

# ThermoSplat: Cross-Modal 3D Gaussian Splatting with Feature Modulation and Geometry Decoupling

Zhaoqi Su, Shihai Chen, Xinyan Lin, Liqin Huang, *Member, IEEE*, Zhipeng Su, and Xiaoqiang Lu, *Senior Member, IEEE*

**Abstract**—Multi-modal scene reconstruction integrating RGB and thermal infrared data is essential for robust environmental perception across diverse lighting and weather conditions. However, extending 3D Gaussian Splatting (3DGS) to multi-spectral scenarios remains challenging. Current approaches often struggle to fully leverage the complementary information of multi-modal data, typically relying on mechanisms that either tend to neglect cross-modal correlations or leverage shared representations that fail to adaptively handle the complex structural correlations and physical discrepancies between spectrums. To address these limitations, we propose ThermoSplat, a novel framework that enables deep spectral-aware reconstruction through active feature modulation and adaptive geometry decoupling. First, we introduce a Cross-Modal FiLM Modulation mechanism that dynamically conditions shared latent features on thermal structural priors, effectively guiding visible texture synthesis with reliable cross-modal geometric cues. Second, to accommodate modality-specific geometric inconsistencies, we propose a Modality-Adaptive Geometric Decoupling scheme that learns independent opacity offsets and executes an independent rasterization pass for the thermal branch. Additionally, a hybrid rendering pipeline is employed to integrate explicit Spherical Harmonics with implicit neural decoding, ensuring both semantic consistency and high-frequency detail preservation. Extensive experiments on the RGBT-Scenes dataset demonstrate that ThermoSplat achieves state-of-the-art rendering quality across both visible and thermal spectrums.

**Index Terms**—3D Gaussian Splatting, RGBT scene reconstruction, multi-modal fusion, neural rendering, feature modulation.

## I. INTRODUCTION

3D scene reconstruction has been widely used in the field of autonomous systems, remote sensing, surveillance, etc. Traditional RGB-based reconstruction methods, while achieving high fidelity in most conditions, often suffer from performance degradation in challenging environments, e.g., low-light conditions, dense smoke, or darkness. To address these limitations, multi-modal reconstruction, especially those integrating RGB and thermal, has emerged as a critical research direction. Unlike RGB sensors, which depend on reflected light, thermal sensors capture long-wave infrared radiation emitted by objects, allowing for the acquisition of stable structural information and heat signatures that are inherently

invariant to illumination changes. This provides a reliable reference for scene geometry under extreme conditions.

The evolution of neural implicit representations has inspired research in multi-modal 3D reconstruction. Previous studies [1], [2] extended the Neural Radiance Fields (NeRF) [3] framework to the infrared spectrum, demonstrating the potential to synthesize thermal views from multi-view observations. However, NeRF-based methods often suffer from high computational cost and slow inference speeds, limiting their abilities in real-time applications. Recently, the emergence of 3D Gaussian Splatting (3DGS) [4] has enabled significantly faster training and rendering, with improved rendering quality. Building on this, several multi-spectral 3DGS frameworks have been proposed. Current state-of-the-art methods can be categorized into two paradigms: either explicitly decomposing 3DGS into modality-specific components to handle property disparities [5], or integrating multi-spectral information into a unified latent space for MLP-based decoding [6], [7]. However, achieving an optimal balance between shared geometry and modality-specific appearance remains non-trivial. The former paradigm often introduces increased model complexity and may face challenges in maintaining cross-modal spatial consistency, while the latter tends to have a limited capacity to precisely model the inherent physical discrepancies and structural variations present across different spectrums.

To bridge these gaps, we present ThermoSplat, a novel cross-modal 3DGS framework that enables deep spectral-aware reconstruction through active feature modulation and adaptive geometry decoupling. Unlike existing methods that either rely on explicit decomposition or unified latent representation, ThermoSplat introduces a FiLM-based [8] feature modulation mechanism. This module dynamically conditions the shared latent representation on thermal structural priors, enabling the model to actively leverage structural infrared features to guide visible texture synthesis. Furthermore, to accommodate the inherent physical discrepancies across different spectral bands, we propose a modality-adaptive geometric decoupling scheme, which allows for independent geometric adjustments in the infrared spectrum, effectively resolving artifacts in regions where transparency or reflectivity varies. Finally, to overcome the detail-loss inherent in pure feature-based decoding, we employ a hybrid rendering pipeline, which integrates explicit Spherical Harmonics (SH) with implicit decoding, achieving high-frequency RGB details while maintaining consistent semantic information across modalities.

Experimental results demonstrate that ThermoSplat achieves state-of-the-art rendering quality across both visible and thermal spectrums. The main contributions of this work are

Zhaoqi Su, Shihai Chen, Xinyan Lin, Liqin Huang, Zhipeng Su, and Xiaoqiang Lu are with the College of Physics and Information Engineering, Fuzhou University, Fuzhou 350108, China.

E-mail: {suzhaoqi, 2501120132, 2501127081, hlq, szp01, luxiaoqiang}@fzu.edu.cn

Corresponding author: Zhipeng Su, Xiaoqiang Lu

This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

summarized as follows:

- **Cross-Modal FiLM Modulation:** We design a Feature-wise Linear Modulation (FiLM) framework to establish deep feature dependencies. By utilizing structural priors to modulate shared latent features, our method enhances texture recovery and cross-modal alignment.
- **Modality-Adaptive Geometric Decoupling:** We introduce a learnable thermal opacity offset and execute an independent rasterization pass that decouples geometric representations between visible and infrared spectrums. This mechanism effectively resolves depth and occlusion misalignments caused by modality-inconsistent physical properties.
- **Hybrid Explicit-Implicit Rendering:** We propose a hybrid rendering pipeline that integrates explicit Spherical Harmonics (SH) with feature-modulated neural decoding. This architecture preserves high-frequency RGB details while maintaining consistent low-frequency semantic information across different modalities.

## II. RELATED WORK

### A. 3D Neural Scene Representation

Recent methods in 3D neural scene representation have shifted from implicit to explicit methods. The implicit NeRF-based methods [3], [9]–[11] represent the scene as a continuous function in 3D space formulated by a shallow network like MLP, achieving view-dependent and photo-realistic scene rendering results. However, these methods suffer from time-consuming training and rendering, limiting their practical use in real-time applications. To address this, 3D Gaussian Splatting (3DGS)-based methods [4], [12]–[14] propose an explicit scene representation paradigm, which leverages 3D Gaussian primitives for explicitly representing the geometric and texture information of the scene, enabling high-fidelity and real-time rendering through a differentiable tile-based rasterization pipeline. Some studies [15]–[17] leverage the idea of both feature-based decoding in NeRF and 3DGS representations to augment the representation capability by distilling high-dimensional latent features into each Gaussian primitive. By integrating these latent features with lightweight decoders, these methods can bypass the limitations of traditional Spherical Harmonics (SH), enabling more complex attribute modeling and cross-modal information interaction.

### B. Neural Thermal and RGBT Scene Reconstruction

Compared to visible light, thermal infrared signals possess distinct physical properties, such as being insensitive to lighting conditions and capable of reflecting the heat distribution of objects. Early attempts mostly extend NeRF-based representations for representing different modalities in a compact manner [1], [2], [18]–[20]. However, due to the volume rendering process in NeRF [3] and its reliance on dense sampling, these implicit methods often face challenges in precisely modeling the high-frequency details. In recent years, the emergence of 3DGS has shifted the focus toward explicit Gaussian-based RGBT (RGB + Thermal) scene modeling. ThermalGaussian [5] pioneered the extension of 3DGS to the

RGBT scene, which optimizes the thermal Gaussian by fine-tuning the pretrained RGB Gaussians and incorporates thermal priors for better scene modeling. Also, it releases the RGBT-Scenes dataset to facilitate benchmarking for multi-modal reconstruction tasks. MS-Splatting [6] formulates the multi-spectral 3D scene using a unified latent space for decoding both RGB and other spectral channels, which is also applied to agricultural NDVI tasks. MS-Splattingv2 [21] uses the optimized joint strategy with RGB initialization to improve rendering quality. MMOne [7] introduces a unified framework that represents multiple modalities, such as RGB, thermal, and language, within a single scene, which designs a multimodal decomposition mechanism for better learning properties of different modalities. Ma et al. [22] decomposes appearance into reflectance and thermal radiance, leveraging the thermal modality as a stable geometric prior to rectify distorted surfaces in low-light RGB inputs. Beyond general multimodal representation, several studies focus on reconstructing thermal infrared signals to tackle ill-posed problems or extreme environmental constraints. Some studies [23], [24] inject physics-based temperature or thermodynamics constraints into thermal 3DGS modeling. Veta-GS [25] introduces a view-dependent deformation field to capture the subtle thermal variations caused by emissivity and transmission effects, effectively reducing artifacts in infrared novel-view synthesis. Others extend the RGBT modeling into more-spectral or hyperspectral scenarios [26], [27]. Despite these advances, existing RGBT frameworks either treat different modalities as independent signals with limited feature interaction, or rely on a shared representation that tends to overlook modality-specific physical discrepancies. These limitations motivate us to explore a more flexible modulation and decoupled modeling for RGBT scene reconstruction.

## III. METHOD

### A. Overview

The overall architecture of ThermoSplat is designed to achieve high-fidelity multi-modal scene reconstruction by addressing spectral-varying properties. As illustrated in Fig. 1, we represent the scene using multi-modal feature-enhanced 3DGS [4]. The pipeline first performs active feature interaction via a **Cross-Modal FiLM Modulation** on the rasterized latent representations, which utilizes thermal structural priors to guide visible texture synthesis. To account for geometric inconsistencies across spectrums, we introduce a **modality-adaptive geometric decoupling** scheme, which uses the learnable offset  $\Delta_t \alpha$  and executes an independent rasterization pass to accommodate modality-specific geometries. Finally, a **hybrid rendering** strategy is employed to combine explicit Spherical Harmonics (SH) with implicit feature-decoded outputs for preserving high-frequency details and view-dependent effects in the visible spectrum.

The remainder of this section provides a detailed formalization of our framework. We first briefly introduce 3D Gaussian Splatting and feature-based rasterization in Section III-B. In Section III-C, we describe the cross-modal feature modulation mechanism, which enables active spectral interaction.

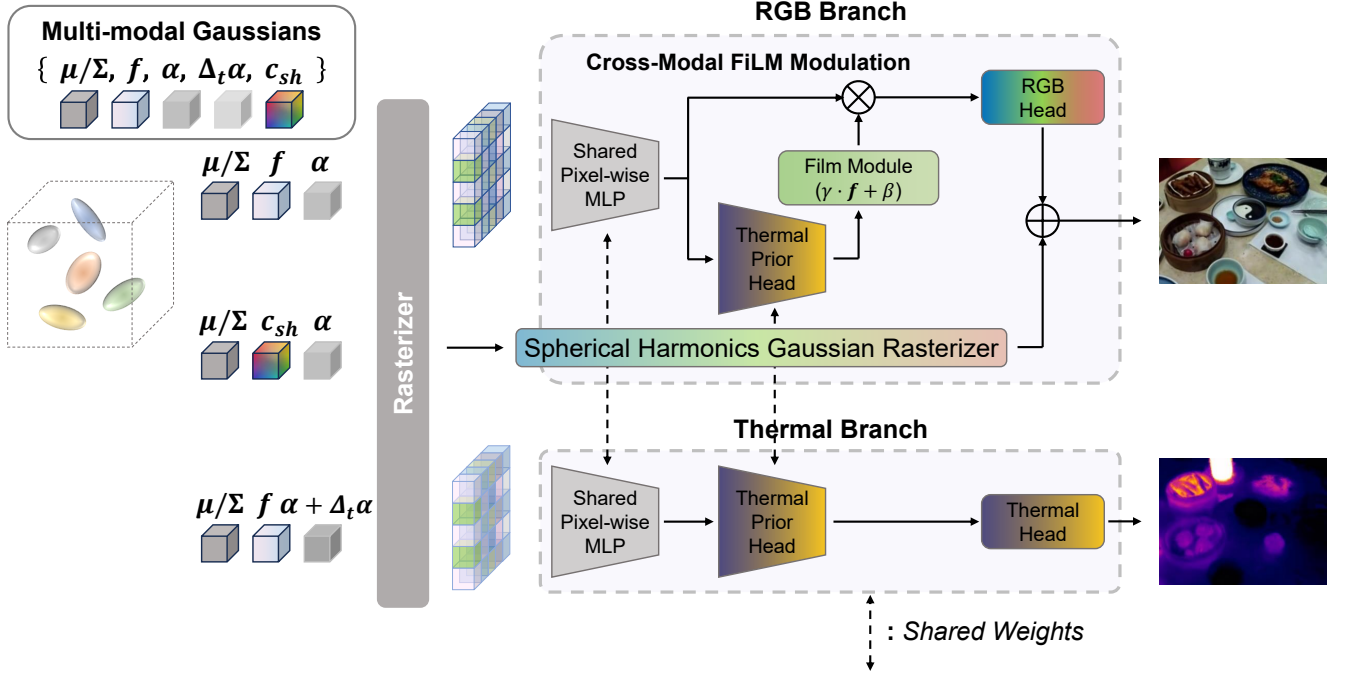


Fig. 1. Overview of the proposed ThermoSplat framework. Given multi-spectral inputs, our method optimizes 3D Gaussian primitives with decoupled properties. (a) Cross-Modal FiLM Modulation dynamically conditions shared latent features on thermal structural priors to guide visible texture synthesis. (b) Modality-Adaptive Geometric Decoupling resolves geometric inconsistencies between visible and infrared spectrums. (c) The Hybrid Rendering pipeline integrates explicit Spherical Harmonics (SH) with implicit neural decoding, ensuring high-frequency detail preservation and cross-modal semantic consistency.

Section III-D presents the Multi-spectral Hybrid Rendering pipeline, where we first detail the modality-adaptive geometric decoupling for modal-specific geometries, followed by the hybrid rendering strategy for RGB synthesis to preserve high-frequency details.

### B. Preliminaries

3D Gaussian Splatting (3DGS) [4] represents a 3D scene as a collection of  $N$  Gaussian primitives. Each Gaussian is characterized by its center position  $\mu \in \mathbb{R}^3$ , an anisotropic covariance  $\Sigma = RSS^T R^T$ , and an opacity value  $\alpha$ . The influence of a Gaussian at a 3D point  $\mathbf{x}'$  is defined as:

$$G(\mathbf{x}'; \mu_i, \Sigma_i) = e^{-\frac{1}{2}(\mathbf{x}' - \mu_i)^T \Sigma_i^{-1}(\mathbf{x}' - \mu_i)}, \quad (1)$$

where  $\mathbf{x}'$  denotes the 3D point in the camera coordinate system. Unlike traditional 3DGS that directly optimizes Spherical Harmonic (SH) coefficients for color, we follow a feature-based splatting paradigm [6] where each Gaussian carries a multi-dimensional latent feature  $f \in \mathbb{R}^d$ . This feature serves as a unified latent representation that can be subsequently decoded into modality-specific signals (e.g., RGB or thermal) via neural networks.

The rendering process follows the point-based  $\alpha$ -blending model. For a specific pixel, the attributes of projected 2D Gaussians are sorted by depth and blended to compute the aggregated pixel value:

$$\mathcal{A} = \sum_{i \in \mathcal{N}} \mathbf{a}_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad (2)$$

where  $\mathbf{a}_i$  denotes the generic attribute of the  $i$ -th Gaussian (such as latent feature  $f_i$  or explicit color attributes) and  $\alpha_i$  is the opacity of the Gaussian at that pixel.

In this work, we leverage this differentiable rasterization to bridge different modalities. By decoupling the feature decoding from the geometric projection, we can perform complex cross-modal modulations in the latent space before generating the final visible and thermal images.

### C. Cross-Modal Feature Modulation

To address the spectral gap between visible and infrared modalities, we propose a cross-modal modulation mechanism. Instead of treating visible and thermal signals as independent entities, our framework leverages structural priors inherent in the thermal spectrum to guide the synthesis of visible textures. As shown in Fig. 1, the proposed Cross-Modal FiLM [8] Modulation integrates feature extraction and conditioning in a unified neural architecture.

**Shared Latent Encoding.** The process begins with the rendered feature map  $\mathcal{A}_f \in \mathbb{R}^{H \times W \times d}$  by rasterizing the per-Gaussian feature  $f_i$  through the 3DGS rendering pipeline [4]. To extract high-level semantic information, we first pass  $\mathcal{A}_f$  through a shared encoder  $\Phi_{shared}$  consisting of multiple pixel-wise linear layers with SiLU activations [28]:

$$h = \Phi_{shared}(\mathcal{A}_f), \quad (3)$$

where  $h$  represents the intermediate feature representation that serves as the common basis for both modalities.

**Thermal-Guided FiLM Modulation.** Distinct from directly applying MLP layers for different spectrums, we intro-

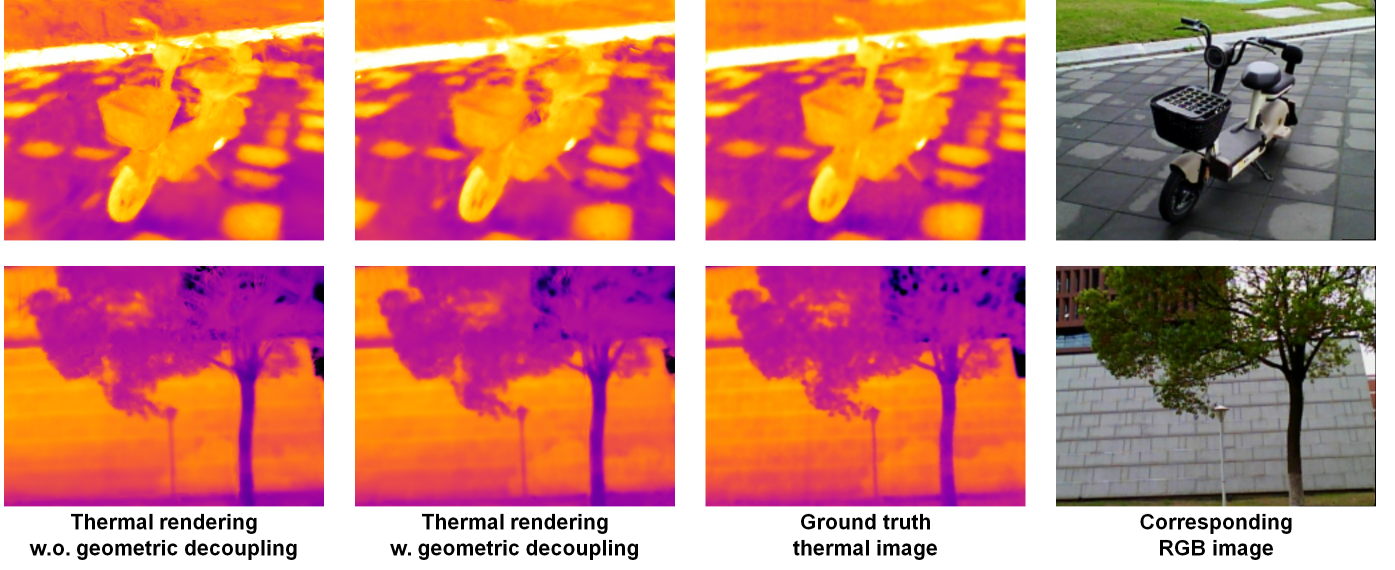


Fig. 2. Thermal rendering results with and without geometric decoupling. The thermal rendering results without geometric decoupling may inherit sharp textures and high-frequency noise from the visible spectrum.

duce a *Thermal Prior Head*  $\Phi_{th}$  and a *FiLM-based modulation* [8] layer. Crucially, as  $h_{th} = \Phi_{th}(h)$  is directly supervised by the thermal rendering task through the subsequent decoding stage (Eq. 6), it is naturally driven to distill structural priors that are most representative of the infrared domain. The feature  $h_{th}$  is subsequently mapped to a set of modulation parameters  $(\gamma, \beta)$  through a linear transformation:

$$[\gamma, \beta] = \text{Linear}_{film}(h_{th}). \quad (4)$$

By treating the infrared information as a conditioning signal, we apply Feature-wise Linear Modulation (FiLM) to the shared representation  $h$ :

$$h_{mod} = \gamma \odot h + \beta, \quad (5)$$

where  $\odot$  denotes element-wise multiplication. This operation dynamically scales and shifts the latent features based on the thermal structural prior, effectively “masking” or “enhancing” regions where visible textures are likely to align with thermal boundaries.

**Modality-Specific Decoding.** Finally, the modulated feature  $h_{mod}$  and the thermal feature  $h_{th}$  are decoded into their respective spectral domains:

$$\begin{aligned} \mathcal{C}_{implicit}^{rgb} &= \text{Sigmoid}(\Phi_{rgb}(h_{mod})), \\ \mathcal{C}^{thermal} &= \text{Sigmoid}(\Phi_{th\_out}(h_{th})), \end{aligned} \quad (6)$$

where  $\mathcal{C}_{implicit}^{rgb}$  provides the base color component for the subsequent hybrid rendering stage. Notably, the thermal-specific feature  $h_{th}$  serves a dual purpose: it acts as the source for FiLM parameter generation while simultaneously being decoded into the infrared signal  $\mathcal{C}^{thermal}$ . This hierarchical modulation ensures that the synthesis of visible images is physically constrained by the cross-modal structural consistency.

#### D. Multi-spectral Hybrid Rendering

Based on the modulated cross-modal features, we develop a dual-branch rendering pipeline to synthesize images in both visible and thermal spectrums. This pipeline addresses the geometric inconsistencies and texture fidelity requirements unique to each modality.

**Modality-Adaptive Geometric Decoupling** Typical multi-modal Gaussian representations assume a shared geometry across all spectrums. However, physical properties such as transparency and reflectivity vary significantly between visible and infrared bands. To accommodate these discrepancies, we introduce a modality-adaptive geometric decoupling scheme.

For the thermal rendering branch, we define a modality-specific opacity  $\alpha_{t,i}$  for each Gaussian  $i$  by adding a learnable offset  $\Delta_t \alpha_i$  to the base opacity  $\alpha_i$ :

$$\alpha_{t,i} = \text{Sigmoid}(\text{Logit}(\alpha_i) + \Delta_t \alpha_i), \quad (7)$$

where  $\Delta_t \alpha_i$  captures the fine-grained geometric deviations. Consequently, the thermal representation is generated via an independent rasterization pass:

$$\mathcal{A}_{f(t)} = \text{Rasterize}(\mu, \Sigma, \alpha_t, f), \quad (8)$$

where  $\mathcal{A}_{f(t)} \in \mathbb{R}^{H \times W \times d}$  represents the thermal-specific feature map generated using the decoupled opacity  $\alpha_t$ . The final thermal image  $\mathcal{C}^{thermal}$  is then decoded from  $\mathcal{A}_{f(t)}$  via the thermal head discussed in Section III-C. This independent pass ensures that occlusions and structural boundaries in the thermal image remain physically consistent with infrared sensors.

As shown in Fig. 2, without the geometric decoupling mechanism, the thermal branch tends to inherit redundant high-frequency textures from the visible spectrum that do not exist in the infrared domain. Our proposed decoupling module effectively filters out these cross-modal artifacts, ensuring that



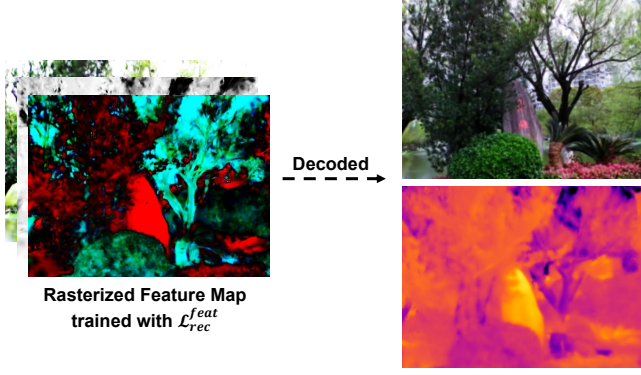


Fig. 3. Feature level reconstruction loss on the rasterized feature maps. Left: rendered feature map, right: reconstructed RGB-thermal scene. Note that  $\mathcal{A}_f$  and  $\mathcal{A}_{f(t)}$  are only different in the opacity used in rasterization.

the thermal rendering preserves its natural smoothness while accurately representing its own structural boundaries.

**Hybrid RGB Synthesis.** While the thermal branch focuses on geometric consistency, the RGB branch requires high-frequency view-dependent details. We propose a hybrid strategy that bridges explicit Gaussian Splatting with implicit neural decoding.

Specifically, the final RGB color  $C^{rgb}$  is formulated as the summation of two components:

$$C^{rgb} = \mathcal{R}_{sh}(\mu, \Sigma, \alpha, \mathbf{c}_{sh}) \oplus C_{implicit}^{rgb}, \quad (9)$$

where  $\mathcal{R}_{sh}$  denotes the explicit color rendered via standard Spherical Harmonic (SH) coefficients  $\mathbf{c}_{sh}$ , capturing view-dependent specular effects. The second term,  $C_{implicit}^{rgb}$ , is the implicit component decoded from the modulated latent features  $h_{mod}$ , providing multi-modal consistent textures. By combining these two components, our hybrid rendering scheme effectively preserves the high-frequency view-dependent properties of explicit rasterization, while simultaneously enriching the visible textures with the structural intelligence of neural-modulated latent features.

### E. Loss Functions

The training objective of ThermoSplat is to optimize the multi-modal Gaussian representation and the neural modulation networks through a composite loss function  $\mathcal{L}$ . This objective ensures that the synthesized visible and thermal images adhere to the ground truth in terms of both pixel intensity and structural topology.

**Spectral Reconstruction Loss** For both the visible and infrared modalities, we employ a combination of  $\ell_1$  loss and Structural Similarity (SSIM) to supervise the final rendered images against the corresponding ground-truth images  $\mathbf{I}_m$ :

$$\mathcal{L}_{rec}^m = (1 - \lambda_s) \|\mathbf{I}_m - \mathbf{C}^m\|_1 + \lambda_s (1 - \text{SSIM}(\mathbf{I}_m, \mathbf{C}^m)), \quad (10)$$

where  $m \in \{rgb, thermal\}$  denotes the spectral modality, and  $\mathbf{C}^m$  are the corresponding output images in our pipeline.

To provide structural guidance during the intermediate stages, we enforce consistency on the rasterized feature maps

$\mathcal{A}_f$  and  $\mathcal{A}_{f(t)}$  by slicing specific channels corresponding to physical properties. As shown in Fig. 3, for the visible branch, we constrain the first three channels of  $\mathcal{A}_f$  to match the RGB appearance. In parallel, for the thermal branch, we supervise the subsequent latent channel (index 3) of  $\mathcal{A}_{f(t)}$  using the transformed thermal map derived from the Ironbow colormap protocol. As this transformed map effectively serves as a proxy for physical temperature and thermal intensity, this constraint encourages the model to learn a compact and structural representation. By applying these latent constraints, we ensure the latent space captures the fundamental visual and thermal distribution before it is decoded. The feature-level reconstruction loss is thus formulated as:

$$\mathcal{L}_{rec}^{feat} = \mathcal{L}(\mathcal{A}_f[:3], \mathbf{I}_{rgb}) + \eta \cdot \mathcal{L}(\mathcal{A}_{f(t)}[3], \mathbf{I}_{th}^{trans}), \quad (11)$$

where  $\mathbf{I}_{th}^{trans}$  represents the temperature-correlated intensity map,  $\mathcal{L}$  denotes the composite  $\ell_1$  and SSIM loss function as defined in Eq. 10.

**Thermal Spatial Regularization** Due to the high-contrast and often sparse nature of infrared signals, we introduce a spatial smoothness constraint on the predicted thermal image:

$$\mathcal{L}_{smooth} = \sum_{p \in \Omega} |\nabla C^{thermal}(p)|, \quad (12)$$

where  $\nabla$  denotes the spatial gradient operator at pixel  $p$ . This term enforces the smooth structural characteristics of the thermal output.

**Total Objective** The final training objective is a weighted summation of the aforementioned reconstruction and regularization terms:

$$\mathcal{L} = \mathcal{L}_{rec} + \lambda_{rf} \mathcal{L}_{rec}^{feat} + \lambda_{sm} \mathcal{L}_{smooth}, \quad (13)$$

where the image-level reconstruction loss is defined as  $\mathcal{L}_{rec} = \mathcal{L}_{rec}^{rgb} + \mathcal{L}_{rec}^{thermal}$ ,  $\lambda_{rf}$  and  $\lambda_{sm}$  are hyper-parameters balancing feature terms and smooth terms. By optimizing this joint objective, our framework ensures that the synthesized modalities satisfy both pixel-level accuracy and the inherent structural characteristics of thermal radiation.

## IV. EXPERIMENTS

### A. Implementation details.

We evaluate our model on the RGBT-Scenes dataset, which comprises over 1,000 calibrated RGB-thermal pairs across ten indoor and outdoor scenes under diverse environmental and lighting conditions. To demonstrate the effectiveness of our approach, we compare our model with state-of-the-art methods, including MMOne [7], MS-Splattingv2 [21] and ThermalGaussian [5]. We also compare our method with the 3DGS [4] baseline trained on both modalities separately. We use Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index (SSIM), and Learned Perceptual Image Patch Similarity (LPIPS) [29] to evaluate the rendering quality of both visible and thermal modalities.

Our framework is implemented on PyTorch and trained with an NVIDIA 3090 GPU. We train our pipeline for 30K iterations, which is the same as the setting used in

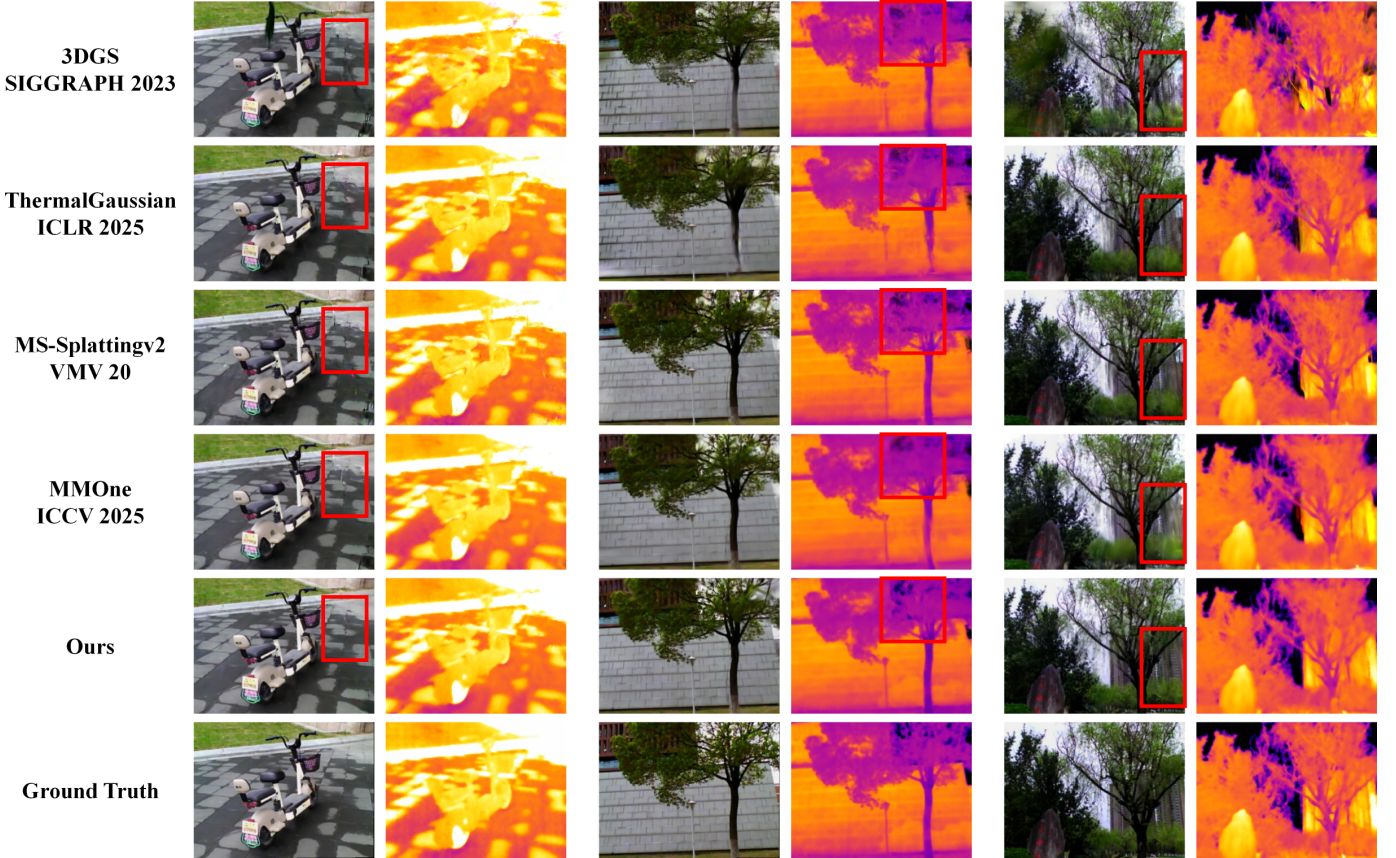


Fig. 4. Qualitative comparison of novel view synthesis results on the RGBT-Scenes dataset. We compare ThermoSplat against state-of-the-art multi-spectral reconstruction methods ThermalGaussian [5], MS-Splattingv2 [21], MMOne [7], and the 3DGS baseline [4]. Our method generates more accurate rendering results and structural details.

3DGS [4], MMOne [7] and ThermalGaussian [5]. For MS-Splattingv2 [21], we follow the training strategy proposed in the paper and train them with 120K iterations. In our experiments, we set per-Gaussian feature dimension  $d = 8$ , and  $\lambda_s = 0.2$ ,  $\eta = 0.5$ ,  $\lambda_{rf} = 1$ ,  $\lambda_{sm} = 0.3$  for loss weights.

### B. Results and Comparisons

We evaluate the novel view synthesis performance on the test set to validate the rendering quality of our method against other state-of-the-art methods. As illustrated in Fig. 4, our method produces results with finer texture details and fewer visual artifacts compared to existing approaches and the 3DGS baseline. Specifically, our model performs better in recovering complex structures that are often blurred or misaligned in the baseline reconstructions, especially on the RGB branch, which is attributed to the proposed cross-modality modulation that effectively leverages structural priors of the scene.

As shown in Tab. I, our method achieves superior performance compared to state-of-the-art baselines across most scenes. Specifically, our model attains the highest average scores in all three metrics for both RGB and thermal modalities. Notably, on the average PSNR, our method outperforms the second-best competitor (MMOne [7]) by 0.34 dB in the RGB spectrum and 0.19 dB in the thermal spectrum. The consistent improvement in SSIM and LPIPS further demonstrates

our model’s capability to reconstruct fine-grained structural details and maintain perceptual fidelity.

While MMOne [7] and MS-Splattingv2 [21] show competitive results in specific scenes (e.g., Dim and Trk), our method demonstrates more robust generalization across diverse environments. These results validate that our modality modulation and geometric decoupling strategy successfully resolves the discrepancies between modalities without compromising the reconstruction quality of the individual branches.

### C. Ablation

To verify the contribution of each design element, we conduct ablation experiments as summarized in Tab. II. First, the comparison between our full model and the “MLP-based” variant demonstrates the advantage of our feature-guided modulation over a standard decoding structure, as the former better leverages spatial-aware features for high-quality appearance synthesis. Second, removing the geometric decoupling mechanism (“w.o. geo. decoup.”) leads to a consistent performance decline in both modalities, confirming that isolating physical geometry from modality-specific radiation is essential for robust RGBT scene modeling. Finally, the exclusion of latent constraints (“w.o. fea(rgb/th)”) results in degradation of perceptual details, which indicates that the feature-level reconstruction loss  $\mathcal{L}_{rec}^{feat}$  is crucial for encouraging the model to capture fundamental structural and intensity distributions.

TABLE I  
QUANTITATIVE EVALUATION OF RGB AND THERMAL (T) RENDERING RESULTS.

M	Metric	Method	Dim	DS	Ebk	RB	Trk	RK	Bldg	II	Pt	LS	Avg.
RGB	PSNR $\uparrow$	3DGS	23.27	21.18	26.17	28.23	22.45	20.74	21.80	24.40	25.65	20.18	23.41
		ThermalGaussian	24.38	21.76	26.85	28.12	24.17	23.14	<b>24.19</b>	24.55	25.48	21.71	24.44
		MS-Splattingv2	24.06	21.18	26.87	28.12	<b>24.54</b>	23.42	23.90	23.77	<u>26.20</u>	<u>22.05</u>	24.41
		MMOne	<b>24.65</b>	<u>22.05</u>	<b>27.43</b>	<b>29.03</b>	23.96	<u>24.12</u>	<u>24.16</u>	<u>25.65</u>	26.01	21.81	<u>24.89</u>
		Ours	<u>24.59</u>	<b>22.12</b>	<u>27.21</u>	<u>28.96</u>	<u>24.31</u>	<b>24.20</b>	24.14	<b>25.98</b>	<b>26.48</b>	<b>24.31</b>	<b>25.23</b>
	SSIM $\uparrow$	3DGS	0.842	0.771	0.902	0.917	0.810	0.765	0.827	0.875	0.867	0.688	0.826
		ThermalGaussian	0.858	0.797	0.905	0.920	0.840	0.822	0.849	0.884	0.855	0.739	0.847
		MS-Splattingv2	0.859	0.788	0.914	<u>0.922</u>	<b>0.859</b>	0.827	<u>0.855</u>	0.877	<u>0.878</u>	<u>0.739</u>	0.852
		MMOne	<u>0.862</u>	<u>0.810</u>	<u>0.918</u>	0.916	0.845	<b>0.842</b>	0.847	<u>0.897</u>	0.876	0.727	<u>0.854</u>
		Ours	<b>0.872</b>	<b>0.818</b>	<b>0.934</b>	<b>0.941</b>	<u>0.859</u>	<u>0.841</u>	<b>0.858</b>	<b>0.911</b>	<b>0.886</b>	<b>0.788</b>	<b>0.871</b>
	LPIPS $\downarrow$	3DGS	0.199	0.271	0.169	0.197	0.244	0.220	0.183	0.193	0.177	0.289	0.214
		ThermalGaussian	0.194	0.253	0.169	0.199	0.211	0.184	0.170	0.186	0.195	0.268	0.203
		MS-Splattingv2	<b>0.150</b>	<u>0.224</u>	<u>0.145</u>	0.197	<u>0.170</u>	<u>0.141</u>	<u>0.145</u>	<u>0.161</u>	<b>0.132</b>	<u>0.211</u>	<u>0.168</u>
		MMOne	0.203	0.254	0.160	0.235	0.226	0.178	0.184	0.183	0.178	0.291	0.209
		Ours	<u>0.155</u>	<b>0.204</b>	<b>0.121</b>	<b>0.164</b>	<b>0.166</b>	<b>0.130</b>	<b>0.131</b>	<b>0.138</b>	<u>0.136</u>	<b>0.180</b>	<b>0.153</b>
T	PSNR $\uparrow$	3DGS	25.99	18.71	20.61	26.55	25.30	26.45	26.83	29.69	24.09	18.48	24.27
		ThermalGaussian	<u>26.46</u>	<b>22.28</b>	23.31	<u>27.17</u>	25.88	26.33	26.72	29.86	26.16	22.27	25.64
		MS-Splattingv2	26.06	21.43	<u>23.32</u>	25.44	<u>26.08</u>	27.24	26.89	<u>29.98</u>	<b>27.01</b>	<u>22.64</u>	25.61
		MMOne	<b>26.90</b>	<u>21.81</u>	<b>23.79</b>	<b>27.39</b>	25.44	<u>27.65</u>	<u>27.06</u>	<b>30.27</b>	26.05	22.52	<u>25.89</u>
		Ours	25.99	21.54	22.95	26.83	<b>26.25</b>	<b>28.48</b>	<b>27.45</b>	29.78	<u>27.00</u>	<b>24.50</b>	<b>26.08</b>
	SSIM $\uparrow$	3DGS	0.889	0.787	0.812	0.914	0.863	0.922	0.896	0.892	0.867	0.768	0.861
		ThermalGaussian	0.886	0.835	0.862	0.919	<u>0.874</u>	0.922	0.888	0.896	0.883	0.850	0.882
		MS-Splattingv2	0.876	0.803	0.855	0.900	0.871	0.927	0.888	0.890	<u>0.903</u>	0.853	0.877
		MMOne	<b>0.894</b>	<b>0.840</b>	<b>0.874</b>	<u>0.926</u>	0.870	<u>0.933</u>	<u>0.902</u>	<u>0.906</u>	0.895	<u>0.861</u>	<u>0.890</u>
		Ours	<u>0.890</u>	<u>0.839</u>	<u>0.865</u>	<b>0.928</b>	<b>0.889</b>	<b>0.941</b>	<b>0.909</b>	<b>0.910</b>	<b>0.912</b>	<b>0.889</b>	<b>0.897</b>
	LPIPS $\downarrow$	3DGS	0.127	0.259	0.307	0.209	0.142	0.126	0.185	0.091	0.227	0.378	0.205
		ThermalGaussian	0.129	0.210	0.203	0.198	0.136	0.124	0.177	0.091	0.181	0.248	0.170
		MS-Splattingv2	<b>0.092</b>	<u>0.164</u>	<b>0.148</b>	<u>0.133</u>	<u>0.107</u>	<u>0.073</u>	<u>0.103</u>	<u>0.064</u>	<b>0.075</b>	<u>0.177</u>	<u>0.114</u>
		MMOne	0.125	0.194	0.201	0.213	0.142	0.127	0.198	0.083	0.205	0.272	0.176
		Ours	<u>0.100</u>	<b>0.149</b>	<u>0.149</u>	<b>0.096</b>	<b>0.094</b>	<b>0.059</b>	<b>0.085</b>	<b>0.057</b>	<u>0.075</u>	<b>0.139</b>	<b>0.101</b>

TABLE II  
ABLATION STUDY. WE CONDUCTED ABLATION EXPERIMENTS ON DIFFERENT MODULES OF OUR PIPELINE.

Method Variant	RGB Modality			Thermal Modality		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
MLP-based	25.14	0.869	0.154	25.88	0.895	0.104
w.o. geo. decoup.	25.07	0.868	0.159	25.82	0.892	0.106
w.o. hybrid rgb	25.07	0.867	0.157	25.86	0.893	0.105
w.o. fea(th)	24.98	0.868	0.155	25.77	0.894	0.107
w.o. fea(rgb)	24.88	0.858	0.190	<u>25.93</u>	<b>0.897</b>	0.104
<b>Ours</b>	<b>25.23</b>	<b>0.871</b>	<b>0.153</b>	<b>26.08</b>	<u>0.897</u>	<b>0.101</b>

Collectively, these results validate that the synergy of the proposed modulation mechanism, geometric decoupling, and latent supervision ensures optimal RGBT reconstruction.

## V. CONCLUSION

In this paper, we present ThermoSplat, a novel cross-modal 3D Gaussian Splatting framework designed for high-fidelity RGBT scene reconstruction. To effectively bridge the gap between visible and thermal modalities, we introduce a Cross-Modal FiLM Modulation mechanism that leverages thermal structural priors to guide visible texture synthesis. Furthermore, to address the inherent geometric inconsistencies caused by disparate physical sensing properties, we propose

a Modality-Adaptive Geometric Decoupling scheme, which enables the model to accurately represent independent spectral characteristics without compromising spatial alignment. Extensive experiments on the RGBT-Scenes dataset demonstrate that our approach achieves state-of-the-art performance in both rendering quality and structural accuracy. By integrating explicit geometric representations with implicit neural feature modulation, ThermoSplat provides a robust and efficient solution for multi-spectral scene understanding in visually degraded environments.

**limitations.** Despite the promising results, ThermoSplat has certain limitations that offer directions for future research. First, the current geometric decoupling scheme primarily focuses on the thermal branch; however, in scenarios with extreme glass reflections or high-transparency surfaces, more complex multi-modal interactions might be required to fully resolve depth ambiguities. Second, the use of latent feature modulation introduces additional memory overhead during the neural decoding phase compared to vanilla 3DGS. Future work will explore more lightweight modulation architectures and investigate the potential of extending this framework to other spectral domains, such as near-infrared or hyperspectral data, to further enhance its versatility and robustness in all-weather environmental perception.

## ACKNOWLEDGMENTS

This paper is supported in part by the National Natural Science Foundation of China (Grant No. 62402274) and the Start-up Funding of Fuzhou University (Grant No. XRC-25164) to Zhaoqi Su; in part by the Education and Scientific Research Project for Middle-aged and Young Teachers of Fujian Province, China (Grant No. JZ250004) to Zhipeng Su; in part by the Special Fund for Promoting High-Quality Development of Marine and Fishery Industries in Fujian Province (Grant No. FJHYF-L-2025-07-005) to Xiaoqiang Lu. The authors would like to acknowledge the use of Gemini to improve the language and readability of the manuscript during the writing process.

## REFERENCES

- [1] Y. Y. Lin, X.-Y. Pan, S. Fridovich-Keil, and G. Wetzstein, “Thermalnerf: Thermal radiance fields,” in *2024 IEEE International Conference on Computational Photography (ICCP)*. IEEE, 2024, pp. 1–12.
- [2] M. Hassan, F. Forest, O. Fink, and M. Mielle, “Thermonerf: A multimodal neural radiance field for joint rgb-thermal novel view synthesis of building facades,” *Adv. Eng. Inform.*, vol. 65, no. PD, May 2025. [Online]. Available: <https://doi.org/10.1016/j.aei.2025.103345>
- [3] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, “Nerf: Representing scenes as neural radiance fields for view synthesis,” *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.
- [4] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, “3d gaussian splatting for real-time radiance field rendering,” *ACM Trans. Graph.*, vol. 42, no. 4, pp. 139–1, 2023.
- [5] R. Lu, H. Chen, Z. Zhu, Y. Qin, M. Lu, L. zhang, C. Yan, and a. xue, “Thermalgaussian: Thermal 3d gaussian splatting,” in *International Conference on Representation Learning*, Y. Yue, A. Garg, N. Peng, F. Sha, and R. Yu, Eds., vol. 2025, 2025, pp. 1105–1117. [Online]. Available: [https://proceedings.iclr.cc/paper\\_files/paper/2025/file/03bdba50e3741ac5e3eaa0e55423587e-Paper-Conference.pdf](https://proceedings.iclr.cc/paper_files/paper/2025/file/03bdba50e3741ac5e3eaa0e55423587e-Paper-Conference.pdf)
- [6] L. Meyer, J. Grün, M. Weiherer, B. Egger, M. Stamminger, and L. Franke, “Multi-spectral gaussian splatting with neural color representation,” *arXiv preprint arXiv:2506.03407*, 2025.
- [7] Z. Gu and B. Wang, “Mmone: Representing multiple modalities in one scene,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2025, pp. 1088–1098.
- [8] E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville, “Film: Visual reasoning with a general conditioning layer,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 32, no. 1, 2018.
- [9] J. T. Barron, B. Mildenhall, M. Tancik, P. Hedman, R. Martin-Brualla, and P. P. Srinivasan, “Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 5855–5864.
- [10] T. Müller, A. Evans, C. Schied, and A. Keller, “Instant neural graphics primitives with a multiresolution hash encoding,” *ACM transactions on graphics (TOG)*, vol. 41, no. 4, pp. 1–15, 2022.
- [11] M. Tancik, E. Weber, E. Ng, R. Li, B. Yi, T. Wang, A. Kristoffersen, J. Austin, K. Salahi, A. Ahuja *et al.*, “Nerfstudio: A modular framework for neural radiance field development,” in *ACM SIGGRAPH 2023 conference proceedings*, 2023, pp. 1–12.
- [12] B. Fei, J. Xu, R. Zhang, Q. Zhou, W. Yang, and Y. He, “3d gaussian splatting as new era: A survey,” *IEEE Transactions on Visualization and Computer Graphics*, 2024.
- [13] Z. Yu, A. Chen, B. Huang, T. Sattler, and A. Geiger, “Mip-splatting: Alias-free 3d gaussian splatting,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2024, pp. 19 447–19 456.
- [14] T. Lu, M. Yu, L. Xu, Y. Xiangli, L. Wang, D. Lin, and B. Dai, “Scaffold-gs: Structured 3d gaussians for view-adaptive rendering,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 20 654–20 664.
- [15] S. Zhou, H. Chang, S. Jiang, Z. Fan, Z. Zhu, D. Xu, P. Chari, S. You, Z. Wang, and A. Kadambi, “Feature 3dgs: Supercharging 3d gaussian splatting to enable distilled feature fields,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 21 676–21 685.
- [16] Z. Dai, T. Liu, and Y. Zhang, “Efficient decoupled feature 3d gaussian splatting via hierarchical compression,” in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 11 156–11 166.
- [17] R.-Z. Qiu, G. Yang, W. Zeng, and X. Wang, “Language-driven physics-based scene synthesis and editing via feature splatting,” in *European Conference on Computer Vision*. Springer, 2024, pp. 368–383.
- [18] J. Xu, M. Liao, R. P. Kathirvel, and V. M. Patel, “Leveraging thermal modality to enhance reconstruction in low-light conditions,” in *European Conference on Computer Vision*. Springer, 2024, pp. 321–339.
- [19] T. Ye, Q. Wu, J. Deng, G. Liu, L. Liu, S. Xia, L. Pang, W. Yu, and L. Pei, “Thermal-nerf: Neural radiance fields from an infrared camera,” in *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2024, pp. 1046–1053.
- [20] M. Özer, M. Weiherer, M. Hundhausen, and B. Egger, “Exploring multi-modal neural scene representations with applications on thermal imaging,” in *European Conference on Computer Vision*. Springer, 2024, pp. 82–98.
- [21] J. Grün, L. Meyer, M. Weiherer, B. Egger, M. Stamminger, and L. Franke, “Towards Integrating Multi-Spectral Imaging with Gaussian Splatting,” in *Vision, Modeling, and Visualization*, B. Egger and T. Günther, Eds. The Eurographics Association, 2025.
- [22] Q. Ma, C. Zou, D. Wang, J. Wang, L. Xiang, and Z. He, “Beyond darkness: Thermal-supervised 3d gaussian splatting for low-light novel view synthesis,” *arXiv preprint arXiv:2511.13011*, 2025.
- [23] Q. Chen, S. Shu, and X. Bai, “Thermal3d-gs: Physics-induced 3d gaussians for thermal infrared novel-view synthesis,” in *European Conference on Computer Vision*. Springer, 2024, pp. 253–269.
- [24] K. Yang, Y. Liu, Z. Cui, Y. Liu, M. Zhang, S. Yan, and Q. Wang, “Ntr-gaussian: Nighttime dynamic thermal reconstruction with 4d gaussian splatting based on thermodynamics,” in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 691–700.
- [25] M. Nam, W. Park, M. Kim, H. Hur, and S. Lee, “Veta-gs: View-dependent deformable 3d gaussian splatting for thermal infrared novel-view synthesis,” in *2025 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2025, pp. 965–970.
- [26] S. N. Sinha, H. Graf, and M. Weinmann, “Spectralgaussians: Semantic, spectral 3d gaussian splatting for multi-spectral scene representation, visualization and analysis,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 227, pp. 789–803, 2025.
- [27] C. Thirgood, O. Mendez, E. Ling, J. Storey, and S. Hadfield, “Hypergs: Hyperspectral 3d gaussian splatting,” in *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 5970–5979.
- [28] S. Elfving, E. Uchibe, and K. Doya, “Sigmoid-weighted linear units for neural network function approximation in reinforcement learning,” *Neural networks*, vol. 107, pp. 3–11, 2018.
- [29] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, “The unreasonable effectiveness of deep features as a perceptual metric,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 586–595.