

Supplementary Material

A Prior on z_{pres}

Here we provide a derivation for the prior on z_{pres} (a vector of length HW made up of all z_{pres}^{ij}), discussed briefly in Section 4.2. Let C be a random variable giving the number of non-zero entries in z_{pres} . Under the prior, z_{pres} is generated by first sampling C , and then drawing z_{pres} uniformly from all binary vectors that have C non-zero entries.

$$p(z_{\text{pres}}) = p(C = nz(z_{\text{pres}}))p(z_{\text{pres}}|C = nz(z_{\text{pres}}))$$

where $nz(\cdot)$ gives the number of non-zero entries in a vector. There are $\binom{HW}{C}$ binary vectors of length HW with C non-zero entries, so:

$$p(z_{\text{pres}}) = p(C = nz(z_{\text{pres}})) \binom{HW}{nz(z_{\text{pres}})}^{-1}$$

For the distribution over C we use a Geometric distribution with parameter s . We use the Geometric interpretation that puts C as the number of failures before a success, and s as the success probability. By setting s to a high value, we put most of the probability mass at low values, thereby putting pressure on the network to explain images using as few objects as possible.

$$p(C = nz(z_{\text{pres}})) = s(1-s)^{nz(z_{\text{pres}})}$$

We then truncate and normalize to ensure C has support $\{0, 1, \dots, HW\}$:

$$\begin{aligned} p(C = nz(z_{\text{pres}})) &= \frac{s(1-s)^{nz(z_{\text{pres}})}}{\sum_{c=0}^{HW} s(1-s)^c} \\ &= \frac{s(1-s)^{nz(z_{\text{pres}})}}{s \frac{(1-(1-s)^{HW+1})}{1-(1-s)}} \\ &= \frac{s(1-s)^{nz(z_{\text{pres}})}}{1-(1-s)^{HW+1}} \end{aligned}$$

Overall, the prior is:

$$p(z_{\text{pres}}) = \frac{s(1-s)^{nz(z_{\text{pres}})}}{(1-(1-s)^{HW+1}) \binom{HW}{nz(z_{\text{pres}})}}$$

B KL Divergence for z_{pres}

Having derived the prior, we now turn to efficient computation of the KL divergence between the distribution over z_{pres}

yielded by the network, namely $q_{\phi}(z_{\text{pres}}|x)$ (embodied in the β_{pres}^{ij} variables from Section 4.3) and $p(z_{\text{pres}})$. Efficient computation of this KL divergence is necessary for computing the second term of the VAE training objective (Equation 3). For convenience, we switch to indexing the entries of z_{pres} using a single index k rather than the grid indices ij . We will use the notation $z_{\text{pres}}^{m:n}$ to mean the sub-vector ranging from index m to n , inclusive.

First, from basic properties of KL divergence we have:

$$D_{KL}(q(z_{\text{pres}}|x) \parallel p(z_{\text{pres}})) = \sum_{k=1}^n E [D_{KL}(q(z_{\text{pres}}^k | z_{\text{pres}}^{1:k-1}, x) \parallel p(z_{\text{pres}}^k | z_{\text{pres}}^{1:k-1}))]$$

where the expectation for index k is taken over variables $z_{\text{pres}}^{1:k-1}$ sampled from the marginal $q(z_{\text{pres}}^{1:k-1}|x)$. We estimate each expectation using a single sample:

$$\approx \sum_{k=1}^n D_{KL}(q(z_{\text{pres}}^k | \hat{z}_{\text{pres}}^{1:k-1}, x) \parallel p(z_{\text{pres}}^k | \hat{z}_{\text{pres}}^{1:k-1}))$$

where \hat{z} indicates values that have been sampled. Note that $q(z_{\text{pres}}^k | \hat{z}_{\text{pres}}^{1:k-1}, x)$ is just Bernoulli(β_{pres}^k) from Section 4.3.

We now turn to computation of $p(z_{\text{pres}}^k | \hat{z}_{\text{pres}}^{1:k-1})$. Recalling that C is a random variable giving the number of non-zero entries in z_{pres} , we have:

$$p(z_{\text{pres}}^k | \hat{z}_{\text{pres}}^{1:k-1}) = \sum_{c=0}^{HW} p(z_{\text{pres}}^k | \hat{z}_{\text{pres}}^{1:k-1}, C) p(C | \hat{z}_{\text{pres}}^{1:k-1}) \quad (\text{B1})$$

We focus first on the factor $p(z_{\text{pres}}^k | \hat{z}_{\text{pres}}^{1:k-1}, C)$. Given that we have sampled $\hat{z}_{\text{pres}}^{1:k-1}$, we still need $C - nz(\hat{z}_{\text{pres}}^{1:k-1})$ non-zeros, and we have $HW - (k-1)$ ‘‘slots’’ to get them with. Recalling that given C , we are assuming uniform sampling of all binary vectors with that number of non-zeros, we have:

$$p(z_{\text{pres}}^k | \hat{z}_{\text{pres}}^{1:k-1}, C) = \frac{C - nz(\hat{z}_{\text{pres}}^{1:k-1})}{HW - (k-1)}$$

The final piece that we need is $p(C | \hat{z}_{\text{pres}}^{1:k-1})$. This can be decomposed recursively as:

$$p(C | \hat{z}_{\text{pres}}^{1:k-1}) \propto p(\hat{z}_{\text{pres}}^{k-1} | C, \hat{z}_{\text{pres}}^{1:k-2}) p(C | \hat{z}_{\text{pres}}^{1:k-2}) \quad (\text{B2})$$

Starting from $k=1$, we first evaluate Equation (B1), then sample z^k from $q_{\phi}(z_{\text{pres}}^k | \hat{z}_{\text{pres}}^{1:k-1})$, and finally update the conditional count distribution $p(C | \hat{z}_{\text{pres}}^{1:k})$ via Equation (B2).

C Object Rendering in the Decoder Network

Here we describe in detail the differentiable rendering algorithm implemented in the decoder network. For each object we have the following values:

1. Object RGB map o^{ij} with shape $(H_{obj}, W_{obj}, 3)$
2. Object transparency map α^{ij} with shape $(H_{obj}, W_{obj}, 1)$
3. Object bounding box b^{ij} derived from z_{where}^{ij} (see Section 4.1 and Figure 1).
4. Relative depth value z_{depth}^{ij}
5. Presence value z_{pres}^{ij}

We also have a background image x_{bg} with shape $(H_{\text{img}}, W_{\text{img}}, 3)$, which may have been predicted from the input image by a neural network, or may be a fixed image calculated from the dataset in some other way.

We first perform two initial computations for each object:

$$\tilde{\alpha}^{ij} = \alpha^{ij} z_{\text{pres}}^{ij}$$

$$\gamma^{ij} = \alpha^{ij} z_{\text{pres}}^{ij} \sigma(-z_{\text{depth}}^{ij})$$

where σ is the sigmoid function. $\tilde{\alpha}^{ij}$ is a transparency map that takes into account whether the object exists, γ^{ij} is an *importance* map which is used to implement a differentiable approximation of relative depth.

For each pixel in the output image, we iterate over all objects, checking whether each object affects the current pixel (i.e. whether the pixel is inside the bounding box for the object). For all objects that affect the pixel, three values are extracted from the object: RGB values o' , a transparency value α' and an importance value γ' . These values are extracted by finding the position of the pixel with respect to the object's coordinate frame, and then using bilinear interpolation to extract a value from the relevant map. The RGB value o' is then mixed with the background using α' . The mixed values from all affecting objects are then mixed with one another using a convex combination consisting of the γ' values normalized to sum to 1 (within the pixel). This mixing step implements the differentiable approximation of relative depth; objects with higher importance values (lower depth) get a larger "share" of the pixels that they affect, and appear on top of objects with lower importance (higher depth). Pseudo-code for this algorithm is given in Algorithm 1.

Algorithm 1 Differentiable rendering algorithm

```

1: function CONTAINS( $b, (y, x)$ )
   Return True iff pixel  $(y, x)$  is inside bounding box  $b$ 
2: end function

3: function INTERPOLATE( $b, v, (y, x)$ )
   Assuming pixel  $(y, x)$  is inside bounding box  $b$ , extract
   a value from map  $v$  for pixel  $(y, x)$  using bilinear inter-
   polation.
4: end function

5: function DIFFERENTIABLERENDERING
6:    $x_{\text{out}} \leftarrow \text{COPY}(x_{\text{bg}})$ 
7:   for pixel with position  $(y, x)$  do
8:     n-writes  $\leftarrow 0$ 
9:     weighted-sum  $\leftarrow 0$ 
10:    importance-sum  $\leftarrow 0$ 
11:    for object with index  $(i, j)$  do
12:      if CONTAINS( $b^{ij}, (y, x)$ ) then
13:         $\alpha' \leftarrow \text{INTERPOLATE}(b^{ij}, \tilde{\alpha}^{ij}, (y, x))$ 
14:         $\gamma' \leftarrow \text{INTERPOLATE}(b^{ij}, \gamma^{ij}, (y, x))$ 
15:         $o' \leftarrow \text{INTERPOLATE}(b^{ij}, o^{ij}, (y, x))$ 
16:         $o' \leftarrow \alpha' o' + (1 - \alpha') x_{\text{bg}}[y, x]$ 
17:        n-writes  $\leftarrow$  n-writes+1
18:        weighted-sum  $\leftarrow$  weighted-sum +  $\gamma' o'$ 
19:        importance-sum  $\leftarrow$  importance-sum +  $\gamma'$ 
20:      end if
21:    end for
22:    if n-writes > 0 then
23:       $x_{\text{out}}[y, x] \leftarrow \frac{\text{weighted-sum}}{\text{importance-sum}}$ 
24:    end if
25:  end for
26:  return  $x_{\text{out}}$ 
27: end function

```
