

Carnegie Mellon

School of Computer Science

Deep Reinforcement Learning and Control

Sim2Real Transfer

Spring 2020, CMU 10-403

Katerina Fragkiadaki



Paradox

The requirement of large number of samples for model-free RL, only possible in simulation, renders model-free RL a model-based framework: we can't do without the simulator.

Choices

We want to learn manipulation and locomotion policies, what do we do?

1. We use a Physics simulator, where Physics rules between objects and/or particles have been hand-coded by engineers. We train our policies there with reinforcement (trial-and-error) and/or demonstrations. We then **transfer** them to the real world.
2. We directly learn policies in the real world.

Choices

We want to learn manipulation and locomotion policies, what do we do?

1. We use a Physics simulator, where Physics rules between objects and/or particles have been hand-coded by engineers. We train our policies there with reinforcement (trial-and-error) and/or demonstrations. We then **transfer** them to the real world.
2. We directly learn policies in the real world.
3. We build better simulators to better model the real world, e.g., food: object deformation, fluids, etc. . Then, GOTO 1.

This Lecture

We want to learn manipulation and locomotion policies, what do we do?

1. We use a Physics simulator, where Physics rules between objects and/or particles have been hand-coded by engineers. We train our policies there with reinforcement (trial-and-error) and/or demonstrations. **We then transfer them to the real world.**
2. We directly learn policies in the real world.
3. We build better simulators to better model the real world, e.g., food: object deformation, fluids, etc.. Then, GOTO 1.

Pros of Simulation

- We can afford many samples
- Safe: we do not want to deploy partially trained policies in the real world
- Avoids wear and tear of the robot
- We can explore creative robot configurations

Cons of Simulation

- Under-modeling: It is hard to exactly replicate the real world and its physics and mechanics
- Large engineering effort into building the environment which we care to manipulate
- Wrong parameters. Even if our physical equations were correct, we would need to estimate the right parameters, e.g., inertia, frictions (system identification).
- Systematic discrepancy w.r.t. the real world regarding:
 1. observations
 2. dynamics

Result: Policies learnt in simulation usually do not directly transfer to the real world

Simulators

MuJoCo: rigid and deformable body simulator on a CPU

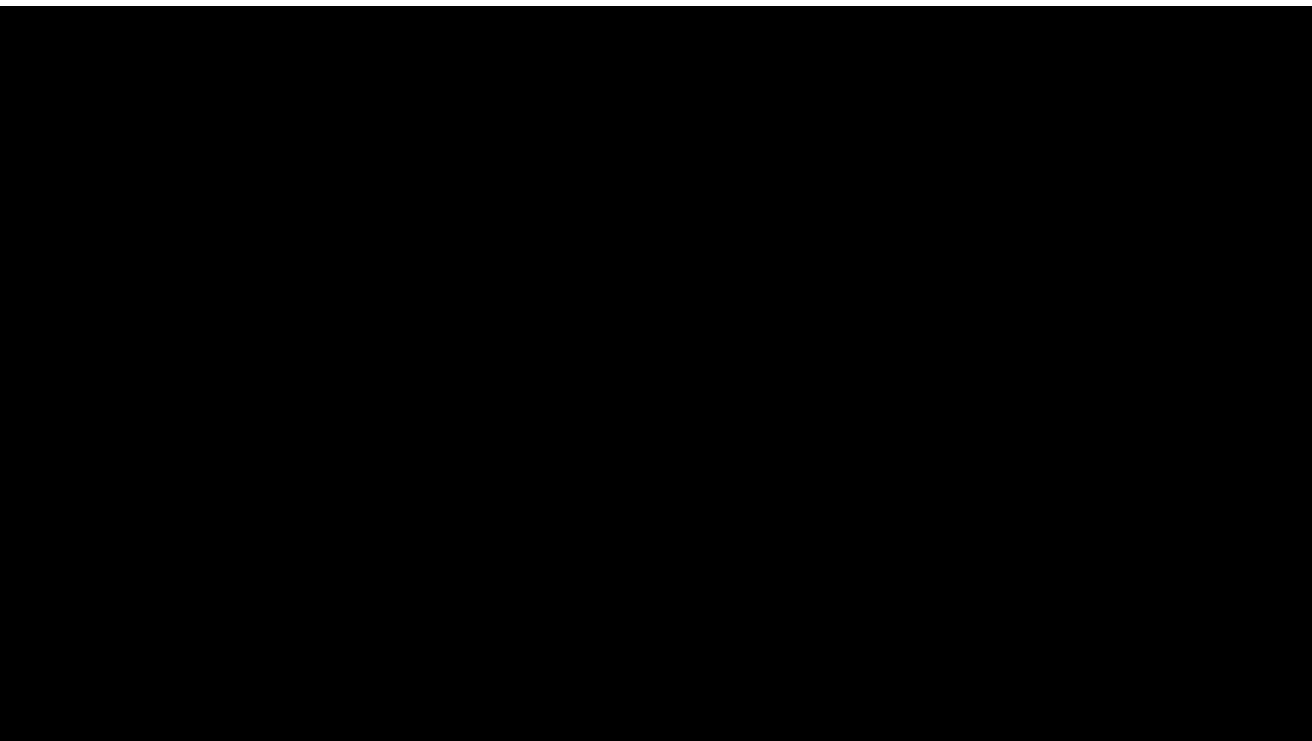
MuJoCo physics

Roboti LLC

www.mujoco.org

<http://www.mujoco.org/image/home/mujocodemo.mp4>

FLEX: particle based simulator on a GPU for rigid / soft bodies, fluids, gas.



<https://www.youtube.com/watch?v=1o0Nuq71gl4>

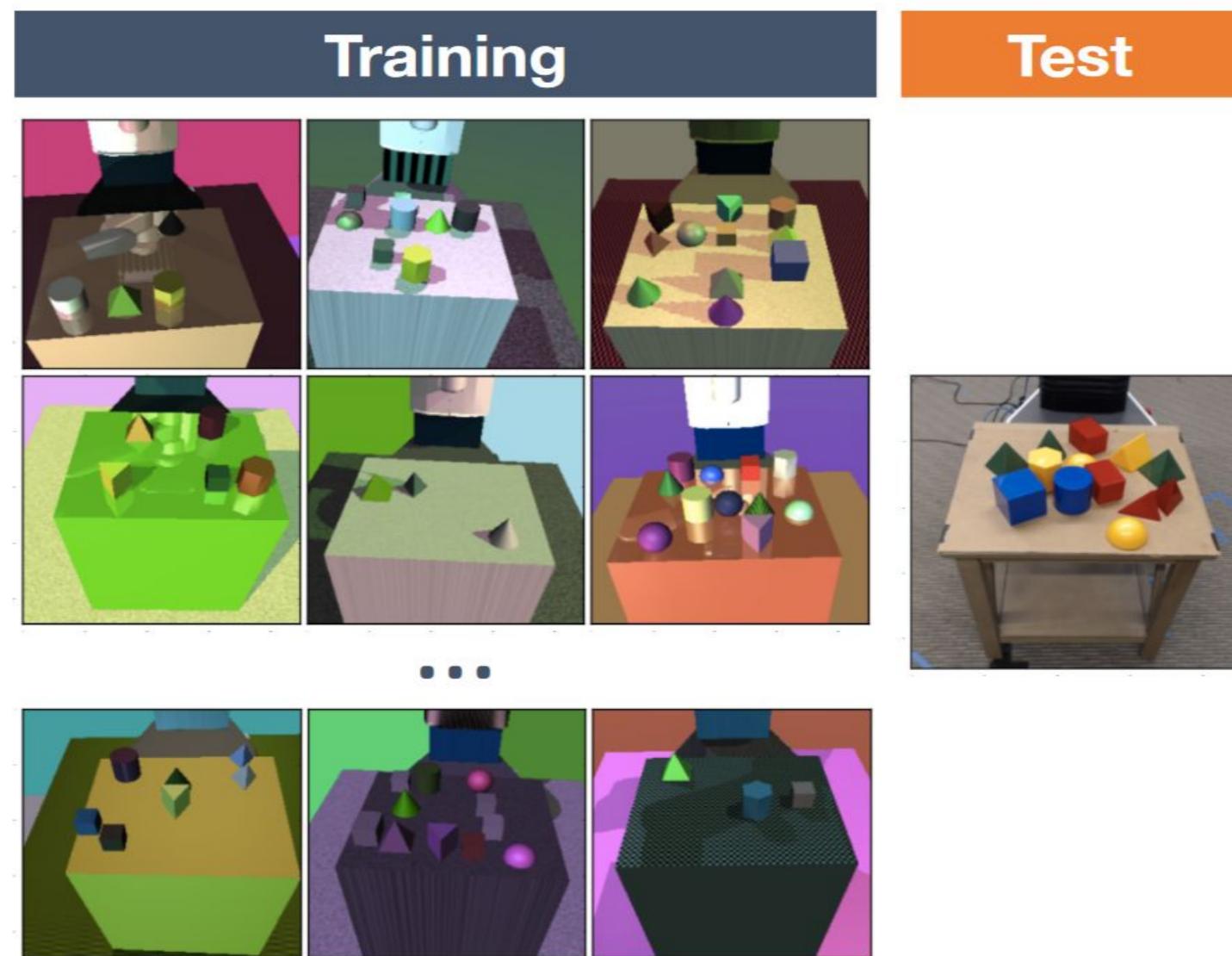
Sim2Real:What has shown to work

- Domain randomization (dynamics, visuals)
- Learning to adapt the textures of the simulator to match the real domain
- Learning to adapt the dynamics of the simulator to match the real domain
- Learning from label images as opposed to pixel images-> semantic maps between simulation and real world are closer than textures
- Abstraction: Learning higher level policies, not low-level controllers because the low level dynamics are very different between Simulation and Reality

What has shown to work

- Domain randomization (dynamics, visuals)
- Learning to adapt the textures of the simulator to match the real domain
- Learning to adapt the dynamics of the simulator to match the real domain
- Learning from label images as opposed to pixel images-> semantic maps between simulation and real world are closer than textures
- Learning higher level policies, not low-level controllers, as the low level dynamics are very different between Sim and REAL

Domain randomization

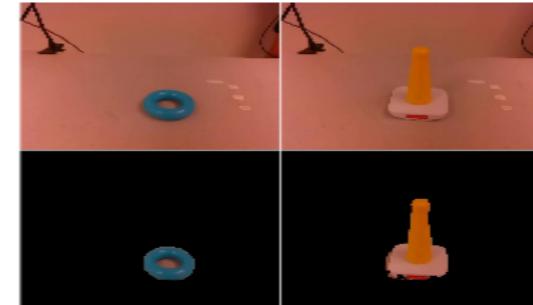


We create (automatically) tons of simulation environments by randomizing textures and camera viewpoints. We use the simulation data to train object detectors

Data dreaming

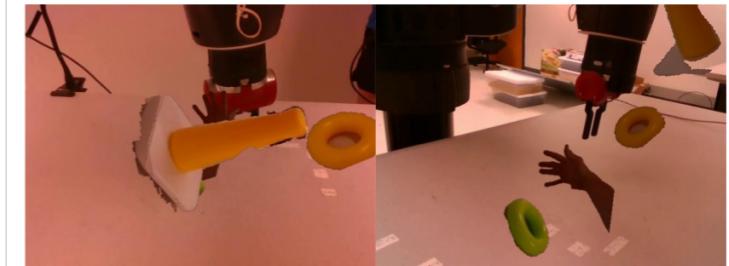
1. Obtaining object masks

- background subtraction gives ground truth object masks



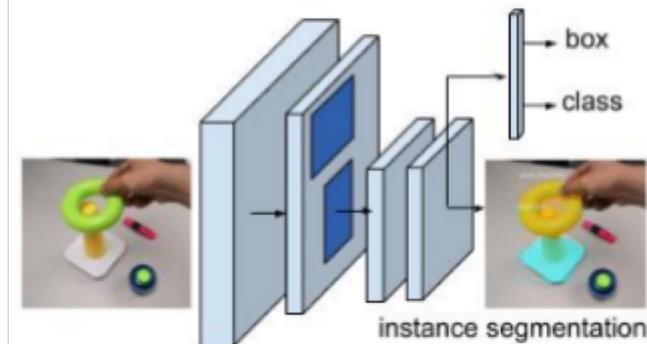
2. Creating synthetic labelled data

- Massive augmentation of ground truth masks by random transformations/occlusions and random backgrounds



3. Training object detectors

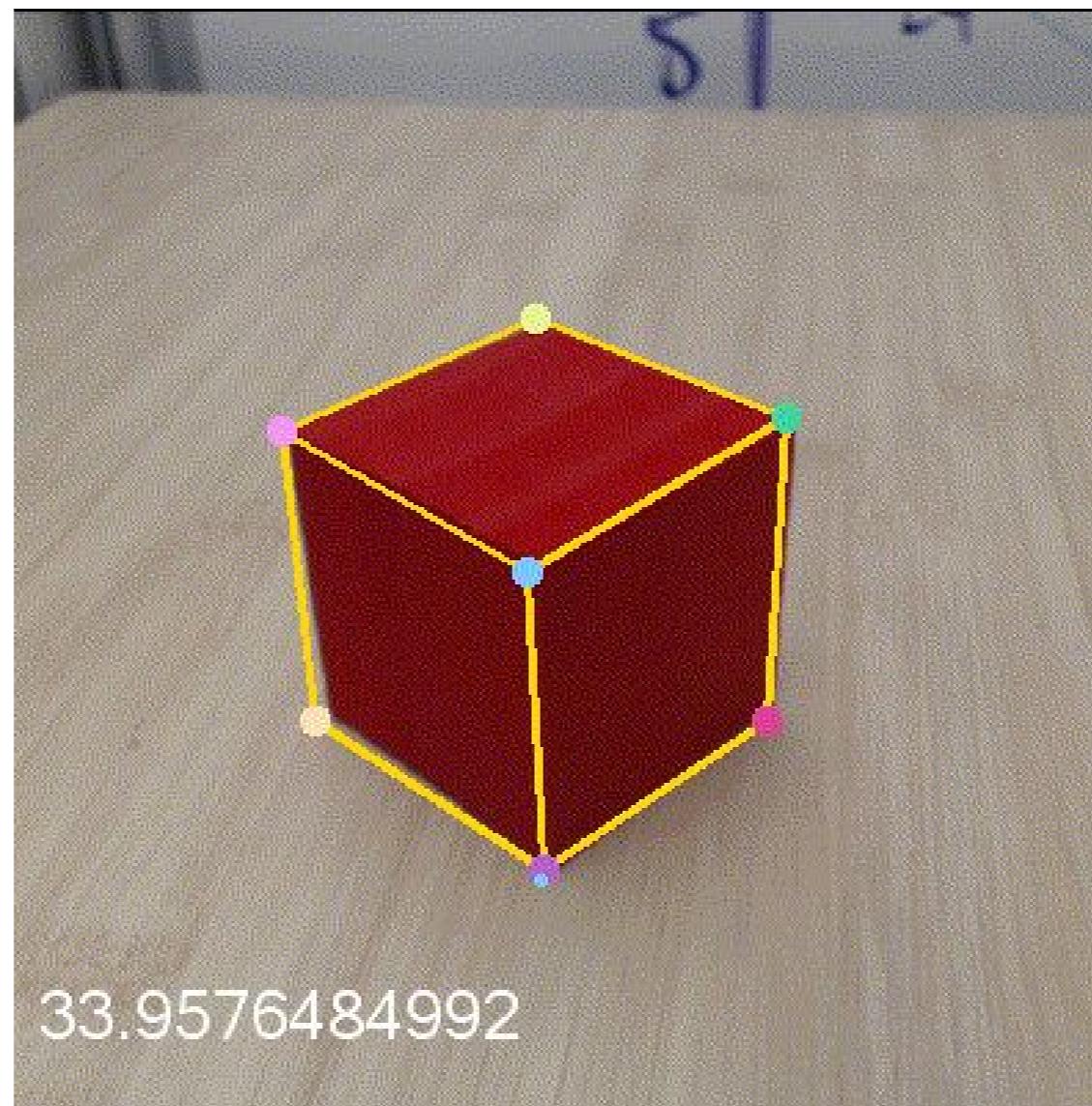
- Mask R-CNN



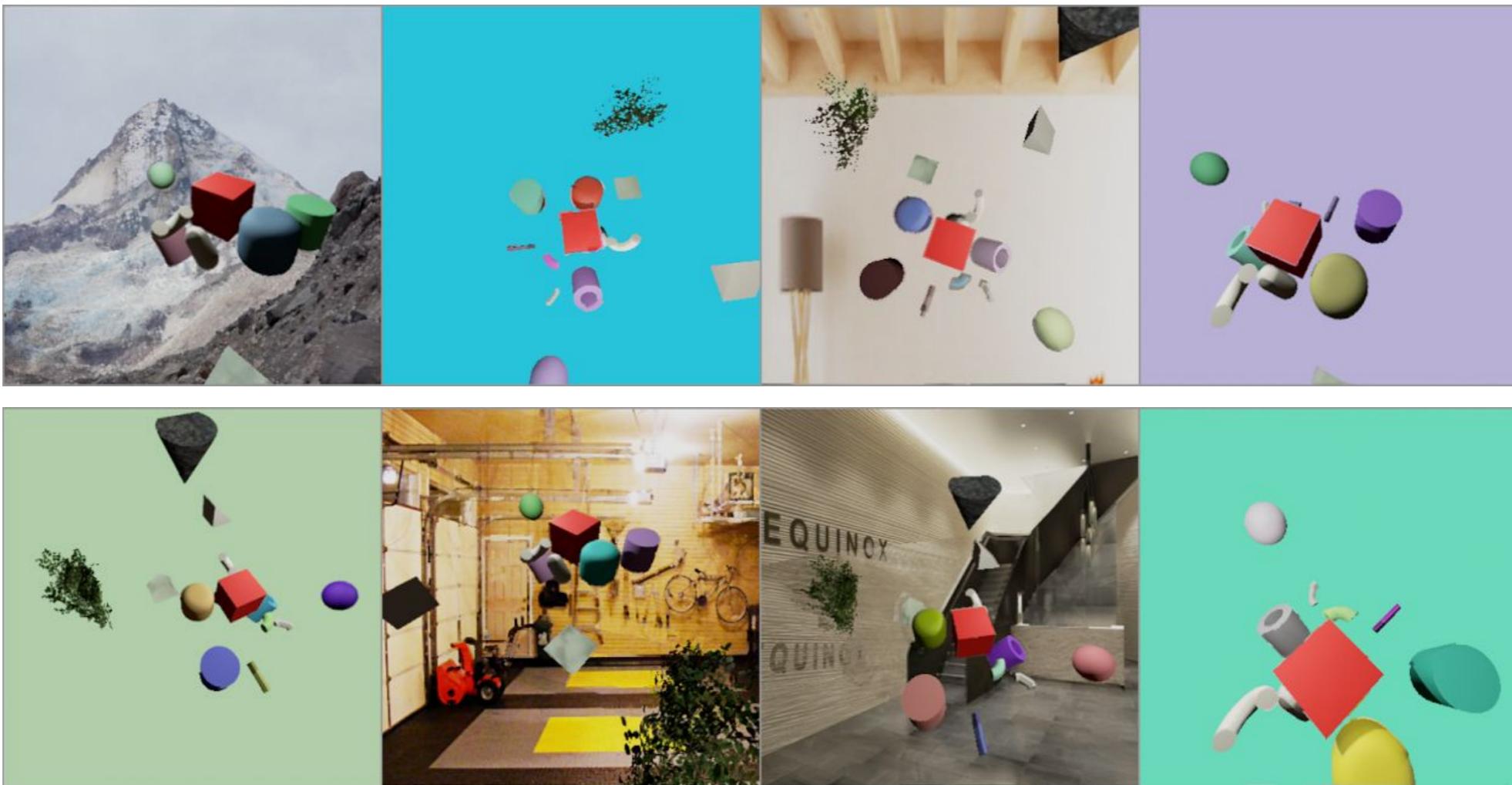
Similar randomization can be used for training object detectors on-the-fly **in the real world** directly.

Let's try a more fine grained task

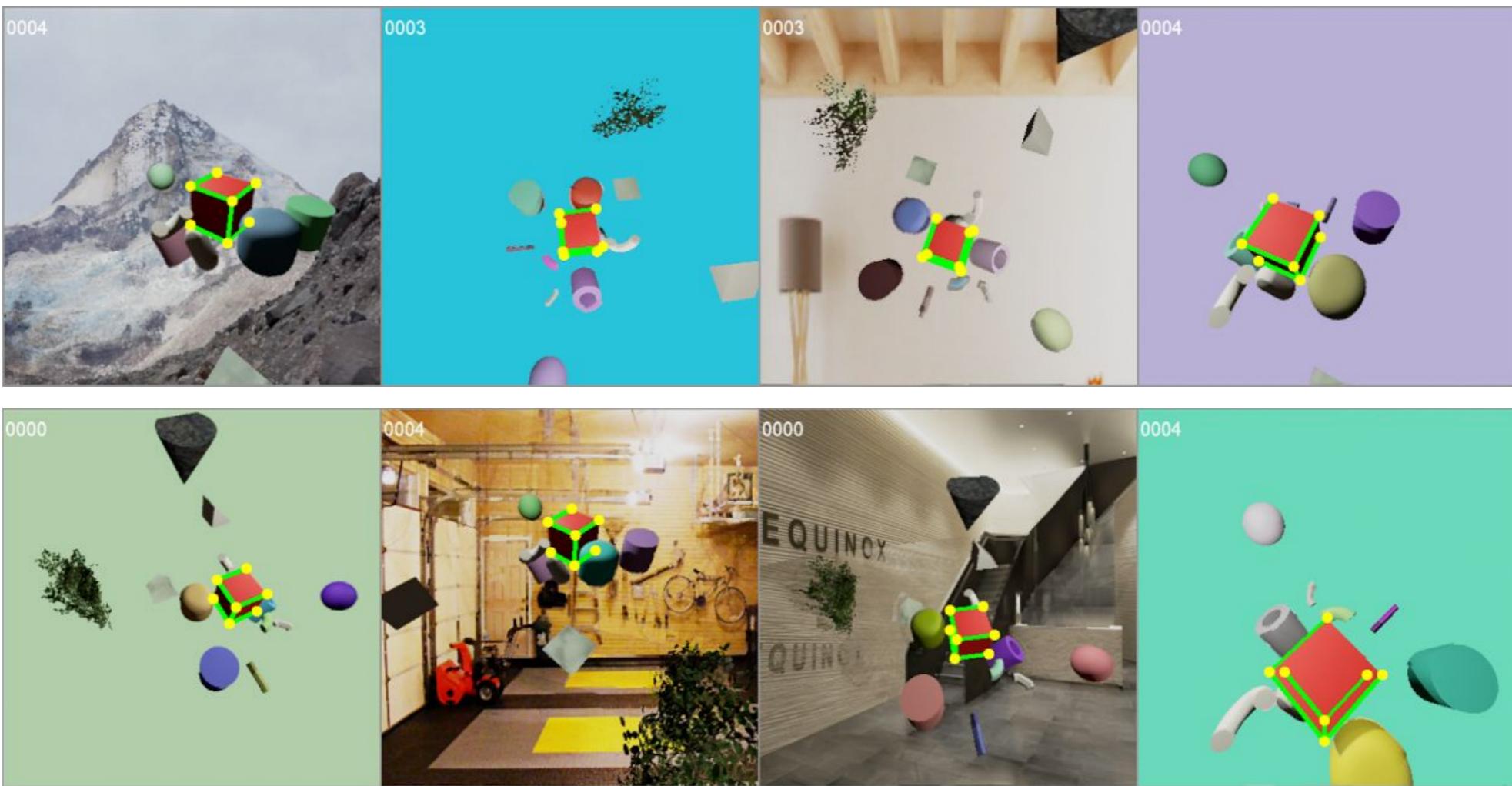
Cuboid Pose Estimation



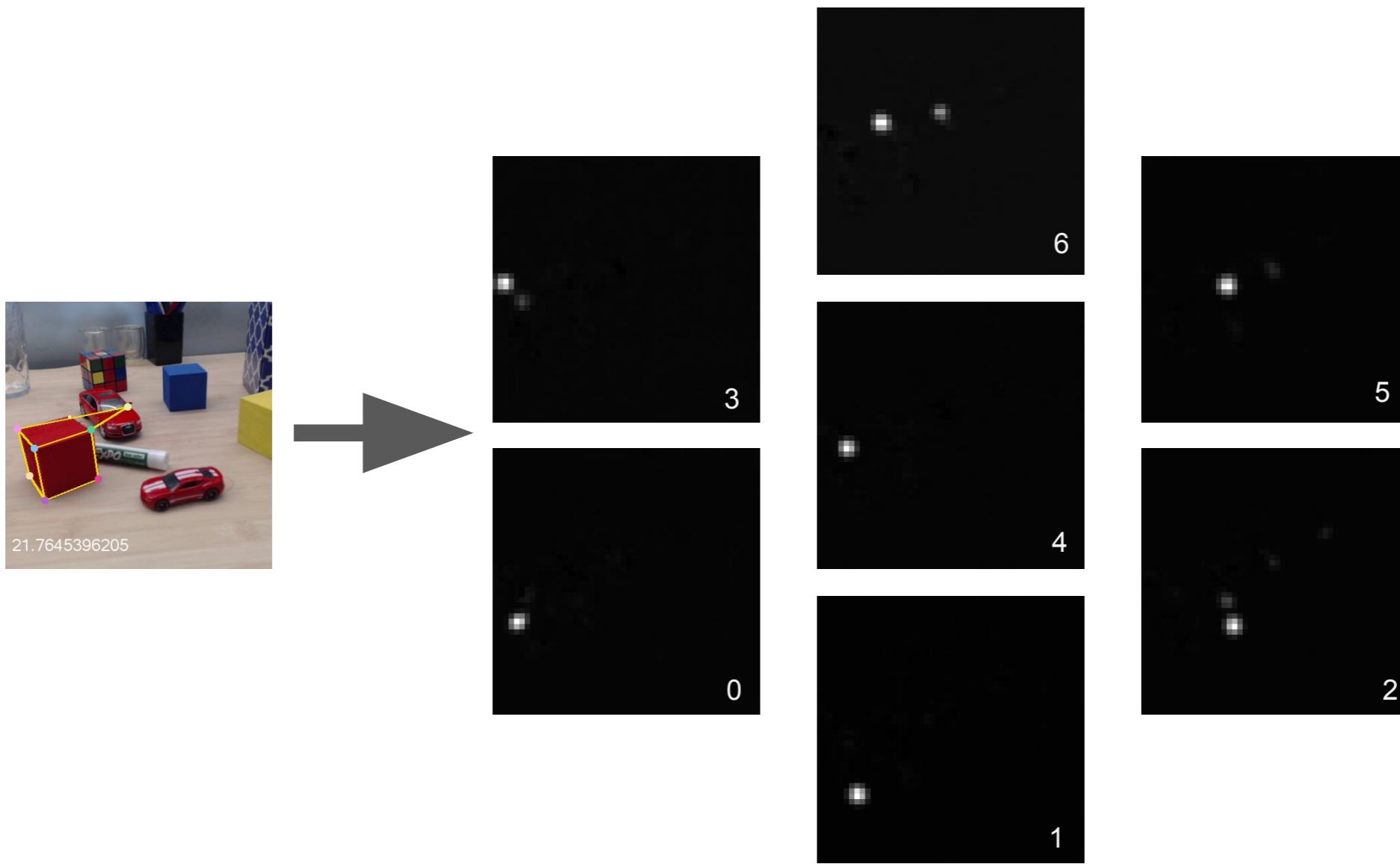
Synthetic data generation



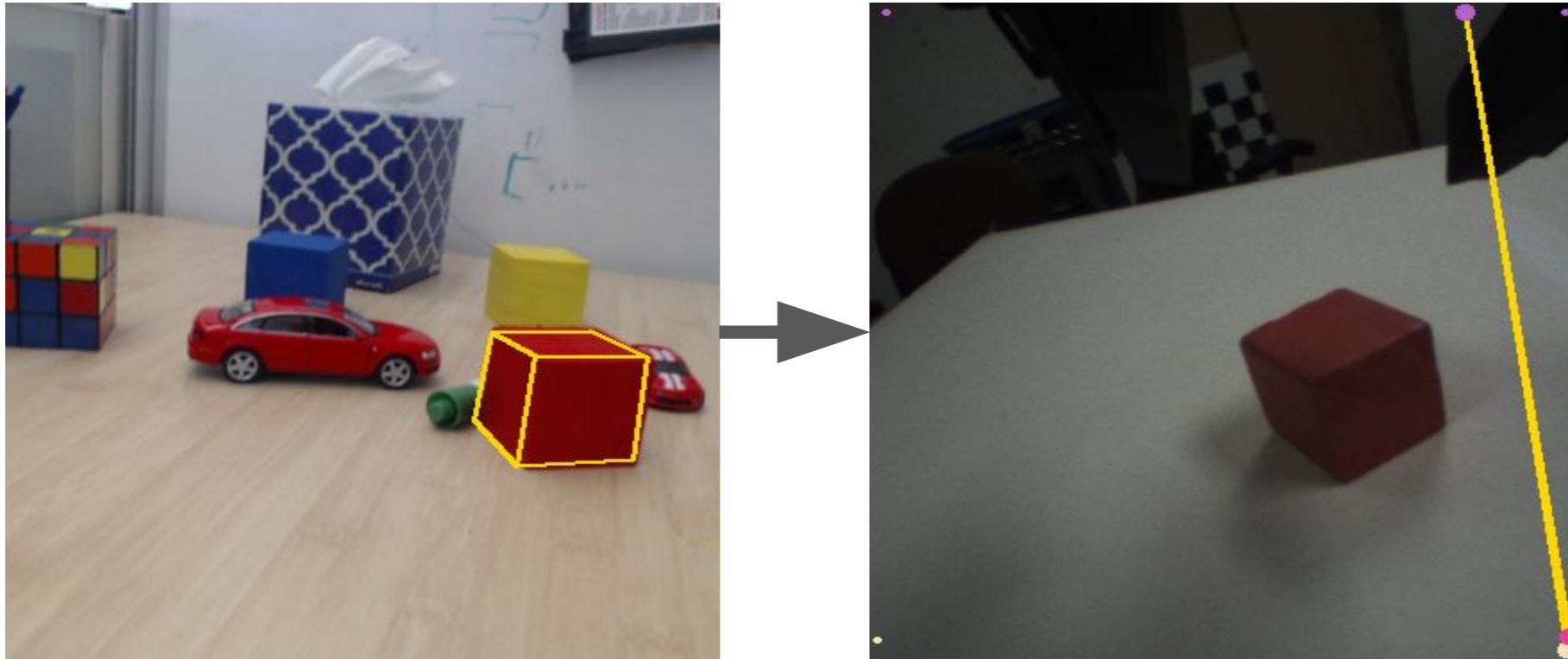
Synthetic data generation



Predicting vertex heatmaps



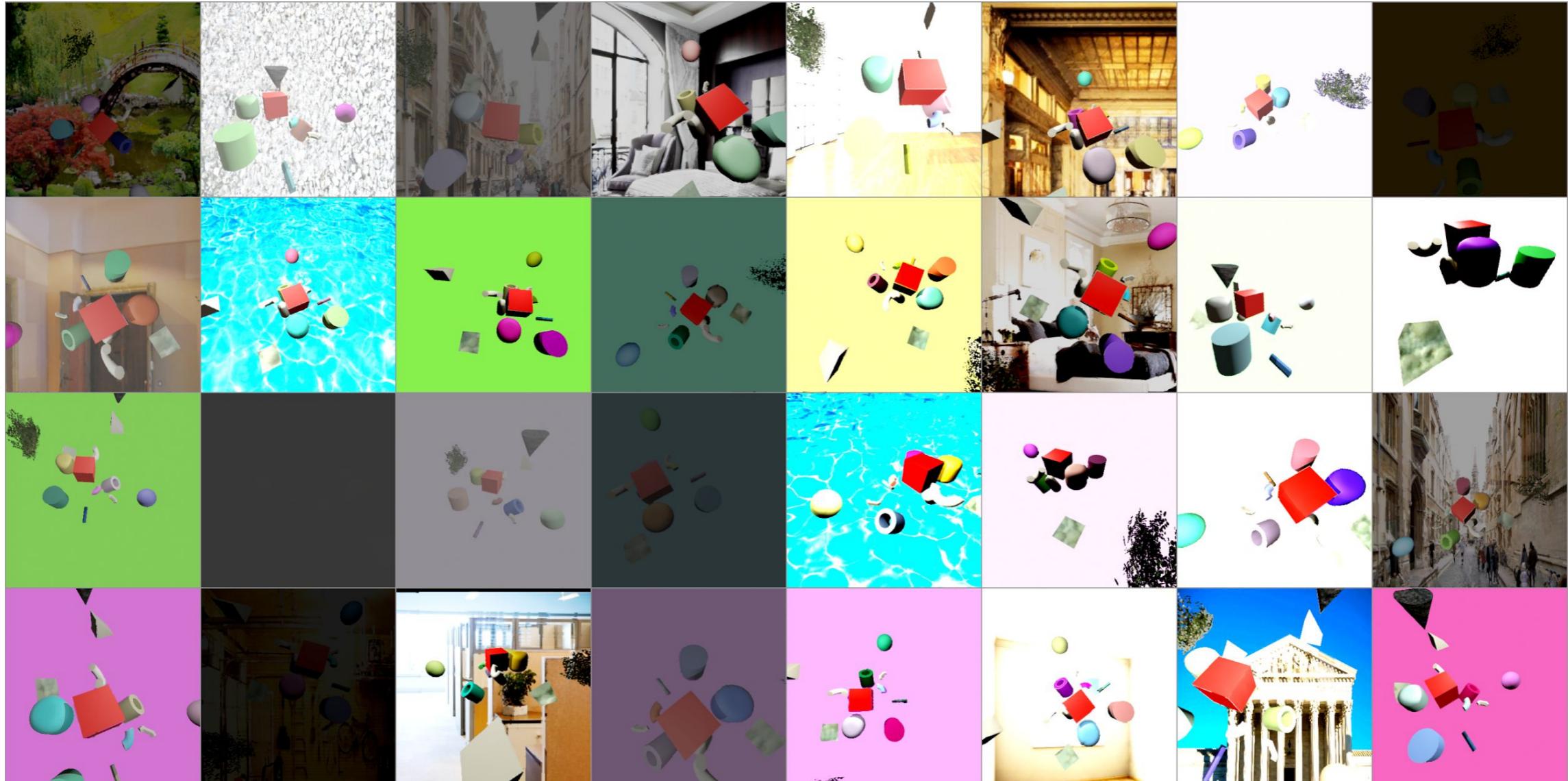
SIM2REAL



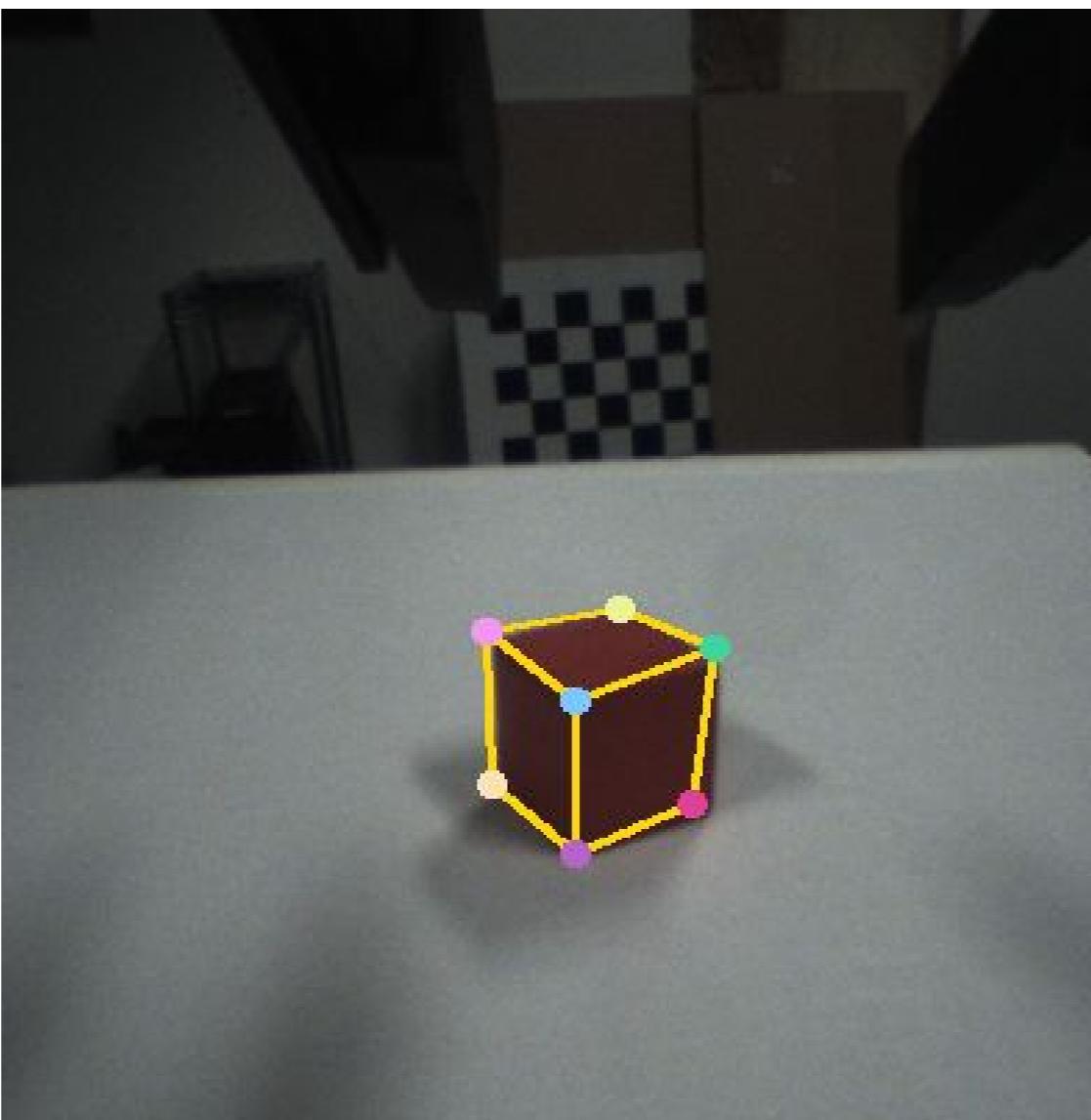
- Pose detector fails when the brightness of the image changes. Solution?
- Randomize also the brightness

Synthetic data generation

Data - Contrast and Brightness



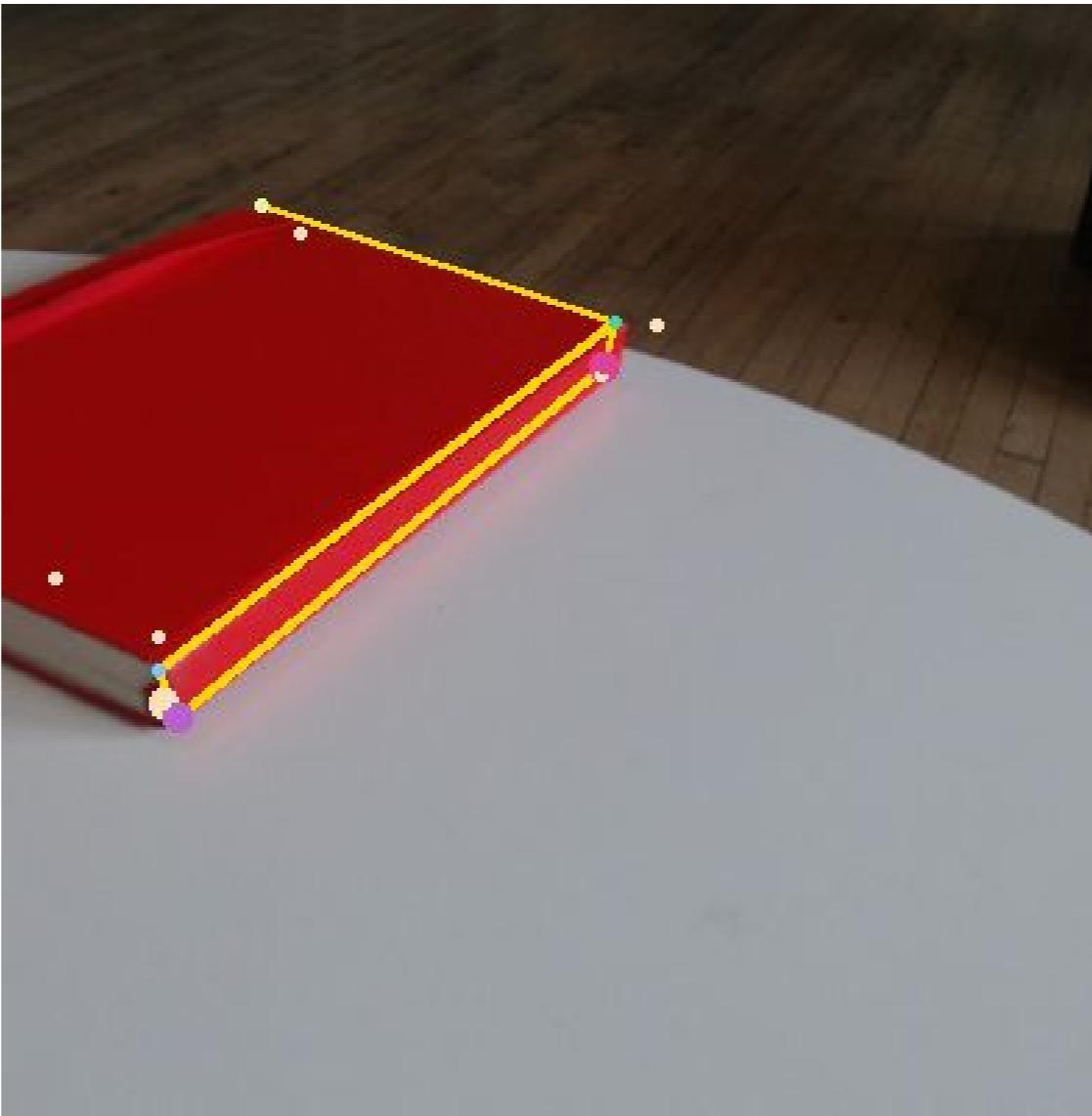
SIM2REAL



- Now it works..

SIM2REAL

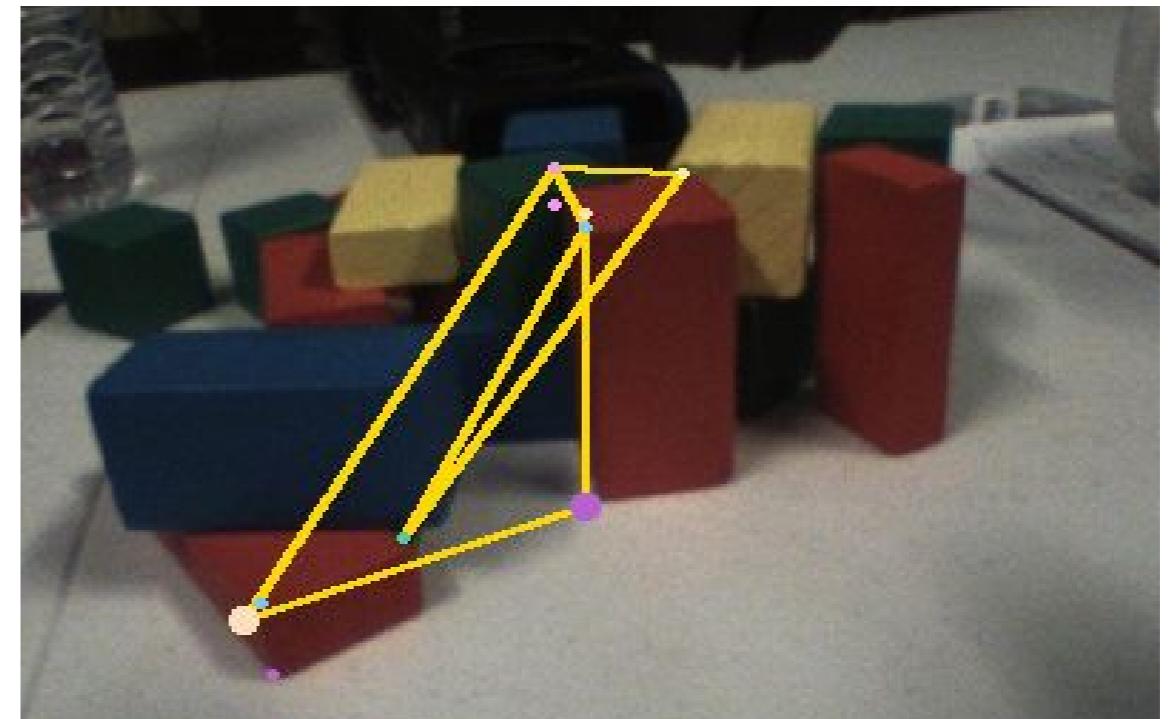
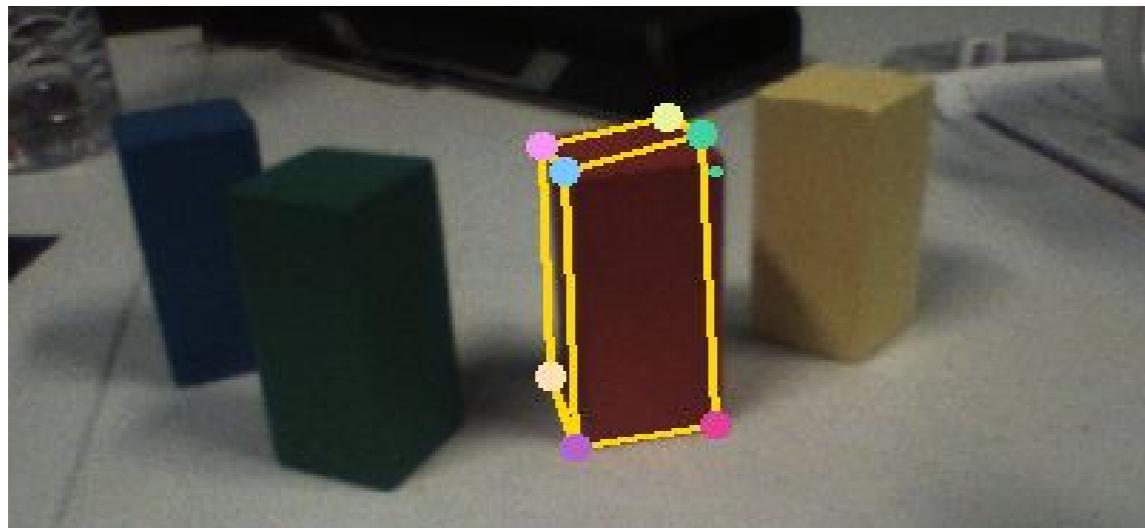
Surprising Result



- Even for non cube objects sometimes

SIM2REAL

Baxter's camera



- It can fail under clutter.
- Solution: use an architecture from computer vision research: combine object detection with vertex heatmap prediction, do not predict vertex heatmaps with the whole image as input

Car detection

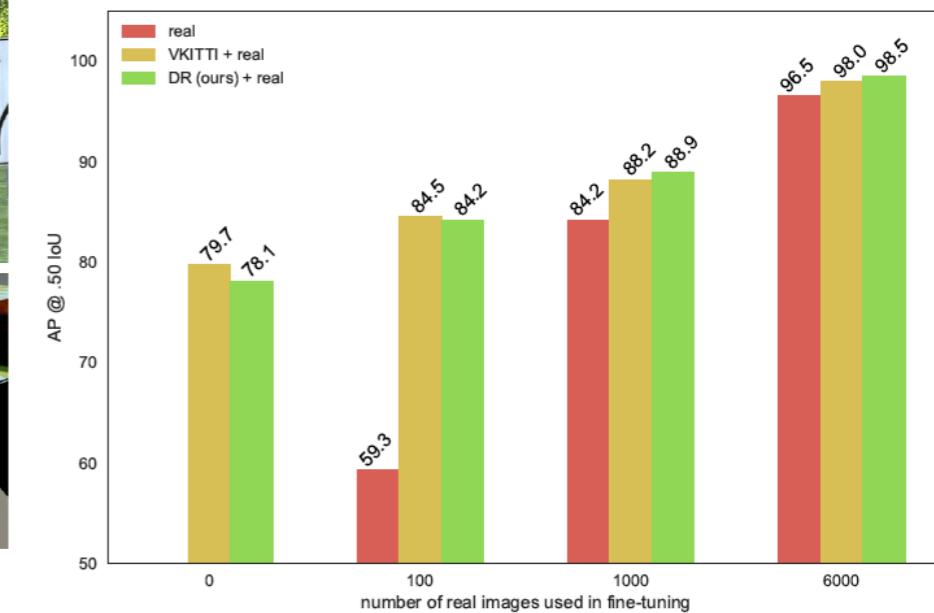
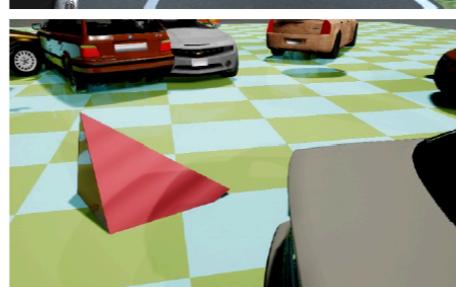
VKITTI (Virtual KITTI): a carefully designed simulation dataset to mimic real driving conditions (large engineering effort)

DR: an automatically created simulation dataset with non-realistic visuals and content (small engineering effort)

VKITTI



DR



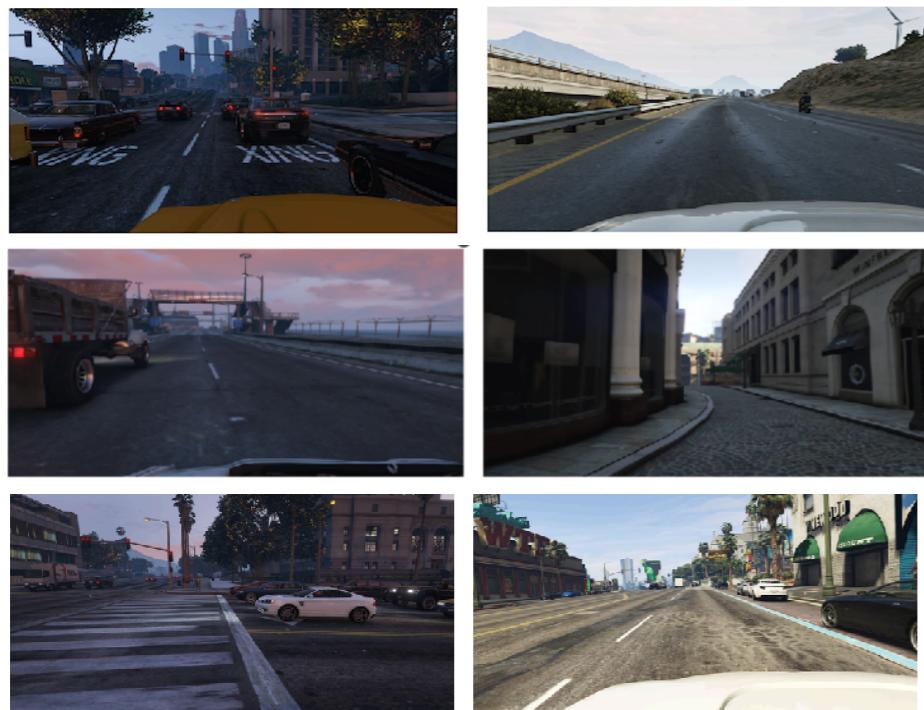
The fewer the real labelled data, the larger the gain from synthetic data

What has shown to work

- Domain randomization (dynamics, visuals)
- Learning to adapt the textures of the simulator to match the real domain
- Learning to adapt the dynamics of the dymulator to match the real domain
- Learning from label (as opposed to pixel) maps-> semantic maps between simulation and real world are closer than textures
- Learning higher level policies, not low-level controllers, as the low level dynamics are very different between Sim and REAL

Domain adaptation for visual observations

GTA: synthetic data of urban scenes from a camera mounted on a car



source

Cityscapes: real data of urban scenes from a camera mounted on a car



target

19 object classes to be detected: people, cars, stop signs, poles, etc.

Our goal: Train detectors and pixel labelers on GTA that generalize to Cityscapes

Baseline

Train a classifier on source and test it on the target, and hope it generalizes

1. Pick a network architecture, e.g. ResNet101 or VGG
2. Download a **pretrained** neural network, e.g., trained for image classification on Imagenet dataset
3. **Finetune** it on the source domain (GTA)
4. Apply it on the target domain (Cityscapes)

Training



Testing

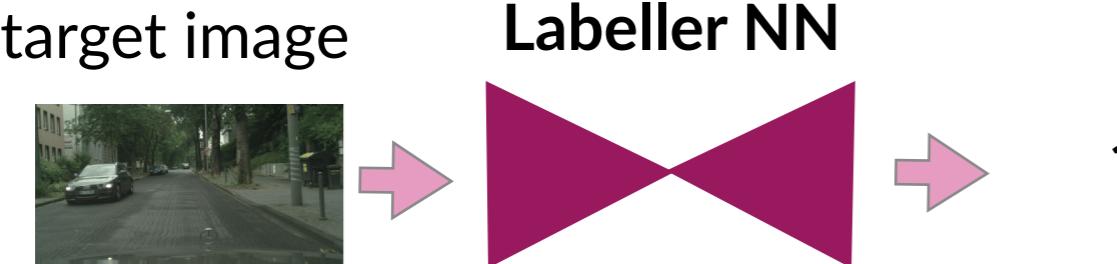


Image classification in Imagenet



Baseline

Train a classifier on source and test it on the target, and hope it generalizes

1. Download a pretrained neural network, e.g., trained for image classification on Imagenet dataset
2. Finetune it on the source domain (GTA)
3. Apply it on the target domain (Cityscapes)

Pretrain on Imagenet -> finetune in GTA->test in GTA: 53% meanIoU

Pretrain on Imagenet -> finetune in GTA->test in Cityscapes: 28% meanIoU

Pretrain on PASCAL -> finetune in GTA->test in GTA: 58.84% meanIoU

Pretrain on PASCAL -> finetune in GTA->test in Cityscapes: 32% meanIoU

Pretrain on PASCAL -> **cotrain** in GTA/PASCAL->test in Cityscapes: 39% meanIoU

Baseline

Train a classifier on source and test it on the target, and hope it generalizes

1. Download a pretrained neural network, e.g., trained for image classification on

Catastrophic forgetting:

- During fine-tuning, the network forgets the general and nicely transferable PASCAL features!
- Finetuning a neural net on a very limited domain is a bad idea for transfer

Pretrain on PASCAL -> finetune in GTA->test in GTA: 58.84% meanIoU
Pretrain on PASCAL -> finetune in GTA -> test in Cityscapes: 22%

Other solutions for catastrophic forgetting?

Baseline

Train a classifier on source and test it on the target, and hope it generalizes

1. Download a **pretrained** neural network, e.g., trained for image classification on Imagenet dataset
2. **Finetune it on the source domain (GTA):**
 1. Adapt only the top layers and keep the earlier frozen.
 2. Cotrain it using both the old task and the new task in the smaller dataset
3. Apply it on the target domain (Cityscapes)

Pretrain on Imagenet -> finetune in GTA->test in GTA: 53% meanIoU

Pretrain on Imagenet -> finetune in GTA->test in Cityscapes: 28% meanIoU

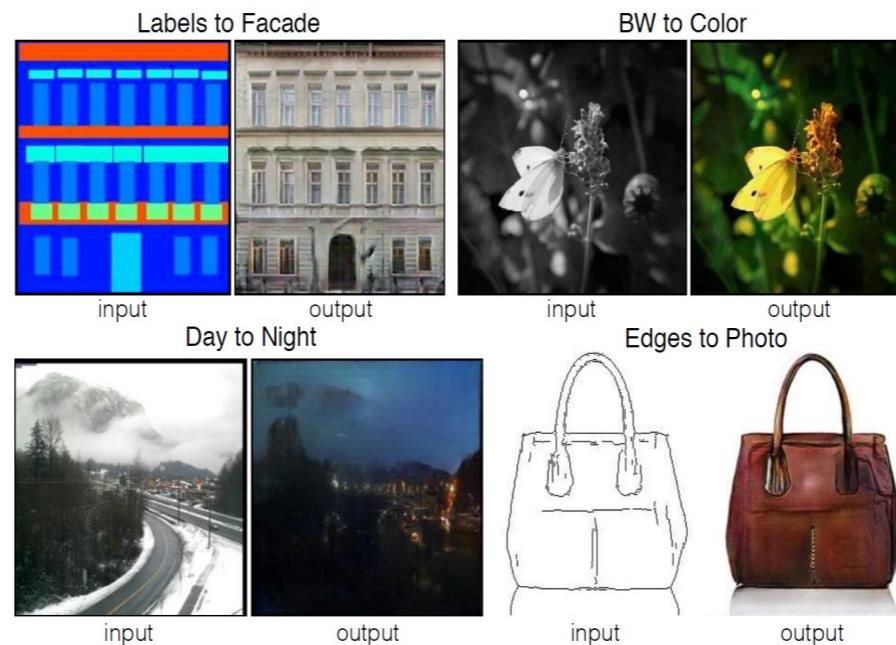
Pretrain on PASCAL -> finetune in GTA->test in GTA: 58.84% meanIoU

Pretrain on PASCAL -> finetune in GTA->test in Cityscapes: **32%** meanIoU

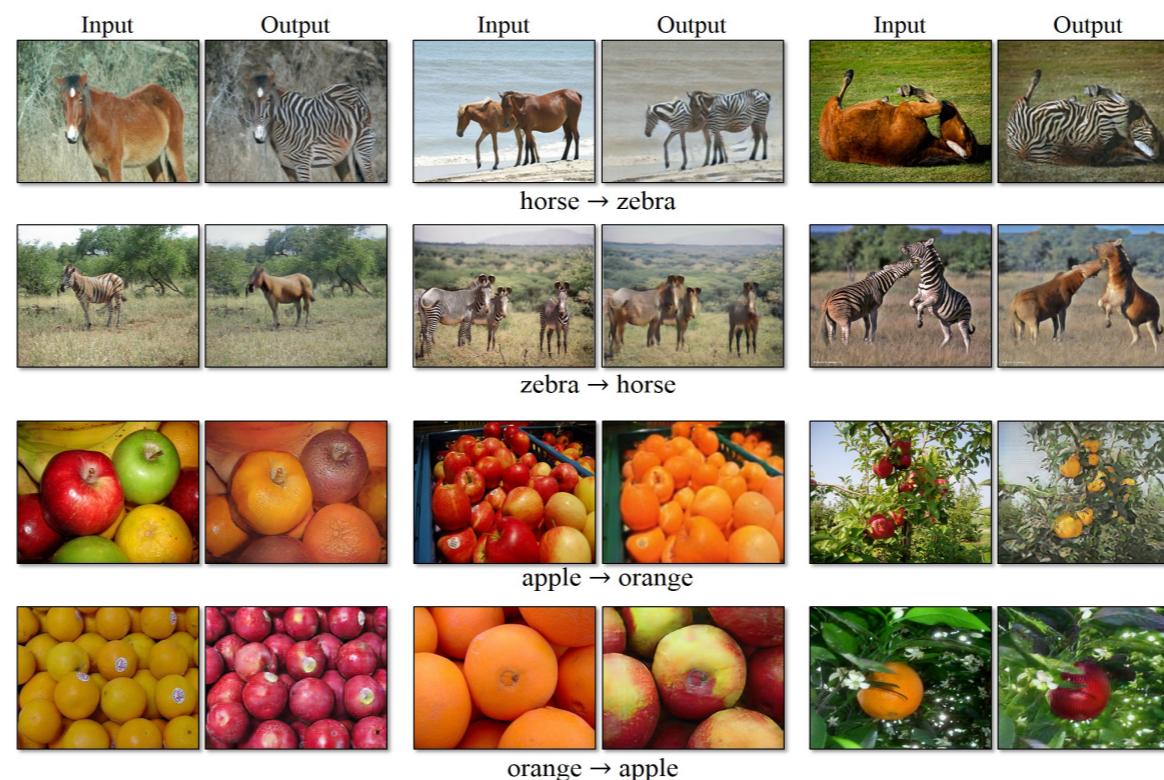
Pretrain on PASCAL -> **cotrain** in GTA/PASCAL->test in Cityscapes: **39%** meanIoU

Learning to translate images across domains (sim and real)

- Paired

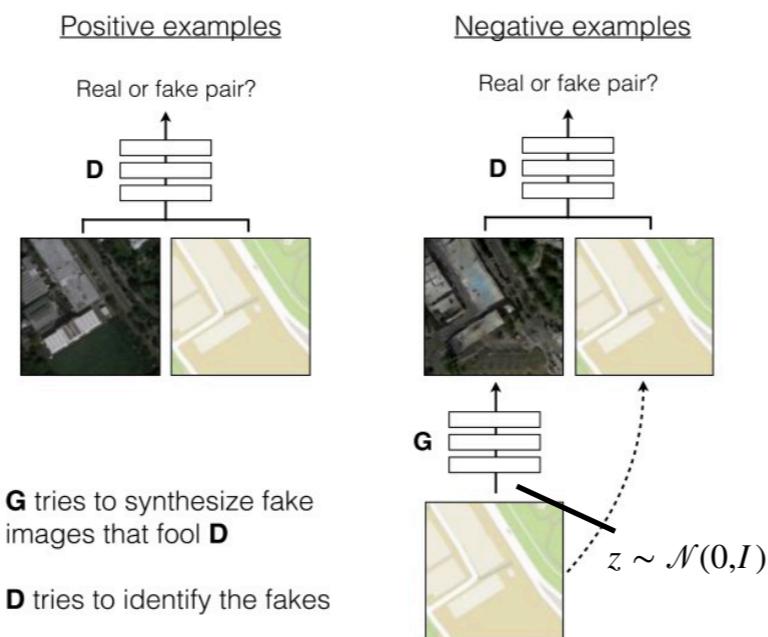


- Unpaired (this is our case)

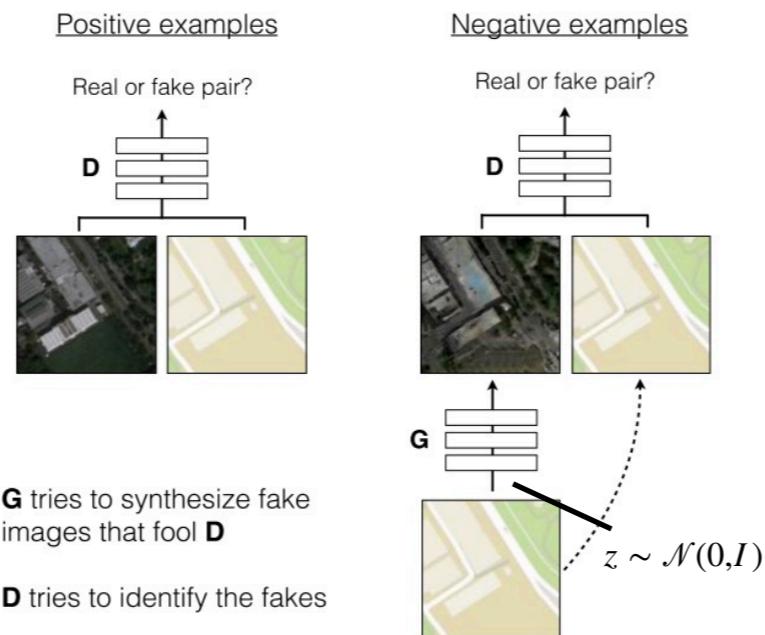


Paired case

- The generator takes the (source) image as input and tries to output the corresponding target image
- Minimize
 - image reconstruction loss
 - adversarial loss: pairs of source-target images as input to discriminator



Paired case

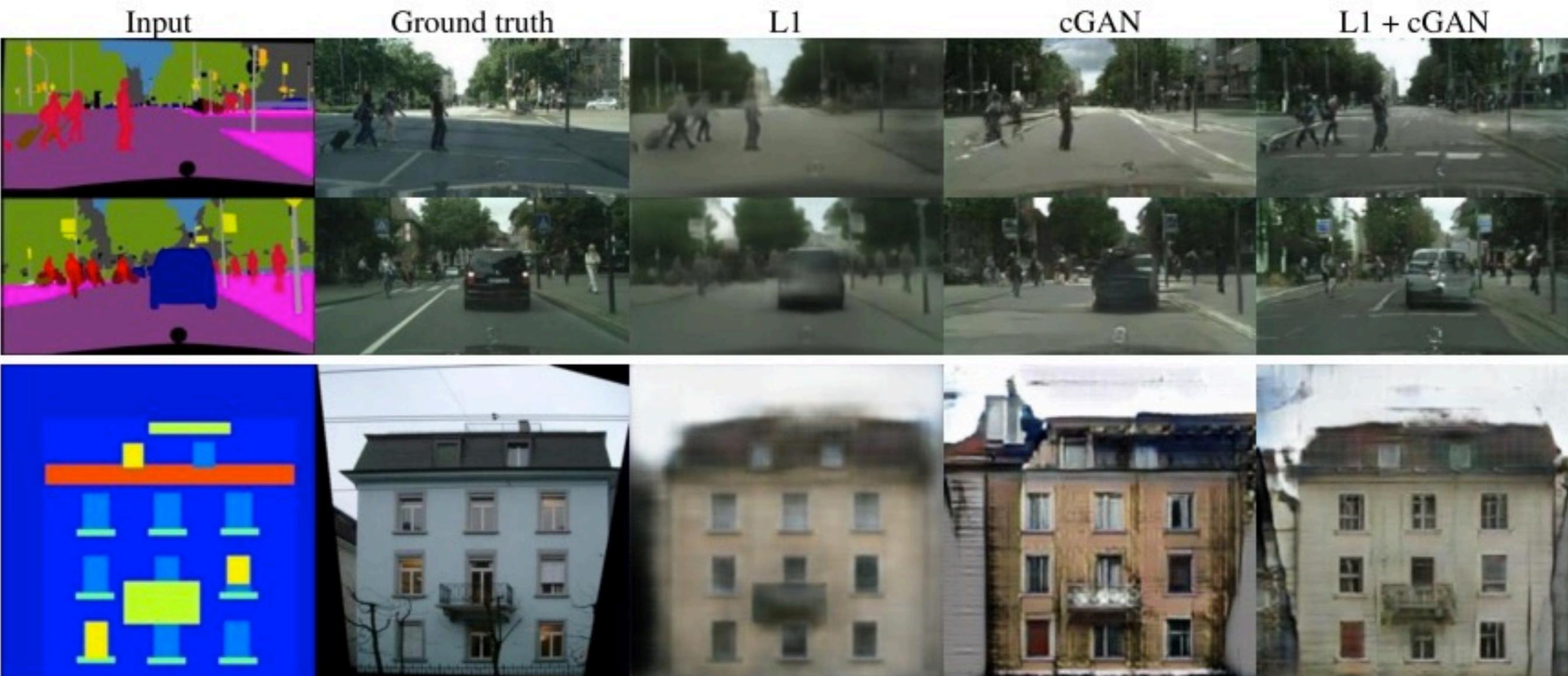


x : source image, y : target image, z : noise

$$\mathcal{L}_{cGAN}(G, D) = \mathbb{E}_{x,y \sim p_{\text{data}}(x,y)}[\log D(x, y)] + \mathbb{E}_{x \sim p_{\text{data}}(x), z \sim p_z(z)}[\log(1 - D(x, G(x, z)))]$$

$$\mathcal{L}_{L1}(G) = \mathbb{E}_{x,y \sim p_{\text{data}}(x,y), z \sim p_z(z)} [\|y - G(x, z)\|_1]$$

Paired case



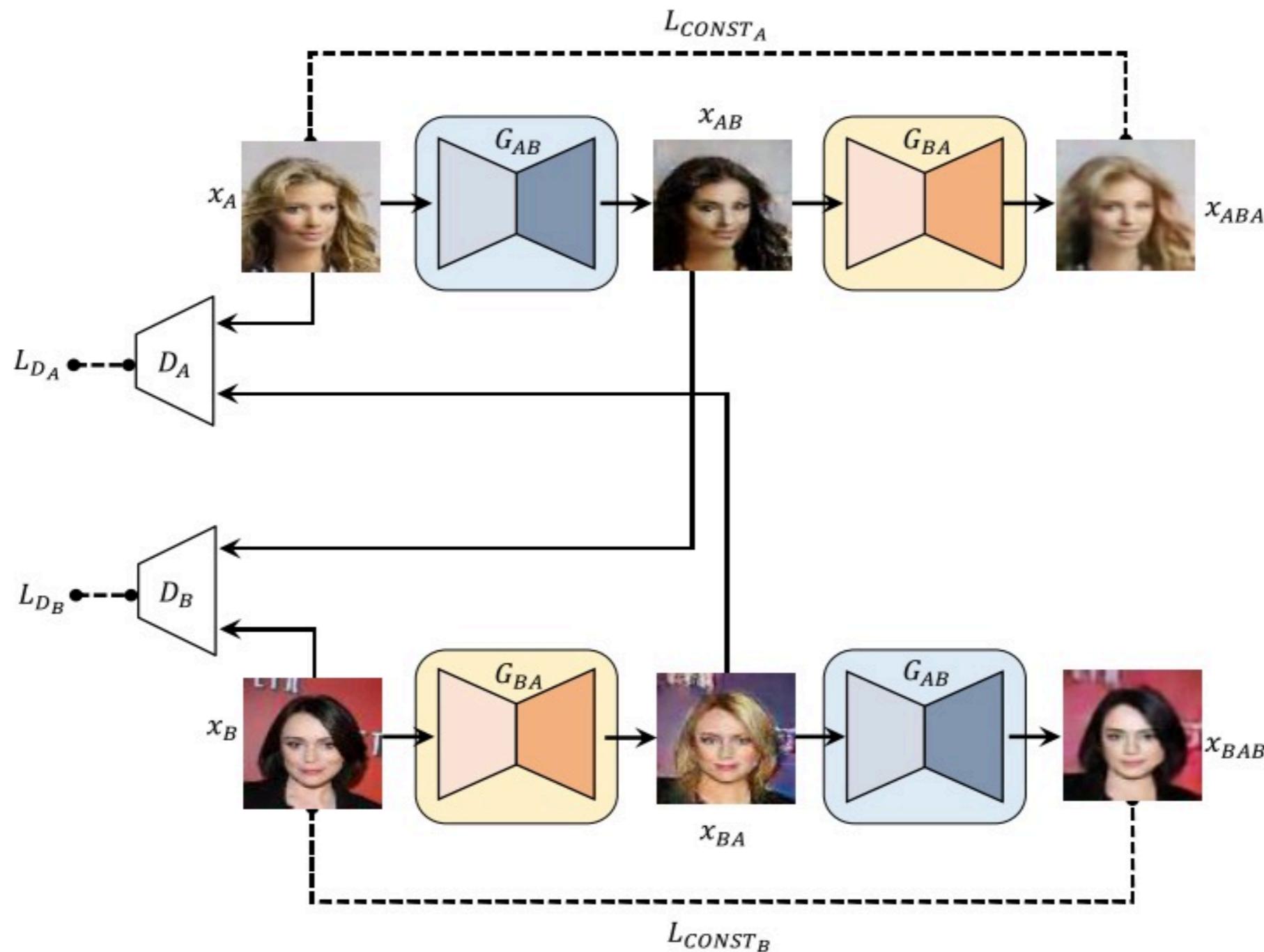
Paired case



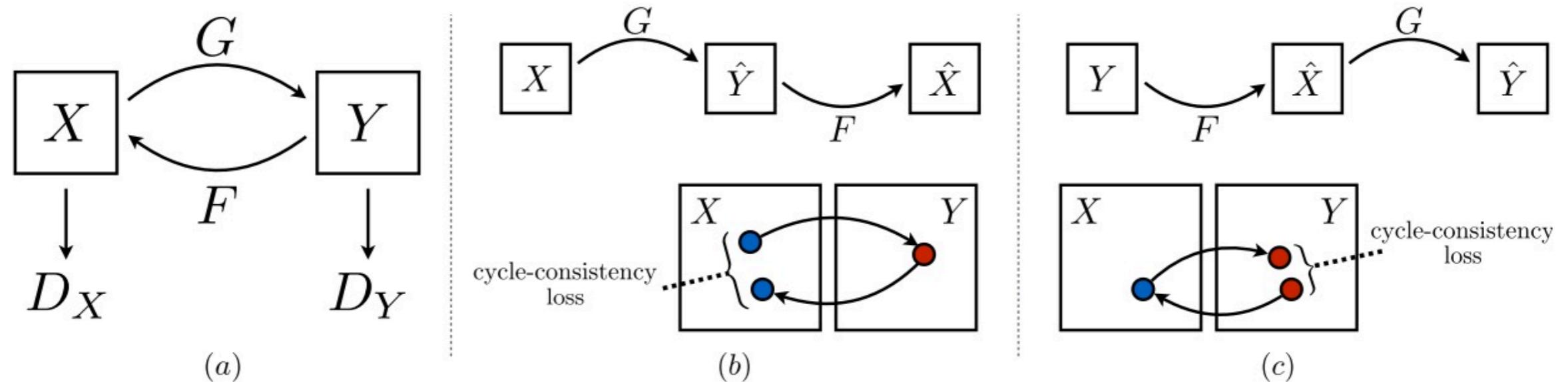
Paired case



Unpaired case



Unpaired case: Cycle GAN / DISCO GAN



Cycle GAN / DISCO GAN

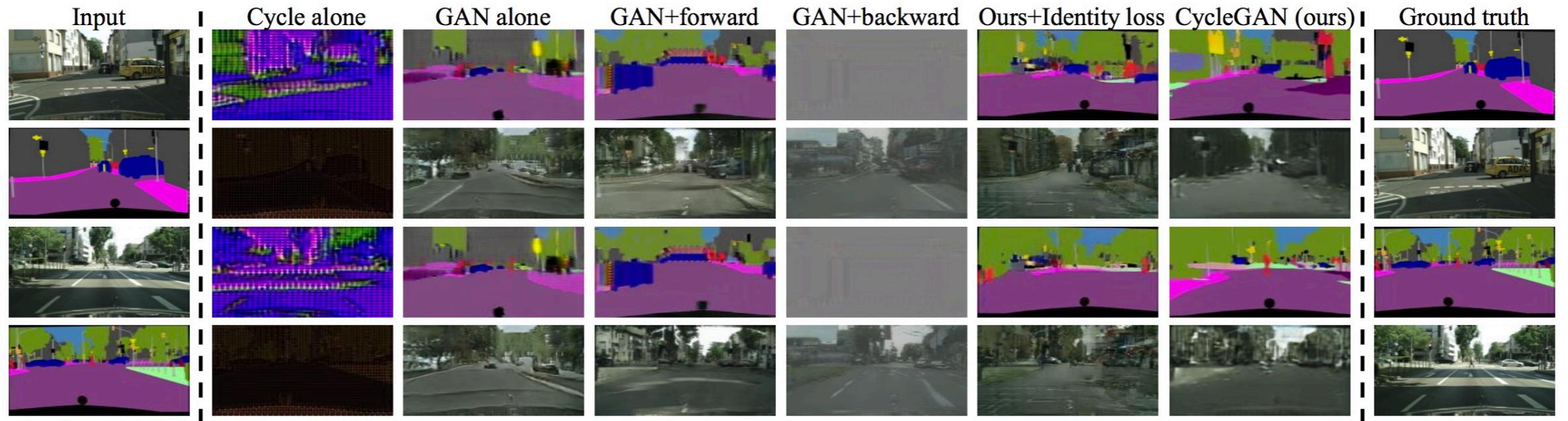
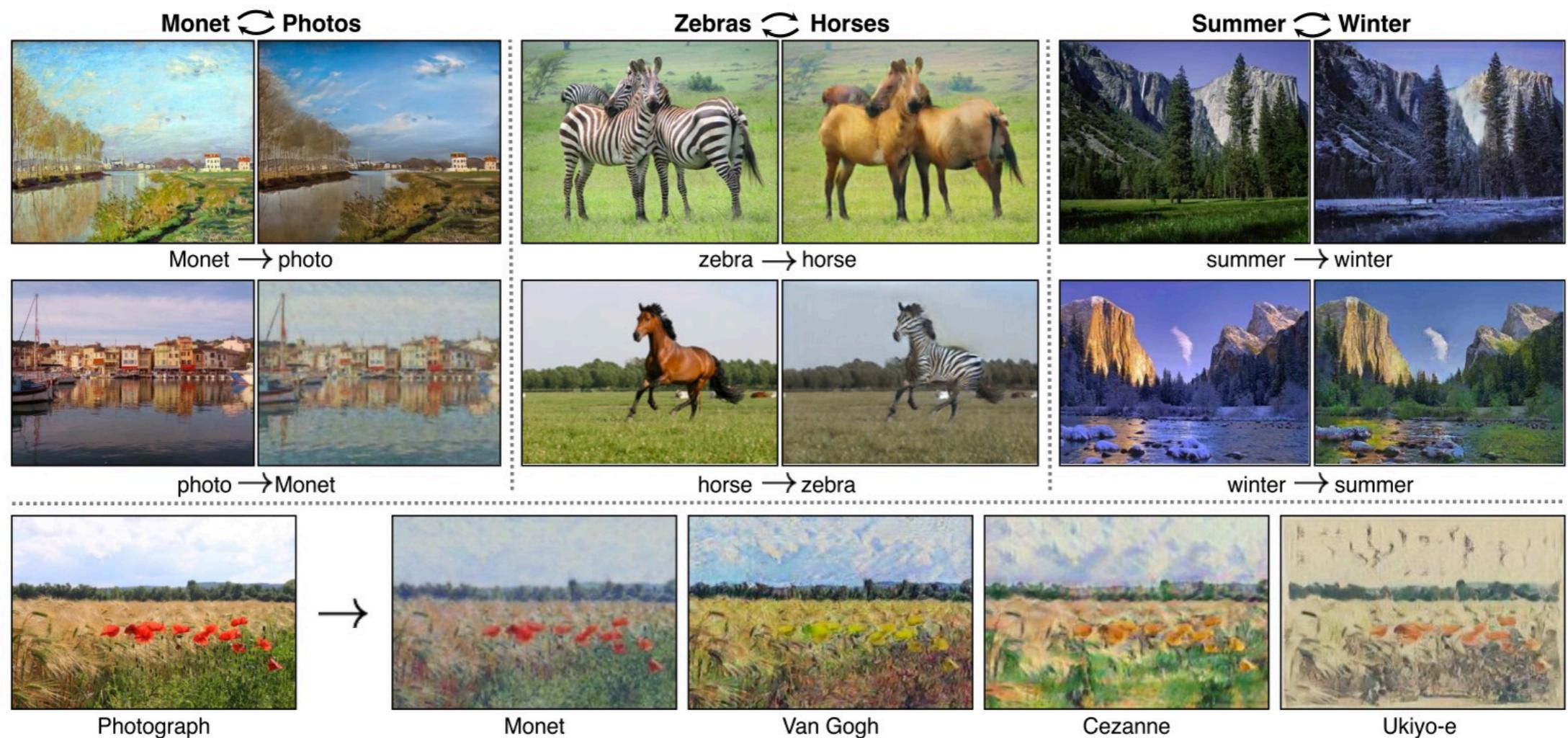
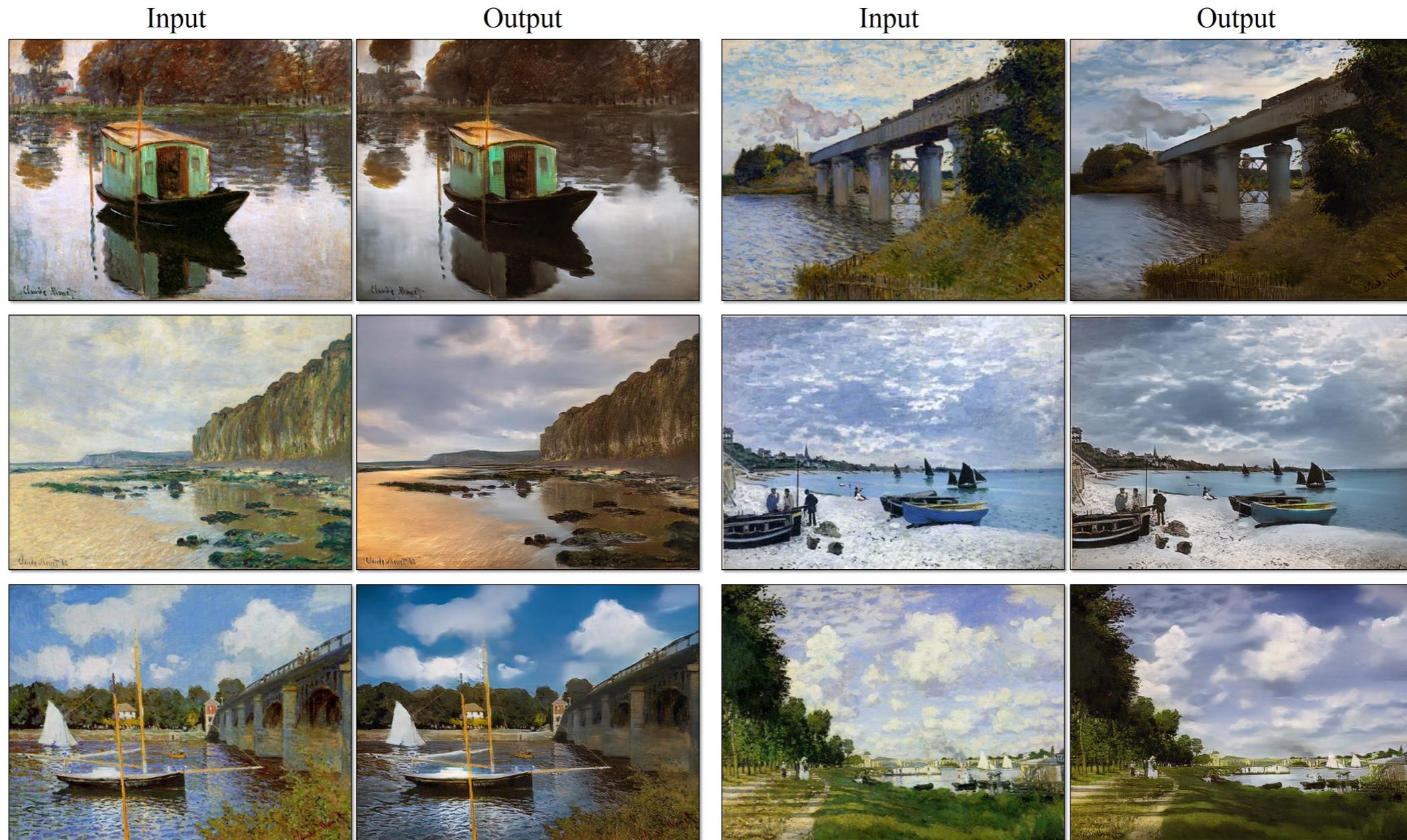


Figure 7: Different variants of our method for mapping labels \leftrightarrow photos trained on cityscapes. From left to right: input, cycle-consistency loss alone, adversarial loss alone, GAN + forward cycle-consistency loss ($F(G(x)) \approx x$), GAN + backward cycle-consistency loss ($G(F(y)) \approx y$), CycleGAN (our full method), and ground truth. Both *Cycle alone* and *GAN + backward* fail to produce images similar to the target domain. *GAN alone* and *GAN + forward* suffer from mode collapse, producing identical label maps regardless of the input photo.

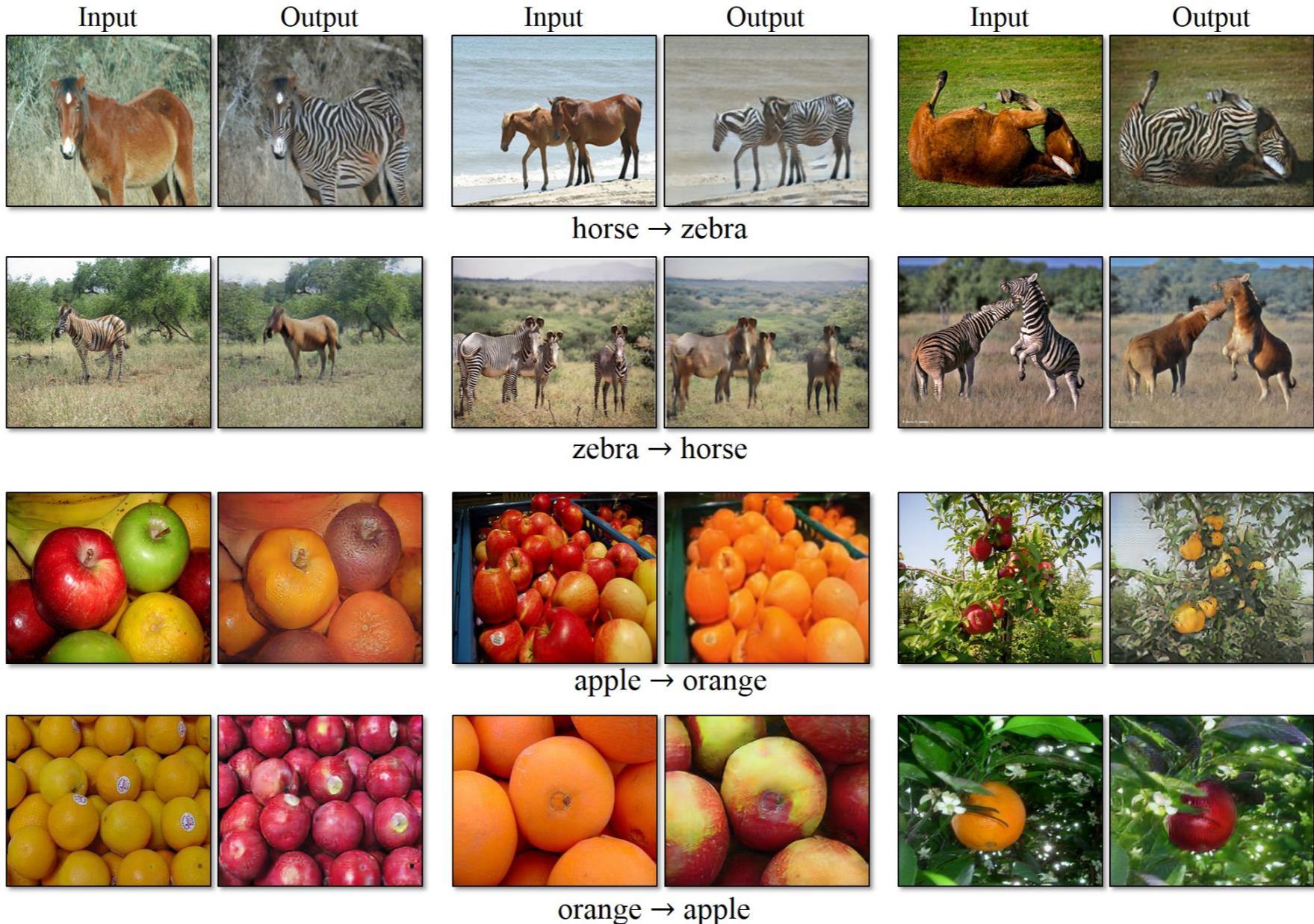
Unpaired case



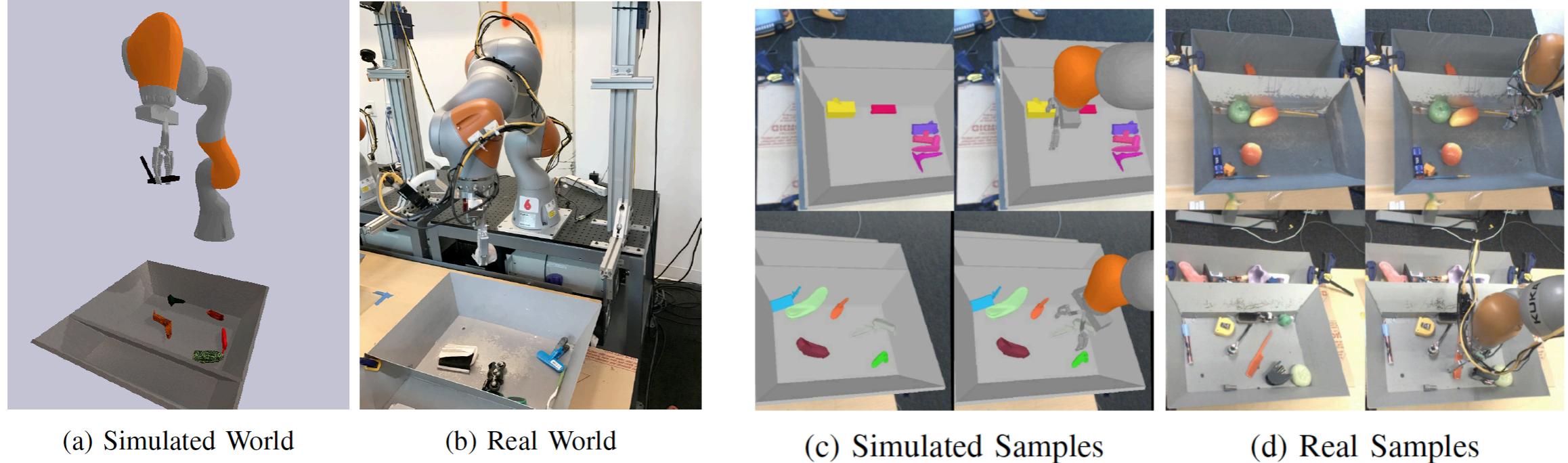
Unpaired case



Unpaired case

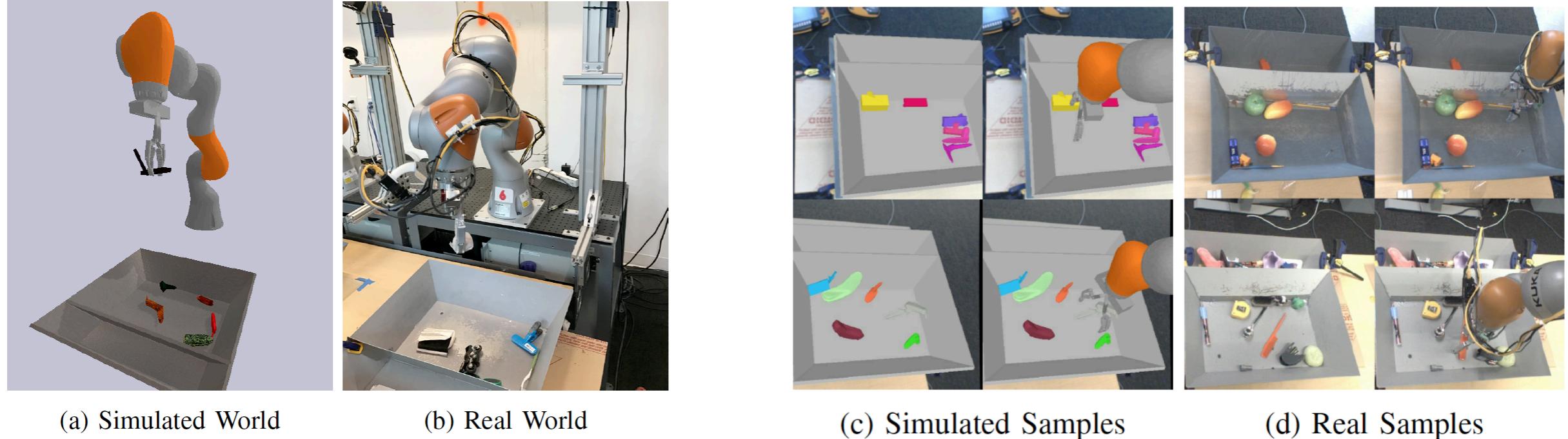


Sim2real for learning to grasp



- I want to learn the function $\text{Grasp}(I, v; \theta)$: given image I and end-effector motion v , will I successfully grasp the object?
- $\text{Grasp}(I, v; \theta)$ can be trained with supervised learning. I want to use a simulated environment to quickly collect lots of samples. I want it to generalize to the real world.

Sim2real for learning to grasp

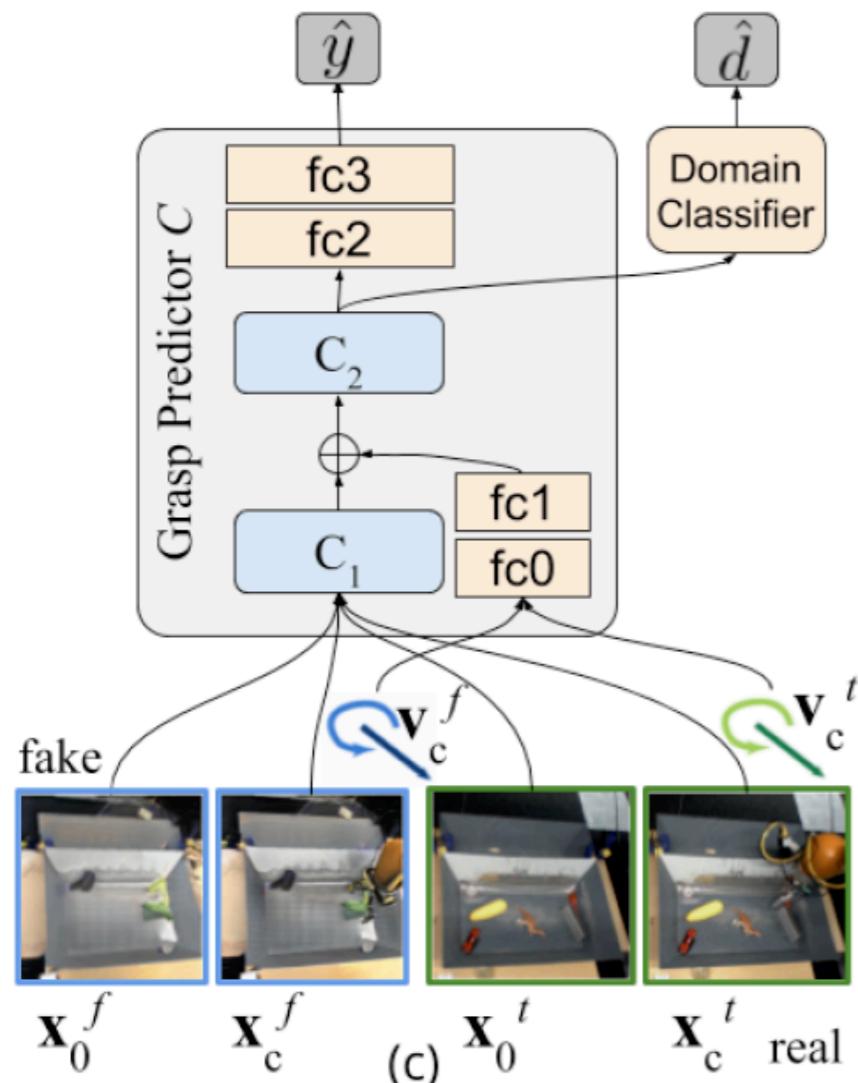


- Use Bullet simulator to emulate the Kuka hardware setup. Camera is mounted over the Kuka shoulder
- 51300 ShapeNet 3D models
- Use progressively better grasping models to collect data
- Randomization: both visuals and dynamics were randomized in simulation: the background image, object masses, textures, coefficients of friction.

Feature adaptation

Two losses: domain confusion loss and grasping prediction loss

Grasping prediction (task loss)



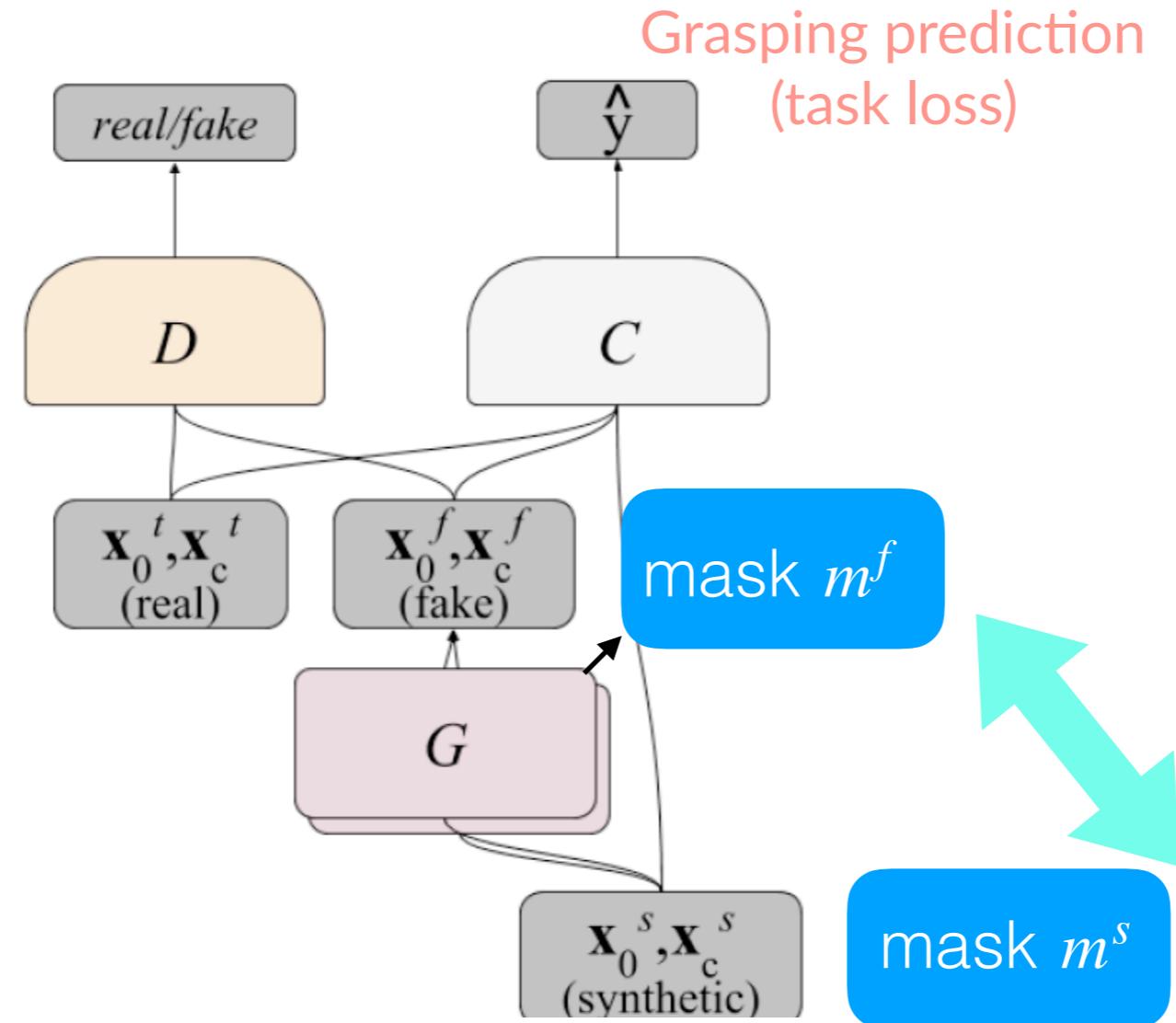
We add a domain classifier, that attempts to classify the domain the features come from

$$\mathcal{L}_{\text{DANN}} = \sum_{i=0}^{N_s+N_t} \{ d_i \log \hat{d}_i + (1 - d_i) \log (1 - \hat{d}_i) \}$$

The shared features C_1, C_2 attempt to confuse the domain classifier (maximize its loss), while the domain classifier features attempts to minimize its loss.

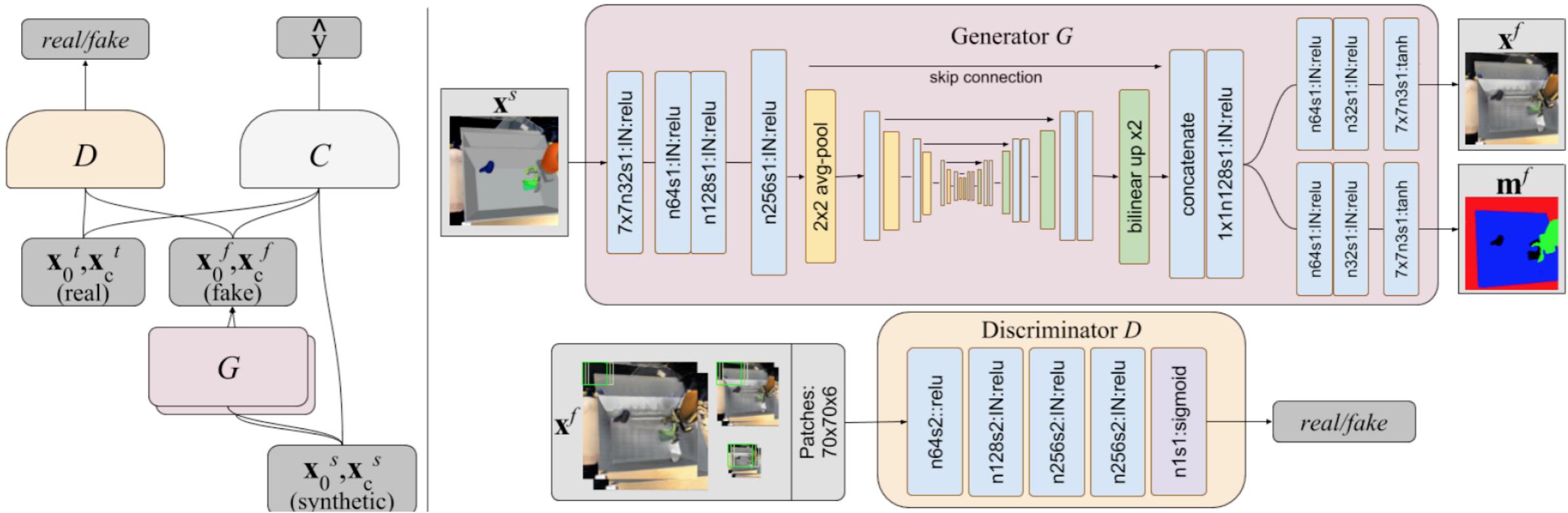
Pixel Adaptation

Three losses: grasping prediction loss, semantic labelling loss, adversarial loss



Pixel Adaptation

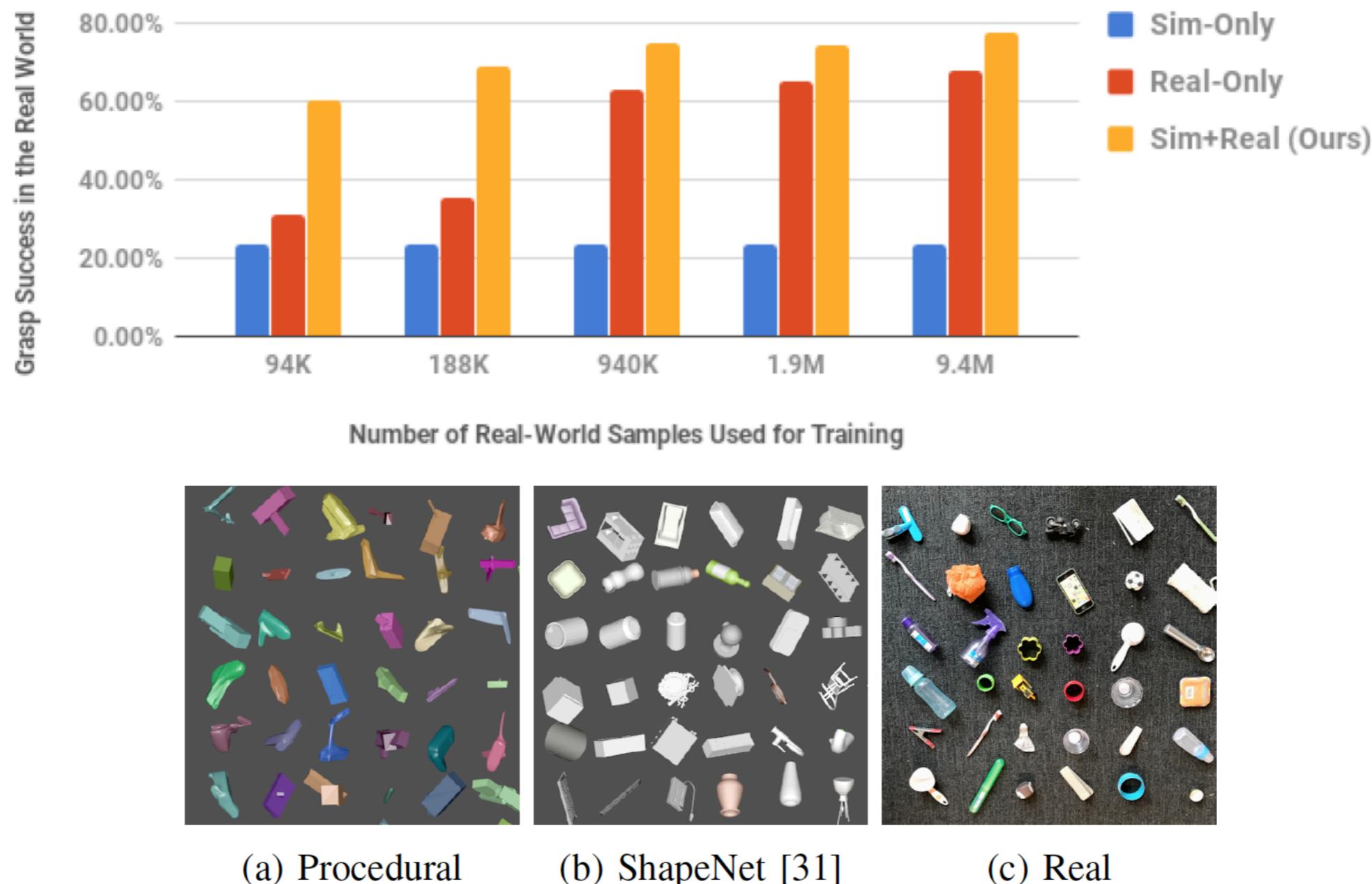
Three losses: grasping prediction loss, semantic labelling loss, adversarial loss



Goal: we want our generator to translate simulated images so that:

1. they do well in the task loss (grasping),
2. look real
3. retain the same semantics as their simulated counterparts

Results



What has shown to work

- Domain randomization (dynamics, visuals)
- Learning to adapt the textures of the simulator to match the real domain
- Learning to adapt the dynamics of the dymulator to match the real domain
- Learning policies not from pixels but rather from label maps-> semantic maps between simulation and real world are closer than textures
- Learning higher level policies, not low-level controllers, as the low level dynamics are very different between Sim and REAL

Driving Policy Transfer via Modularity and Abstraction

Matthias Müller
Visual Computing Center
KAUST, Saudi Arabia

Alexey Dosovitskiy
Intelligent Systems Lab
Intel Labs, Germany

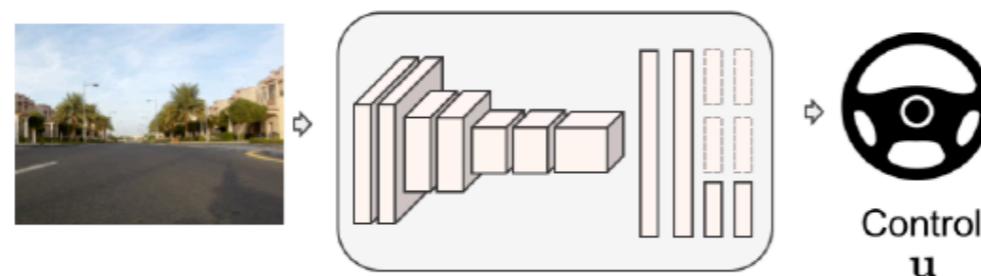
Bernard Ghanem
Visual Computing Center
KAUST, Saudi Arabia

Vladlen Koltun
Intelligent Systems Lab
Intel Labs, USA

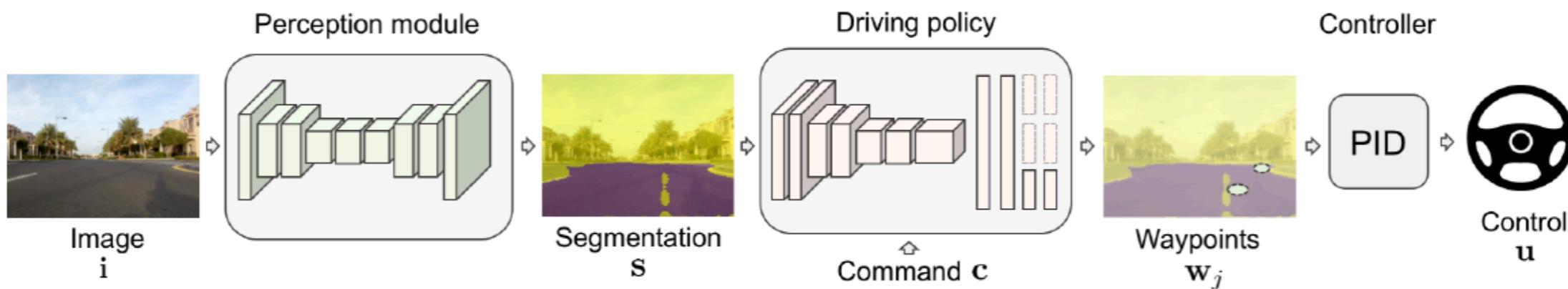
Idea: the driving policy is not directly exposed to raw perceptual input or low-level vehicle dynamics.

Main idea

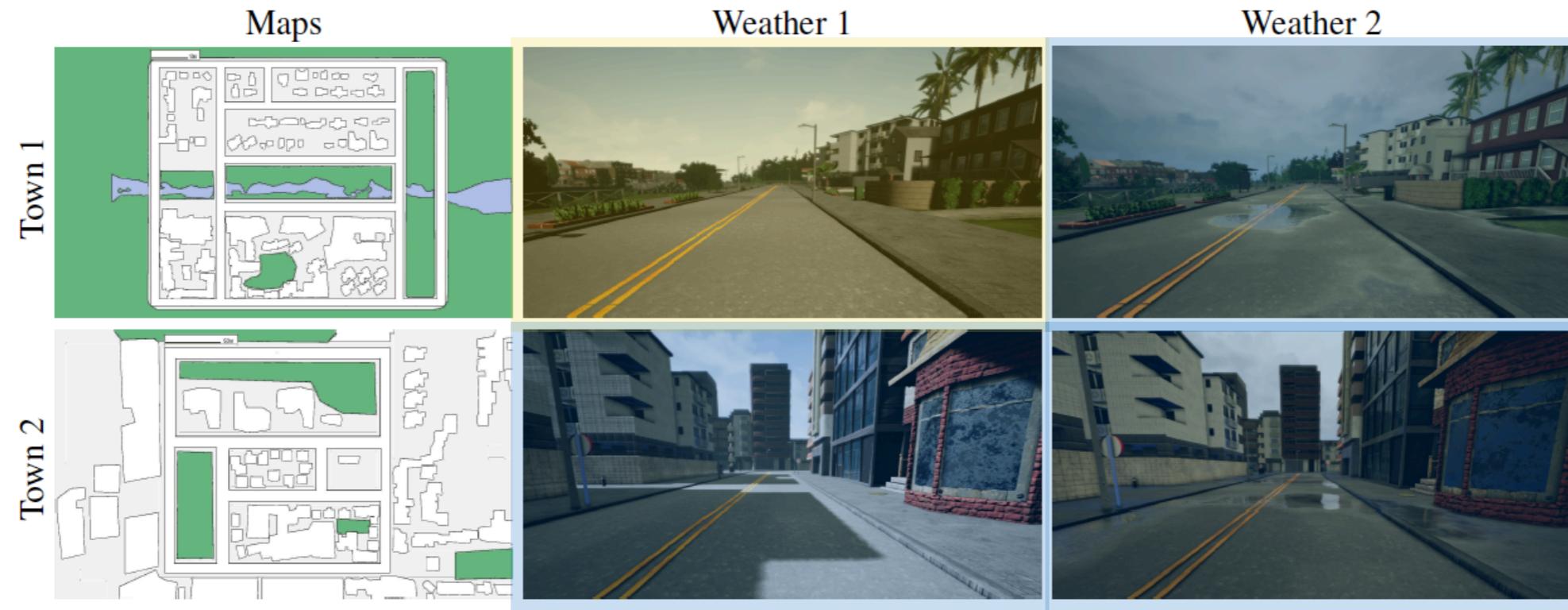
Pixels to steering wheel mapping is not SIM2REAL transferable: image textures and car dynamics mismatch



Instead: **label maps to waypoint** mapping is better SIM2REAL transferable:
label maps and waypoints are similar across SIM and REAL. A low-level
controller will take the car from waypoint to waypoint in the real world



Results: Train/Test



We train policies via behaviour cloning (standard regression loss) in Town1/ Weather1 dataset, and evaluate them on all four.

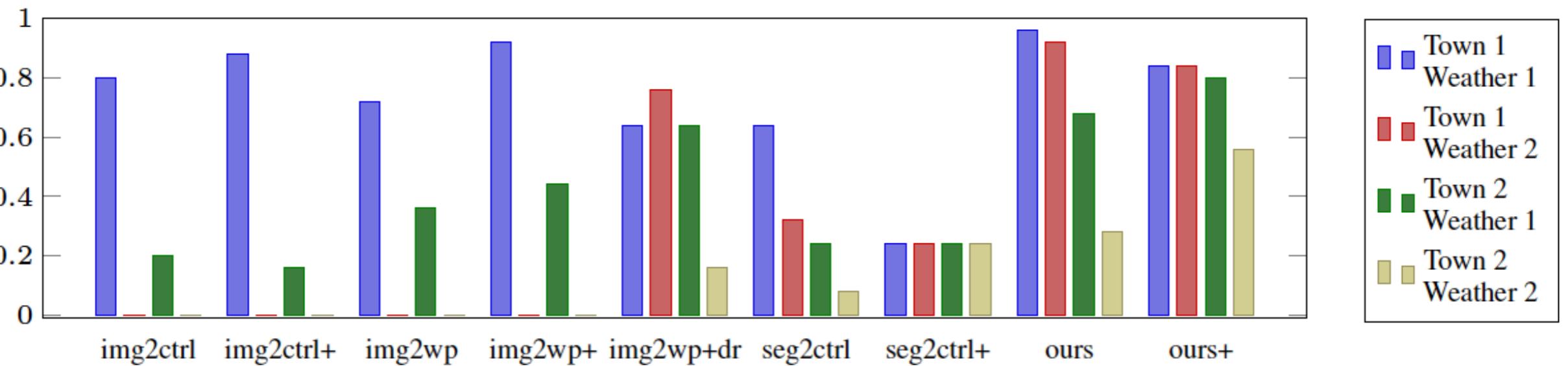


Figure 4: Quantitative evaluation of goal-directed navigation in simulation. We report the success rate over 25 navigation trials in four town-weather combinations. The models have been trained in Town 1 and Weather 1. The evaluated models are: *img2ctrl* – predicting low-level control from color images; *img2wp* – predicting waypoints from color images; *seg2ctrl* – predicting low-level control from the segmentation produced by the perception module; *ours* – predicting waypoints from the segmentation produced by the perception module. Suffix ‘+’ denotes models trained with data augmentation, and ‘+dr’ denotes the model trained with domain randomization.