

Appendix

A. Model Training Details

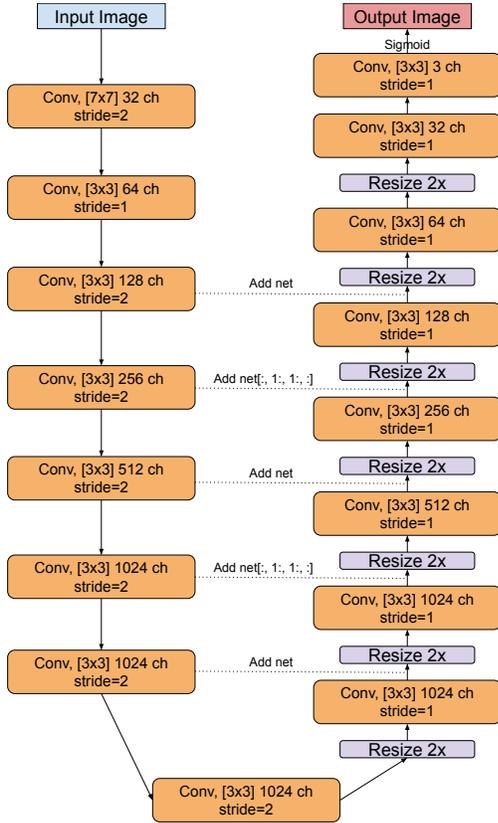


Figure 5. The U-net architecture used for G and F generator networks. Each convolutional block is shown with kernel size, number of filters and stride. Spectral normalization is applied to all convolutions. Images are resized to be twice as high and wide using nearest neighbor interpolation. Intermediate outputs from the down-convolutions (left) are added during the up-convolutions (right) as shown by the dotted lines, in two cases the first row and column are dropped during addition to match sizes. Instance normalization [32] is applied to all convolutions except the final output convolution.

We use data-sets consisting of grasping episodes from simulation and real robots. For Robot 1 and Robot 2 the data-sets are 580,000 and 80,000 real robot episodes respectively. Both data-sets are collected by starting with a human-designed scripted policy, which succeeds a small fraction of the time. Models are trained with this data, and periodically, those models are deployed to the robot to collect data from a better policy. When collecting data, random exploration noise is added to collect more diverse data. For this paper, we randomly subsample smaller datasets from these larger sets, to study the performance when using varying amounts of real episodes. For both setups several mil-

lion simulated episodes are also generated during on-policy training.

Typical episodes contain 6-10 states represented by a 512 pixel high, 640 pixel wide RGB image. To increase data diversity images are randomly cropped to 472x472 during training. At inference time, the center 472x472 square from the image is used. The generator for the GAN is a convolutional neural network with a U-Net architecture [28] as shown in 5. The discriminator is smaller convolutional neural network that operates on three scales of the input image. Both networks are described in detail in [3]. The robotic grasping task is trained via QT-Opt with the Q-function represented as a convolutional neural network (see [20] for architecture). RL-CycleGAN jointly trains the Cycle-GAN along with the Q-function during QT-Opt. Models are trained on Google TPUv3 Pod as in [5] and required *bfloat16* precision training to fit in memory. Each batch had 8 real images and 8 simulated images. We use Adam optimizer [21] with $\beta_1 = 0.1$ and $\beta_2 = 0.999$ and a constant learning rate of 0.0001. We employ Spectral Normalization [36] in the GAN generator networks and find that it improves stability.

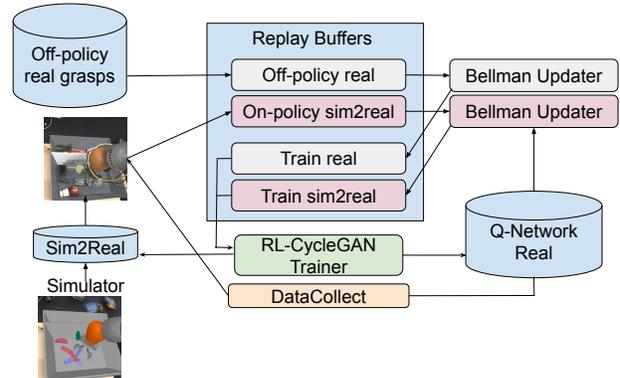


Figure 6. Training RL-CycleGAN via QT-Opt.

Figure 6 shows how we train RL-CycleGAN via QT-Opt. Images from a simulator are transformed by Sim2Real generator G and then passed to Q_{real} to generate an action. In this way, on-policy (w.r.t Q_{real}) episodes are generated in the simulation-to-real environment. Off-policy real grasps are read from disk. Separate replay buffers and bellman update instances are used for the off-policy real and off-policy simulation-to-real data. RL-CycleGAN is trained with batches with equal parts from real and simulation-to-real data. During training we evaluate the performance of both Q_{sim} and Q_{real} in the simulator, with simulation-to-real applied prior to evaluating Q_{real} . Training converges when certain conditions are met, the simulation-to-real images look realistic, the cycled images look reasonable with a reasonably low \mathcal{L}_{cyc} , both Q_{sim} and Q_{real} per-



Figure 8. Robot 2: simulated versions of 51 common real world objects are used when training the RL-CycleGAN , including plastic bottles (8), coffee cups (18), plastic utensils (3), drink cans (6), mugs (15), and wine glass (1).

jects are created: plastic bottles (8), coffee cups (18), plastic utensils (3), drink cans (6), mugs (15), and wine glass (1), shown in Figure 8. These do not cover all the real world objects used evaluation.