



浙江大學  
ZHEJIANG UNIVERSITY

ByteDance

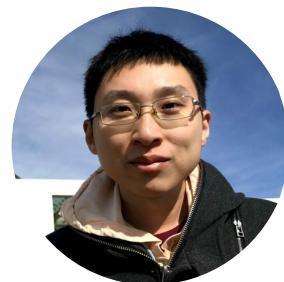


MONASH  
University



NEURAL INFORMATION  
PROCESSING SYSTEMS

# ZPressor: Bottleneck-Aware Compression for Scalable Feed-Forward 3DGs



Weijie Wang<sup>1✉</sup> Donny Y. Chen<sup>2✉</sup>

Zeyu Zhang<sup>3</sup>

Duochao Shi<sup>1</sup>

Akide Liu<sup>3</sup>

Bohan Zhuang<sup>1</sup>

✉ Corresponding authors.

<sup>1</sup>Zhejiang University <sup>2</sup>ByteDance Seed <sup>3</sup>Monash University  
**Weijie Wang (王伟杰)**

College of Computer Science and Technology, Zhejiang University

2025/6/16

# About Me

## »» Weijie Wang (王伟杰)

First-year Ph.D. student @ **ZIP Lab, State Key Lab of CAD&CG, ZJU**

Supervised by **Prof. Bohan Zhuang**

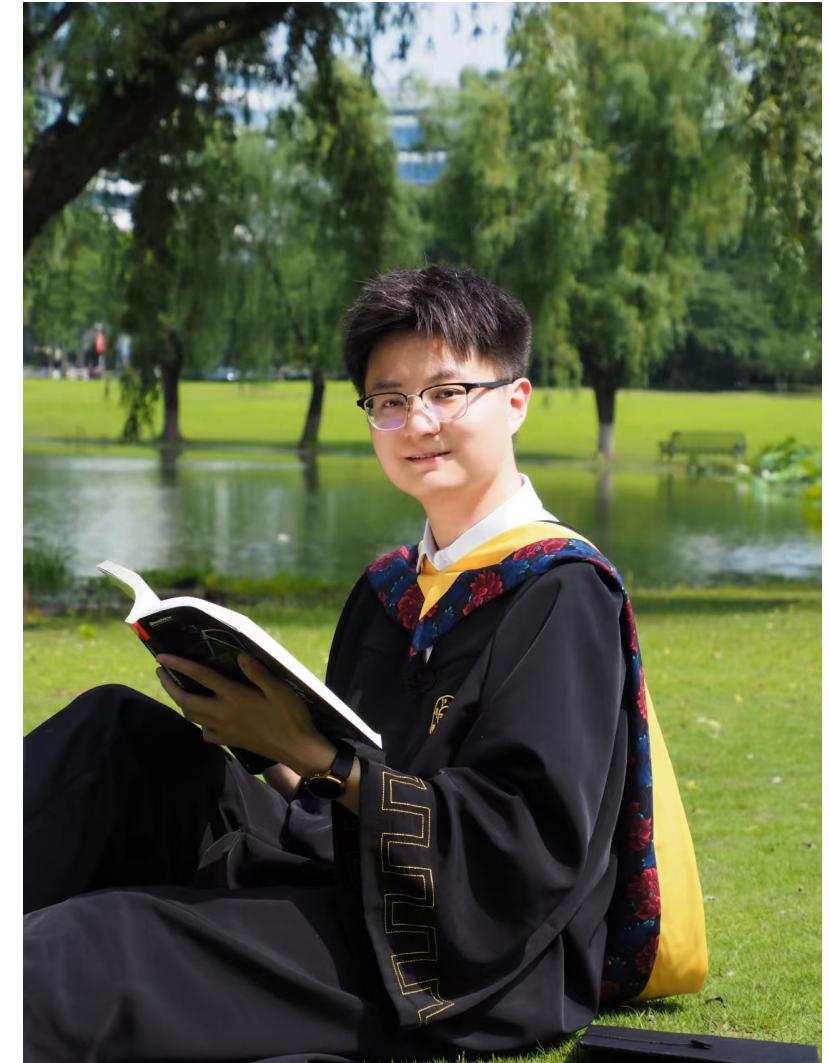
### Research Interest:

- **Feed-Forward Reconstruction:** [ZPressor](#), [PM-Loss](#), [VolSplat](#)
- **Dynamic Reconstruction:** [Street Gaussians](#), [DriveGen3D](#)
- **Interactive Generation:** [WonderTurbo](#)

**Webpage:** <https://lhmd.top>

**Email:** wangweijie@zju.edu.cn

**Wechat:** zju-lhmd

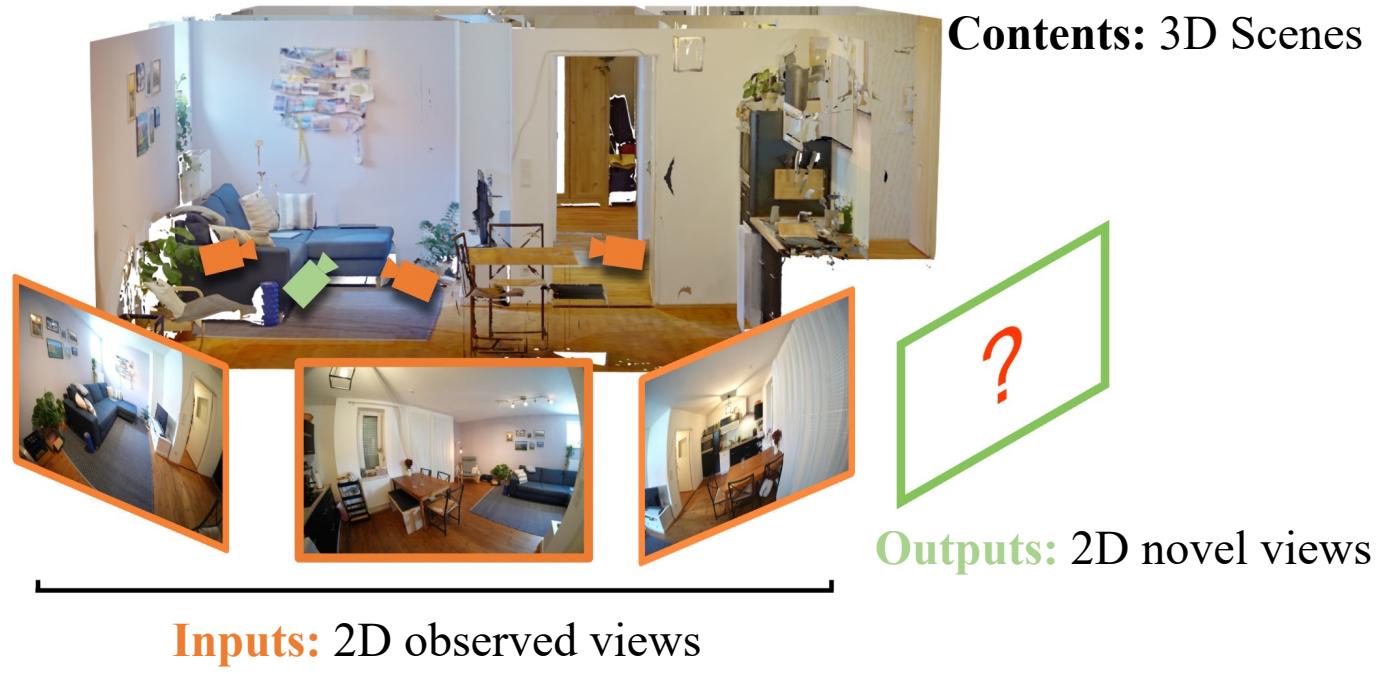


# Background

# Tasks

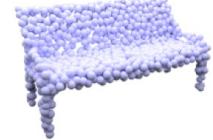
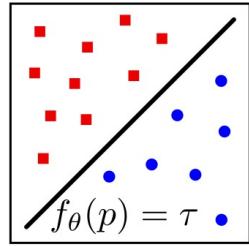
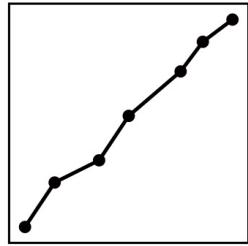
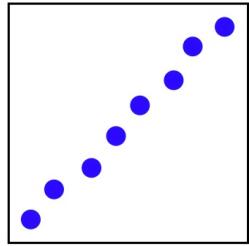
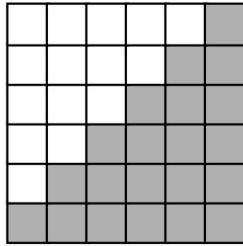


3D Reconstruction



Novel View Synthesis

# 3D Representations

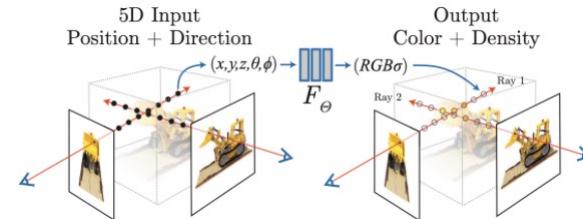


Voxel

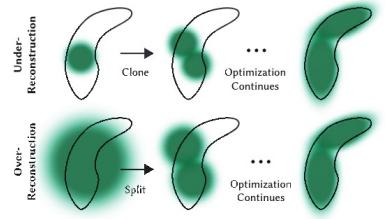
Point Cloud

Mesh

Occupancy  
Networks



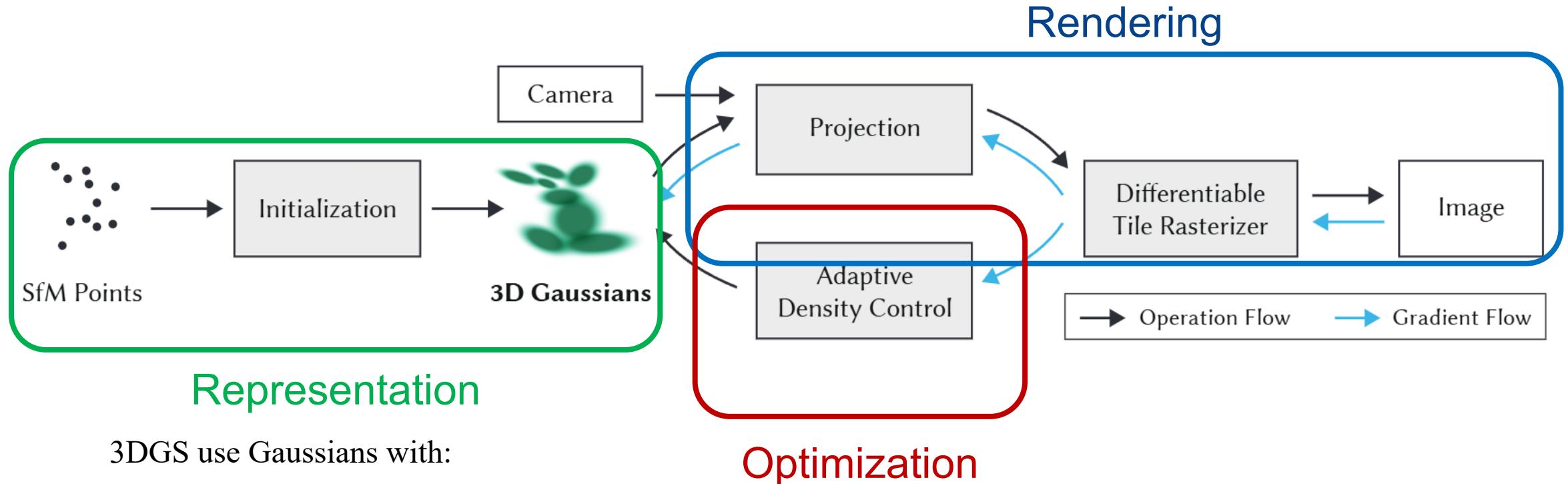
Neural Radiance Field (NeRF)



3D Gaussian Splatting (3DGS)

There is no canonical representation in 3D. We chose 3DGS since it performs the best for NVS in general.

# 3D Gaussian Splatting (3DGS)



3DGS use Gaussians with:

- $\mu$ : Gaussian center position (xyz)
- $\alpha$ : opacity; (how transparent)
- $\Sigma$ : covariance; (scale, rotation)
- $c$ : color; (spherical harmonic)

# Limitations of Per-Scene based 3DGS

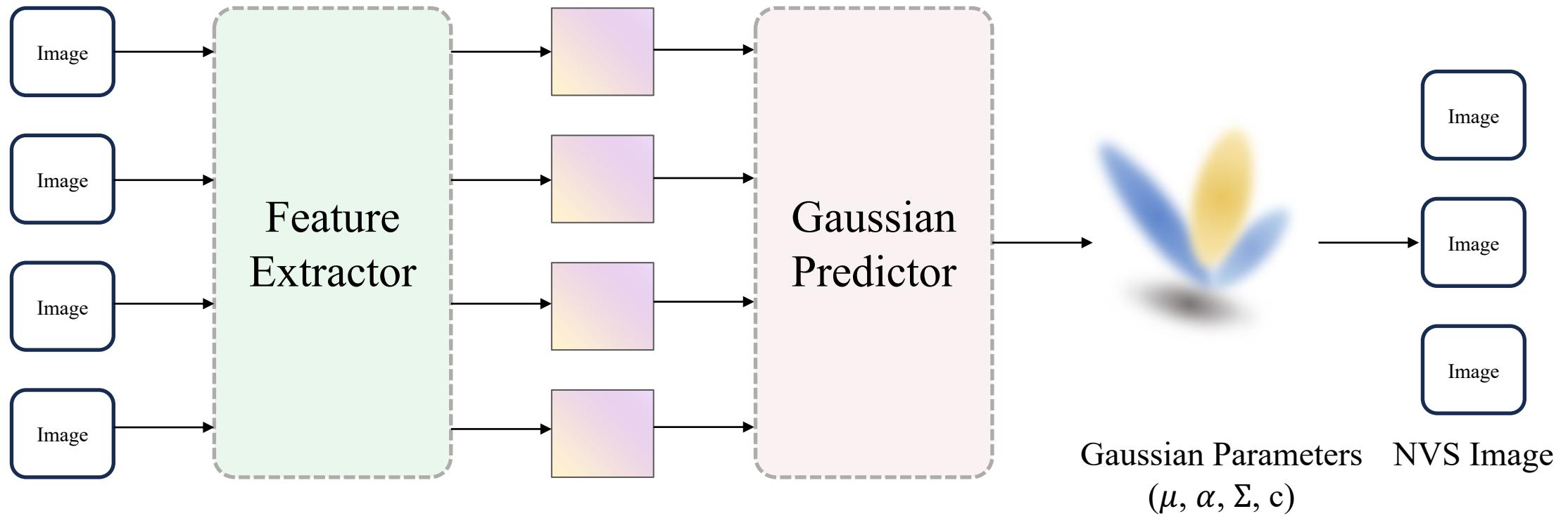
1. **Time:** requires applying the optimization process to *each scene* (20+ mins)
2. **Space:** requires additional permanent storage for the 3D representation of *each scene* (10+ M)



The bicycle scene takes: ~50 mins, ~100 M

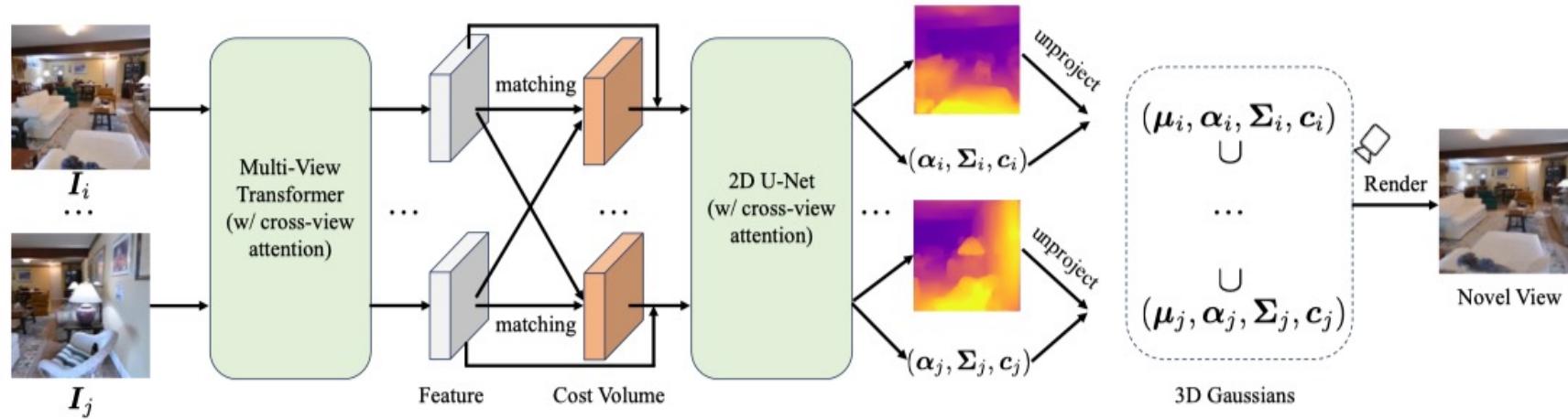
**Note:** Here , we refer to the inria's version of 3DGS;  
NOT those improved models such as sparse-view 3DGS, fast-training 3DGS, 3DGS compression, *etc.*

# Pipeline of Feed-Forward 3DGS



Almost all feed-forward 3DGS networks use this paradigm.

# Example: MVSplat



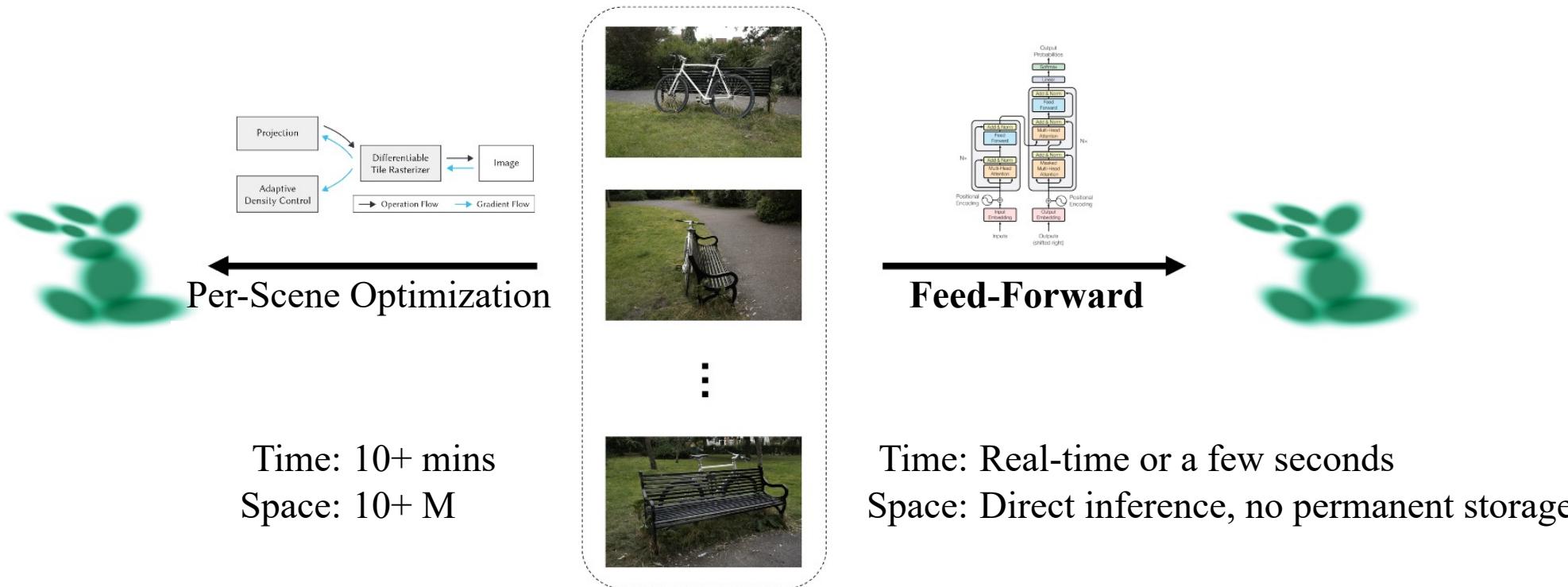
$$f_{\theta} : \{(\mathbf{I}^i, \mathbf{P}^i)\}_{i=1}^K \mapsto \{(\boldsymbol{\mu}_j, \alpha_j, \boldsymbol{\Sigma}_j, \mathbf{c}_j)\}_{j=1}^{H \times W \times K}$$

**Inputs:** Multi-view images, with corresponding camera poses

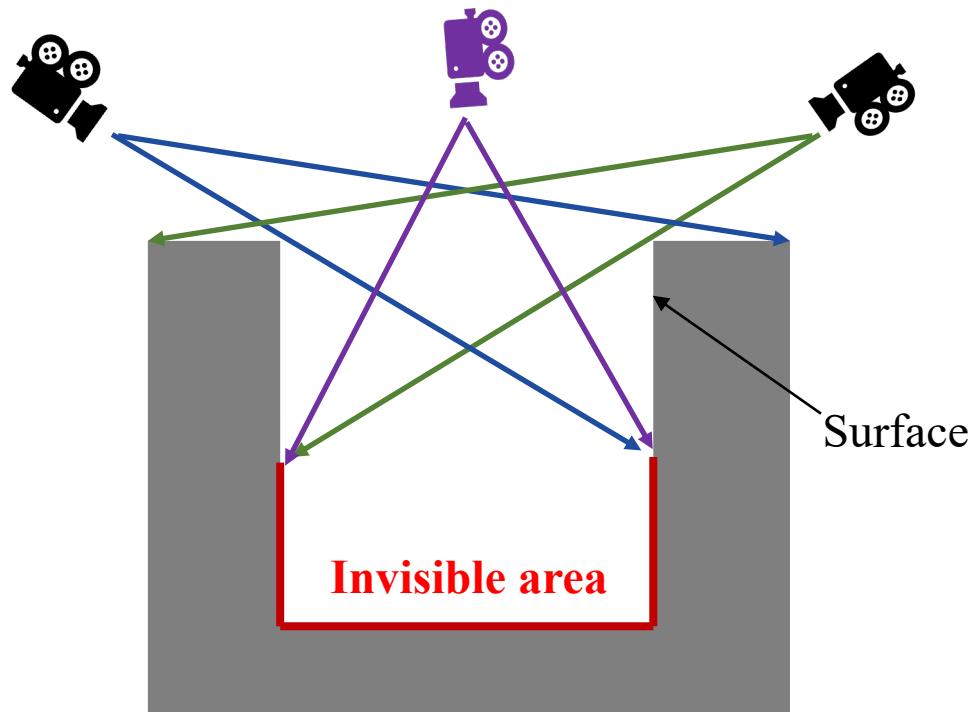
**Outputs:** Pixel-align 3D Gaussians for the scenes

**NVS:** Render the predicted 3DGS from novel viewpoints

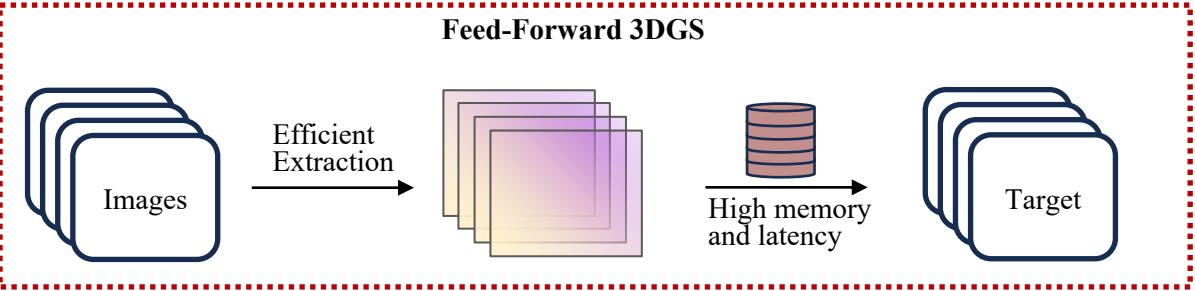
# Per-Scene VS Feed-Forward



# Challenges in Feed-Forward 3DGS



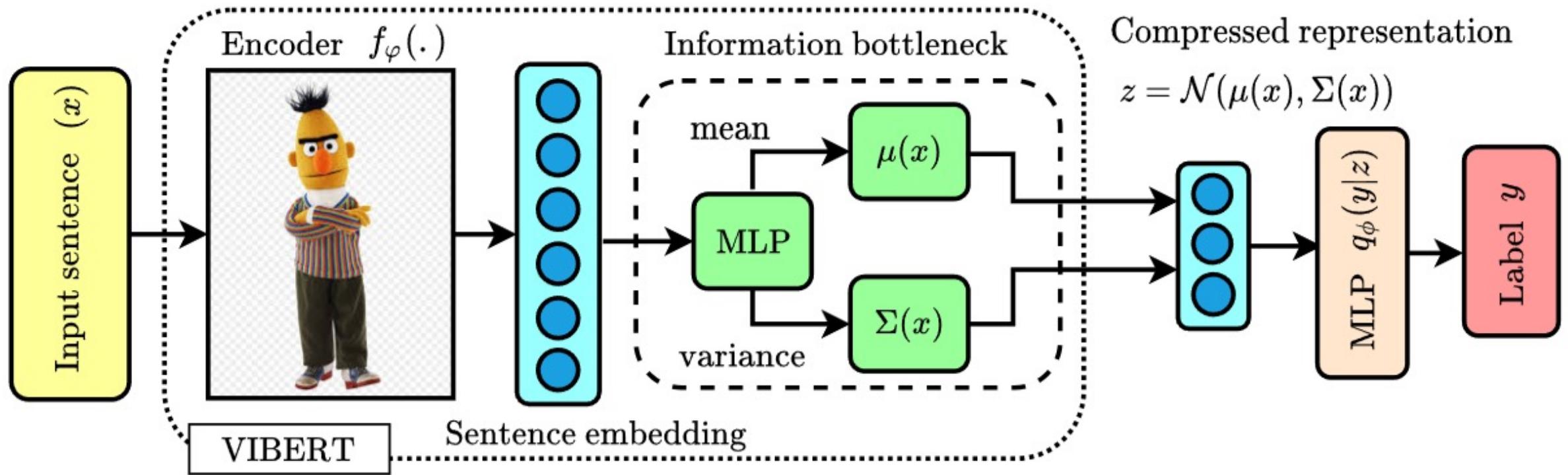
We need denser views to **provide more information**,  
but at the same time not be influenced by  
**redundancy**.



The scalability of feed-forward 3DGS is  
fundamentally constrained by the **limited capacity** of  
their encoders.

# **ZPressor: Bottleneck-Aware Compression for Scalable Feed- Forward 3DGS**

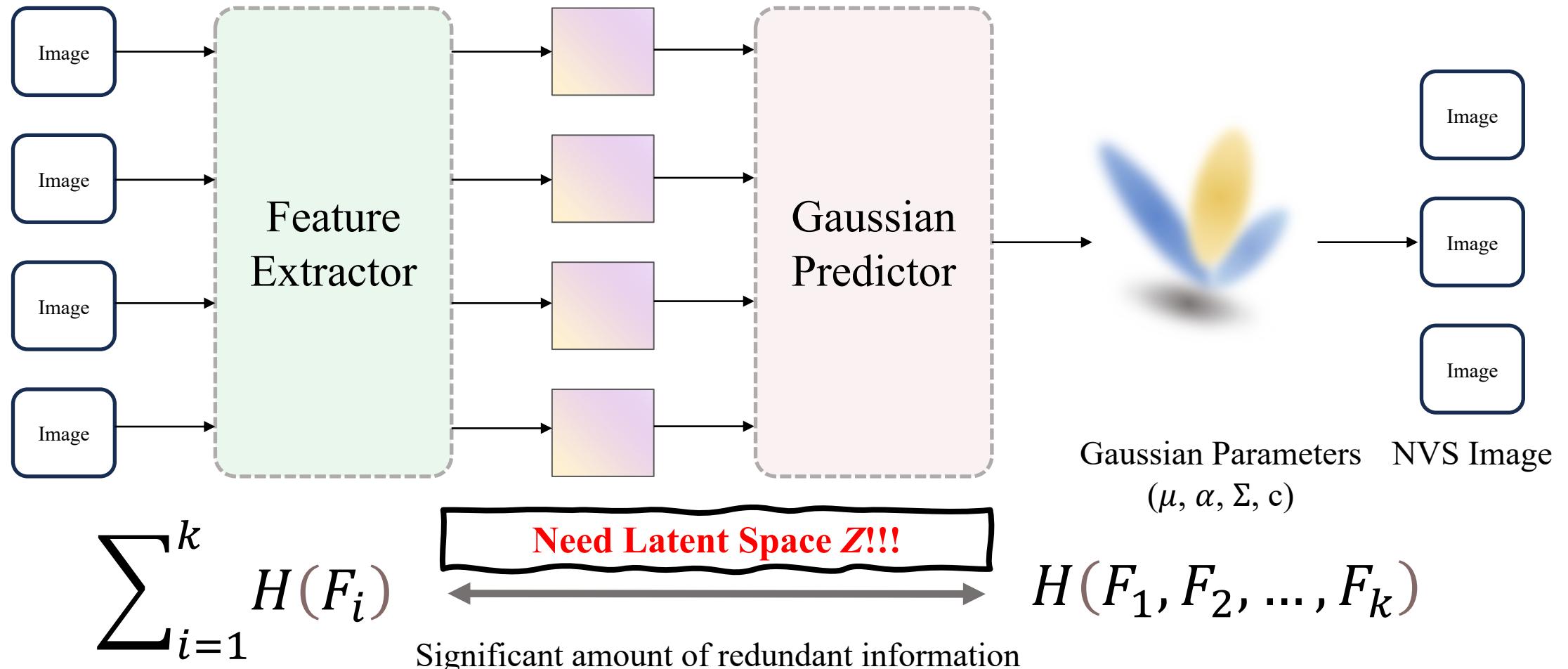
# Information Bottleneck Theory



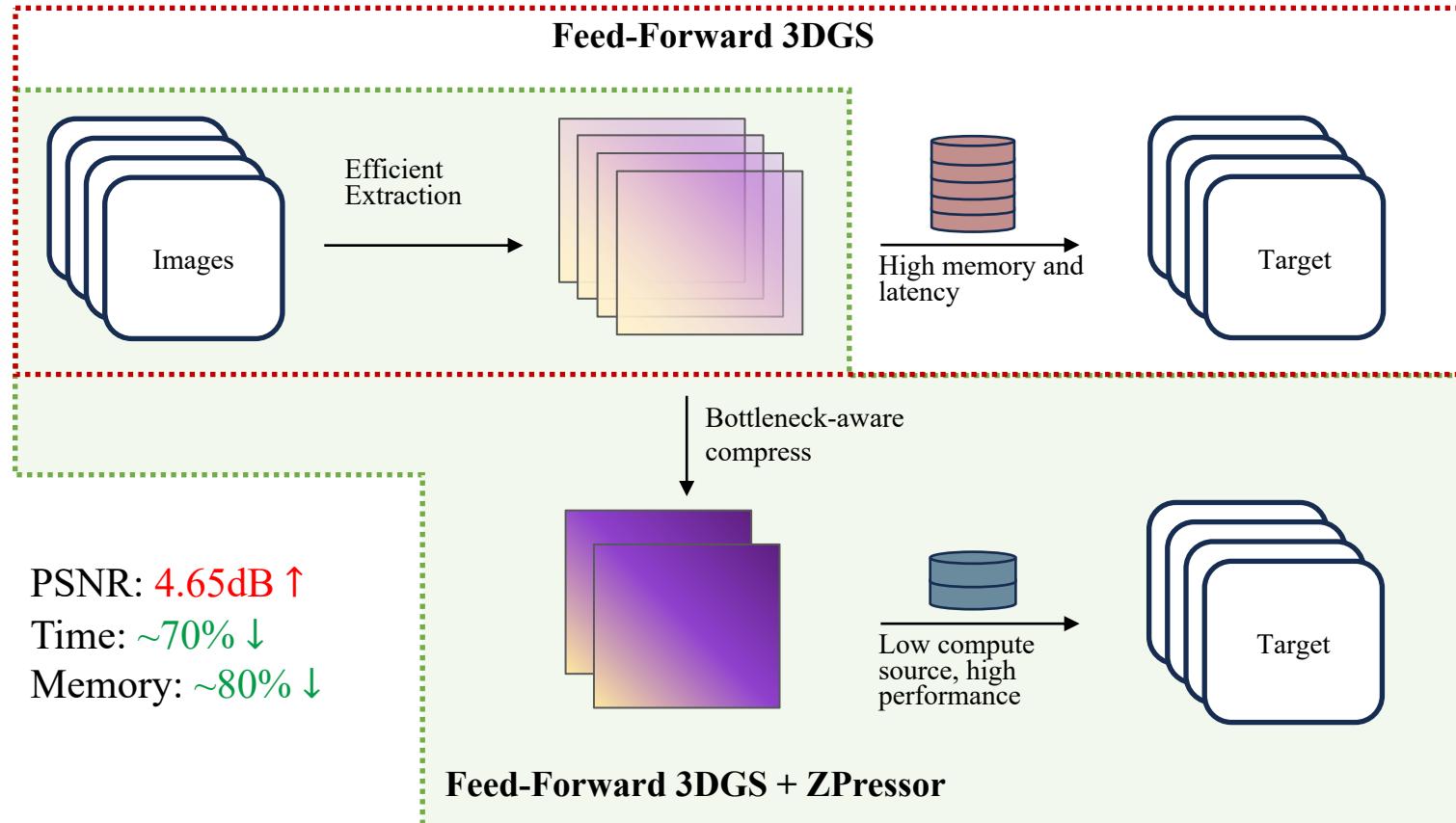
$$I(Z, Y; \theta) = \int dx dy p(z, y|\theta) \log \frac{p(z, y|\theta)}{p(z|\theta)p(y|\theta)}.$$

$$\min_z IB = \underbrace{\beta I(\mathcal{X}, \mathcal{Z})}_{\text{Compression Score}} - \underbrace{I(\mathcal{Z}, \mathcal{Y})}_{\text{Prediction Score}}$$

# Information Flow in FF 3DGS



# Bottleneck-Aware Compression

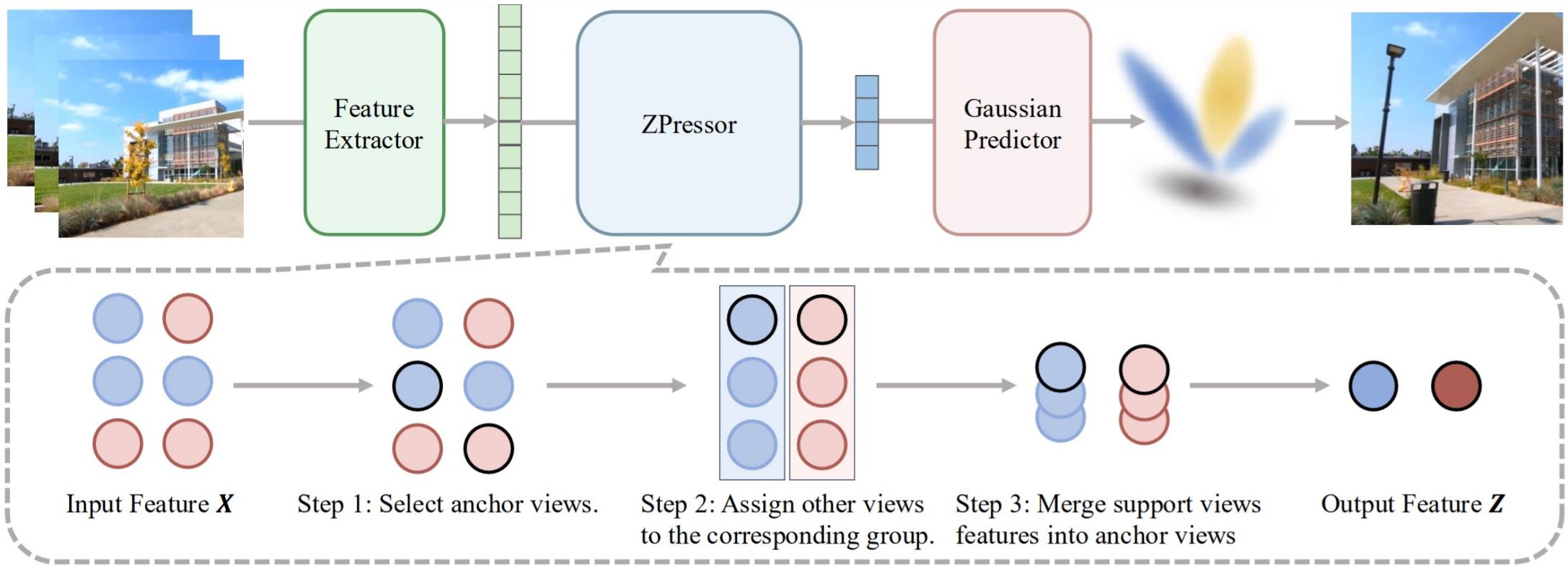


$$\min_{\mathcal{Z}} IB = \underbrace{\beta I(\mathcal{X}, \mathcal{Z})}_{\text{Compression Score}} - \underbrace{I(\mathcal{Z}, \mathcal{Y})}_{\text{Prediction Score}}$$

- 1. Compression Score:** Minimizing  $I(\mathcal{X}, \mathcal{Z})$
- 2. Prediction Score:** Maximizing  $I(\mathcal{Z}, \mathcal{Y})$

Note: The mutual information (MI) of two random variables  $I(\cdot, \cdot)$  is a measure of the mutual dependence between the two variables.

# Zpressor: Overview



**Anchor View Selection**

**Support-to-anchor Assignment**

**Views Information Fusion**

# Anchor View Selection

---

**Algorithm 2** Farthest Point Sampling for Anchor View Selection

---

**Input:** Set of view camera positions  $\mathcal{T} = \{\mathbf{T}_1, \mathbf{T}_2, \dots, \mathbf{T}_K\}$ , Number of anchor views  $N$

**Output:** Indices of the selected anchor views  $\mathcal{S} = \{\mathbf{T}_{a_1}, \mathbf{T}_{a_2}, \dots, \mathbf{T}_{a_n}\}$

Initialize the set of anchor view indices  $\mathcal{S} \leftarrow \emptyset$

Randomly select a random anchor view  $\mathbf{T}_{a_1} \in \mathcal{T}$ , where  $\mathbf{T}_{a_1} \sim \text{Uniform}(\mathcal{T})$

Add  $\mathbf{T}_{a_1}$  to  $\mathcal{S}$ :  $\mathcal{S} \leftarrow \{\mathbf{T}_{a_1}\}$

**for**  $j \leftarrow 2$  to  $N$  **do**

    Initialize a dictionary to store minimum distances  $D \leftarrow \{\}$

**for**  $k \leftarrow 1$  to  $K$  **do**

**if**  $k \notin \mathcal{S}$  **then**

            Calculate the minimum distance  $d_k \leftarrow \min_{i \in \mathcal{S}} \|\mathbf{T}_k - \mathbf{T}_i\|_2$

            Store the distance:  $D[k] \leftarrow d_k$

**end if**

**end for**

    Find the view position  $T_{a_j}$  with the maximum minimum distance:  $T_{a_j} \leftarrow \arg \max_{k \notin \mathcal{S}} D[k]$

    Add  $a_j$  to  $\mathcal{S}$ :  $\mathcal{S} \leftarrow \mathcal{S} \cup \{T_{a_j}\}$

**end for**

**return**  $\mathcal{S}$

---

# Support-to-anchor Assignment

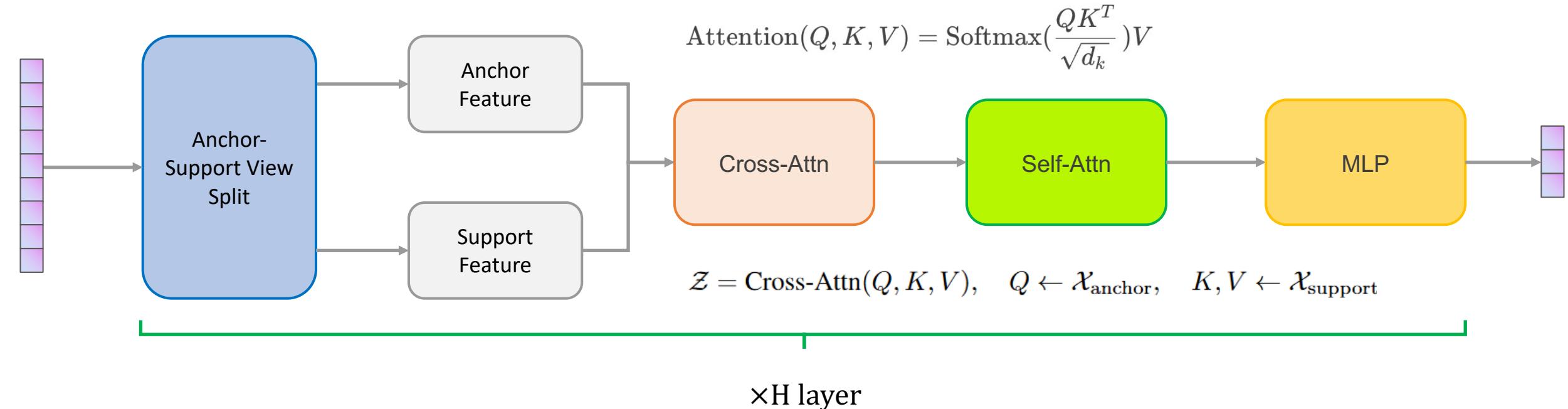


View Groups after Step 1 and Step 2

- Once anchor views are selected, each support view is assigned to its nearest anchor based on **camera position**.
- This grouping ensures that support views, which capture complementary scene details, are paired with **the most spatially relevant** anchor views.
- This pairing thereby ensures the effectiveness of information fusion.
- Formally, the cluster assignment to the  $i$ -th anchor view can be denoted as:

$$\mathcal{C}_i = \{f(\mathbf{T}) \in \mathcal{X}_{\text{support}} \mid \|\mathbf{T} - \mathbf{T}_{a_i}\| \leq \|\mathbf{T} - \mathbf{T}_{a_j}\|, \forall j \neq i\}$$

# Views Information Fusion



Design of Feature Fusion Networks. Feature Fusion  
by Cross-Attention, Self-Attention and MLP.

# Results on DL3DV with DepthSplat

Views	Methods	PSNR↑	SSIM↑	LPIPS↓
36 views	DepthSplat	19.23	0.666	0.286
	DepthSplat + ZPressor	<b>23.88</b> <sub>+4.65</sub>	<b>0.815</b> <sub>+0.149</sub>	<b>0.150</b> <sub>-0.136</sub>
24 views	DepthSplat	20.38	0.711	0.253
	DepthSplat + ZPressor	<b>24.26</b> <sub>+3.88</sub>	<b>0.820</b> <sub>+0.109</sub>	<b>0.147</b> <sub>-0.106</sub>
16 views	DepthSplat	22.07	0.773	0.195
	DepthSplat + ZPressor	<b>24.25</b> <sub>+2.18</sub>	<b>0.819</b> <sub>+0.046</sub>	<b>0.147</b> <sub>-0.047</sub>
12 views	DepthSplat	23.32	0.807	0.162
	DepthSplat + ZPressor	<b>24.30</b> <sub>+0.97</sub>	<b>0.821</b> <sub>+0.014</sub>	<b>0.146</b> <sub>-0.017</sub>

# Results on RE10K with MVsplat

Views	Methods	PSNR↑	SSIM↑	LPIPS↓
36 views	pixelSplat	OOM	OOM	OOM
	pixelSplat + ZPressor	<b>26.59</b>	<b>0.849</b>	<b>0.225</b>
	MVsplat	24.19	0.851	0.155
	MVsplat + ZPressor	<b>27.34<sub>+3.15</sub></b>	<b>0.893<sub>+0.042</sub></b>	<b>0.113<sub>-0.042</sub></b>
24 views	pixelSplat	OOM	OOM	OOM
	pixelSplat + ZPressor	<b>26.72</b>	<b>0.851</b>	<b>0.223</b>
	MVsplat	25.00	0.871	0.137
	MVsplat + ZPressor	<b>27.49<sub>+2.49</sub></b>	<b>0.895<sub>+0.024</sub></b>	<b>0.111<sub>-0.026</sub></b>
16 views	pixelSplat	OOM	OOM	OOM
	pixelSplat + ZPressor	<b>26.81</b>	<b>0.853</b>	<b>0.221</b>
	MVsplat	25.86	0.888	0.120
	MVsplat + ZPressor	<b>27.60<sub>+1.74</sub></b>	<b>0.896<sub>+0.008</sub></b>	<b>0.110<sub>-0.010</sub></b>
8 views	pixelSplat	26.19	0.852	<b>0.215</b>
	pixelSplat + ZPressor	<b>26.86<sub>+0.67</sub></b>	<b>0.854<sub>+0.002</sub></b>	<b>0.219<sub>+0.004</sub></b>
	MVsplat	26.94	<b>0.902</b>	<b>0.107</b>
	MVsplat + ZPressor	<b>27.72<sub>+0.78</sub></b>	<b>0.897<sub>-0.005</sub></b>	<b>0.109<sub>+0.002</sub></b>

# Qualitative comparison

## Visualization on DL3DV (36 Input Views)



a62c330f5403e2e41a82a74c4e865b705c5706843b992fae2fe2e538b122d984



63798f5c6fbfc4eb686268248b8ecbc8d87d920b2bcce967eeaedfd3b3b6d82

# Analysis of model efficiency

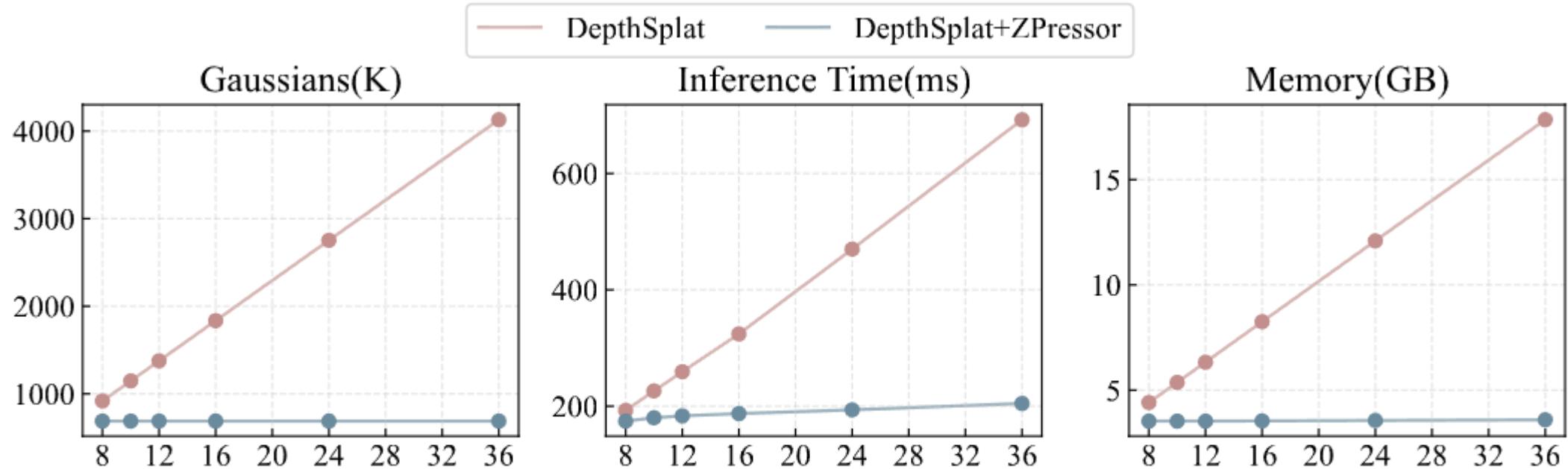
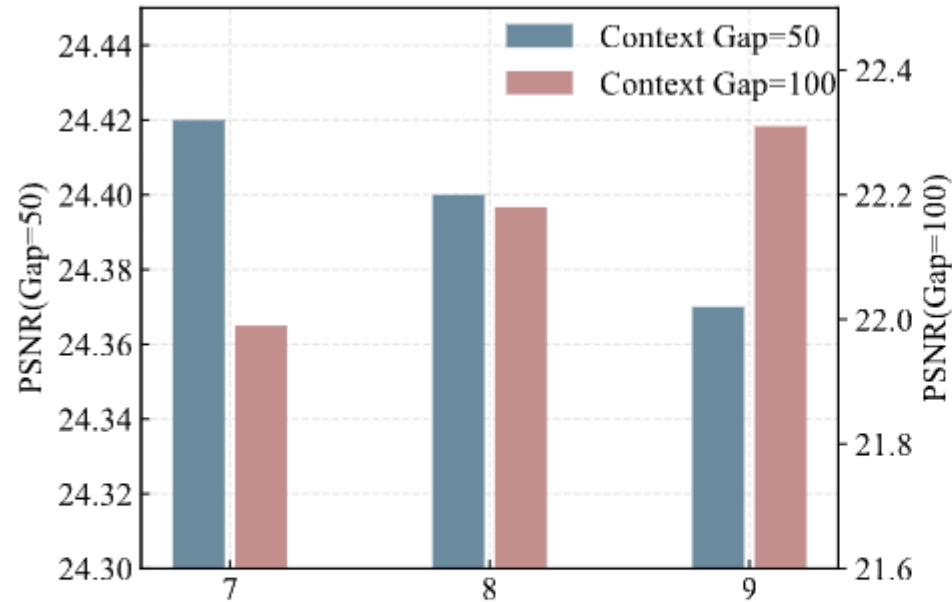


Figure 5: **Efficiency analysis.** We report the number of Gaussians (K), inference time (ms) and peak memory (GB) of DepthSplat [12] and DepthSplat with ZPressor.

# Analysis of the Information Bottleneck

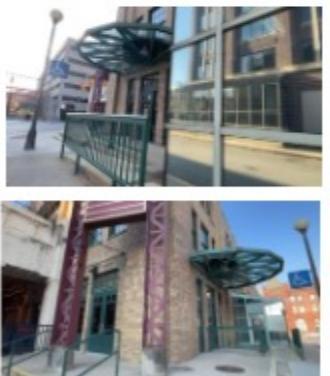


**Figure 6: Analysis of the bottleneck constraint.**  
We compare the performance of ZPressor in different scale of scene coverage.

# Limitations



⋮



⋮



Inputs (~500 views)

DepthSplat + ZPressor

ZPressor exhibits limitations when processing scenarios with an **extremely high** density of input views.

# More Information



Paper, code and model  
will be available on our  
project page



ZIP Lab. We are currently  
recruiting research  
assistants for 3D LM topic



Weijie Wang's homepage.  
Actively seeking  
cooperation

# THANK YOU

**Weijie Wang (王伟杰)**

College of Computer Science and Technology, Zhejiang University

2025/6/16

# Cross Dataset Generalization on ACID

Views	Methods	PSNR↑	SSIM↑	LPIPS↓
36 views	pixelSplat	OOM	OOM	OOM
	pixelSplat + Ours	<b>27.78</b>	<b>0.823</b>	<b>0.238</b>
	MVSplat	24.89	0.812	0.179
	MVSplat + Ours	<b>28.16<sub>+3.27</sub></b>	<b>0.853<sub>+0.041</sub></b>	<b>0.145<sub>-0.034</sub></b>
24 views	pixelSplat	OOM	OOM	OOM
	pixelSplat + Ours	<b>27.91</b>	<b>0.825</b>	<b>0.235</b>
	MVSplat	25.46	0.829	0.167
	MVSplat + Ours	<b>28.33<sub>+2.87</sub></b>	<b>0.856<sub>+0.027</sub></b>	<b>0.142<sub>-0.025</sub></b>
16 views	pixelSplat	OOM	OOM	OOM
	pixelSplat + Ours	<b>27.97</b>	<b>0.826</b>	<b>0.234</b>
	MVSplat	26.08	0.844	0.156
	MVSplat + Ours	<b>28.42<sub>+2.34</sub></b>	<b>0.858<sub>+0.014</sub></b>	<b>0.141<sub>-0.015</sub></b>
8 views	pixelSplat	26.69	0.807	0.260
	pixelSplat + Ours	<b>28.05<sub>+1.36</sub></b>	<b>0.828<sub>+0.021</sub></b>	<b>0.234<sub>-0.026</sub></b>
	MVSplat	27.89	<b>0.864</b>	<b>0.140</b>
	MVSplat + Ours	<b>28.60<sub>+0.71</sub></b>	0.860 <sub>-0.004</sub>	<b>0.140<sub>-0.000</sub></b>

# Ablation Studies

Table 4: **Ablation study of our method with DepthSplat [12] on the DL3DV dataset [17]**. Models are evaluated by rendering eight novel views using 12 input views.

Methods	PSNR↑	SSIM↑	LPIPS↓	Time (s)	Peak Memory (GB)
DepthSplat + ZPressor	<b>24.30</b>	<b>0.821</b>	<b>0.146</b>	0.184	3.80
w/o multi-blocks	24.18	0.817	0.149	<b>0.140</b>	<b>3.79</b>
w/o self-attention	23.85	0.810	0.156	0.183	3.80
DepthSplat	23.32	0.808	0.162	0.260	6.80

**Note:** All ablation models and training settings will be available on our GitHub project.