



#### Introduction

#### **Background:**

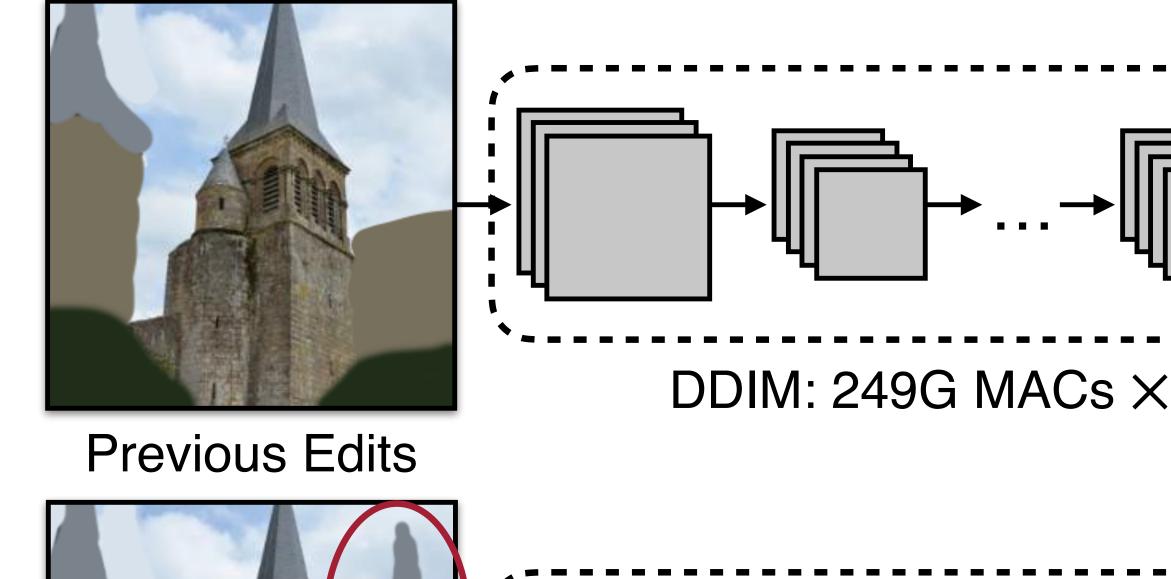
Carnegie

University

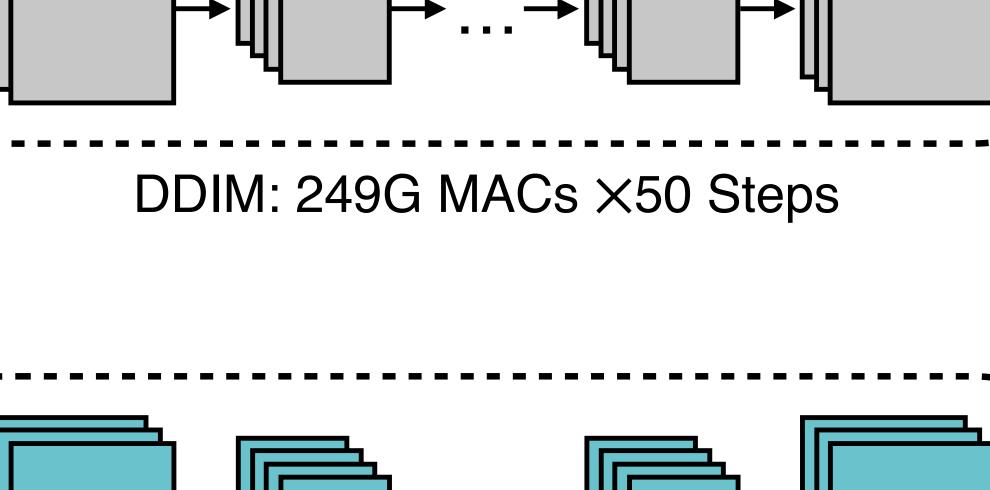
Mellon

Recent models often synthesize the entire image, even for a minor edit, wasting significant computation.

Computed Activations for Previous Edits







DDIM: 249G MACs ×50 Steps



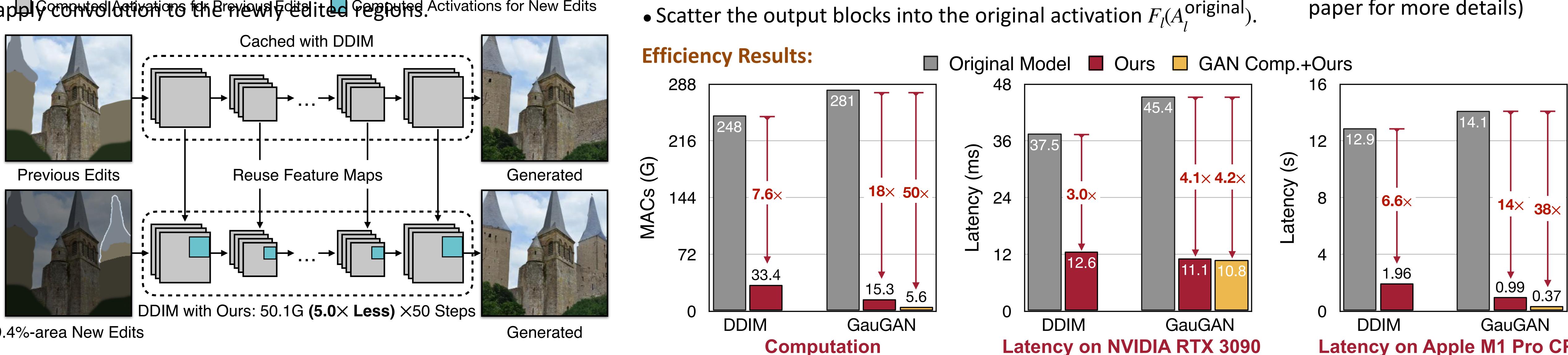
Generated



9.4%-area New Edits

### Idea — Spatially Sparse Inference (SSI):

Cache the features of previous edits. Then reuse them and only apply conved at one for Breviewal Edited Cogregated Activations for New Edits



Active Indices

**Tiling-based Sparse Convolution:** 

• Precompute the features of original image.

• Gather the active blocks along batch dimension.

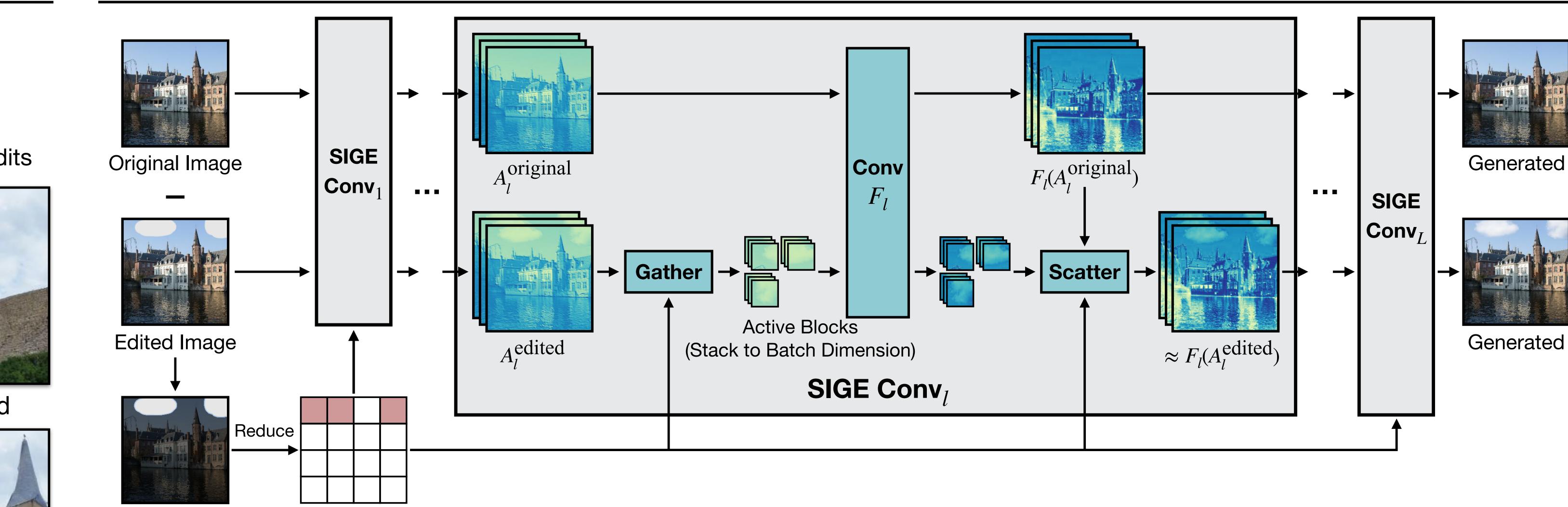
• Compute a difference mask between the original and edited image.

• Divide mask to small blocks and reduce it to active block indices.

## **Efficient Spatially Sparse Inference for Conditional GANs and Diffusion Models**

Muyang Li<sup>1</sup>, Ji Lin<sup>2</sup>, Chenlin Meng<sup>3</sup>, Stefano Ermon<sup>3</sup>, Song Han<sup>2</sup>, and Jun-Yan Zhu<sup>1</sup>

<sup>1</sup>Carnegie Mellon University, <sup>2</sup>Massachusetts Institute of Technology, and <sup>3</sup>Stanford University

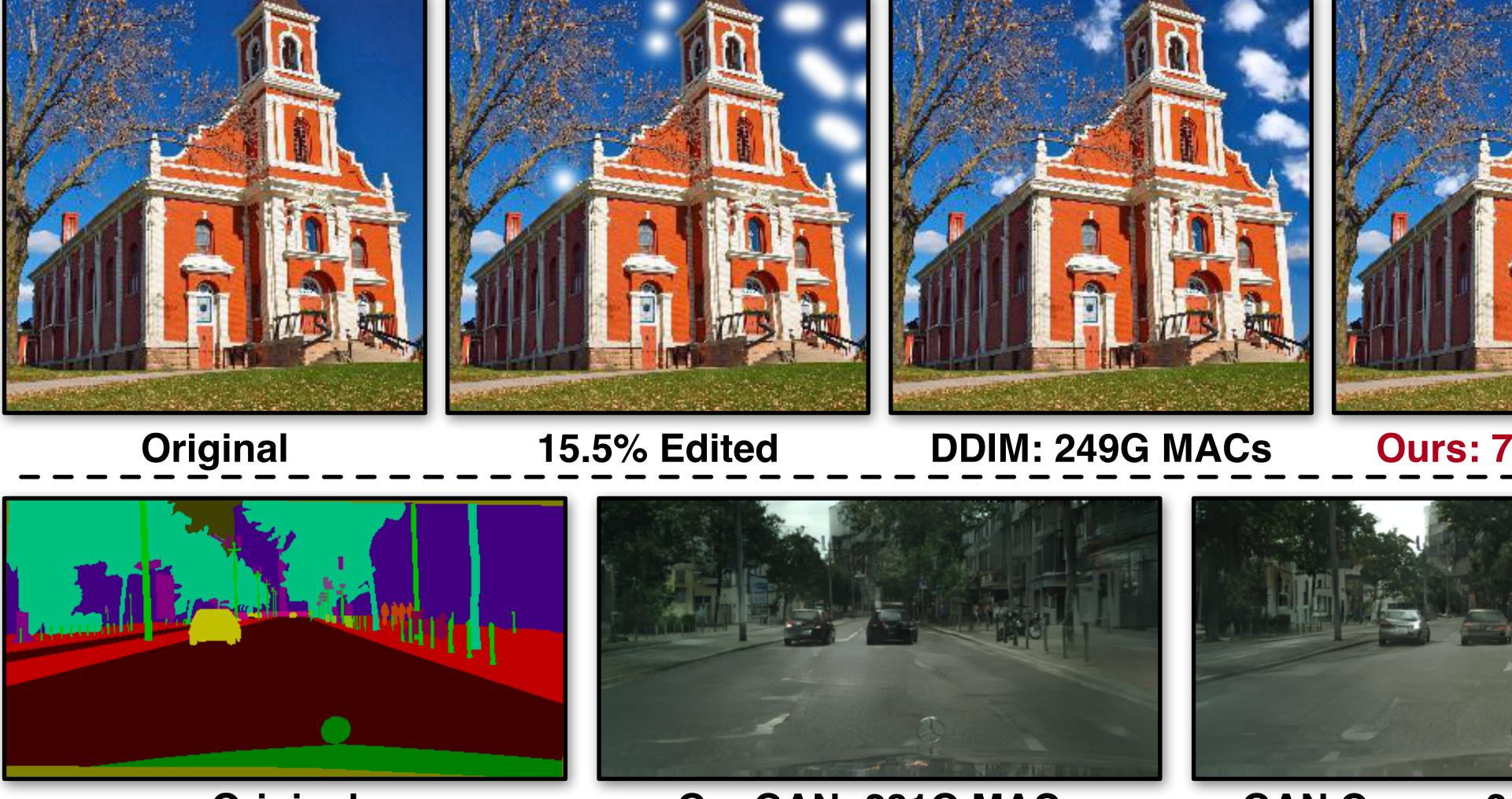


Generated

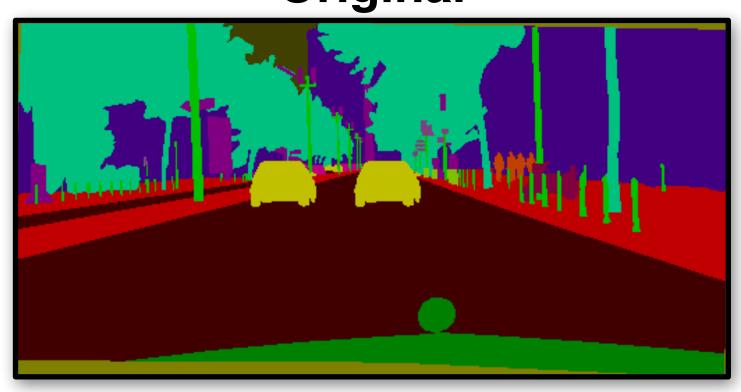
Difference Mask

## **Sparse Incremental Generative Engine (SIGE)**

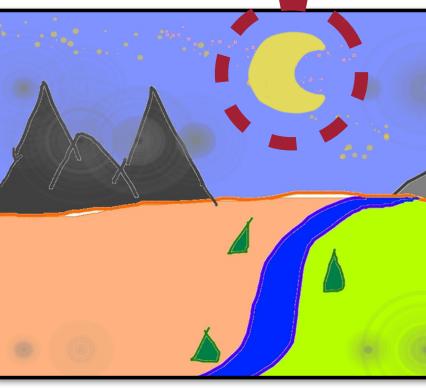




#### Original



1.18% Edited



2.1% Edited

#### **References:**

- [1] SBNet: Sparse Blocks Network for Fast Inference, Ren et al., CVPR 2018
- [2] Semantic Image Synthesis with Spatially-Adaptive Normalization (GauGAN), Park et al., CVPR 2019
- [3] GAN Compression: Efficient Architectures for Interactive Conditional GANs, Li et al., CVPR 2020 [4] Denoising Diffusion Implicit Model (DDIM), Song et al., ICLR 2021
- [5] High-resolution image synthesis with latent diffusion models, Rombach et al., CVPR 2022
- [6] Progressive Distillation for Fast Sampling of Diffusion Models, Salimans et al., ICLR 2022
- [7] SDEdit: Guided Image Synthesis and Editing with Stochastic Differential Equations, Meng et al., ICLR 2022

- Refuger the original statistics to replace normalization
- with scale+shift.
- Kernel fusion. (Refer to our paper for more details)

Latency on Apple M1 Pro CPU

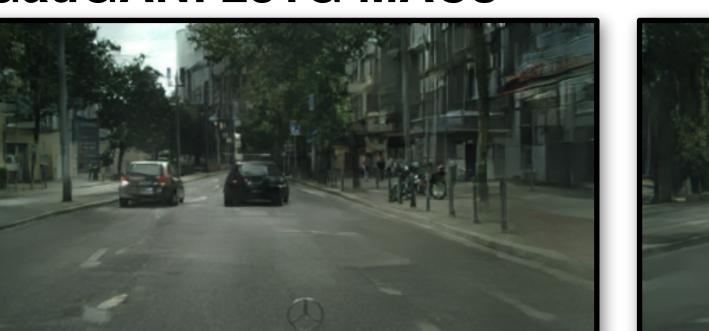






#### **Quality Results on DDIM and GauGAN**

# GauGAN: 281G MACs



#### GAN Comp.: 31.2G (9.0×)



GAN Comp.+Ours: 5.59G (50×)

#### **Extension to Stable-Diffusion**

A fantasy landscape, trending on artstation

Ours: 15.3G (18x)



Full: 1358GMACs 369ms



Ours: 189GMACs (7.2x) 55ms(6.8x)