.419	0.735	0.791	0.749	0.647	0.83	0.90

Implementation Details To render objects that are affected by rigid-body physics, we store the rigid transformations at each timestep and apply them to the 3D Gaussians during rendering. For animating objects based on keypoints, we use Bézier interpolation to produce a smooth trajectory. Additional implementation details regarding the various VFX modules can be found in the supplementary.

4.2. Qualitative evaluation

Qualitative comparison In Fig. 5, we compare the visual quality of static scene editing across different methods. Our approach outperforms baselines in object insertion and manipulation, delivering realistic and accurate edits while preserving scene structure. In contrast, Instruct-N2N struggles with localized editing, FRESCO fails at structural preservation, and DGE, while producing realistic videos, cannot ensure instruction alignment. Additionally, AutoVFX provides richer capabilities, such as precise material editing (“*make it mirror-like*”), accurate object counting (“*drop five basketballs*”), and advanced visual effects (“*make it on fire*”).

Dynamic video simulation We present additional results for dynamic VFX video in Fig. 4. These highlight our method’s ability to generate a wide range of realistic, physically plausible dynamic simulations from text instructions, using modules like rigid body simulations, object animation, smoke and fire, and object fracturing. None of the generative editing baselines support this feature. We also conduct experiments on autonomous driving simulation using the Waymo dataset [76], as shown in Fig. 6. AutoVFX enables both realistic rendering and realistic physical interaction between cars in collision scenarios.

4.3. Quantitative evaluation

We also provide a quantitative evaluation of our method. The evaluation is based on nine metrics categorized into three groups: “*Semantic Consistency Measures*”, “*Multimodal LLM Quality Evaluation*” and “*User Study*”. These metrics collectively provide a comprehensive assessment of the quality and effectiveness of text-guided visual edits. The quantitative results are presented in Table 2.

Semantic Consistency Measures We incorporate the “Multiple Objects” metric from VBench [37] to verify the

presence of objects after editing. This metric assesses whether multiple objects are correctly composed within the edits using a detection module, ensuring that the desired semantic content has been successfully modified. Instead of using GRiT [82] as the detection module, we employ Grounded-SAM [69] for this task. We assess the success rate of visual content editing across all frames and for all possible edits. We also adopt “CLIP Similarity” and “CLIP Directional Similarity” metrics as proposed in DGE [10] for evaluation. “CLIP Similarity” measures the alignment between the text instructions and each edited frame, while “CLIP Directional Similarity” evaluates the temporal consistency of the edits across frames. Both metrics operate in CLIP space. As shown in Table 2, our method significantly outperforms other approaches in object detection score and CLIP directional similarity score, while achieving comparable results in CLIP similarity score. In particular, we improve the object detection score by a large margin, suggesting that our video edits reflect the goal of object-level changes. We notice that CLIP similarity is less discriminative among all methods and conjecture that this might be because global CLIP is not sensitive to capturing local changes, such as the insertion of small objects or dynamics. These results indicate that our method effectively aligns the edited outcomes with the provided text instructions.

Multimodal LLM Quality Evaluation Inspired by [84], we utilize multimodal LLMs as a powerful, interpretable, text-driven model for image quality assessment. Specifically, we prompt GPT-4o to evaluate and compare the “Overall Perceptual Quality” of four different methods based on three criteria: “Text Alignment”, “Photorealism”, and “Structural Preservation”. We also ask GPT-4o to assign a quality score of each criterion to each method, ranging from 0 to 1, with 1 being the highest (detailed prompts could be found in supplementary). From Table 2, our method outperforms other approaches across all four metrics by a significant margin. In particular, AutoVFX creates more realistic videos, preserves structure better, and excels in “Text Alignment” by a most prominent margin. This demonstrates that our video editing produce high-quality and reasonable image edits that fully reflect the desired text instructions, which is desirable in downstream VFX applications.



Figure 5. **Qualitative comparison on static editing.**



Figure 6. **Dynamic simulation of AutoVFX on driving scenes.**

User Study We conduct a user study to evaluate “Text Alignment” and “Overall Realism” of performed edits. To address the potential bias where minimal changes to visual content are perceived as more realistic, we structured the survey as follows: first, users are asked to evaluate which edited videos best aligned with the given text instructions, allowing for multiple selections. In the second question, users are required to choose the most realistic video from the set of choices they previously selected. We collected a total of 36 user samples. For detailed information on the user study methodology, please refer to the supplementary materials. As shown in Table 2, our method receives a higher preference from users in both “Text Alignment” and “Overall Realism” categories, showing that the edits are not

only accurate but also appealing to human judgment.

5. Conclusion

We presented AutoVFX, a system that automatically creates physically-grounded VFX given a monocular video and natural language instructions. AutoVFX combines neural scene modeling, LLM-based code generation, and physical simulation to allow realistic and easily controllable VFX creation. Experimental results demonstrate that AutoVFX outperforms existing scene editing methods based on a variety of practical criteria. We envision that AutoVFX will facilitate both the acceleration and democratization of visual content creation, helping both experienced artists and everyday users create the high-quality VFX that they desire.

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmerschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 5
- [2] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021. 3
- [3] Autodesk, INC. Maya, 2019. <https://autodesk.com/maya>. 1
- [4] Omer Bar-Tal, Dolev Ofri-Amar, Rafail Fridman, Yoni Kasten, and Tali Dekel. Text2live: Text-driven layered image and video editing. In *European Conference on Computer Vision*, pages 707–723. Springer, 2022. 1, 3
- [5] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022. 5, 6
- [6] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023. 3
- [7] Tim Brooks, Aleksander Holynski, and Alexei A. Efros. Instructpix2pix: Learning to follow image editing instructions. In *CVPR*, 2023. 3
- [8] Duygu Ceylan, Chun-Hao Huang, and Niloy J. Mitra. Pix2video: Video editing using image diffusion. In *ICCV*, 2023. 1, 3
- [9] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021. 3
- [10] Minghao Chen, Iro Laina, and Andrea Vedaldi. Dge: Direct gaussian 3d editing by consistent multi-view editing. *arXiv preprint arXiv:2404.18929*, 2024. 1, 3, 6, 7, 8
- [11] Yiwen Chen, Zilong Chen, Chi Zhang, Feng Wang, Xiaofeng Yang, Yikai Wang, Zhongang Cai, Lei Yang, Huaping Liu, and Guosheng Lin. Gaussianeditor: Swift and controllable 3d editing with gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 21476–21485, 2024. 1, 3
- [12] Ho Kei Cheng, Seoung Wug Oh, Brian Price, Alexander Schwing, and Joon-Young Lee. Tracking anything with decoupled video segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1316–1326, 2023. 4
- [13] Yisol Choi, Sangkyung Kwak, Kyungmin Lee, Hyungwon Choi, and Jinwoo Shin. Improving diffusion models for virtual try-on. *arXiv preprint arXiv:2403.05139*, 2024. 1
- [14] Chong Bao and Bangbang Yang, Zeng Junyi, Bao Hujun, Zhang Yinda, Cui Zhaopeng, and Zhang Guofeng. Neumesh: Learning disentangled neural mesh-based implicit field for geometry and texture editing. In *European Conference on Computer Vision (ECCV)*, 2022. 3
- [15] Mark Christiansen. *Adobe After Effects CC Visual Effects and Compositing Studio Techniques*. Adobe Press, 2013. 1
- [16] Blender Online Community. *Blender - a 3D modelling and rendering package*. Blender Foundation, Stichting Blender Foundation, Amsterdam, 2018. 1
- [17] Yuren Cong, Mengmeng Xu, Christian Simon, Shoufa Chen, Jiawei Ren, Yanping Xie, Juan-Manuel Perez-Rua, Bodo Rosenhahn, Tao Xiang, and Sen He. Flatten: optical flow-guided attention for consistent text-to-video editing. *arXiv preprint arXiv:2310.05922*, 2023. 1, 3
- [18] Erwin Coumans. Bullet physics simulation. In *ACM SIGGRAPH 2015 Courses*, 2015. 5
- [19] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13142–13153, 2023. 5
- [20] Jiahua Dong and Yu-Xiong Wang. Vica-nerf: View-consistency-aware 3d editing of neural radiance fields. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 1, 3
- [21] Iddo Drori, Sarah Zhang, Reece Shuttleworth, Leonard Tang, Albert Lu, Elizabeth Ke, Kevin Liu, Linda Chen, Sunny Tran, Newman Cheng, et al. A neural network solves, explains, and generates university math problems by program synthesis and few-shot learning at human level. *Proceedings of the National Academy of Sciences*, 119(32):e2123433119, 2022. 3
- [22] Jiemin Fang, Junjie Wang, Xiaopeng Zhang, Lingxi Xie, and Qi Tian. Gaussianeditor: Editing 3d gaussians delicately with text instructions. In *CVPR*, 2024. 1, 3
- [23] Shuangkang Fang, Yufeng Wang, Yi-Hsuan Tsai, Yi Yang, Wenrui Ding, Shuchang Zhou, and Ming-Hsuan Yang. Chat-edit-3d: Interactive 3d scene editing via text prompts. *arXiv preprint arXiv:2407.06842*, 2024. 1, 3, 4
- [24] Yutao Feng, Yintong Shang, Xuan Li, Tianjia Shao, Chenfanfu Jiang, and Yin Yang. Pie-nerf: Physics-based interactive elastodynamics with nerf, 2023. 3
- [25] Yutao Feng, Xiang Feng, Yintong Shang, Ying Jiang, Chang Yu, Zeshun Zong, Tianjia Shao, Hongzhi Wu, Kun Zhou, Chenfanfu Jiang, and Yin Yang. Gaussian splashing: Unified particles for versatile motion synthesis and rendering. *arXiv preprint arXiv:2401.15318*, 2024. 3
- [26] Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith, Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but not perceive. *arXiv preprint arXiv:2404.12390*, 2024. 3
- [27] Gege Gao, Weiyang Liu, Anpei Chen, Andreas Geiger, and Bernhard Schölkopf. Graphdreamer: Compositional 3d scene synthesis from scene graphs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21295–21304, 2024. 4

- [28] Ruiyuan Gao, Kai Chen, Enze Xie, Lanqing Hong, Zhenguo Li, Dit-Yan Yeung, and Qiang Xu. Magicdrive: Street view generation with diverse 3d geometry control. *arXiv preprint arXiv:2310.02601*, 2023. 1
- [29] Michal Geyer, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Tokenflow: Consistent diffusion features for consistent video editing. *arXiv preprint arXiv:2307.10373*, 2023. 1, 3
- [30] Antoine Guédon and Vincent Lepetit. Sugar: Surface-aligned gaussian splatting for efficient 3d mesh reconstruction and high-quality mesh rendering. *arXiv preprint arXiv:2311.12775*, 2023. 4
- [31] Tanmay Gupta and Aniruddha Kembhavi. Visual programming: Compositional visual reasoning without training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14953–14962, 2023. 3, 4, 5
- [32] Ayaan Haque, Matthew Tancik, Alexei Efros, Aleksander Holynski, and Angjoo Kanazawa. Instruct-nerf2nerf: Editing 3d scenes with instructions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023. 1, 3, 6, 7, 8
- [33] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *arXiv preprint arXiv:2006.11239*, 2020. 3
- [34] Ziniu Hu, Ahmet Iscen, Aashi Jain, Thomas Kipf, Yisong Yue, David A Ross, Cordelia Schmid, and Alireza Fathi. Scenecraft: An llm agent for synthesizing 3d scenes as blender code. In *Forty-first International Conference on Machine Learning*, 2024. 4
- [35] Ian Huang, Guandao Yang, and Leonidas Guibas. Blender-alchemy: Editing 3d graphics with vision-language models. In *European Conference on Computer Vision*, pages 297–314. Springer, 2025. 3
- [36] Wenlong Huang, Chen Wang, Ruohan Zhang, Yunzhu Li, Jiajun Wu, and Li Fei-Fei. Voxposer: Composable 3d value maps for robotic manipulation with language models. *arXiv preprint arXiv:2307.05973*, 2023. 4
- [37] Ziqi Huang, Yanan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, Yaohui Wang, Xinyuan Chen, Limin Wang, Dahua Lin, Yu Qiao, and Ziwei Liu. VBench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 7
- [38] Side Effects Software Inc. *SideFX Houdini FX*. SideFX, 2018. 1
- [39] Ying Jiang, Chang Yu, Tianyi Xie, Xuan Li, Yutao Feng, Huamin Wang, Minchen Li, Henry Lau, Feng Gao, Yin Yang, and Chenfanfu Jiang. Vr-gs: A physical dynamics-aware interactive gaussian splatting system in virtual reality. *arXiv preprint arXiv:2401.16663*, 2024. 3
- [40] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023. 4
- [41] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics (ToG)*, 36(4):1–13, 2017. 6
- [42] Max Ku, Cong Wei, Weiming Ren, Huan Yang, and Wenhui Chen. Anyv2v: A plug-and-play framework for any video-to-video editing tasks. *arXiv preprint arXiv:2403.14468*, 2024. 1, 3
- [43] Xuan Li, Yi-Ling Qiao, Peter Yichen Chen, Krishna Murthy Jatavallabhula, Ming Lin, Chenfanfu Jiang, and Chuang Gan. PAC-nerf: Physics augmented continuum neural radiance fields for geometry-agnostic system identification. In *The Eleventh International Conference on Learning Representations*, 2023. 3
- [44] Xiaofan Li, Yifu Zhang, and Xiaoqing Ye. Drivingdiffusion: Layout-guided multi-view driving scene video generation with latent diffusion model. *arXiv preprint arXiv:2310.07771*, 2023. 1
- [45] Yuan Li, Zhi-Hao Lin, David Forsyth, Jia-Bin Huang, and Shenlong Wang. Climatenerf: Extreme weather synthesis in neural radiance field. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 1, 3
- [46] Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as policies: Language model programs for embodied control. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9493–9500. IEEE, 2023. 4, 5
- [47] Ruofan Liang, Zan Gojcic, Merlin Nimier-David, David Acuna, Nandita Vijaykumar, Sanja Fidler, and Zian Wang. Photorealistic object insertion with diffusion-guided inverse rendering. *arXiv preprint arXiv:2408.09702*, 2024. 1
- [48] Zhi-Hao Lin, Bohan Liu, Yi-Ting Chen, David Forsyth, Jia-Bin Huang, Anand Bhattad, and Shenlong Wang. Urbanir: Large-scale urban scene inverse rendering from a single video. *arXiv preprint arXiv:2306.09349*, 2023.
- [49] Zhi-Hao Lin, Jia-Bin Huang, Zhengqin Li, Zhao Dong, Christian Richardt, Tuotuo Li, Michael Zollhöfer, Johannes Kopf, Shenlong Wang, and Changil Kim. Iris: Inverse rendering of indoor scenes from low dynamic range images. *arXiv preprint arXiv:2401.12977*, 2024. 1
- [50] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 4
- [51] Shaowei Liu, Zhongzheng Ren, Saurabh Gupta, and Shenlong Wang. Physgen: Rigid-body physics-grounded image-to-video generation. In *European Conference on Computer Vision*, pages 360–378. Springer, 2025. 3
- [52] Jiayi Lv, Yi Huang, Mingfu Yan, Jiancheng Huang, Jianzhuang Liu, Yifan Liu, Yafei Wen, Xiaoxin Chen, and Shifeng Chen. Gpt4motion: Scripting physical motions in text-to-video generation via blender-oriented gpt planning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 1430–1440, 2024. 3, 4
- [53] Yunsheng Ma, Can Cui, Xu Cao, Wenqian Ye, Peiran Liu, Juanwu Lu, Amr Abdelraouf, Rohit Gupta, Kyungtae Han,

- Aniket Bera, et al. Lampilot: An open benchmark dataset for autonomous driving with language model programs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15141–15151, 2024. [4](#)
- [54] 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. In *International Conference on Learning Representations*, 2022. [3](#)
- [55] Ashkan Mirzaei, Tristan Aumentado-Armstrong, Marcus A Brubaker, Jonathan Kelly, Alex Levinstein, Konstantinos G Derpanis, and Igor Gilitschenski. Watch your steps: Local image and scene editing by text instructions. *arXiv preprint arXiv:2308.08947*, 2023. [1, 3](#)
- [56] Chong Mou, Mingdeng Cao, Xintao Wang, Zhaoyang Zhang, Ying Shan, and Jian Zhang. Revideo: Remake a video with motion and content control, 2024. [1, 3](#)
- [57] Hao Ouyang, Qiuyu Wang, Yuxi Xiao, Qingyan Bai, Juntao Zhang, Kecheng Zheng, Xiaowei Zhou, Qifeng Chen, and Yujun Shen. Codef: Content deformation fields for temporally consistent video processing. *arXiv preprint arXiv:2308.07926*, 2023. [1, 3](#)
- [58] Yicong Peng, Yichao Yan, Shengqi Liu, Yuhao Cheng, Shanyan Guan, Bowen Pan, Guangtao Zhai, and Xiaokang Yang. Cagenerf: Cage-based neural radiance field for generalized 3d deformation and animation. *Advances in Neural Information Processing Systems*, 35:31402–31415, 2022. [3](#)
- [59] Pakkapon Phongthawee, Worameth Chinchuthakun, Nontaphat Sinsunthithet, Amit Raj, Varun Jampani, Pramook Khungurn, and Supasorn Suwajanakorn. Diffusionlight: Light probes for free by painting a chrome ball. *arXiv preprint arXiv:2312.09168*, 2023. [5](#)
- [60] Chenyang Qi, Xiaodong Cun, Yong Zhang, Chenyang Lei, Xintao Wang, Ying Shan, and Qifeng Chen. Fatezero: Fusing attentions for zero-shot text-based video editing. *arXiv:2303.09535*, 2023. [1, 3](#)
- [61] Yi-Ling Qiao, Alexander Gao, and Ming C. Lin. Neu-physics: Editable neural geometry and physics from monocular videos. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2022. [3](#)
- [62] Yi-Ling Qiao, Alexander Gao, Yiran Xu, Yue Feng, Jia-Bin Huang, and Ming C. Lin. Dynamic mesh-aware radiance fields. *ICCV*, 2023. [1, 3](#)
- [63] Lingteng Qiu, Guanying Chen, Xiaodong Gu, Qi Zuo, Mutian Xu, Yushuang Wu, Weihao Yuan, Zilong Dong, Liefeng Bo, and Xiaoguang Han. Richdreamer: A generalizable normal-depth diffusion model for detail richness in text-to-3d. *arXiv preprint arXiv:2311.16918*, 2023. [5](#)
- [64] Ri-Zhao Qiu, Ge Yang, Weijia Zeng, and Xiaolong Wang. Language-driven physics-based scene synthesis and editing via feature splatting. In *European Conference on Computer Vision (ECCV)*, 2024. [3](#)
- [65] 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. [5](#)
- [66] Alexander Raistrick, Lahav Lipson, Zeyu Ma, Lingjie Mei, Mingzhe Wang, Yiming Zuo, Karhan Kayan, Hongyu Wen, Beining Han, Yihan Wang, et al. Infinite photorealistic worlds using procedural generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12630–12641, 2023. [4](#)
- [67] Alexander Raistrick, Lingjie Mei, Karhan Kayan, David Yan, Yiming Zuo, Beining Han, Hongyu Wen, Meenal Parakh, Stamatis Alexandropoulos, Lahav Lipson, et al. Infinigen indoors: Photorealistic indoor scenes using procedural generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21783–21794, 2024. [4](#)
- [68] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019. [5](#)
- [69] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, Zhaoyang Zeng, Hao Zhang, Feng Li, Jie Yang, Hongyang Li, Qing Jiang, and Lei Zhang. Grounded sam: Assembling open-world models for diverse visual tasks, 2024. [7](#)
- [70] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021. [3](#)
- [71] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. [6](#)
- [72] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016. [4](#)
- [73] Yuan Shen, Bhargav Chandaka, Zhi-Hao Lin, Albert Zhai, Hang Cui, David Forsyth, and Shenlong Wang. Sim-on-wheels: Physical world in the loop simulation for self-driving. *IEEE Robotics and Automation Letters*, 2023. [1](#)
- [74] Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. Progprompt: Generating situated robot task plans using large language models. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11523–11530. IEEE, 2023. [4](#)
- [75] Liangchen Song, Liangliang Cao, Jiatao Gu, Yifan Jiang, Junsong Yuan, and Hao Tang. Efficient-nerf2nerf: Streamlining text-driven 3d editing with multiview correspondence-enhanced diffusion models. *arXiv preprint arXiv:2312.08563*, 2023. [3](#)
- [76] Pei Sun, Henrik Kretschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2446–2454, 2020. [6, 7](#)
- [77] Shuo Sun, Zekai Gu, Tianchen Sun, Jiawei Sun, Chengran Yuan, Yuhang Han, Dongen Li, and Marcelo H Ang.

- Drivescenenet: Generating diverse and realistic driving scenarios from scratch. *IEEE Robotics and Automation Letters*, 2024. 1
- [78] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky. Resolution-robust large mask inpainting with fourier convolutions. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2149–2159, 2022. 5
- [79] Nils Thuerey and Tobias Pfaff. MantaFlow, 2018. <http://mantaflow.com>. 5
- [80] Yuxi Wei, Zi Wang, Yifan Lu, Chenxin Xu, Changxing Liu, Hao Zhao, Siheng Chen, and Yanfeng Wang. Editable scene simulation for autonomous driving via collaborative llm-agents. *arXiv preprint arXiv:2402.05746*, 2024. 3
- [81] Bichen Wu, Ching-Yao Chuang, Xiaoyan Wang, Yichen Jia, Kapil Krishnakumar, Tong Xiao, Feng Liang, Licheng Yu, and Peter Vajda. Fairy: Fast parallelized instruction-guided video-to-video synthesis. *arXiv preprint arXiv:2312.13834*, 2023. 1, 3
- [82] Jialian Wu, Jianfeng Wang, Zhengyuan Yang, Zhe Gan, Zicheng Liu, Junsong Yuan, and Lijuan Wang. Grit: A generative region-to-text transformer for object understanding. *arXiv preprint arXiv:2212.00280*, 2022. 7
- [83] 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. 1, 3
- [84] Tianhe Wu, Kede Ma, Jie Liang, Yujiu Yang, and Lei Zhang. A comprehensive study of multimodal large language models for image quality assessment. *arXiv preprint arXiv:2403.10854v3*, 2024. 7
- [85] Hongchi Xia, Zhi-Hao Lin, Wei-Chiu Ma, and Shenlong Wang. Video2game: Real-time interactive realistic and browser-compatible environment from a single video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4578–4588, 2024. 3
- [86] Tianyi Xie, Zeshun Zong, Yuxing Qiu, Xuan Li, Yutao Feng, Yin Yang, and Chenfanfu Jiang. Physgaussian: Physics-integrated 3d gaussians for generative dynamics. *arXiv preprint arXiv:2311.12198*, 2023. 3
- [87] Tianhan Xu and Tatsuya Harada. Deforming radiance fields with cages. In *ECCV*, 2022. 3
- [88] Teng Xu, Jiamin Chen, Peng Chen, Youjia Zhang, Junqing Yu, and Wei Yang. Tiger: Text-instructed 3d gaussian retrieval and coherent editing. *arXiv preprint arXiv:2405.14455*, 2024. 1, 3
- [89] Jinggang Yang, Yuhao Dong, Shuai Liu, Bo Li, Ziyue Wang, Chencheng Jiang, Haoran Tan, Jiamu Kang, Yuanhan Zhang, Kaiyang Zhou, et al. Octopus: Embodied vision-language programmer from environmental feedback. *arXiv preprint arXiv:2310.08588*, 2023. 4
- [90] Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan Tompson, Dale Schuurmans, and Pieter Abbeel. Learning interactive real-world simulators. *arXiv preprint arXiv:2310.06114*, 2023. 1
- [91] Shuai Yang, Yifan Zhou, Ziwei Liu, , and Chen Change Loy. Fresco: Spatial-temporal correspondence for zero-shot video translation. In *CVPR*, 2024. 1, 3, 6, 7, 8
- [92] Lior Yariv, Peter Hedman, Christian Reiser, Dor Verbin, Pratul P Srinivasan, Richard Szeliski, Jonathan T Barron, and Ben Mildenhall. Baked sdf: Meshing neural sdfs for real-time view synthesis. In *ACM SIGGRAPH 2023 Conference Proceedings*, pages 1–9, 2023. 4
- [93] Mingqiao Ye, Martin Danelljan, Fisher Yu, and Lei Ke. Gaussian grouping: Segment and edit anything in 3d scenes. In *ECCV*, 2024. 3
- [94] Chandan Yeshwanth, Yueh-Cheng Liu, Matthias Nießner, and Angela Dai. Scannet++: A high-fidelity dataset of 3d indoor scenes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12–22, 2023. 4, 6
- [95] Yu-Jie Yuan, Yang-Tian Sun, Yu-Kun Lai, Yuewen Ma, Rongfei Jia, and Lin Gao. Nerf-editing: geometry editing of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18353–18364, 2022. 3
- [96] Albert J Zhai, Yuan Shen, Emily Y Chen, Gloria X Wang, Xinlei Wang, Sheng Wang, Kaiyu Guan, and Shenlong Wang. Physical property understanding from language-embedded feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 28296–28305, 2024. 3
- [97] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models, 2023. 3
- [98] Tianyuan Zhang, Hong-Xing Yu, Rundui Wu, Brandon Y. Feng, Changxi Zheng, Noah Snavely, Jiajun Wu, and William T. Freeman. PhysDreamer: Physics-based interaction with 3d objects via video generation. *arxiv*, 2024. 3
- [99] Mengqi Zhou, Jun Hou, Chuanchen Luo, Yuxi Wang, Zhaoxiang Zhang, and Junran Peng. Scenex: Procedural controllable large-scale scene generation via large-language models. *arXiv preprint arXiv:2403.15698*, 2024. 4
- [100] Jingyu Zhuang, Chen Wang, Liang Lin, Lingjie Liu, and Guanbin Li. Dreameditor: Text-driven 3d scene editing with neural fields. In *SIGGRAPH Asia 2023 Conference Papers*, pages 1–10, 2023. 1, 3
- [101] Jingyu Zhuang, Di Kang, Yan-Pei Cao, Guanbin Li, Liang Lin, and Ying Shan. Tip-editor: An accurate 3d editor following both text-prompts and image-prompts. *arXiv preprint arXiv:2401.14828*, 2024. 1, 3