



# LoBE-GS: Load-Balanced and Efficient 3D Gaussian Splatting for Large-Scale Scene Reconstruction

Sheng-Hsiang Hung<sup>1</sup>, Jeng Wen Joshua Lean<sup>1</sup>, Ting-Yu Yen<sup>1</sup>, Wei-Fang Sun<sup>2</sup>, Simon See<sup>2</sup>, Shih-Hsuan Hung<sup>1</sup>, Chih-Yuan Yao<sup>3</sup>, Hung-Kuo Chu<sup>1</sup>

<sup>1</sup>National Tsing Hua University <sup>2</sup>NVIDIA AI Technology Center (NVAITC) <sup>3</sup>National Taiwan University of Science and Technology



## Abstract

We present **LoBE-GS**, a load-balanced and efficient framework for large-scale 3D Gaussian Splatting. Unlike prior large-scale 3DGS methods such as CityGS, DOGS and VastGS, which focus on reconstruction and partitioning heuristics but ignore runtime imbalance, LoBE-GS uses a lightweight runtime-correlated proxy to steer partitioning. Guided by the proxy, we apply (i) load balance-aware scene partitioning to produce blocks with similar predicted runtime, (ii) fast camera selection to minimize partition overhead, and (iii) visibility cropping and selective densification to reduce per-block optimization. Across large-scale scenes, LoBE-GS achieves up to **2x** end-to-end training speedup over previous methods while preserving rendering quality. After training, we prune and merge the block outputs, export the unified model into USDZ format, and stream it via NVIDIA Omniverse Kit for real-time visualization, collaborative review, and interactive playback.

## Motivation

**Runtime goal.** Partition-merge with parallel fine stage:

$$T_{E2E} = T_{coarse} + T_{partition} + \max_b T_{fine}.$$

⇒ **Minimize the slowest block.**

**What predicts runtime?** Correlation study across 5 datasets shows:

- Weak: block area  $A^{(b)}$ , in-block Gaussian count  $G_{blk}^{(b)}$ , camera count  $C^{(b)}$ , and average visible per camera  $G_{avg\_blk}^{(b)}$ .

- Strong** : total visible Gaussians per block  $G_{vis}^{(b)}$  before fine optimization.

•••••  $A$  ( $r=0.50$ ) •••••  $C^{(b)}$  ( $r=0.12$ ) •••••  $G_{vis}^{(b)}$  ( $r=0.88$ ) •••••  $G_{blk}^{(b)}$  ( $r=0.52$ ) •••••  $G_{avg}^{(b)}$  ( $r=0.11$ )

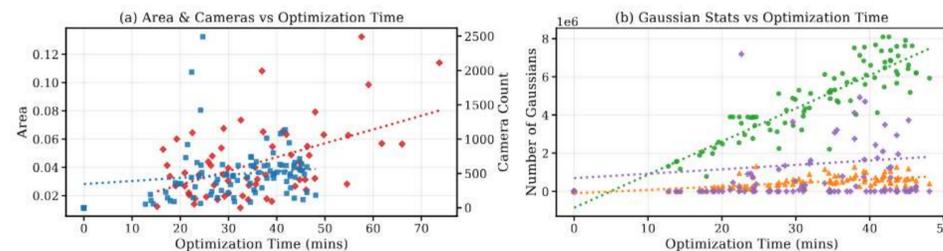


Figure 1. Correlation between per-block training time and block-level statistics under CityGS's partitioning.

## Methodology

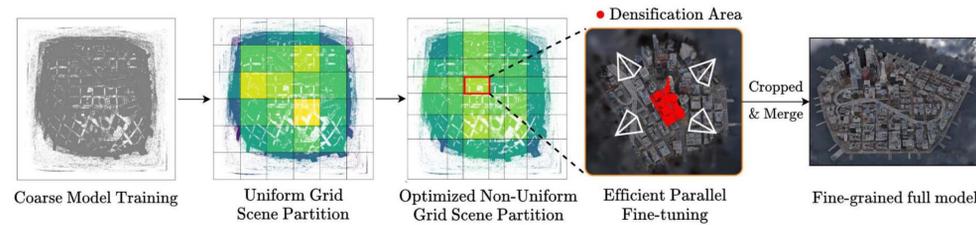


Figure 2. Overview of our pipeline

We train a coarse 3DGS model, optimize grid cuts with load-balance-aware partitioning, apply visibility cropping and selective densification during parallel fine-tuning, then prune outside regions and merge the blocks into a single high-quality model.

### Load Balance-Aware Scene Partition (LB-SP)

**Goal:** Minimize the slowest block's cost by balancing initial visible Gaussians per block (a strong runtime proxy).

**How:** Bayesian Optimization (BO) over vertical/horizontal cuts with ordering constraints; evaluate maximum visible Gaussians each round; pick best cuts.

**Implementation note:** visible Gaussians per block computed in **NVIDIA Warp** for near-CUDA speed across BO iterations.

### Fast Camera Selection (FCS)

**Key idea:** Only **N** projections: render per-view depth once from the coarse model, back-project to a dense point cloud, then count points inside each block to get a visibility ratio

**Assign views:** For each block, select cameras whose visibility of that block exceeds a chosen threshold, and reuse those per-view point clouds across BO iterations.

### Visibility Cropping (VC) & Selective Densification (SD)

**VC:** Before fine-training block  $b$ , keep only Gaussians visible from its assigned **cameras** instead of the full coarse model → big cut in optimizer updates and memory.

**SD:** During fine-training, densify only inside the block; visible-but-out-of-block Gaussians are optimized but not densified → fewer new points, faster training.

## Experiments

All Experiments ran on 10 nodes on **OVX server** (8 **NVIDIA L40 GPUs/node**, 80 GPUs total) via **Run:ai** platform.

Methods	MatrixCity-Aerial			Mill19			UrbanScene3D			Methods	MatrixCity-Aerial			Mill19			UrbanScene3D		
	PSNR	SSIM	LPIPS	PSNR	SSIM	LPIPS	PSNR	SSIM	LPIPS		PSNR	SSIM	LPIPS	PSNR	SSIM	LPIPS	PSNR	SSIM	LPIPS
3DGS <sup>†</sup>	23.67	0.735	0.384	22.97	0.749	0.291	21.25	0.814	0.239	VastGS <sup>†</sup>	28.33	0.835	0.220	23.50	0.735	0.245	21.83	0.730	0.261
CityGS	27.46	0.865	0.204	23.66	0.796	<b>0.237</b>	<b>21.70</b>	0.825	0.221	DOGS	28.58	0.847	0.219	24.26	0.762	<b>0.231</b>	23.18	0.772	0.232
Ours	<b>27.74</b>	<b>0.875</b>	<b>0.186</b>	<b>23.87</b>	<b>0.797</b>	0.240	21.33	<b>0.826</b>	<b>0.213</b>	Ours	<b>28.91</b>	<b>0.879</b>	<b>0.187</b>	<b>24.68</b>	<b>0.795</b>	0.241	<b>23.55</b>	<b>0.832</b>	<b>0.213</b>

Table 1. Quantitative Comparison

Table 2. Quantitative Comparison on color-correct metrics

Methods	MatrixCity-Aerial				Mill19				UrbanScene3D			
	$T_{coarse}$	$T_{partition}$	$\max T_{fine}$	$T_{E2E}$	$T_{coarse}$	$T_{partition}$	$\max T_{fine}$	$T_{E2E}$	$T_{coarse}$	$T_{partition}$	$\max T_{fine}$	$T_{E2E}$
3DGS <sup>†</sup>	01:50	-	-	01:50	01:20	-	-	01:20	01:01	-	-	01:01
VastGS <sup>†</sup>	-	00:48	01:13	02:01	-	<b>00:05</b>	00:42	<b>00:47</b>	-	00:17	00:40	00:57
CityGS	00:52	01:39	01:00	03:31	01:03	00:15	01:10	02:28	00:43	00:20	01:04	02:07
Ours	<b>00:38</b>	<b>00:16</b>	<b>00:30</b>	<b>01:24</b>	<b>00:07</b>	<b>00:36</b>	01:07	<b>00:21</b>	<b>00:07</b>	<b>00:28</b>	<b>00:55</b>	

Table 3. End-to-End time comparison

## Interactive 3DGS Map via Omniverse Streaming

**Challenge:** Large Scale 3DGS (e.g., 50M Gaussians, 8 GB) cannot run on phones/laptops — client rendering is 1–4 FPS or unusable.

**How:** Browser-based rendering using **RTX** and **NuRec** via **Omniverse Streaming** on an RTX 5090, with USDZ assets converted by **3DGRUT**. WebRTC streams 30 FPS video to a Three.js frontend. UE5/OSM metadata (building boxes, paths) and invisible Three.js overlay objects enable clickable scene interactions and road-level navigation. Low-latency camera transform sync keeps frontend/backend views aligned.

**Result:** Stable **30 FPS** on laptops & phones ( $\approx$ **7–30x** speedup vs client-side Three.js).



Figure 3, 4. Interactive 3DGS map of NTHU campus