

Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, Bryan Catanzaro

GENERATIVE ADVERSARIAL NETWORKS

Unconditional GANs



Image credit: Celebrity dataset, Jensen Huang, Founder and CEO of NVIDIA, Ian Goodfellow, Father of GANs.

After training for a while using NVIDIA DGX1 machines

Fun sampling time begin

$$z_1, z_2, z_3, \dots$$
 Generator

Image credit: NVIDIA StyleGAN



CONDITIONAL GANS Allow user more control on the sampling process



Sampling $z \sim Z, y \sim Y$ (testing) $f \sim Z, y \sim Y$ output style Given info (e.g. image, text)

SKETCH-CONDITIONAL GANS



IMAGE-CONDITIONAL GANS







MASK-CONDITIONAL GANS

Semantic Image Synthesis



LIVE DEMO

- I need to get an RTX Ready Laptop (<u>https://www.nvidia.com/en-us/geforce/gaming-laptops/20-series/</u>)
- It is running live in GTC
- Will be online for everyone to try out in NVIDIA AI Playground website (<u>https://www.nvidia.com/en-</u> <u>us/research/ai-playground/</u>)









0 x =x =0 β 0 same output! affine transform de-normalization normalization

Batch Norm (loffe et al. 2015)

removes label information



- Do not feed the label map directly to network
- Use the label map to generate normalization layers instead



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



SPADE SPatially Adaptive DE-normalization

SPADE RESIDUAL BLOCKS



SPADE GENERATOR

































IMAGE RESULTS





IMAGE RESULTS













IMAGE-TO-IMAGE SYNTHESIS















MOTIVATION

• Al-based rendering



Traditional graphics

Geometry, texture, lighting



Machine learning graphics

Data
MOTIVATION

- Al-based rendering
- High-level semantic manipulation



PREVIOUS WORK

Image translation



pix2pixHD [2018], CRN [2017], pix2pix [2017]

Video style transfer





MoCoGAN [2018], TGAN [2017], VGAN [2016]



PREVIOUS WORK: FRAME-BY-FRAME RESULT



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

Sequential Generator



Sequential Generator



Multi-scale Discriminators

Image Discriminator





Video Discriminator

 D_3

Spatio-temporally Progressive Training

Spatially progressive



Temporally progressive



Alternating training



RESULTS

- Semantic \rightarrow Street view scenes
- Edges \rightarrow Human faces
- Poses \rightarrow Human bodies

RESULTS

- Semantic \rightarrow Street view scenes
- Edges \rightarrow Human faces
- Poses \rightarrow Human bodies

STREET VIEW: CITYSCAPES



Semantic map



pix2pixHD



COVST (video style transfer)



Ours

STREET VIEW: BOSTON



STREET VIEW: NYC



RESULTS

- Semantic \rightarrow Street view scenes
- Edges \rightarrow Human faces
- Poses \rightarrow Human bodies

FACE SWAPPING (FACE \rightarrow EDGE \rightarrow FACE)



input

edges

output

FACE SWAPPING (SLIMMER FACE)

input



(slimmed) edges (slimmed) output

FACE SWAPPING (SLIMMER FACE)



input (slimmed) edges (slimmed) output

MULTI-MODAL EDGE \rightarrow FACE





Style 1

Style 2

Style 3

RESULTS

- Semantic \rightarrow Street view scenes
- Edges \rightarrow Human faces
- Poses \rightarrow Human bodies



9 🛛 💿 NIDIA



input

poses





51 💿 💿 🕺 💿



input

poses



MOTION TRANSFER





EXTENSION: FRAME PREDICTION

- Goal: predict future frames given past frames
- Our method: decompose prediction into two steps
 - 1. predict the semantic map for next frame
 - 2. synthesize the frame based on the semantic map

EXTENSION: FRAME PREDICTION





Ground truth







MCNet

INTERACTIVE GRAPHICS

- Real-time inference
- Combining with existing graphics pipeline
- Domain gap between real input and synthetic input

- Real-time inference
- Combining with existing graphics pipeline
- Domain gap between real input and synthetic input

- Real-time inference
 - FP16 + TensorRT \rightarrow ~5 times speed up
 - 36ms (27.8 fps) for 1080p inference
 - Overall: 15~25 fps

- Real-time inference
- Combining with existing graphics pipeline
 - CARLA: open-source simulator for autonomous driving research
 - Make game engine render semantic maps
 - Pass the maps to the network and display the inference result

- Real-time inference
- Combining with existing graphics pipeline
- Domain gap between *real* input and *synthetic* input
 - Network trained on real data but tested on synthetic data
 - Things that differ: Object shapes/edges, density of objects, camera viewpoints, etc
 - On-going work

ORIGINAL CARLA IMAGE



RENDERED SEMANTIC MAPS



RECORDED DEMO RESULTS



RECORDED DEMO RESULTS




- What can we achieve?
- What can it be used for?

- What can we achieve?
 - Synthesize high-res realistic images



- What can we achieve?
 - Synthesize high-res realistic images
 - Produce temporally-smooth videos



- What can we achieve?
 - Synthesize high-res realistic images
 - Produce temporally-smooth videos
 - Reinvent interactive graphics



- What can we achieve?
- What can it be used for?
 - AI-based rendering
 - High-level semantic manipulation



THANK YOU



https://github.com/NVIDIA/vid2vid

