



Video Synthesis and Editing

Jun-Yan Zhu

16-726 Learning-based Image Synthesis, Spring 2022

Logistics

- HW5 award on Wed.
- 5 min (including Q & A) x 23 groups / 60 min ~ 2 hours
- Strict timing.
- Preliminary results.
- Common failure cases:
 - You didn't show up.
 - You gave a 10 min talk.
 - You spent 4.5 min on related work.
 - We cannot open your link.

Image Editing and Synthesis

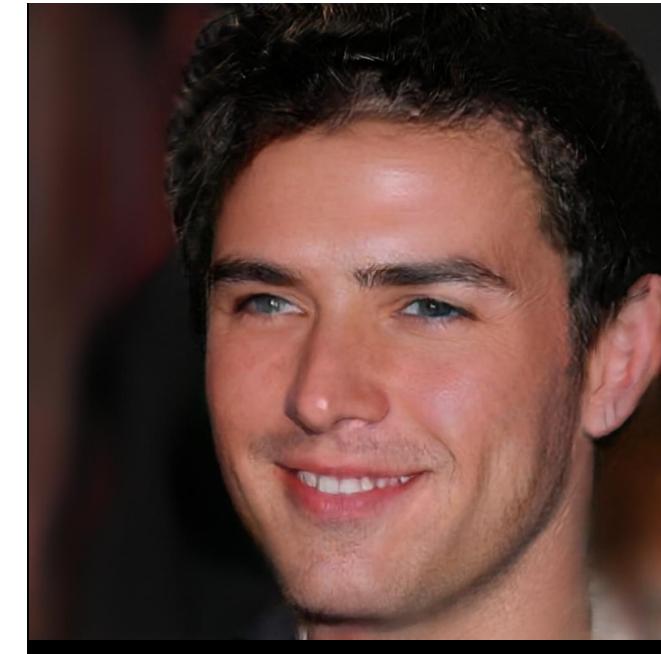


Image Editing and Synthesis

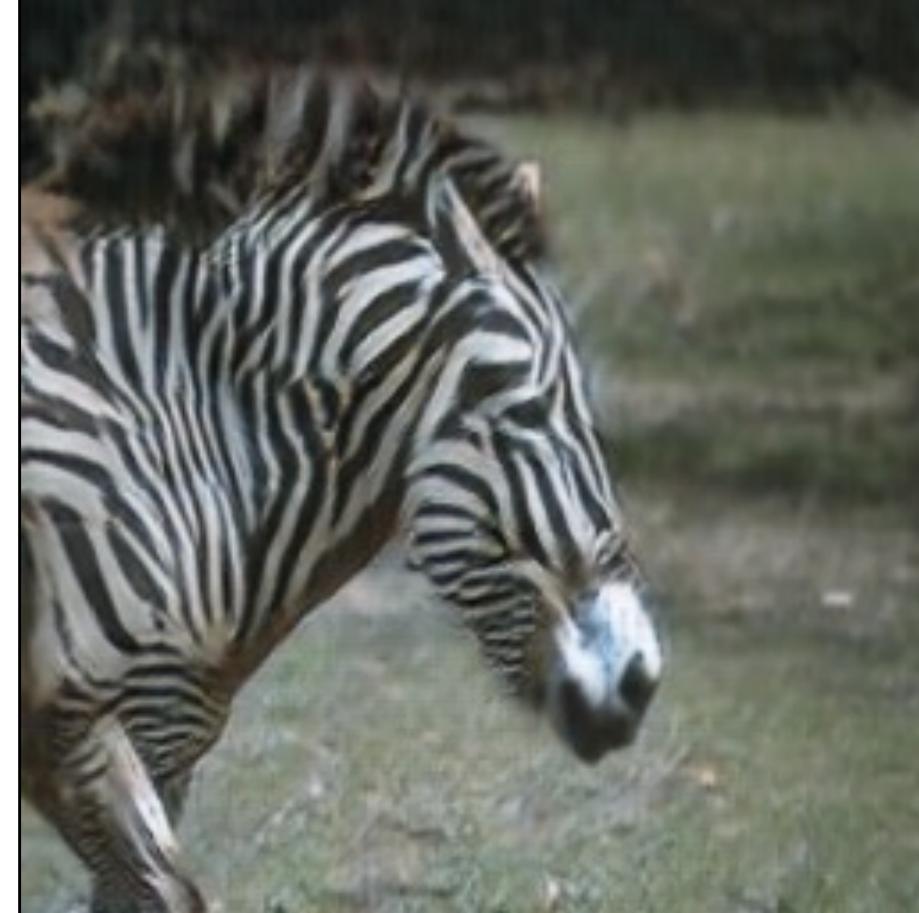
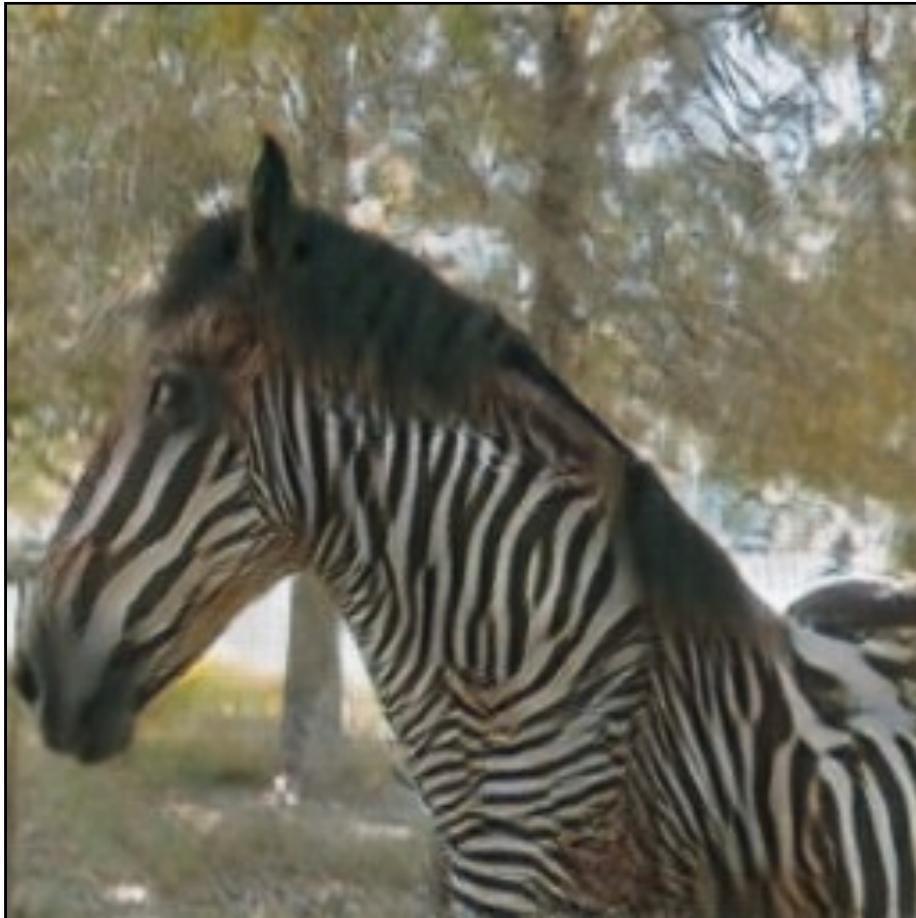
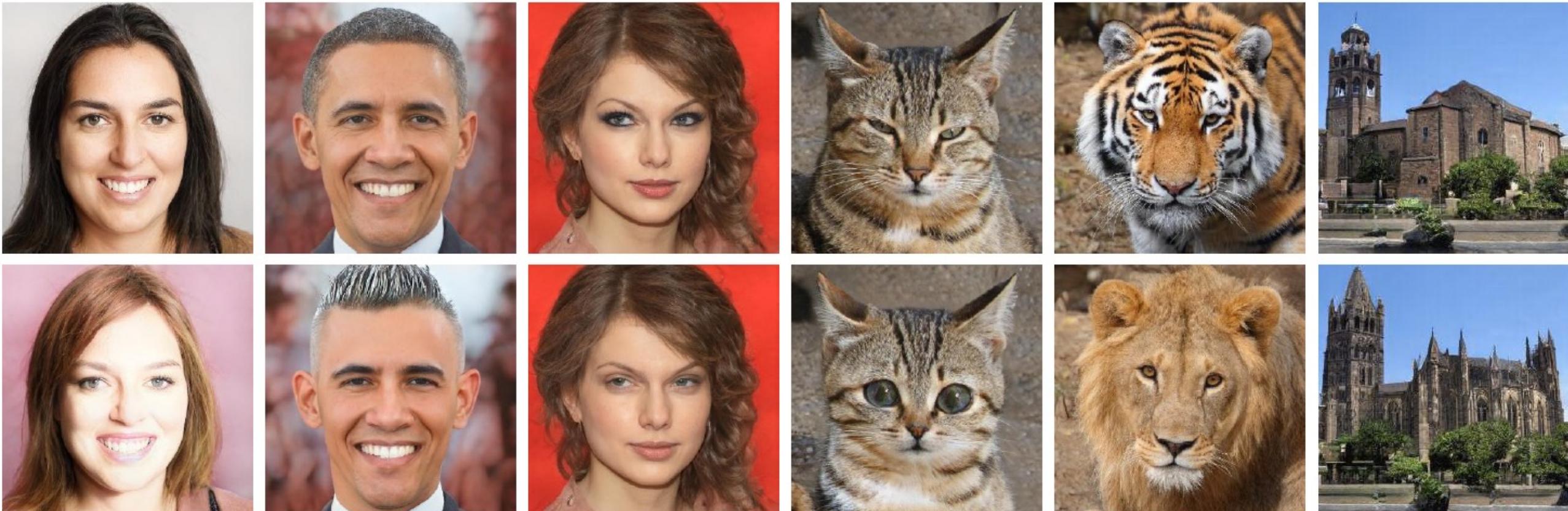


Image Editing and Synthesis



“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

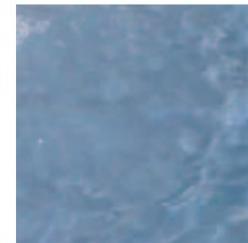
Spatial-temporal Constraints

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

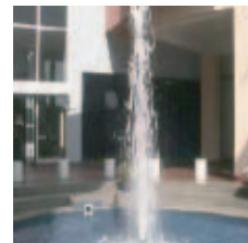
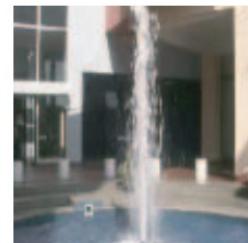
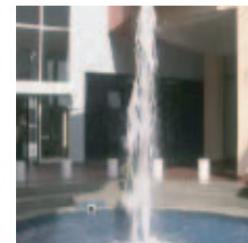
Output
Image

Guidance Gradient

Background



Foreground

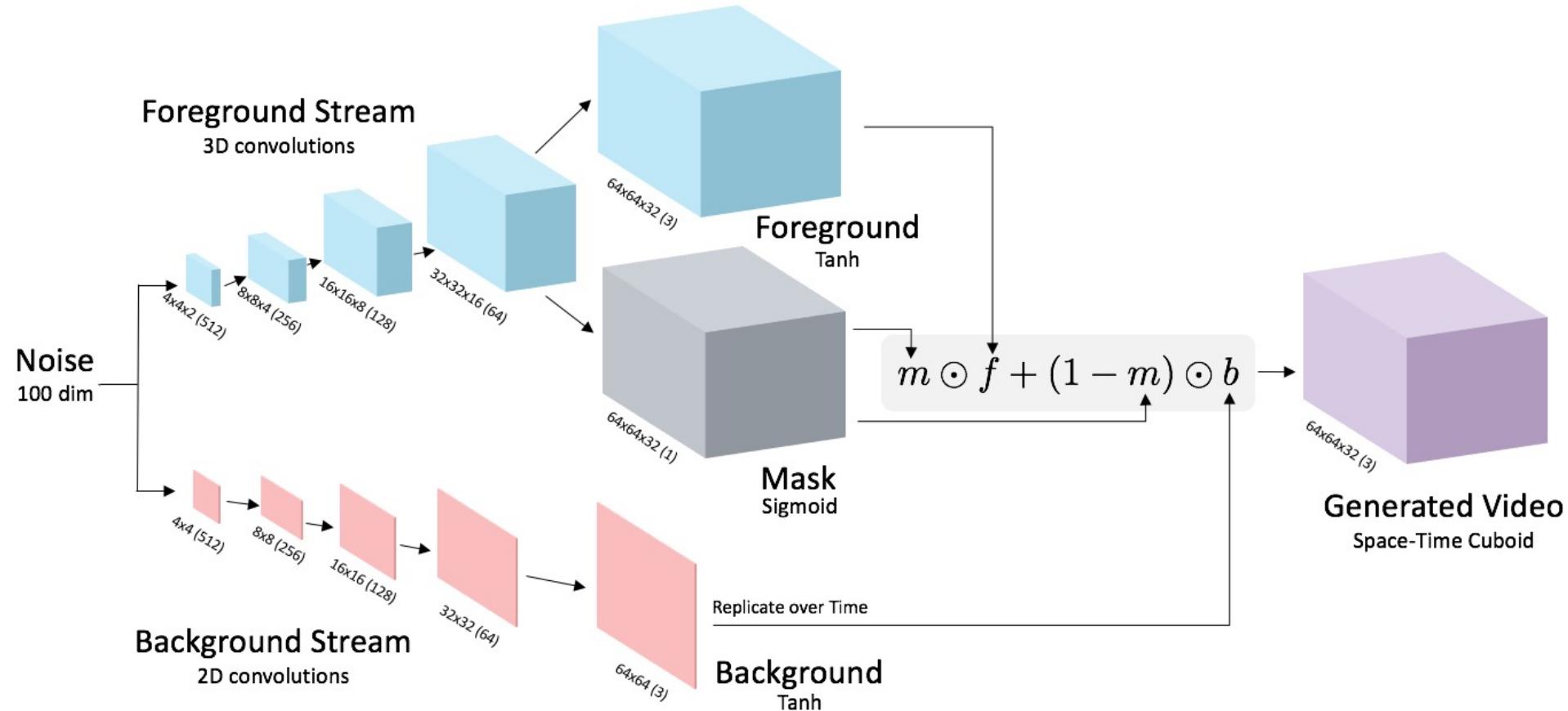


Output

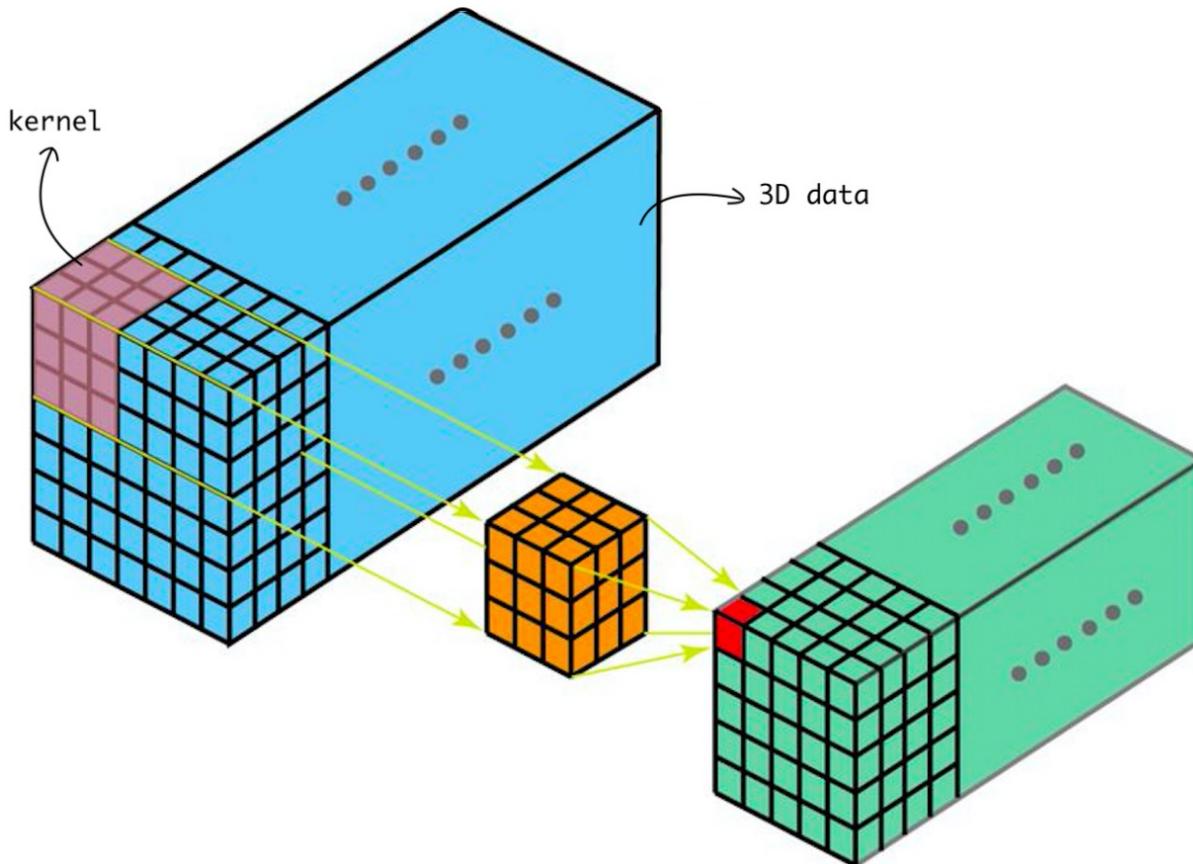


Time

Case Study: Video GANs



Recap: 3D Conv



Easy to implement:

- Replace 2D by 3D in your code

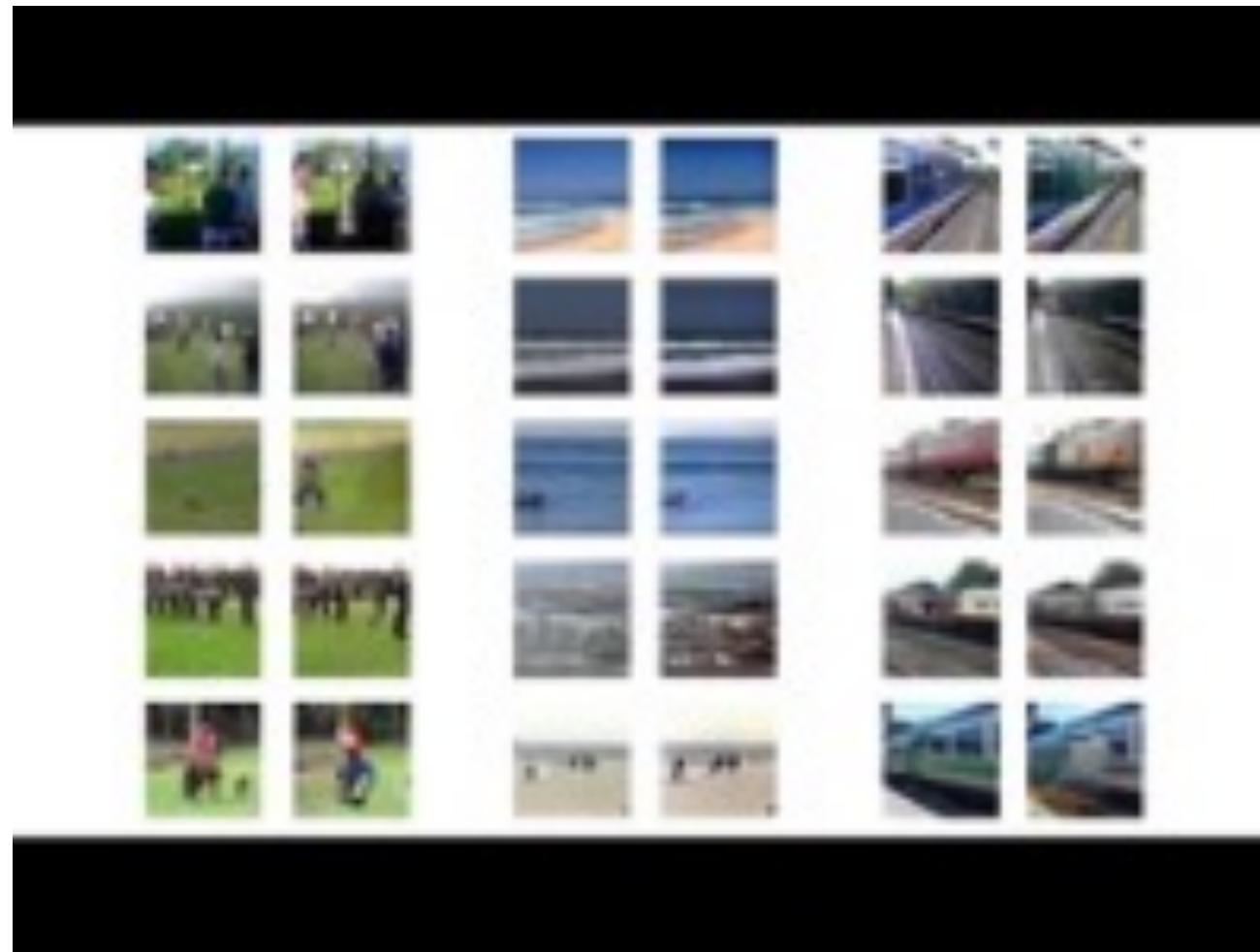
e.g., Conv2D -> Conv3D

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



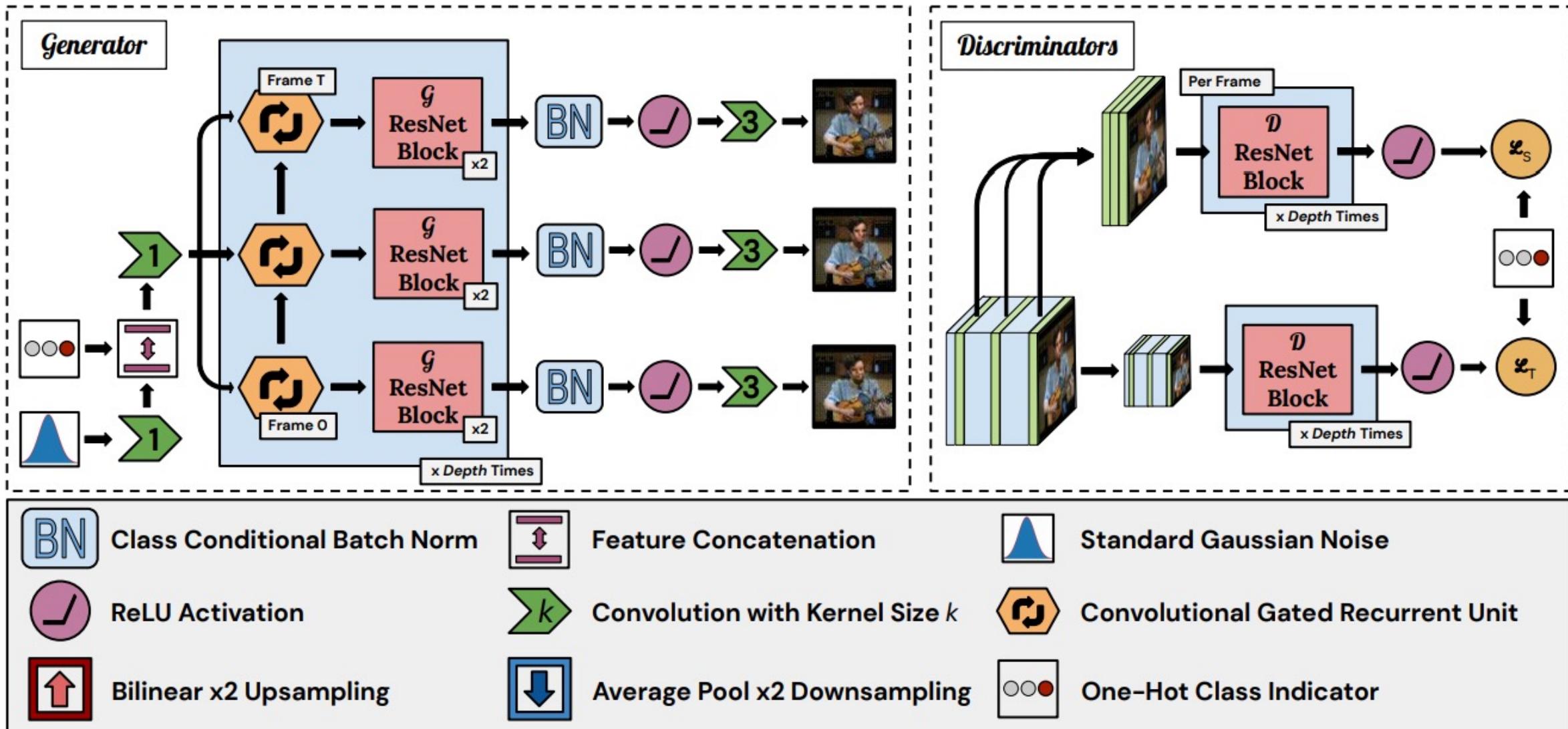
Case Study: DVDGANs



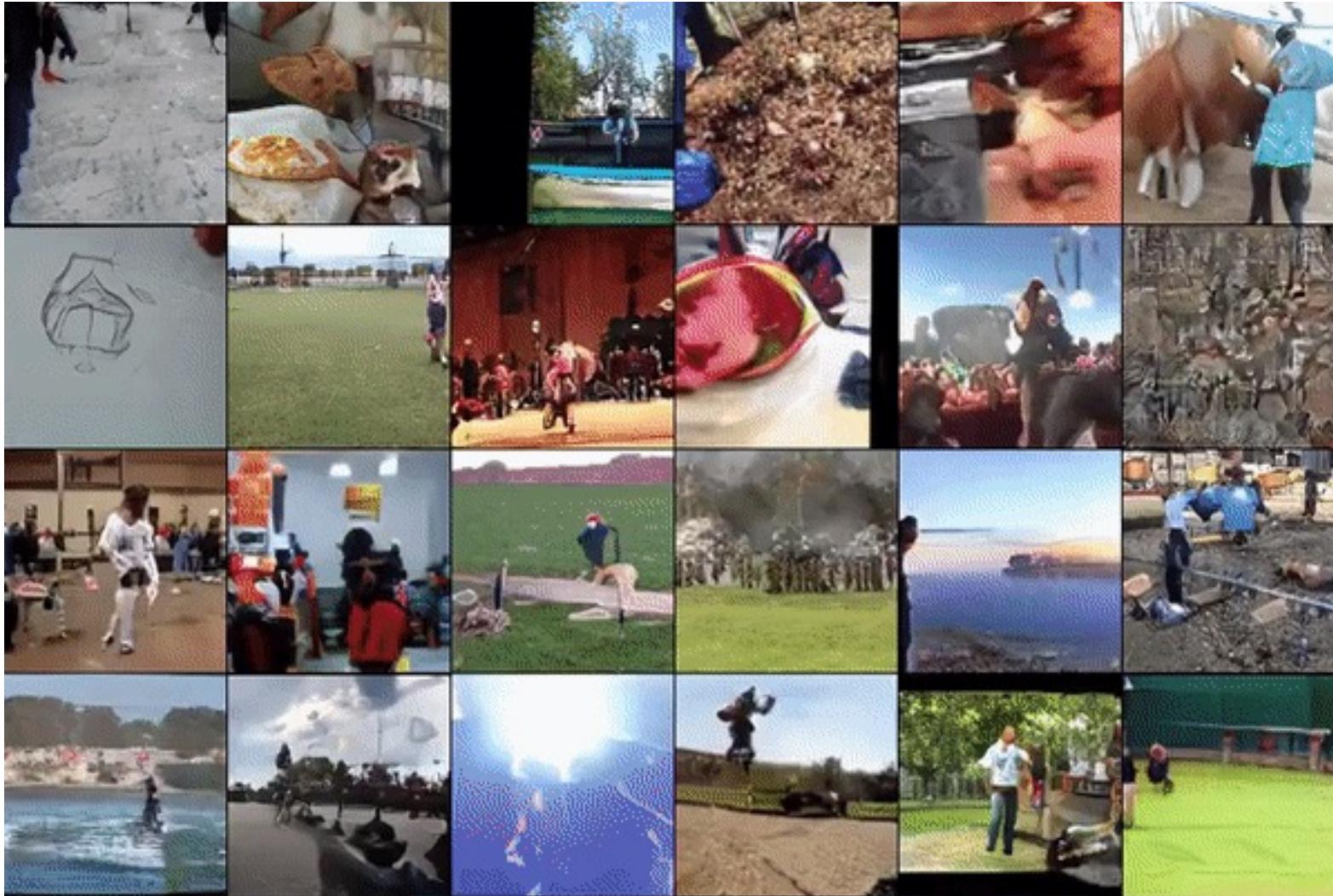
BigGAN-based generator + Spatial Discriminator + Temporal Discriminator

16
Clark et all, arXiv 2019

Case Study: DV DGANs



Case Study: DVDGANs



How to Generalize to Videos

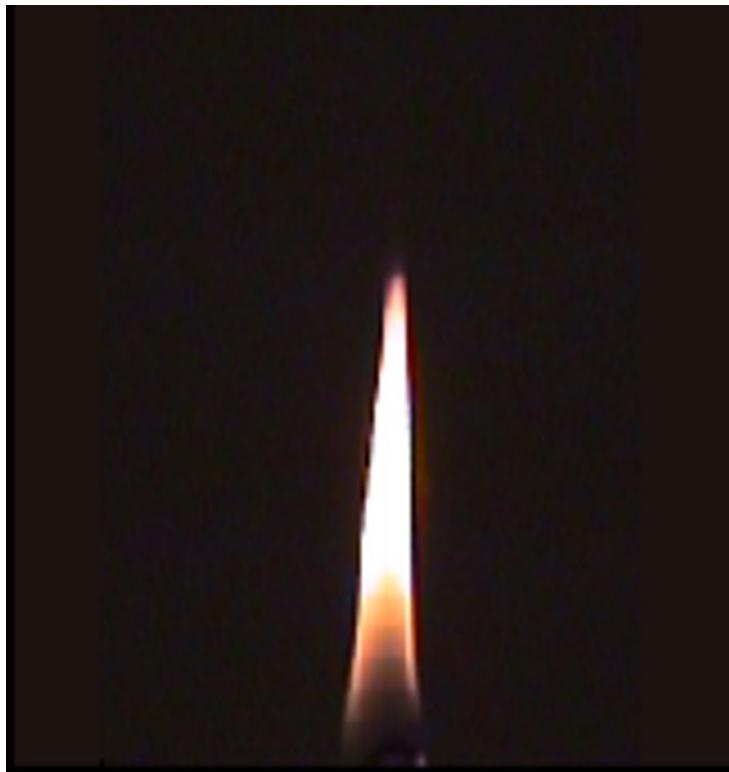
- Idea 1: Frame-by-Frame
 - Temporal inconsistency (flicking, color drift)
- Idea 2: Video as 3D data (height x width x time)
 - memory-intensive and time-consuming
 - only work for a short video at low resolution
- Idea 3: recurrent (autoregressive) synthesis
 - Generate 1st frame, generate 2nd frame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

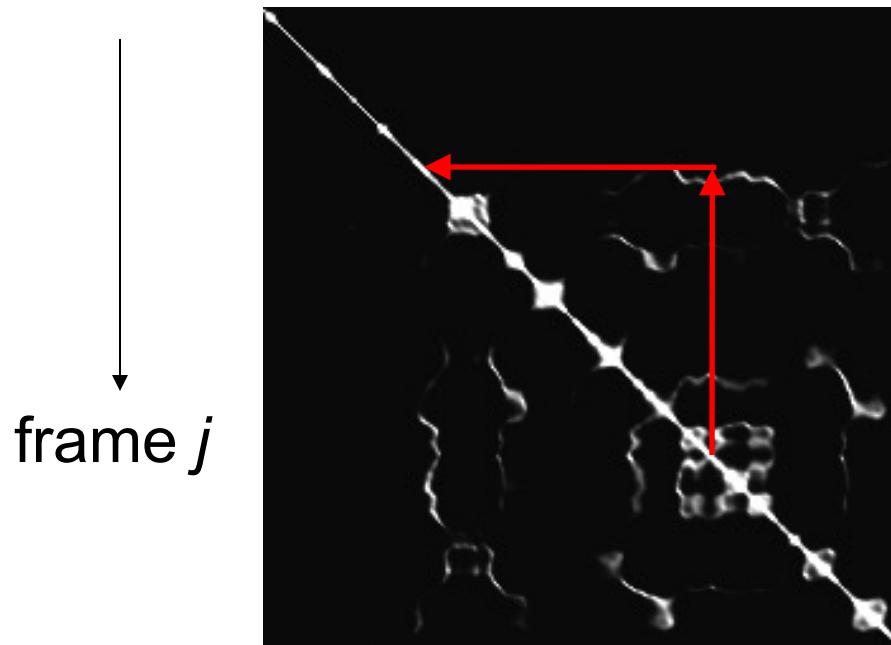


Finding good transitions

- Compute L_2 distance $D_{i,j}$ between all frames

vs.

→ frame i

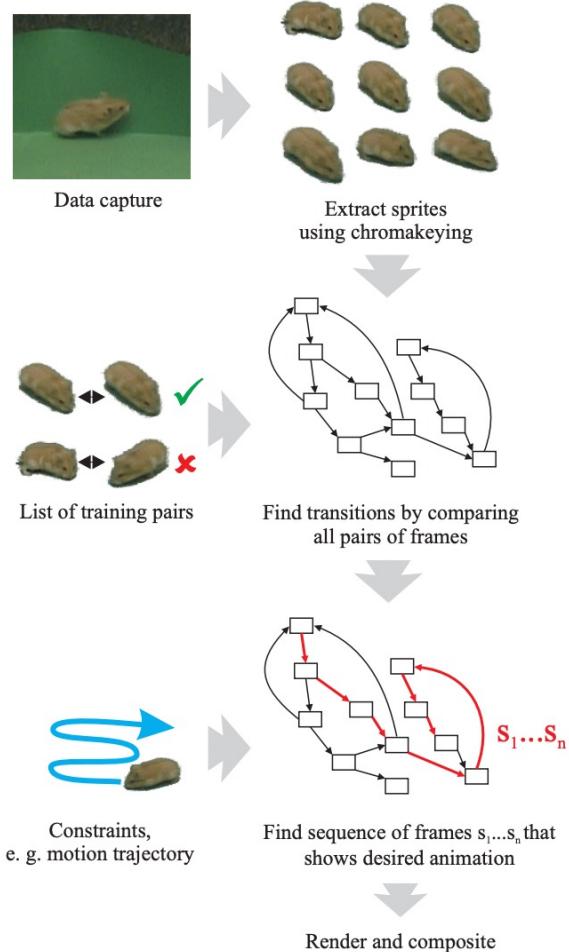


Similar frames make good transitions

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

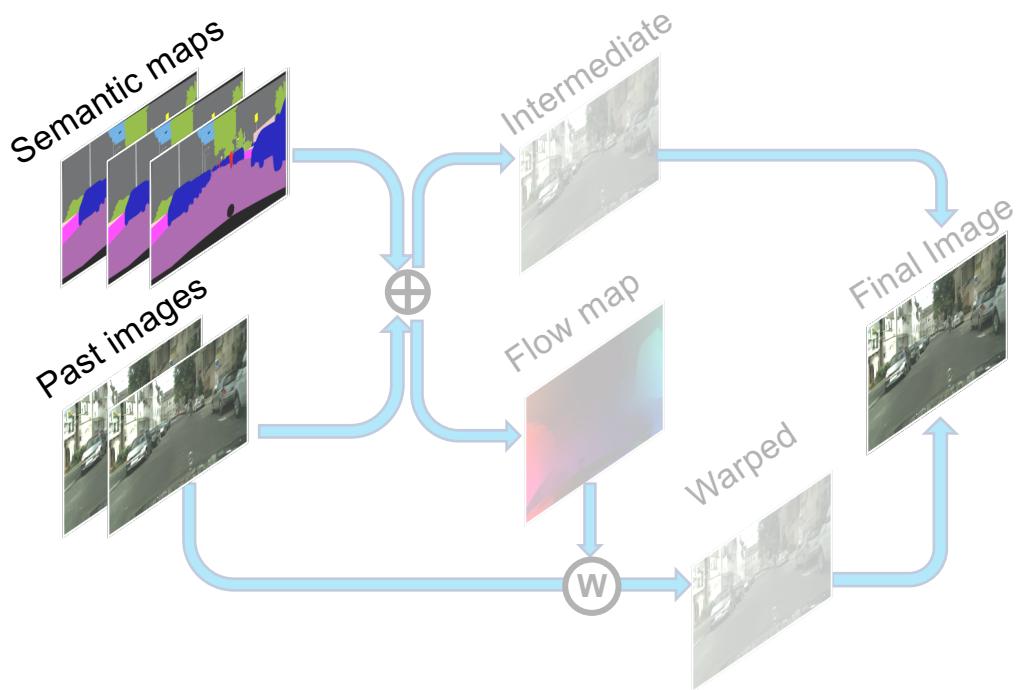


vid2vid

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

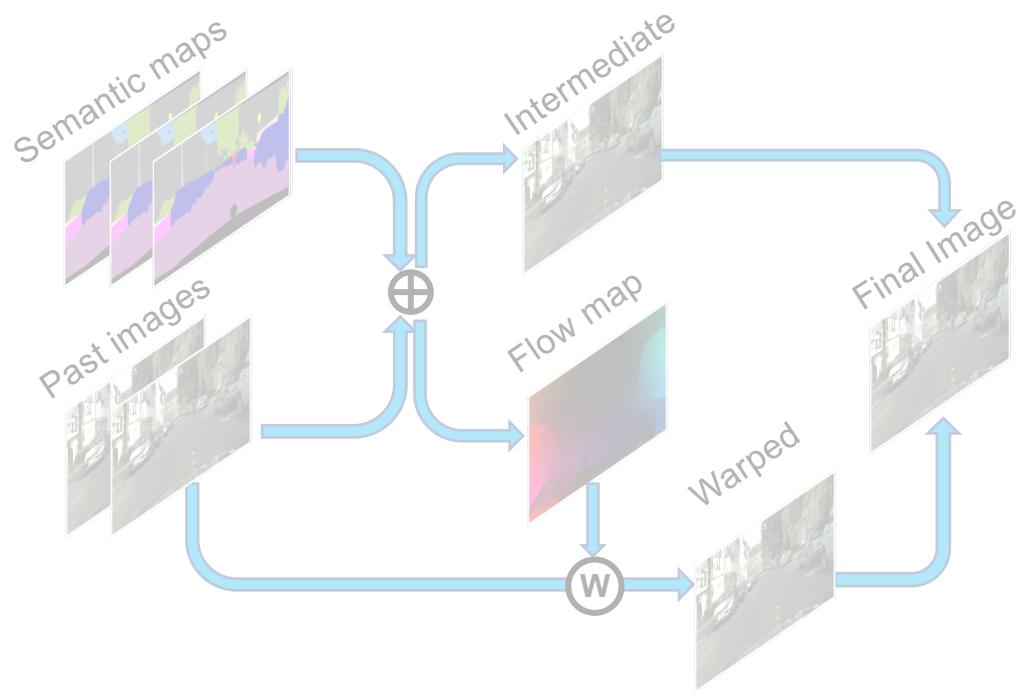
vid2vid

Sequential Generator



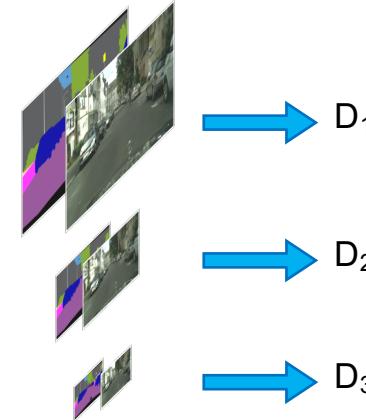
vid2vid

Sequential Generator

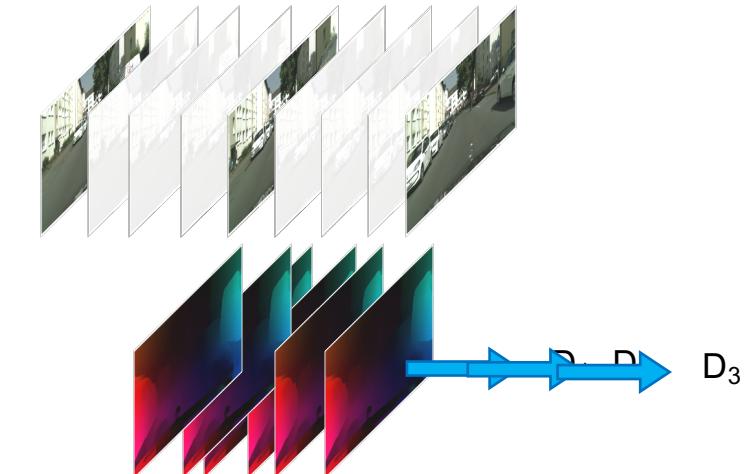


Multi-scale Discriminators

Image Discriminator



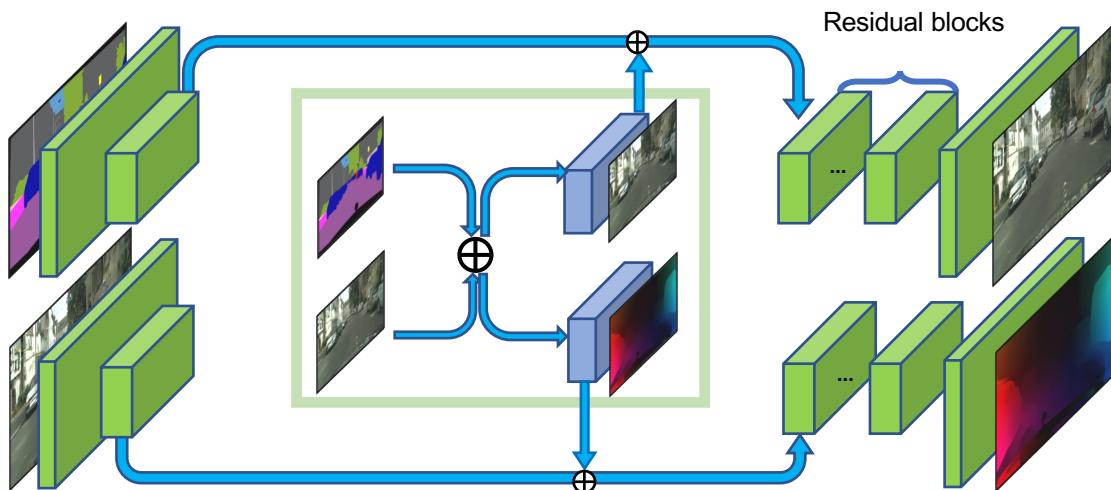
Video Discriminator



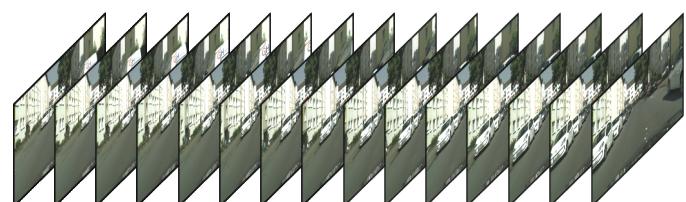
vid2vid

Spatio-temporally Progressive Training

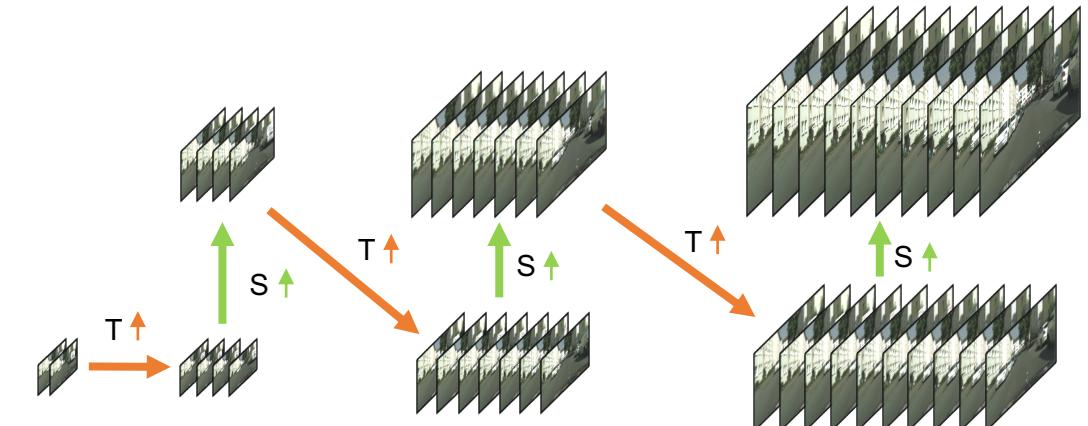
Spatially progressive



Temporally progressive



Alternating training



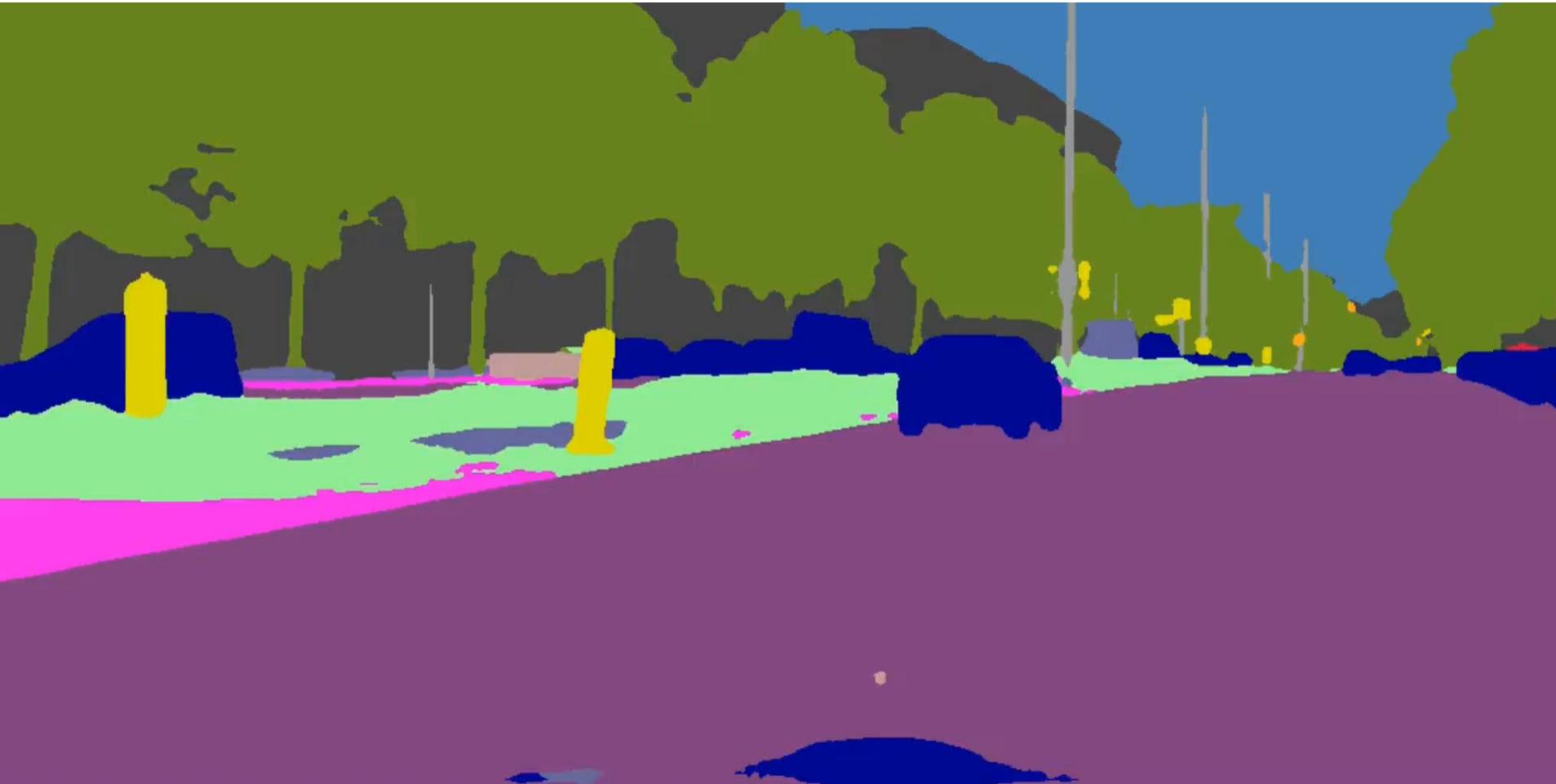
vid2vid Results

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

vid2vid Results

- Semantic → Street view scenes

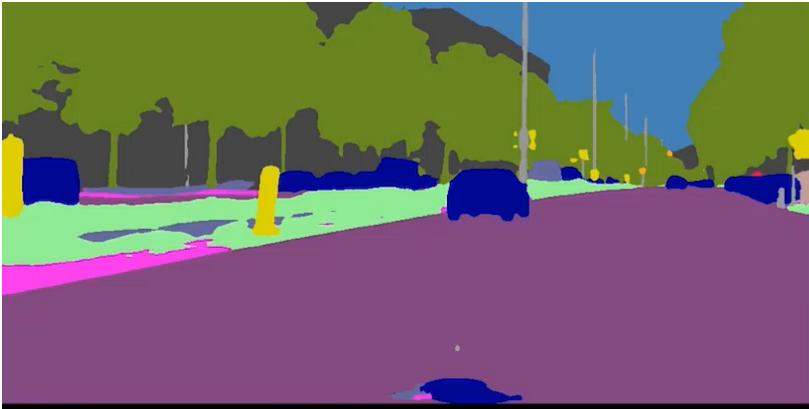
Street View: Cityscapes



Street View: Cityscapes



Street View: Cityscapes



Labels



pix2pixHD



COVST



Ours

Street View: Boston



Street View: NYC



Results

- Edges → Human faces

Face Swapping (face → edge → face)

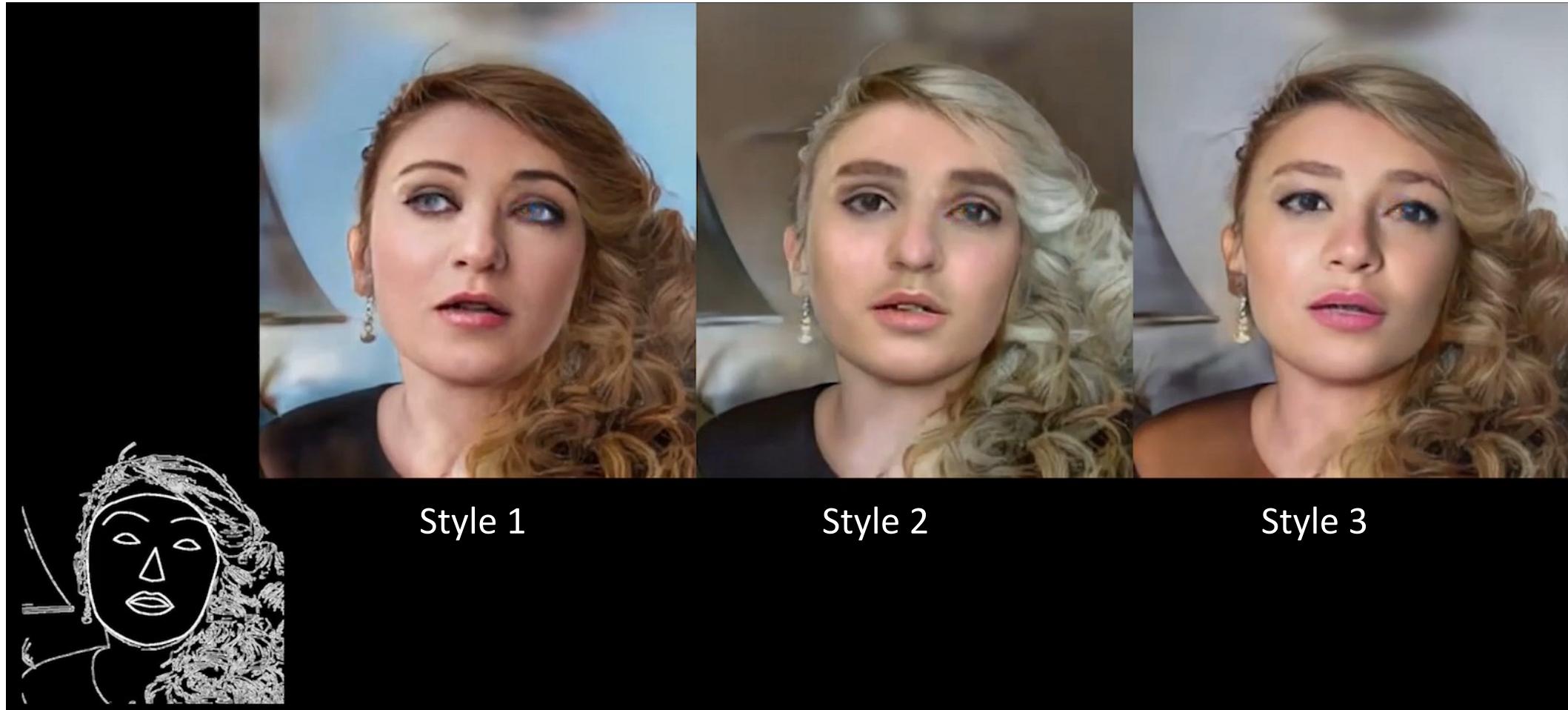


input

edges

output

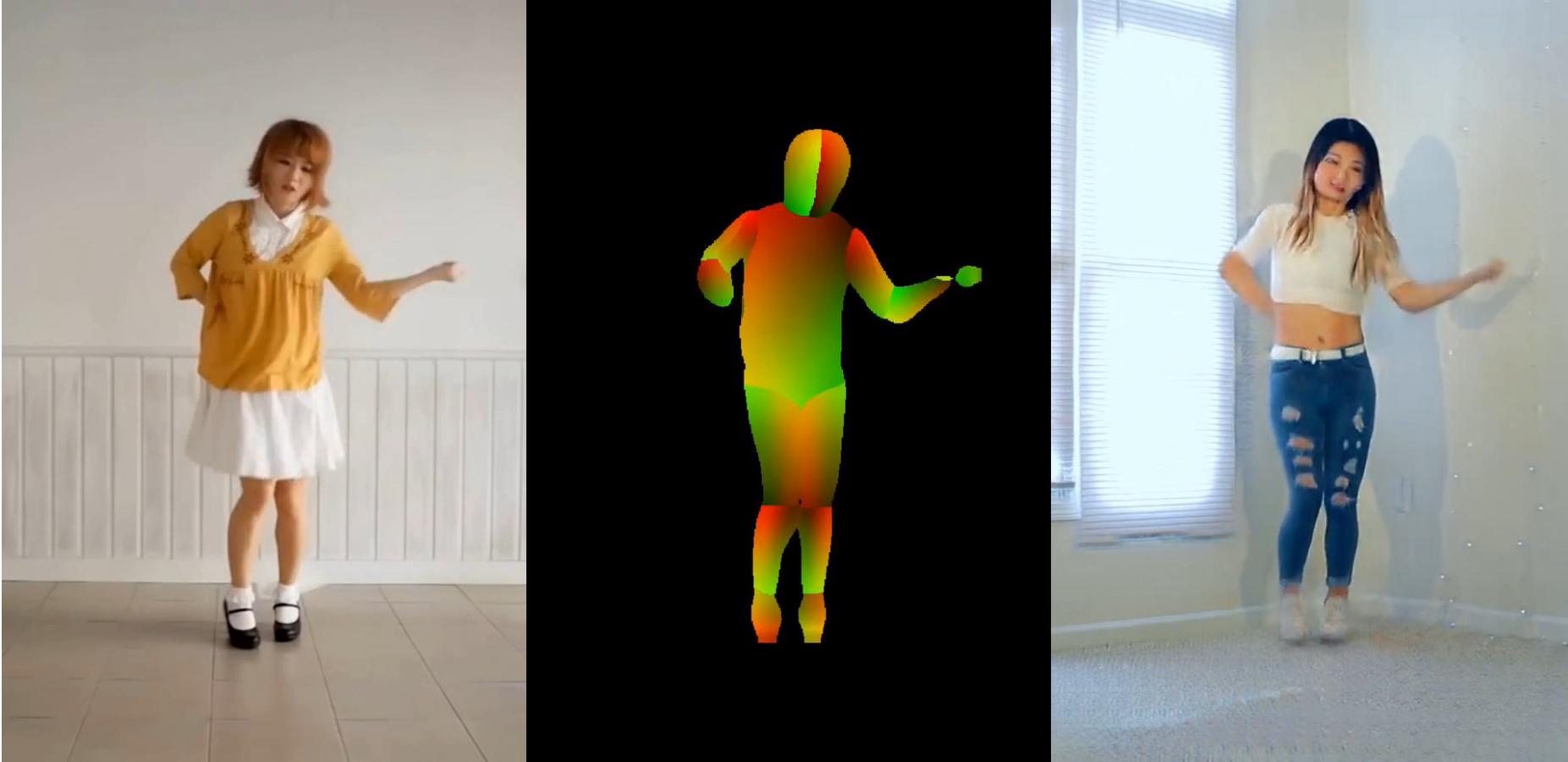
Multi-modal Edge → Face



Results

- Poses → Human bodies

Motion Transfer (body → pose → body)

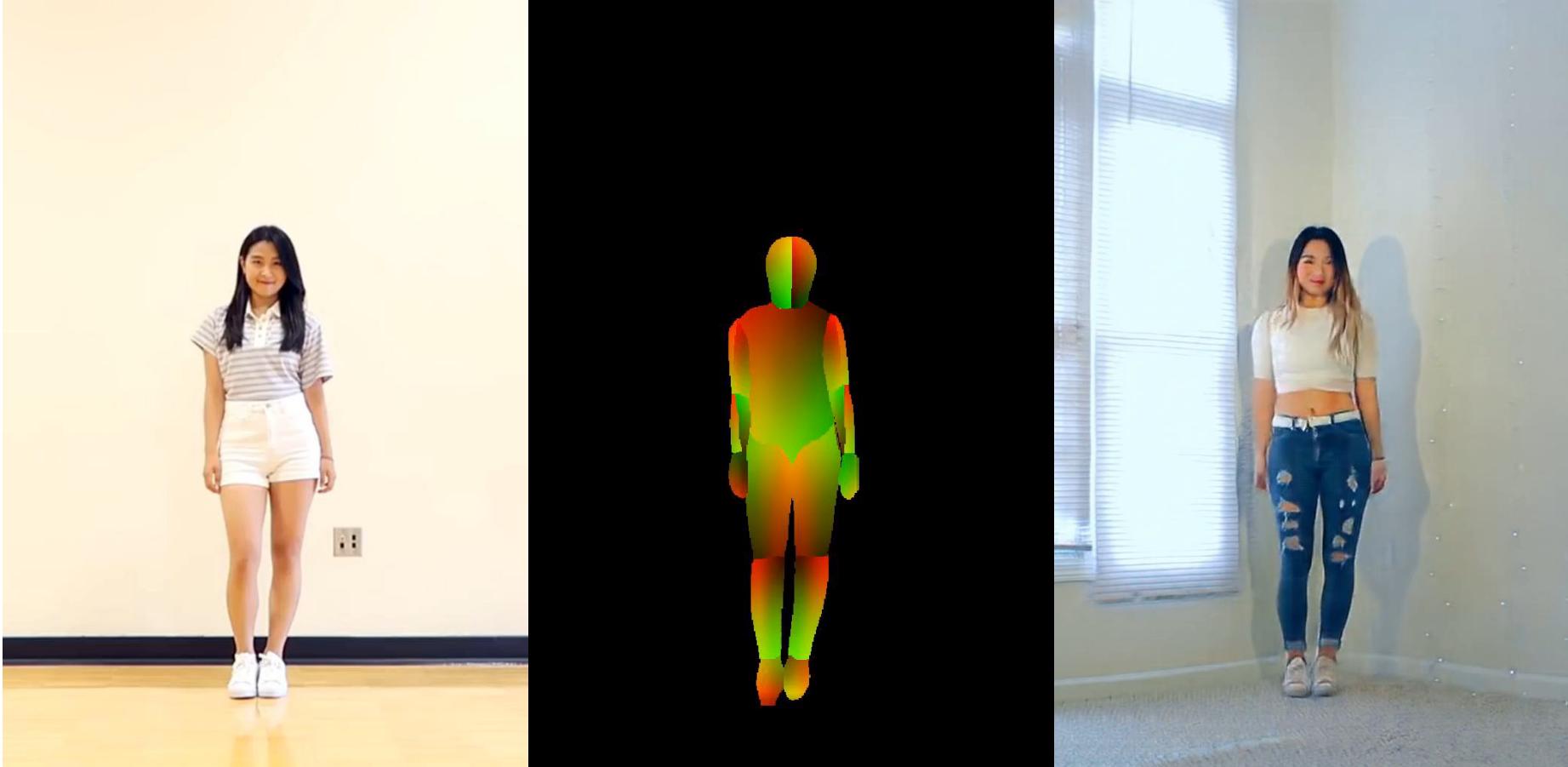


input

poses

output

Motion Transfer (body → pose → body)



input

poses

output

Motion Transfer (body → pose → body)



input

poses

output

Motion Transfer (body → pose → body)



input

poses

output

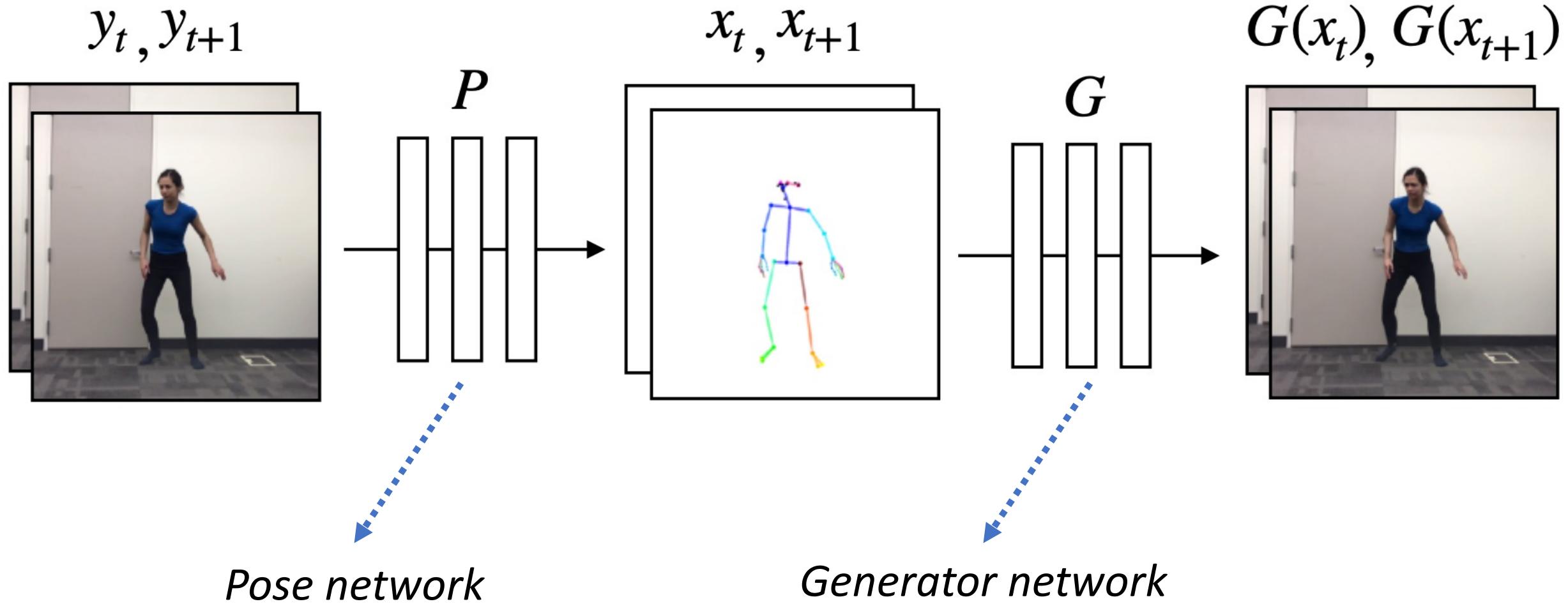
More Dancing ...

Source Subject

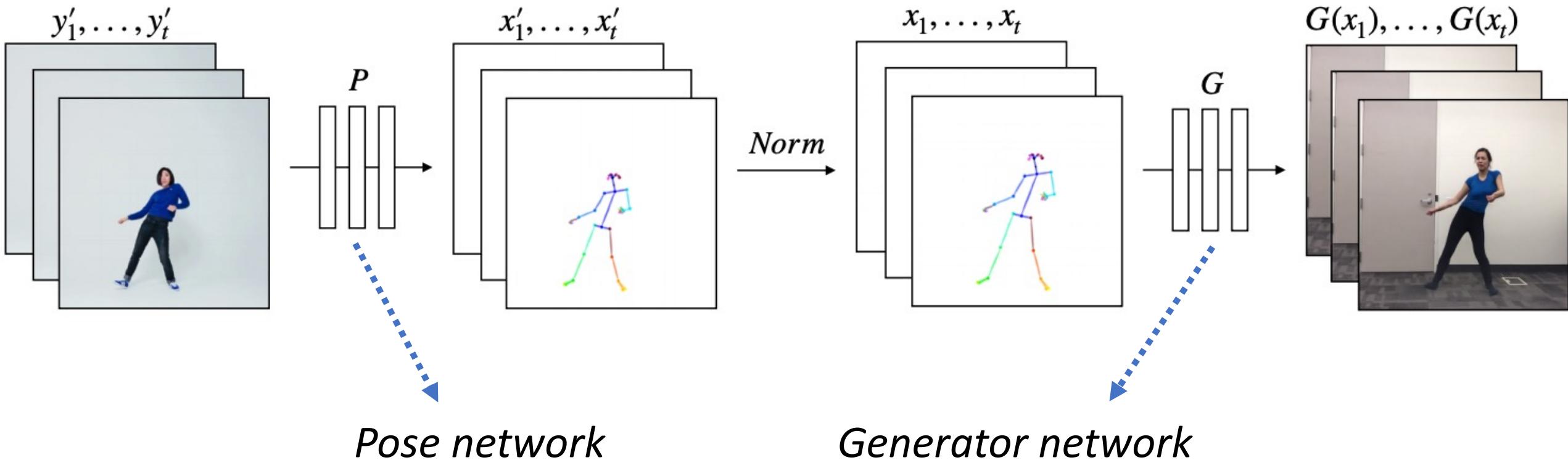
*Challenging due to missed detections



Pose-guided Synthesis



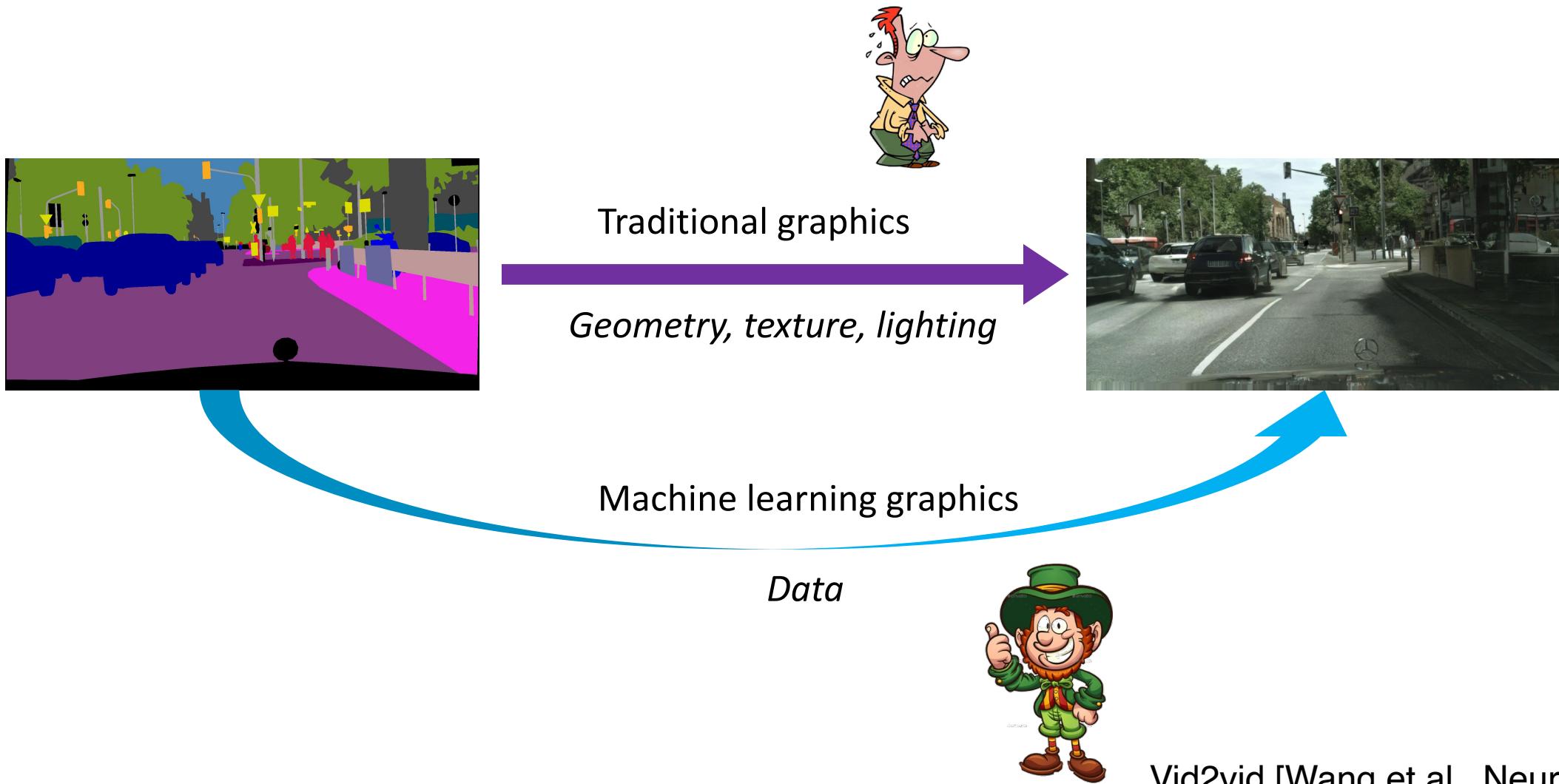
Transfer Phase



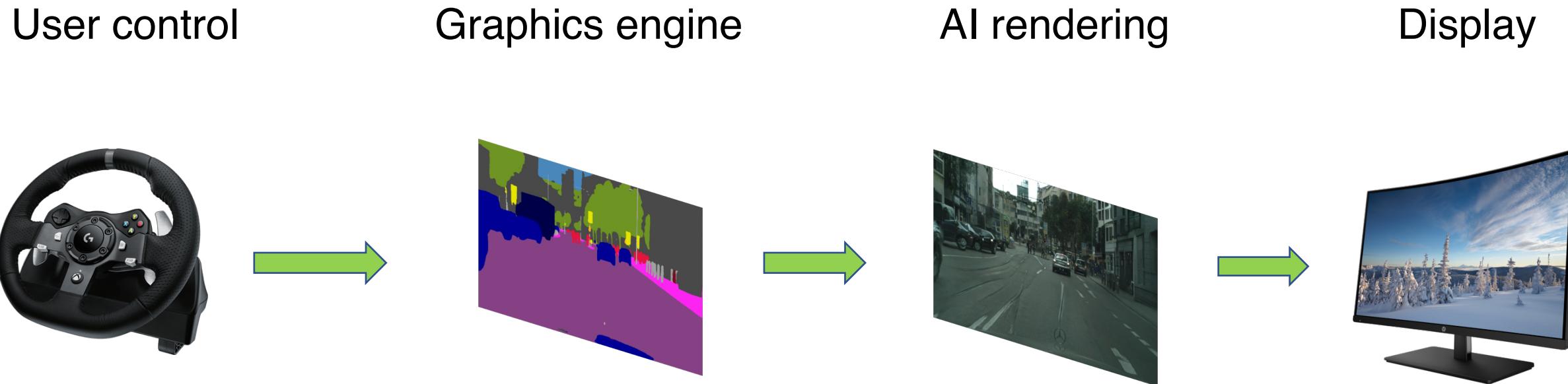


51
Everybody Dance Now [Chan et al., ICCV 2019]

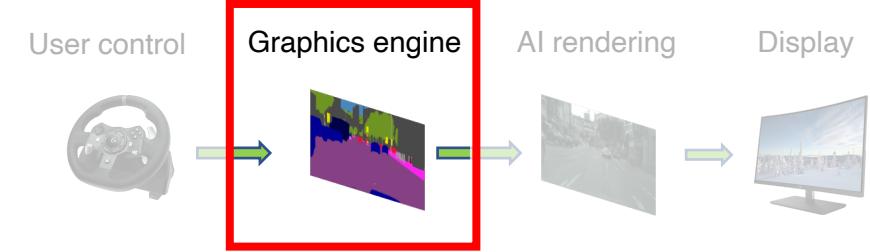
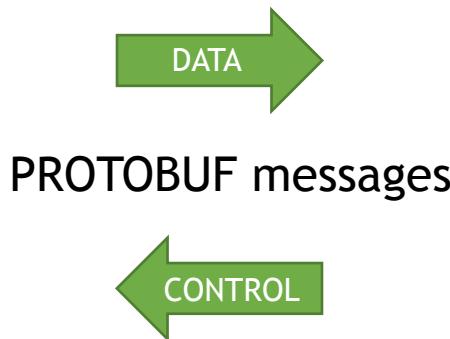
ML-based Rendering



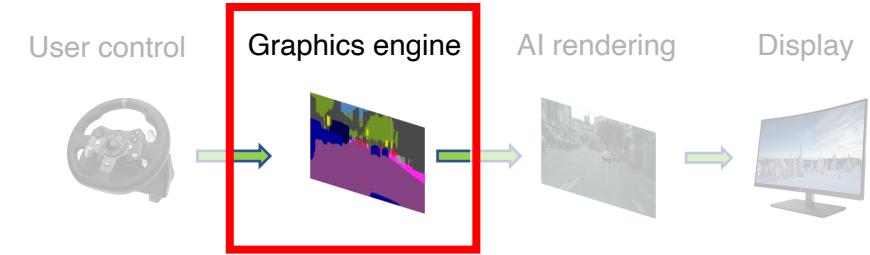
vid2vid Extensions: Interactive Graphics



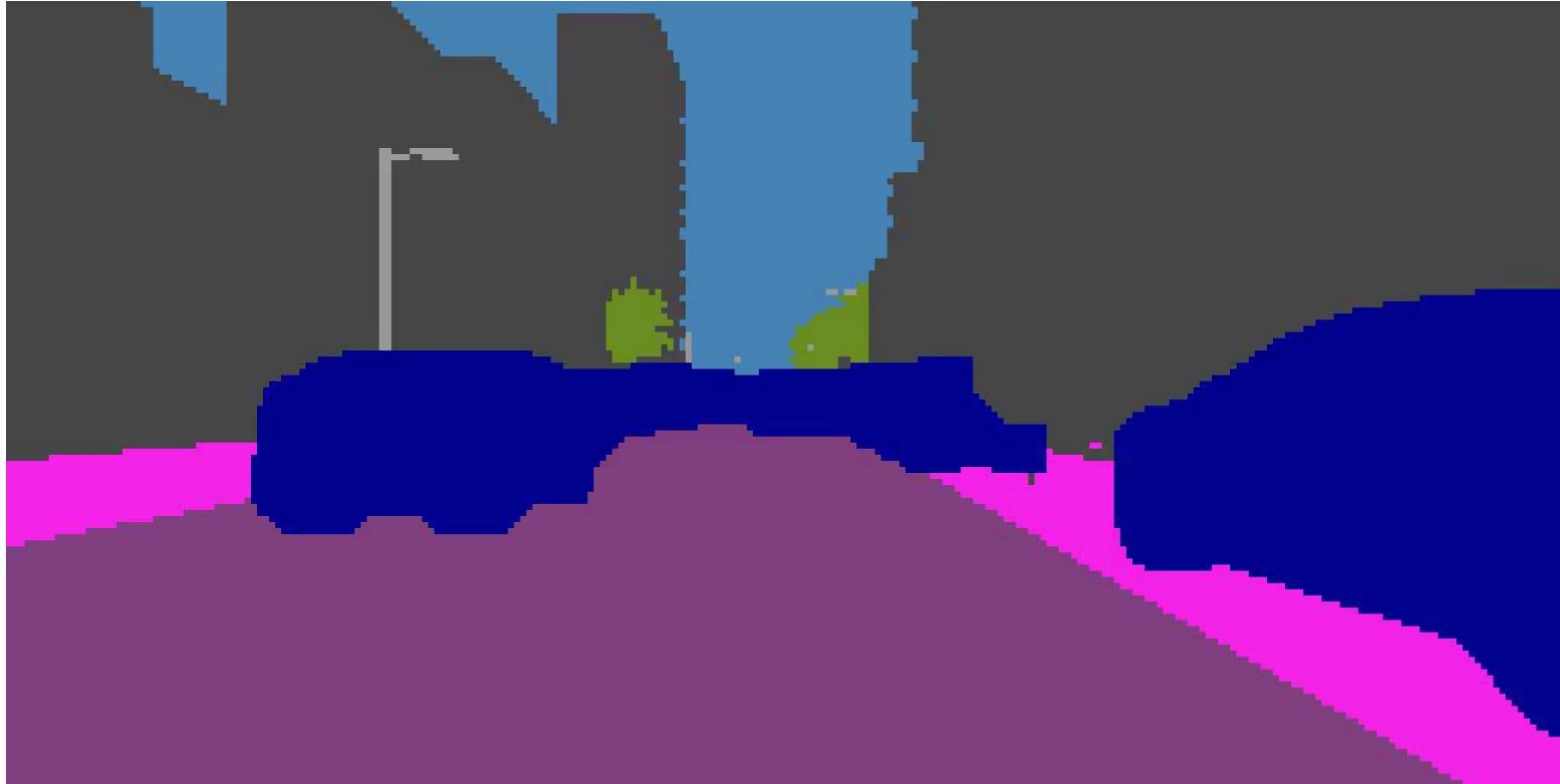
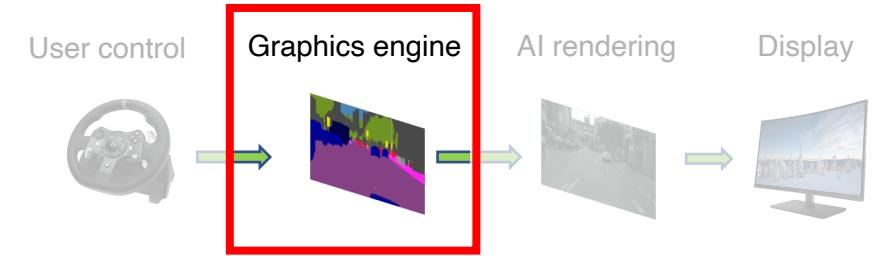
Graphics Engine: CARLA



Original CARLA Sequence

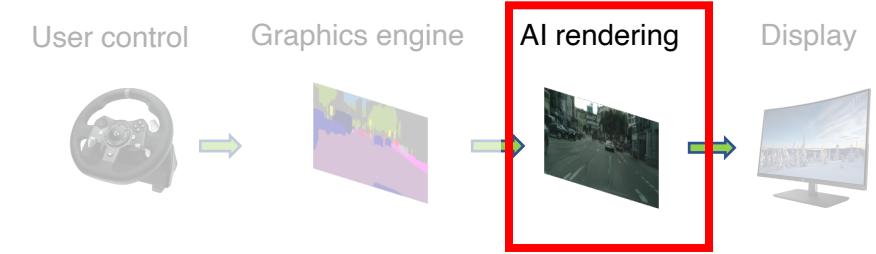
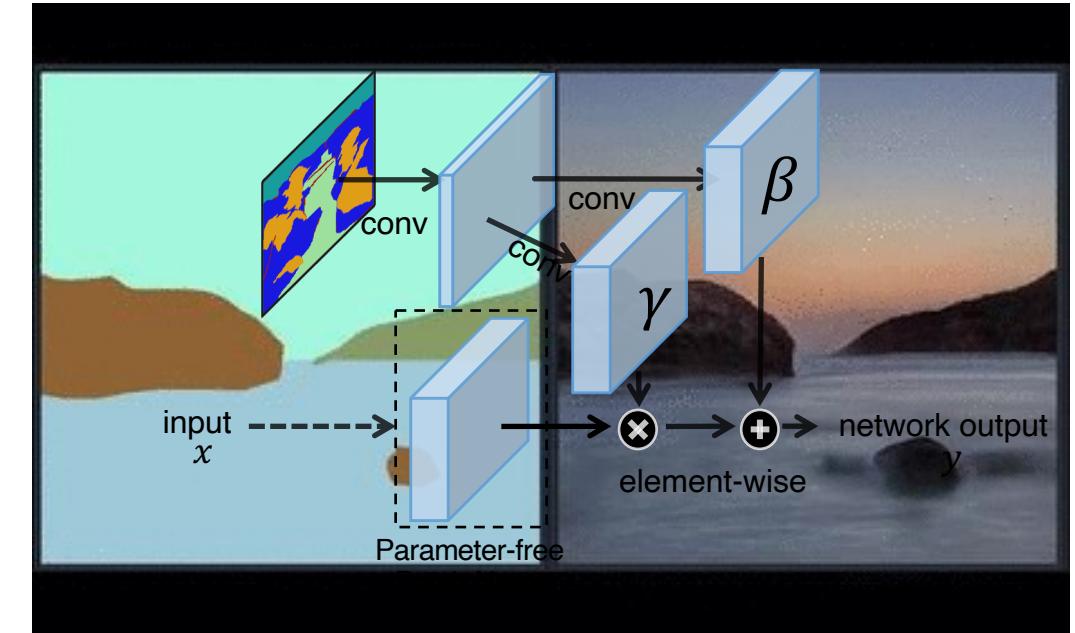
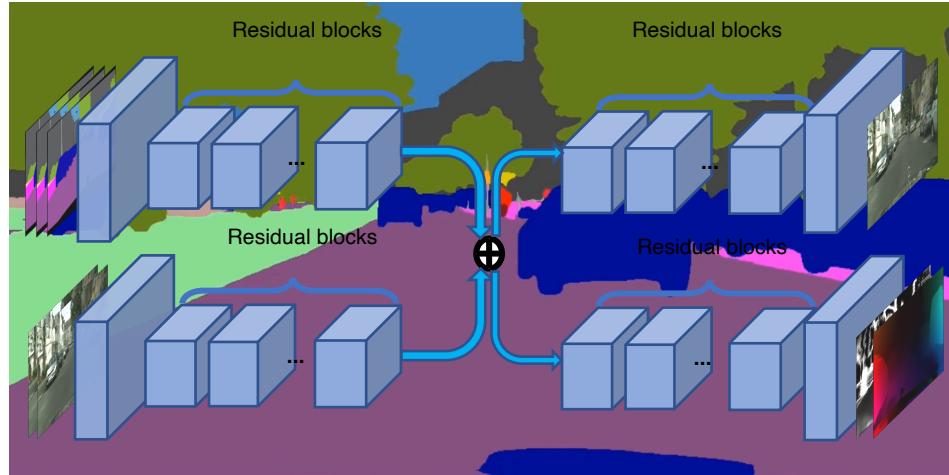


CARLA Semantic Maps

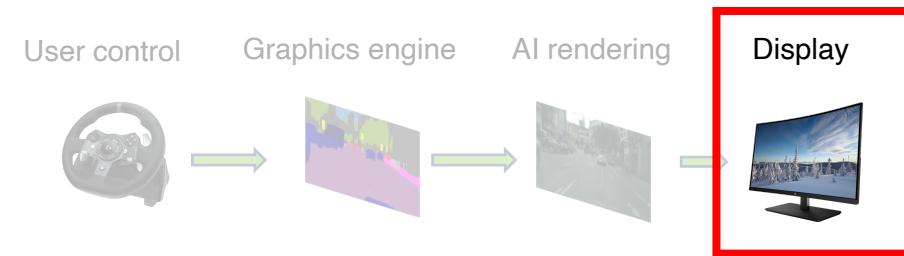


Methodology

- 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



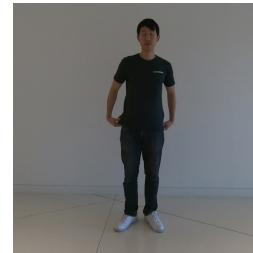
Adaptive Video Translation

Disadvantages of vid2vid

- Separate models for each dataset



model 1



model 2



model 3

- Generalizing to new persons requires

Collecting new data

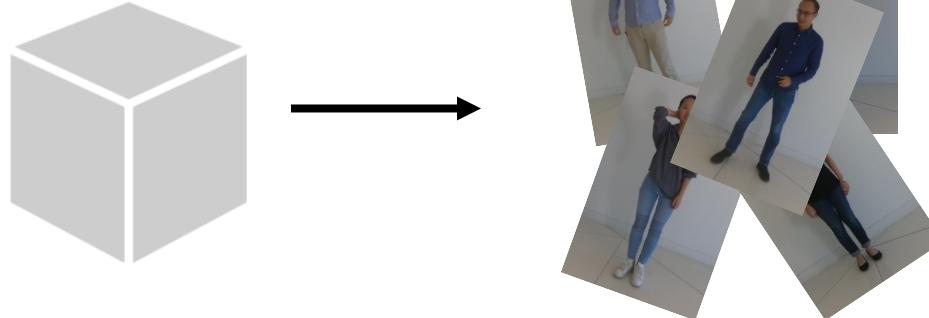


Training

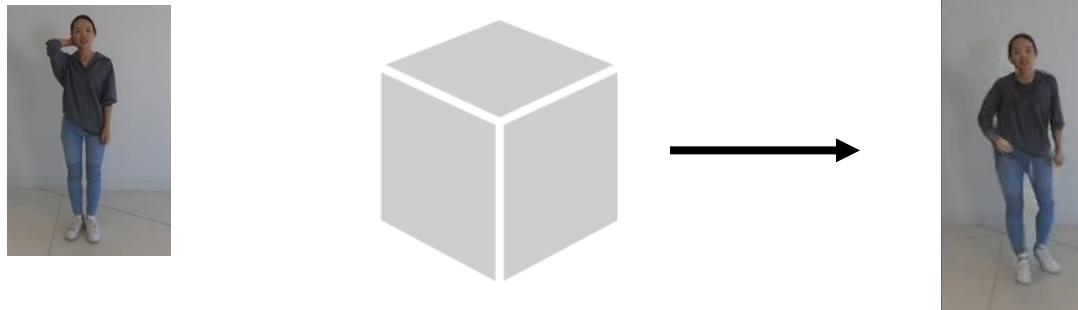


Wouldn't it be great if...

- One model for all

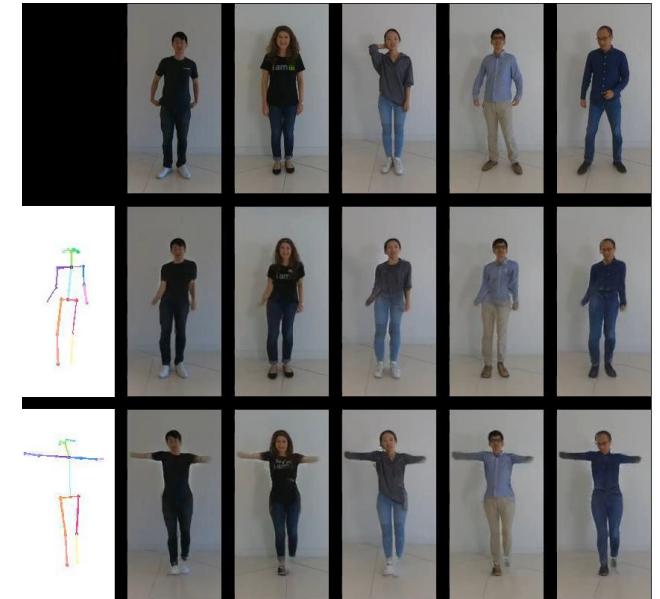


- Dynamically determine the style at run time
 - based on an *exemplar image*



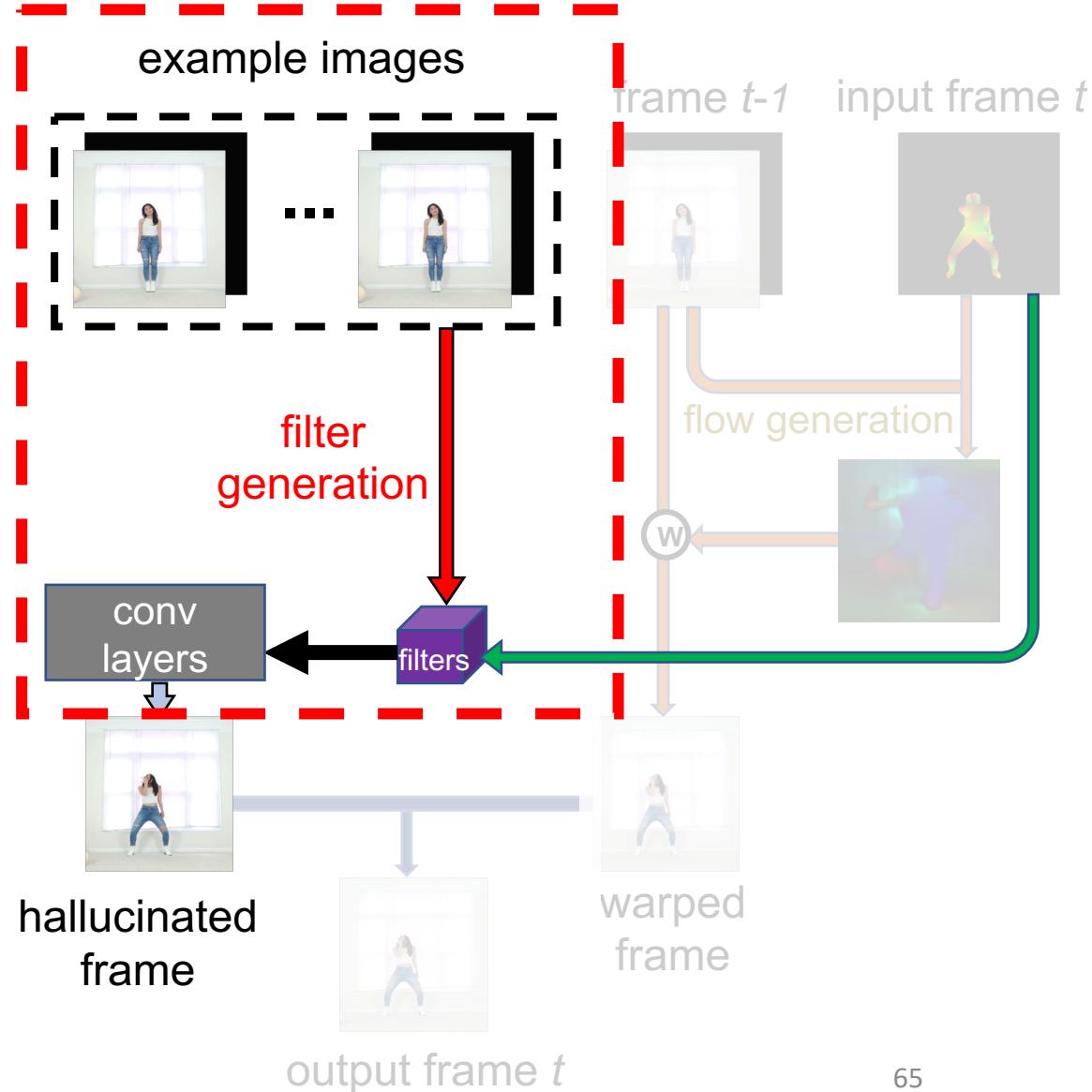
Adaptive Video-to-Video Translation

T.-C. Wang, M.-Y. Liu, A. Tao, G. Liu, J. Kautz, B. Catanzaro, “Few-shot Adaptive Video-to-Video Synthesis,” NeurIPS 2019.



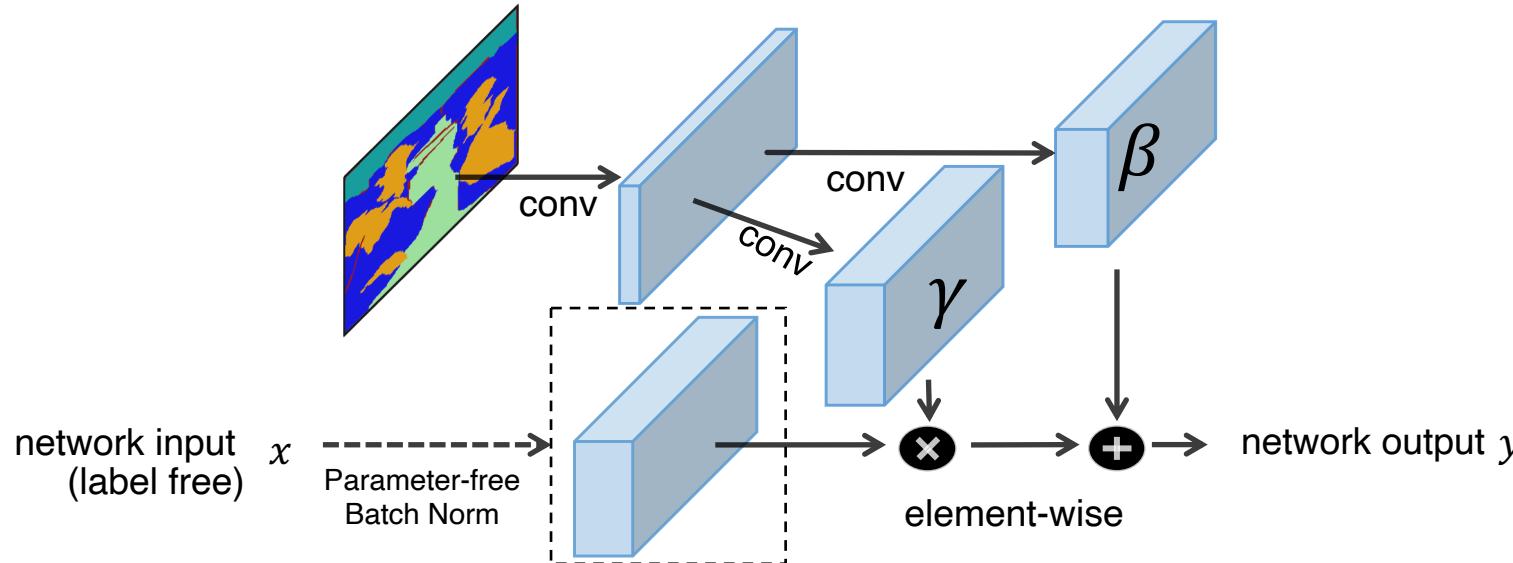
Adaptive vid2vid: overflow

- Original vid2vid
 - Output frame = Hallucinated frame + Warped frame
- Adaptive vid2vid
 - Hallucinated frames
 - generated based on example images
 - Using a filter generation scheme



Adaptive vid2vid

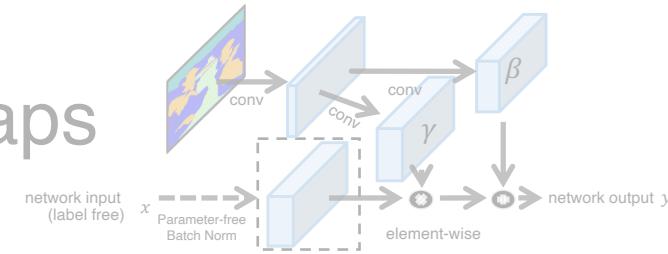
- Based on SPADE (GauGAN)
 - Prior work: ~~input semantics \rightarrow encoder-decoder \rightarrow output image~~
 - Instead: input semantics
 - \rightarrow ***spatially-varying*** normalization maps
 - \rightarrow used in every BatchNorm



$$y = \frac{x - \mu}{\sigma} \cdot \gamma + \beta$$

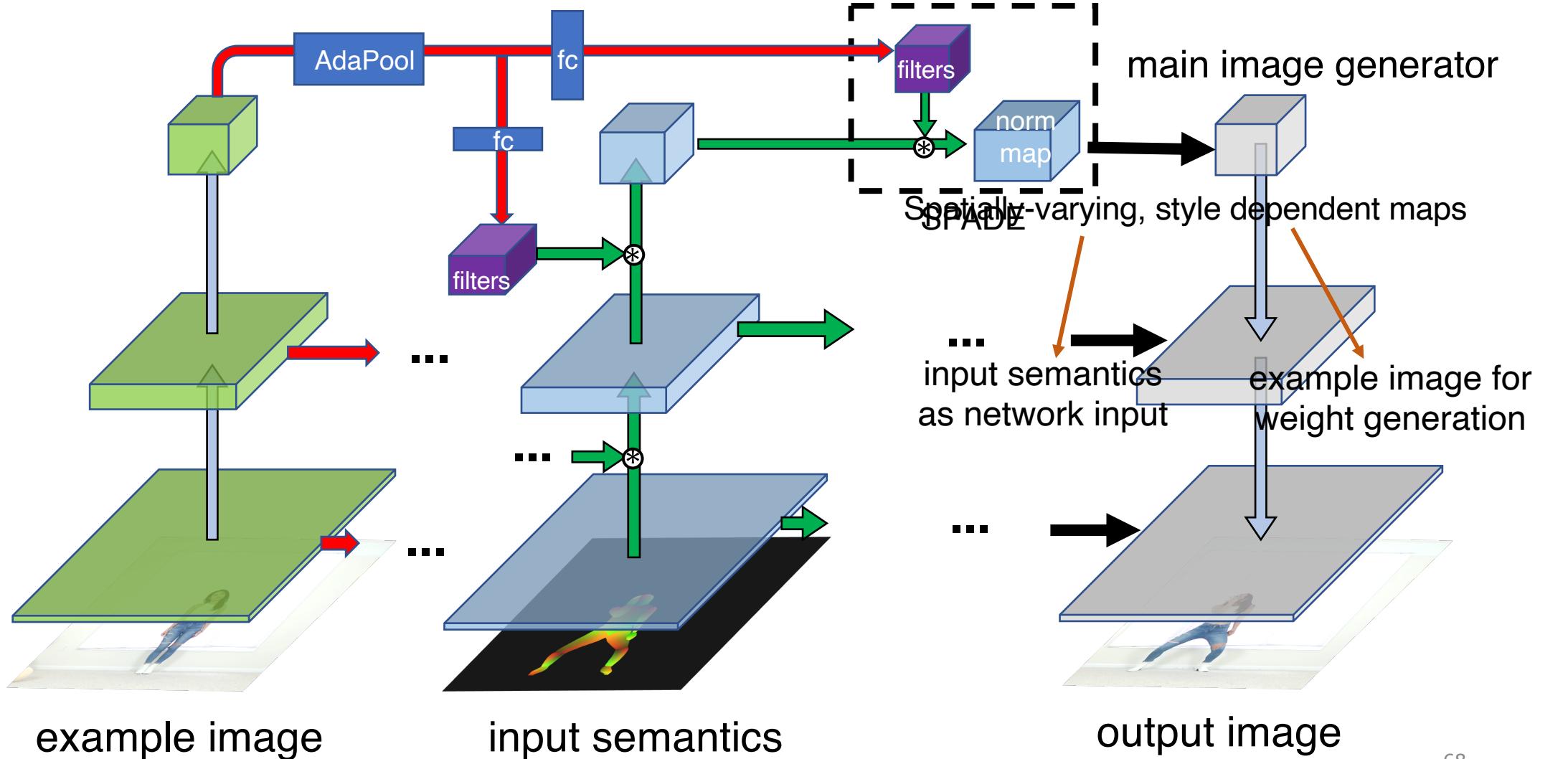
Adaptive vid2vid

- Based on SPADE (GauGAN)
 - Prior work: input semantics \rightarrow encoder-decoder \rightarrow output image
 - Instead: input semantics
 \rightarrow ***spatially-varying*** normalization maps
 \rightarrow used in every BatchNorm
- Given an additional exemplar image
 - Dynamically configure the ***network weights*** in SPADE
 - Generate ***spatially-varying, style-dependent*** normalization maps
 - Spatial info \leftarrow input semantics
 - Style info \leftarrow exemplar images



Dynamic Weight Generation

- filter generation
- normal convolution
- dynamic convolution
- normalization
- cube → convolution filters



Adaptive vid2vid: Training

- From a video
 - Randomly sample a clip
 - Randomly sample another reference frame(s)
- Make the network generate the clip
 - Based on the reference frame

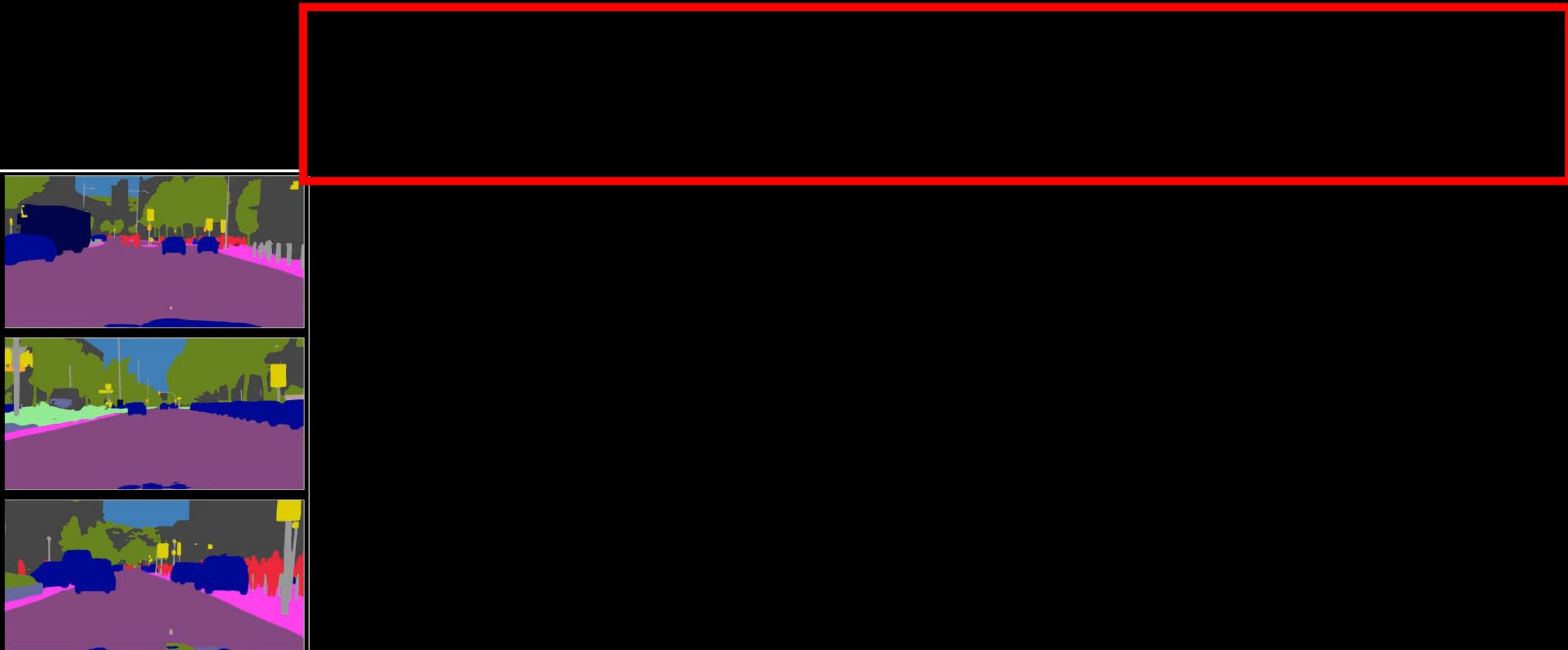
Adaptive vid2vid: Testing

- Given an example image
- Finetune on the example image
 - Network output should be the same as the example
 - Only finetune for a few iterations
- For faces: normalize keypoints
 - To the same as example image
 - To better preserve identity

Results

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

Street View Scenes



Input segmentations

Edges → Faces

Example images



Edges → Faces



Example image



Input videos



Extracted edges



Synthesized result

Poses → Body

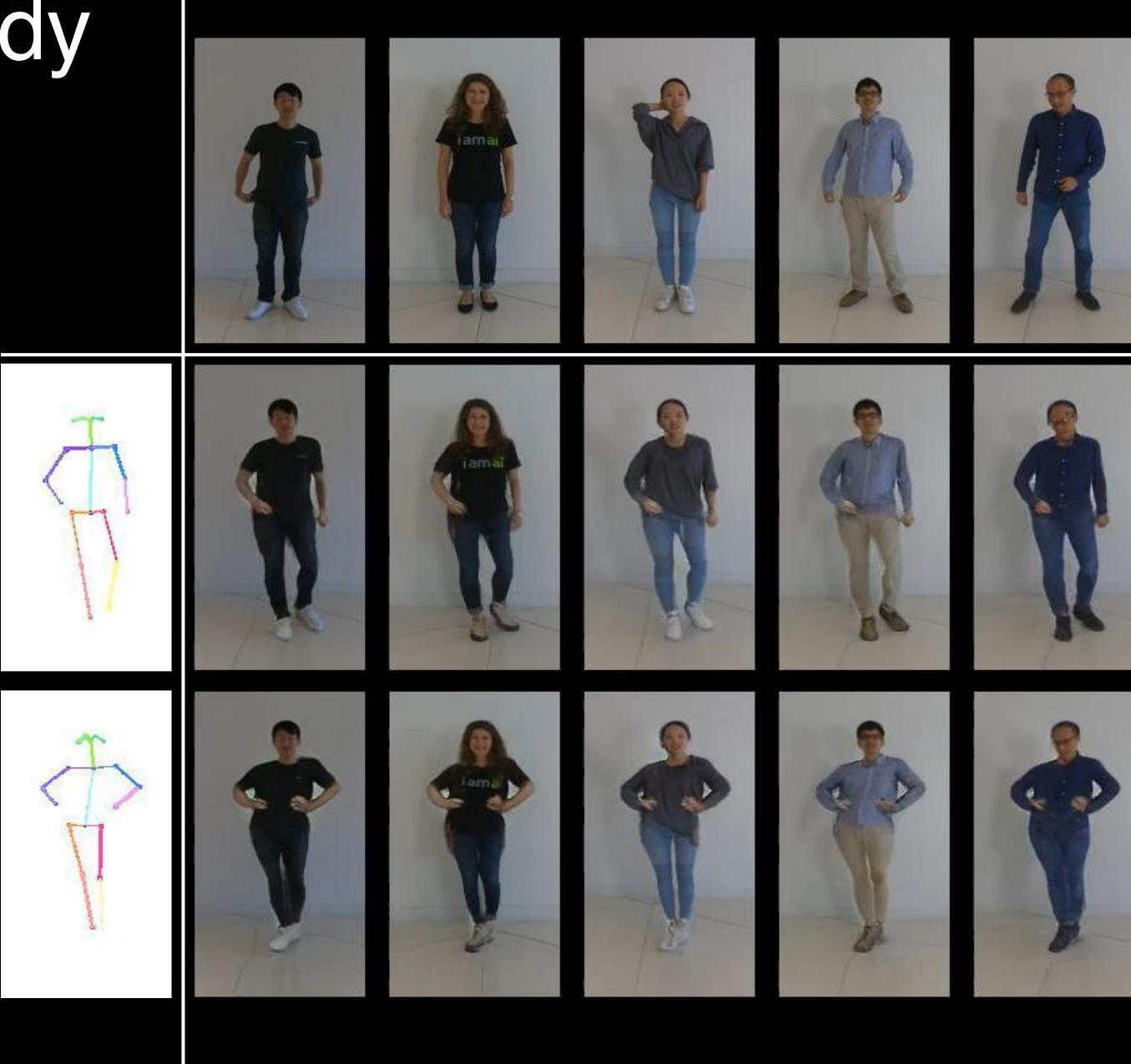
Input poses



Example
images

Synthesized
videos

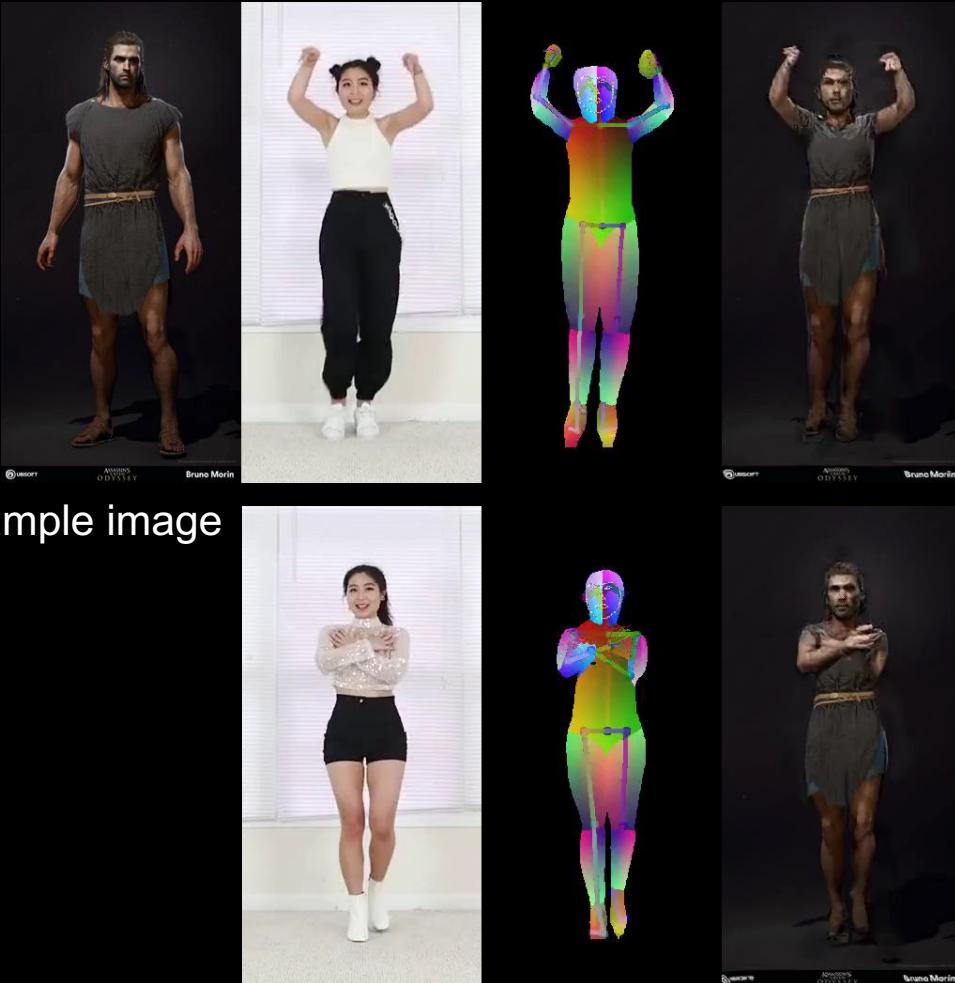
Poses → Body



Poses → Body



Poses → Body



Face Video Compression

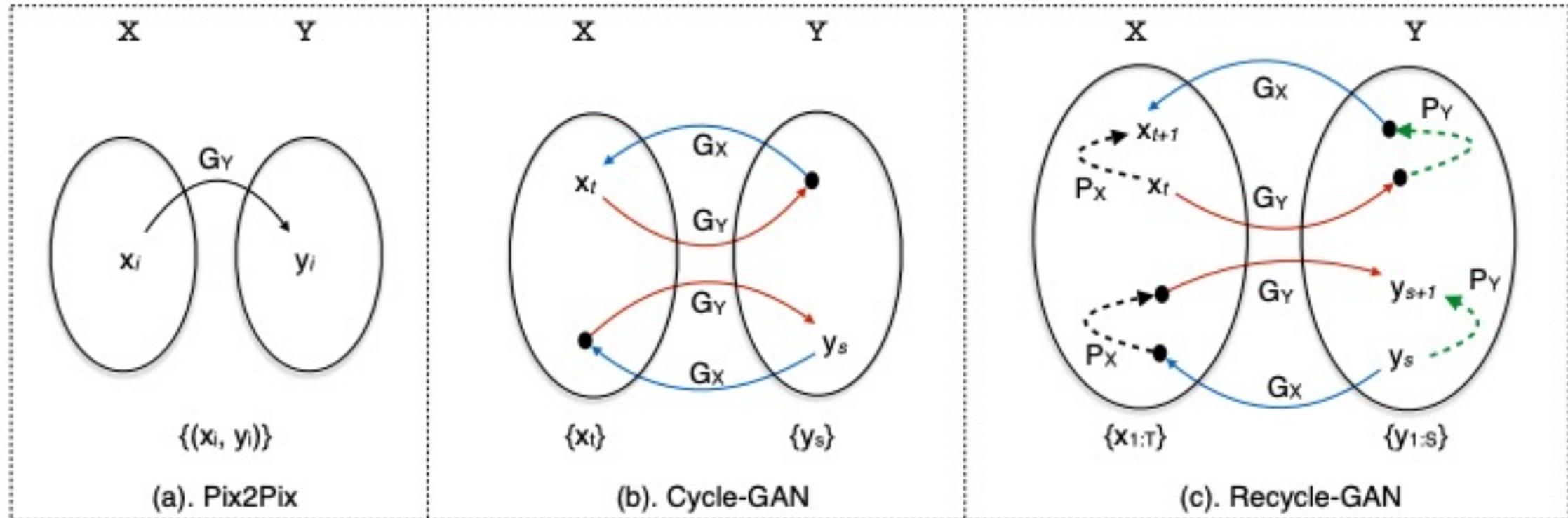


Paired vs. Unpaired

Unpaired Video Learning with RecycleGAN



Unpaired Video Learning with RecycleGAN



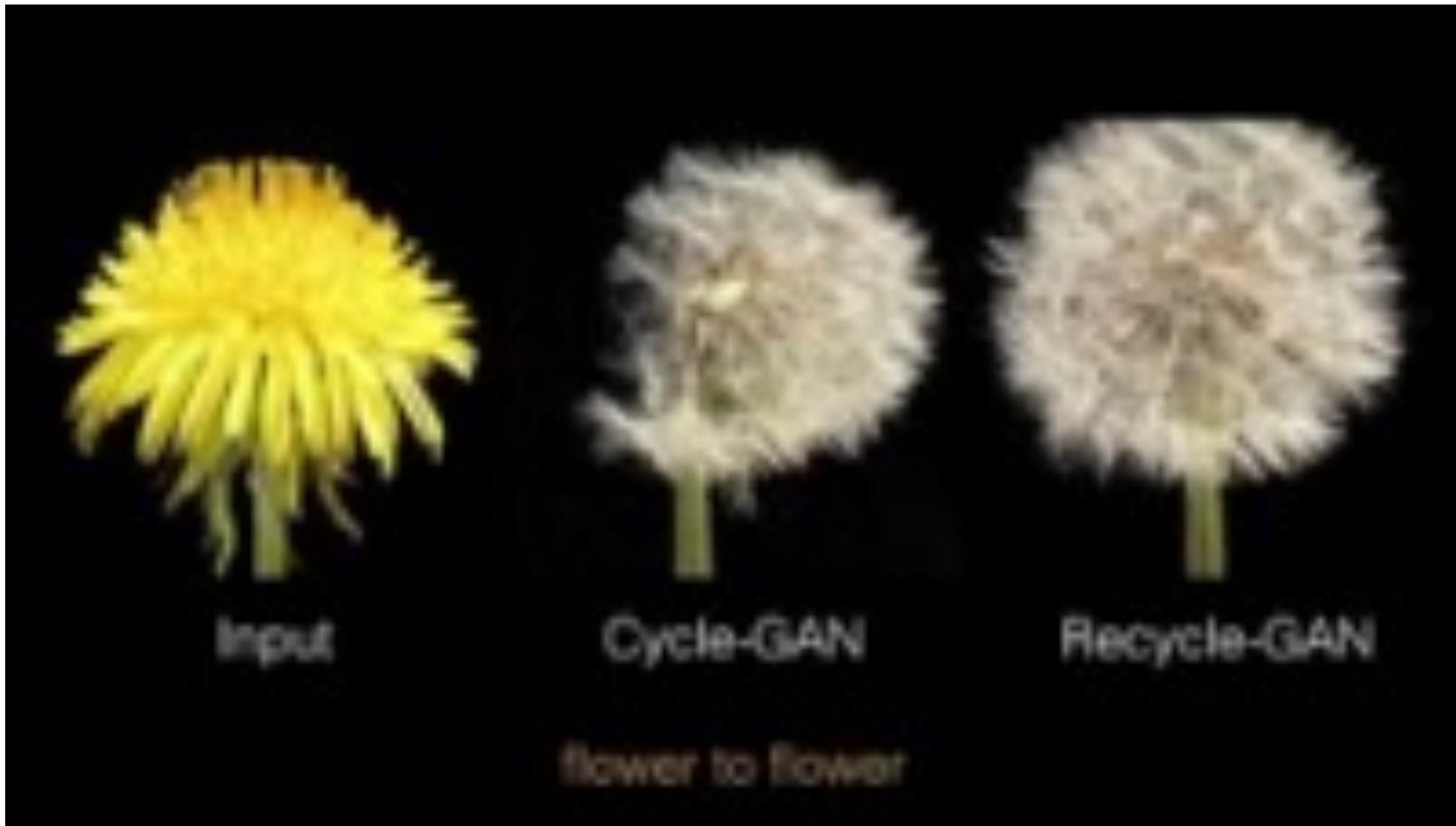
Prediction Loss
$$L_\tau(P_X) = \sum_t \|\mathbf{x}_{t+1} - P_X(\mathbf{x}_{1:t})\|^2,$$

Recycle Loss
$$L_r(G_X, G_Y, P_Y) = \sum_t \|\mathbf{x}_{t+1} - G_X(P_Y(G_Y(\mathbf{x}_{1:t})))\|^2,$$

Unpaired Video Learning with RecycleGAN



Unpaired Video Learning with RecycleGAN



Thank You!



16-726, Spring 2022

<https://learning-image-synthesis.github.io/sp22/>