

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

DeepView: View synthesis with learned gradient descent

Supplemental Materials

Anonymous CVPR submission

Paper ID 2439

1. Reducing Memory Usage

A naive implementation of our model requires a prohibitive amount of RAM for both inference and training. Internally, convolutional layers within the CNNs compute activations with as many as 128 channels. The required RAM for the activations for even a single one of these layers would be $K \times D \times N_r \times N_c \times 128$ — for even moderately sized MPis this would require hundreds of gigabytes of GPU RAM, and this is for a single network layer within a single iteration of the network.

Reducing RAM during inference is relatively simple as we can tile the inference of each per-iteration CNN across the $D \times N_r \times N_c$ MPI volume, recomputing the gradient components for each view as needed. The tiled inference of the per-iteration CNNs can then be run on GPU. Using this strategy inference takes 50s for a 12 input $64 \times 980 \times 580$ MPI on a P100 GPU.

Reducing RAM at training time is more complex. The inference method described above is not directly applicable since we need to back-propagate the gradients. More importantly, a deep network typically requires tens or hundreds of thousands of iterations to converge and a per step time of 50s would make the training time impractical. Instead, during training the network is trained to produce a small 32×32 crop within a target image. By carefully considering both the CNN padding and the MPI and view geometry at each iteration we can compute both the needed crops from each input view and the minimal volume of the MPI that needs to be computed at each iteration in order to produce a given target image patch without border effects. Note that this calculation is complex since in order to produce the gradient components for a specific input view crop at iteration n we need to have available the MPI volume visible to that crop at the previous iteration $n - 1$, but in order to compute that MPI volume in $n - 1$ we need to have more of the input view's area etc. The required MPI volume that needs to be computed for a given target crop thus increases as the number of LGD iterations increases.

To further reduce RAM we discard activations within

each of the per-iteration networks, as described in [3] and tile computation along the depth plane dimension. Even with these memory optimizations, at higher resolutions we are limited to training a single example patch at a time and rely on synchronized replicas to get large effective batch sizes.

2. Generating training samples from the Spaces dataset

To generate a sample we first randomly select a scene and a random rig position. We then randomly select the input views from the set of possible view sets within the rig (see Figure 4). These views form the input to the network. The target view is then randomly selected from the remaining views across all rigs that are within 6cm of the convex hull of the input views, and a maximum of 7cm from the plane of the input views. A 32×32 crop is then chosen from the target view.

3. Training hyperparameters

We used distributed training with synchronized replicas to increase the effective batch size. We typically used 16 replicas, but for some experiments we used less replicas but computed gradients from two examples in each replica before accumulating the batch.

As in Adler *et al.* [1] we used global gradient clipping, with a threshold of 8.0.

We use feature similarity [2, 5] as our training loss \mathcal{L}_f , specifically the *conv1_2*, *conv2_2* and *conv3_3* layers of a pre-trained VGG-16 network [4]. The *conv3_3* layer has a receptive field of 40×40 , so we first reflect padded the 32×32 target and training crop to be 40×40 . Following Chen *et al.* [2] we then computed the loss by summing the L_1 difference of the layers, weighted by the empirically determined weights of [11.17, 35.04, 29.09]. Finally, we divided the loss by the area of the crop 32×32 , producing a loss in a more reasonable range which improved the numerical stability of global gradient clipping.

108	References	162
109		163
110	[1] J. Adler and O. Öktem. Learned primal-dual reconstruction.	164
111	<i>IEEE Transactions on Medical Imaging</i> , 37:1322–1332, 2018.	165
112	1	166
113	[2] Q. Chen and V. Koltun. Photographic image synthesis with	167
114	cascaded refinement networks. <i>CoRR</i> , abs/1707.09405, 2017.	168
115	1	169
116	[3] T. Chen, B. Xu, C. Zhang, and C. Guestrin. Training deep nets	170
117	with sublinear memory cost. <i>CoRR</i> , abs/1604.06174, 2016. 1	171
118	[4] K. Simonyan and A. Zisserman. Very deep convolutional	172
119	networks for large-scale image recognition. <i>CoRR</i> ,	173
120	abs/1409.1556, 2014. 1	174
121	[5] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang.	175
122	The unreasonable effectiveness of deep features as a perceptual	176
123	metric. <i>CoRR</i> , abs/1801.03924, 2018. 1	177
124		178
125		179
126		180
127		181
128		182
129		183
130		184
131		185
132		186
133		187
134		188
135		189
136		190
137		191
138		192
139		193
140		194
141		195
142		196
143		197
144		198
145		199
146		200
147		201
148		202
149		203
150		204
151		205
152		206
153		207
154		208
155		209
156		210
157		211
158		212
159		213
160		214
161		215