# FaithFusion: Harmonizing Reconstruction and Generation via Pixel-wise Information Gain

**YuAn Wang**[1*] **Xiaofan Li**[1*,†] **Chi Huang**[1] **Wenhao Zhang**[1,2]
**Hao Li**[1] **Bosheng Wang**[1] **Xun Sun**[1] **Jun Wang**[1]

[1]Baidu Inc. [2]Nanjing University

Project page: https://shalfun.github.io/faithfusion

## Abstract

*In controllable driving-scene reconstruction and 3D scene generation, maintaining geometric fidelity while synthesizing visually plausible appearance under large viewpoint shifts is crucial. However, effective fusion of geometry-based 3DGS and appearance-driven diffusion models faces inherent challenges, as the absence of pixel-wise, 3D-consistent editing criteria often leads to over-restoration and geometric drift. To address these issues, we introduce **FaithFusion**, a 3DGS-diffusion fusion framework driven by pixel-wise Expected Information Gain (EIG). EIG acts as a unified policy for coherent spatio-temporal synthesis: it guides diffusion as a spatial prior to refine high-uncertainty regions, while its pixel-level weighting distills the edits back into 3DGS. The resulting plug-and-play system is free from extra prior conditions and structural modifications. Extensive experiments on the Waymo dataset demonstrate that our approach attains SOTA performance across NTA-IoU, NTL-IoU, and FID, maintaining an FID of 107.47 even at 6 meters lane shift. Our code is available at https://github.com/wangyuanbiubiubiu/FaithFusion.*

## 1. Introduction

Building a controllable driving world for closed-loop simulation [26, 43, 52] requires jointly achieving geometric fidelity in reconstruction and controllability in appearance generation. Neural rendering, particularly Neural Radiance Fields (NeRF) [1, 12, 33] and 3D Gaussian Splatting (3DGS) [23, 24, 30], has transformed 3D representations and novel view synthesis. However, under sparse observations, heavy occlusions, or viewpoints far from the training trajectory, NeRF and 3DGS often yield geometric inconsistencies and artifacts [5, 55, 61]. Diffusion mod-
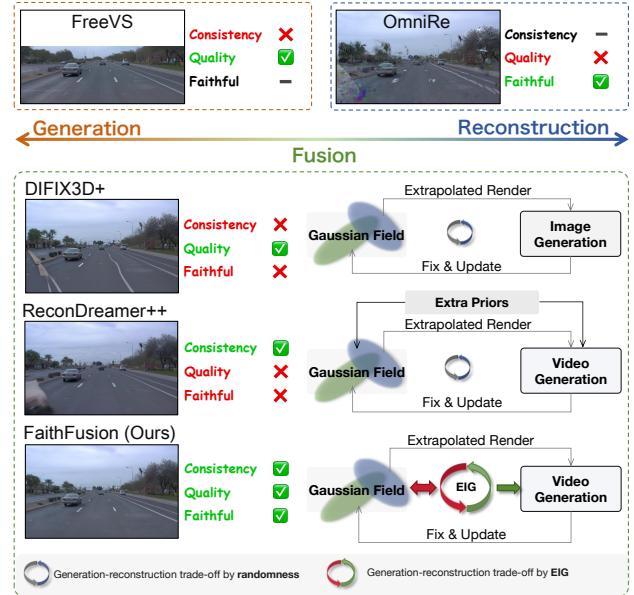


Figure 1. **Comparative overview.** Comparison of FreeVS [46], OmniRe [7], the fusion-based methods DIFIX3D+ [47] and ReconDreamer++ [64], and our EIG-integrated *FaithFusion*, which simultaneously achieves consistency, quality, and faithfulness.

els [16, 28, 40] excel at image and video generation and restoration, yet without pixel-level, geometry-consistent guidance, they tend to cause over-restoration and introduce geometric drift. Thus, balancing faithful reconstruction with controllable generation remains a central challenge for unified world modeling [59, 63, 64].

Unlike existing reconstruction or generation methods, fusion reconstruction–generation approaches primarily adopt an online progressive loop paradigm, following a "render–restoration–feedback" workflow designed to minimize information loss during 3DGS rendering, as shown in Fig. 1. Building upon this foundation, current efforts

---

*Equal contribution.
†Corresponding author: Shalfunnn@gmail.com

focus on optimizing both the generation and reconstruction branches: the former aims to improve temporal consistency via extra prior conditions [35, 46], while DI-FIX3D+ [47] achieves a balance between efficiency and quality through single-step, single-frame restoration. The reconstruction branch involves structural modifications to the 3DGS architecture or the introduction of robust geometry priors [64]. However, these explorations face common limitations: they predominantly rely on view-level heuristics to decide "where, when, and how much to edit". This dependence on coarse-grained guidance directly leads to insufficient control over generation; once diffusion is activated, it often overwrites already correct regions, causing over-restoration and geometric drift. Therefore, a principled decision mechanism that determines which regions to generate and which to preserve is still absent.

Motivated by these limitations, we introduce **FaithFusion**, a 3DGS-diffusion fusion paradigm driven by pixel-wise Expected Information Gain (EIG). Our core insight is to reformulate the decision of whether and how much to edit a pixel into a forward-looking information-theoretic metric—how much the edit reduces posterior uncertainty. To this end, we tightly couple a Laplace-approximated EIG with the differentiable 3DGS renderer to derive a pixel-level estimation specifically tailored for it. As illustrated in Fig. 2, *FaithFusion* builds upon this foundation by:

- **Generation branch:** EIG serves as a spatial weighting function that guides diffusion to generate content only in high-uncertainty regions, effectively suppressing over-restoration and geometric drift.
- **Reconstruction branch:** EIG functions as a pixel-wise loss weight, progressively distilling high-value edits back into 3DGS and forming a unified loop of edit triggering, strength modulation, and knowledge feedback.
- **System properties:** Without relying on external priors or modifying the 3DGS architecture, *FaithFusion* preserves trajectory fidelity while significantly improving spatio-temporal consistency and perceptual quality under large viewpoint shifts such as lane changes.

We conduct comprehensive evaluations on the Waymo dataset [41], achieving substantial improvements on NTA-IoU, NTL-IoU, and FID metrics. Extensive visualizations demonstrate precise corrections in under-constrained regions, minimal disturbance to original trajectories, and consistent behavior across diverse viewpoints. In summary, *FaithFusion* replaces heuristic decisions about "where, when, and how much to edit" with a principled, information-theoretic formulation: offering a concise, interpretable, and generalizable framework for unified, controllable 3D scene modeling.

## 2. Related Work

**Driving Scene Reconstruction.** The introduction of NeRF [2, 33, 34] and increasingly influential 3DGS [23, 30, 58] has brought high-quality novel view synthesis (NVS) to the forefront of autonomous driving research. A primary challenge is coping with the ubiquitous dynamics in driving scenes. Approaches differ in how they handle dynamic content: self-supervised or weakly supervised methods, whether NeRF-based [44, 51] or 3DGS-based [6, 20, 53], aim to avoid annotation but struggle to robustly disentangle dynamics under motion blur, illumination changes, and other complex cues. 3DGS, thanks to its explicit primitives, naturally supports foreground-background decoupling, and annotation-driven methods [7, 49, 65, 66] are widely adopted. Beyond dynamics, enhancing novel view rendering quality from sparse views remains a key challenge, often tackled through two main strategies. First, researchers introduce geometric priors like LiDAR depth [42, 62], pretrained depth maps [9, 45, 57, 60], and driving-specific ground parameterizations [12, 32, 39, 64, 65]. However, these priors can inject noise and bias. Second, methods use virtual view augmentation [5, 8] to enrich observations, which risks local minima due to dependence on accurate warping. Despite progress, large viewpoint shifts (such as lane-change) remain particularly challenging, as shape-radiance ambiguity and unobserved regions still lead to geometric inconsistencies and artifacts.

**Information Gain in Radiance Fields.** EIG quantifies the value of new observations based on model uncertainty, measuring the uncertainty reduction via information acquisition [25]. Radiance-field uncertainty estimation typically falls into three categories: *variational inference* [18, 36, 38, 48], which learns probabilistic distributions but demands architectural changes and high training cost; *Monte Carlo sampling* [31, 38], which captures uncertainty via sample dispersion but often relies on specific low-dimensional assumptions; and the *Laplace approximation* [11, 22], a post-hoc, architecture-agnostic, efficient option we adopt for its plug-and-play nature. In novel-view optimization and active mapping, ActiveNeRF [36] and FisherRF [22] already leverage EIG or its heuristics [29] for view selection, maximizing knowledge gain by measuring uncertainty reduction. Crucially, prior work largely remains at the view level; we extend FisherRF's theory to the pixel level, enabling fine-grained, interpretable diffusion guidance.

**Driving Scenes with Diffusion Priors.** Diffusion models have recently revolutionized image [16, 28, 37] and video [3, 17, 54] generation, offering a new paradigm for handling large viewpoint changes in autonomous driving. Coupling 3DGS reconstruction priors with diffusion is often preferred over direct pose-conditioned generation [27, 63], as it reduces the task to restoring degraded novel-view renderings [13, 21, 46, 59], which is more tractable. How-
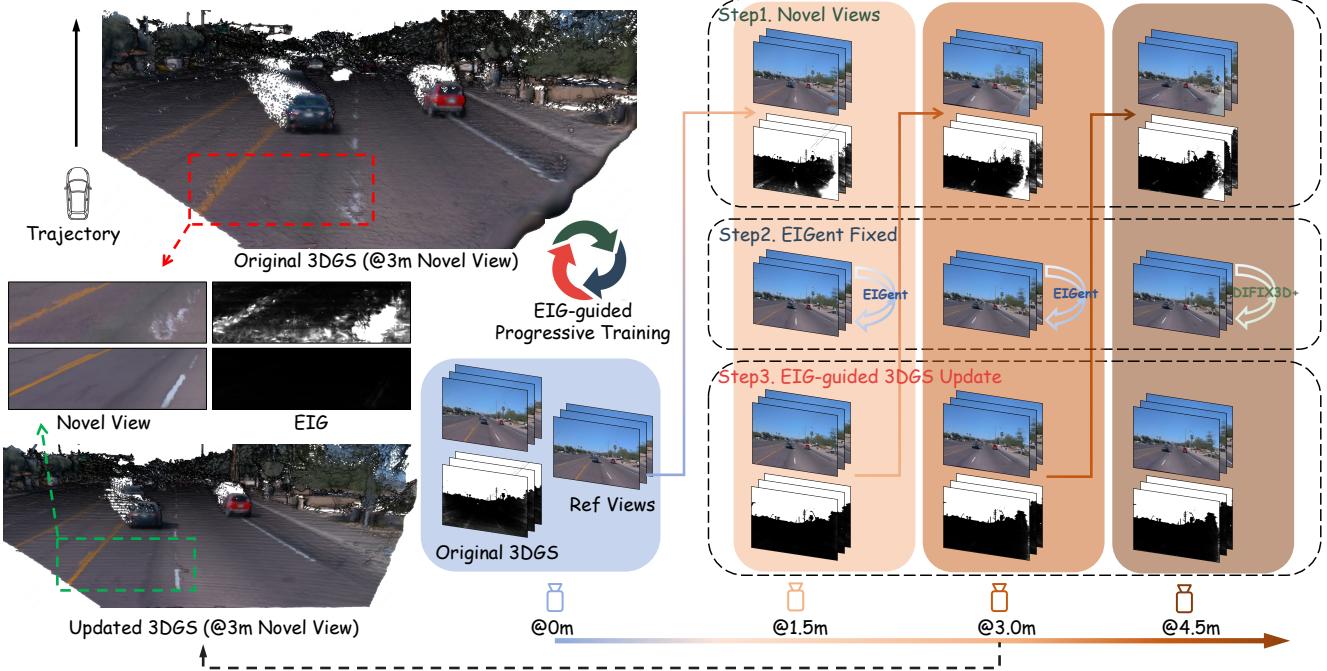
Figure 2. **FaithFusion pipeline.** The EIG-guided progressive training loop with three steps: **Step 1:** *Novel-view synthesis.* Render laterally offset novel views and their pixel-level EIG maps from the original 3DGS. **Step 2:** *EIGent Fixed.* Feed the renders and EIG maps into EIGent to repair high-EIG regions—using Video DiT early for spatio-temporal consistency and DIFIX3D+ later for per-frame perceptual refinement. **Step 3:** *EIG-guided 3DGS Update.* Fine-tune the 3DGS model with the EIGent-restored views and EIG maps.

ever, these approaches often depend on extra prior conditions, such as LiDAR [46, 50], sparse annotations [4, 35], or require 3DGS-specific adaptations to shrink the gap between generated results and reconstruction [64]. This dependence limits their generality. In contrast, we propose a novel reconstruction–generation fusion paradigm driven by the intrinsic EIG of 3DGS. This information-theoretic metric replaces heuristics and strong geometric conditions, significantly enhancing spatio-temporal consistency and perceptual quality under large viewpoint shifts.

## 3. Method

Our primary goal is to develop a method capable of synthesizing high-fidelity, temporally and spatially consistent 4D representations for challenging, far-field novel viewpoints (e.g., lane changes), where traditional 3DGS-based methods often fail. To this end, we propose *FaithFusion*, a controllable 3DGS–diffusion fusion framework driven by pixel-wise Expected Information Gain (EIG). Under a progressive fusion scheme [14, 35, 47], EIG serves as a unified spatial policy: it first directs the diffusion model to synthesize high-information regions in novel views (e.g., unseen areas), and then guides selective 3DGS fine-tuning to assimilate the generated content, yielding coherent generation–reconstruction coupling.

We describe pixel-wise EIG computation (Sec. 3.1), the dual-branch *EIGent* guidance (Sec. 3.2), and the EIG-driven progressive knowledge integration (Sec. 3.3); the overall pipeline is shown in Fig. 2. The framework is *plug-and-play* and can be seamlessly integrated into mainstream street-scene 3DGS systems [6, 7, 49, 53].

### 3.1. Expected Information Gain in 3DGS

**3D Gaussian Splatting Parameterization.** Original 3DGS [23] represents a static scene with a set of anisotropic Gaussians parameterized by world-space position $\boldsymbol{\mu}_w \in \mathbb{R}^3$, rotation $\mathbf{q}_w \in \mathbb{R}^4$, and scale $\mathbf{s} \in \mathbb{R}^3$. For dynamic street scenes, existing approaches further augment object parameterizations [6, 7, 49, 53] to capture motion and deformation. To model view-dependent appearance, each Gaussian maintains spherical harmonic coefficients $\mathbf{c} \in \mathbb{R}^k$ and opacity $o \in \mathbb{R}$. Together, these parameters $\omega$ yield photorealistic renderings at target timestamps after $\alpha$-blending and projection onto the image plane:

$$\mathbf{C} = \sum_{i \in \mathcal{M}} \mathbf{c}_i \alpha_i' \prod_{j=1}^{i-1} \left(1 - \alpha_j'\right), \tag{1}$$

where $\mathcal{M}$ denotes the set of Gaussians intersected by the ray, ordered by depth, and $\alpha_i'$ is determined by the opacity $o$ and the 2D Gaussian after linear projection.

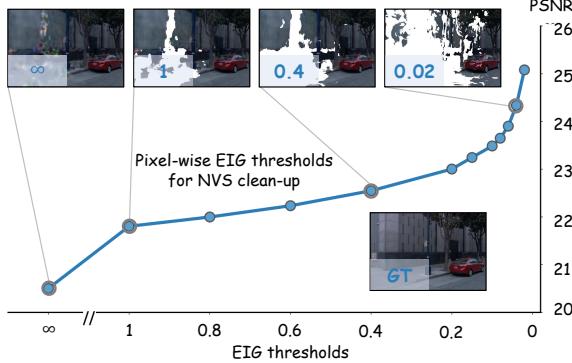**Pixel-wise Expected Information Gain with 3DGS.**

3

Figure 3. **Image quality vs. EIG mask threshold.** We validate pixel-level EIG as a proxy for novel-view synthesis quality by progressively retaining high-EIG regions and evaluating PSNR. The consistent decrease in PSNR as high-EIG regions are retained confirms that higher EIG marks lower-quality rendering.

Given a training set $D_{train}$ containing views $X_{train}$ (camera parameters) and corresponding image sequences $Y_{train}$, the differentiable renderer $\mathcal{F}$ of 3DGS is optimized over parameters $\omega$. The goal is to obtain the point estimate $\omega^*$ that minimizes the negative log-likelihood between rendered images $\hat{Y}_i^{train} = \mathcal{F}(X_i^{train}, \omega)$ and real observations $Y_i^{train}$. Assuming Gaussian observation noise, this is equivalent to minimizing the reconstruction error:

$$\omega^* = \arg\min_{\omega} \sum_{(X_i, Y_i) \in \mathcal{D}_{train}} \|Y_i^{train} - \mathcal{F}(X_i^{train}, \omega)\|_2^2 \tag{2}$$

To quantify the uncertainty of $\omega^*$ and estimate potential information gain from new observations, we apply the Laplace approximation [10], modeling the posterior with a Gaussian distribution $\Omega$:

$$\Omega \approx \mathcal{N}(\omega^*, H''[\omega^*]^{-1}), \tag{3}$$

where $H''[\omega^*]$ is the Hessian of the negative log-likelihood at $\omega^*$. The expectation of this Hessian under the predictive distribution corresponds to the Fisher information, quantifying how strongly the observations constrain the parameters. FisherRF [22] treats the Fisher information as an uncertainty proxy, integrating it into the $\alpha$-blending process to compute pixel-level uncertainty at training views, but its application is limited to the observed training data.

For novel views $X_{NVS}$, the optimal 3DGS parameters $\omega^*$ render novel-view sequences $Y_{NVS}$. The corresponding Expected Information Gain, defined as $EIG(\Omega; Y_i^{NVS}|X_i^{NVS})$, is the difference between the prior entropy of $\Omega$ and the expected posterior entropy after observing $Y_i^{NVS}$ [19]:

$$EIG = \mathbb{H}[\Omega] - \mathbb{E}_{p(Y_i|X_i)}[\mathbb{H}[\Omega|Y_i^{NVS}, X_i^{NVS}]]. \tag{4}$$

where $\mathbb{H}[\cdot]$ denotes the differential entropy of the distribution, and $\mathbb{E}_{p(\cdot)}[\cdot]$ denotes the expectation with respect to the

predictive distribution $p(Y_i^{NVS}|X_i^{NVS})$.

Leveraging the Laplace approximation and properties of the Fisher information, this quantity can be efficiently estimated. To obtain a compact and computable upper bound, we apply the inequality $\log \det(A + I_d) \leq \operatorname{tr}(A)$ (Lemma 5.1) [25] and exploit the additive property of Fisher information, leading to the trace form[*]:

$$EIG \leq \frac{1}{2} \sum_i \operatorname{tr}\left(H''[Y_i^{NVS}|X_i^{NVS}, \omega^*]H''[\omega^*]^{-1}\right). \tag{5}$$

While Eq. (5) yields the aggregate EIG for a full novel view, EIGent requires pixel-level granularity to guide local edits precisely. We extend EIG to each pixel by accumulating the Fisher information contributions of Gaussians intersected along each rendering ray. Algorithm 1 outlines the computation process: accumulating global Fisher information during training, and computing and mapping pixel-wise EIG for novel views. Following the methodology of BayesRays [11], and as illustrated in Fig. 3, our cross-camera evaluations on the Waymo dataset [41] demonstrate that this pixel-level EIG strongly correlates with novel view synthesis quality[*].

---

**Algorithm 1** Pixel-wise Expected Information Gain Computation in 3DGS

---

**Input:** Trained 3DGS model parameters $\omega^*$, Total number of 3D Gaussians $N$, Train views $X_{train}$, Novel views $X_{NVS}$, Differentiable rasterization $\mathcal{F}$, Hessian computation function(Laplace diag, per-Gaussian) $\mathcal{H}$

**Output:** Pixel-wise Expected Information Gain $EIG_{NVS}$

1: $H''[\omega^*] \leftarrow \mathbf{0} \in \mathbb{R}^N$
2: **for** each train view $X_i^{train}$ in $X_{train}$ **do**
3:      $H''[\omega^*] \leftarrow H''[\omega^*] + \mathcal{H}[\mathcal{F}(X_i^{train}, \omega^*)]$
4: **end for**
5: **for** each novel view $X_j^{NVS}$ in $X_{NVS}$ **do**
6:      $H''[Y_j^{NVS}|X_j^{NVS}, \omega^*] \leftarrow \mathcal{H}[\mathcal{F}(X_j^{NVS}, \omega^*)]$
7:      $EIG_j^{GS} \leftarrow H''[Y_j^{NVS}|X_j^{NVS}, \omega^*] \odot (H''[\omega^*])^{-1}$
8:      $EIG_j^{NVS} = \mathcal{F}((X_j^{NVS}, EIG_j^{GS}), \omega^*)$
9: **end for**

---

### 3.2. EIGent: EIG Guided Dual-Branch Controllable Generation

**Dataset Generation.** To train the EIG-guided video restoration task, we construct datasets following the cross-camera referencing strategy of [47]. As illustrated in Fig. 4, we first train a street-scene 3DGS model on forward-camera sequences, then render it from right-front view. This process simultaneously produces novel-view renderings and corresponding EIG maps computed via Algorithm 1, which are paired with real observations to form triplets for the restoration task. To ensure data validity, we filter the raw

---

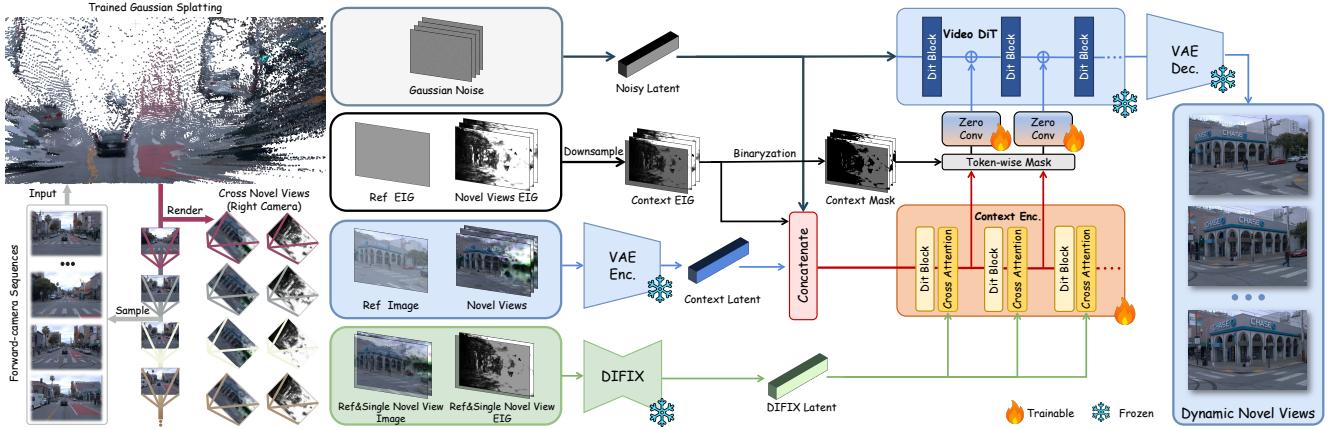[*]Please refer to the supplemental material for more details.

Figure 4. **Overview of EIGent.** *Data:* Cross-view pairing: a forward-camera–trained 3DGS renders right-front views to produce artifact-prone novel-view renders and per-pixel EIG (Alg. 1), temporally aligned with real right-front videos. *Architecture:* EIGent is a dual-branch model with coarse-to-fine EIG guidance: downsampled $E$, noise latent $L_N$, and VAE latent $L$ feed a lightweight context encoder $\mathcal{G}$; a mask $M$ suppresses high–EIG regions. Via cross-attention with a DIFIX branch, cues are injected into a pretrained DiT backbone, enabling EIG-aware controllable repair and foreground spatio-temporal consistency.

sequences. We remove nearly stationary clips and compute inter-camera overlap (based on FOV) to guarantee sufficient reconstruction signals and adequate cross-view coverage. Furthermore, we mitigate EIG estimation distortion caused by large floaters near unseen regions by imposing explicit scale constraints on 3D Gaussians, confining artifacts within controllable bounds.

**EIGent architecture.** EIGent employs EIG as a spatial prior for controllable restoration, providing interpretable, pixel-wise priorities:

- **High-gain regions.** Areas of low rendering quality or missing information that require focused restoration and content generation.
- **Low-gain regions.** Backgrounds with reliable content where the model should preserve original structures.

To incorporate this prior effectively, we design a dual-branch control architecture, shown in Fig. 4. A lightweight EIG-guided context encoder operates alongside a pretrained DiT backbone,

decoupling stable background preservation from temporally consistent foreground generation.

Given an input video $V$, a VAE encoder $\mathcal{E}$ maps it to a latent $L = \mathcal{E}(V)$, and the pixel-wise EIG map is downsampled to $E$. The multi-scale guidance strategy fuses these signals through EIG-guided context injection:

$$\epsilon_\theta(z_t, t, C)_k = \epsilon\theta(z_t, t, C)_k + M \odot \mathcal{G}(L_N, L, E)_k, \quad (6)$$

where $\epsilon_\theta$ denotes the DiT denoiser, $z_t$ the noisy latent, $t$ the timestep, $C$ the conditioning input, $\odot$ the Hadamard product, $L_N$ the noise latent, and $k$ the feature layer index. $\mathcal{G}$ represents the lightweight EIG-guided context encoder, cloned from the first four layers of the pretrained DiT.

Furthermore, to enhance per-frame quality, we fuse external repair cues (e.g., DIFIX latent) with the context branch $\mathcal{G}$ via cross-attention, regulating the fusion with spatial weights from $E$ and a binary mask $M$ that filters regions of extreme uncertainty (e.g., EIG above a threshold) and selectively incorporates reliable cues across multiple scales. This dual control injects coarse spatial metadata via $E$ while ensuring only trustworthy, background-relevant information reaches the DiT backbone, preventing contamination of stable contexts.

Overall, this EIG-driven coarse-to-fine guidance and fusion strategy enables full exploitation of the DiT architecture, improving both perceptual quality and spatio-temporal consistency in restored video.

### 3.3. Progressive EIG-Aware Diffusion-to-3DGS Knowledge Integration

To fully exploit the restoration capability of EIGent and inject high-quality generated content into 3DGS in an orderly and controllable manner, we adopt a progressive knowledge integration strategy inspired by prior work [35, 47]. Different from existing approaches that rely on heuristics or view-level control, our key idea is to use **pixel-wise EIG** as the guiding signal, enabling finer-grained and more interpretable fusion.

Taking a standard street-scene 3DGS pipeline [7, 49] as an example: the overall loss is composed of the original-trajectory term and the novel-trajectory term. The loss for the original trajectory is

$$\mathcal{L}_{ori}(\omega) = \lambda_r \mathcal{L}1^{ori} + (1 - \lambda_r)\mathcal{L}_{SSIM}^{ori} + \lambda_d \mathcal{L}_{depth}^{ori}, \quad (7)$$

where $\mathcal{L}_1$ and $\mathcal{L}_{SSIM}$ denote the L1 and SSIM losses be-

Figure 5. **Qualitative comparison on Waymo [41].** Novel-view renderings for the same trajectory across representative methods [35, 46, 47, 64]. Orange boxes highlight regions where our approach yields noticeably better results.

tween rendered images and ground-truth views, $\mathcal{L}_{depth}$ is the depth supervision against sparse LiDAR, and $\lambda_r, \lambda_d$ are the corresponding weights.

To incorporate the restored novel views into the optimization, we introduce a novel-view loss $\mathcal{L}_{novel}(\omega)$ with two components:

- **EIG pixel-wise weighting.** The normalized EIG map is used as a pixel-level weight matrix $\lambda_{EIG}$ to modulate the image loss, so that 3DGS focuses optimization on the regions with the highest information gain (i.e., the most under-constrained areas):

$$\mathcal{L}_{img}^{novel} = \lambda_{EIG} \odot \left( \lambda_r \mathcal{L}_1^{novel} + (1 - \lambda_r) \mathcal{L}_{SSIM}^{novel} \right) \quad (8)$$

- **Sparse depth supervision.** We further employ point-cloud projections aggregated from neighboring frames as sparse depth supervision $\mathcal{L}_{depth}^{novel}$ to preserve geometric consistency in novel views.

The novel-trajectory loss is thus

$$\mathcal{L}_{novel}(\omega) = \mathcal{L}_{img}^{novel} + \lambda_d \mathcal{L}_{depth}^{novel}, \quad (9)$$

and is used to fine-tune the existing 3DGS model. This fine-tuning relies on EIGent-restored views in the early stage to prioritize spatial structure and cross-frame coherence, followed by DIFIX3D+ refined views once the expansion reaches its maximum range and stabilizes.

## 4. Experiments

### 4.1. Experiment Setup

**3DGS Training and Evaluation Protocol.** We conduct our experiments on the Waymo dataset [41], strictly following the experimental protocol established by Recon-Dreamer [35]. Specifically, the 3DGS model is trained on 8 distinct clips (40 frames per clip) using only the forward-camera data. During evaluation, we assess the cross-lane rendering quality. The training trajectory is expanded progressively by 1.0 m every 2,000 iterations starting from step 3,000 for synthesizing novel viewpoints.

**EIGent Model Fine-Tuning and Data Preparation.** The video generation components, the external DIFIX model and the dual-branch video generator within EIGent, are fine-tuned using a dedicated dataset. We preprocess the first 200 training clips of the Waymo dataset [41] for this purpose. Applying the data filtering strategy detailed in Sec. 3.2, this process yields 936 video triplets. The image-generation branch (DIFIX) is fine-tuned using LoRA under the default DIFIX3D+ [47] configuration. Meanwhile, for the video branch (EIGent), we freeze the base diffusion model (CogVideo-5B-I2V [54]) and train only the context encoder for 14,336 steps with a learning rate of $1 \times 10^{-5}$.

**Baselines.** To comprehensively evaluate *FaithFusion*, we

6

| Method | Extra Condition | | | Lane Shift @ 3m | | | Lane Shift @ 6m | | |
|---|---|---|---|---|---|---|---|---|---|
| | LiDAR | Box | HDMap | NTA-IoU ↑ | NTL-IoU ↑ | FID↓ | NTA-IoU ↑ | NTL-IoU ↑ | FID↓ |
| OmniRe [7] | × | × | × | 0.424 | 51.73 | 188.42 | 0.423 | 49.08 | 191.00 |
| FreeVS [46] | ✓ | × | × | 0.505 | 56.84 | 104.23 | 0.465 | 55.37 | 121.44 |
| ReconDreamer [35] | × | ✓ | ✓ | 0.539 | 54.58 | 93.56 | 0.467 | 52.58 | 149.19 |
| ReconDreamer++ [64]* | × | ✓ | ✓ | 0.572 | <u>57.06</u> | <u>72.02</u> | 0.489 | **56.57** | <u>111.92</u> |
| DIFIX3D+ [47] | × | × | × | <u>0.578</u> | 56.94 | 84.12 | <u>0.504</u> | 53.77 | 120.24 |
| FaithFusion | × | × | × | **0.581** | **57.67** | **71.51** | **0.517** | <u>55.78</u> | **107.47** |

Table 1. Comparison of different lane shifts on the Waymo dataset [41], highlighting key methodology requirements. Extra Condition indicates the reliance on additional data injected as a condition to guide the synthesis process (e.g., LiDAR, 3D boxes, HDMap). ∗ denotes that the method requires significant architectural or geometrical modifications, including decomposed modeling and new trajectory field.

integrate it into the general 3DGS reconstruction framework OmniRe [7] and compare it with FreeVS [46] as a representative generation method, as well as three fusion-based novel view synthesis methods, including Recon-Dreamer [35], ReconDreamer++ [64], and DIFIX3D+ [47]. **Evaluation Metrics.** Following DriveDreamer4D [63], we report Novel Trajectory Agent IoU (NTA-IoU), Novel Trajectory Lane IoU (NTL-IoU), and FID as the primary evaluation metrics. We also introduce two EIG-partitioned metrics in our ablation study: FID-UCR, assessing Under-Constrained Regions (UCR), and FID-HPR, for High-Confidence Regions (HPR).

## 4.2. Comparison with State-of-the-Art Methods

Comparison among representative SOTA methods is shown in Tab. 1. For a fair comparison, we crop the outputs of FreeVS [46] to exclude regions without LiDAR coverage when computing metrics. Benefiting from the plug-and-play nature of DIFIX3D+ [47], we integrate it into OmniRe [7] using the official gsplat [56] interface.

At a lane shift of 3 meters, *FaithFusion* approaches the best IoU scores and achieves the lowest FID (71.51), demonstrating strong 3D semantic stability and visual generalization. This robustness stems from EIG's dual role in the video-generation and DIFIX3D+ [47] restoration stages, and from the differentiated constraints of the EIG-weighted reconstruction loss. High-confidence regions that are co-visible with the original trajectory are refined for details while maintaining structural coherence, whereas under-constrained regions are guided to generate plausible geometry and semantics. As a result, progressive updates yield stable and high-quality novel-view generation even at moderate trajectory deviations. Visualization results (upper half of Fig. 5) show that road structures from the original view are well preserved under EIG guidance; in synthesized regions, the EIG map clarifies repair needs, enabling semantically correct completion (e.g., building facades) and reducing 3D semantic inconsistencies.

| Method | FID ↓ | | |
|---|---|---|---|
| | Total | UCR | HPR |
| DIFIX3D+ (Baseline) | 120.24 | 147.97 | 152.66 |
| + EIG Guided DIFIX3D+ | 119.01 | 143.80 | 149.82 |
| ++ EIGent Dual-Stage Fusion | 113.94 | 137.58 | 153.69 |
| **+++ EIG Recon (Full FaithFusion)** | **107.47** | **137.02** | **147.75** |

Table 2. **Ablation Study: Incremental Contributions of Faith-Fusion's Core Components.** Results on the most challenging 6-meter lateral-shift novel-view synthesis task, showing the gain from sequentially adding the three proposed EIG-guided modules to the DIFIX3D+ baseline.

At a lane shift of 6 meters, most methods suffer severe performance degradation due to accumulated errors, whereas *FaithFusion* maintains stable performance (NTA-IoU: 0.517, NTL-IoU: 55.78, FID: 107.47). Here, EIG guides the model to focus on critical unseen structures during generation to ensure geometric fidelity, while constraining high-confidence regions during repair to seamlessly integrate new content with existing areas, thus achieving a balance between global coherence and fine-detail fidelity. Visualization results (lower half of Fig. 5) illustrate these effects: while competing methods often produce blurred or semantically inconsistent results, EIGent—guided explicitly by EIG and enhanced by robust video generation—precisely localizes and restores missing regions. Our results exhibit state-of-the-art global coherence, avoiding spurious deformations such as the erroneous ground bending in the third row, and preventing unrealistic artifacts often observed when using DIFIX3D+ [47] alone.

## 4.3. Ablation Study

We conduct comprehensive ablations to quantify the contribution of each component in *FaithFusion*. We argue that traditional global FID is insufficient to accurately mea-
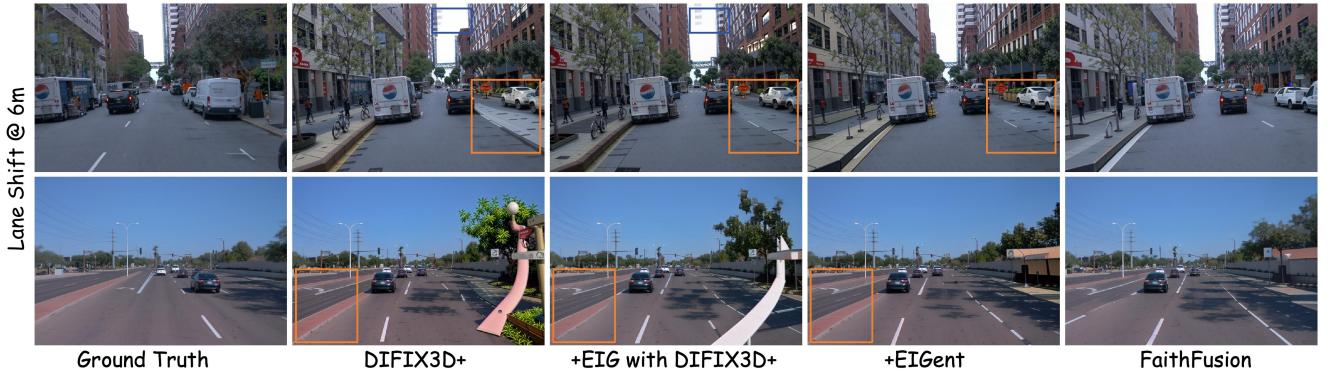
Figure 6. **Ablation overview.** We incrementally integrate EIG-guided components into the OmniRe baseline. The results highlight the incremental contributions of EIG guidance in resolving over-restoration and geometric drift by acting as a unified pixel-wise editing policy.
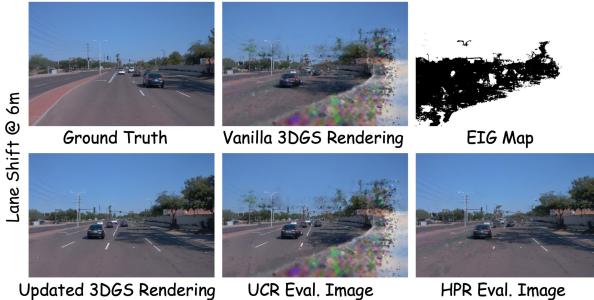


Figure 7. **Region partition for EIG-based evaluation.** We use EIG as a proxy for rendering quality. With a threshold $\tau = 0.4$, renderings are partitioned into UCR and HPR, and we report region-specific metrics: FID-UCR and FID-HPR.

sure fine-grained performance differences across regions of varying confidence. Given that pixel-wise EIG effectively reflects local rendering quality, we introduce two EIG-partitioned metrics to distinguish performance across uncertainty levels. Specifically, we define Under-Constrained Regions (UCR) and High-Confidence Regions (HPR), with the corresponding metrics being FID-UCR and FID-HPR. These regions are defined by the EIG threshold detailed in Fig. 7. During evaluation, the non-evaluated region is filled with vanilla 3DGS renderings to ensure complete inputs.

Overall, the improvements of our full system stem from three complementary components operating at different stages: As shown in Tab. 2, introducing EIG guidance consistently reduces FID, with a total drop of about 1.23. Compared with pure DIFIX3D+ [47], EIG focuses the repair on low-EIG (well-reconstructed) regions and suppresses unnecessary hallucinations, thereby improving overall 3DGS performance. The visualization results this trend: DIFIX3D+ [47] with EIG guidance alleviates semantic mismatches in low-EIG areas, as visually confirmed by Fig. 6,

although high-EIG regions may still exhibit deviations due to the lack of strong priors.

Adding EIGent further enhances generation quality across both partitioned regions, reducing the overall FID by 5.07 and lowering FID-UCR by 6.22. This indicates that EIGent effectively compensates for weaknesses in novel-view synthesis and makes the outputs closer to real images.

A slight increase in FID-HPR is observed, which mainly comes from enforcing temporal consistency in video diffusion, where stronger consistency can dampen fine-grained appearance details. Nonetheless, the visualization results, as shown in the lower half of Fig. 6, show that EIGent produces semantically more coherent and plausible content, especially in regions that require completion or semantic reasoning.

Finally, the progressive EIG-aware diffusion-to-GS integration reduces the total FID to **107.47**, improving by 12.77 over the baseline. FID-UCR and FID-HPR decrease by 10.95 and 4.91, respectively. This strategy retains EIGent's benefits while preventing over-restoration in low-EIG regions, ensuring structural and appearance consistency. Visualization results further confirm these gains, showing a balanced trade-off between detail restoration and global coherence.

## 5. Conclusion

We introduce *FaithFusion*, a 3DGS–diffusion fusion paradigm driven by pixel-wise Expected Information Gain (EIG), which unifies faithful reconstruction and controllable generation by converting heuristic editing decisions into an information-theoretic quantity. This cross-modal EIG guidance employs EIG as a spatial weight for content suppression on the generation side and as a loss weight for selective knowledge distillation on the reconstruction side, offering strong generality and interpretability. Systematic experiments on the Waymo dataset [41] demonstrate that

*FaithFusion* significantly improves spatio-temporal consistency and perceptual quality under large viewpoint shifts, achieving state-of-the-art results across major metrics.

**Limitations and Future Work.** While EIG effectively slows the accumulation of errors inherent in the 3DGS–diffusion fusion paradigm, the issue is not fully eliminated; this suggests that customizing the 3DGS model architecture could be key to further reduction. Furthermore, EIG is widely utilized in active exploration and mapping, and its incorporation into *FaithFusion* provides a natural bridge for future work, enabling active mapping strategies to significantly enhance overall efficiency.

# References

[1] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5855–5864, 2021. 1

[2] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Zip-nerf: Anti-aliased grid-based neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19697–19705, 2023. 2

[3] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22563–22575, 2023. 2

[4] Anthony Chen, Wenzhao Zheng, Yida Wang, Xueyang Zhang, Kun Zhan, Peng Jia, Kurt Keutzer, and Shangbang Zhang. Geodrive: 3d geometry-informed driving world model with precise action control. *arXiv preprint arXiv:2505.22421*, 2025. 3

[5] Danpeng Chen, Hai Li, Weicai Ye, Yifan Wang, Weijian Xie, Shangjin Zhai, Nan Wang, Haomin Liu, Hujun Bao, and Guofeng Zhang. Pgsr: Planar-based gaussian splatting for efficient and high-fidelity surface reconstruction. *IEEE Transactions on Visualization and Computer Graphics*, 2024. 1, 2

[6] Yurui Chen, Chun Gu, Junzhe Jiang, Xiatian Zhu, and Li Zhang. Periodic vibration gaussian: Dynamic urban scene reconstruction and real-time rendering. *arXiv preprint arXiv:2311.18561*, 2023. 2, 3

[7] Ziyu Chen, Jiawei Yang, Jiahui Huang, Riccardo de Lutio, Janick Martinez Esturo, Boris Ivanovic, Or Litany, Zan Gojcic, Sanja Fidler, Marco Pavone, et al. Omnire: Omni urban scene reconstruction. *arXiv preprint arXiv:2408.16760*, 2024. 1, 2, 3, 5, 7

[8] Kai Cheng, Xiaoxiao Long, Wei Yin, Jin Wang, Zhiqiang Wu, Yuexin Ma, Kaixuan Wang, Xiaozhi Chen, and Xuejin Chen. Uc-nerf: Neural radiance field for under-calibrated multi-view cameras in autonomous driving. *arXiv preprint arXiv:2311.16945*, 2023. 2

[9] Jaeyoung Chung, Jeongtaek Oh, and Kyoung Mu Lee. Depth-regularized optimization for 3d gaussian splatting in few-shot images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 811–820, 2024. 2

[10] Erik Daxberger, Agustinus Kristiadi, Alexander Immer, Runa Eschenhagen, Matthias Bauer, and Philipp Hennig. Laplace redux-effortless bayesian deep learning. *Advances in neural information processing systems*, 34:20089–20103, 2021. 4

[11] Lily Goli, Cody Reading, Silvia Sellán, Alec Jacobson, and Andrea Tagliasacchi. Bayes' rays: Uncertainty quantification for neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20061–20070, 2024. 2, 4

[12] Jianfei Guo, Nianchen Deng, Xinyang Li, Yeqi Bai, Botian Shi, Chiyu Wang, Chenjing Ding, Dongliang Wang, and Yikang Li. Streetsurf: Extending multi-view implicit surface reconstruction to street views. *arXiv preprint arXiv:2306.04988*, 2023. 1, 2

[13] Huasong Han, Kaixuan Zhou, Xiaoxiao Long, Yusen Wang, and Chunxia Xiao. Ggs: Generalizable gaussian splatting for lane switching in autonomous driving. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 3329–3337, 2025. 2

[14] Ayaan Haque, Matthew Tancik, Alexei A Efros, Aleksander Holynski, and Angjoo Kanazawa. Instruct-nerf2nerf: Editing 3d scenes with instructions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19740–19750, 2023. 3

[15] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 1

[16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 1, 2

[17] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646, 2022. 2

[18] Matthew D Hoffman, Tuan Anh Le, Pavel Sountsov, Christopher Suter, Ben Lee, Vikash K Mansinghka, and Rif A Saurous. Probnerf: Uncertainty-aware inference of 3d shapes from 2d images. In *International Conference on Artificial Intelligence and Statistics*, pages 10425–10444. PMLR, 2023. 2

[19] Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. Bayesian active learning for classification and preference learning. *arXiv preprint arXiv:1112.5745*, 2011. 4, 2

[20] Nan Huang, Xiaobao Wei, Wenzhao Zheng, Pengju An, Ming Lu, Wei Zhan, Masayoshi Tomizuka, Kurt Keutzer, and Shanghang Zhang. $s^3$gaussian: Self-supervised

street gaussians for autonomous driving. *arXiv preprint arXiv:2405.20323*, 2024. 2

[21] Sungwon Hwang, Min-Jung Kim, Taewoong Kang, Jayeon Kang, and Jaegul Choo. Vegs: View extrapolation of urban scenes in 3d gaussian splatting using learned priors. In *European Conference on Computer Vision*, pages 1–18. Springer, 2024. 2

[22] Wen Jiang, Boshu Lei, and Kostas Daniilidis. Fisherrf: Active view selection and mapping with radiance fields using fisher information. In *European Conference on Computer Vision*, pages 422–440. Springer, 2024. 2, 4

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

[24] Shakiba Kheradmand, Daniel Rebain, Gopal Sharma, Weiwei Sun, Yang-Che Tseng, Hossam Isack, Abhishek Kar, Andrea Tagliasacchi, and Kwang Moo Yi. 3d gaussian splatting as markov chain monte carlo. *Advances in Neural Information Processing Systems*, 37:80965–80986, 2024. 1

[25] Andreas Kirsch and Yarin Gal. Unifying approaches in active learning and active sampling via fisher information and information-theoretic quantities. *arXiv preprint arXiv:2208.00549*, 2022. 2, 4, 3

[26] Xiaofan Li, Yifu Zhang, and Xiaoqing Ye. Drivingdiffusion: Layout-guided multi-view driving scenarios video generation with latent diffusion model. In *European Conference on Computer Vision*, pages 469–485. Springer, 2024. 1

[27] Xiaofan Li, Chenming Wu, Zhao Yang, Zhihao Xu, Yumeng Zhang, Dingkang Liang, Ji Wan, and Jun Wang. Driverse: Navigation world model for driving simulation via multi-modal trajectory prompting and motion alignment. In *Proceedings of the 33rd ACM International Conference on Multimedia*, pages 9753–9762, 2025. 2

[28] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022. 1, 2

[29] Xi Liu, Chaoyi Zhou, and Siyu Huang. 3dgs-enhancer: Enhancing unbounded 3d gaussian splatting with view-consistent 2d diffusion priors. *Advances in Neural Information Processing Systems*, 37:133305–133327, 2024. 2

[30] Tao Lu, Mulin Yu, Linning Xu, Yuanbo Xiangli, Limin Wang, Dahua Lin, and Bo Dai. Scaffold-gs: Structured 3d gaussians for view-adaptive rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20654–20664, 2024. 1, 2

[31] Linjie Lyu, Ayush Tewari, Marc Habermann, Shunsuke Saito, Michael Zollhöfer, Thomas Leimkühler, and Christian Theobalt. Manifold sampling for differentiable uncertainty in radiance fields. In *SIGGRAPH Asia 2024 Conference Papers*, pages 1–11, 2024. 2

[32] Ruohong Mei, Wei Sui, Jiaxin Zhang, Xue Qin, Gang Wang, Tao Peng, Tao Chen, and Cong Yang. Rome: Towards large scale road surface reconstruction via mesh representation. *IEEE Transactions on Intelligent Vehicles*, 2024. 2

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

[34] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM transactions on graphics (TOG)*, 41(4):1–15, 2022. 2

[35] Chaojun Ni, Guosheng Zhao, Xiaofeng Wang, Zheng Zhu, Wenkang Qin, Guan Huang, Chen Liu, Yuyin Chen, Yida Wang, Xueyang Zhang, et al. Recondreamer: Crafting world models for driving scene reconstruction via online restoration. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 1559–1569, 2025. 2, 3, 5, 6, 7, 4

[36] Xuran Pan, Zihang Lai, Shiji Song, and Gao Huang. Activenerf: Learning where to see with uncertainty estimation. In *European Conference on Computer Vision*, pages 230–246. Springer, 2022. 2

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

[38] Luca Savant, Diego Valsesia, and Enrico Magli. Modeling uncertainty for gaussian splatting. *arXiv preprint arXiv:2403.18476*, 2024. 2

[39] Yedong Shen, Xinran Zhang, Yifan Duan, Shiqi Zhang, Heng Li, Yilong Wu, Jianmin Ji, and Yanyong Zhang. Oggaussian: Occupancy based street gaussians for autonomous driving. *arXiv preprint arXiv:2502.14235*, 2025. 2

[40] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. In *International Conference on Machine Learning (ICML)*, 2023. 1

[41] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2446–2454, 2020. 2, 4, 6, 7, 8, 3

[42] Shanlin Sun, Bingbing Zhuang, Ziyu Jiang, Buyu Liu, Xiaohui Xie, and Manmohan Chandraker. Lidarf: Delving into lidar for neural radiance field on street scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19563–19572, 2024. 2

[43] Adam Tonderski, Carl Lindström, Georg Hess, William Ljungbergh, Lennart Svensson, and Christoffer Petersson. Neurad: Neural rendering for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14895–14904, 2024. 1

[44] Haithem Turki, Jason Y Zhang, Francesco Ferroni, and Deva Ramanan. Suds: Scalable urban dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12375–12385, 2023. 2

[45] Guangcong Wang, Zhaoxi Chen, Chen Change Loy, and Ziwei Liu. Sparsenerf: Distilling depth ranking for few-shot novel view synthesis. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9065–9076, 2023. 2

[46] Qitai Wang, Lue Fan, Yuqi Wang, Yuntao Chen, and Zhaoxiang Zhang. Freevs: Generative view synthesis on free driving trajectory. *arXiv preprint arXiv:2410.18079*, 2024. 1, 2, 3, 6, 7, 4, 5

[47] Jay Zhangjie Wu, Yuxuan Zhang, Haithem Turki, Xuanchi Ren, Jun Gao, Mike Zheng Shou, Sanja Fidler, Zan Gojcic, and Huan Ling. Difix3d+: Improving 3d reconstructions with single-step diffusion models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 26024–26035, 2025. 1, 2, 3, 4, 5, 6, 7, 8

[48] Ziyang Xie, Junge Zhang, Wenye Li, Feihu Zhang, and Li Zhang. S-nerf: Neural radiance fields for street views. *arXiv preprint arXiv:2303.00749*, 2023. 2

[49] Yunzhi Yan, Haotong Lin, Chenxu Zhou, Weijie Wang, Haiyang Sun, Kun Zhan, Xianpeng Lang, Xiaowei Zhou, and Sida Peng. Street gaussians: Modeling dynamic urban scenes with gaussian splatting. In *European Conference on Computer Vision*, pages 156–173. Springer, 2024. 2, 3, 5, 1

[50] Yunzhi Yan, Zhen Xu, Haotong Lin, Haian Jin, Haoyu Guo, Yida Wang, Kun Zhan, Xianpeng Lang, Hujun Bao, Xiaowei Zhou, et al. Streetcrafter: Street view synthesis with controllable video diffusion models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 822–832, 2025. 3

[51] Jiawei Yang, Boris Ivanovic, Or Litany, Xinshuo Weng, Seung Wook Kim, Boyi Li, Tong Che, Danfei Xu, Sanja Fidler, Marco Pavone, et al. Emernerf: Emergent spatial-temporal scene decomposition via self-supervision. *arXiv preprint arXiv:2311.02077*, 2023. 2

[52] Ze Yang, Yun Chen, Jingkang Wang, Sivabalan Manivasagam, Wei-Chiu Ma, Anqi Joyce Yang, and Raquel Urtasun. Unisim: A neural closed-loop sensor simulator. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1389–1399, 2023. 1

[53] Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin. Deformable 3d gaussians for high-fidelity monocular dynamic scene reconstruction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 20331–20341, 2024. 2, 3

[54] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024. 2, 6

[55] Zesong Yang, Ru Zhang, Jiale Shi, Zixiang Ai, Boming Zhao, Hujun Bao, Luwei Yang, and Zhaopeng Cui. Gurecon: Learning detailed 3d geometric uncertainties for neural surface reconstruction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 9363–9372, 2025. 1

[56] Vickie Ye, Ruilong Li, Justin Kerr, Matias Turkulainen, Brent Yi, Zhuoyang Pan, Otto Seiskari, Jianbo Ye, Jeffrey Hu, Matthew Tancik, et al. gsplat: An open-source library for gaussian splatting. *Journal of Machine Learning Research*, 26(34):1–17, 2025. 7

[57] Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sattler, and Andreas Geiger. Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction. *Advances in neural information processing systems*, 35:25018–25032, 2022. 2

[58] Zehao Yu, Anpei Chen, Binbin Huang, Torsten Sattler, and Andreas Geiger. Mip-splatting: Alias-free 3d gaussian splatting. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19447–19456, 2024. 2

[59] Zhongrui Yu, Haoran Wang, Jinze Yang, Hanzhang Wang, Zeke Xie, Yunfeng Cai, Jiale Cao, Zhong Ji, and Mingming Sun. Sgd: Street view synthesis with gaussian splatting and diffusion prior. *arXiv preprint arXiv:2403.20079*, 2024. 1, 2

[60] Baowen Zhang, Chuan Fang, Rakesh Shrestha, Yixun Liang, Xiaoxiao Long, and Ping Tan. Rade-gs: Rasterizing depth in gaussian splatting. *arXiv preprint arXiv:2406.01467*, 2024. 2

[61] Jiawei Zhang, Jiahe Li, Xiaohan Yu, Lei Huang, Lin Gu, Jin Zheng, and Xiao Bai. Cor-gs: sparse-view 3d gaussian splatting via co-regularization. In *European Conference on Computer Vision*, pages 335–352. Springer, 2024. 1

[62] Cheng Zhao, Su Sun, Ruoyu Wang, Yuliang Guo, JunJun Wan, Zhou Huang, Xinyu Huang, Yingjie Victor Chen, and Liu Ren. Tclc-gs: Tightly coupled lidar-camera gaussian splatting for autonomous driving. *arXiv preprint arXiv:2404.02410*, 2024. 2

[63] Guosheng Zhao, Chaojun Ni, Xiaofeng Wang, Zheng Zhu, Xueyang Zhang, Yida Wang, Guan Huang, Xinze Chen, Boyuan Wang, Youyi Zhang, et al. Drivedreamer4d: World models are effective data machines for 4d driving scene representation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 12015–12026, 2025. 1, 2, 7

[64] Guosheng Zhao, Xiaofeng Wang, Chaojun Ni, Zheng Zhu, Wenkang Qin, Guan Huang, and Xingang Wang. Recondreamer++: Harmonizing generative and reconstructive models for driving scene representation. *arXiv preprint arXiv:2503.18438*, 2025. 1, 2, 3, 6, 7

[65] Hongyu Zhou, Longzhong Lin, Jiabao Wang, Yichong Lu, Dongfeng Bai, Bingbing Liu, Yue Wang, Andreas Geiger, and Yiyi Liao. Hugsim: A real-time, photo-realistic and closed-loop simulator for autonomous driving. *arXiv preprint arXiv:2412.01718*, 2024. 2

[66] Xiaoyu Zhou, Zhiwei Lin, Xiaojun Shan, Yongtao Wang, Deqing Sun, and Ming-Hsuan Yang. Drivinggaussian: Composite gaussian splatting for surrounding dynamic autonomous driving scenes. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 21634–21643, 2024. 2

# Supplementary Material

## A. Additional Implementation Details

### A.1. Driving Scene Representation

FaithFusion's scene representation is built upon the decoupled 3D Gaussian Splatting (3DGS) structure widely adopted for driving scenes [7, 49, 64]. This representation is explicitly decomposed into three core components: a *Static Background* ($\mathcal{G}^{\text{bg}}$), *Dynamic Rigid Objects* ($\bar{\mathcal{G}}_v^{\text{rigid}}$ in local canonical space), and a *Sky Model* ($C_{\text{sky}}$), to precisely model geometry and motion through distinct semantic representations.

**Static Background.** The static background component is represented by a set of static Gaussian primitives $\mathcal{G}^{\text{bg}}$. These Gaussians are defined and optimized directly in the global world coordinate system, and all their attributes remain invariant over time.

**Rigid Body Representation and Transformation.** Gaussians belonging to a rigid object $v$ are defined in its local canonical space. The set of Gaussians in this space, $\bar{\mathcal{G}}_v^{\text{rigid}}$, do not change their internal attributes (mean $\bar{\mu}$, rotation $\bar{q}$, scale $\bar{s}$, opacity $\bar{o}$, and SH coefficients $\bar{c}$) over time $t$. The object's motion is captured entirely by an external rigid transformation $\mathbf{T}_v(t) \in \text{SE}(3)$, which transforms the Gaussians into world space $\mathcal{G}_v^{\text{rigid}}(t)$:

$$\mathcal{G}_v^{\text{rigid}}(t) = \mathbf{T}_v(t) \otimes \bar{\mathcal{G}}_v^{\text{rigid}}. \tag{S1}$$

The transformation operator $\otimes$ specifically updates the mean position $\mu(t)$ and rotation $\mathbf{q}(t)$ of the Gaussians when moved to the world coordinate system, while other attributes (scale, opacity, and SH coefficients) remain unchanged. We decompose the rigid pose as $\mathbf{T}_v(t) = (\mathbf{R}_v(t), \mathbf{t}_v(t))$, and the world-space mean position $\mu(t)$ is obtained by:

$$\mu(t) = \mathbf{R}_v(t)\bar{\mu} + \mathbf{t}_v(t), \tag{S2}$$

the rotation $\mathbf{q}(t)$ is updated by composing the object's rotational component $\mathbf{R}_v(t)$ with the canonical rotation $\bar{q}$:

$$\mathbf{q}(t) = \text{Rot}(\mathbf{R}_v(t), \bar{q}), \tag{S3}$$

where $\text{Rot}(\cdot)$ denotes rotating the quaternion $\bar{q}$ by the rotation matrix $\mathbf{R}_v(t)$. In this manner, the motion of dynamic objects is accurately modeled, ensuring geometric consistency across the time sequence.

**Sky Model Compositing.** The sky model is treated as a separate optimizable environmental texture map $C_{\text{sky}}$ to fit large-scale appearance. The final pixel color $C$ is obtained by $\alpha$-blending the rendered Gaussian image $C_{\mathcal{G}}$ with the sky image $C_{\text{sky}}$, where $C_{\mathcal{G}}$ is rendered from all Gaussians ($\mathcal{G}^{\text{bg}}$ and $\{\mathcal{G}_v^{\text{rigid}}\}$):

$$C = C_{\mathcal{G}} + (1 - O_{\mathcal{G}})C_{\text{sky}}, \tag{S4}$$

where $O_{\mathcal{G}}$ is the rendered opacity mask accumulated from all Gaussian primitives. This strategy addresses the challenge of reconstructing unbounded distant scenes.

### A.2. Evaluation Metrics

We conduct comprehensive quantitative evaluations following the protocol established by DriveDreamer4D [63], utilizing metrics that jointly assess the 3DGS novel view synthesis capability and the resulting spatiotemporal consistency under complex conditional shifts (e.g., lane shifts).

**Spatiotemporal Coherence Metrics.** (↑) To rigorously evaluate the coherence of dynamic elements and static scene structure, we employ two core metrics. Both quantify accuracy by comparing features detected in the rendered image with ground truth features that are geometrically projected from the original 3D scene onto the new trajectory view.

- **Novel Trajectory Agent IoU (NTA-IoU):** Measures the spatiotemporal accuracy of foreground dynamic agents (vehicles). Its computation involves detecting 2D bounding boxes on the rendered frames and comparing them to the projected ground-truth 3D bounding boxes. High NTA-IoU ensures accurate agent placement and adherence to the underlying 3D structure, leading to precise corrections in under-constrained regions.

- **Novel Trajectory Lane IoU (NTL-IoU):** Measures the geometric fidelity and spatiotemporal coherence of background lane lines. By comparing lane lines detected in the synthesized image against projected ground truth (often derived from the HDMap), it specifically verifies the integrity of the environment's static geometry. This metric reflects minimal disturbance to the original scene structure and guarantees environmental consistency.

**Perceptual Quality Metric.** (↓) We use Fréchet Inception Distance (FID) [15] to evaluate the overall visual realism and distributional quality of the 3DGS rendered novel view frames. FID calculates the distance between two multivariate Gaussian distributions fitted to the deep feature representations (from an Inception network) of generated frames and real frames. This score reflects the distribution-level similarity in a high-level perceptual space. A lower FID score indicates superior visual quality and consistent behavior across diverse viewpoints.

## A.3. Expected Information Gain Derivations

We provide the detailed derivation for the Expected Information Gain (EIG) approximation utilized in the main paper (Equation 4). This derivation strictly follows the unified framework for Bayesian optimal experimental design and information-theoretic approximations [19, 25].

**Motivation.** The analytical computation of the EIG definition in Equation 4 is intractable, requiring evaluation of complex posterior parameter distributions and expectations over the observation space. To obtain a highly efficient and differentiable acquisition function for 3DGS, our goal is to derive a *computable upper bound* of the EIG. This is achieved by combining the *Laplace approximation* (to simplify the entropy terms, Prop. 3.2/3.5 in [25]) and the *log-determinant inequality* (to obtain the final trace form upper bound, Lemma 5.1 in [25]).

**EIG Definition.** The EIG quantifies the *predicted reduction in uncertainty of 3DGS model parameters* $\mathbf{\Omega}$ if a new observation $(Y_{NVS}^{\text{gt}})$ at the novel view $X_{NVS}$ is acquired. Following the mutual information definition of EIG in [25] (Section 5.1), this uncertainty reduction is formally the mutual information between $\mathbf{\Omega}$ and $Y_{NVS}^{\text{gt}}$ (conditioned on $X_{NVS}$):

$$\text{EIG} = I\left[\mathbf{\Omega}; Y_{NVS}^{\text{gt}} \mid X_{NVS}\right] = \mathbb{H}[\mathbf{\Omega}] - \mathbb{E}_{p(Y_{NVS}^{\text{gt}}|X_{NVS})}\left[\mathbb{H}\left[\mathbf{\Omega} \mid Y_{NVS}^{\text{gt}}, X_{NVS}\right]\right], \tag{S5}$$

where $\mathbb{H}[\mathbf{\Omega}]$ denotes the *prior entropy*, and $\mathbb{H}[\mathbf{\Omega} \mid Y_{NVS}^{\text{gt}}, X_{NVS}]$ denotes the *posterior entropy*.

**! Note on Observation Index $i$ and EIG Decomposition.** The index $i$ in the main paper's equations is used to denote different scopes of observation:

- In the optimization objective (Eq. 2), $i$ denotes an individual training view.
- In the EIG definition (Eq. 4), $i$ denotes an individual view $Y_i^{NVS}$ within the novel view sequence $Y_{NVS}$. This formula calculates the information gain from a single such view.

The full EIG for the entire novel view sequence is obtained via view additivity. The final computable bound (Eq. 5) is then derived by applying Fisher information additivity principle (Prop. 4.2 in [25]) across all pixels $j$ in every view $i$. Therefore, the summation index $i$ in the final trace form (Eq. 5) is implicitly a flattened sum over all pixels in novel view sequence.

**Laplace Proxy Justification.** In Eq. (S5), $p(Y_{NVS}^{\text{gt}} \mid X_{NVS})$ is the true predictive distribution of real observations. We use the deterministic 3DGS rendered result $Y_{NVS}$ as a computationally tractable proxy for $Y_{NVS}^{\text{gt}}$. This is justified by the Laplace approximation (Prop. 3.2 in [25]), which models 3DGS parameters $\mathbf{\Omega}$ as a Gaussian posterior around the MAP parameters $\omega^*$ ($\mathbf{\Omega} \sim \mathcal{N}(\omega^*, (H''[\omega^*])^{-1})$). Here, $H''[\omega^*]$ is the Hessian of the negative log-posterior, serving as the inverse prior covariance. Since $Y_{NVS}$ is a deterministic function of $\mathbf{\Omega}$, $Y_{NVS}$ (conditioned on $\omega^*$) approximates $p(Y_{NVS}^{\text{gt}} \mid X_{NVS})$ for this Gaussian parameter posterior.

**Laplace Approximation of Entropy.** To compute the entropy terms in Eq. (S5), we apply the Gaussian differential entropy formula: $\mathbb{H}[\mathcal{N}(\mu, \mathbf{\Sigma})] = \frac{1}{2}\log\det(2\pi e\mathbf{\Sigma})$, where the covariance $\mathbf{\Sigma}$ is the inverse of the *observed information matrix* (Hessian of the negative log-posterior, Prop. 3.2 in [25]). For EIG, we distinguish two key observed information matrices:

- **Prior observed information:** $H''[\omega^*]$ (Hessian of the negative log-posterior of $\mathbf{\Omega}$ evaluated at $\omega^*$, corresponding to $\mathbb{H}[\mathbf{\Omega}]$);
- **Posterior observed information:** $H''[\mathbf{\Omega} \mid Y_{NVS}] = H''[\omega^*] + H''[Y_{NVS} \mid \omega^*]$ (sum of prior information and novel-view information, via the information additivity principle in Prop. 4.2 of [25]), where $H''[Y_{NVS} \mid \omega^*]$ is the Hessian of the negative log-likelihood of $Y_{NVS}$ (conditioned on $\omega^*$).

Substituting these Gaussian entropy approximations into Eq. (S5) yields:

$$\text{EIG} \approx \frac{1}{2}\log\det\left(2\pi e(H''[\omega^*])^{-1}\right) - \mathbb{E}_{p(Y_{NVS}|X_{NVS}, \omega^*)}\left[\frac{1}{2}\log\det\left(2\pi e\left(H''[\omega^*] + H''[Y_{NVS} \mid \omega^*]\right)^{-1}\right)\right] \tag{S6}$$

$$= \frac{1}{2}\left[\log\det\left((H''[\omega^*])^{-1}\right) - \mathbb{E}_{p(Y_{NVS}|X_{NVS}, \omega^*)}\left[\log\det\left((H''[\omega^*] + H''[Y_{NVS} \mid \omega^*])^{-1}\right)\right]\right] \tag{S7}$$

$$= \frac{1}{2}\mathbb{E}_{p(Y_{NVS}|X_{NVS}, \omega^*)}\left[\log\det\left(H''[\omega^*] + H''[Y_{NVS} \mid \omega^*]\right) - \log\det\left(H''[\omega^*]\right)\right] \tag{S8}$$

$$= \frac{1}{2}\mathbb{E}_{p(Y_{NVS}|X_{NVS}, \omega^*)}\left[\log\det\left(\mathbf{I} + H''[Y_{NVS} \mid \omega^*](H''[\omega^*])^{-1}\right)\right]. \tag{S9}$$

**Trace Form Upper Bound and Pixel-Level Decomposition.** We apply the log determinant inequality ($\log\det(\mathbf{I} + \mathbf{A}) \leq \text{tr}(\mathbf{A})$) (Lemma 5.1 in [25]) to Eq. (S9). By the linearity of the trace operator, and substituting the expectation of the novel-view Hessian $\mathbb{E}_{p(Y_{NVS}|X_{NVS}, \omega^*)}[H''[Y_{NVS} \mid \omega^*]]$ with its Fisher information (Prop. 4.1 in [25]), we get:

$$\text{EIG} \leq \frac{1}{2}\mathbb{E}_{p(Y_{NVS}|X_{NVS},\omega^*)}\left[\text{tr}\left(H''[Y_{NVS}\mid\omega^*](H''[\omega^*])^{-1}\right)\right] \tag{S10}$$

$$= \frac{1}{2}\text{tr}\left(\mathbb{E}_{p(Y_{NVS}|X_{NVS},\omega^*)}\left[H''[Y_{NVS}\mid\omega^*]\right](H''[\omega^*])^{-1}\right) \tag{S11}$$

$$= \frac{1}{2}\text{tr}\left(H''[Y_{NVS}\mid X_{NVS},\omega^*](H''[\omega^*])^{-1}\right). \tag{S12}$$

Finally, leveraging the Fisher information additivity principle (Prop. 4.2 in [25]), the total Fisher information of $Y_{NVS}$ is the sum of pixel-level Fisher information $H''[Y_{i,NVS}\mid X_{i,NVS},\omega^*]$ (one per pixel $i$). Substituting these pixel-wise Fisher information into Eq. (S12) yields the final trace-form approximation (main paper Equation 5):

$$\text{EIG} \leq \frac{1}{2}\sum_i \text{tr}\left(H''[Y_{i,NVS}\mid X_{i,NVS},\omega^*](H''[\omega^*])^{-1}\right). \tag{S13}$$
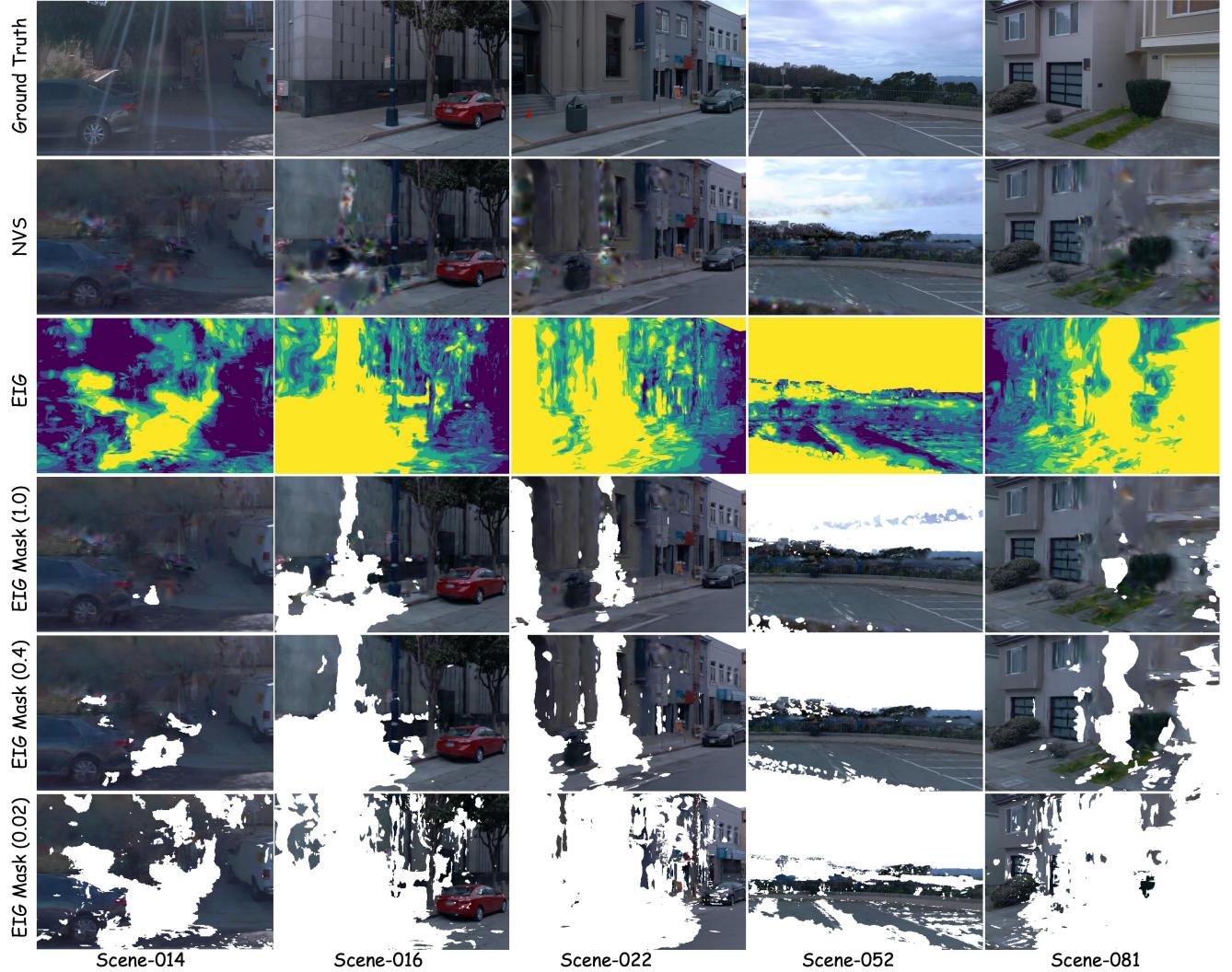


Figure S1. Visualization of EIG as a novel view synthesis quality proxy on representative Waymo [41] scenes. **Rows (Top to Bottom):** (1) Ground Truth, (2) Novel View Synthesis, (3) Pixel-wise EIG map (yellow = high EIG), (4-6) NVS masked by EIG thresholds ($\tau = 1.0, 0.4, 0.02$). White areas are excluded high-EIG pixels, confirming high EIG aligns with NVS artifacts. Sky regions are handled by a separate model and are excluded from this EIG analysis.

Figure S2. **Extended Qualitative Comparison on Waymo [41].** This figure provides additional novel view renderings for the same trajectory across representative methods, complementing the results shown in Fig. 5 of the main paper. Our method (last column) consistently maintains superior detail and fidelity across challenging regions, highlighted by the orange boxes, compared to methods [35, 46, 47].

## B. Additional Visualization Results

### B.1. EIG Correlation Validation and Evaluation Protocol

As quantified in Fig. 3 of the main paper, our cross-camera evaluation validates that pixel-level EIG is highly correlated with NVS quality. We detail the specific evaluation protocol here.

We first compute co-visible frame sequences between target cameras based on their Field of View (FoV) to ensure sufficient multi-view observation redundancy—a factor critical for optimizing 3D geometry and rendering fidelity. Considering that large low-frequency regions (e.g., solid colors, low-light scenes, or smooth surfaces) disproportionately inflate PSNR scores in standard NVS evaluations, which fails to reflect the model's ability to capture complex 3D structure, we filtered out frames predominantly containing such low-frequency information when assessing EIG-NVS correlation. This procedure ensures our validation efforts focus on EIG's efficacy in high-frequency detail and intricate geometry, effectively eliminating the inherent PSNR bias. Following this rigorous filtering process, we identified $4,245$ evaluation pairs for correlation validation.

For qualitative validation, Fig. S1 visualizes the EIG map and subsequent masking results on representative Waymo [41] scenes. The figure confirms the correlation: high EIG consistently aligns with NVS artifacts, and progressively masking these high-EIG pixels yields substantially improved perceptual clarity. This direct visual evidence not only validates the intuitive link between EIG and synthesis quality but also underscores EIG's unique advantage as a reliable, pixel-level proxy: it requires no manual annotation of artifacts, operates in a fully unsupervised manner, and provides fine-grained spatial guidance for targeted synthesis refinement.

### B.2. More Qualitative Results

To complement the qualitative analysis in the main paper (Figs. 5 and 6), extended visualization results are provided in Fig. S2 and Fig. S3, where our key conclusions regarding scene synthesis quality and consistency are further validated. All high-resolution visuals and frame-by-frame trajectory comparisons (covering original and novel trajectories with 3 meters/6 meters lane shifts) are available in the accompanying *qualitative_supplement* folder.

4

Figure S3. **Detailed Ablation Study of EIG-Guided Components.** This figure provides an extended analysis by quantitatively measuring the incremental performance of integrating EIG-guided components into the OmniRe baseline (Sec. 4.3 in the main paper for the overview). The comprehensive results further confirm the significant role of EIG guidance in coherently integrating diffusion edits and distilling them back into the 3DGS structure, thus mitigating over-restoration and geometric drift.

Fig. S2 extends the main paper comparisons, showing *FaithFusion* significantly outperforms baselines ( [35, 46, 47]) in preserving fine details (e.g., lane markings, building facades) and 3D coherence under large viewpoint shifts. Notably, the black regions in the upper part of the FreeVS [46] results were manually padded to align with the visualization scale of other methods, as its original output crops the sky region. Our method avoids spurious artifacts (ground bending, semantic mismatches) even at 6 meters lane offsets, aligning with our conclusion that EIG-guided control enables precise "generate-preserve" decisions.

Fig. S3 details our ablation results, validating EIG's role as a unified policy that harmonizes diffusion and 3DGS. Consistent with our core insight of replacing heuristics with information-theoretic guidance, EIG suppresses over-restoration in high-confidence regions (preserving 3DGS fidelity) and refines under-constrained areas (enhancing quality). This mechanism resolves the reconstruction-generation trade-off, successfully delivering the three key goals: consistency, quality, and faithfulness.