## A. Supplementary Material

First, we show a list of hyper-parameters and implementation details for all of our adaptation methods in Section A.1. Then we show some samples of images and their corresponding per-pixel masks, along with the verification algorithm for counting and occlusions in Section A.2. Then we show some graph samples from our pool of manually designed scenes for the W-VQA dataset and describe their functionality for our automatic triplet (IQA) generation in Section A.3. Finally in Sections A.4 and A.5 we show some samples from W-VQA and H-VQA we randomly select from a diverse set of scenes, with different backgrounds, camera position and illumination.

#### A.1. Hyper-parameter Selection

The following are all the hyper-parameter selection for all of our algorithms: lr refers to learning rate, E to the number of training epochs, O to the optimizer type,  $O_{wd}$  is the optimizer weight decay,  $O_{\epsilon}$  is the term added to the denominator to improve numerical stability,  $O_{\beta}$  are a tuple of coefficients used for computing running averages of gradient and its square. For the Adversarial and MMD methods, the auto-encoder network (AE) is trained separately, in a 2 step format following Zhang et al. [67] Two-stage DA; in both cases the first number in E refers to the training epoch parameter for the AE. For Domain Independent,  $di_{tokens}$  is the additional output we use for the synthetic answer tokens.

Adversarial		MMD	
lr = 15e - 4	E = 100 + 13	lr = 1e - 3	E = 150 + 13
$O_{wd} = 1e - 6$	O = Adam	$O_{wd} = 1e - 4$	O = Adam
$O_{\epsilon} = 1e - 4$	$O_{\beta} = (0.8, 0.8)$	$O_{\epsilon} = 1e - 4$	$O_{\beta} = (0.8, 0.8)$
$\alpha = \frac{2}{(1 + exp(-10*p)) - 1}$		$\alpha = 0.4$	$\beta = 0.6$
Domain Inde	pendent	F-S	SWAP
Domain Indep $lr = 15e - 4$	pendent $E = 13$	F-S = lr = 15e - 4	SWAP $E = 13$
$\begin{array}{c} \text{Domain Inder}\\ lr = 15e-4\\ O_{wd} = 0.2 \end{array}$	E = 13 $O = Adam$	$F-S$ $lr = 15e - 4$ $O_{wd} = 1e - 1$	E = 13 $O = Adam$
$\begin{array}{c} \mbox{Domain Index}\\ lr = 15e-4\\ O_{wd} = 0.2\\ O_{\epsilon} = 1e-9 \end{array}$	$E = 13$ $O = Adam$ $O_{\beta} = (0.9, 0.9)$	$F-S$ $lr = 15e - 4$ $O_{wd} = 1e - 1$ $O_{\epsilon} = 1e - 9$	

Table 5. Hyper-parameter selection details for all methods.

### A.2. RGB and Mask Samples

ThreeDWorld (TDW) [14] allows to capture the RGB images from the camera view along with the id and category per-pixel semantic masks, which we later use to verify the number of objects in the image and avoid object occlusions. Figure 6 shows some samples we randomly select from our generated W-VQA set. The first column correspond to the RGB image, the second and third columns correspond to the category and id masks respectively. We verify if an object overlaps to another and assess the object counts by computing the intersection over union.



Figure 6. Random samples from the images we generate using TDW along with their category masks (second row) and id masks (third row).

#### A.3. Scene-Graph Samples

Let E denote the set of scene entities and consider the set of binary relations R. Then a scene graph  $SG \in E \times R \times E$ is a collection of ordered triplets (o, p, o) = object, position, and object. For example, as shown in the first sample in Figure 7, with A=lamp, B=table, C=backpack, the triplet (A, position, B) indicates that a lamp is on top of the table, or the table is under the lamp. Similarly, the triplet (B, position, C) indicates that the backpack is to the left of the table, or the table is to the right of the backpack. In this way, from a relationship, there are at least two possible positions,  $p \wedge p^{-1}$ , e.g., p = left and  $p^{-1} = \text{right}$ . When sampling from these graphs, each node in E could also be assigned three different attributes: the number of objects to appear in the same scene n = randrange(20), the color, and material type which are selected from a list of available materials and colors from the set of Records in TDW  $^{3}$ .



Figure 7. Some of the scene graphs designed for our automated synthetic dataset generation. While generating images, we select one graph and randomly select the number of objects per position node := [A, B, C, D], it's color and materials. Then we use the grammar introduced in Section 3.2 to generate the questions and corresponding answers.

<sup>&</sup>lt;sup>3</sup>https://github.com/threedworld-mit/tdw/blob/ master / Documentation / misc \_ frontend / materials \_ textures\_colors.md

# A.4. W-VQA Generated Samples

We show some random samples we generate for our W-VQA dataset in Figure 8, following Section 3.2.



Figure 8. Additional samples of our W-VQA dataset. The first row showcase simple configurations using the same background. The second row shows diverse compositions using indoor scenes. The third row shows compositions of challenging counting questions. The fourth row shows outdoor objects and scenes. Finally, the fifth row shows materials and color related questions using the same object in different camera positions. Best viewed in color.

## A.5. H-VQA Generated Samples

We show some random samples we generate for our H-VQA dataset in Figure 9.



Q/ How many sofas are in the image? A/ 3 Q/ What is on top of the table? A/ a notebook Q/ What is next to the table? A/ a red sofa Q/ Is there a chair in the room? A/ no



Q/ How many plants are in the image? A/ 2 Q/ How many chairs are in the room? A/ 3 Q/ How many paintings are in the room? A/ 1 Q/ Is there a sofa in the living room? A/ yes



Q/ How many paintings are in the image? A/ 2 Q/ How many books are in the image? A/ 9 Q/ How many sofas are in the image? A/ 3 Q/ Is there a chair in the picture? A/ no



Q/ How many televisions are in the bathroom? A/ 0 Q/ How many plants are in the room? A/ 0 Q/ How many towels are in the bathroom? A/ 1 Q/ Is there a trash can in the image? A/ yes



Q/ How many bottles are in the image? A/ 2 Q/ How many plants are in the image? A/ 1 Q/ Is there a sofa in the room? A/ yes Q/ Is there a statue in the waiting room? A/ no



Q/ How many frames are in the image? A/ 8 Q/ How many lamps are in the bedroom? A/ 2 Q/ How many tables are in the bedroom? A/ 2 Q/ Is there a bicycle in the room? A/ no



Q/ How many chairs are in the room? A/ 3 Q/ How many shelves are in the image? A/ 3 Q/ How many plants are in the image? A/ 0 Q/ is there a computer in the office? A/ yes



Q/ How many glasses are in the kitchen? A/ 3 Q/ How many chairs are in the kitchen? A/ 2 Q/ How many plants are in the kitchen? A/ 0 Q/ Is there a refrigerator in the picture? A/ yes



Q/ How many books are in the room? A/ 20 Q/ How many pillows are in the image? A/ 1 Q/ How many sofas are in the picture? A/ 1 Q/ Is there a plant in the image? A/ yes



Q/ How many plants are in the room? A/ 1 Q/ How many pillows are in the image? A/ 1 Q/ Is there a table in the room? A/ yes Q/ Is there a desk in the room? A/ no







Q/ How many laptops are in the image? A/ 1 Q/ How many speakers are in the image? A/ 5 Q/ How many vases are in the image? A/ 0 Q/ Is there a table in the room? A/ yes



Q/ How many sofas are in the image? A/ 1 Q/ How many bags are in the image? A/ 0 Q/ How many tables are in the image? A/ 1 Q/ Is there a chair can in the picture? A/ yes



Q/ How many chairs are in the room? A/ 4 Q/ How many sofas are in the image? A/ 0 Q/ How many beds are in the image? A/ 0 Q/ Is there a table in the picture? A/ yes



Q/ How many chairs are in the room? A/ 5 Q/ How many plants are in the image? A/ 0 Q/ Is there a picture frame in the room? A/ yes Q/ Is there a magazine in the waiting room? A/ no



Q/ How many doors are in the image? A/ 0 Q/ What is on top of the shelf? A/ a frame Q/ Is there a paint in the room? A/ yes Q/ Is there a sofa in the room? A/ no



Q/ How many towels are in the image? A/ 1 Q/ How many lamps are in the bathroom? A/ 0 Q/ How many ducks are in the image? A/ 1 Q/ Is there a soap in the room? A/ yes



Q/ How many pears are in the image? A/ 8 Q/ How many plants are in the office? A/ 1 Q/ How many chairs are in the office? A/ 0 Q/ Is there a chair in the picture? A/ no



Q/ How many towels are in the bathroom? A/ 1 Q/ How many toothbrushes are in the image? A/ 2 Q/ How many soaps are in the bathroom? A/ 0 Q/ Is there a bathtub in the image? A/ yes



Q/ How many lemons are in the picture? A/ 3 Q/ How many sofas are in the image? A/ 2 Q/ Is there a water dispenser in the room? A/ yes Q/ Is there a laptop in the room? A/ yes

Figure 9. Additional samples of our H-VQA dataset. We generate questions and answers from manual and existing semantic annotations from Hypersim [45]. Best viewed in color.