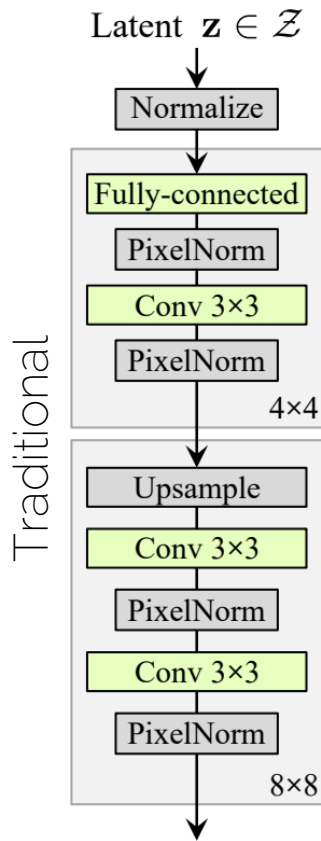
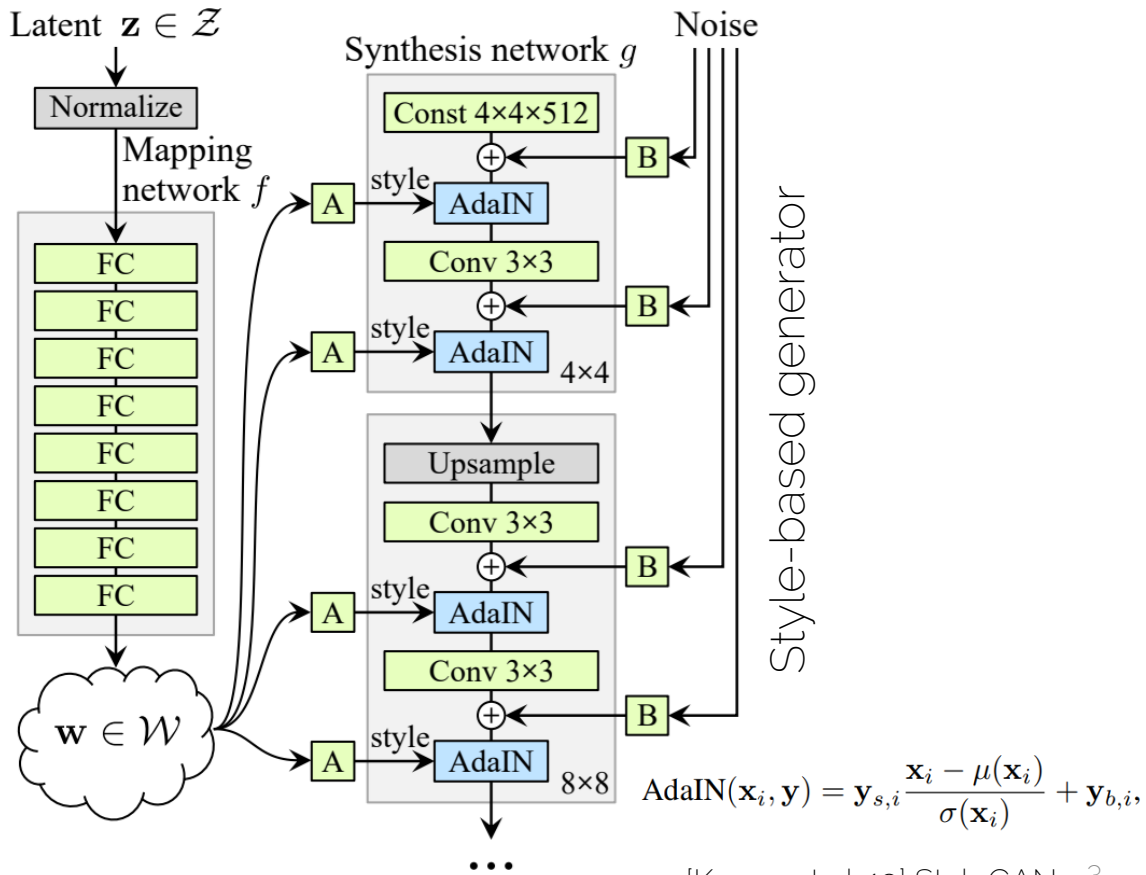
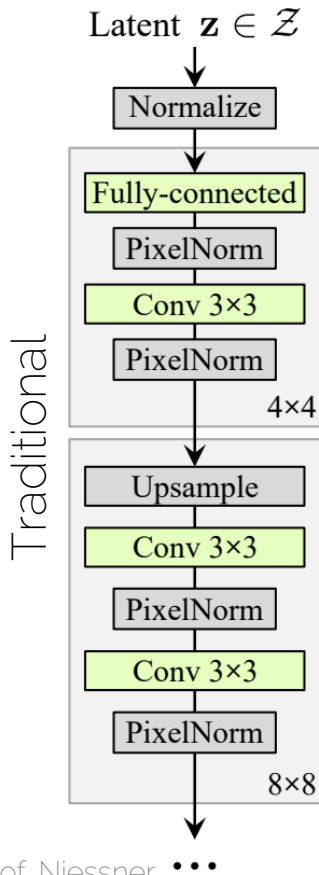


Style GAN

StyleGAN



StyleGAN



StyleGAN

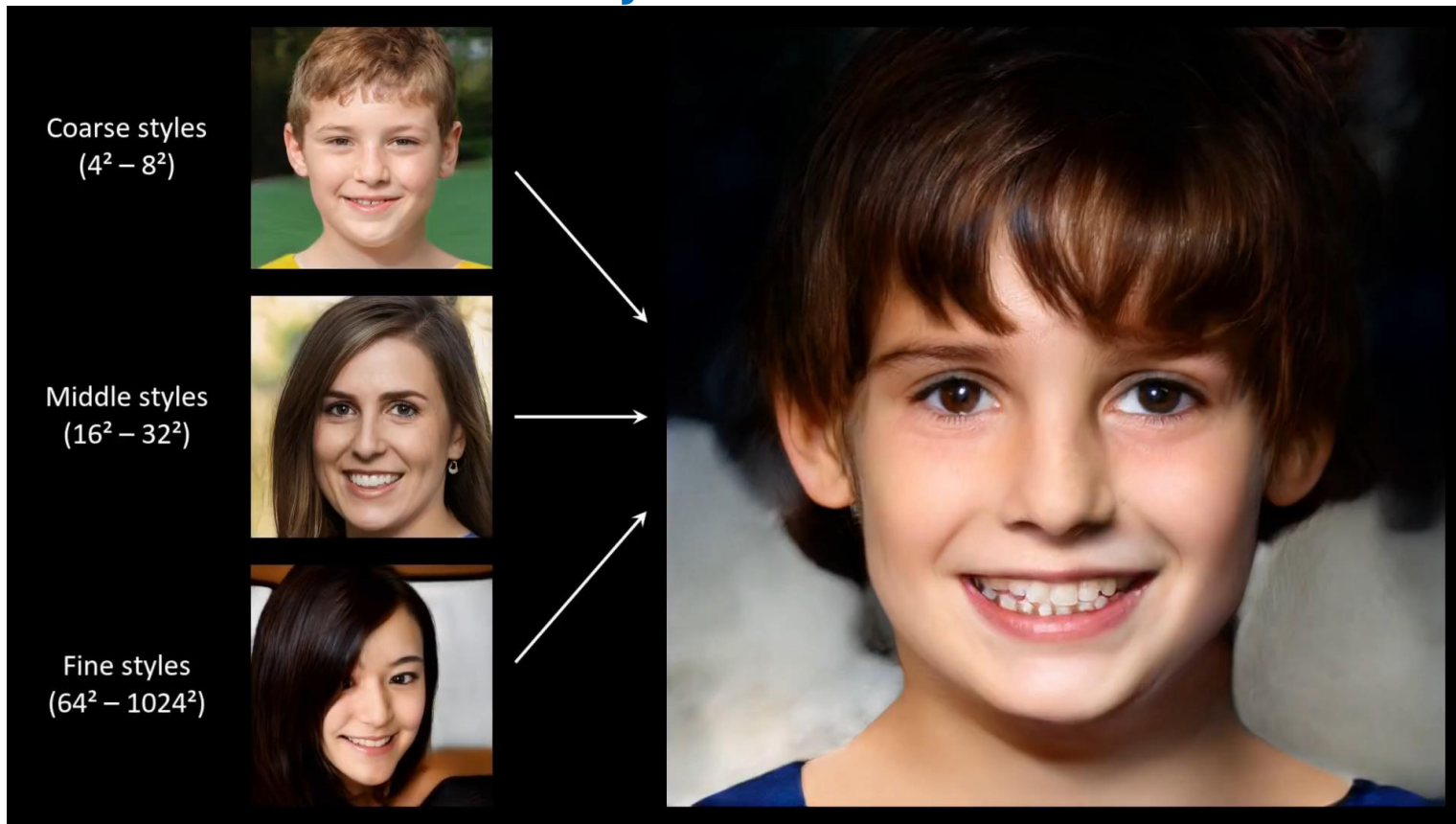
Method	CelebA-HQ	FFHQ
A Baseline Progressive GAN [30]	7.79	8.04
B + Tuning (incl. bilinear up/down)	6.11	5.25
C + Add mapping and styles	5.34	4.85
D + Remove traditional input	5.07	4.88
E + Add noise inputs	5.06	4.42
F + Mixing regularization	5.17	4.40

FID (Frechet inception distance) on 50k gen. images
-> Architecture is similar to Progressive Growing GAN

StyleGAN



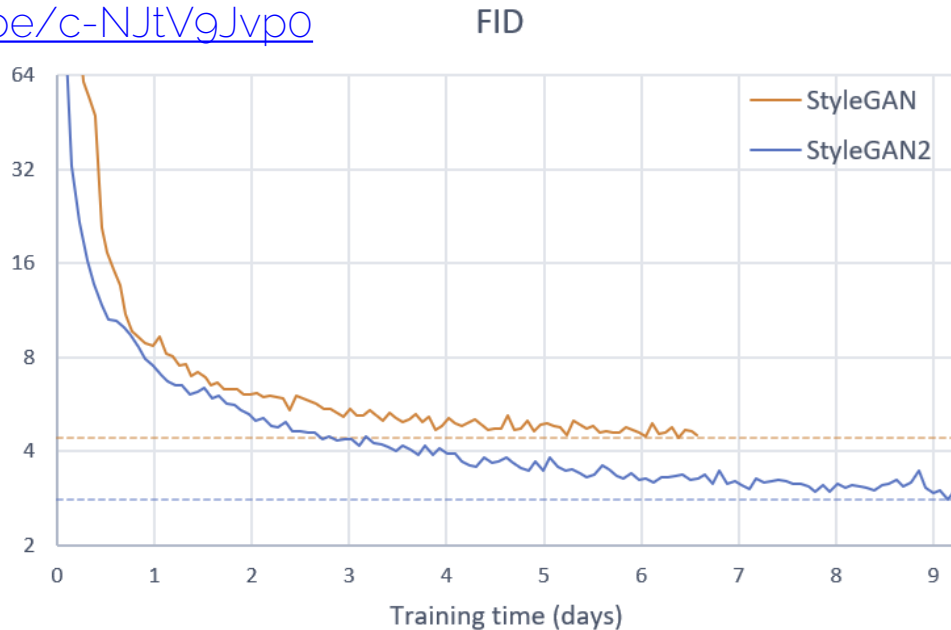
StyleGAN



StyleGAN2

Interesting analysis about design choices!

- <https://arxiv.org/pdf/1912.04958.pdf>
- <https://github.com/NVlabs/stylegan2>
- <https://youtu.be/c-NJtVgJvp0>



Autoregressive Models

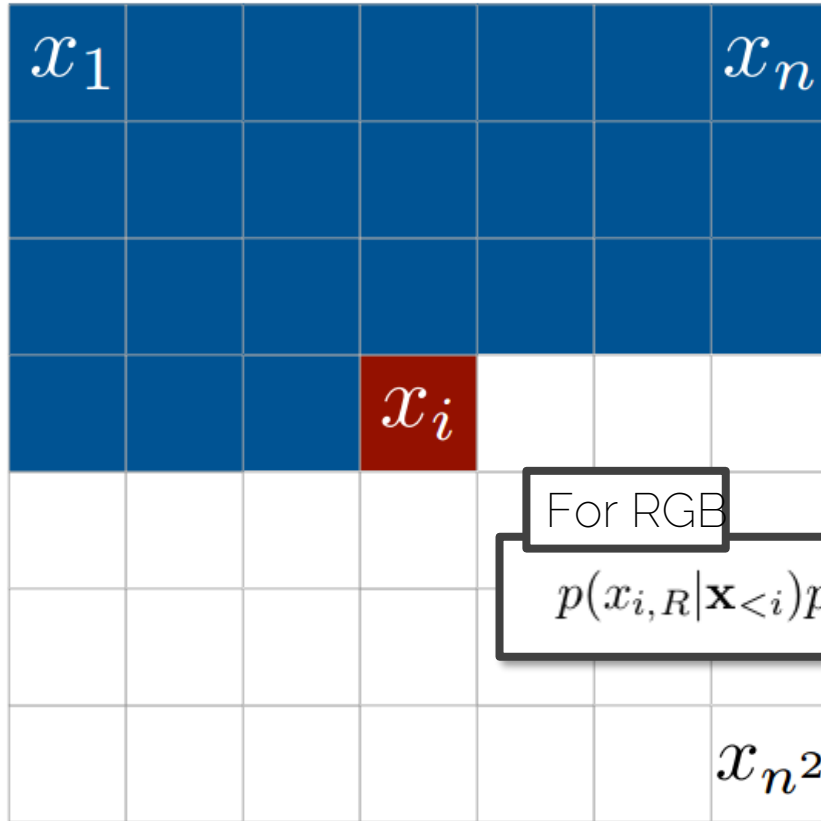
Autoregressive Models vs GANs

- GANs learn implicit data distribution
 - i.e., output are samples (distribution is in model)
- Autoregressive models learn an explicit distribution governed by a prior imposed by model structure
 - i.e., outputs are probabilities (e.g., softmax)

PixelRNN

- Goal: model distribution of natural images
- Interpret pixels of an image as product of conditional distributions
 - Modeling an image \rightarrow sequence problem
 - Predict one pixel at a time
 - Next pixel determined by all previously predicted pixels
 - Use a Recurrent Neural Network

PixelRNN

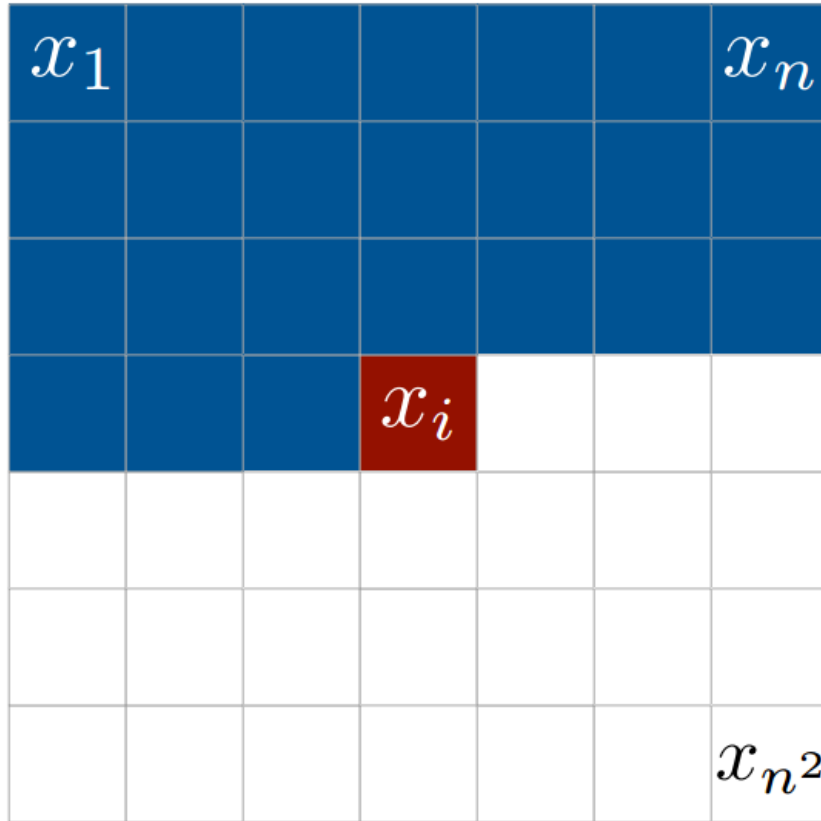


$$p(\mathbf{x}) = \prod_{i=1}^{n^2} p(x_i | x_1, \dots, x_{i-1})$$

For RGB

$$p(x_{i,R} | \mathbf{x}_{<i}) p(x_{i,G} | \mathbf{x}_{<i}, x_{i,R}) p(x_{i,B} | \mathbf{x}_{<i}, x_{i,R}, x_{i,G})$$

PixelRNN

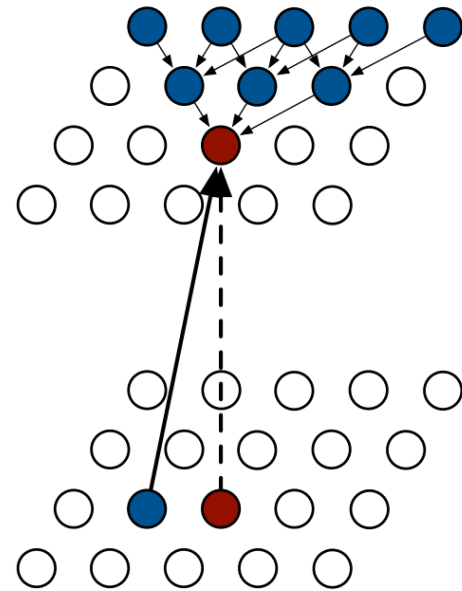


$$p(\mathbf{x}) = \prod_{i=1}^{n^2} p(x_i | x_1, \dots, x_{i-1})$$

$x_i \in [0, 255]$
→ 256-way softmax

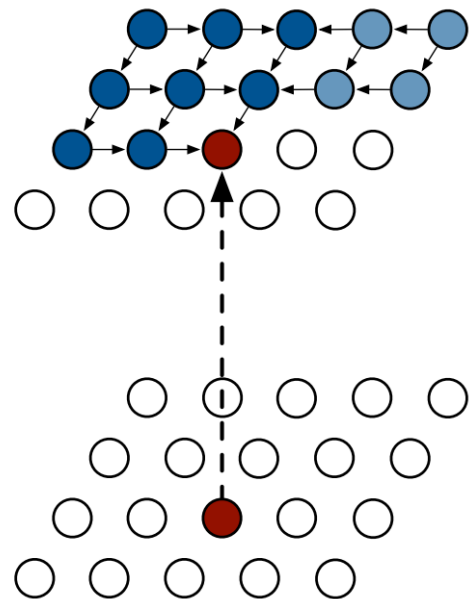
PixelRNN

- Row LSTM model architecture
- Image processed row by row
- Hidden state of pixel depends on the 3 pixels above it
 - Can compute pixels in row in parallel
- Incomplete context for each pixel



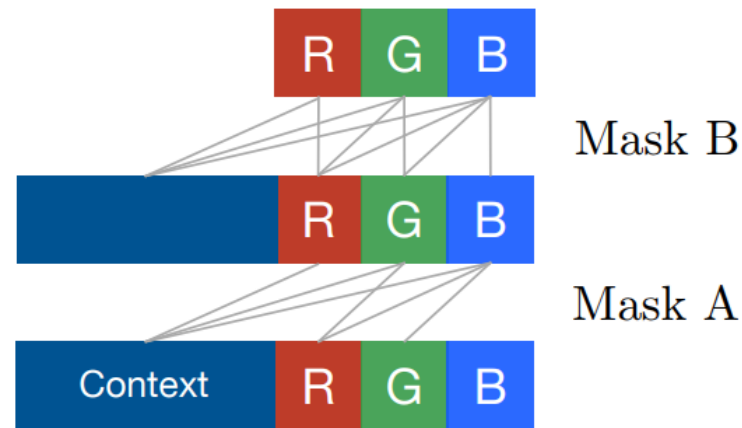
PixelRNN

- Diagonal BiLSTM model architecture
- Solve incomplete context problem
- Hidden state of pixel $p_{i,j}$ depends on $p_{i,j-1}$ and $p_{i-1,j}$
- Image processed by diagonals



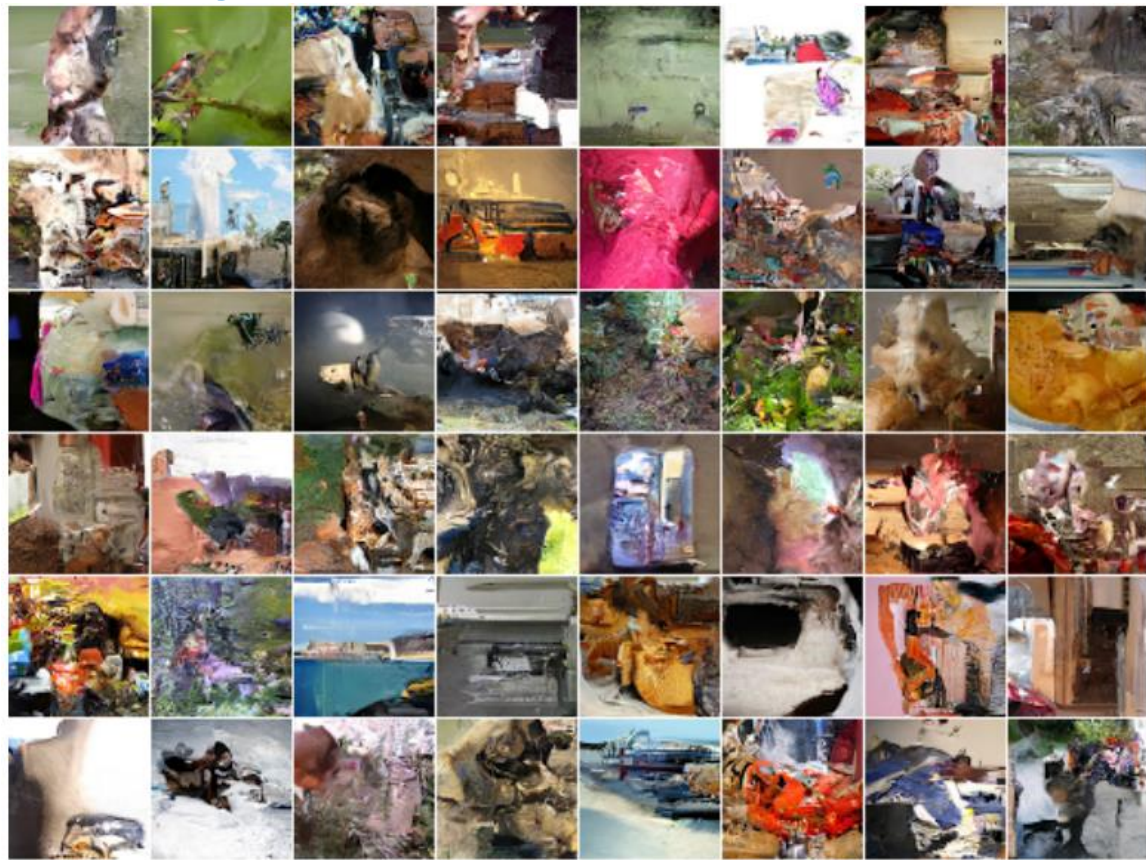
PixelRNN

- Masked Convolutions
- Only previously predicted values can be used as context
- Mask A: restrict context during 1st conv
- Mask B: subsequent convs
- Masking by zeroing out values



PixelRNN

- Generated 64x64 images, trained on ImageNet



PixelCNN

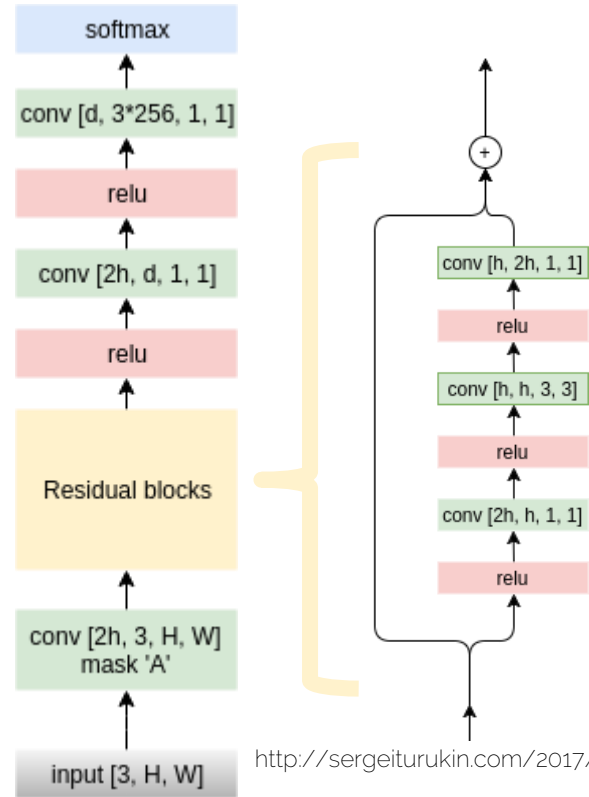
- Row and Diagonal LSTM layers have potentially unbounded dependency range within the receptive field
 - Can be very computationally costly
- PixelCNN:
 - standard convs capture a bounded receptive field
 - All pixel features can be computed at once (during training)

PixelCNN

- Model preserves spatial dimensions
- Masked convolutions to avoid seeing future context

1	1	1	1	1
1	1	1	1	1
1	1	0	0	0
0	0	0	0	0
0	0	0	0	0

Mask A



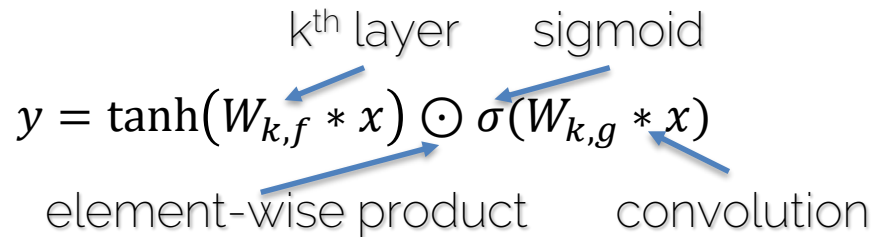
Gated PixelCNN

- Gated blocks
- Imitate multiplicative complexity of PixelRNNs to reduce performance gap between PixelCNN and PixelRNN
- Replace ReLU with gated block of sigmoid, tanh

$$y = \tanh(W_{k,f} * x) \odot \sigma(W_{k,g} * x)$$

kth layer sigmoid

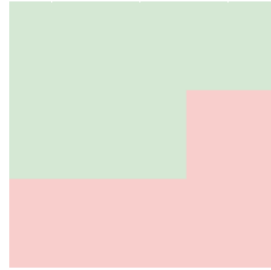
element-wise product convolution

The diagram shows the equation $y = \tanh(W_{k,f} * x) \odot \sigma(W_{k,g} * x)$. Above the equation, 'kth layer' has arrows pointing to $W_{k,f}$ and $W_{k,g}$, and 'sigmoid' has an arrow pointing to σ . Below the equation, 'element-wise product' has an arrow pointing to \odot , and 'convolution' has an arrow pointing to $*$.

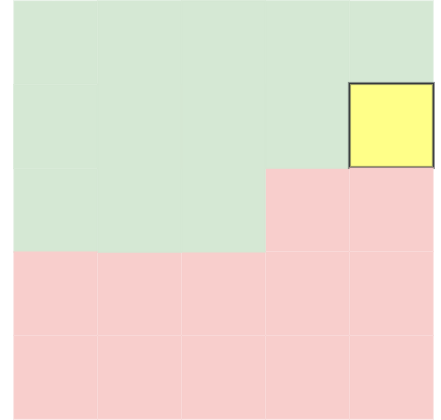
PixelCNN Blind Spot

1	1	1	1	1
1	1	1	1	1
1	1	0	0	0
0	0	0	0	0
0	0	0	0	0

5x5 image / 3x3 conv



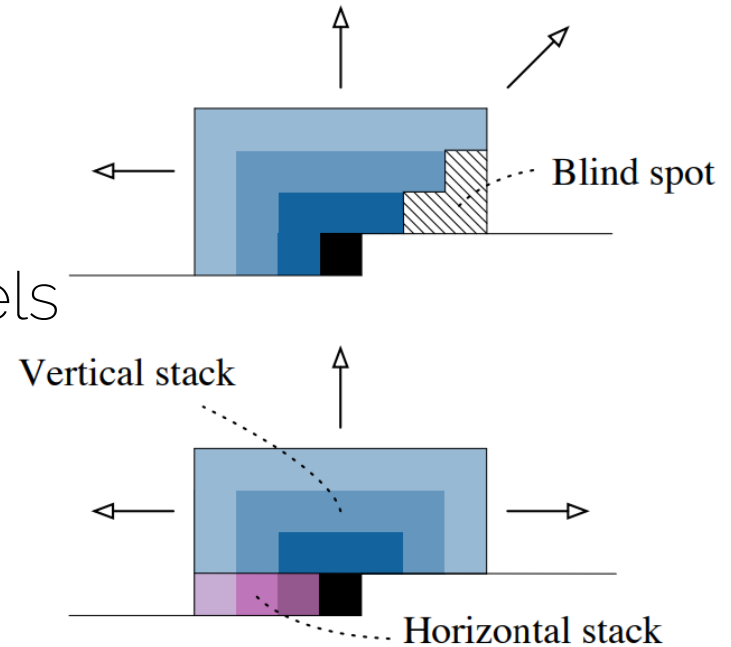
Receptive Field



Unseen context

PixelCNN: Eliminating Blind Spot


- Split convolution to two stacks
- Horizontal stack conditions on current row
- Vertical stack conditions on pixels above



Conditional PixelCNN

- Conditional image generation
- E.g., condition on semantic class, text description

latent vector to be conditioned on

$$y = \tanh(W_{k,f} * x + V_{k,f}^T h) \odot \sigma(W_{k,g} * x + V_{k,g}^T h)$$


Conditional PixelCNN



Coral Reef



Sorrel horse

Autoregressive Models vs GANs

- Advantages of autoregressive:
 - Explicitly model probability densities
 - More stable training
 - Can be applied to both discrete and continuous data
- Advantages of GANs:
 - Have been empirically demonstrated to produce higher quality images
 - Faster to train

Autoregressive Models

- State of the art is pretty impressive 😊

Vector Quantized Variational AutoEncoder



Generating Diverse High-Fidelity Images with VQ-VAE-2

<https://arxiv.org/pdf/1906.00446.pdf> [Razavi et al. 19]

Generative Models on Videos

GANs on Videos

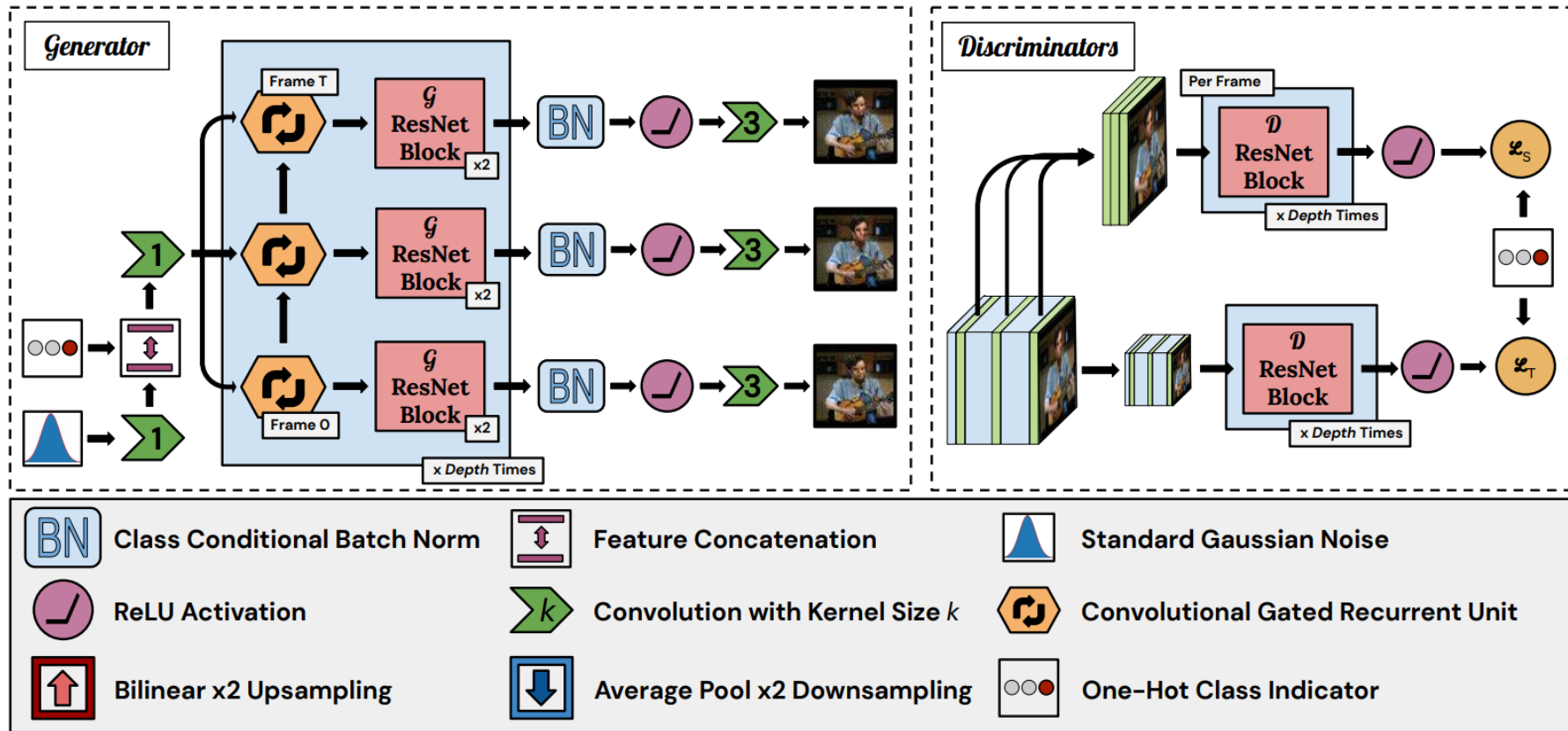
Two options

- Single random variable z seeds entire video (all frames)
 - Very high dimensional output
 - How to do for variable length?
 - Future frames deterministic given past
- Random variable z for each frame of the video
 - Need conditioning for future from the past
 - How to get combination of past frames + random vectors during training

General issues

- Temporal coherency
- Drift over time (many models collapse to mean image)

GANs on Videos: DVD-GAN



GANs on Videos: DVD-GAN



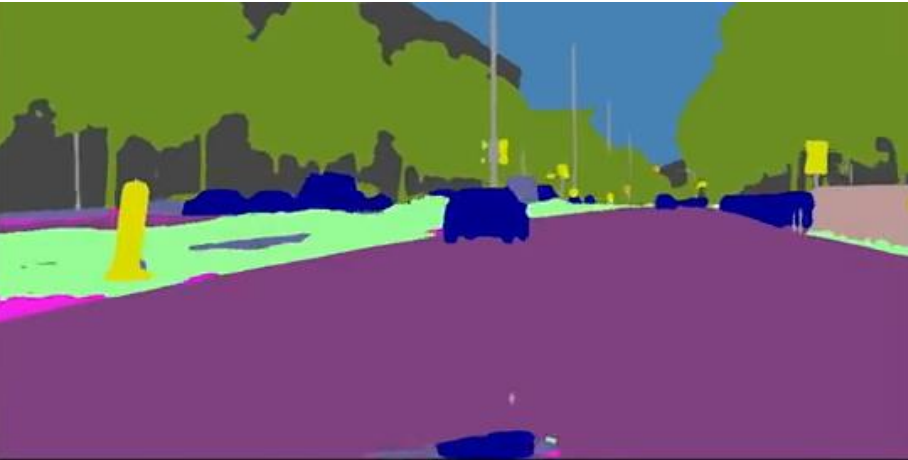
Time

GANs on Videos: DVD-GAN

- Trained on Kinetics-600 dataset
 - 256×256 , 128×128 , and 64×64
 - Lengths of up 48 frames
- > This is state of the art!
- > Videos from scratch still incredibly challenging

Conditional GANs on Videos

- Challenge:
 - Each frame is high quality, but temporally inconsistent



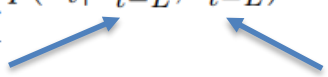
Labels



pix2pixHD

Video-to-Video Synthesis

- Sequential Generator:

$$p(\tilde{\mathbf{x}}_1^T | \mathbf{s}_1^T) = \prod_{t=1}^T p(\tilde{\mathbf{x}}_t | \tilde{\mathbf{x}}_{t-L}^{t-1}, \mathbf{s}_{t-L}^t).$$


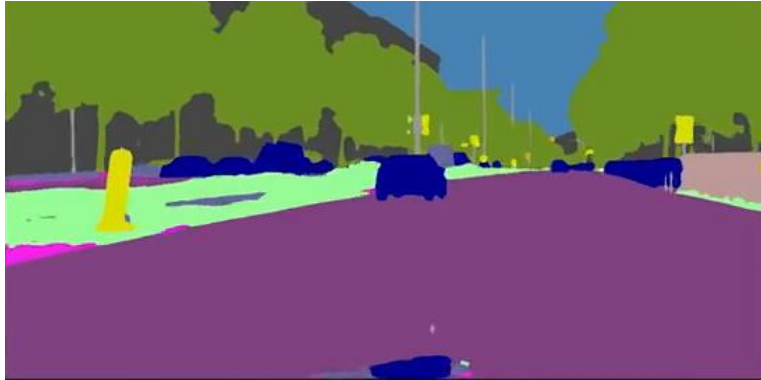
past L generated frames past L source frames
(set L = 2)

- Conditional Image Discriminator D_I (is it real image)
- Conditional Video Discriminator D_V (temp. consistency via flow)

Full Learning Objective:

$$\min_F \left(\max_{D_I} \mathcal{L}_I(F, D_I) + \max_{D_V} \mathcal{L}_V(F, D_V) \right) + \lambda_W \mathcal{L}_W(F),$$

Video-to-Video Synthesis



Labels



pix2pixHD

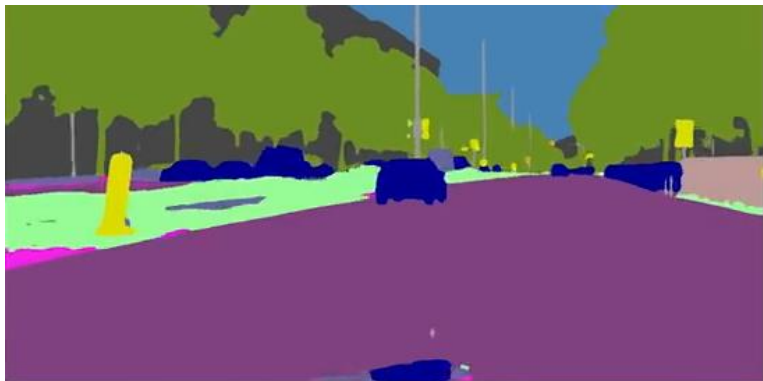


COVST



Ours

Video-to-Video Synthesis



Labels



pix2pixHD



COVST



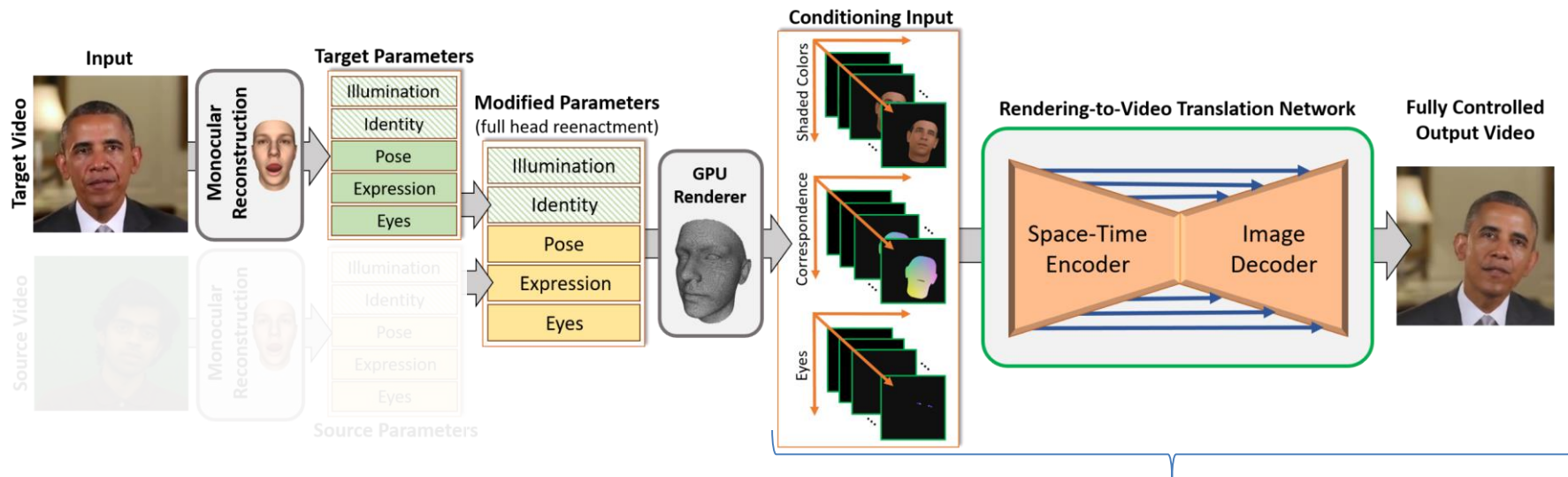
Ours

Video-to-Video Synthesis

- Key ideas:
 - Separate discriminator for temporal parts
 - In this case based on optical flow
 - Consider recent history of prev. frames
 - Train all of it jointly

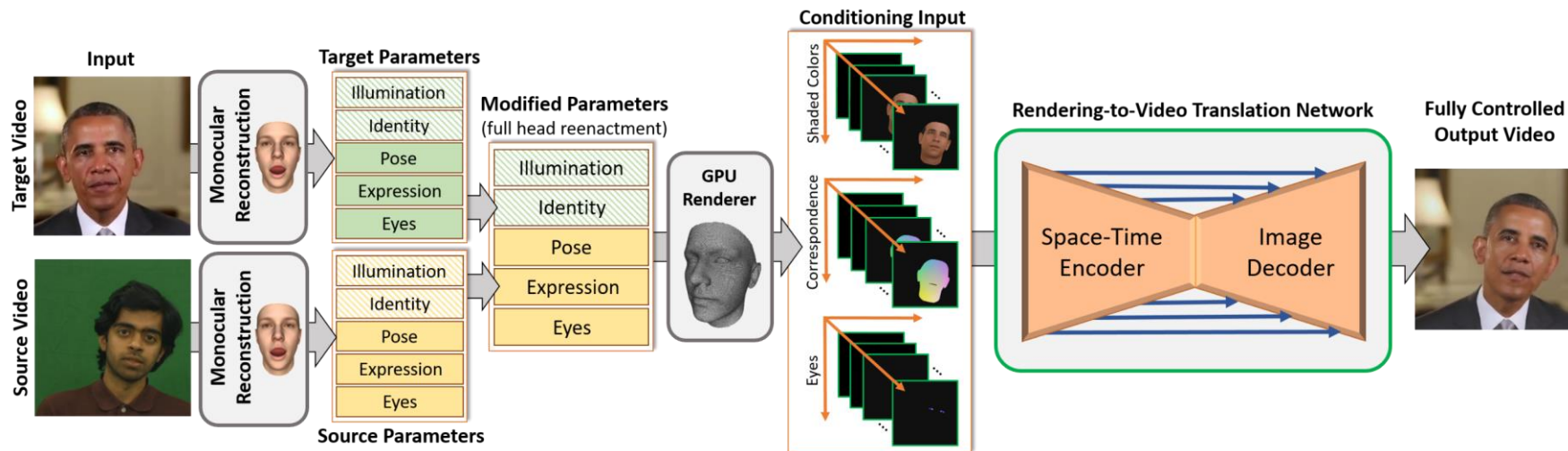
Deep Video Portraits

Deep Video Portraits

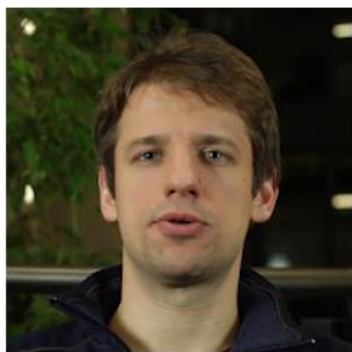


Similar to “Image-to-Image Translation” (Pix2Pix) [Isola et al.]

Deep Video Portraits



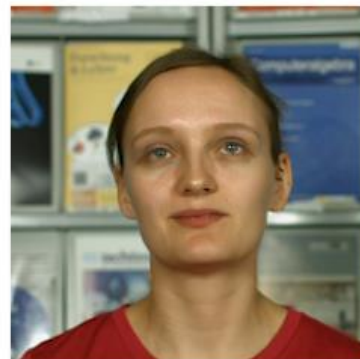
Deep Video Portraits



Source Sequence



Conditioning Images

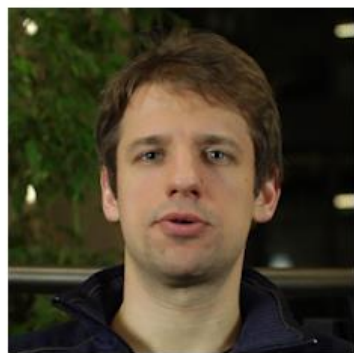


Result

Neural Network converts synthetic data to realistic video



Deep Video Portraits



Source Sequence

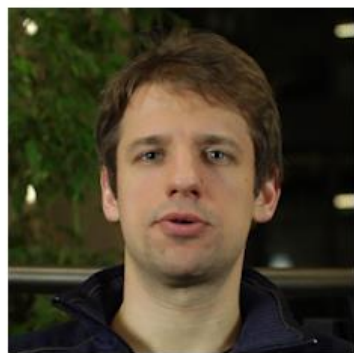


Conditioning Images



Result

Deep Video Portraits



Source Sequence



Conditioning Images



Result

Deep Video Portraits



Deep Video Portraits



Interactive Video Editing

2x speed

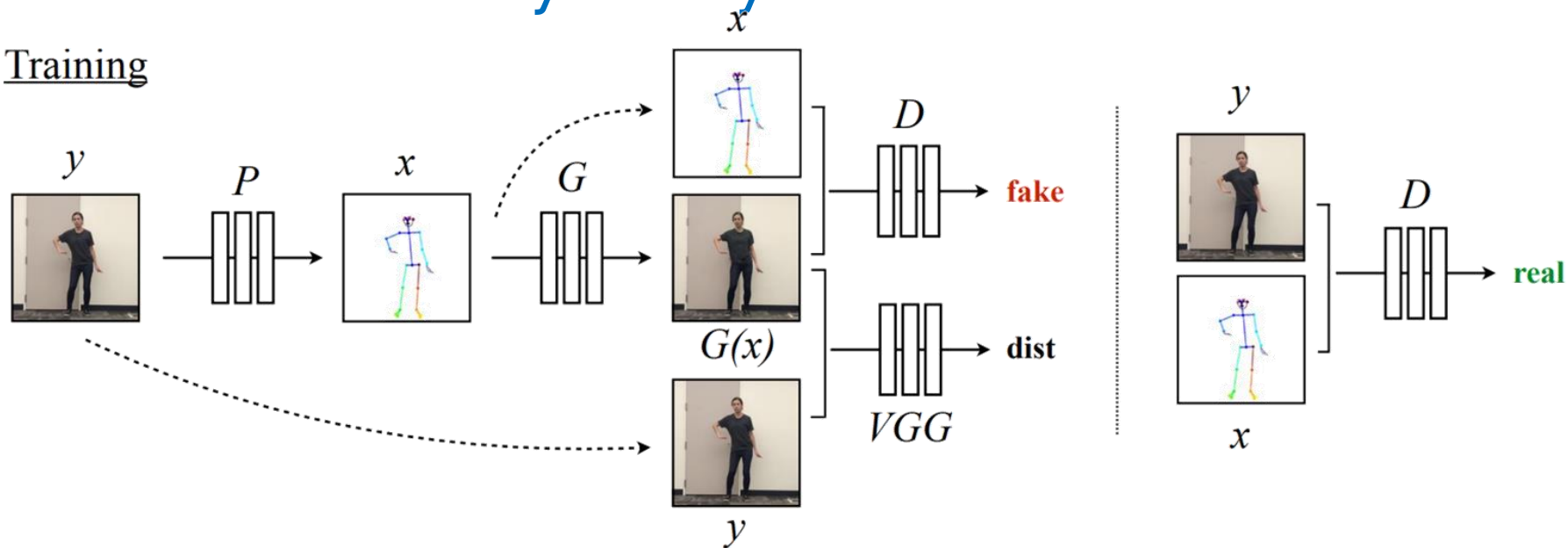
Deep Video Portraits: Insights

- Synthetic data for tracking is great anchor / stabilizer
- Overfitting on small datasets works pretty well
- Need to stay within training set w.r.t. motions
- No real learning; essentially, optimizing the problem with SGD
 - > should be pretty interesting for future directions

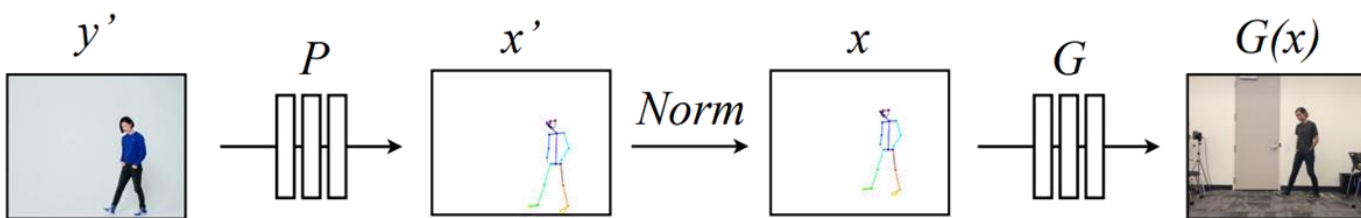
Everybody Dance Now

Everybody Dance Now

Training



Transfer



Everybody Dance Now

Source Subject



Everybody Dance Now

- cGANs work with different input
- Requires consistent input i.e., accurate tracking
- Network has no explicit 3D notion



Everybody Dance Now: Insights

- Conditioning via tracking seems promising!
 - Tracking quality translates to resulting image quality
 - Tracking human skeletons is less developed than faces
 - Temporally it's not stable... (e.g., OpenPose etc.)
 - Fun fact, there were like 4 papers with a similar idea that appeared around the same time...

Next Lectures

- Next Lectures:
 - Neural Rendering
 - 3D Deep Learning
- Keep working on the projects!

See you next week 😊