# Generative Adversarial Transformers

**Drew A. Hudson & Larry Zitnick**

ICML 2021

**facebook** AI    **Stanford**

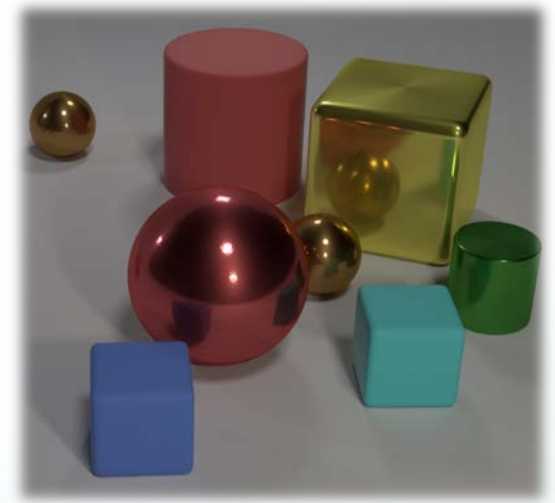Most GANs use Convolution

fake

# Convolutional GANs struggle with long-range dependencies



**Consistency**   **Global Structure**   **Compositionality**

# Transformers

O(n²)

use Self-Attention
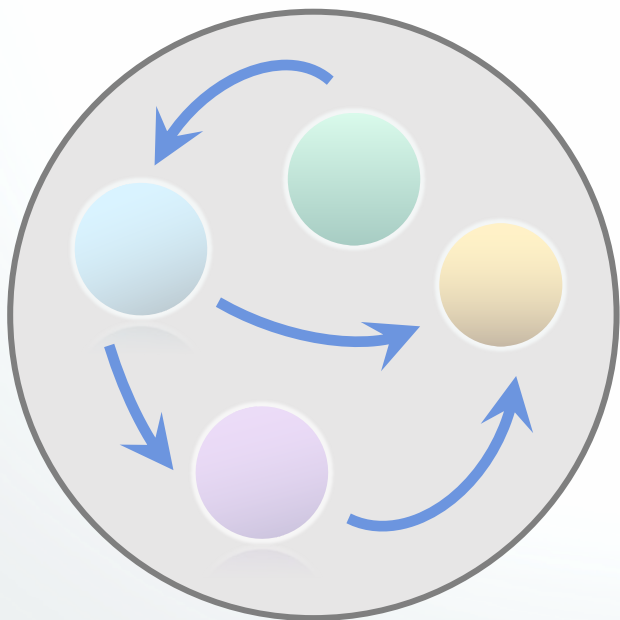
to model long-range

interactions

but they don't scale

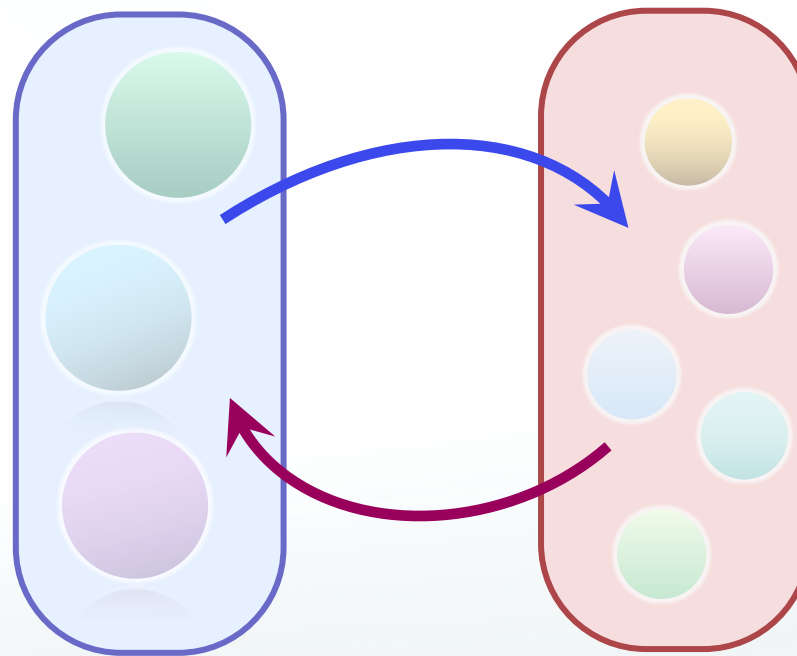# Can we efficiently apply transformers for image generation?

# Generative Adversarial Transformers
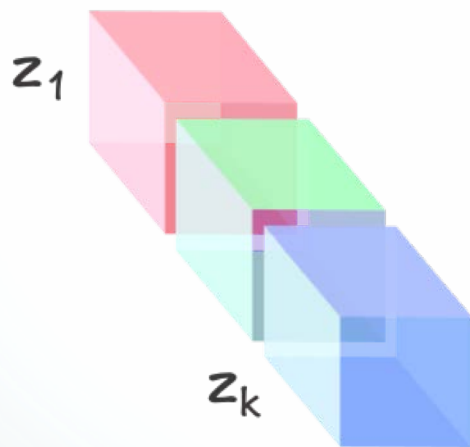
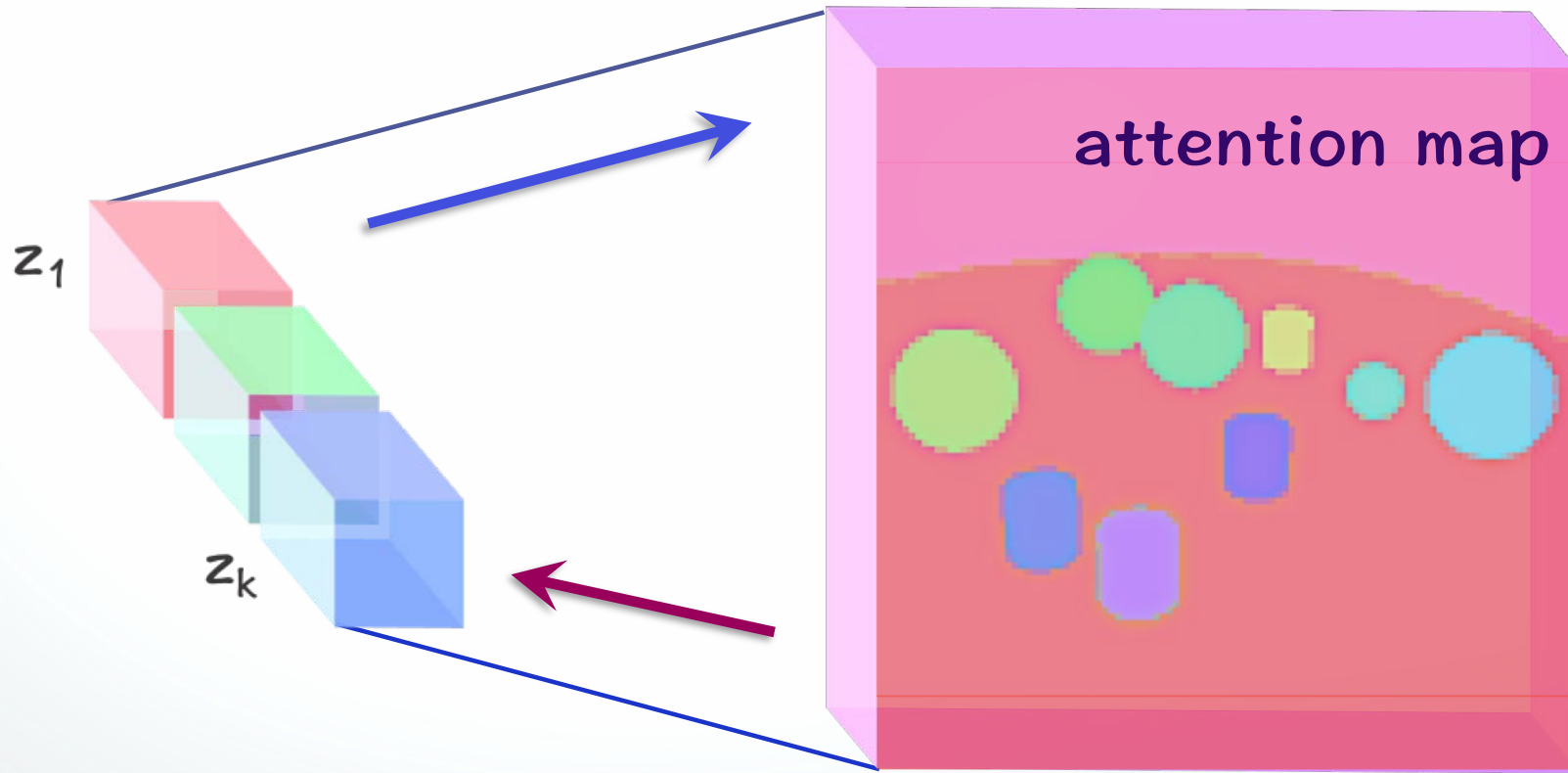Self-Attention

Latents

Image

Bipartite Attention

# From Transformer to GANformer



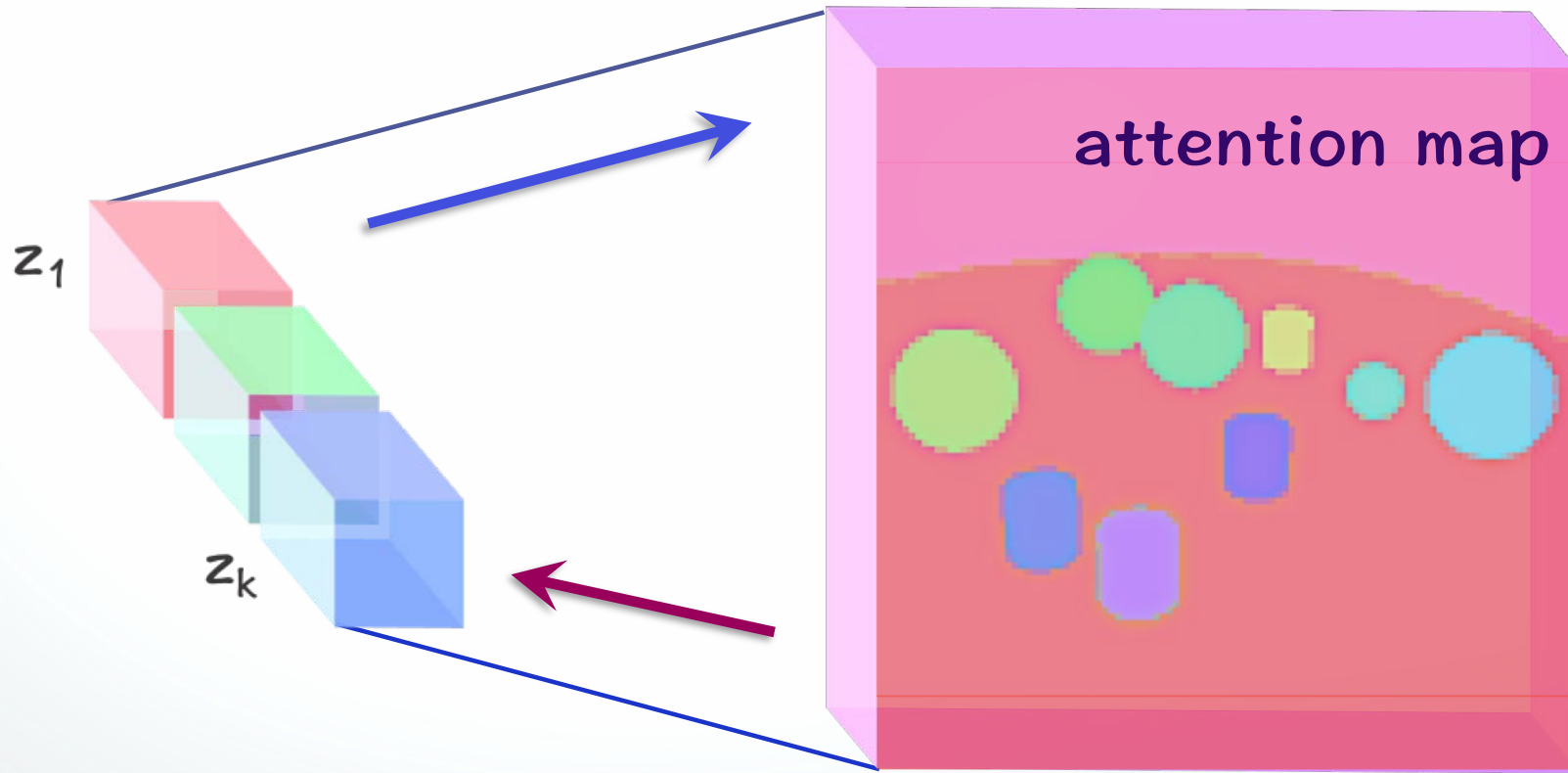**Compositional latent space with multiple variables**

# From Transformer to GANformer



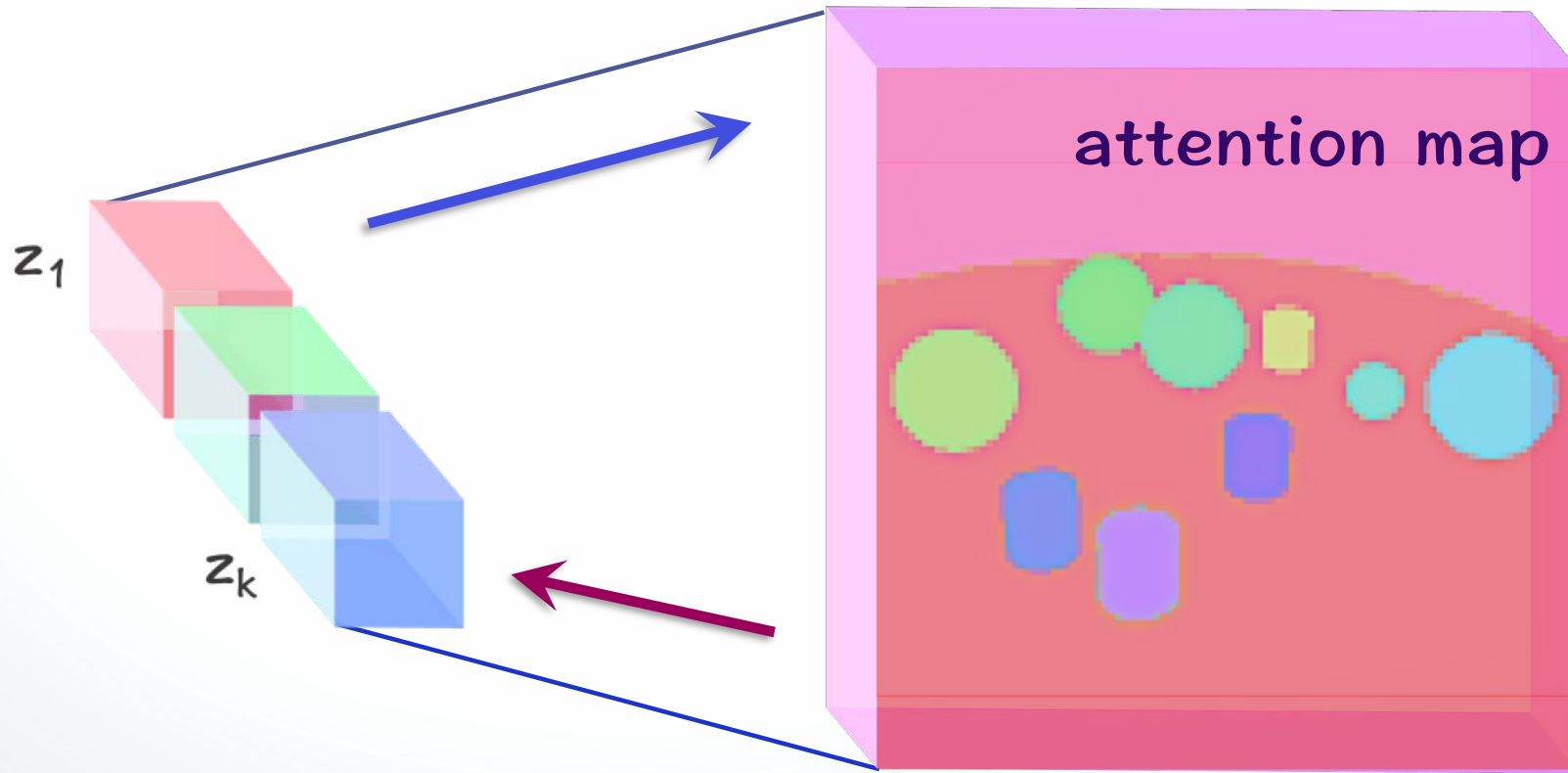**Compositional** **latent space** **with** **multiple** **variables**

# From Transformer to GANformer

attention map

$z_1$

$z_k$

**Bidirectional interaction** between the **latents** and the **image** enables **bottom-up & top-down processing**.
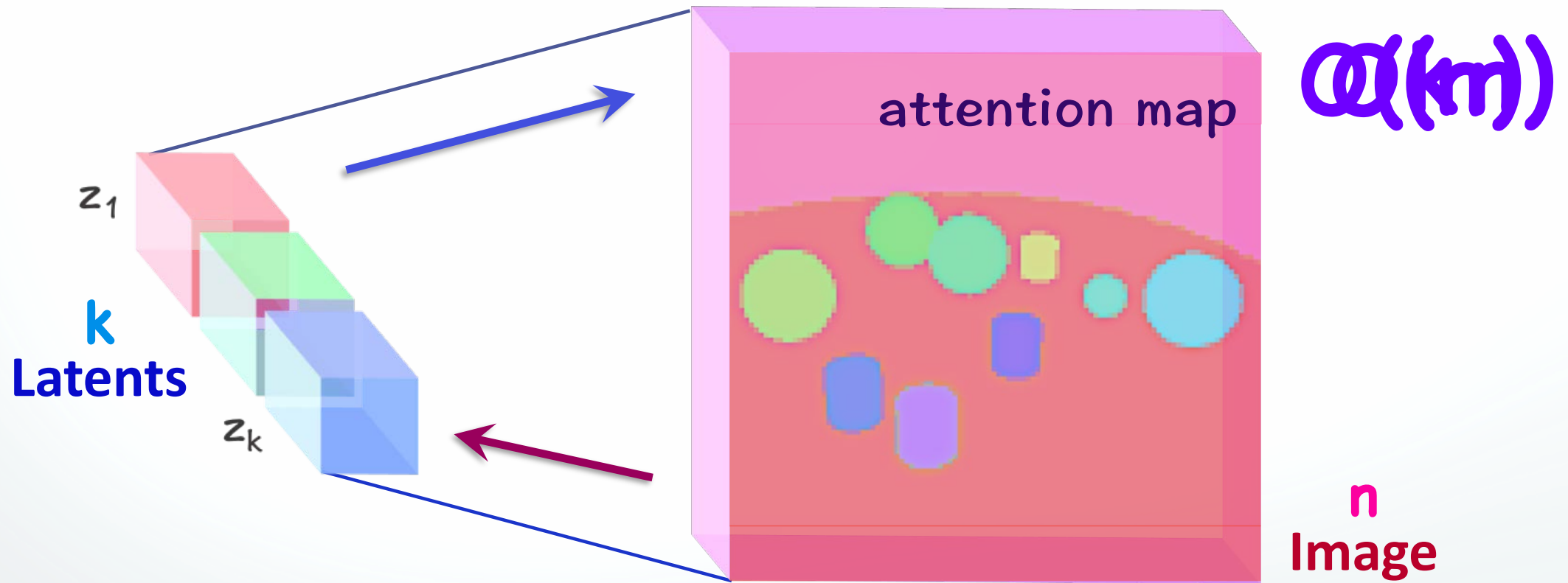
# From Transformer to GANformer



The **latents** guide the image synthesis **cooperatively**, **attending** and **modulating** different **objects or entities**.

# From Transformer to GANformer



$O(km)$

attention map

$z_1$

**k**
**Latents**

$z_k$

**n**
**Image**

**Has linear efficiency that scales to high-resolutions while capturing long-range dependencies.**

# GANformer's Attention Maps

**The latents attend to semantic entities and cooperatively generate compositional scenes.**

# GANformer's Attention Maps

**The latents attend to semantic entities and cooperatively generate compositional scenes.**

# GANformer's Attention Maps

The **latents attend** to **semantic entities** and **cooperatively** generate **compositional** scenes.

# GANformer's Attention Maps

The **latents attend** to **semantic entities** and **cooperatively** generate **compositional** scenes.
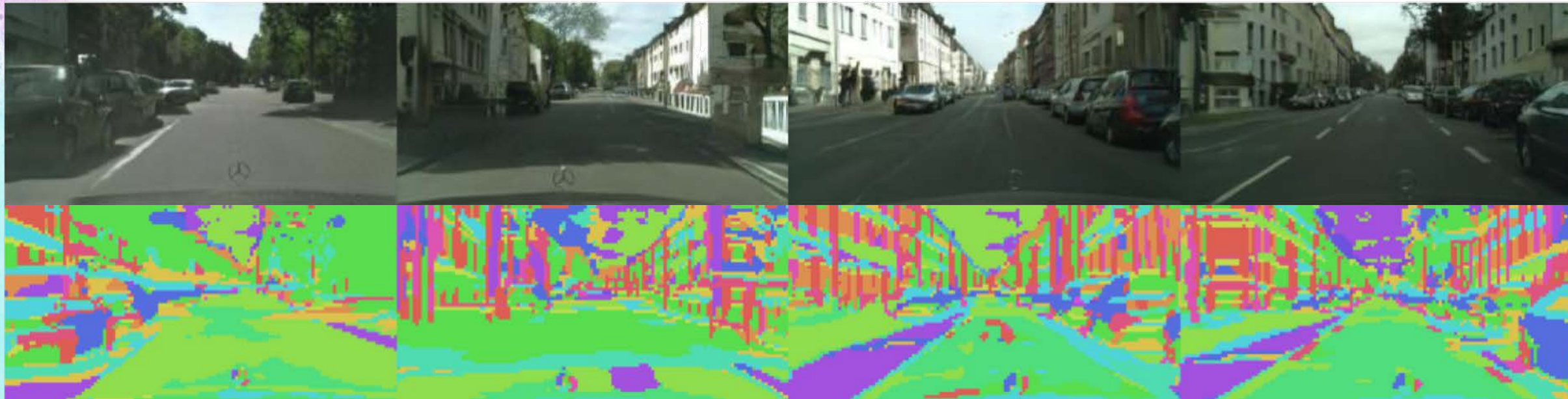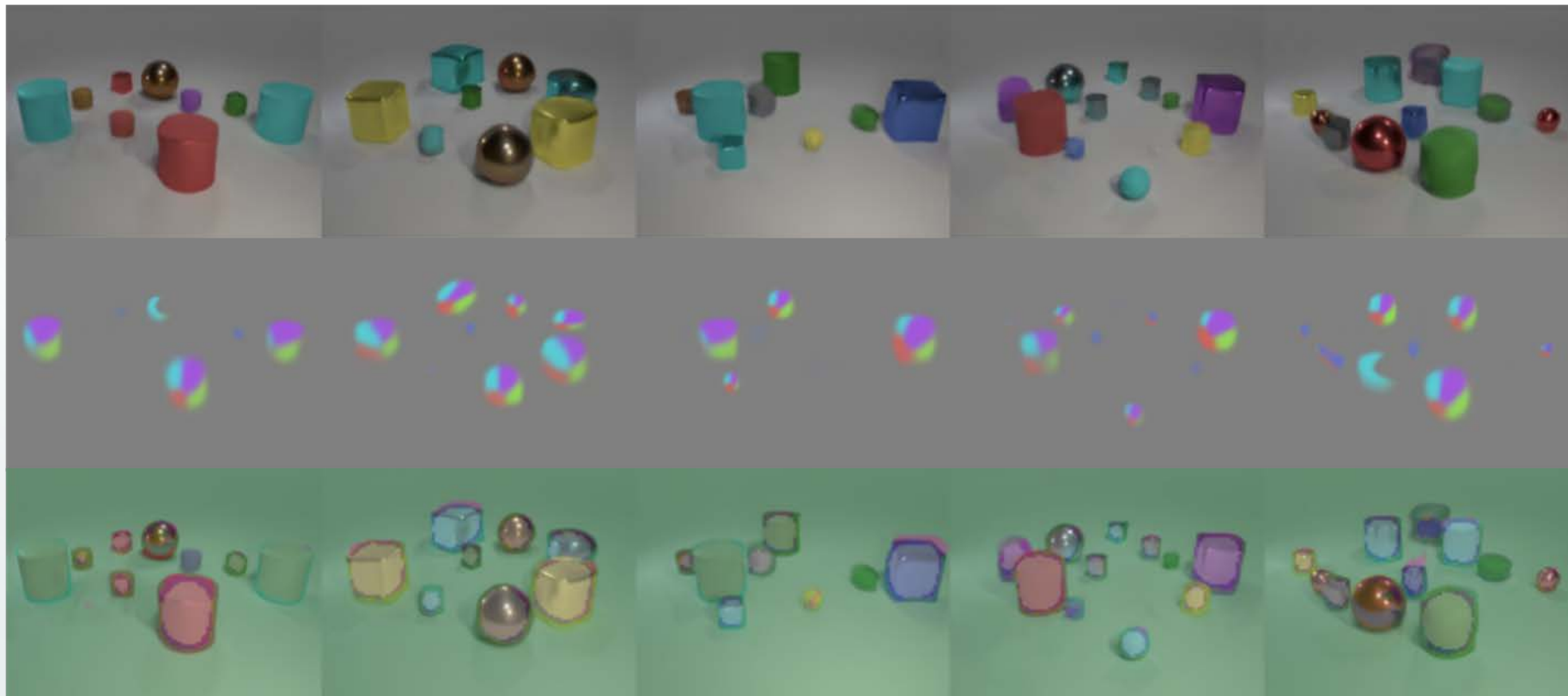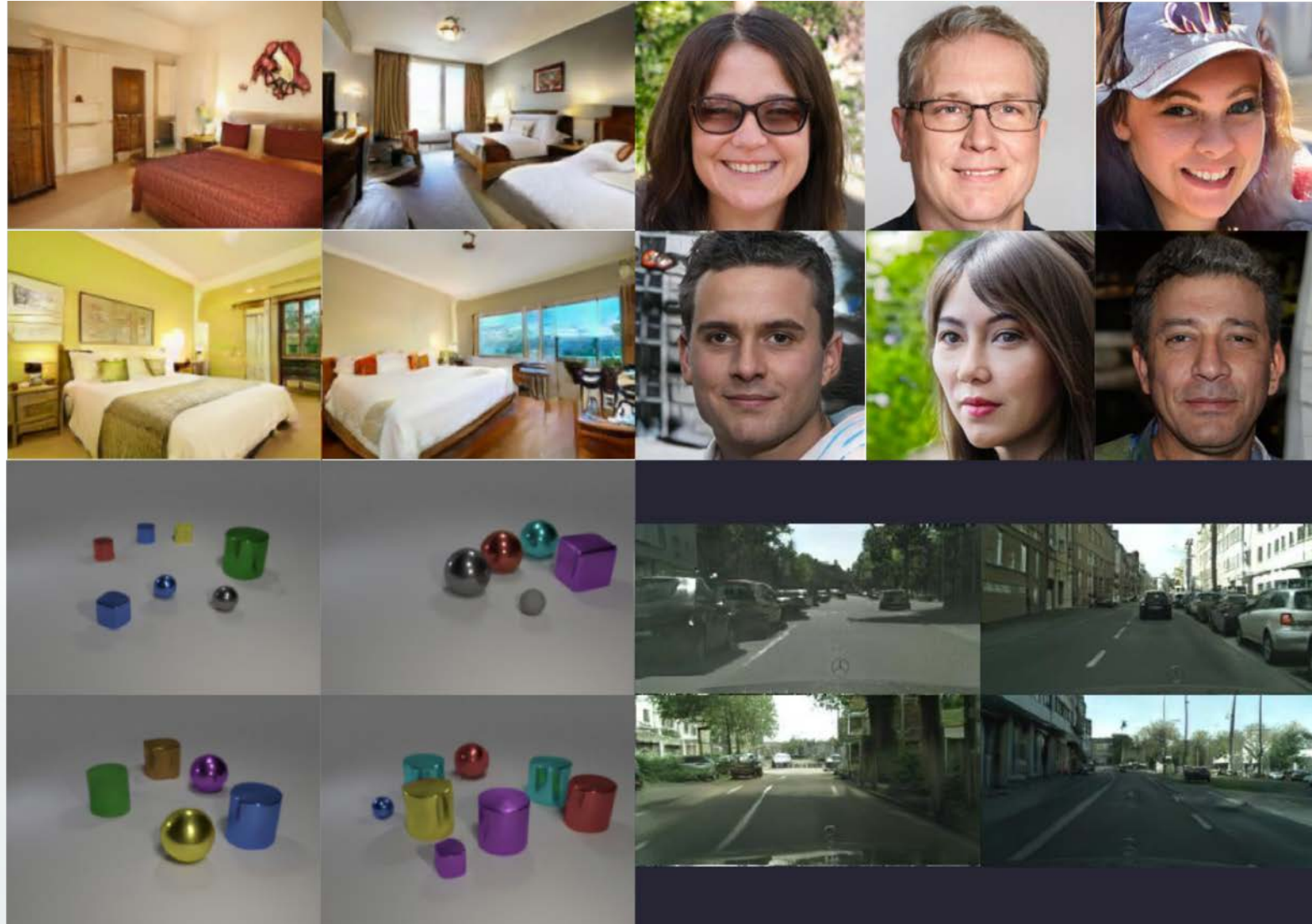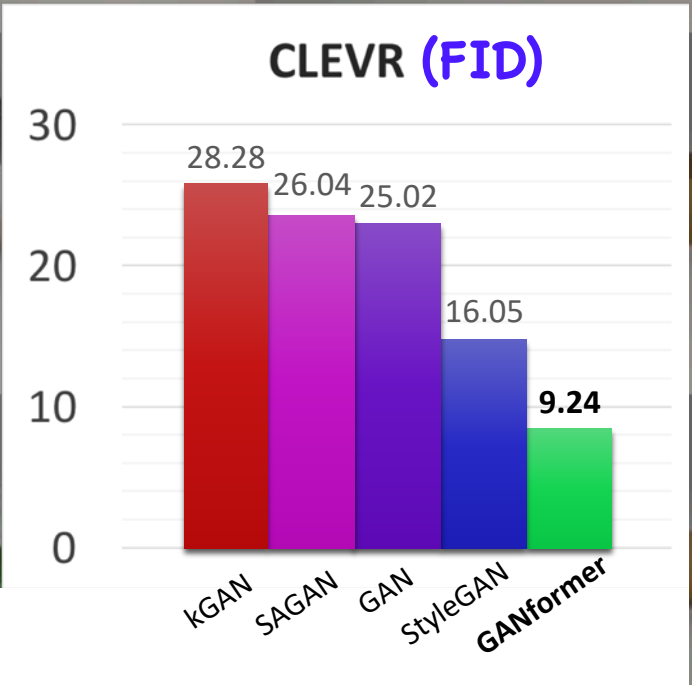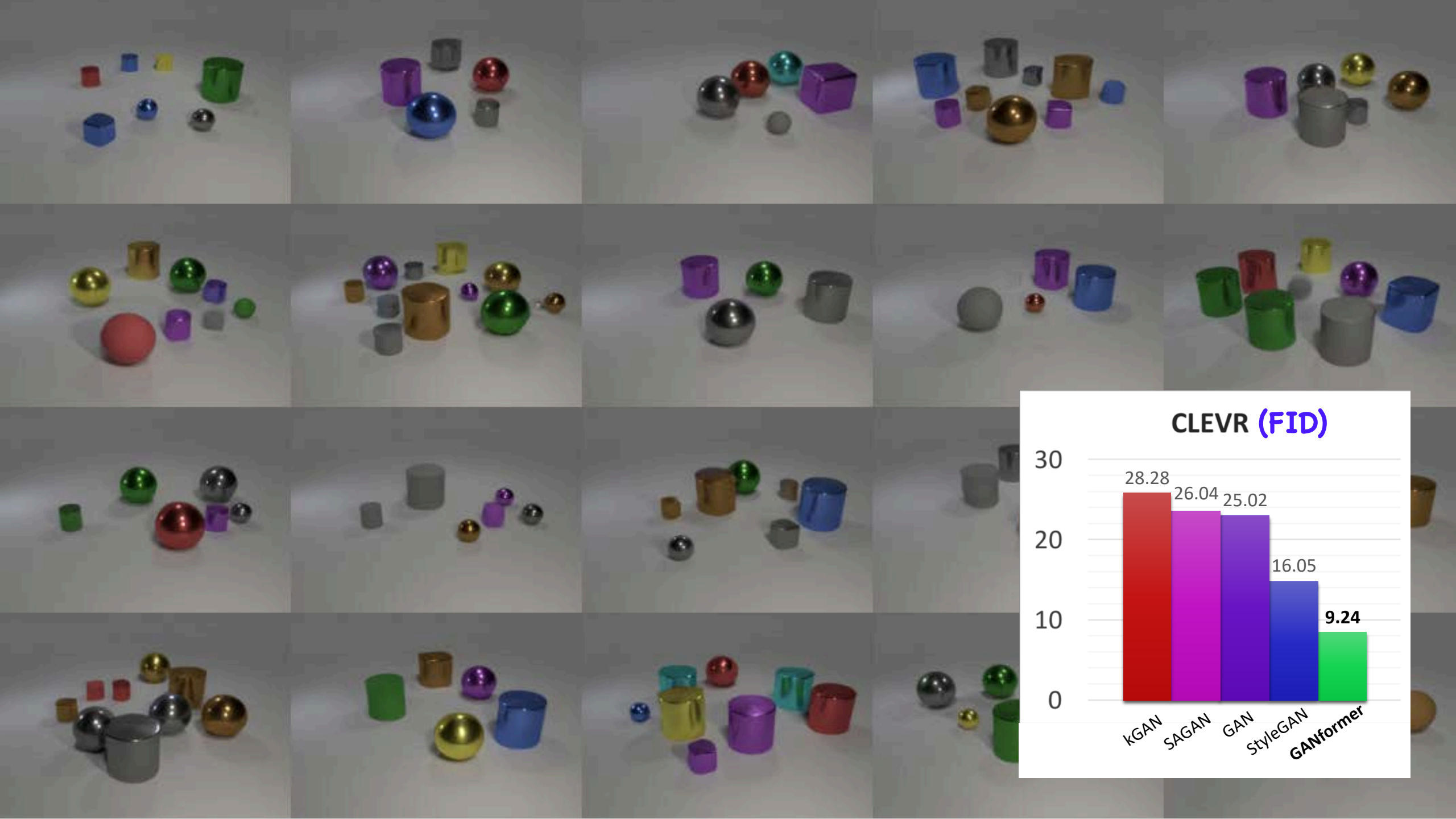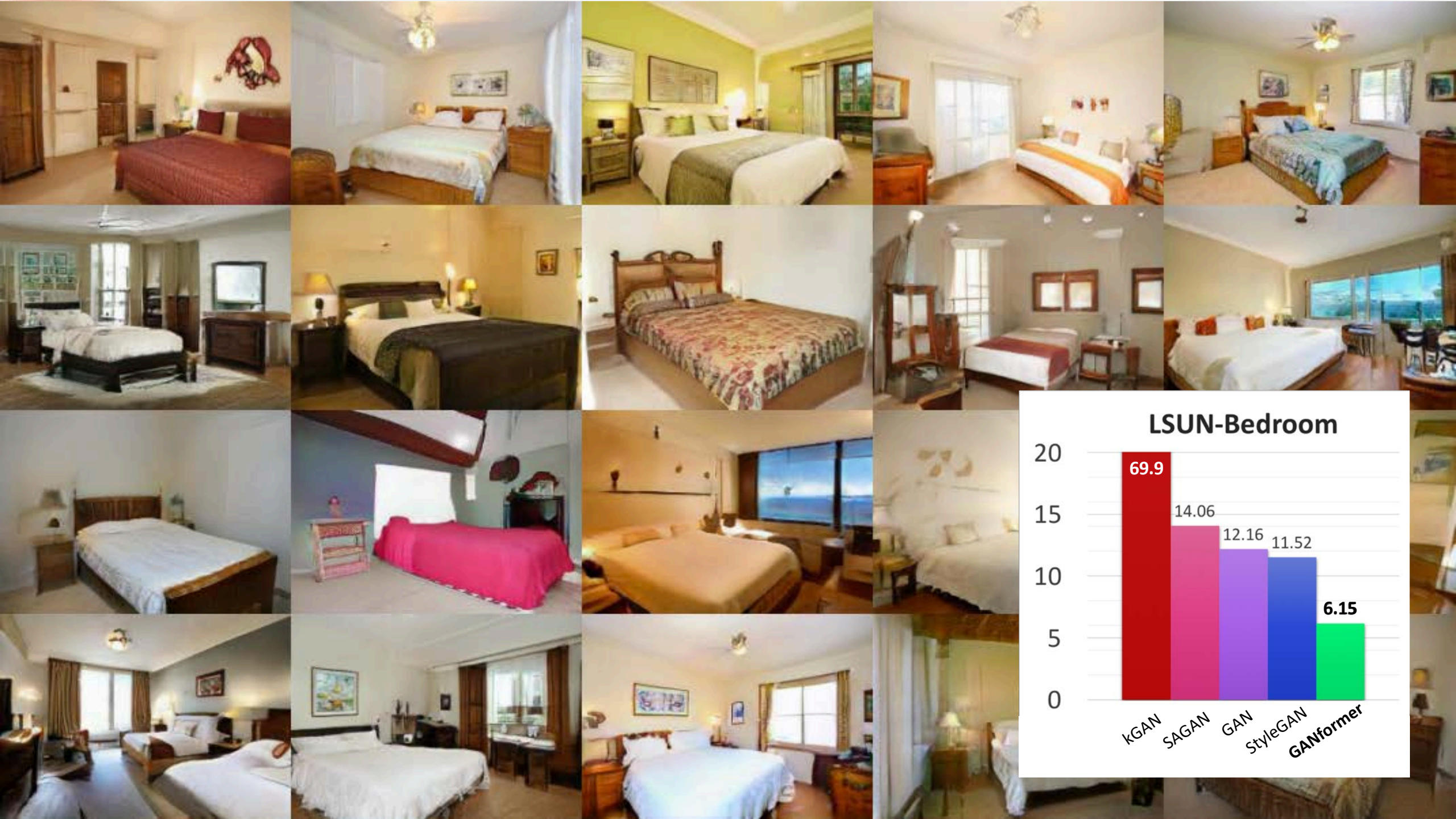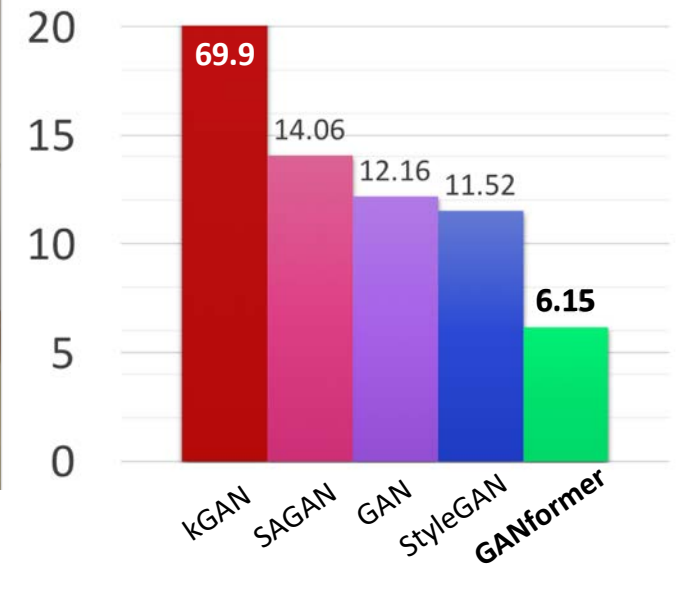
# Excels on highly-structured scenes

CLEVR (FID)

| kGAN | SAGAN | GAN | StyleGAN | GANformer |
|------|-------|-----|----------|-----------|
| 28.28 | 26.04 | 25.02 | 16.05 | **9.24** |

LSUN-Bedroom

| kGAN | SAGAN | GAN | StyleGAN | GANformer |
|------|-------|-----|----------|-----------|
| 69.9 | 14.06 | 12.16 | 11.52 | 6.15 |

Cityscapes

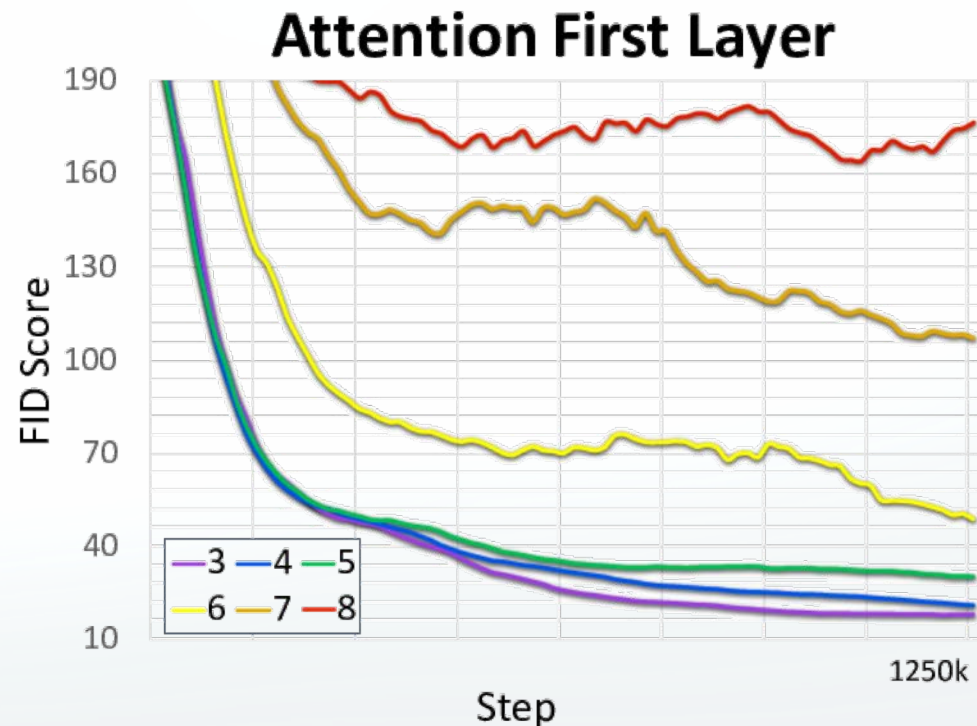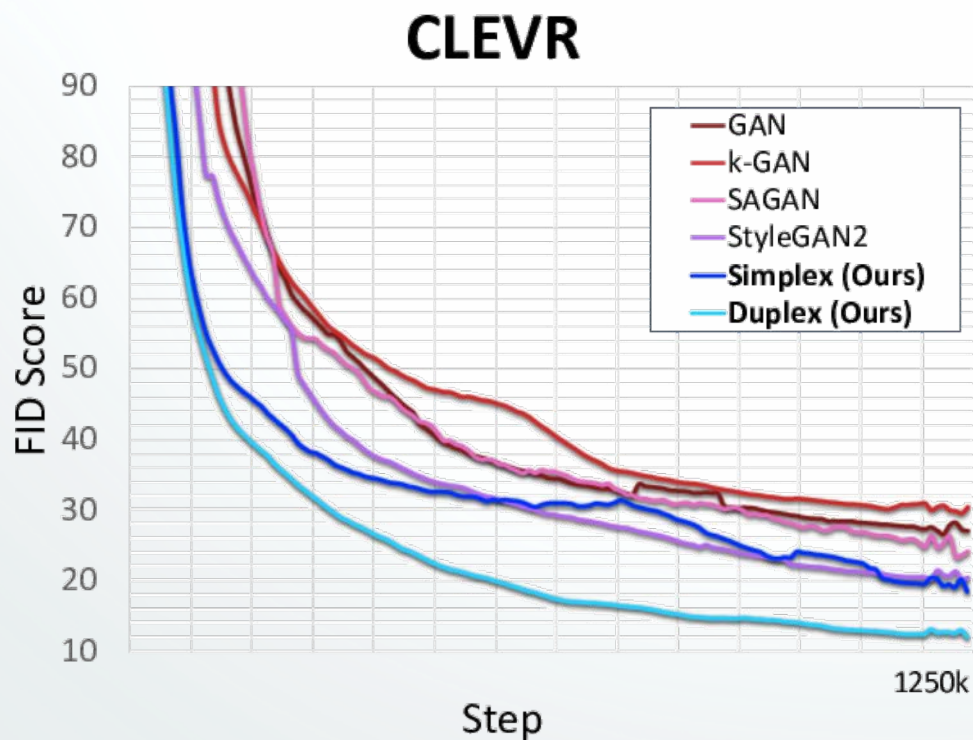| kGAN | SAGAN | GAN | StyleGAN | **GANformer** |
|------|-------|-----|----------|--------------|
| 51.08 | 12.81 | 11.57 | 8.35 | 5.23 |

# GANformer Model Analysis

**The model learns faster and enjoys higher data-efficiency.**

**Latents are more disentangled; Images are more diverse.**

`|ıₗₗₗ| State of the Art` `Image Generation on CLEVR`   `|ıₗₗₗ| State of the Art` `Image Generation on Cityscapes`

`|ıₗₗₗ| State of the Art` `Image Generation on LSUN Bedroom 256 x 256 (FID-10k-training-steps metric)`

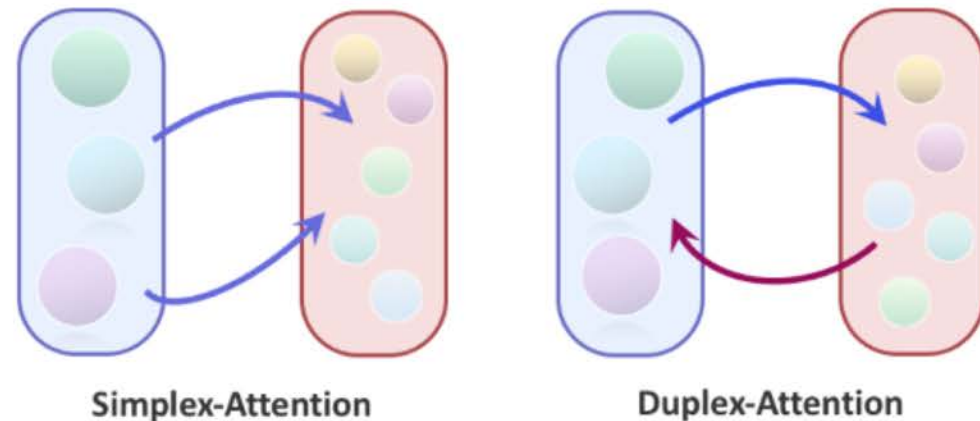`python 3.7`  `tensorflow 1.14`  `cudnn 10.0`  `license MIT`

# GANformer: Generative Adversarial Transformers

**github.com / dorarad / gansformer**

**The Bipartite Transformer**

Simplex-Attention                    Duplex-Attention

github.com / dorarad / gansformer

Thank you! ☺