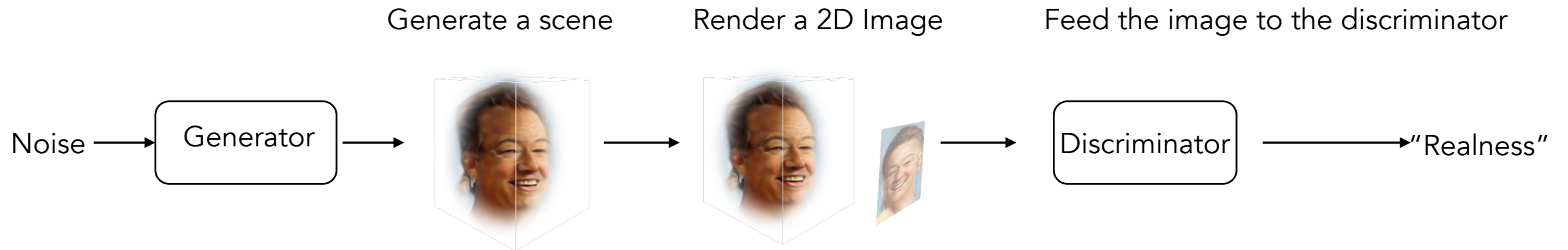


Training a 3D-Aware GAN

Training Steps

1. Generate a representation of a *scene*
2. Render the scene from a random camera pose
3. Feed the image to a 2D discriminator
4. Backpropagate through the discriminator and differentiable rendering



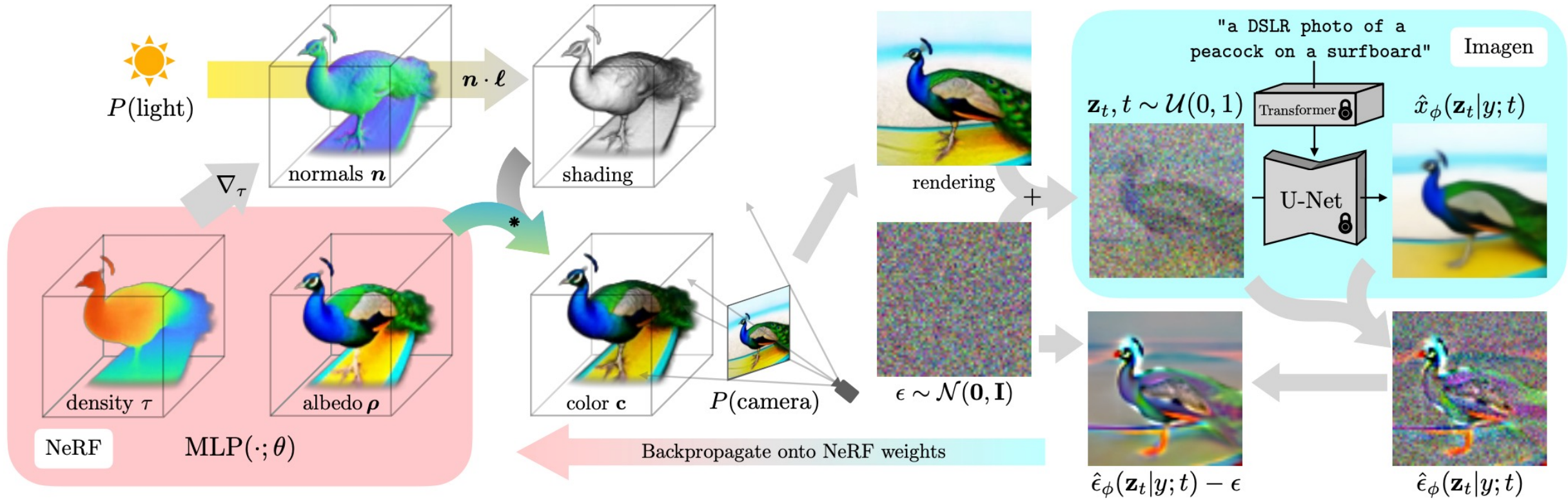
Slide credit: Eric Chan

Text-based Editing

a DSLR photo of a squirrel
wearing a purple hoodie
reading a book



Text-based Editing



FOR loop

Step 1. Render a view using existing NeRF

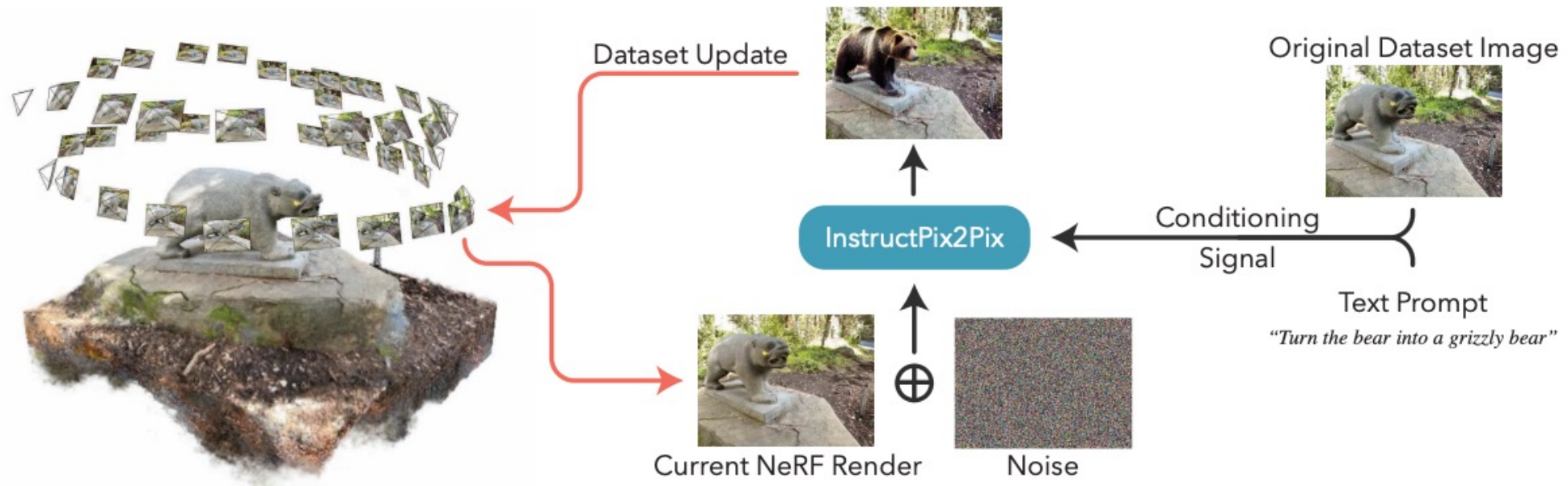
Step 2. Add noise and denoise using a pre-trained Stable Diffusion model

Step 3. Update NeRF parameters with the gradient (difference between added and predicted noises)

Instruct NeRF2NeRF



Instruct NeRF2NeRF



FOR loop

Step 1. Render a view using existing NeRF

Step 2. Use InstructPix2Pix to produce output images

Step 3. Update NeRF parameters with the generated result from Step 2

InstructPix2pix: image-conditional diffusion model (<https://www.timothybrooks.com/instruct-pix2pix/>)



Video Synthesis and Editing

Jun-Yan Zhu

16-726 Learning-based Image Synthesis, Spring 2023

Image Editing and Synthesis

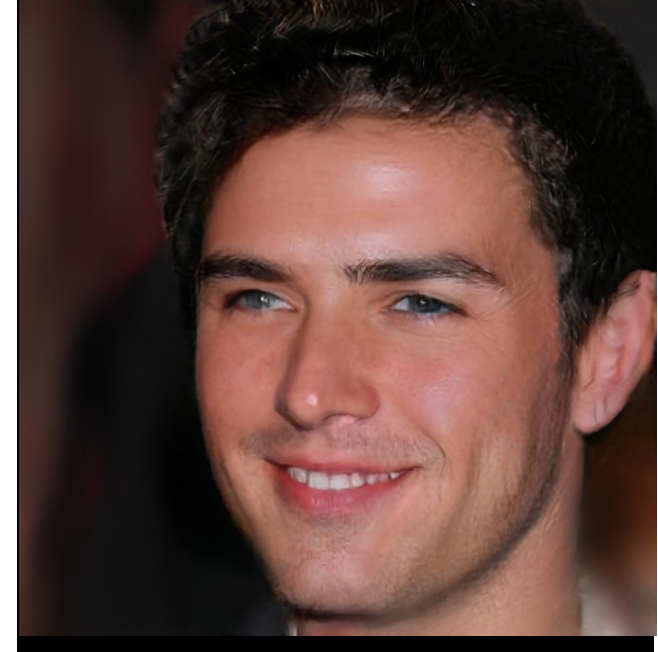


Image Editing and Synthesis



Image Editing and Synthesis



“Emma Stone”

“Mohawk hairstyle”

“Without makeup”

“Cute cat”

“Lion”

“Gothic church”

Images vs. Videos

- What are the major differences?
 - 2D vs. 3D
 - Spatial vs. spatial + temporal
- Differences between 3D vs. videos?
 - (x, y, z) vs. (x, y, t)
- Why are videos much more challenging?
 - Understanding motions/actions/intensions.
 - Long-term dependence.
 - Hard to annotate.
 - Computationally-expensive (e.g., memory, training time)

Why not apply it to video?

Someone gave you an image synthesis model

asked you to build a video synthesis application

How to Generalize to Videos

- Idea 1: Frame-by-Frame

Frame-by-Frame Result (pix2pixHD)



Frame-by-Frame Result (CycleGAN)



How to Generalize to Videos

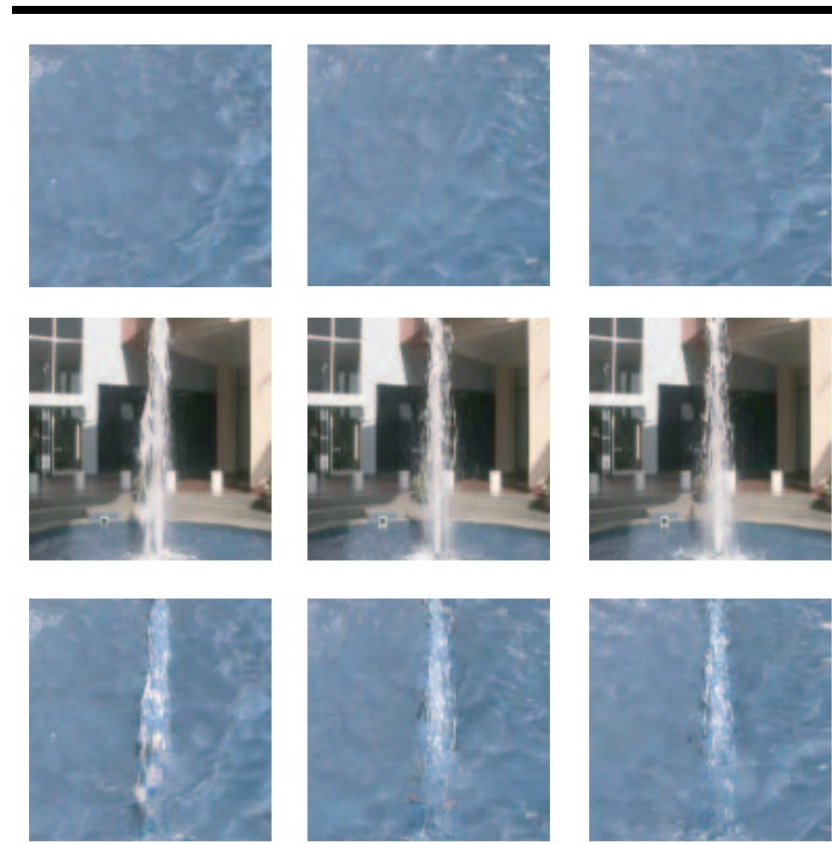
- Idea 1: Frame-by-Frame
 - Temporal inconsistency (flicking, color drift)
- Idea 2: Video as 3D data (height x width x time)

Case Study: 3D Poisson blending

Spatial-temporal Constraints

Background

Time



$$\begin{aligned}
 F(\nabla I, G) &= \|\nabla I - G\|^2 \\
 &= \left(\frac{\partial I}{\partial x} - G_x\right)^2 + \left(\frac{\partial I}{\partial y} - G_y\right)^2 + \left(\frac{\partial I}{\partial t} - G_t\right)^2
 \end{aligned}$$

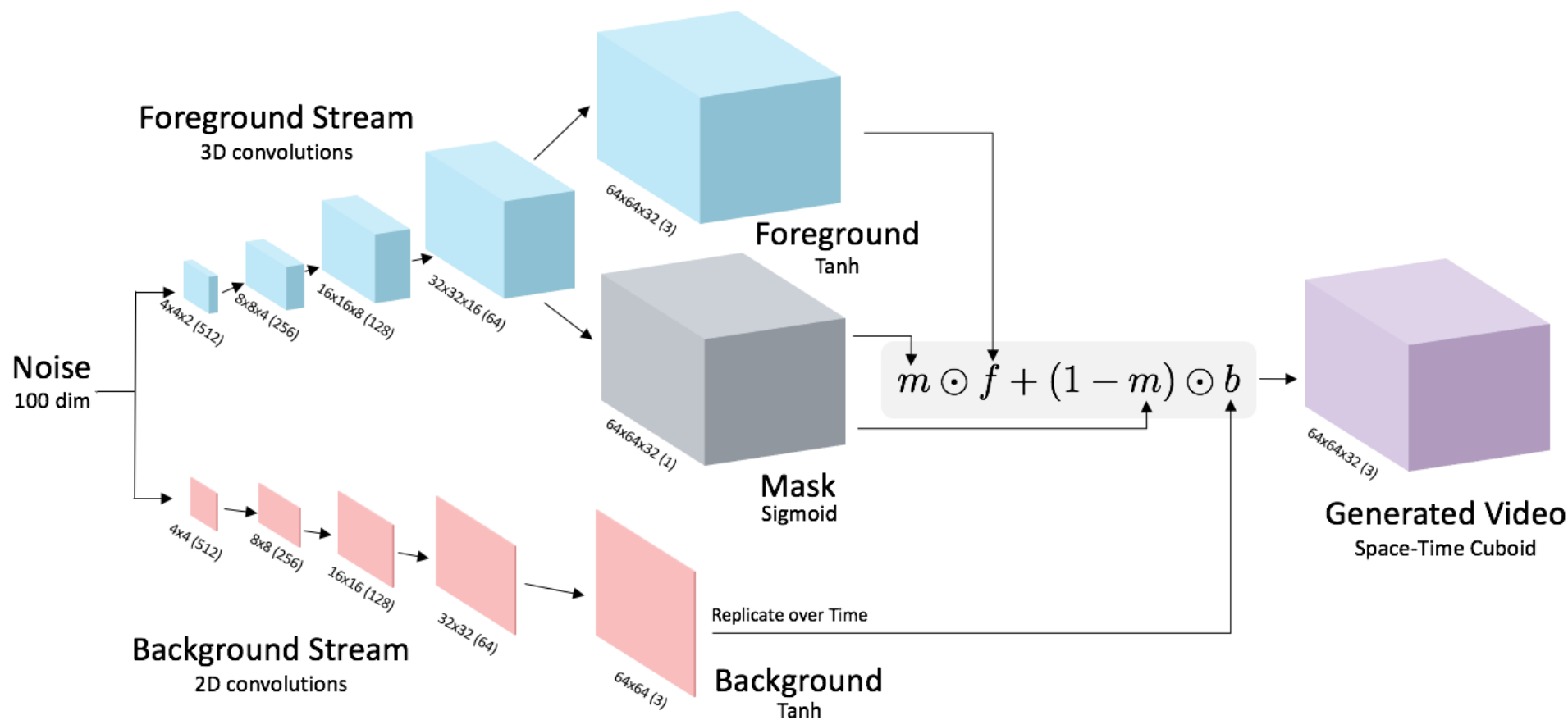
Foreground

Output Image

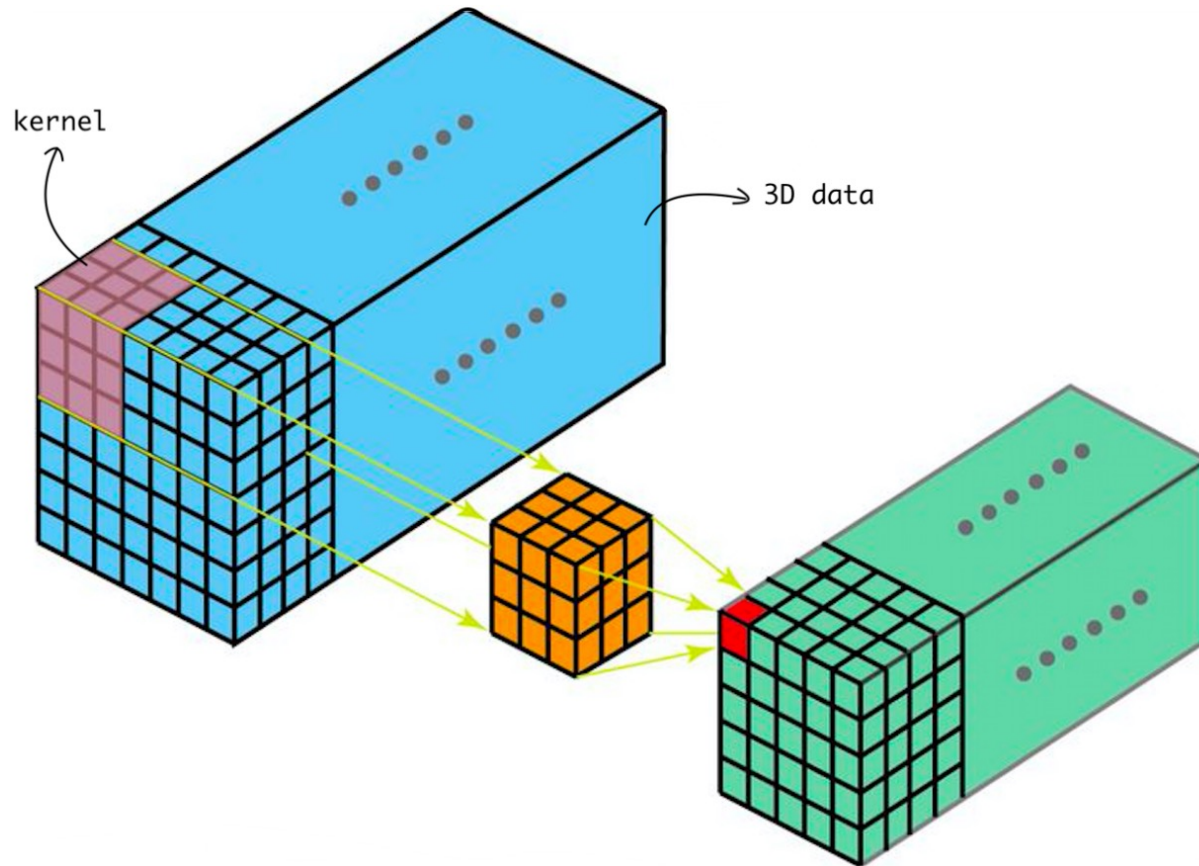
Guidance Gradient

Output

Case Study: Video GANs



Recap: 3D Conv



Easy to implement:

- Replace 2D by 3D in your code

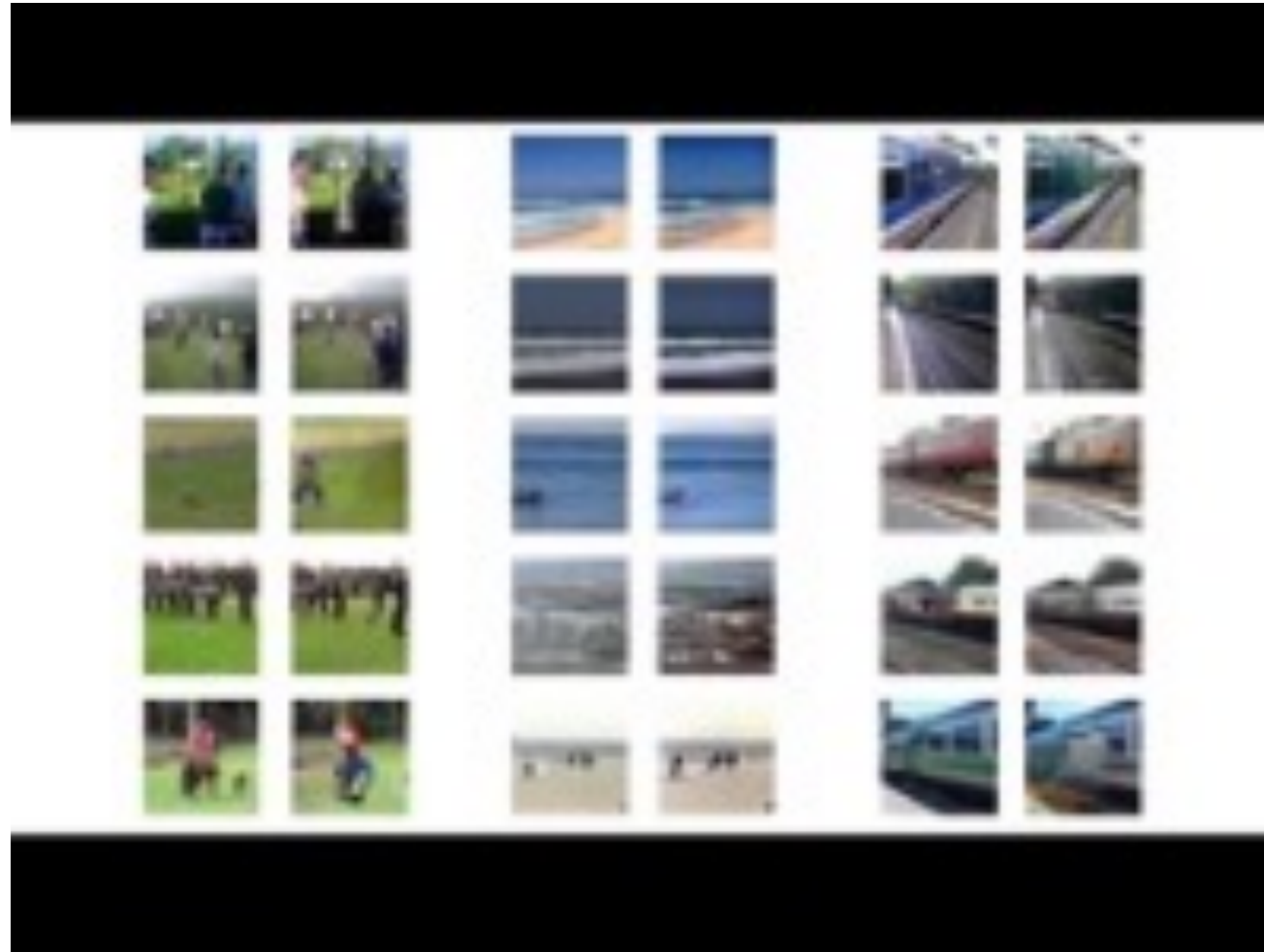
e.g., Conv2D -> Conv3D

ConvTranspose2d -> ConvTranspose3d

MaxPool2d -> MaxPool3d

```
CLASS torch.nn.Conv3d(in_channels, out_channels, kernel_size, stride=1, padding=0,  
dilation=1, groups=1, bias=True, padding_mode='zeros', device=None, dtype=None) [SOURCE]
```

Case Study: Video GANs

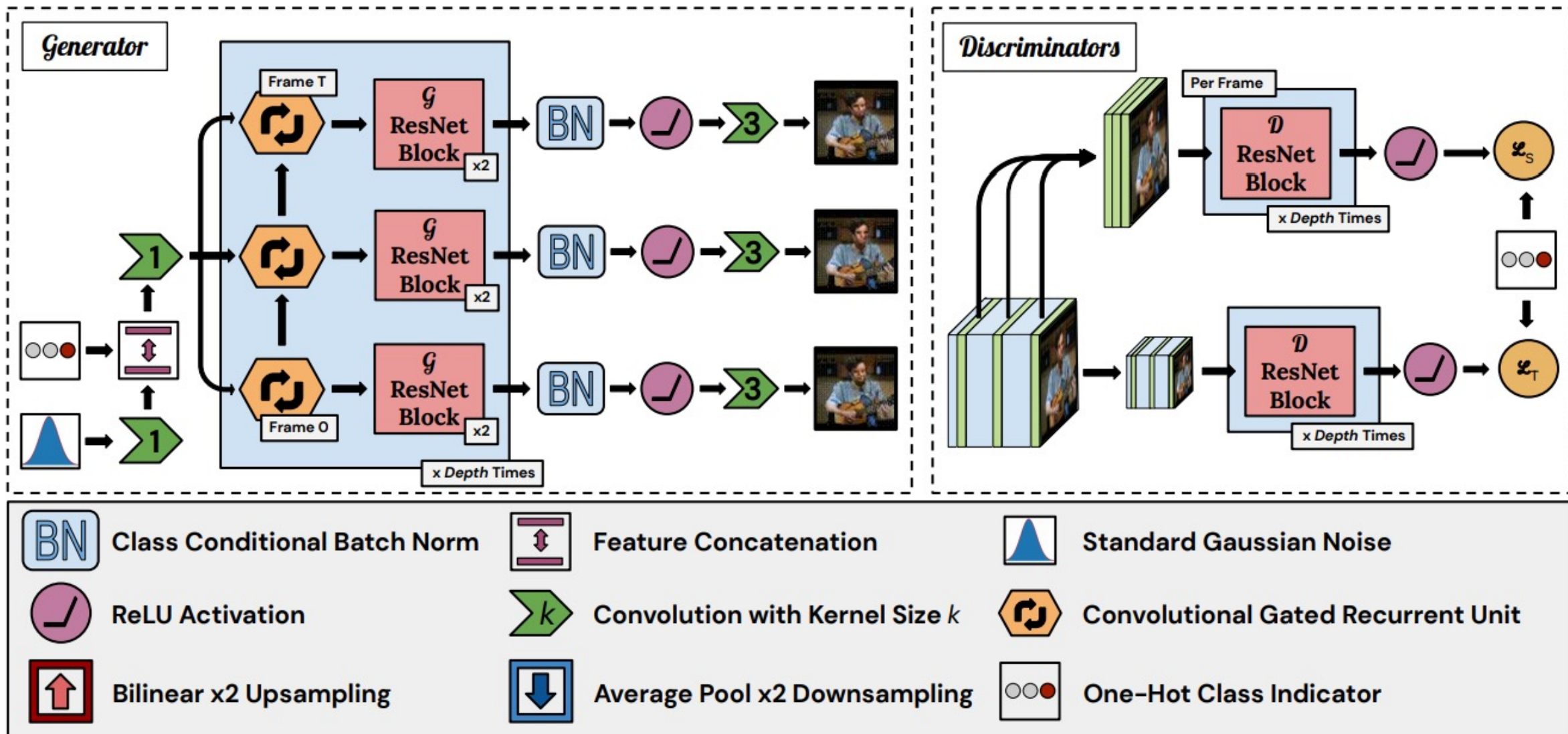


Case Study: DVDGANs



BigGAN-based generator + Spatial Discriminator + Temporal Discriminator

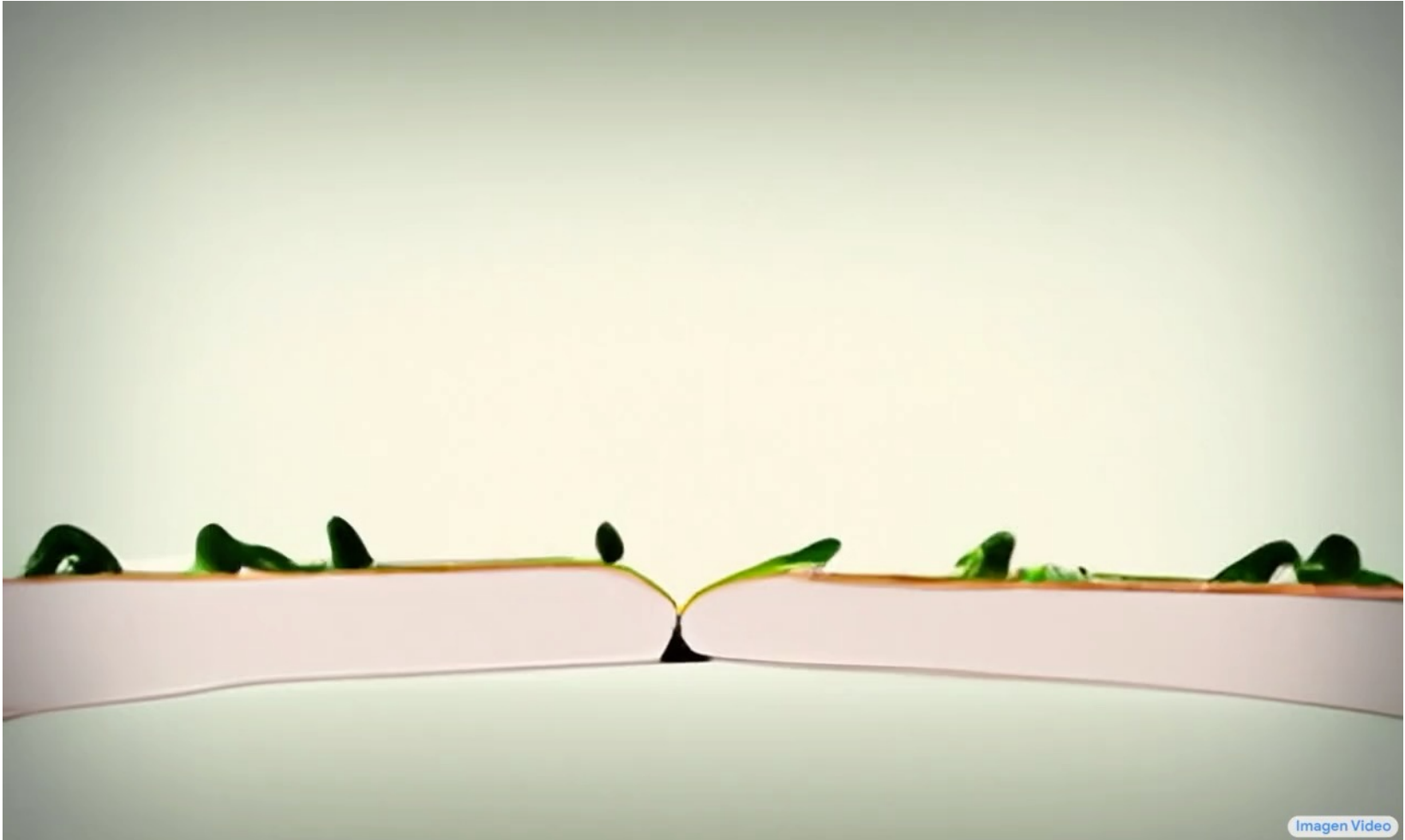
Case Study: DVDGANs



Case Study: DVDGANs



Case Study: Imagen Video



Case Study: Imagen Video



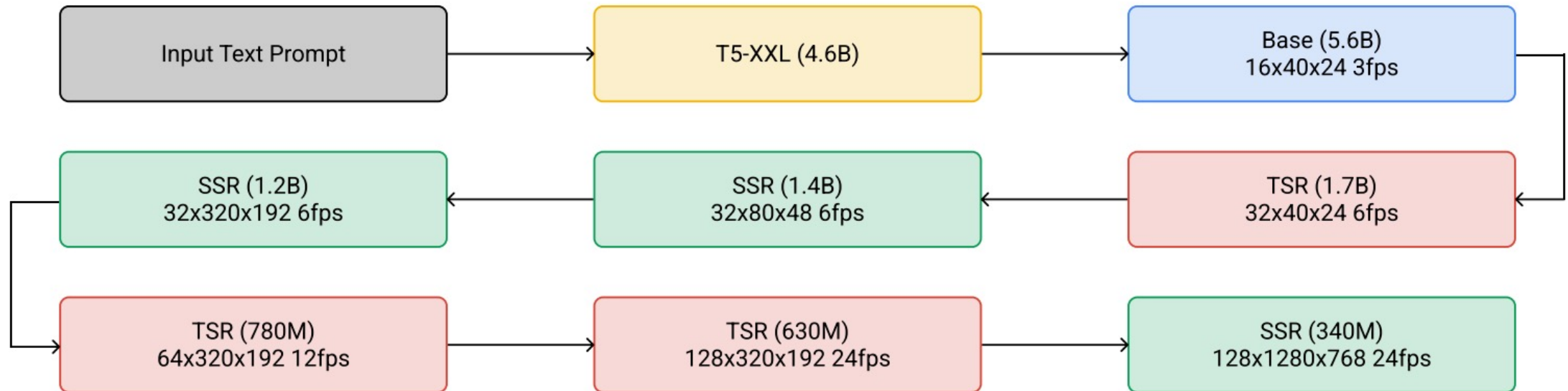
An astronaut riding a horse

Case Study: Imagen Video

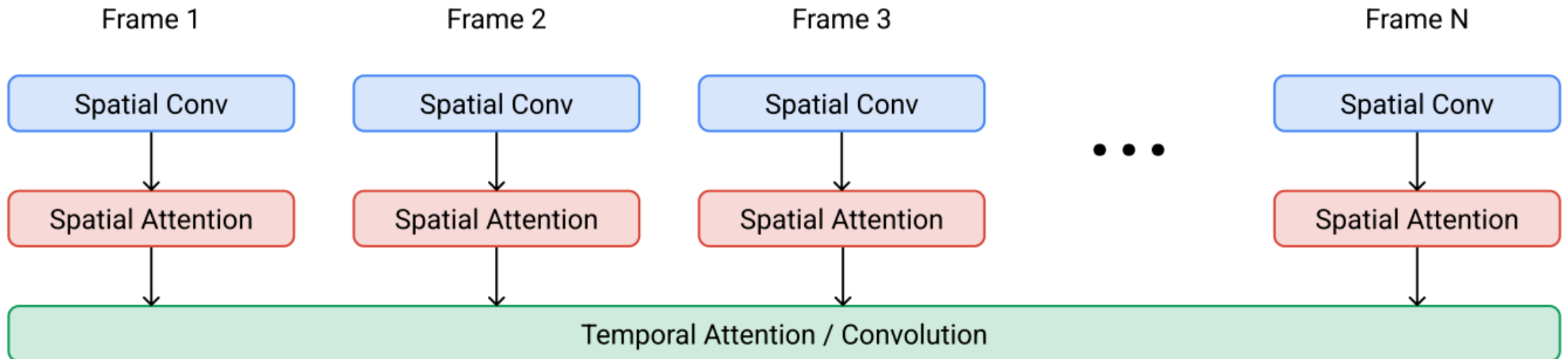


Wooden figure surfing on a surfboard in space

Case Study: Imagen Video



Case Study: Imagen Video



How to Generalize to Videos

- Idea 1: Frame-by-Frame
 - Temporal inconsistency (flicking, color drift)
- Idea 2: Video as 3D data (height x width x time)
 - memory-intensive and time-consuming
 - only work for a short video at low resolution
- Idea 3: recurrent (autoregressive) synthesis
 - Generate 1st frame, generate 2nd frame based on 1st one, ...
 - Using optical flow (optional): warp 1st frame to 2nd frame.

Text Synthesis

- [Shannon, '48] proposed a way to generate English-looking text using N-grams:
 - Assume a generalized Markov model
 - Use a large text to compute prob. distributions of each letter given N-1 previous letters
 - Starting from a seed repeatedly sample this Markov chain to generate new letters
 - Also works for whole words

WE NEED TO EAT CAKE

Case Study: Video Textures



Still image



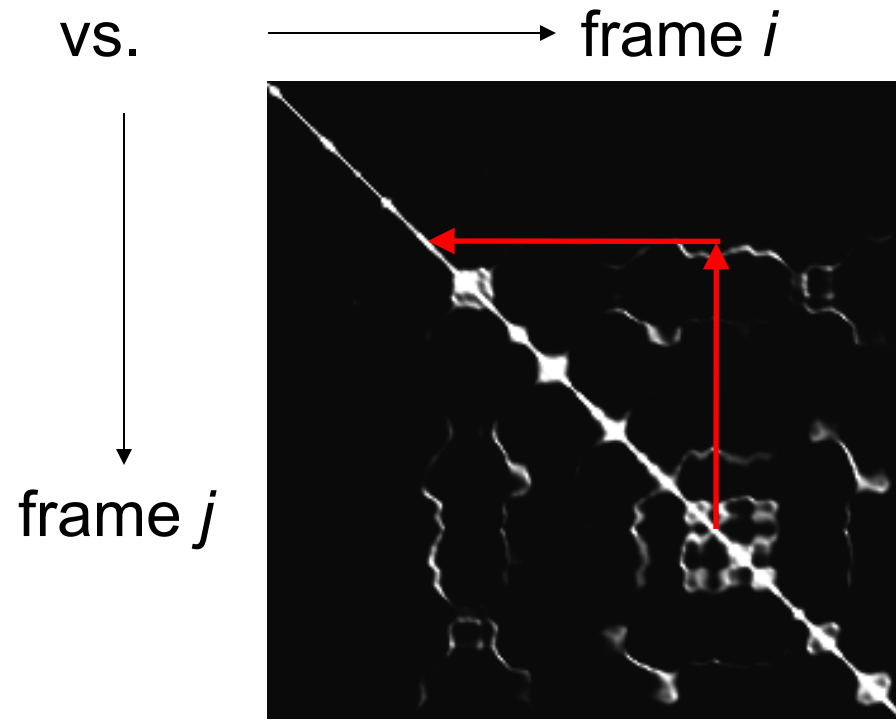
Loopy video

Case Study: Video Textures



Finding good transitions

- Compute L₂ distance $D_{i,j}$ between all frames

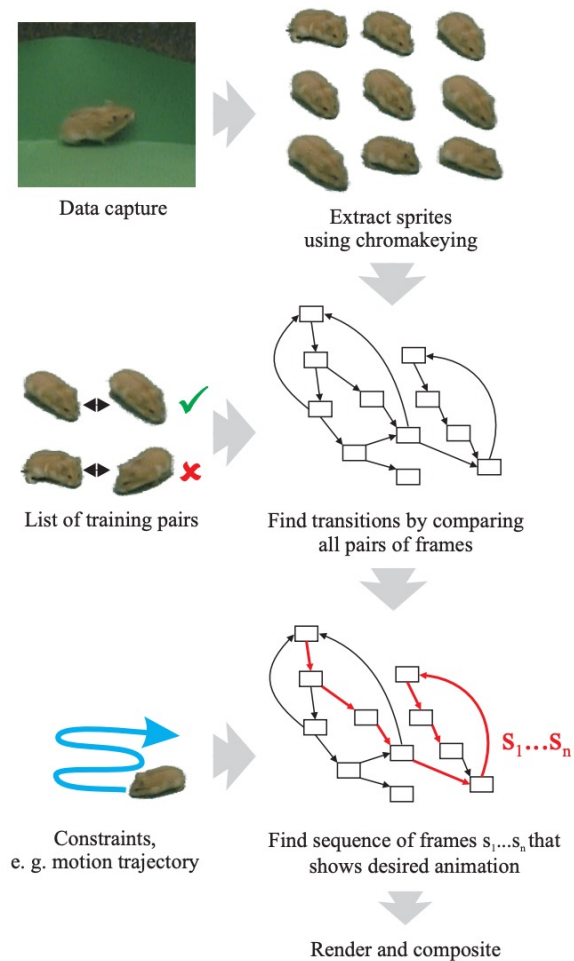


Similar frames make good transitions

Case Study: Video Textures



Case Study: Controlled Animation of Video Sprites



Case Study: Video-to-Video Translation



T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, G. Liu, A. Tao, J. Kautz, B. Catanzaro,
“Video-to-Video Synthesis,” NeurIPS 2018.

<https://github.com/NVIDIA/vid2vid>

Previous Work: Frame-by-Frame Result

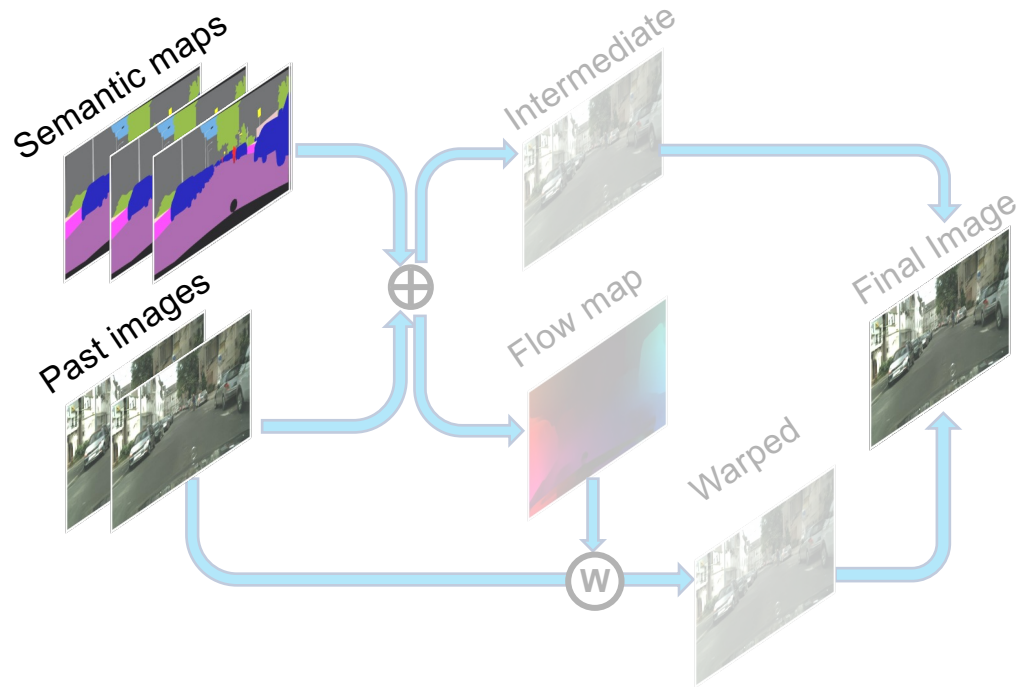


vid2vid

- Sequential generator
- Multi-scale temporal discriminator
- Spatio-temporal progressive training procedure

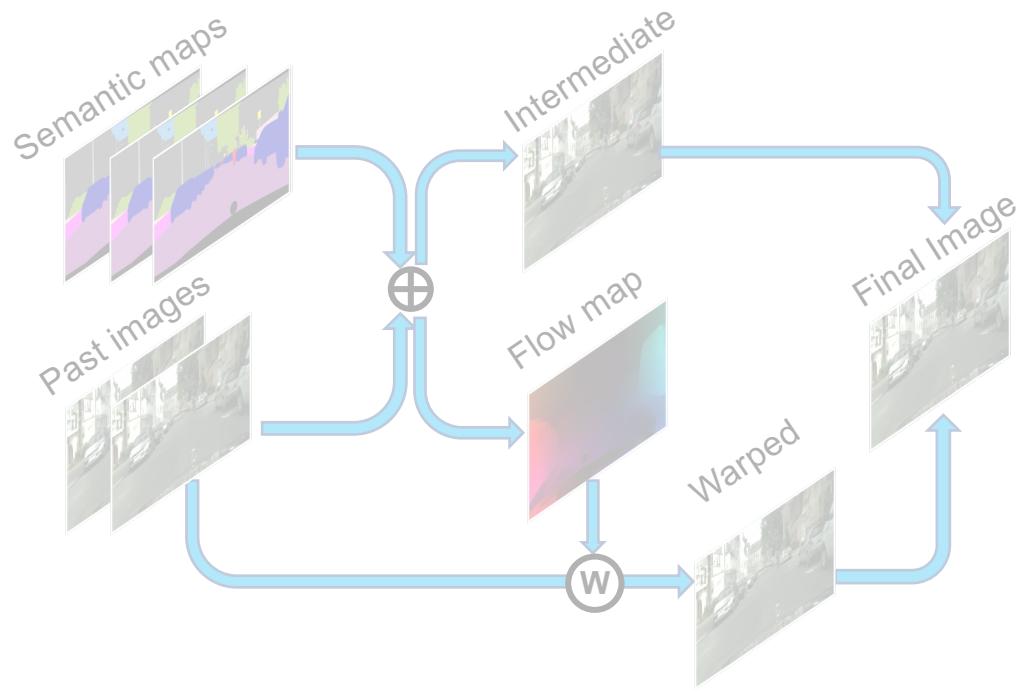
vid2vid

Sequential Generator



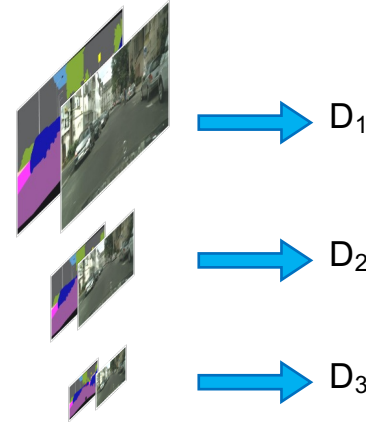
vid2vid

Sequential Generator

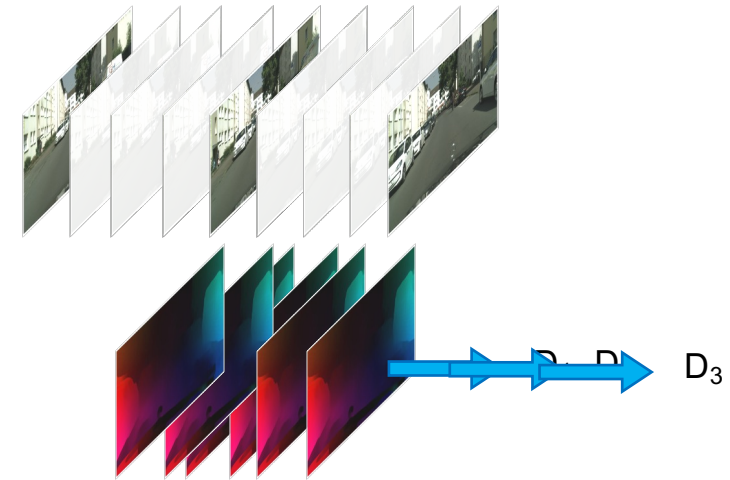


Multi-scale Discriminators

Image Discriminator



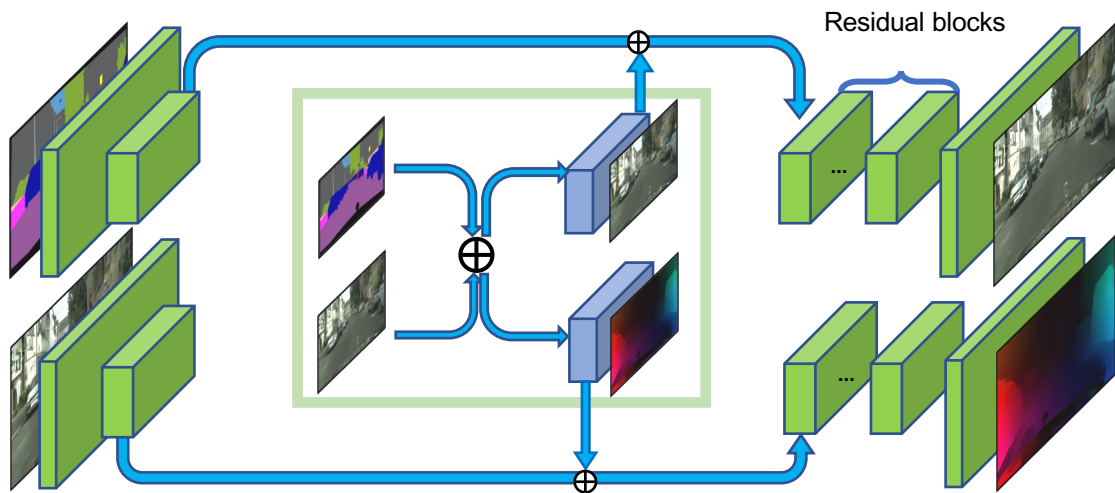
Video Discriminator



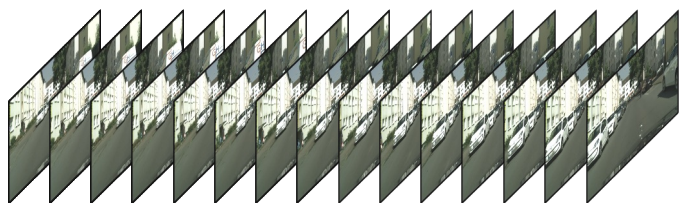
vid2vid

Spatio-temporally Progressive Training

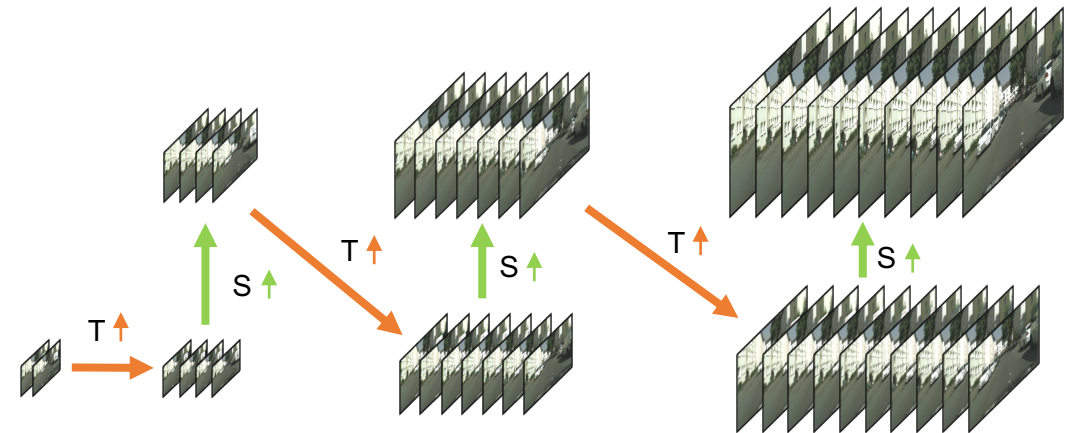
Spatially progressive



Temporally progressive



Alternating training



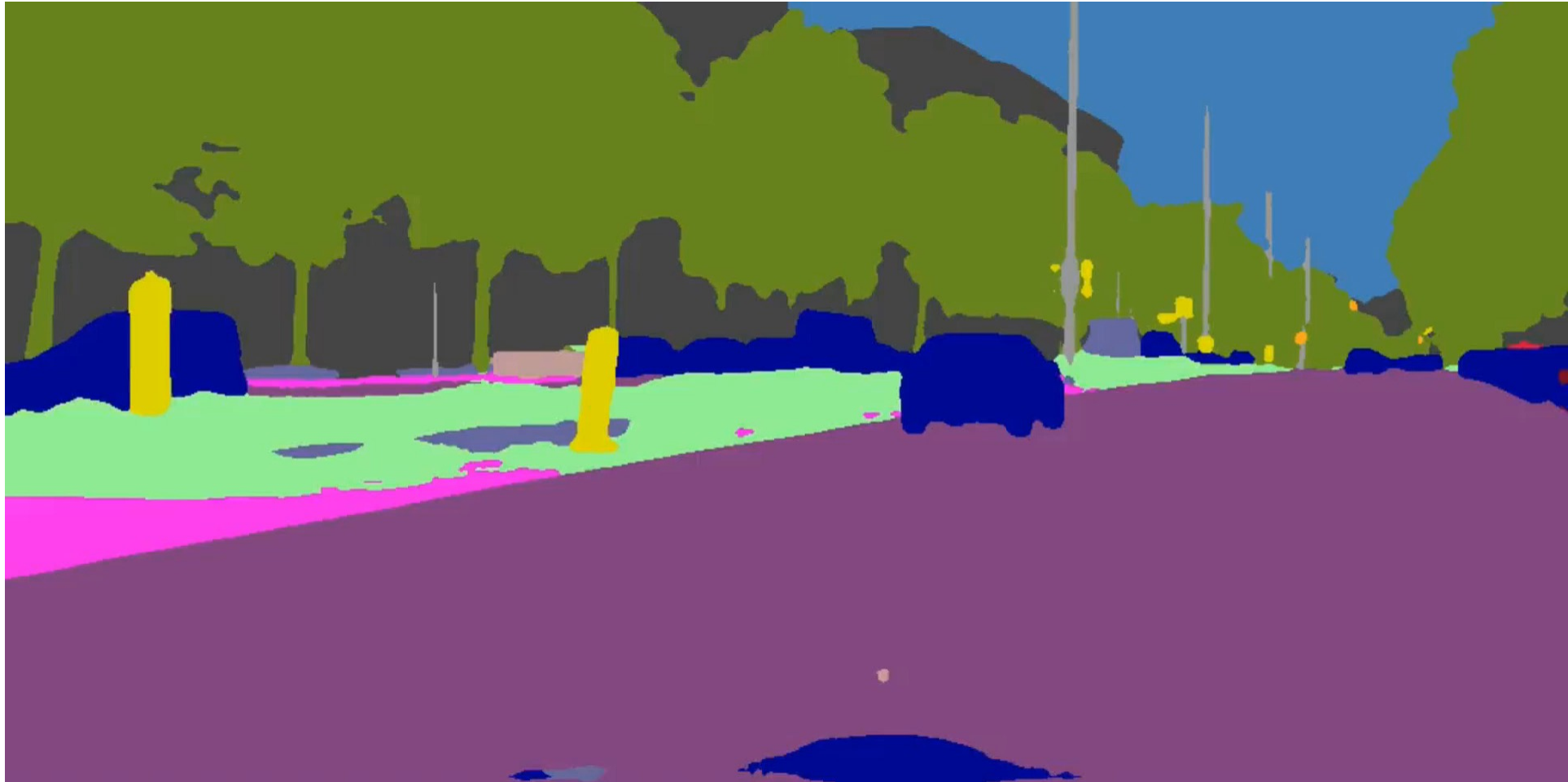
vid2vid Results

- Semantic → Street view scenes
- Edges → Human faces
- Poses → Human bodies

vid2vid Results

- Semantic → Street view scenes

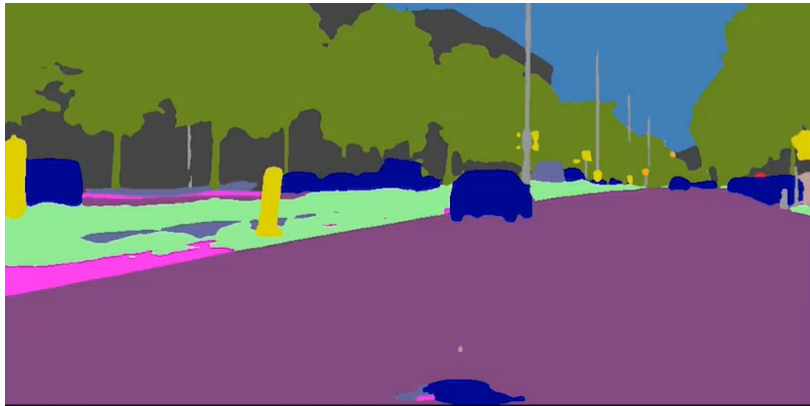
Street View: Cityscapes



Street View: Cityscapes



Street View: Cityscapes



Labels



pix2pixHD

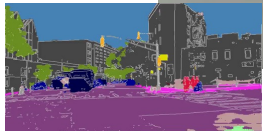


COVST



Ours

Street View: Boston



Street View: NYC



Results

- Edges → Human faces

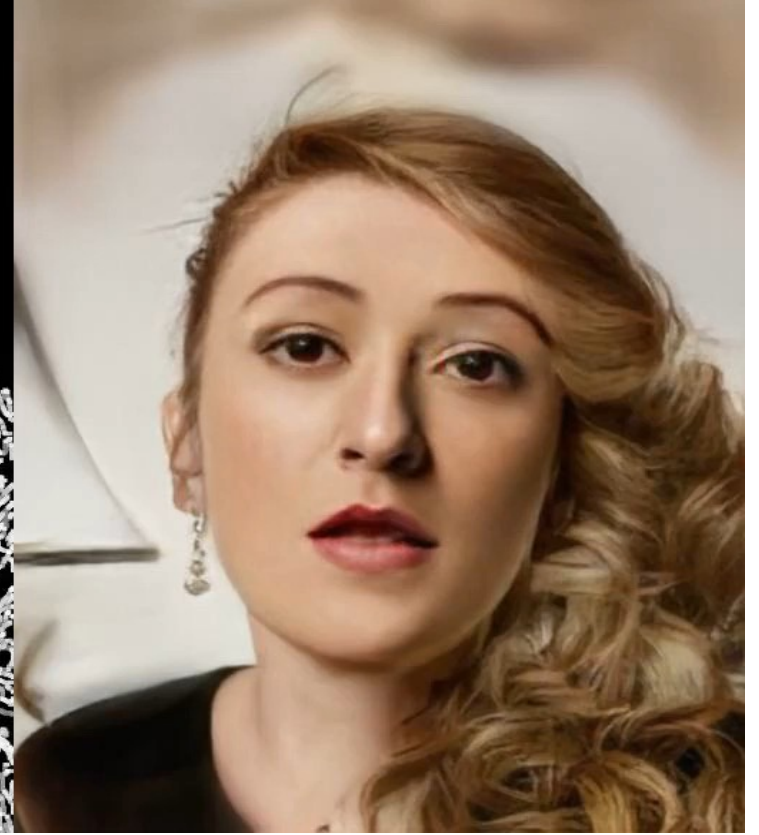
Face Swapping (face \rightarrow edge \rightarrow face)



input

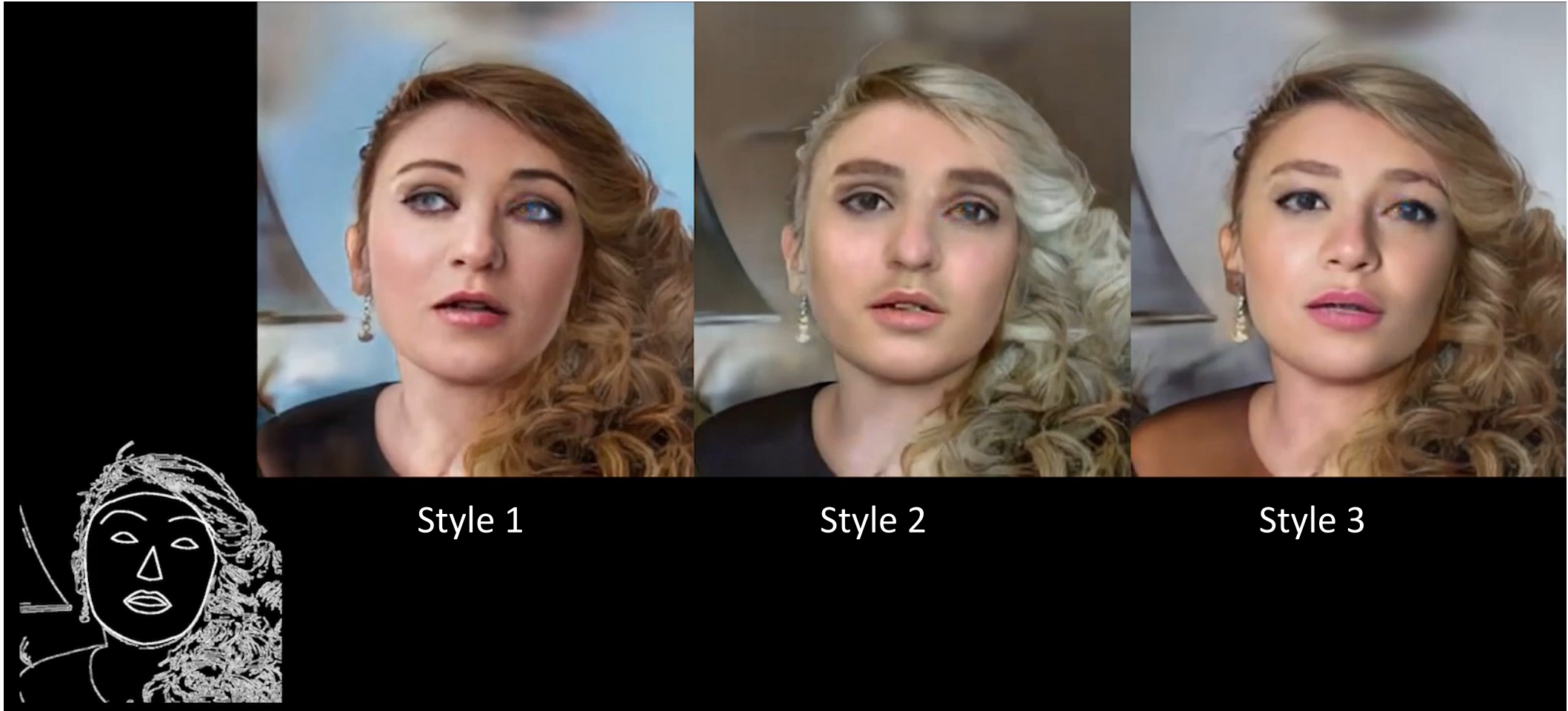


edges



output

Multi-modal Edge \rightarrow Face



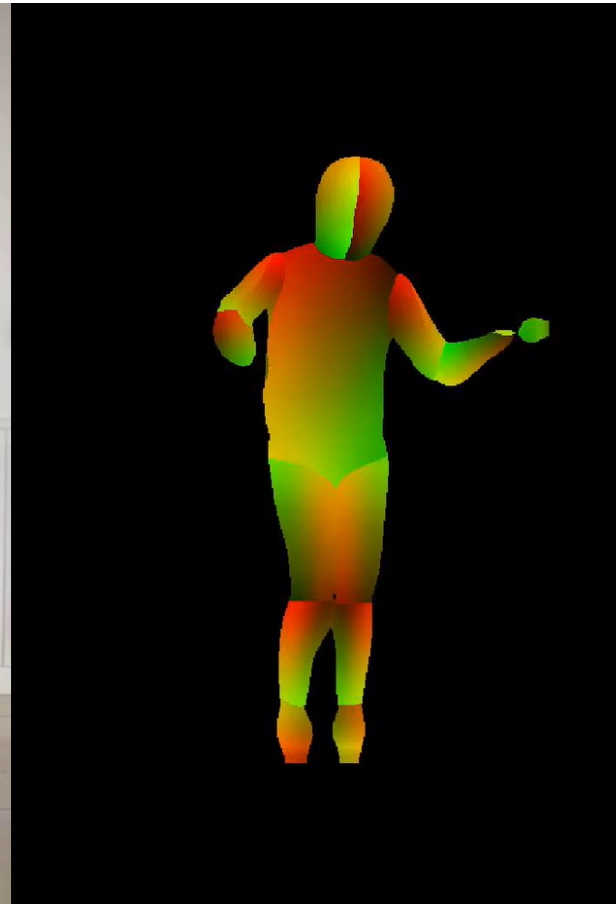
Results

- Poses → Human bodies

Motion Transfer (body \rightarrow pose \rightarrow body)



input



poses

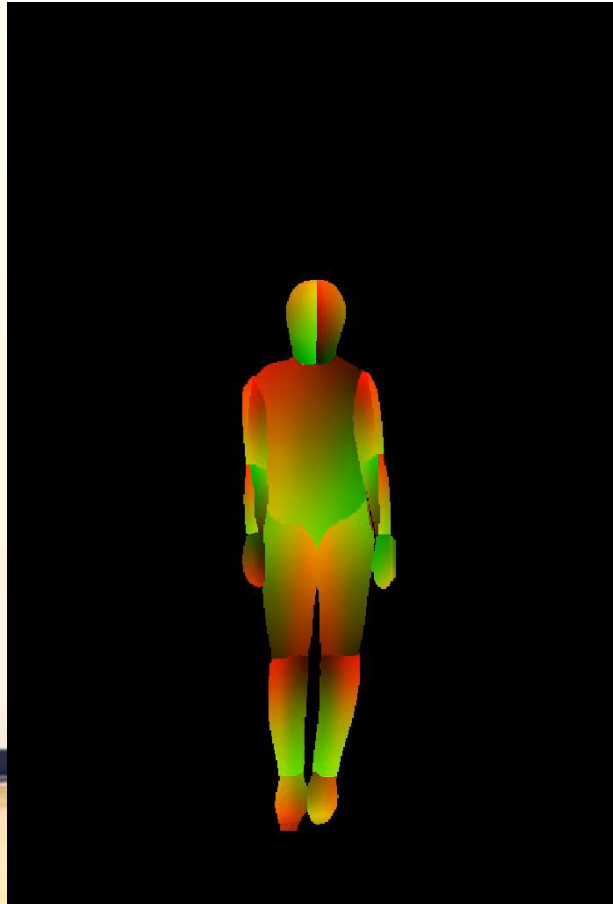


output

Motion Transfer (body \rightarrow pose \rightarrow body)



input



poses

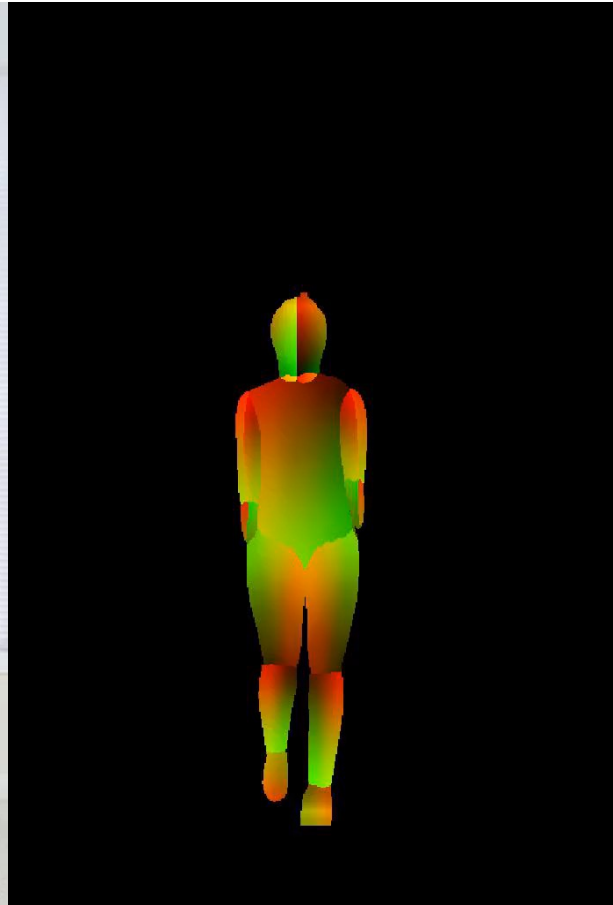


output

Motion Transfer (body \rightarrow pose \rightarrow body)



input



poses



output

Motion Transfer (body \rightarrow pose \rightarrow body)



input

poses

output

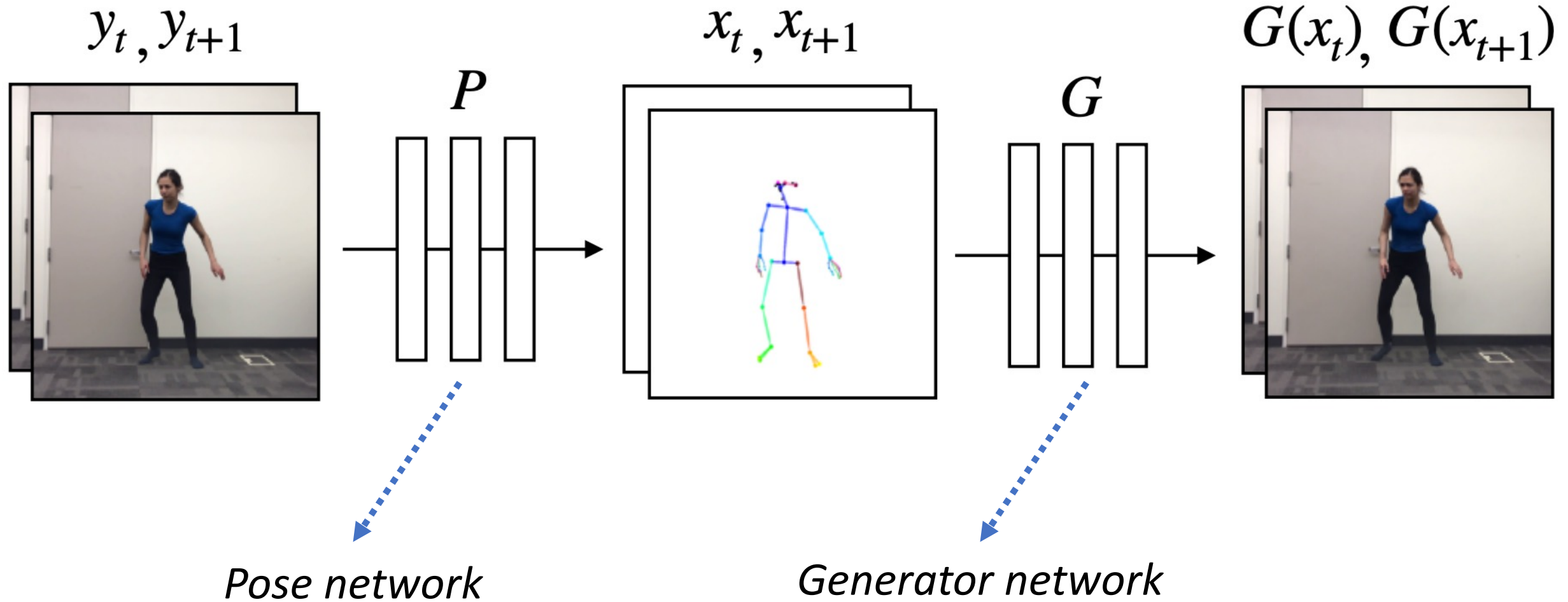
More Dancing ...

Source Subject

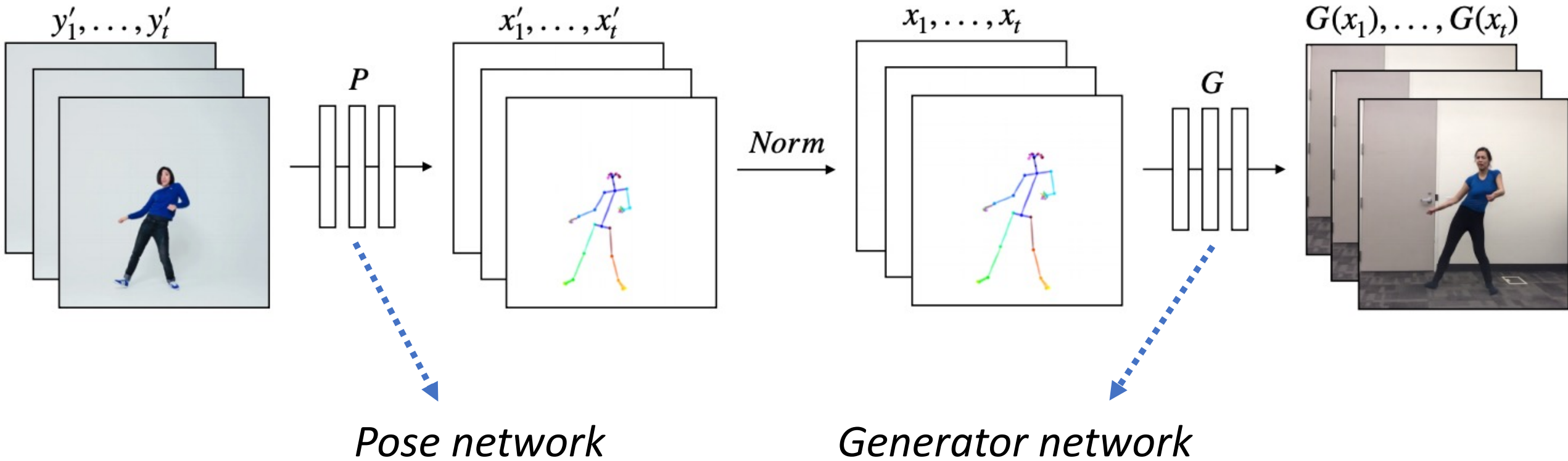
*Challenging due to missed detections



Pose-guided Synthesis

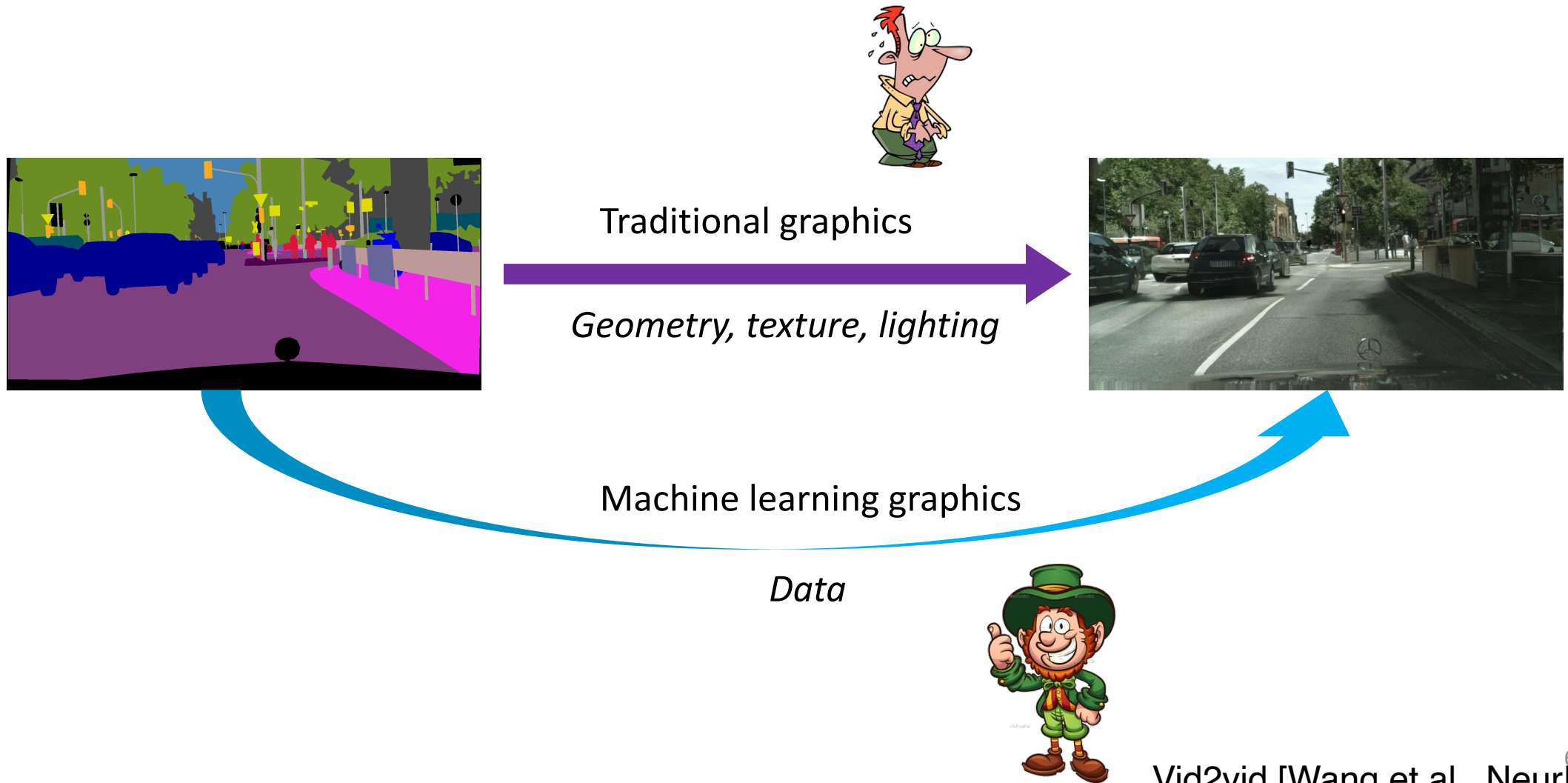


Transfer Phase





ML-based Rendering



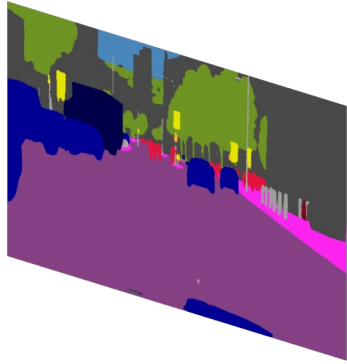
vid2vid Extensions: Interactive Graphics

User control

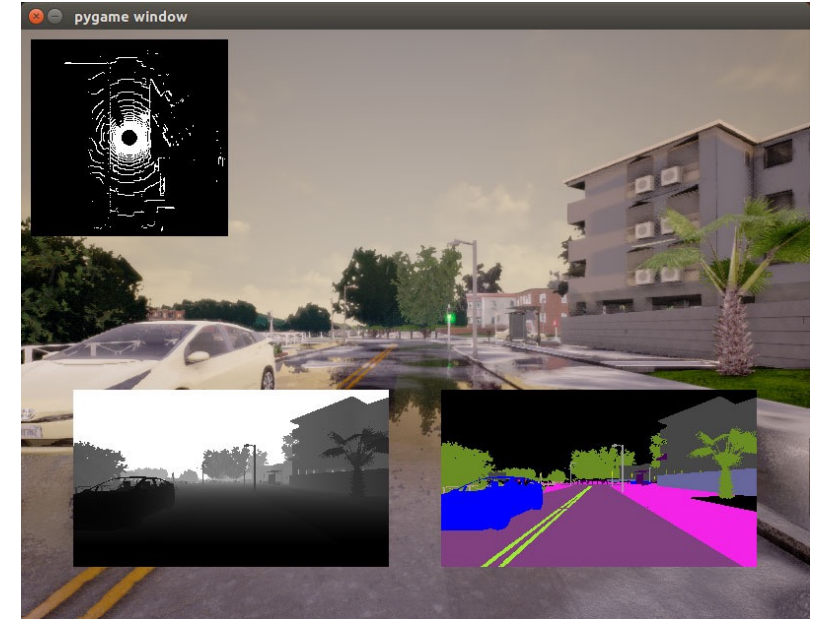
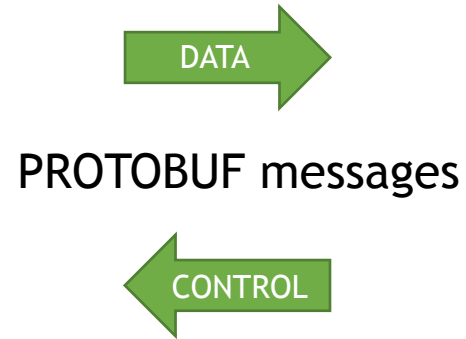
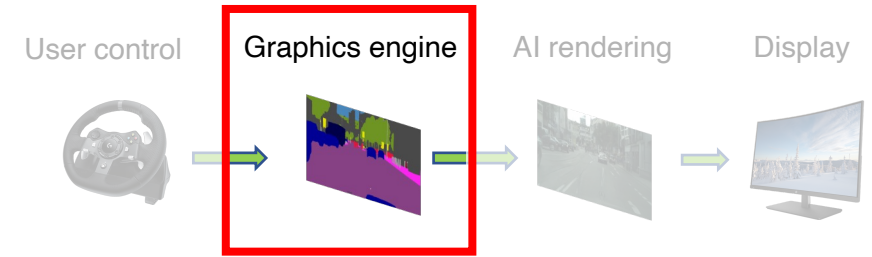
Graphics engine

AI rendering

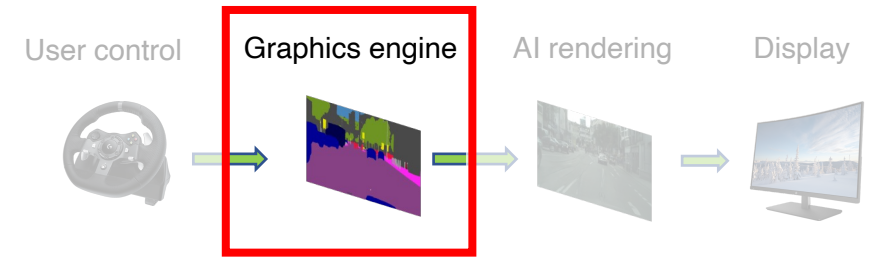
Display



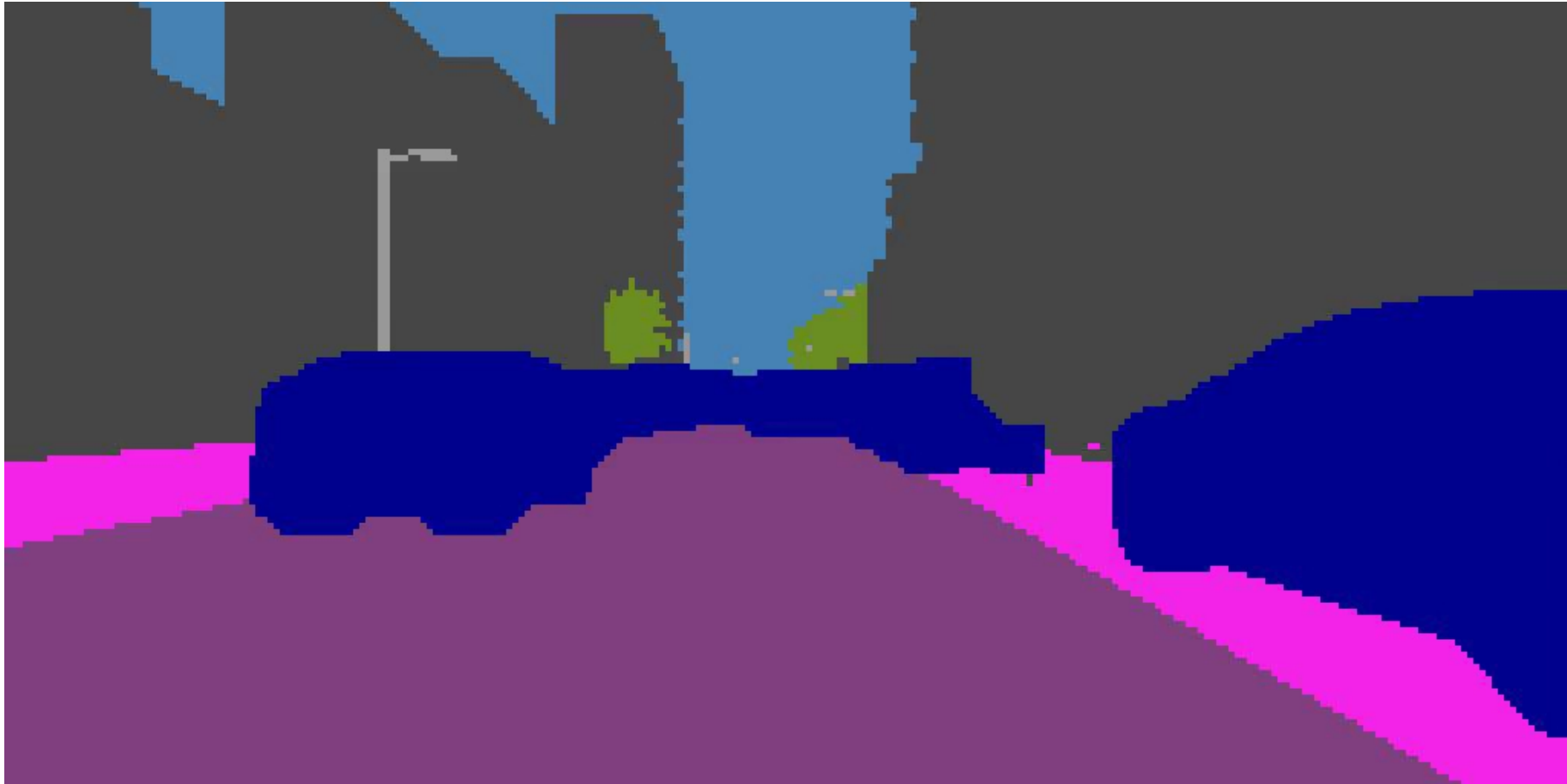
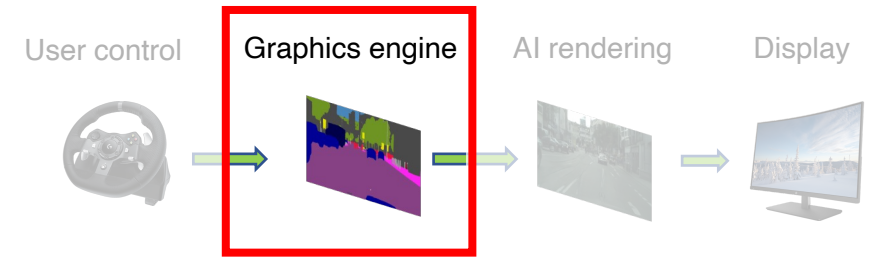
Graphics Engine: CARLA



Original CARLA Sequence

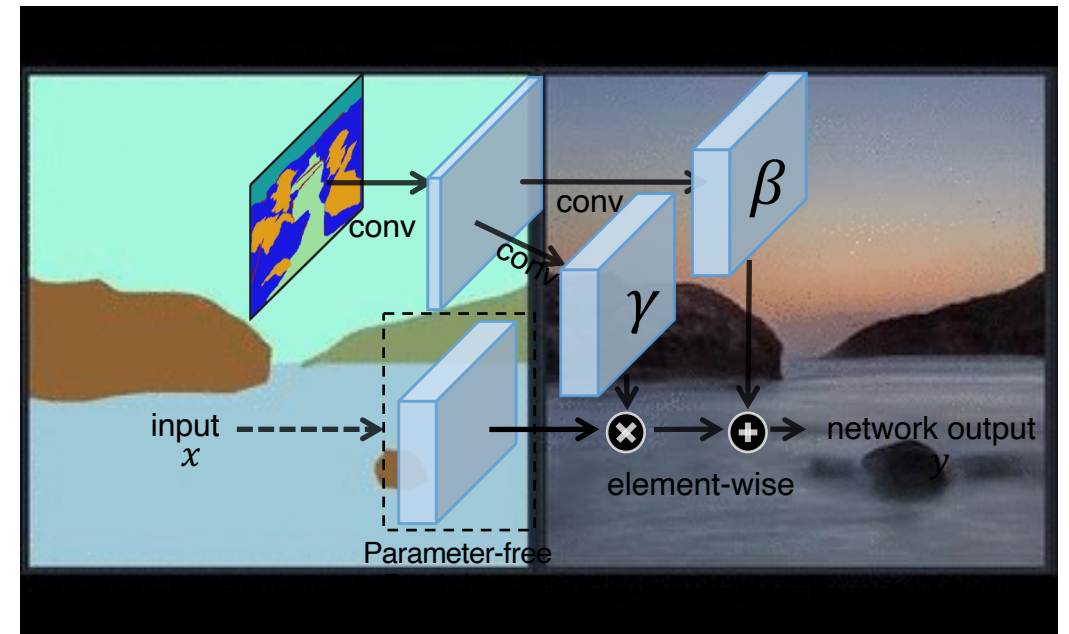
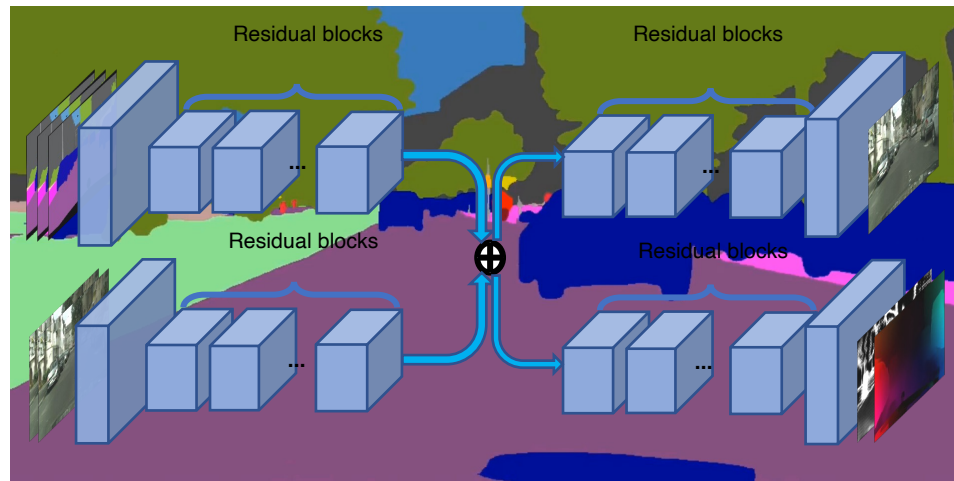
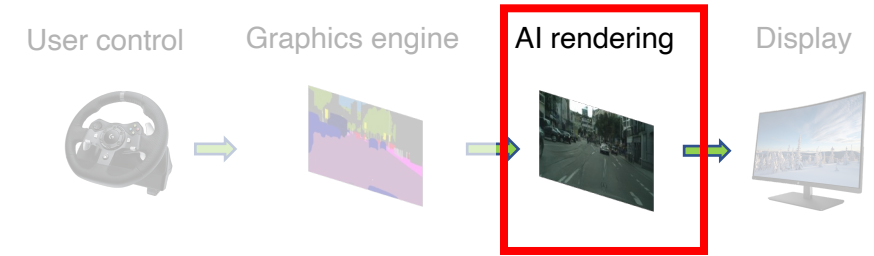


CARLA Semantic Maps

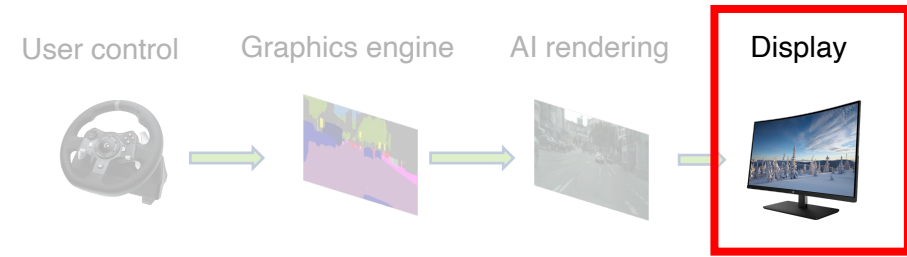


Methodology

- Combine vid2vid with SPADE



Demo Result



Driving Game



vid2game by FAIR

O. Gafni, L. Wolf, Y. Taigman. "Vid2Game: Controllable Characters Extracted from Real-World Videos," 2019



Thank You!



16-726, Spring 2023

<https://learning-image-synthesis.github.io/>