

SwinGS: Sliding Window Gaussian Splatting for Volumetric Video Streaming with Arbitrary Length

Bangya Liu

University of Wisconsin-Madison
bangya@cs.wisc.edu

Suman Banerjee

University of Wisconsin-Madison
suman@cs.wisc.edu

Abstract

Recent advances in 3D Gaussian Splatting (3DGS) have garnered significant attention in computer vision and computer graphics due to its high rendering speed and remarkable quality. While extant research has endeavored to extend the application of 3DGS from static to dynamic scenes, such efforts have been consistently impeded by excessive model sizes, constraints on video duration, and content deviation. These limitations significantly compromise the streamability of dynamic 3D Gaussian models, thereby restricting their utility in volumetric video streaming.

This paper introduces **SwinGS**, a streaming-friendly paradigm representing volumetric video as a per-frame-update 3D Gaussian model that could be easily scaled to arbitrary video length. Specifically, we incorporate a sliding-window based incremental optimization during the train stage as well as a straightforward rendering at client side. We implement a prototype of **SwinGS** and demonstrate its streamability across various datasets and scenes. Additionally, we develop an interactive WebGL viewer enabling real-time volumetric video playback on most devices with modern browsers, including smartphones and tablets. Experimental results show that **SwinGS** reduces transmission costs by 83.6% compared to previous work.

ACM Reference Format:

Bangya Liu, and Suman Banerjee. 2023 SwinGS: Sliding Window Gaussian Splatting for Volumetric Video Streaming with Arbitrary Length. In *XXXXX (XXXXX)*, October 2–6, 2023, Madrid, Spain. 1 pages. <https://doi.org/10.1145/1234567890>

1. Introduction

Volumetric video, also known as free view video (FVV), represents a revolutionary media format that enables viewers to experience content as if physically present in the scene. Unlike traditional video captured from a single perspective, volumetric video encapsulates the depth, shape, and motion of objects as well as people within a scene. This 3D representation can be viewed from any position or perspective in virtual reality (VR), augmented reality (AR), or on flat screens with user interaction. Beyond entertainment, volumetric video plays crucial roles in autonomous vehicles [39], robotics vision, and teleoperation [29].

Historically, volumetric video has relied on point clouds [8, 15] and meshes as foundational elements. However, these approaches have struggled to balance video quality with storage and bandwidth efficiency. Recent advances in computer graphics have introduced a new family of 3D scene representations: neural rendering. This includes Neural Radiance Fields (NeRF) [25] and the emerging 3D Gaussian Splatting (3DGS) [13]. While NeRF achieves superior rendering quality with compact storage, it suffers from intensive computational costs due to its sampling process, resulting in low frame rates. The 3D Gaussian model, a “volumetric and differentiable” variant of point clouds, has emerged as a promising alternative.

Industry demos of 3DGS, including Teleport from Varjo [32] and Gracia [1], have garnered significant attention upon release. More recently, static 3D Gaussian Splatting has been showcased on the Pico XR headset [2] and Apple Vision Pro via MetalSplat [3], offering impressive immersive experiences.

Recent research works [6, 9, 12, 18, 23, 30, 33, 35–37] have demonstrated the potential of 3DGS in representing dynamic 3D scenes. However, significant gaps remain between dynamic 3DGS and a fully realized 3D Gaussian-based volumetric video. Previous attempts have fallen short in three key areas: (i) excessive model sizes, (ii) limited video duration, and (iii) lack of mechanisms to handle content deviation across extended time spans. All three are crucial for streaming volumetric video from server to client.

To address these challenges, we propose a novel paradigm representing volumetric video as a dynamic 3D Gaussian model, **SwinGS**, yet in a streaming-friendly style compared with previous work. To rendering a volumetric video, the model retires a subset of 3D Gaussians from previous frames and introduces new Gaussians at the beginning of the next frame. Through assigning each Gaussian a clear lifespan indicating when it joins and leaves the model, the model can be easily adapted to new content in the subsequent frames. This solves the content deviation issue (iii). On the other hand, the transmission of a single bulky model is broken into consecutive transmissions of small pieces of data encapsulating per-frame update, making streaming of arbitrary length videos theoretically viable and addressing concerns (i) and (ii).

To facilitate such paradigm, we innovatively propose a sliding-window based incremental optimization approach. Within the window, Gaussians contributing to later frames will be optimized together with Gaussians contributing to earlier frames, while partial of the Gaussians are frozen. In this setup, Gaussians are naturally shared among frames, in the meantime that rendering quality of earlier frames are unaffected during the optimization for later frames.

Specifically, we train the 3D Gaussian model using Stochastic Gradient Langevin Dynamics (SGLD) and Gaussian relocation, as proposed in 3DGS-MCMC [14]. This approach allows the model to adapt to various 3D scenes across different frames, while keeping a constant number of Gaussians throughout training.

We have implemented **SwiNGS** on top of the 3DGS-MCMC codebase and evaluated it using various scenes from the ActorsHQ dataset [10] and DyNeRF dataset [17].

Our contributions can be summarized as follows:

- To the best of our knowledge, we are the first to extend 3DGS to the field of long volumetric video streaming, identifying its unique challenges compared to 3D Gaussian splatting tailored for short scenes.
- Targeting on the challenge, we propose **SwiNGS**, a new paradigm to represent volumetric video with a per-frame-update 3D Gaussian model. We also designed and developed a sliding-window based incremental optimization approach to facilitate the model training and convenient rendering at client side.
- We implement and evaluate the proposed **SwiNGS**, demonstrating its feasibility with diverse datasets. We also develop a WebGL-based viewer¹ that enables easy playback of volumetric video on portable devices.

2. Background

2.1. 3D Gaussian Splatting

3D Gaussian model is an emerging graphic primitives to represent a 3D scene, and 3D Gaussian splatting (3DGS) [13] is the technique rendering a model into 2D images with given camera poses. In a 3D Gaussian model, the scene is represented with a set of Gaussian points parameterized by covariance Σ , center position μ , color c , and opacity α . The intensity of a 3D gaussian at a given location x in the 3D space could be defined as:

$$G(x, \mu, \Sigma) = e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)} \quad (1)$$

In practice, Σ is decomposed into a rotation matrix R and a scaling matrix S , to guarantee semi-definiteness, as shown in Equation 2. Usually, 3D Gaussians could be visualized as

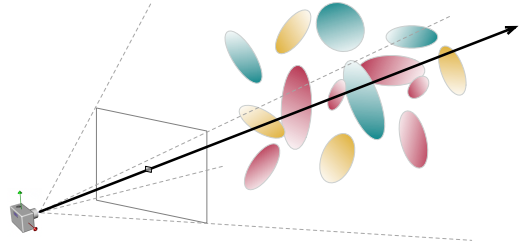


Figure 1: Rasterization Rendering of 3DGS

ellipsoids in the 3D space, in which case R could be interpreted as the orientation and S as the length of axes of visualized ellipsoid.

$$\Sigma = RSS^T R^T \quad (2)$$

3DGS is a rasterization-based rendering technique. When it comes to rendering, as shown by Figure 1, ray r emitting from camera center to the queried pixel on the rendering image will traverse a subset of 3D Gaussians, and the final color C of that pixel is computed by alpha blending this subset with the order of z-depth from close to far:

$$C(r) = \sum_{i \in N} c_i \alpha'_i \prod_{j=1}^{i-1} (1 - \alpha'_j) \quad (3)$$

Here α'_i is determined by multiplying opacity α_i with the integration of 3D Gaussian G 's intensity along the ray. In practice, integration will be performed by sampling from a 2D Gaussian projected from the original 3D Gaussian.

2.2. Dynamic 3DGS

Recent works have extended vanilla 3D Gaussian Splatting (3DGS) to dynamic scenes, primarily following two type of approaches: deformation field-based and 4D primitive-based methods.

Deformable 3DGS [37] pioneered the use of multilayer perceptrons (MLPs) to implement deformation fields, allowing 3D Gaussians to exhibit different properties across time frames, following which, [35] enhanced the deformation field's fitting capacity by introducing a hexplane [5] ahead of the MLP. Another early stage work, DynamicGS [23] incorporated rigidity loss to improve tracking accuracy between frames. Several recent works, including SC-GS [9], HiFi4G [12], and Superpoint Gaussian [33], proposed hierarchical structures where higher-layer Gaussians act as deformable skeletons, while lower-layer Gaussians bound to them fit appearance. Taking a different approach, Space-timeGS [18] used MLPs to encode appearance and parameterized polynomials to represent deformation.

¹A live demo is publicly available at <https://swingsplat.github.io/demo/>. The **SwiNGS** codebase will be made available following the paper's acceptance, in compliance with the double-blind review process.

Primitive	Rendering Quality	Rendering Speed	Accessibility	Storage
3DGS [13]	high	high (>100fps)	relatively easy	middle
NeRF [25]	high	low (<10fps)	relatively easy	low
point cloud	depends	middle	easy	high
mesh	depends	high	poor	middle

Table 1: Comparison between different 3D primitives

4D primitive-based approaches incorporate time as the fourth dimension of a Gaussian, pioneered by 4DGS [36] and 4D-Roter GS [6]. During rendering, these methods project (also called slice or condition) 4D primitives into 3D Gaussians before processing them through the alpha blending pipeline.

2.3. Volumetric Video Streaming

Traditional primitives for streamable volumetric video include point clouds and 3D meshes, where codecs play a crucial role due to limited bandwidth budget. Groot [15] presents a system that leverages an advanced octree-based codec to efficiently compress point cloud payloads. Another work, ViVo [8], focusing on point clouds, reduces bandwidth usage by considering visibility based on user viewport. More recently, MetaStream [7] introduced a comprehensive, point cloud-based system for capturing, creating, delivering, and rendering volumetric videos in an end-to-end fashion.

Transcoding offers another popular paradigm for volumetric video streaming. Vues [22], for instance, offloads rendering tasks to an edge server instead of rendering received 3D primitives on the user side, achieving a balance between rendering quality and bandwidth utilization.

There are some preliminary attempt trying to represent and stream volumetric video with neural radiance field (NeRF) [20, 21]. Yet the use of 3D Gaussians as primitives for volumetric video streaming remains a largely unexplored field, with most dynamic 3DGS research focusing on short video clips lasting less than a dozen seconds. The research most similar to our proposed **SwinGS** is 3DGStream [30], which introduces a Neural Transformation Cache (NTC) as a per-frame deformation field. A recent work, MGA [31], approaches the problem from a different angle. Treating network bandwidth as a constraint, it optimizes bitrate allocation by tuning various encoding parameters for 3DGS frames to achieve optimal overall rendering quality.

3. Motivation

3.1. Pitfall of Previous Methods

Among various primitives representing 3D scenes, including NeRF [25], point cloud, and mesh, 3D Gaussian achieves a balance between rendering quality, speed, accessibility,

and storage cost, as shown in Table 1. While NeRF delivers high-quality renderings with reasonable storage, its computational demands limit rendering speed. Works including StreamRF [16] and NeRFPlayer [28] attempted to increase rendering speed, but at a cost of dramatically increased model size. Meshes, although relatively compact and rendering-friendly, require extensive manual labor to create, especially for high-quality volumetric videos. Point clouds offer decent frame rates and are easily accessible from 3D dense reconstruction, but their rendering quality correlates with point count, leading to high storage costs and challenges like hole filling. Recent research on dynamic Gaussian splatting [6, 9, 12, 18, 23, 30, 33, 35–37] demonstrates the potential of 3D Gaussian models for representing 3D scenes. However, these approaches remain incompatible with volumetric video streaming due to several limitations:

Excessive Model Size: A naive approach could be constructing static 3D Gaussian models for each frame, which leads to substantial traffic overhead. To tackling with this, one of the baselines, 3DGStream [30], employs per-frame Neural Transformation Codes (NTC) to transform Gaussians between frames. However, this still incurs a storage overhead of approximately 7.8MB/frame, requiring a minimum bandwidth of 200MB/s for 30fps video. Other works [9, 12, 33, 35, 37] utilize a shared neural network for Gaussian deformation or appearance encoding, along with initial Gaussians for the first frame. While this reduces storage costs, the entire model still occupies hundreds of megabytes [18] and must be transmitted at once. Any packet loss during transmission could block the entire rendering pipeline, making it unsuitable for streaming applications.

Limited Video Length: Previous approaches are typically constrained to dynamic scenes lasting only a few seconds, primarily due to the limited capacity of a single neural network. [18] confirms that increasing clip length from 50 to 300 frames results in a PSNR drop from 29.48 to 29.17 for the Flame Salmon dataset of DyNeRF [17]. Scenes with more substantial motion and longer duration are likely to experience more severe degradation, as the neural network must learn and remember increasingly diverse and distinct Gaussian deformations across frames.

Long-term Content Deviation: Most current methods lack mechanisms for introducing new Gaussians when ob-

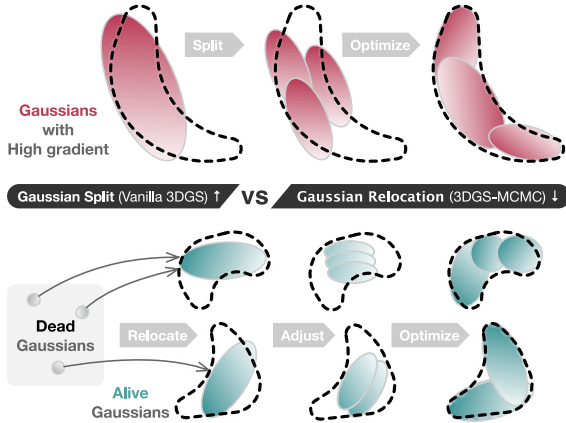


Figure 2: Different Gaussian densification methods (dotted line referring to the shape to fit)

jects enter the scene or removing them when objects exit. This limitation not only reduces rendering quality but also exacerbates the model size issue, as all Gaussians must be transmitted simultaneously. While fragmentation with multiple dynamic Gaussian models could potentially address this concern, it introduces new challenges such as visual discontinuities between fragments and does not alleviate the burst nature of model transmission, in addition to extra storage cost coming with multiple deformation networks.

3.2. 3DGS-MCMC for Bandwidth Shaping

For streamable content, maintaining a uniform data volume across time is crucial. Significant variations in the number of Gaussians between different time frames of a volumetric video can compromise the quality of service (QoS) on the client side, especially when the user is in a high mobility scenario with constrained network bandwidth.

Recent work, 3DGS-MCMC [14], introduces a novel approach to Gaussian densification during model optimization, as shown in Figure 2. Instead of simply splitting one Gaussian into several, 3DGS-MCMC relocates “dead” Gaussians (those with low opacity) to the positions of “alive” Gaussians (those with high opacity hence high presence in the scene). After that, parameters of both “dead” and “alive” Gaussians are adjusted in a way that the Gaussians distribution keeps approximately consistent. With total number of Gaussians keeps constant after densification, this method allows for precise control over the number of Gaussians.

Notably, this relocation operation is not limited to densification but also facilitates smooth transitions of Gaussian distributions between consecutive frames. Gaussians representing exiting objects are optimized into “dead” Gaussians with diminishing opacity and scale. In subsequent iterations, these “dead” Gaussians can be repurposed to represent newly appearing objects.

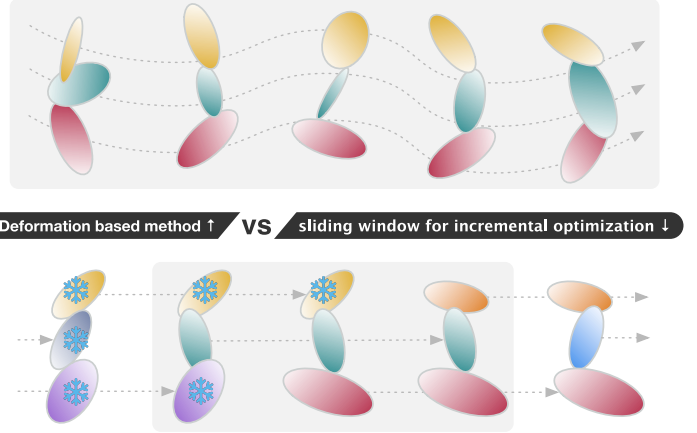


Figure 3: Different paradigm for dynamic 3D Gaussian model (dotted arrow referring to the lifespan of each Gaussian)

Building on this concept, the intuition could be formalized as such: 3DGS-MCMC continuously transitions the 3D Gaussian distribution using a combination of SGLD [34] and Gaussians relocation. This process shapes the model to represent different 3D scenes at various time frames of a volumetric video. Concurrently, a sliding window moves along the timeline, capturing snapshots of the Gaussian distributions at different times for volumetric video playback.

3.3. Sliding Window for Incremental Optimization

The key difference between proposed **SwiNGS** and previous work is visualized in Figure 3. Previous methods [4, 9, 18, 36, 37] dominantly adopt a deformation-based paradigm where a fixed set of Gaussians in the canonical space is deformed by a carefully designed deformation network at different frame to fitting the dynamics in the video. Yet the limited fitting capacities of deformation network and such tight coupling between frames makes the model hard to train. Hence, complicated modules including hexplane [5] and resfields [24] are integrated to compensate for that, usually with an extra overhead of storage and transmission.

Yet our proposed paradigm assigns a clear lifespan for each Gaussian so that the optimization is always focusing on a short snippet instead of the whole video. We could conveniently optimize those “temporal-local” Gaussians using a sliding window, then derive a streamable model in an incremental style. During the window optimization, Gaussians who also contributes to out-of-window frames are frozen to avoid degradation of previous frames’ rendering quality. This combination of per-frame-update paradigm and incremental optimization approach, makes training on long video sequence feasible.

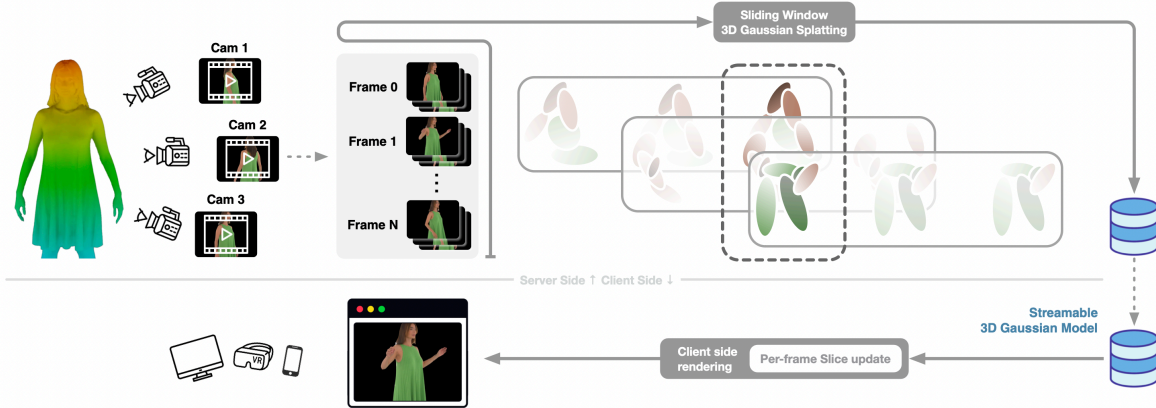


Figure 4: Overview of **SwinGS**

4. Design of SwinGS

4.1. Overview

In this paper, we introduce **SwinGS**, a streaming-friendly paradigm representing volumetric video as a per-frame-update 3D Gaussian model. Figure 4 provides an overview.

SwinGS begins with multi-view video input, accompanied by corresponding camera poses. The initial step involves decomposing these videos into individual frames, which are then clustered based on their frame index. This process yields a sequence of folders, each corresponding to a specific frame in the volumetric video.

Following frame clustering, we train a 3D Gaussians model using SGLD and Gaussian relocation to fit various frames within the sliding window. For incremental optimization, we freeze and archive a small portion of Gaussians at the end of each window training, in the meantime, inject some new optimizable Gaussians for the future frames. Consequently, we allow each Gaussian to contribute to image rendering across a small spanning of frames.

When it comes to video streaming, similarly, we only need to stream and update the small portion of Gaussians to update the model, which helps substantially reduce bandwidth requirements and makes efficient streaming of volumetric video possible.

The following two subsections will delve into the details of the training process and real-time streaming implementation, respectively.

4.2. Incremental Optimization

The general workflow of model training is outlined in Algorithm 1. The process involves **an outer loop** that shifts a sliding window across the video, iterating from $[0, \text{swin_size})$ to the video’s end, and **an inner loop** trains on frames randomly sampled within this sliding window. Here, constant swin_size represents the window length, which

also defines the maximum lifespan of a Gaussian and num_gs represents the maximum of Gaussians that will coexist in the model to render an image.

Each of the Gaussians uses two integers to indicate its lifespan: “start” and “expire”. The Gaussian will participate in rendering only if the current frame falls into its lifespan. During model training, we have two sets of Gaussians in the GPU memory, gs and matured , as shown in Algorithm 1 as global variables. To make it simple, gs are those Gaussians we are currently optimizing, usually helping fit the content in the new frames, while matured are those Gaussians that have been snapshotted already and are not optimizable, contributing to the rendering of previous frames as well as current frames. Whenever there is a frame training, both gs and matured will be used to derive loss yet there is no gradient for matured Gaussians.

Upon completing training within a sliding window, we increment both the start and end frames by one. We then check if any Gaussian’s lifespan begins earlier than the window’s start frame. If so, we **mature** that Gaussian and save its parameters.

It might be helpful to explain the difference between two concepts: **expire** and **mature**. **Expire** refers to the relationship between a Gaussian and the current rendering frame, determining whether the Gaussian contributes to that frame’s rendering. A Gaussian is considered **active** if the current frame falls within its lifespan. On the other hand, **mature** relates to optimization: if the first contributing frames of a Gaussian fall out after window sliding, that Gaussian should no longer be optimized, in order to preserve the training result for that fallen frame.

Notably, a matured Gaussian may still contribute to rendering any of its involving frames that remain within the current training sliding window. Thus, the **active** Gaussians for a current rendering frame comprise both optimizable

Algorithm 1: SwinGS Training Procedure

```

1 global gs, matured, stream, trainset
2 const swin_size, num_gs, relocate_period, iterations
3 procedure train_swin(st, ed):
4   for iter in range(iterations):
5     gt = sample_frame_between(trainset, st, ed)
6     frame = ground_truth.frame
7     active_idx = gs.filter(start ≤ frame < expire)
8     active_ma_idx = matured.filter(start ≤ frame < expire)
9     active_gs = gs[active_idx] + matured[active_ma_idx]
10    pred = render(gt.cam, active_gs)
11    loss = loss_func(gt.image, pred) + reg(active_gs)
12    loss.backward()
13    with no_grad:
14      gs[active_idx].param += λnoise * ε
15      if iter % relocate_period == 0 then
16        | relocate(gs[active_idx])
17  procedure mature(st):
18    mature_idx = gs.filter(start < st)
19    stream.write(gs[mature_idx].detach())
20    matured += gs[mature_idx].detach()
21    matured = matured[-num_gs:]
22    gs[mature_idx].birth = gs[mature_idx].expire
23    gs[mature_idx].start = gs[mature_idx].expire
24    gs[mature_idx].expire = gs[mature_idx].start + swin_size
25  procedure main:
26    gs[num_gs].param = random()
27    gs[num_gs].birth = 0
28    gs[num_gs].start = 0
29    gs[num_gs].expire = swin_size
30    train_swin(0, swin_size)
31    schedule_expire(gs)
32    for st in range(1, trainset.total_frames):
33      | mature(st)
34      | train_swin(st, st + swin_size)
35    mature(trainset.total_frames)

```

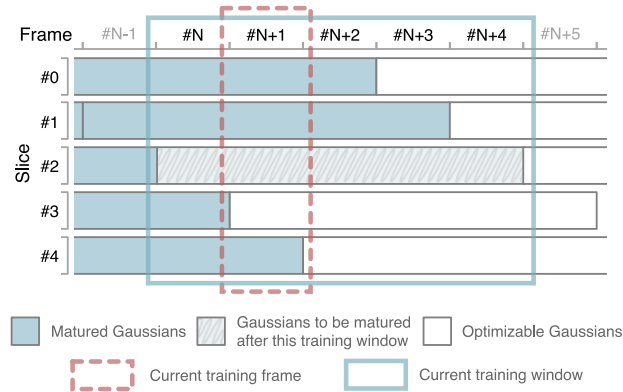


Figure 5: Sliding window with Gaussians

and matured Gaussians. We append matured Gaussians to a streamable model on disk as they reach maturity.

In addition to maturing partial Gaussians, the `mature()` procedure in Algorithm 1 also appends new optimizable Gaussians to maintain the number of Gaussians within a model. Those new Gaussians, together with rest of optimizable Gaussians in the model, will help fitting the content in new frames. In practice, we directly update lifespan-related variables as a bypass, avoiding the memory overhead of “release-then-allocate” operations.

To further accelerate this process, we evenly divide the entire 3D Gaussians model into several smaller slices. Gaussians within the same slice share identical lifespans and mature together, while different slices have misaligned lifespans. This arrangement ensures that exactly one slice matures in each frame. Figure 5 illustrates this slicing approach, with `swin_size` set to 5. Each bar represents a group of Gaussians with identical lifespans. Within a slice, multiple bars are positioned in a bumper-to-bumper style across different frames, as older Gaussians trained for previous frames mature and new optimizable Gaussians are introduced. Darker bars denote matured Gaussians, while white bars represent optimizable Gaussians. For training frame `#N+1`, matured Gaussians in slices `#0`, `#1`, `#4`, and optimizable Gaussians in slices `#2` and `#3` participate in image rendering. After model training within the `[N, N+5)` window, the sketchy bar of slice `#2` will mature, because there will be no further chance to optimize model on frame `#N`.

The function `schedule_expire()` mentioned in the Algorithm 1 is the one that performs such slicing. For the first sliding window `[0, swin_size)`, each Gaussian is initialized with full lifespan cross the window. After the first sliding window, `schedule_expire()` will bring forward the expire field of all the Gaussians in a way that after the first frame, $\text{num_update} = \frac{\text{num_gs}}{\text{swin_size}}$ Gaussians will expire at each further frame. In practice, we let Gaussians with higher opacity expire later and those with low opacity expire sooner.

Under this setup, each new frame in the volumetric video is represented as a small slice of data in the final streamable model. And the original streaming payload is reduced from `num_gs` Gaussians per frame to $\frac{\text{num_gs}}{\text{swin_size}}$ Gaussians per frame, which dramatically reduce the required bandwidth. Further because the original rendering interface is unchanged, the only extra cost here during the rendering stage is the memory operations within GPU.

4.3. Real Time Streaming and Rendering

Algorithm 2 illustrates the rendering process of the 3D Gaussian model on the client’s device. At beginning, the first `swin_size` slices consisting of total `num_gs` 3D Gaussians is streamed from the remote server to the client device. Then it is loaded to GPU to render the first frame.

Algorithm 2: SwinGS Streaming and Rendering

```
1 global frame, buffer, events, stream, user
2 const swin_size, num_gs, FPS
3 procedure render_thread:
4   while true:
5     active_idx = buffer.filter(start ≤ frame < expire)
6     render(user.cam, buffer[active_idx])
7 procedure update_thread:
8   while true:
9     sleep( $\frac{1000}{FPS}$ )
10    for (target_frame, slice, update) in events:
11      if target_frame == frame then
12        update_first = slice * num_update
13        update_last = (slice+1) * num_update
14        buffer[update_first:update_last] = update
15        break
16    frame += 1
17 procedure main:
18   frame = 0
19   buffer[:num_gs] = stream.read(num_gs)
20   num_update =  $\frac{\text{num\_gs}}{\text{swin\_size}}$ 
21   threading render_thread
22   threading update_thread
23   while true:
24     update = stream.read(num_update)
25     target_frame = update[0].birth
26     slice = target_frame % swin_size
27     events.append([target_frame, slice, update])
```

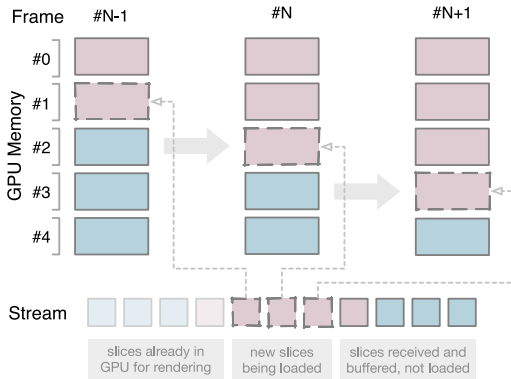


Figure 6: Per-frame slice update in client device's GPU

Subsequent slices will be received, processed, and buffered by the client. Those slices will be inserted into the GPU memory replacing expired slices. The `num_update` in Algorithm 2 represents the number of Gaussians in each slice, while `buffer` abstracts the GPU memory. The GPU memory is divided into `swin_size` slices, with each slice containing `num_update` Gaussians.

Figure 6 depicts how the client device's GPU memory is updated, in a slice by slice manner, as new streaming data from the server is received. For frame#N-1 slice #1 gets updated, while slice #2 and #3 get updated at frame#N and frame#N+1 respectively.

To decouple GPU rendering from streaming data receiving, which is IO intensive, two dedicated threads are instantiated. Received slice will first be buffered in CPU memory, specifically in the events queue. When it is the time to rendering its corresponding frame, `update_thread` will migrate the slice from CPU to GPU memory. `render_thread` keeps rendering images at the same FPS as vanilla 3DGS with current Gaussians in the GPU.

Compared to 3DGStream [30], which recalculates and refreshes all Gaussians for every frame, our approach significantly reduces both streaming traffic and GPU operations. These optimizations are crucial for mobile applications where resources are constrained.

4.4. Implementation

Our demo is built upon the foundation of 3DGS-MCMC [14]. We extended the original codebase by transitioning from a per-frame training approach to a per-window training strategy, as proposed in Algorithm 1. While preserving the CUDA-based differentiable Gaussian rasterization module, we refactored the `GaussianModel` and `Scene` classes to accommodate Gaussians. A new `SwinManager` class was implemented to handle sliding window.

To incorporate the perturbation required by Stochastic Gradient Langevin Dynamics (SGLD) [34], we introduced scaled noise for optimizable active Gaussians post-training for each frame. Our loss function adheres to the practice established in [14], encompassing image quality measurements and regularization terms for opacity α and scaling S .

In the following subsections, we discuss the key engineering challenges encountered during implementation and introduce our WebGL-based viewer.

4.5. Other Challenges

Adapting 3DGS-MCMC, originally designed for static 3D reconstruction, to our cross-frame Gaussians presented significant challenges. We focus on two primary obstacles: the relocation mechanism for Gaussians and the implementation of an efficient training data loader.

4.5.1. Relocation for Gaussians. Relocation is a crucial mechanism in 3DGS-MCMC that facilitates model transition between frames while maintaining a constant number of Gaussians. As outlined in [14], the relocation operation must preserve the probability of the sample state (i.e. Gaussians' distribution) before and after the move to prevent the collapse of the Markov chain Monte Carlo sampling process.

However, relocating Gaussians is complicated by the fact that active Gaussians for a given frame may originate from different slices or have varying statuses, as illustrated in Figure 5. Simply relocating “dead” Gaussians to the positions of “alive” Gaussians would disrupt the counting uniformity among slices, significantly complicating streaming and data loading processes. We considered restricting relocation to Gaussians within the same slice but determined this approach could potentially limit the fitting capacity of our Gaussian model. Instead, we propose a novel policy for relocating Gaussians across various slices:

- We permit the relocation of any optimizable “dead” Gaussian from an earlier slice to any optimizable “alive” Gaussian from a later slice. The terms “earlier” and “later” are defined by the “birth” field.
- Post-relocation, the “dead” Gaussian adopts the “start” field from the “alive” Gaussian. However, it retains its original “birth” and “expire” values. This approach allows the Gaussian to remain in its original slice while participating in rendering at a later start frame.

This strategy maintains the integrity of the slice structure while allowing for more flexible Gaussian optimization across the spacetime volume.

4.5.2. Data Loader for Trainset. The second major challenge stemmed from the limited GPU memory available in our testing environment. The original 3DGS codebase loads all training set images into GPU memory at the initiation of model training. This approach becomes infeasible when dealing with volumetric video containing numerous image sequences from different cameras, each comprising hundreds or thousands of frames. To address this constraint, we applied two key modifications:

- We developed a `LazyCamera` class to load images in a lazy manner, significantly reducing initial memory requirements and improving dataset loading speed.
- To manage memory constraints when training with longer sliding window sizes, we maintain a maximum number of frames in GPU memory. This approach involves dynamically unloading and reloading different frames as needed during the training process.

These modifications allow our system to handle volumetric video datasets efficiently with limited GPU resources.

4.6. WebGL Viewer

To visually demonstrate the feasibility of our proposed paradigm, we implemented a web-based viewer building upon the open-source project `antimatter15/splat2`, which was originally designed for rendering static 3D Gaussian models in browsers supporting WebGL 1.0. We extended this viewer to support rendering the streamable models as proposed by

SwinGS, while maintaining full interactivity for user view-point manipulation.

The vanilla viewer leverages the `ReadableStream` interface from JavaScript’s Web API to fetch data from the server in a streaming fashion. The fetched binary sequence is parsed into Gaussian, which are then transmitted to a rendering worker. Because there is no native support for Gaussian Splatting in WebGL, the worker will convert the Gaussian objects into textures that can be directly fed to the shader for real-time rendering.

To enable rendering our streamable model as volumetric video, our modified version incorporated the per-frame slice updates as proposed in Algorithm 2. The streamable model is hosted on Hugging Face, a popular online platform for sharing and distributing machine learning models. On the client side, whenever a new slice arrives, the rendering set of Gaussians will be scheduled to be updated with that new slice. However, due to the limitations of web applications in directly manipulating GPU memory, we perform slice updates on the Gaussians buffer within the rendering worker.

5. Evaluation

5.1. Dataset

We comprehensively evaluate **SwinGS** using two distinct datasets: `ActorsHQ` [10] and `DyNeRF` [17], each offering unique characteristics suitable for our volumetric video streaming task.

5.2. Setup

For model training in **SwinGS**, we carefully tuned the hyper-parameters to achieve optimal performance. The key parameters were set as follows: `scale_reg` at $1e-2$, `opacity_reg` at $2e-2$, and `noise_lr` at $5e4$ ($5e5$ for `ActorsHQ`). The degree of sphere harmonics function for viewpoint dependent coloring is set to 1 to reduce storage cost. These values were determined through extensive experimentation to balance model accuracy and computational efficiency.

We initialize our model using SfM points derived from COLMAP [27]. This approach provides a strong initial geometry estimate, significantly improving convergence speed and final model quality. We also explored random initialization, but found it prone to overfitting in our experimental setup. By default, our training setup is adapted to the specific characteristics of each dataset as follows:

Scene	Genesis Iters	Iters	Total Num	Swin Size
ActorsHQ [10]	30K	20K	100K	5
DyNeRF [17]	30K	10K	200K	5

²<https://github.com/antimatter15/splat>

Method	Cook Spinach	Cut Beef	Sear Steak	Mean PSNR	FPS	Storage per Frame	Streamable
4DGS [35]	32.46	32.90	32.49	32.62	30	0.3MB	✗
SpacetimeGS [18]	33.18	33.72	33.89	33.60	140	1MB	✗
3DGStream [30]	33.31	33.21	33.01	33.17	215	7.8MB	✓
Ours	32.43	32.02	33.00	32.48	300+	1.2MB	✓

Table 3: Overall comparison with other neural rendering methods Evaluate on DyNeRF Dataset. *per-scene PSNR is not reported

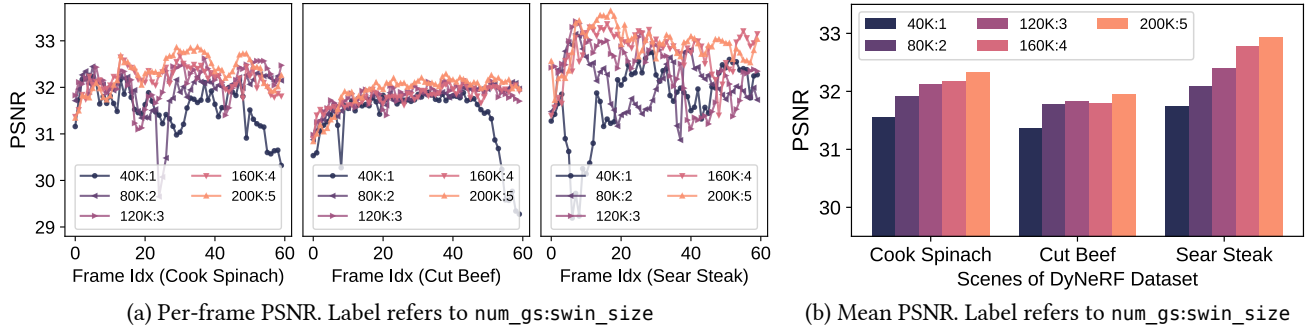


Figure 7: Streaming with various setup but same bandwidth on DyNeRF dataset [17]

Geneiss `iters` refers to the training iterations for the very first sliding window, i.e. $[0, \text{swin_size})$ and `Iters` refers to training iterations for the following sliding windows. Total Number indicates the maximum active Gaussians for a frame.

The choice of Gaussian counts (100K for ActorsHQ and 200K for DyNeRF) was made to balance model expressiveness with computational efficiency as well as storage cost. DyNeRF scenes benefit from a higher Gaussian count due to the presence of both dynamic foreground elements and static background details. The sliding window size of 5 frames was selected as an optimal trade-off between temporal coherence and computational resources.

For testset and trainset split, we adopts the convention from 3DGS [13] where `l1ffhold` parameter set as 8, which means, we will hold the video frames of one camera per eight cameras, leaving a 7:1 trainset/testset ratio. PSNR is chosen as the rendering quality metrics and FPS is adopted as rendering speed metrics.

5.3. Result

We evaluate the performance of **SwinGS** across multiple dimensions, including rendering quality, speed, and traffic cost. We also explore the impact of key parameters such as `swin_size` and `num_gs`. Finally, we present an analysis of rendering latency.

5.3.1. Overall comparison.

We compare **SwinGS** with previous works that utilize neural rendering to reconstruct dynamic 3D scenes, as shown in Table 3. We select three scenes from the DyNeRF

[17] dataset as benchmarks, training with the first 60 frames of each scene. We set `swin_size`=5 with `num_gs`=200K, updating 40K Gaussians per frame.

SwinGS surpasses NeRF-based methods like [16] and achieves comparable rendering quality to 4DGS [35]. While our PSNR is slightly lower than SpacetimeGS [18] and 3DGStream [30], **SwinGS** excels in balancing high streamability with low per-frame storage requirements. Unlike baseline methods that require compulsory neural network forwarding for numerous Gaussians, our approach primarily incurs costs from GPU memory operations, enabling our model to achieve the best FPS.

5.3.2. Impact of `swin_size` and `num_gs`.

The key design parameters in **SwinGS** are `swin_size` and `num_gs`. The former determines how many consecutive frames a Gaussian participates in, while the latter defines the total number of Gaussians used to fit a single frame. A larger `swin_size` typically results in higher content sharing among frames, saving bandwidth. Conversely, more `num_gs` can potentially provide better detail in the rendered image, albeit at the cost of increased traffic and storage. The bandwidth required for streaming a model can be formatted as:

$$\text{BW} = \text{FPS}_{\text{video}} \times \frac{\text{num_gs}}{\text{swin_size}} \times N_{\text{bytes/GS}} \quad (4)$$

This implies that the actual traffic cost is proportional to $\text{num_update} = \frac{\text{num_gs}}{\text{swin_size}}$, given a fixed video FPS and storage for each individual Gaussian (32 byte for WebGL viewer).

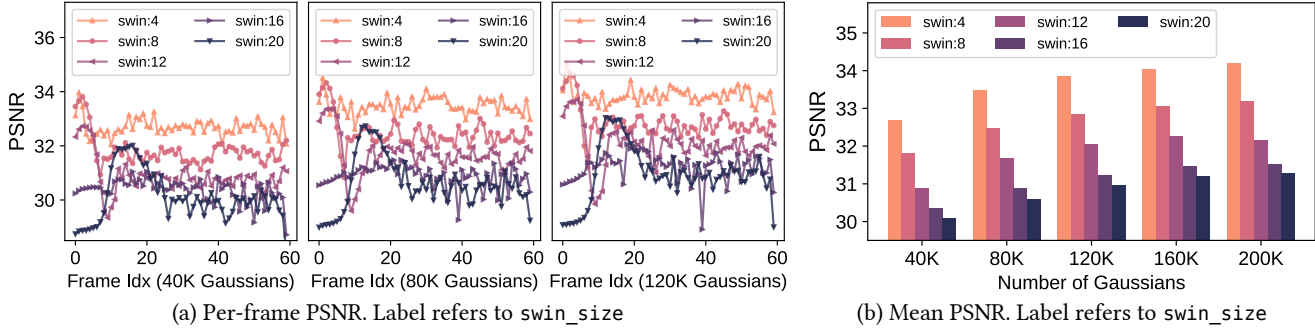
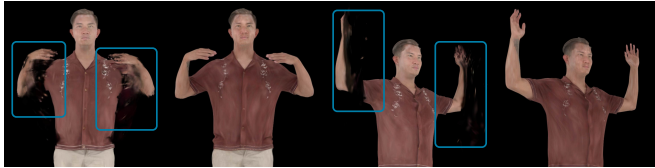
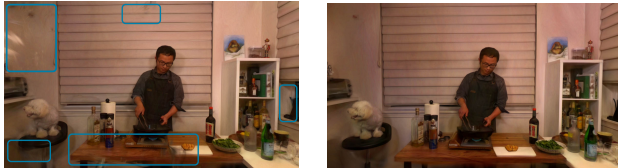


Figure 8: Streaming with various setup but wider swin_size
On ActorsHQ-Actor02-Sequence01 [10]



(a) Artifact due to over-sized swin_size
swin_size from left to right are: 16, 4, 16, 4



(b) Random point initialization might cause more floaters
left: random point cloud init, right: SfM init

Figure 9: Two types of artifact in **SwINGS**

Due to the content sharing feature of Gaussians in our method, we observed that, under limited bandwidth constraints, our approach consistently achieves higher rendering quality compared to transmitting each frame’s model in its entirety.

We trained the model with 5 different setups, varying swin_size from 1 to 5, where 1 indicates no common Gaussians among frames. The bandwidth was constrained to 40K Gaussians per frame. As shown in Figure 7(b), with limited bandwidth, the rendering quality of the current frame benefits from borrowing Gaussians from previous frames, especially for the DyNeRF dataset, which consists of dynamic human movement and relatively static background.

Figure 7(a) illustrates how per-frame PSNR differs across different setups and scenes. At some individual time frames, lower swin_size values can achieve comparable PSNR to swin_size=5, but they suffer from fluctuating PSNR over a larger scope. In other words, higher swin_size values result in more robust rendering quality due to increased redundancy of Gaussians across the time domain.

We also trained the model using the ActorsHQ [10] dataset with various setups, exploring the impact of swin_size and num_gs on the final rendering quality separately. We adopted wider swin_size values including 8, 12, 16, and 20. The results are shown in Figure 8. Figure 8(b) demonstrates that with a fixed number of Gaussians for rendering each frame, a larger sliding window typically decreases the rendering quality noticeably. This observation holds for both large and small num_gs models. Examining

per-frame PSNR, as shown in Figure 8(a), we find that smaller swin_size values display higher robustness in rendering quality. This distinct difference compared to the DyNeRF dataset results can be attributed to ActorsHQ being a fully dynamic volumetric video with limited sharable static content between frames. Consequently, when an oversized sliding window is set, Gaussians must be “torn” and to fit contents belonging to a wide range of frames.

5.3.3. Visual Artifacts.

Figure 9(a) provides a qualitative visualization of the aforementioned tearing from the Actor02-Sequence01 scene. The setup with swin_size=4 renders the actors’ movements perfectly, while training the model with swin_size=16 causes large areas of artifacts. In the first image, the actor’s hands exhibit a high degree of blurring, as Gaussians attempt to fit this waving-hand motion. Since per-pixel L1 loss is an important component of the final loss function, these Gaussians are optimized to the “average” value of all frames.

Figure 9(b) illustrates another type of artifact: floaters. Gaussian models in 3DGS-MCMC can be initialized either with cloud points extracted from Structure-from-Motion (SfM) or with points randomly sampled within a fixed range. Random initialization is more likely to produce floaters, often indicating model overfitting.

These floaters may also exacerbate the artifacts. In the third image of Figure 9(a), bulky floaters can be seen blocking the actor’s arms. This can be explained by some randomly initialized Gaussian points being optimized into

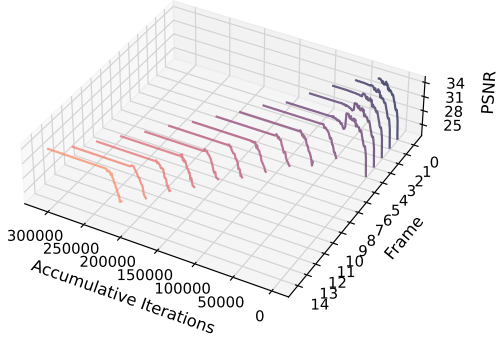


Figure 10: Test PSNR for the first 15 frames of Actor02-Sequence01

black floaters, which are barely distinguishable from the black background.

In summary, `swin_size` and `num_gs` are crucial parameters when applying **SwingS**. While increasing `num_gs` generally enhances PSNR, the final rendering quality is more sensitive to `swin_size`. To best benefit from shared Gaussians among frames, the sliding window size should match the degree of dynamics within the training data. For scenes with moderate motion and static backgrounds, larger `swin_size` values are worth exploring. For highly dynamic scenes, shorter sliding windows are preferable.

5.3.4. Streaming video with arbitrary length.

A key feature of **SwingS** is its ability to render very long volumetric videos without compromising rendering quality. This is attributed to the **mature** mechanism introduced in Section 4.2. When a frame has completed all its training tasks for the Gaussian model and is no longer being trained on, all its **active** Gaussians are matured and no longer optimized. When the sliding window shifts to the second half of the volumetric video and trains on new frames, the rendering quality of frames in the first half remains unaffected.

Figure 10 visualizes this process by plotting the test PSNR for the first 15 frames, with `swin_size` set to 5. Accumulative iteration refers to the total iterations the model has been trained on since the very first frame.

The sliding window [0,5) is trained first to create a preliminary 3D Gaussian model for frame#0. After this initial window, one-fifth of the Gaussians are matured, and the same number of new Gaussians join the model. These new Gaussians, along with the remaining four-fifths, are optimized by training on the sliding window [1,6).

As frame-by-frame training progresses, the plateau line at the tail of each curve indicates that one frame’s PSNR does not degrade with more following frames training. In fact, a

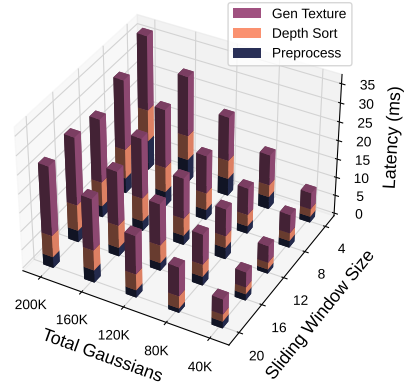


Figure 11: Rendering latency decomposition for WebGL viewer

Device (SoC)	preproc	sort	texture	overall
MacBook pro (M3 pro)	3.00	5.81	18.46	27.27
iPhone (A18 pro)	6.00	4.44	19.90	30.34
iPad (M1)	6.48	4.94	22.29	33.71
Pixel (Snapdragon 765)	13.02	17.37	49.59	79.98

Table 4: Rendering latency decomposition on different devices (ms)

frame’s PSNR becomes relatively converged and stable after the first fifth of its active Gaussians are matured.

5.3.5. Overhead at client side.

5.3.5.1. Differentiable rasterizer rendering. The primary cost introduced by **SwingS**, compared to vanilla 3DGS, is the memory operation to replace expired Gaussians with new ones in the GPU. This allows for a rendering speed of over 300 FPS when we use the vanilla differentiable Gaussian rasterization module to render the 3D scene.

5.3.5.2. WebGL rendering. However, when streaming volumetric video with the WebGL viewer, the scenario becomes more complex. Typically, after chunks of raw byte streams are read from the server, several steps are required before rendering the streaming data into images: raw data preprocessing, depth sorting, and texture generation. The first step involves parsing the binary raw data into Gaussian objects. Depth sorting arranges all Gaussians according to their distance from the camera center, from close to far. The third step, texture generation, is necessary because 3D Gaussians cannot be directly rendered by WebGL; they must be converted into texture data before being delivered to the shader for rendering. All steps introduces additional overhead to the rendering latency.

Figure 11 visualizes the latency composition when rendering volumetric video on a laptop. Raw data processing generally takes less than 1 ms, while depth sorting and texture generation take longer, ranging from 5ms to 18ms. There is a clear correlation between the number of Gaus-

sians and the time required for sorting and texture generation. This is reasonable considering our WebGL implementation does not manipulate the GPU directly and retransmits the texture data corresponding to all active Gaussians to the shader as a whole for every new frame. When configured with `swin_size` as 4 and `num_gs` as 200K, it takes approximately 34ms to complete the full pipeline for one frame, resulting in a worst-case video frame rate of around 30 FPS for serialized computation.

Considering the three stages could be fully paralleled among video frames, we expect an optimized version achieving over 60fps for video frame rate with texture generation as bottleneck stage for 18ms for our WebGL viewer.

Table 4 further profiles the WebGL viewer on a wide range of mobile devices from performance laptop to portable smartphones. Most latest devices are capable of video playback with 30fps.

6. Discussion

6.1. 3DGS and Point Cloud-Based Methods

While Table 1 presents 3D Gaussian Splatting (3DGS) and point cloud methods as alternatives, 3D Gaussians can actually be viewed as an extension of 3D points. In addition to the position (`xyz`) and color (`rgba`) properties inherent to points, 3D Gaussians incorporate rotation (`R`) and scaling (`S`). This relationship allows for seamless adaptation of point cloud-based volumetric video streaming innovations to the 3D Gaussian representation.

For instance, GROOT [15] employs octrees for efficient geometry data compression, a data structure recently also has been applied to 3D Gaussians in Octree-GS [26] to reduce the total number of Gaussians in a rendering scene. Similarly, RTGS [19] adopts an approach analogous to ViVo [8], utilizing user camera poses to optimize rendering resource allocation across different scene sections.

On the other hand, current off-the-shelf point cloud codec, MPEG Point Cloud Compression (PCC) [11] or a general data compressor like arithmetic entropy encoding could also be integrated to further reduce the transmission load during video streaming, at a cost of longer decoding time on the client side.

We anticipate further extensions of point cloud techniques to 3DGS, for example point cloud super-resolution [38], which could potentially reduce data streaming traffic costs in future implementations.

6.2. Adaptive Bitrate Streaming

As illustrated in Figure 8(b), `num_gs` and `swin_size` are controllable parameters that directly impact both bandwidth usage and rendering quality. This characteristic enables

adaptive bitrate control, facilitating a smooth streaming experience for users.

For example, during network congestion, we can reduce the number of Gaussians per frame by transmitting only a subset of Gaussians sampled from each slice with reset marked as empty padding. This approach helps prevent network overload while maintaining acceptable video quality.

6.3. Limitations and Future Work

As the first effort to adapt 3DGS from short 3D dynamic scenes to long volumetric videos, **SwiNGS** faces several challenges that warrant further investigation.

6.3.1. Inflexible Maturation Schedule. Current Gaussian maturation process follows a fixed schedule, where Gaussians mature after exactly `swin_size` frames to facilitate efficient batch operations in GPU memory. Although the Gaussian relocation mechanism and “birth” field allow some inter-slice migration, we have yet to fully optimize bandwidth usage with the most informative content. An ideal scenario would involve replacing the least useful Gaussians with the most informative new ones, rather than updating a preassigned fixed subset. For instance, within 40K Gaussians per frame, we could prioritize updates for Gaussians representing human motion while less frequently updating those depicting static backgrounds.

6.3.2. Prolonged Training Time. A significant challenge we face is the extended model training time. Using a single NVIDIA RTX 4090 with 24GB memory, training each sliding window of the DyNeRF dataset takes approximately 2 minutes. For a one-minute volumetric video (60 seconds, 30FPS, 1800 sliding windows), this translates to 60 GPU hours. To address this issue, we might need to either enhance training process parallelism or reduce the training time for individual sliding windows.

7. Conclusion

Drawing inspiration from recent advancements in neural rendering, our work, **SwiNGS**, adapts 3D Gaussian Splatting (3DGS) techniques to the challenging domain of volumetric video streaming. We first identify the unique challenges inherent in this task compared to previous dynamic 3DGS task. In response, we propose a novel method that employs a sliding window technique for training 3D Gaussian models and captures Gaussian snapshots for each frame in a slice-by-slice manner.

To demonstrate the efficacy of **SwiNGS**, we conduct experiments using various scenes from two distinct datasets. Furthermore, we develop a WebGL application capable of streaming volumetric video onto mobile devices, showcasing the practical applicability of our approach.

Our work represents a significant step forward in the realm of volumetric video streaming, leveraging the strengths of 3DGS: compact representation, high rendering quality, and rapid rendering speed. We believe **SwiNGS** opens up new avenues for research and development in this exciting field. As the demand for immersive and interactive visual experiences continues to grow, we anticipate that our contribution will catalyze further innovations in real-time volumetric video streaming.

Acknowledgements

Our **SwiNGS** viewer is built on the basis of Kevin Kwok's <https://antimatter15.com/splat/>

References

- [1] 2024. Gracia. Retrieved from <https://store.steampowered.com/app/2802520/Gracia/>
- [2] 2024. unlocking next-gen rendering: 3d gaussian splatting on pico 4 ultra. Retrieved from <https://developer.picoxr.com/news/3dgs-pico4ultra/>
- [3] 2024. metalsplatter for apple vision pro. Retrieved from <https://radiancefields.com/metalsplatter-for-apple-vision-pro>
- [4] Jeongmin Bae, Seoha Kim, Youngsik Yun, Hahyun Lee, Gun Bang, and Youngjung Uh. 2025. Per-gaussian embedding-based deformation for deformable 3d gaussian splatting. In *European Conference on Computer Vision*, 2025. 321–335.
- [5] Ang Cao and Justin Johnson. 2023. Hexplane: A fast representation for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. 130–141.
- [6] Yuanxing Duan, Fangyin Wei, Qiyu Dai, Yuhang He, Wenzheng Chen, and Baoquan Chen. 2024. 4D-Rotor Gaussian Splatting: Towards Efficient Novel View Synthesis for Dynamic Scenes. In *ACM SIGGRAPH 2024 Conference Papers*, 2024. 1–11.
- [7] Yongjie Guan, Xueyu Hou, Nan Wu, Bo Han, and Tao Han. 2023. Metastream: Live volumetric content capture, creation, delivery, and rendering in real time. In *Proceedings of the 29th Annual International Conference on Mobile Computing and Networking*, 2023. 1–15.
- [8] Bo Han, Yu Liu, and Feng Qian. 2020. ViVo: Visibility-aware mobile volumetric video streaming. In *Proceedings of the 26th annual international conference on mobile computing and networking*, 2020. 1–13.
- [9] Yi-Hua Huang, Yang-Tian Sun, Ziyi Yang, Xiaoyang Lyu, Yan-Pei Cao, and Xiaojuan Qi. 2024. Sc-gs: Sparse-controlled gaussian splatting for editable dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 4220–4230.
- [10] Mustafa Işık, Martin Rünz, Markos Georgopoulos, Taras Khakhulin, Jonathan Starck, Lourdes Agapito, and Matthias Nießner. 2023. Humanrf: High-fidelity neural radiance fields for humans in motion. *ACM Transactions on Graphics (TOG)* 42, 4 (2023), 1–12.
- [11] Euee S Jang, Marius Preda, Khaled Mammou, Alexis M Tourapis, Jungsun Kim, Danillo B Graziosi, Sungryeul Rhyu, and Madhukar Budagavi. 2019. Video-based point-cloud-compression standard in MPEG: From evidence collection to committee draft [standards in a nutshell]. *IEEE Signal Processing Magazine* 36, 3 (2019), 118–123.
- [12] Yuheng Jiang, Zhehao Shen, Penghao Wang, Zhuo Su, Yu Hong, Yingliang Zhang, Jingyi Yu, and Lan Xu. 2024. Hifi4g: High-fidelity human performance rendering via compact

- gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 19734–19745.
- [13] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 2023. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM Trans. Graph.* 42, 4 (2023), 139–131.
- [14] Shakiba Kheradmand, Daniel Rebain, Gopal Sharma, Weiwei Sun, Jeff Tseng, Hossam Isack, Abhishek Kar, Andrea Tagliasacchi, and Kwang Moo Yi. 2024. 3D Gaussian Splatting as Markov Chain Monte Carlo. *arXiv preprint arXiv:2404.09591* (2024).
- [15] Kyungjin Lee, Juheon Yi, Youngki Lee, Sunghyun Choi, and Young Min Kim. 2020. GROOT: a real-time streaming system of high-fidelity volumetric videos. In *Proceedings of the 26th Annual International Conference on Mobile Computing and Networking*, 2020. 1–14.
- [16] Lingzhi Li, Zhen Shen, Zhongshu Wang, Li Shen, and Ping Tan. 2022. Streaming radiance fields for 3d video synthesis. *Advances in Neural Information Processing Systems* 35, (2022), 13485–13498.
- [17] Tianye Li, Mira Slavcheva, Michael Zollhoefer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, Richard Newcombe, and others. 2022. Neural 3d video synthesis from multi-view video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 5521–5531.
- [18] Zhan Li, Zhang Chen, Zhong Li, and Yi Xu. 2024. Spacetime gaussian feature splatting for real-time dynamic view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 8508–8520.
- [19] Weikai Lin, Yu Feng, and Yuhao Zhu. 2024. RTGS: Enabling Real-Time Gaussian Splatting on Mobile Devices Using Efficiency-Guided Pruning and Foveated Rendering. *arXiv preprint arXiv:2407.00435* (2024).
- [20] Junhua Liu, Yuanyuan Wang, Yan Wang, Yufeng Wang, Shuguang Cui, and Fangxin Wang. 2023. Mobile volumetric video streaming system through implicit neural representation. In *Proceedings of the 2023 Workshop on Emerging Multimedia Systems*, 2023. 1–7.
- [21] Kaiyan Liu, Ruizhi Cheng, Nan Wu, and Bo Han. 2023. Toward next-generation volumetric video streaming with neural-based content representations. In *Proceedings of the 1st ACM Workshop on Mobile Immersive Computing, Networking, and Systems*, 2023. 199–207.
- [22] Yu Liu, Bo Han, Feng Qian, Arvind Narayanan, and Zhi-Li Zhang. 2022. Vues: Practical mobile volumetric video streaming through multiview transcoding. In *Proceedings of the 28th Annual International Conference on Mobile Computing And Networking*, 2022. 514–527.
- [23] Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. 2023. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. *arXiv preprint arXiv:2308.09713* (2023).
- [24] Marko Mihajlovic, Sergey Prokudin, Marc Pollefeys, and Siyu Tang. 2023. Resfields: Residual neural fields for spatiotemporal signals. *arXiv preprint arXiv:2309.03160* (2023).
- [25] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM* 65, 1 (2021), 99–106.
- [26] Kerui Ren, Lihan Jiang, Tao Lu, Mulin Yu, Linning Xu, Zhangkai Ni, and Bo Dai. 2024. Octree-gs: Towards consistent real-time rendering with lod-structured 3d gaussians. *arXiv preprint arXiv:2403.17898* (2024).
- [27] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. 2016. Pixelwise View Selection for Unstructured Multi-View Stereo. In *European Conference on Computer Vision (ECCV)*, 2016.
- [28] Liangchen Song, Anpei Chen, Zhong Li, Zhang Chen, Lele Chen, Junsong Yuan, Yi Xu, and Andreas Geiger. 2023. Nerfplayer: A streamable dynamic scene representation with decomposed neural radiance fields. *IEEE Transactions on Visualization and Computer Graphics* 29, 5 (2023), 2732–2742.
- [29] Patrick Stotko, Stefan Krumpfen, Max Schwarz, Christian Lenz, Sven Behnke, Reinhard Klein, and Michael Weinmann. 2019. A VR system for immersive teleoperation and live exploration with a mobile robot. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2019. 3630–3637.
- [30] Jiakai Sun, Han Jiao, Guangyuan Li, Zhanjie Zhang, Lei Zhao, and Wei Xing. 2024. 3dstream: On-the-fly training of 3d gaussians for efficient streaming of photo-realistic free-viewpoint videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 20675–20685.
- [31] Yuan-Chun Sun, Yuang Shi, Wei Tsang Ooi, Chun-Ying Huang, and Cheng-Hsin Hsu. 2024. Multi-frame Bitrate Allocation of Dynamic 3D Gaussian Splatting Streaming Over Dynamic Networks. In *Proceedings of the 2024 SIGCOMM Workshop on Emerging Multimedia Systems*, 2024. 1–7.
- [32] Varjo Technologies. 2024. Varjo Demonstrates Teleport, a Powerful New Service for Turning Real-World Places into Virtual Experiences. Retrieved from <https://varjo.com/press-release/varjo-demonstrates-teleport-a-powerful-new-service-for-turning-real-world-places-into-virtual-experiences/>
- [33] Diwen Wan, Ruijie Lu, and Gang Zeng. 2024. Superpoint Gaussian Splatting for Real-Time High-Fidelity Dynamic Scene Reconstruction. *arXiv preprint arXiv:2406.03697* (2024).
- [34] Max Welling and Yee W Teh. 2011. Bayesian learning via stochastic gradient Langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, 2011. 681–688.
- [35] Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Xinggang Wang. 2024. 4d gaussian splatting for real-time dynamic scene rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 20310–20320.

- [36] Zeyu Yang, Hongye Yang, Zijie Pan, Xiatian Zhu, and Li Zhang. 2023. Real-time photorealistic dynamic scene representation and rendering with 4d gaussian splatting. *arXiv preprint arXiv:2310.10642* (2023).
- [37] Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin. 2024. Deformable 3d gaussians for high-fidelity monocular dynamic scene reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 20331–20341.
- [38] Anlan Zhang, Chendong Wang, Bo Han, and Feng Qian. 2021. Efficient volumetric video streaming through super resolution. In *Proceedings of the 22nd International Workshop on Mobile Computing Systems and Applications*, 2021. 106–111.
- [39] Pengyuan Zhou, Jinjing Zhu, Yiting Wang, Yunfan Lu, Zixiang Wei, Haolin Shi, Yuchen Ding, Yu Gao, Qinglong Huang, Yan Shi, and others. 2022. Vetaverse: A survey on the intersection of Metaverse, vehicles, and transportation systems. *arXiv preprint arXiv:2210.15109* (2022).