

VeGaS: Video Gaussian Splatting

Weronika Smolak-Dyżewska*, Dawid Malarz*, Kornel Howil*,
 Jan Kaczmarczyk, Marcin Mazur, Przemysław Spurek
 Jagiellonian University
 Faculty of Mathematics and Computer Science
 weronika.smolak@doctoral.uj.edu.pl

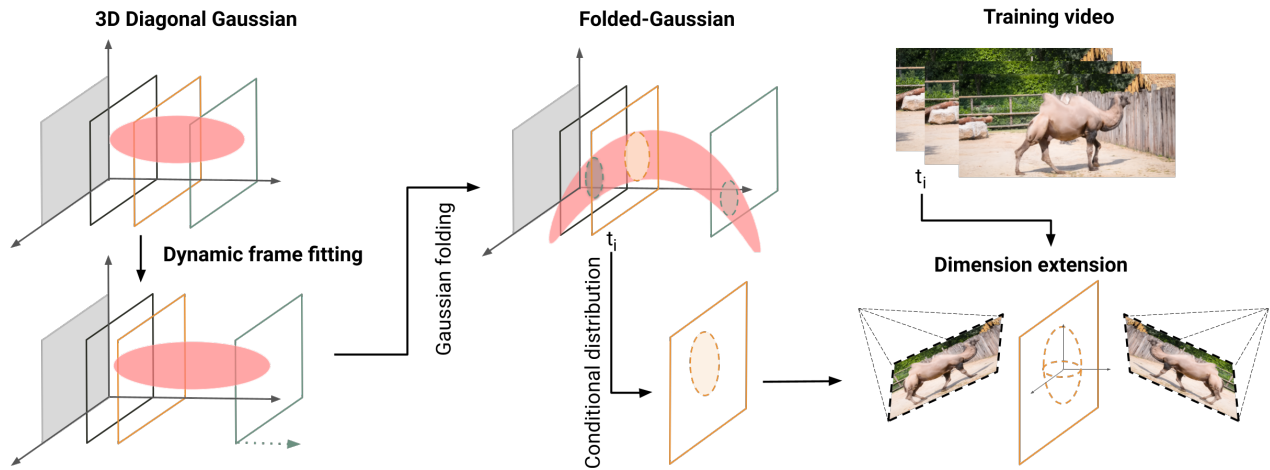


Figure 1. Graphical summary of our Video Gaussian Splatting (VeGaS) model. The initial step involves the use of diagonal 3D Gaussians and frames with equal distances. Then, dynamic frame fitting and Gaussian folding are employed to approximate nonlinear structures within a video stream. Each frame is modeled by 2D Gaussians obtained by conditioning of 3D Folded-Gaussians at frame occurrence time t_i . This representation allows for the creation of high-quality renderings of video data and facilitates a wide range of modifications.

Abstract

Implicit Neural Representations (INRs) employ neural networks to approximate discrete data as continuous functions. In the context of video data, such models can be utilized to transform the coordinates of pixel locations along with frame occurrence times (or indices) into RGB color values. Although INRs facilitate effective compression, they are unsuitable for editing purposes. One potential solution is to use a 3D Gaussian Splatting (3DGS) based model, such as the Video Gaussian Representation (VGR), which is capable of encoding video as a multitude of 3D Gaussians and is applicable for numerous video processing operations, including editing. Nevertheless, in this case, the capacity for modification is constrained to a limited set of basic transformations. To address this issue, we introduce

the Video Gaussian Splatting (VeGaS) model, which enables realistic modifications of video data. To construct VeGaS, we propose a novel family of Folded-Gaussian distributions designed to capture nonlinear dynamics in a video stream and model consecutive frames by 2D Gaussians obtained as respective conditional distributions. Our experiments demonstrate that VeGaS outperforms state-of-the-art solutions in frame reconstruction tasks and allows realistic modifications of video data. The code is available at: <https://github.com/gmum/VeGaS>.

1. Introduction

Implicit Neural Representations (INRs) [27] employs neural networks to describe discrete data as a smooth, continuous function. They have emerged as a promising method for continuously encoding a variety of signals, including

*These authors contributed equally to this work

images [14], videos [5], audio [30], and 3D shapes [24]. INRs are frequently utilized in the context of 2D imagery by training networks that map pixel coordinates to RGB color values, thus encoding a structure of images in neural network weights. This approach offers several benefits, including applications in compression [5], hyper-resolution [14], or as an integral part of generative models [11, 28].

On the other hand, the 3D Gaussian Splatting (3DGS) framework [13], initially proposed for the modeling of 3D scenes, has recently been adapted with 2D images. In particular, the GaussianImage method [40] has demonstrated promising results in image reconstruction by efficiently encoding images in the 2D space, with a strong focus on model efficiency and reduced training time. Furthermore, the MiraGe representation [33] has demonstrated the feasibility of generating realistic modifications of 2D images.

Similarly to 2D images, INR produces a continuous representation of videos [5]. In such a case, neural networks transform the pixel coordinates and time frames into RGB color. Such models provide good reconstruction quality and compression ratios. Unfortunately, INR ultimately failed with the editing of videos. To solve such a problem, we can use the Gaussian Splatting solution. Video Gaussian Representation (VGR) [29] uses Gaussians in the canonical position and deformation function, which transfers such a Gaussian to each time frame. This model is capable of handling a variety of video processing tasks, such as video editing. Nevertheless, these changes are restricted to linear transformations and translations.

This paper introduces the Video Gaussian Splatting (VeGaS) model, which presents evidence that the 3DGS approach can be adapted for use with 2D video data. In particular, video frames are treated as parallel planes within a 3D space, and a 3D Gaussian Splatting is employed to model the transitions observed between the content of subsequent frames. By conditioning 3D Gaussian components at specified time points, 2D Gaussians are tailored to the selected frames. It is crucial to emphasize that our solution surpasses the classical Gaussian Splatting model by enhancing its capacity to integrate intricate distributions, thus facilitating more exact modeling of swift alterations in video sequences. In particular, we introduce Folded-Gaussians, a family of functions that model nonlinear structures and produce classical 2D Gaussian distributions after conditioning. It should be noted that the employment of the 3DGS framework for video modeling allows for the utilization of both extensive Gaussians to represent the background, which remains largely static over time, and brief Gaussians to represent elements present in only a few frames. Furthermore, VeGaS employs a MiraGe-based representation to model individual frames, which allows one to modify both the entire video and selected frames, resulting in high-quality renderings, as shown in Figure 2.

The following represents a comprehensive account of our significant contributions:

- we introduce Folded-Gaussians, a novel family of functions that model nonlinear structures and can be readily incorporated into the 3D Gaussian Splatting framework,
- we propose the VeGaS model, which allows for the processing of 2D video data using the Folded-Gaussians,
- we conduct experiments that demonstrate the superiority of VeGaS for reconstruction tasks and show its efficiency in producing realistic modifications of video data.

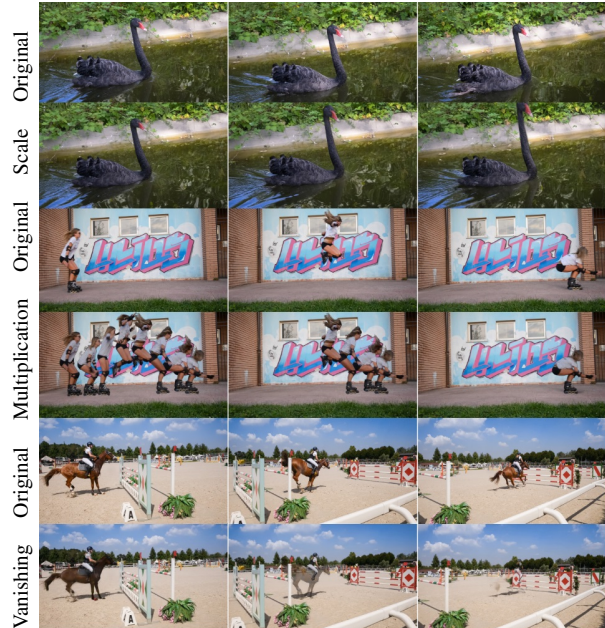


Figure 2. *Video edition*. Note that VeGaS enables modification of selected objects on a global scale, including operations such as multiplication and scaling. The model was trained on the DAVIS dataset [25].

2. Related Works

The decomposition of videos into layered representations enables the utilization of sophisticated video editing techniques. In [12], the authors decompose an image into textured layers and learn a corresponding deformation field, which allows for efficient video editing. The method outlined in [38] involves the partitioning of videos into discrete motion groups, with each group being driven by an MLP-based representation. INVE [9] employs a bidirectional warping field to facilitate extensive video tracking and editing over extended periods of time. In [2], the authors propose an improvement to the rendering of lighting and color details. This is achieved by incorporating additional layers and residual color maps, which serve to enhance the representation of illumination effects in the video. CoDeF [23] employs a multi-resolution hash grid and a shallow MLP to model frame-by-frame deformations relative to a canon-

ical image. This approach enables editing in the canonical space, with changes effectively propagated across the entire video. A comparable representation is utilized in GenDeF [35] for the generation of controllable videos.

The generative potential of latent diffusion models has been harnessed in various research endeavors within the context of video editing [26]. In [41], the authors integrate control signals into the network during video reconstruction, thereby guiding the editing process. A related technique involves frame interpolation to generate edited videos from specifically edited keyframes [22], while another method employs a token merging approach to incorporate control signals during [18]. Moreover, some works investigate inversion techniques for video editing [17, 42].

3D Gaussian Splatting (3DGS) [13] models 3D static scenes using a set of Gaussian components. Recently, numerous generalizations have been proposed for the representation of dynamic scenes. In [21], the authors employ a multiview dynamic dataset coupled with an incremental, frame-based strategy. Nevertheless, this method does not account for inter-frame correlation and requires a considerable amount of storage space for extended sequences. The approach presented in [15, 37] employs an MLP to represent temporal alterations in Gaussians. In contrast, in [36] the authors utilize an MLP in conjunction with a decomposed neural voxel encoding technique to enhance training and storage effectiveness. In [20], dynamic scenes are divided into dynamic and static segments, which are optimized independently and then combined to facilitate decoupling. Other research has sought to enhance the reconstruction of dynamic scenes by incorporating external priors. For instance, diffusion priors have been shown to serve as effective regularization terms in the optimization process [39]. In [6], the authors propose 4DRotorGS which employs a four-dimensional Gaussian, with the fourth dimension dedicated to time.

Furthermore, 3DGS was utilized to modify scene geometry based on the underlying meshes. In [7], the positioning of 3D Gaussians on an explicit mesh enables the utilization of mesh rendering to facilitate adaptive refinement. This method is contingent upon the use of an extracted mesh as a proxy; therefore, it is inoperable in the event of a failure in the mesh extraction process. On the other hand, in [8] explicit meshes are derived from 3DGS representations through the regularization of Gaussians over surfaces. This process, which involves a significant optimization and refinement phase, is particularly resource-intensive. Another example is given in [10], where sparse control points are used for 3D scene dynamics. However, this approach encounters difficulties with extensive edit movements and requires precise static node selection. In turn, GaMeS [32] integrates 3DGS with mesh extraction, though this approach is only effective for static scenes. In contrast, D-MiSo [31]

is a mesh-based approach, specifically designed for dynamic scenes, which employs a simple 3DGS technique pipeline to allow real-time editing of dynamic scenes.

In [29], the authors introduce the Video Gaussian Representation (VGR), which employs 3D Gaussian Splatting to model video data. This approach is closely related to ours and therefore represents the most reasonable baseline for the VeGaS model. VGR uses Gaussians in the canonical position, transferring them further to each frame occurrence time, and a deformation function. This model is capable of handling a variety of video processing tasks, such as video editing. Nevertheless, the possible changes are restricted to linear transformations and translations, which represents a limitation relative to our proposed solution.

3. Folded-Gaussians

In this section, we introduce the concept of a Folded-Gaussian distribution, which can be seen as a generalization of a classical Gaussian distribution in order to capture nonlinear structures. It should be noted that Folded-Gaussians constitute a novel family of distributions that we further employ to represent video data (see the next section). Accordingly, in order to emphasize this relationship, a terminology based on the concept of a space-time variable is employed, given that each video can be regarded as a sequence of successive frames occurring at discrete time points. For the reader’s convenience, we begin with a simple two-dimensional toy example, before extending our presentation to the multidimensional case.

Toy Example in \mathbb{R}^2 Our toy example starts with a two-dimensional Gaussian distribution $\mathcal{N}(m, \Sigma)$ of a space-time random variable $x = (s, t) \in \mathbb{R} \times \mathbb{R}$, which is given by a mean vector $m = (m_s, m_t)$ and a covariance matrix

$$\Sigma = \begin{bmatrix} \sigma_s^2 & 0 \\ 0 & \sigma_t^2 \end{bmatrix}. \quad (1)$$

In this case, the density function is defined by the following formula:

$$N(m, \Sigma)(x) = N(m_s, \sigma_s^2)(s) \cdot N(m_t, \sigma_t^2)(t), \quad (2)$$

where

$$N(m, \sigma^2)(x) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{|x - m|^2}{2\sigma^2}\right). \quad (3)$$

Using such a distribution, we can model ellipses, which can be considered as a simple linear structure spanned along the coordinate axes. Therefore, we propose a generalization of a classical 2D Gaussian that allows us to deal with nonlinear patterns. Specifically, we are looking for a two-dimensional distribution for which conditioning on the time

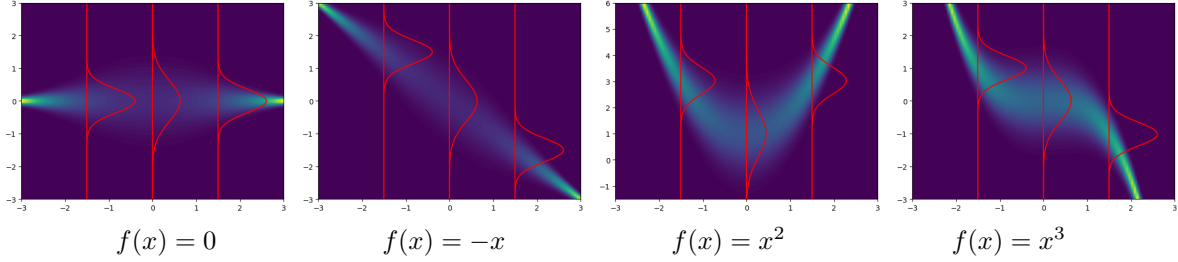


Figure 3. Folded-Gaussian distribution is capable of capturing both linear and nonlinear patterns. It is crucial to highlight that the conditional distributions (marked in red) are classical Gaussians.

variable would produce one-dimensional Gaussians aligned along an arbitrary curve (not necessarily linear), as shown in Figure 3. A possible solution is to ensure that a conditional distribution of $s|t$ is a Gaussian distribution

$$\mathcal{N}(m_s + f(m_t - t), a(t)\sigma_s^2), \quad (4)$$

where $f: \mathbb{R} \rightarrow \mathbb{R}$ and $a: \mathbb{R} \rightarrow [0, 1]$ are functions designed to capture a desired time-dependent shift and rescaling of the space variable, respectively. In practice, we use a polynomial of a given order as f and the likelihood of the time variable (scaled to the unit interval) as a , i.e.,

$$a(t) = \frac{\mathcal{N}(m_t, \sigma_t^2)(t)}{\mathcal{N}(m_t, \sigma_t^2)(m_t)}. \quad (5)$$

Finally, to recover a joint distribution of a space-time variable, we can apply the standard chain rule for random variables, which leads to a density function given by the following formula:

$$\mathcal{N}(m_s + f(m_t - t), a(t)\sigma_s^2)(s|t) \cdot \mathcal{N}(m_t, \sigma_t^2)(t). \quad (6)$$

It is important to note that while marginal distributions of t and conditional distributions of $s|t$ are both Gaussian, the resulting joint distribution is not a Gaussian anymore. This makes it an interesting object not only for applications, but also for further theoretical studies. To ensure the applicability of our contribution to a wider range of contexts, we extend the discussion to an arbitrary dimension in the following paragraph.

Folded-Gaussians in \mathbb{R}^d We begin with a multivariate Gaussian distribution $\mathcal{N}(\mathbf{m}, \Sigma)$, which is defined for a space-time random variable $\mathbf{x} = (s, t) \in \mathbb{R}^{d-1} \times \mathbb{R}$ by a mean vector $\mathbf{m} = (m_s, m_t)$ and a covariance matrix

$$\Sigma = \begin{bmatrix} \Sigma_s & 0 \\ 0 & \sigma_t^2 \end{bmatrix}, \quad (7)$$

where Σ_s denotes the diagonal covariance matrix of the space component of \mathbf{x} . The probability density function

(PDF) is then factorized into two independent normal densities, i.e.,

$$\mathcal{N}(\mathbf{m}, \Sigma)(\mathbf{x}) = \mathcal{N}(m_s, \Sigma_s)(s) \cdot \mathcal{N}(m_t, \sigma_t^2)(t). \quad (8)$$

Note that such a distribution permits the modeling of only simple linear structures that are spanned along coordinate axes (see the leftmost picture in Figure 3 for an appropriate example in the two-dimensional case).

To address the aforementioned limitation, we propose incorporating non-trivial conditioning, which allows for a more flexible representation. In particular, we apply the following time-dependent transformation on the space variable:

$$s \rightarrow \sqrt{a(t)}(s - m_s) + m_s + f(m_t - t), \quad (9)$$

where $a: \mathbb{R} \rightarrow (0, 1]$ and $f: \mathbb{R} \rightarrow \mathbb{R}^{d-1}$ are suitably chosen functions. This yields the new space variable with the conditional normal density

$$\mathcal{N}(m_{s|t}, \Sigma_{s|t}, a, f)(s|t) = \mathcal{N}(m_s + f(m_t - t), a(t)\Sigma_s)(s|t), \quad (10)$$

which in turn gives rise to a novel Folded-Gaussian distribution with the PDF defined as¹:

$$\mathcal{FN}(\mathbf{m}, \Sigma, a, f)(\mathbf{x}) = \mathcal{N}(m_{s|t}, \Sigma_{s|t}, a, f)(s|t) \cdot \mathcal{N}(m_t, \sigma_t^2)(t), \quad (11)$$

where $s|t$ represents the new time-conditioned space random variable and $\mathbf{x} = (s|t, t)$ (we do not change the notation as we do not believe it would cause confusion). A notable benefit of Folded-Gaussians is their ability to effectively capture a range of relationships present in the data. This is due to the inherent flexibility in selecting the functions f and a . In the context of the VeGaS model, the polynomial function f (with trained coefficients) is employed in conjunction with the likelihood-based function a (as in Equation 5). Consequently, our approach allows

¹Note that in this case we are using the standard chain rule for random variables.

us to encompass both linear and non-linear patterns, as illustrated in Figure 3, which refers to a simplified case of two-dimensional distributions. Furthermore, the incorporation of likelihood-based time-dependent rescaling leads to the disappearance of the tails of Folded-Gaussians, thus facilitating the capture of elements present in only a portion of a video stream (an optimal scenario would entail these elements initially approaching the camera and subsequently receding from view).

The following paragraph offers further theoretical insight into the Folded-Gaussian distribution, including the formal arguments that underpin Equations (10) and (11).

Theoretical Study It should first be noted that the transformation given in Equation (9) is of an affine form $As + b$. Consequently, the distribution of the random variable $s|t$ is also Gaussian with parameters $Am_s + b$ and $A\Sigma_s A^T$. Given that in our case we have

$$A = \sqrt{a(t)}I_d, b = \left(1 - \sqrt{a(t)}\right) m_s + f(m_t - t), \quad (12)$$

we can conclude that

$$s|t \sim \mathcal{N}(m_s + f(m_t - t), a(t)\Sigma_s), \quad (13)$$

which justifies the assertion made in Equation (10). Moreover, a straightforward assessment of the correctness of the definition of the Folded-Gaussian distribution, as given by Equation (11), can be conducted through the following calculation:

$$\begin{aligned} & \int \mathcal{FN}(m, \Sigma, a, f)(x) dx = \\ & \int \left(\int \mathcal{N}(m_{s|t}, \Sigma_{s|t}, a, f)(s|t) ds |t \right) \mathcal{N}(m_t, \sigma_t^2)(t) dt \quad (14) \\ & = \int \mathcal{N}(m_t, \sigma_t^2)(t) dt = 1. \end{aligned}$$

Similarly, based on Equation (10), the PDF of the conditional distribution of the random variable $s|t$ can be computed as follows:

$$\begin{aligned} & \frac{\mathcal{FN}(m, \Sigma, a, f)(s|t, t)}{\int \mathcal{FN}(m, \Sigma, a, f)(s|t, t) ds} = \\ & \frac{\mathcal{N}(m_t, \sigma_t^2)(t) \cdot \mathcal{N}(m_{s|t}, \Sigma_{s|t}, a, f)(s|t)}{\mathcal{N}(m_t, \sigma_t^2)(t) \cdot \int \mathcal{N}(m_{s|t}, \Sigma_{s|t}, a, f)(s|t) ds |t} \quad (15) \\ & = \mathcal{N}(m_{s|t}, \Sigma_{s|t}, a, f)(s|t), \end{aligned}$$

which corroborates our previous assertion.

4. Video Gaussian Splatting

This section introduces our Video Gaussian Splatting (VeGaS) model. We start with a brief overview of the 3D

Gaussian Splatting (3DGS) [13] method and then proceed to the MiraGe [33] approach for 2D images, while at the same time tailoring the presentation for straightforward integration into our proposed solution. The section concludes with a detailed description of the VeGaS model.

3D Gaussian Splatting The 3D Gaussian Splatting (3DGS) [13] method employs a family of three-dimensional Gaussian distributions

$$\mathcal{G}_{3DGS} = \{(\mathcal{N}(m, \Sigma), \rho, c)\}, \quad (16)$$

characterized by a set of attributes, including location (mean) m , covariance matrix Σ , opacity ρ , and color c . In practice, the covariance matrix Σ is factorized as $\Sigma = RSSR^T$, where R is the rotation matrix and S is a diagonal matrix containing the scaling parameters. Therefore, it is also possible to use the notation $N(m, R, S)$ instead of $N(m, \Sigma)$.

The efficiency of the 3DGS technique is primarily attributable to its rendering process, which involves the projection of 3D Gaussians onto a two-dimensional space. Throughout the training process, all parameters are optimized according to the mean square error (MSE) cost function. Since such a procedure often results in local minima, 3DGS can employ supplementary training methods that include component creation, removal, and repositioning based on the proposed heuristic, which is both a fast and effective strategy. In addition, the GS training process is executed within the CUDA kernel, allowing for rapid training and real-time rendering.

Gaussian Splatting for 2D images The MiraGe [33] approach employs the 3DGS technique to accommodate 2D images. This is accomplished by employing flat Gaussians positioned on the plane spanned by the canonical vectors $e_1 = (1, 0, 0)$ and $e_2 = (0, 1, 0)$, which gives rise to a specific type of parametrization. In essence, this method deals with a family of 3D Gaussian components of the form

$$\mathcal{G}_{MiraGe} = \{(\mathcal{N}(m, R, S), \rho, c)\}, \quad (17)$$

where $m = (m_1, m_2, 0)$, $S = \text{diag}(s_1, s_2, \varepsilon)$, and

$$R = [r_1, r_2, e_3] = \begin{bmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix}. \quad (18)$$

(ε is a small positive constant used to ensure compatibility with the three-dimensional framework.) Subsequently, utilizing the parametrization proposed by the GaMeS [32] model, such flat Gaussians can be represented by three points (triangle face)

$$V = [m, v_1, v_2] = \mathcal{T}(m, R, S), \quad (19)$$

with the vertices defined as $v_1 = m + s_1 r_1$, and $v_2 = m + s_2 r_2$. On the other hand, given a face representation $V = [m, v_1, v_2]$, the Gaussian component

$$\mathcal{N}(m, R, S) = \mathcal{N}(\mathcal{T}^{-1}(V)) \quad (20)$$

can be reconstructed through the mean m , the rotation matrix $R = [r_1, r_2, e_3]$, and the scaling matrix $S = \text{diag}(s_1, s_2, \varepsilon)$, where the parameters are defined by the following formulas:

$$r_1 = \frac{v_1 - m}{\|v_1 - m\|}, \quad r_2 = \text{orth}(v_2 - m; r_1, e_2), \quad (21)$$

$$s_1 = \|v_1 - m\|, \quad s_2 = \langle v_2 - m, r_2 \rangle. \quad (22)$$

In this context, $\text{orth}(\cdot)$ represents one iteration of the Gram-Schmidt process [1]. We would like to highlight that the formulas presented above have been adjusted to align with our framework and may therefore differ slightly from those provided in [32, 33].

The use of the GaMeS parameterization allows for the modification of the position, scale, and rotation of Gaussians by altering the underlying triangle face. Furthermore, the MiraGe extension facilitates the manipulation of 2D images within 3D space, thereby creating the illusion of three-dimensional effects.



Figure 4. *Video edition*. Note that VeGaS permits selection of a single frame and modification of some of its elements. The model was trained on the DAVIS dataset [25].

Video Gaussian Splatting Consider a video comprising a sequence of frames $[I_{t_1}, \dots, I_{t_n}]$, indexed by their occurrence times scaled to the unit interval $[0, 1]$. In this context, the MiraGe model may be employed for each consecutive frame, as it can be treated as a separate 2D image. Consequently, this would result in a joined family of 3D Gaussian distributions

$$\mathcal{G}_{\text{MiraGe}}^{t_1} \cup \dots \cup \mathcal{G}_{\text{MiraGe}}^{t_n}, \quad (23)$$

where each $\mathcal{G}_{\text{MiraGe}}^{t_i}$ is given by Equation (17), which could be considered an adequate representation of the entire data. However, such an approach completely ignores the relationships that naturally exist within a video stream.

To overcome the mentioned limitation, we propose to construct Gaussians related to successive frames by

conditioning of corresponding three-dimensional Folded-Gaussian distributions at frames occurrence times. The resulting Video Gaussian Splatting (VeGaS) model is thus formally defined as a collection on 3D Folded-Gaussians

$$\mathcal{G}_{\text{VeGaS}} = \{(\mathcal{FN}(m, \Sigma, a, f), \rho, c)\}, \quad (24)$$

where for each component, three-dimensional extensions (see the preceding paragraph) around two-dimensional conditional distributions

$$\mathcal{N}(m_{s|t_1}, \Sigma_{s|t_1}, a, f), \dots, \mathcal{N}(m_{s|t_n}, \Sigma_{s|t_n}, a, f) \quad (25)$$

are built using Equation (10), with respect to the common opacity ρ and color c .

It is important to highlight our method does not use the fixed frames occurrence times t_1, t_2, \dots, t_n , but learns them through an optimization procedure that guarantees superior reconstruction quality. Specifically, we use the dynamic frame fitting function $f_t: \mathbb{Z}_+ \rightarrow [0, 1]$, which maps the frame number k to its scaled occurrence time t_k as follows

$$t_k = f_t(k) = \sum_{i=1}^k \sigma(w)_i = \sum_{i=1}^k \frac{e^{w_i}}{\sum_{j=1}^{n-1} e^{w_j}}, \quad (26)$$

where w_1, w_2, \dots, w_{n-1} are trainable parameters. (For interpolated frames between I_k and I_{k+1} , we use evenly spaced times between t_k and t_{k+1}). Additionally, it should be noted that in the VeGaS model we are dealing with two types of spatial distributions. First, we have Folded-Gaussians, which represent the dynamics in the video stream. Second, we generate flat conditional distributions, which are then extended to 3D Gaussians by adding a small orthogonal component (controlled by the value of ε). In this case, two opposing cameras reconstruct a single 2D image – one produces the original image, while the other produces its mirrored version.

5. Experiments

This section presents an extensive experimental study of the VeGaS model in a variety of settings, which compares its efficiency with respect to a diverse range of state-of-the-art solutions.

Datasets The efficacy of our method was evaluated on two datasets: the Bunny dataset [16] and the DAVIS dataset [25]. The Bunny dataset comprises 132 frames with a resolution of 720×1280 . In accordance with the specifications outlined in [4], the video is cropped to a resolution of 640×1280 . The DAVIS dataset is a high-quality and high-resolution collection of videos utilized for the purpose of video object segmentation. It encompasses a multitude of videos, each comprising a total of less than 100 frames. This dataset is accessible in two distinct versions: a full-resolution version and a more compact 480p version.

Table 1. *Frame reconstruction*. Performance of the VeGaS model on the evaluation setting proposed in [29], using various videos from the DAVIS dataset [25], in terms of the PSNR metric. Note that, in each situation, VeGaS obtains the best metric scores.

Model	Bear	Cows	Elephant	Breakdance-Flare	Train	Camel	Kite-surf	Average
Omnimotion [34]	22.96	23.93	26.59	24.45	22.85	23.98	23.72	24.07
CoDeF [23]	29.17	28.82	30.50	25.99	26.53	26.10	27.17	27.75
VGR [29]	30.17	28.24	29.82	27.18	28.09	27.74	27.82	28.44
VeGaS-Full (our)	31.79	27.64	30.93	29.37	31.20	30.76	35.84	31.08
VeGaS-480p (our)	33.23	30.27	33.28	33.19	32.86	32.23	38.23	33.31

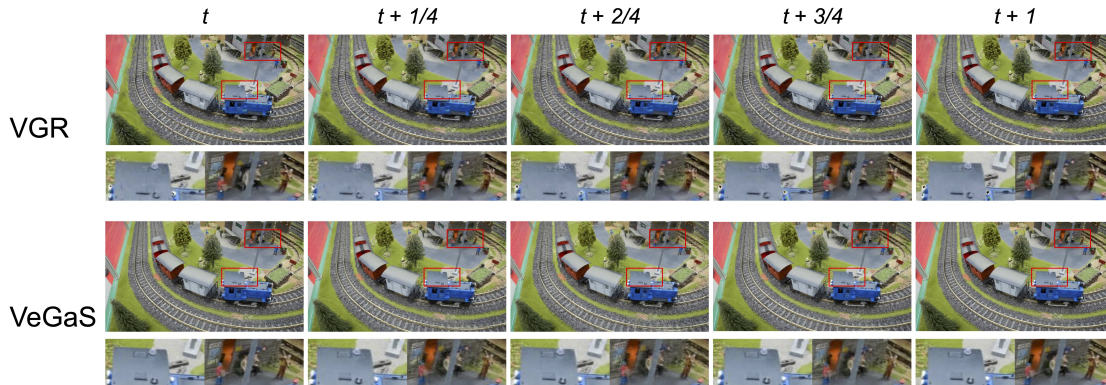


Figure 5. *Frame interpolation*. Qualitative results obtained by VeGaS and VGR [29] on a selected video object from the DAVIS dataset [25]. Frames at times t and $t + 1$ are reconstructions of two consecutive original frames, while frames at times $t + 1/4$, $t + 2/4$, and $t + 3/4$ are interpolated. Note that VeGaS produces outcomes that are slightly more favorable.

Implementation Details The initialization process entails the uniform sampling of Gaussian means m_1 and m_2 , which are positioned as points within a two-dimensional bounding box. The activation function utilized for m_t is the sigmoid function, which is initialized in a manner that ensures the resulting values from the activation process are uniformly distributed between 0 and 1. Similarly, the exponential function is employed as the activation function for σ_t , with the initial values distributed uniformly between 0.01 and 1. The coefficients of the polynomial function f are sampled uniformly between -1 and 1 . Furthermore, the rotation matrix is parameterized as a single number (rotation angle) and the activation function employed is sigmoid multiplied by 2π . The angle is initialized to be uniformly distributed between 0 and 2π .

The model is trained for 30,000 steps at a batch size of 3, utilizing a polynomial function of degree 7 and 500,000 initial Gaussians, unless otherwise specified. The learning rates, densification, pruning, and opacity reset settings are all consistent with the 3DGS [13] framework. In accordance with the MiraGe approach, two cameras are utilized: the initial camera generates the original image, while the second camera produces its mirrored version.

Frame Reconstruction We conducted a series of experiments to assess the efficacy of our method in frame recon-

struction tasks. The first experimental setup was adapted from that proposed in [29], where the authors introduce the VGR model and evaluate its performance compared to two state-of-the-art baselines, namely Omnimotion [34] and CoDeF [23]. Table 1 presents the values of rendering quality metrics for various videos from the DAVIS dataset. As there is no information in [29] on the resolution used for the evaluation, we report our results on both cases. It should be noted that, in each situation, VeGaS obtains the best metric scores. Furthermore, our model is capable of reconstructing videos with high quality and fidelity, as illustrated in Figure 6, which provides visualizations based on a selected object from the DAVIS dataset [25].

In the other experiments, the scene setting proposed by the authors of [43] was used to evaluate the DNeRV model against various NeRF-based baselines, namely [3], E-NeRV [19], HNeRV [4], and DNeRV [43]. In accordance with this protocol, the full-resolution version of each scene is center-cropped to a resolution of 960×1920 . The results (i.e., the values of the rendering quality metrics) are presented in Table 2, which shows the values obtained for various video objects from the DAVIS dataset. It should be noted that VeGaS outperforms all the considered NeRF-based models.

Frame Interpolation In subsequent experiments, we utilized the continuous representation of the video data pro-

Table 2. *Frame reconstruction*. Performance of the VeGaS model on setting proposed in [43], using various videos from the DAVIS dataset [25], in terms of the PSNR and SSIM metrics. Note that VeGaS outperforms all baseline models.

	NeRV [3]		E-NeRV [19]		HNeRV [4]		DNeRV [43]		VeGaS (our)	
	PSNR↑	SSIM↑	PSNR↑	SSIM↑	PSNR↑	SSIM↑	PSNR↑	SSIM↑	PSNR↑	SSIM↑
Blackswan	28.48	0.812	29.38	0.867	30.35	0.891	30.92	0.913	34.92	0.932
Bmx-bumps	29.42	0.864	28.90	0.851	29.98	0.872	30.59	0.890	33.01	0.915
Bmx-trees	26.24	0.789	27.26	0.876	28.76	0.861	29.63	0.882	31.78	0.896
Breakdance	26.45	0.915	28.33	0.941	30.45	0.961	30.88	0.968	32.27	0.950
Camel	24.81	0.781	25.85	0.844	26.71	0.844	27.38	0.887	31.12	0.886
Car-round	24.68	0.857	26.01	0.912	27.75	0.912	29.35	0.937	32.75	0.941
Car-shadow	26.41	0.871	30.41	0.922	31.32	0.936	31.95	0.944	36.41	0.956
Car-turn	27.45	0.813	29.02	0.888	29.65	0.879	30.25	0.892	31.44	0.852
Cows	22.55	0.702	23.74	0.819	24.11	0.792	24.88	0.827	27.97	0.834
Dance-twirl	25.79	0.797	27.07	0.864	28.19	0.845	29.13	0.870	30.45	0.850
Dog	28.17	0.795	30.40	0.882	30.96	0.898	31.32	0.905	34.52	0.914
Average	26.40	0.818	27.85	0.879	28.93	0.881	29.66	0.901	32.42	0.902

Table 3. *Ablation study*. Effect of a batch size and a degree of the polynomial function f on the performance of VeGaS in terms of the PSNR metric and the final number of Gaussians (in parentheses). The model was trained on the Bunny dataset [16].

Batch size	Polynomial degree					Mean training time
	1	3	5	7	9	
1	36.73 (1.26M)	37.31 (1.77M)	37.30 (1.73M)	37.36 (1.72M)	37.42 (1.83M)	31m20s
3	38.15 (0.57M)	38.31 (0.58M)	38.39 (0.59M)	38.53 (0.62M)	38.24 (0.59M)	56m30s
5	37.84 (0.33M)	37.94 (0.32M)	37.95 (0.32M)	37.92 (0.31M)	37.94 (0.31M)	1h15m29s

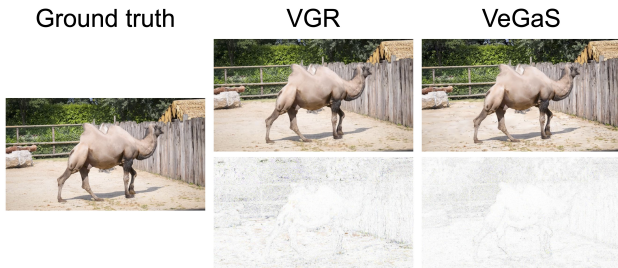


Figure 6. *Frame reconstruction*. Qualitative results obtained by VeGaS and VGR [29] on a selected video object from the DAVIS dataset [25]. The bottom row illustrates the discrepancy between the ground truth and the reconstructed image. Note that VeGaS is capable of reconstructing videos with high quality and fidelity.

Table 4. *Ablation study*. Effect of the initial number of Gaussians on the performance of VeGaS in terms of the PSNR metric, final number of Gaussians, and training time. The model was trained on the Bunny dataset [16] with batch size 3 and polynomial degree 7.

Initial Gaussians	PSNR↑	Final Gaussians	Training time
0.10M	38.53	0.62M	57m00s
0.20M	38.85	0.62M	57m29s
0.30M	38.99	0.62M	58m03s
0.40M	38.96	0.64M	59m36s
0.50M	39.02	0.65M	58m58s
0.60M	38.86	0.66M	1h00m53s

vided by our model to examine the potential for frame interpolation at a desired upsampling rate. To generate addi-

tional frames, the Folded-Gaussians were sliced at uniform intervals between each pair of consecutive frames. Figure 5 compares the results obtained by VeGaS and VGR [29] on a selected video object from the DAVIS dataset. The qualitative study reveals that interpolations using our method yield superior results. However, it should be noted that the source code for VGR has not been publicly released by the authors of [29], preventing a direct comparison of the respective rendering metrics scores.

Video Edition To illustrate the adaptability of the VeGaS model in editing video data, a series of experiments were carried out on entire scenes and specific objects from the DAVIS dataset. The results, presented in Figures 2 and 4, confirm that our method allows for both global modification (e.g., multiplication or scaling) of selected objects and for the choice of a single frame to modify some of its elements.

Ablation Study In our ablation study, we examined the impact of different hyperparameters of the VeGaS model trained on the Bunny dataset [16]. Table 3 presents the final rendering quality metric scores and the number of Gaussians obtained for various batch sizes and degrees of the polynomial function f . In turn, Table 4 presents the final metric values, numbers of Gaussians, and training times with reference to the initial numbers of Gaussians. As can be observed, VeGaS attains superior results when applied with a batch size of 3 and a polynomial degree of 7, with a

starting number of 0.50M Gaussian components.

6. Conclusions

In this paper, we propose the VeGaS model, which has been designed for the processing of video. To construct VeGaS, we have introduced a novel family of Folded-Gaussian distributions, which allow for the capture of nonlinear patterns in the video stream. The results of the conducted experiments demonstrate that our method enables superior reconstructions and realistic modifications within video frames.

References

- [1] Ake Björck. Numerics of Gram-Schmidt orthogonalization. *Linear Algebra and Its Applications*, 197:297–316, 1994. 6
- [2] Cheng-Hung Chan, Cheng-Yang Yuan, Cheng Sun, and Hwann-Tzong Chen. Hashing neural video decomposition with multiplicative residuals in space-time. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7743–7753, 2023. 2
- [3] Hao Chen, Bo He, Hanyu Wang, Yixuan Ren, Ser Nam Lim, and Abhinav Shrivastava. Nerv: Neural representations for videos. *Advances in Neural Information Processing Systems*, 34:21557–21568, 2021. 7, 8
- [4] Hao Chen, Matthew Gwilliam, Ser-Nam Lim, and Abhinav Shrivastava. Hnerv: A hybrid neural representation for videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10270–10279, 2023. 6, 7, 8
- [5] Zeyuan Chen, Yinbo Chen, Jingwen Liu, Xingqian Xu, Vidit Goel, Zhangyang Wang, Humphrey Shi, and Xiaolong Wang. Videoinr: Learning video implicit neural representation for continuous space-time super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2047–2057, 2022. 2
- [6] Yuanxing Duan, Fangyin Wei, Qiyu Dai, Yuhang He, Wenzheng Chen, and Baoquan Chen. 4d-rotor gaussian splatting: towards efficient novel view synthesis for dynamic scenes. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–11, 2024. 3
- [7] Lin Gao, Jie Yang, Bo-Tao Zhang, Jia-Mu Sun, Yu-Jie Yuan, Hongbo Fu, and Yu-Kun Lai. Mesh-based gaussian splatting for real-time large-scale deformation, 2024. 3
- [8] 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. 3
- [9] Jiahui Huang, Leonid Sigal, Kwang Moo Yi, Oliver Wang, and Joon-Young Lee. Inve: Interactive neural video editing. *arXiv preprint arXiv:2307.07663*, 2023. 2
- [10] Yi-Hua Huang, Yang-Tian Sun, Ziyi Yang, Xiaoyang Lyu, Yan-Pei Cao, and Xiaojuan Qi. Sc-gs: Sparse-controlled gaussian splatting for editable dynamic scenes. *arXiv preprint arXiv:2312.14937*, 2023. 3
- [11] Adam Kania, Artur Kasymov, Maciej Zięba, and Przemysław Spurek. Hypernerf-gan: Hypernetwork approach to 3d nerf gan. *arXiv preprint arXiv:2301.11631*. 2
- [12] Yoni Kasten, Dolev Ofri, Oliver Wang, and Tali Dekel. Layered neural atlases for consistent video editing. *ACM Transactions on Graphics (TOG)*, 40(6):1–12, 2021. 2
- [13] 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. 2, 3, 5, 7
- [14] Sylwester Kłoczek, Łukasz Maziarka, Maciej Wołczyk, Jacek Tabor, Jakub Nowak, and Marek Śmieja. Hypernetwork functional image representation. In *International Conference on Artificial Neural Networks*, pages 496–510. Springer, 2019. 2
- [15] Agelos Kratimenos, Jiahui Lei, and Kostas Daniilidis. Dynmf: Neural motion factorization for real-time dynamic view synthesis with 3d gaussian splatting. *arXiv*, 2023. 3
- [16] Janus B. Kristensen. Big buck bunny. 2010. 6, 8
- [17] Maomao Li, Yu Li, Tianyu Yang, Yunfei Liu, Dongxu Yue, Zhihui Lin, and Dong Xu. A video is worth 256 bases: Spatial-temporal expectation-maximization inversion for zero-shot video editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7528–7537, 2024. 3
- [18] Xirui Li, Chao Ma, Xiaokang Yang, and Ming-Hsuan Yang. Vidtime: Video token merging for zero-shot video editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7486–7495, 2024. 3
- [19] Zizhang Li, Mengmeng Wang, Huaijin Pi, Kechun Xu, Jianbiao Mei, and Yong Liu. E-nerv: Expedite neural video representation with disentangled spatial-temporal context. In *European Conference on Computer Vision*, pages 267–284. Springer, 2022. 7, 8
- [20] Yiqing Liang, Numair Khan, Zhengqin Li, Thu Nguyen-Phuoc, Douglas Lanman, James Tompkin, and Lei Xiao. Gafre: Gaussian deformation fields for real-time dynamic novel view synthesis, 2023. 3
- [21] Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. In *3DV*, 2024. 3
- [22] Haoyu Ma, Shahin Mahdizadehghadam, Bichen Wu, Zhipeng Fan, Yuchao Gu, Wenliang Zhao, Lior Shapira, and Xiaohui Xie. Maskint: Video editing via interpolative non-autoregressive masked transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7403–7412, 2024. 3
- [23] 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. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8089–8099, 2024. 2, 7
- [24] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 165–174, 2019. 2
- [25] Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams, Luc Van Gool, Markus H. Gross, and Alexander Sorkine-Hornung. A benchmark dataset and evaluation methodology

- for video object segmentation. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016*, 2016. 2, 6, 7, 8
- [26] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 3
- [27] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. In *Advances in Neural Information Processing Systems*, pages 7462–7473. Curran Associates, Inc., 2020. 1
- [28] Ivan Skorokhodov, Savva Ignatyev, and Mohamed Elhoseiny. Adversarial generation of continuous images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10753–10764, 2021. 2
- [29] Yang-Tian Sun, Yi-Hua Huang, Lin Ma, Xiaoyang Lyu, Yan-Pei Cao, and Xiaojuan Qi. Splatter a video: Video gaussian representation for versatile processing. *arXiv preprint arXiv:2406.13870*, 2024. 2, 3, 7, 8
- [30] Filip Szatkowski, Karol J Piczak, Przemysław Spurek, Jacek Tabor, and Tomasz Trzcziński. Hypernetworks build implicit neural representations of sounds. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 661–676. Springer, 2023. 2
- [31] Joanna Waczyńska, Piotr Borycki, Joanna Kaleta, Sławomir Tadeja, and Przemysław Spurek. D-miso: Editing dynamic 3d scenes using multi-gaussians soup. *arXiv preprint arXiv:2405.14276*, 2024. 3
- [32] Joanna Waczyńska, Piotr Borycki, Sławomir Tadeja, Jacek Tabor, and Przemysław Spurek. Games: Mesh-based adapting and modification of gaussian splatting. *arXiv preprint arXiv:2402.01459*, 2024. 3, 5, 6
- [33] Joanna Waczyńska, Tomasz Szczepanik, Piotr Borycki, Sławomir Tadeja, Thomas Bohné, and Przemysław Spurek. Mirage: Editable 2d images using gaussian splatting. *arXiv preprint arXiv:2410.01521*, 2024. 2, 5, 6
- [34] Qianqian Wang, Yen-Yu Chang, Ruojin Cai, Zhengqi Li, Bharath Hariharan, Aleksander Holynski, and Noah Snavely. Tracking everything everywhere all at once. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19795–19806, 2023. 7
- [35] Wen Wang, Kecheng Zheng, Qiuyu Wang, Hao Chen, Zifan Shi, Ceyuan Yang, Yujun Shen, and Chunhua Shen. Gendef: Learning generative deformation field for video generation. *arXiv preprint arXiv:2312.04561*, 2023. 3
- [36] Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Wang Xinggang. 4d gaussian splatting for real-time dynamic scene rendering. *arXiv preprint arXiv:2310.08528*, 2023. 3
- [37] 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. 3
- [38] Vickie Ye, Zhengqi Li, Richard Tucker, Angjoo Kanazawa, and Noah Snavely. Deformable sprites for unsupervised video decomposition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2657–2666, 2022. 2
- [39] Tingyang Zhang, Qingzhe Gao, Weiyu Li, Libin Liu, and Baoquan Chen. Bags: Building animatable gaussian splatting from a monocular video with diffusion priors, 2024. 3
- [40] Xinjie Zhang, Xingtong Ge, Tongda Xu, Dailan He, Yan Wang, Hongwei Qin, Guo Lu, Jing Geng, and Jun Zhang. Gaussianimage: 1000 fps image representation and compression by 2d gaussian splatting. In *European Conference on Computer Vision*, 2024. 2
- [41] Yabo Zhang, Yuxiang Wei, Dongsheng Jiang, Xiaopeng Zhang, Wangmeng Zuo, and Qi Tian. Controlvideo: Training-free controllable text-to-video generation. *arXiv preprint arXiv:2305.13077*, 2023. 3
- [42] Youyuan Zhang, Xuan Ju, and James J Clark. Fastvideoedit: Leveraging consistency models for efficient text-to-video editing. *arXiv preprint arXiv:2403.06269*, 2024. 3
- [43] Qi Zhao, M Salman Asif, and Zhan Ma. Dnerv: Modeling inherent dynamics via difference neural representation for videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2031–2040, 2023. 7, 8