

Learning Where to Look: Data-Driven Viewpoint Set Selection for 3D Scenes

Kyle Genova, Manolis Savva, Angel X. Chang, and Thomas Funkhouser

Princeton University

Abstract

The use of rendered images, whether from completely synthetic datasets or from 3D reconstructions, is increasingly prevalent in vision tasks. However, little attention has been given to how the selection of viewpoints affects the performance of rendered training sets. In this paper, we propose a data-driven approach to view set selection. Given a set of example images, we extract statistics describing their contents and generate a set of views matching the distribution of those statistics. Motivated by semantic segmentation tasks, we model the spatial distribution of each semantic object category within an image view volume. We provide a search algorithm that generates a sampling of likely candidate views according to the example distribution, and a set selection algorithm that chooses a subset of the candidates that jointly cover the example distribution. Results of experiments with these algorithms on SUNCG indicate that they are indeed able to produce view distributions similar to an example set from NYUDv2 according to the earth mover's distance. Furthermore, the selected views improve performance on semantic segmentation compared to alternative view selection algorithms.

1. Introduction

Rendering of 3D scenes, both synthetic and reconstructed, is a promising way to generate training data for deep learning methods in computer vision. For example, several methods have trained on rendered images to improve the performance of semantic segmentation networks [40]. The extent to which the rendered images can be useful as training data depends on the quality of the match between the rendered data and real-world images. For example, when rendering images for training semantic segmentation algorithms, the statistics of object occurrences, sizes, and placements influence the outcome.

Yet, little previous work has investigated algorithms to select camera views for the purposes of generating training sets for vision tasks. Most previous work on view se-

lection has focused on optimizing the aesthetic properties of rendered images for applications in computer graphics [1, 2, 5, 10, 26, 19] or visible surface coverage for surveillance [21] and surface reconstruction [41, 31].

In this paper, we investigate a new view set selection problem motivated by generating training sets for computer vision tasks. We pose the following view selection problem: select a set of views whose “distribution of image content” best matches a set of example images. This problem statement defines the objective as matching a latent distribution generating the example set, rather than optimizing a particular function and thus is quite different than previous work. This requires choosing a representative example set, defining a “distribution of image content,” and searching the infinite space of view sets for the best match.

In this paper, we investigate one concrete instance of the problem motivated by generating training data for semantic segmentation: *to select a set of views where the pixel distribution of object observations for specific semantic categories matches that of an example image set.* By matching the spatial distributions of object categories in example images, we may be able to train deep networks more effectively than is possible now with cameras placed heuristically [40], along trajectories [23], or randomly [13].

We propose an algorithm with three components to address the problem. The first reduces a candidate image into a low-dimensional representation that models the distances between the spatial distributions of pixel sets for each semantic category. The second suggests candidate positions in a 3D computer graphics scene that are likely to yield rendered images with a given pixel distribution of semantic categories. The third uses submodular optimization to select a set amongst the candidate views to match the overall latent distribution of the example set.

Results of experiments with this algorithm demonstrate that it is practical, efficient, and more effective than traditional approaches at selecting camera views for the SUNCG dataset [36] when asked to match the distribution of views from NYUDv2 [34]. We find that our selected views are quantitatively more similar to NYUDv2's than previous ap-

proaches attempting to model the distribution heuristically, and we show that training an FCN network [20] on them provides better performance than for other view sets. We make the following contributions:

- We are the first to focus on the problem of viewpoint selection for synthetic dataset generation in the context of data-hungry tasks in computer vision.
- We present a model for viewpoint selection that matches a data prior over object occurrence patterns specified through example images of the real world.
- We introduce an algorithm based on submodular optimization for selecting a set of images approximately optimizing a function measuring deviation from an statistical distribution derived from examples.
- We present results demonstrating that our agent-agnostic algorithms for view set selection outperform previous alternatives that are manually tuned to match a particular agent.

2. Related work

View selection has been studied in several contexts, including camera placement in graphics, next-best view prediction and sensor placement in robotics, canonical views in perceptual psychology, among many others.

Camera optimization in computer graphics. Early work in graphics has used information entropy as a measure of view quality to select good views of objects [38]. Another line of work encodes image composition principles to optimize virtual camera placements such that they focus on specified objects in the frame and are judged to be aesthetically pleasing by people [1, 2, 5, 10, 26]. More recent work has taken a similar approach in the more specific scenario of product images [19]. Freitag et al. [9] compare several viewpoint quality estimation metrics on two virtual scenes, though none of the metrics involve matching real-world data semantics. A related and rich body of work studies automated camera path planning for real-time graphics systems — a good survey is given by Jankowski and Hachet [15]. Generally, these systems use manually encoded heuristics to create smooth camera trajectories. Unlike our work, the focus is not selecting a set of static views, or on matching real-world view statistics based on object semantics. In general, related work in computer graphics has not considered modeling camera priors based on real-world image data, and it has not evaluated the impact of camera viewpoints on computer vision tasks.

Next-best view prediction and camera placement optimization. Next-best view prediction has been addressed in the context of various problems in robotics and computer vision: 3D reconstruction [17], volumetric occupancy prediction [39], and object pose estimation [7] among others.

However, the next-best view problem is predominantly addressed in an active sensing context where planning for the next most informative view is desired. In contrast, we focus on the offline problem of selecting a set of static viewpoints given a 3D indoor scene dataset as input. Therefore, we are more closely related to work in camera placement optimization. The camera placement problem is typically cast as a version of the art gallery problem, seeking to minimize the number of cameras needed to ensure visual coverage of a target environment. There is much prior work in this area — a recent survey is provided by Mavrinac and Chen [21]. Similar problem formulations have been used in robotics for view planning of 3D sensing [3]. Our problem statement is distinct since we match the distributions of semantically meaningful objects instead of simply optimizing for coverage of a single input environment.

Canonical views of objects and scenes. Internet image collections of particular objects were shown to mirror canonical viewpoint preferences by people [24]. Other work has conditioned the camera viewpoint on object keypoints to jointly predict the 3D object pose and the viewpoint [37]. Judgments of preferred views collected through crowdsourcing have been used to train a predictor for preferred viewpoints of objects [32]. Ehinger and Oliva [8] ask people to select a “desirable” (i.e., canonical) orientation within 2D panoramas and show that chosen views correlate highly with directions that maximize the visible volume and also coincide with prominent navigatable paths. Our approach assumes that images taken by people for inclusion in image datasets exhibit a prior over natural viewpoints that we can extract and leverage to select views in 3D scenes.

Rendering synthetic 3D scenes for data generation. There is a recent explosion in generation of synthetic training data for many vision tasks including tracking, object recognition, semantic segmentation, pose estimation, and optical flow among others [4, 27, 12, 14, 13, 22, 30, 33, 29]. Work that rendered synthetic 3D images has focused on domains where the camera viewpoint is highly constrained (e.g., driving [30, 33, 29]), has used manually specified camera trajectories (e.g., by recording first-person user trajectories in the virtual scene, or from camera paths in the real world [12, 14]), or has purportedly randomly sampled the space of camera views [13]. No previous paper has focused on view set selection or investigated its impact on trained models.

3. Data-driven viewpoint selection

In this paper, we investigate a new data-driven approach to view selection in synthetic scenes. The input to our system is a 3D computer graphics model of an indoor scene (e.g., from SUNCG [36]) and a set of example RGB-D images containing semantic segmentations (e.g., from

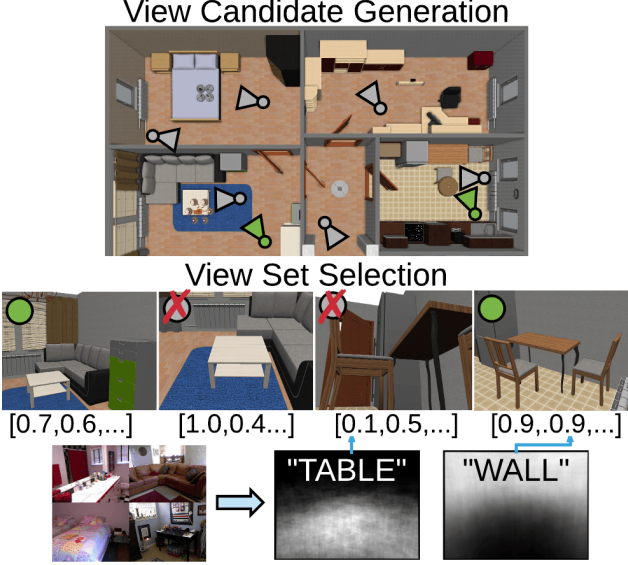


Figure 1: **Top:** candidate views in a target 3D scene with the selected output view set in green. **Bottom:** category pdfs (right) from example set (left) are used to estimate image likelihood along multiple semantic axes. **Middle:** a set of views is selected to jointly match the object distribution statistics of the input.

NYUDv2 [34]). The output is a set of views (6 DoF camera poses), where the spatial pixel distribution of semantic objects in images rendered from the views ideally matches the distribution observed in the example set (Figure 1).

We model the views of the example set as random samples drawn from a latent distribution. To model that distribution, we represent it as a set of n probability density functions (pdfs) representing the x , y , and depth positions of pixels for each of n semantic classes observed in the example set. This pixel-level representation captures the spatial layout of different object classes in the example images, as depicted in Figure 1 (bottom right shows a 2D representation of the 3D histogram associated with “table”).

We select a set of views approximately covering an estimate of that distribution with the two step process depicted in Algorithm 1. During the first step, we generate candidate camera positions by sampling viewpoints according to pdfs dependent only on depths and categories of visible objects (i.e., the original pdfs with rotations and pixel xy marginalized). During the second step, we choose a set of candidate views that approximately match the entire pdf as a set using an algorithm based on submodular maximization. The following three subsections describe the pdf representation and the two algorithmic steps in detail.

3.1. Candidate representation

The first issue to address with this data-driven approach is to choose a suitable representation for a view. Ideally, the

Algorithm 1 Viewpoint Generation with RGBD Priors

```

1: function GENERATE(Scene,Room)
2:   Voxelize Room into 1000-10000 voxels  $\mathcal{V}$ 
3:   for  $v \in \mathcal{V}$  do
4:      $W_{v,:}^{\mathcal{V}} \leftarrow 0$ 
5:      $\mathcal{K}_c \leftarrow 0$ 
6:     for  $i \in [\#Samples]$  do
7:       Sample Point  $p^o \sim \mathcal{U}(v)$ 
8:       Sample Direction  $d^o \sim \mathcal{A}$ 
9:       Ray  $r \leftarrow (p^o, d^o)$ 
10:      if  $r$ .Intersects(Scene) then
11:         $p \leftarrow \text{IntersectionPoint}()$ 
12:         $c \leftarrow \text{ObjectAt}(p)$ .Category()
13:         $d \leftarrow \|p - r_p\|_2$ 
14:         $W_{v,c}^{\mathcal{V}} \leftarrow W_{v,c}^{\mathcal{V}} + \sum_x \sum_y f_c(x, y, d)$ 
15:         $\mathcal{K}_c \leftarrow \mathcal{K}_c + 1$ 
16:       $W_v^{\mathcal{V}} \leftarrow \sum_{c \in \mathcal{C}} \frac{w_c}{\mathcal{K}_c} W_{v,c}^{\mathcal{V}}$ 
17:      for  $i \in [\#CandidateCameras]$  do
18:        Select Voxel  $v \sim W^{\mathcal{V}}$ 
19:        Sample Eye  $e \sim \mathcal{U}(v)$ 
20:        Sample Gaze  $g \sim \mathcal{A}$ 
21:        Camera  $c_i \leftarrow (e, g)$ 
22:        Image  $I \in \mathcal{I} \leftarrow \text{Render}(\text{Scene}, c_i)$ 
23:         $W_i^{\mathcal{I}} \leftarrow \sum_{(x,y,d,c) \in \mathcal{I}} f_c(x, y, d)$ 
24:        for  $c \in \mathcal{C}$  do
25:           $W_{i,c}^{\mathcal{I}} \leftarrow \frac{w_c \sum_{(x,y,d,c') \in \mathcal{I}} f_c(x, y, d) \mathbb{1}(c'=c)}{\sum_{(x,y,d,c') \in \mathcal{I}} \mathbb{1}(c'=c)}$ 
26:        Order  $I_i$  by descending  $\sum_{c \in \mathcal{C}} W_{i,c}^{\mathcal{I}}$ 
27:        return  $I_1 \dots I_{threshold}$ 
28: function SELECT( $I \in \mathcal{I}, W_{I,c}^{\mathcal{I}}, h_c, k$ )
29:   for  $c \in \mathcal{C}$  do
30:      $F_c \leftarrow \text{MinHeap}()$ 
31:     for  $i \in [h_c]$  do
32:       Push 0 onto  $F_c$ 
33:    $S \leftarrow \{\}$ 
34:   for  $I \in \mathcal{I}$  do
35:      $\Delta(I|S) \leftarrow \sum_{c \in \mathcal{C}} W_{i,c}^{\mathcal{I}}$ 
36:   for  $i \in [k]$  do
37:     Order  $\mathcal{I}$  by descending  $\Delta(I|S)$  as  $\mathcal{I}_1 \dots \mathcal{I}_{|\mathcal{I}|}$ 
38:     for  $i \in |\mathcal{I}|$  do
39:        $\Delta(\mathcal{I}_i|S) \leftarrow \sum_{c \in \mathcal{C}} \max(W_{i,c}^{\mathcal{I}} - F_c.\text{top}(), 0)$ 
40:       if  $\Delta(\mathcal{I}_i|S) > \Delta(\mathcal{I}_{i+1}|S)$  then
41:         for  $c \in \mathcal{C}$  do
42:           Pop the top element from  $F_c$  as  $T$ 
43:           Push  $\max(T, W_{i,c}^{\mathcal{I}})$  onto  $F_c$ 
44:         Move  $\mathcal{I}_i$  from  $\mathcal{I}$  to  $S$ .
45:       break
46:   return  $S$ 

```

representation should be small, while retaining the essential information about what makes a view representative of ones

in an example set – i.e., it should be both concise and informative. Of course, we could compute any feature vector to represent a view in our framework, including a single number based on an embedding, features from a deep network, or all the pixels of a rendered image. However, motivated by applications of semantic segmentation and scene understanding, we choose to represent each view $I \in \mathcal{I}$ as a n -dimensional vector $w_{I,1\dots n}$, where each dimension encodes information about the likelihood of a single semantic category’s contribution to a rendered image for the view. Since our examples come from the training set of NYUDv2, we select n to be 40, where each value $w_{I,c}$ corresponds to an NYU40 category.

More formally, we represent the contribution from a category $c \in \mathcal{C}$ to an image $I \in \mathcal{I}$ as the average likelihood of the pixels $p = (x, y, d, c') \in I$ where $c' = c$. Categories not present in an image are assigned a score of zero. We compute the likelihood of an observation (x, y, d) conditioned on class c as $f_c(x, y, d)$, the value of the pdf of category c defined over the three dimensional view volume. It represents the likelihood of observing category c at pixel locations x and y and depth d in the view volume relative to all other points in the view volume, and is normalized with ℓ_1 normalization. We approximate each function $f_c(x, y, d)$ as a 3-dimensional histogram. For each category c , we compute its contribution to w_I for any view I by looking up the likelihood of $f_c(x, y, d)$ for every pixel (x, y, d) in I and taking the average.

Intuitively this representation encodes the likelihood of observing a particular spatial distribution within each semantic category. This general approach could be applied to other properties of scene observations or other data. For RGB images, where depth is not available, we can use a two-dimensional histogram, which is equivalent to eliminating the uncertainty by integrating over the depth dimension. Alternatively, one could define $f_c(x, y, d) := f_c(x, y)p_c(d)$ where $p_c(d)$ is a specified prior over depth, i.e. for the semantic category sofa and the goal of generating NYUDv2-like images a gaussian centered 1.5 meters away might be reasonable. The details of $f_c(x, y, d)$ do not affect the rest of the algorithm, as long as it is nonnegative and larger values represent more desirable contributions, which are requirements for submodular maximization.

3.2. Candidate generation

Our next step is to generate a discrete set of candidate views. Given a CAD model of an indoor scene and a method to map any particular view I to an n -dimensional vector w describing the likelihood of its scene observations with respect to all n semantic categories (as described in the previous subsection), the goal is to generate a set of candidate views \mathcal{I} that will form the input to the view set selection algorithm in the next step.

Ideally, the output candidate set \mathcal{I} will contain all views likely to be selected for the final output (high recall), but with as few extras as possible (high precision). Finding the optimal set is not tractable: the space of views is infinite (with 6 continuous degrees of freedom representing translations and rotations), and evaluating the representation for any given view requires rendering an image and counting pixels in each semantic category. Thus, we utilize approximations leveraging the specific properties of our view representation. As our feature vector w for each view measures its likelihood in the example distribution along each of several axes, we can approximate overall likelihood for a view or set of views by averaging the feature vectors over all in the set. In particular, we can estimate the likelihood for particular view positions by averaging over all views at that position, marginalizing rotations and pixel locations (x, y) . This allows us to first sample view positions with high estimated likelihood and then later choose rotations, which saves having to render specific images in the first step.

Our algorithm first voxelizes the 3D scene, weighting the voxels $v \in \mathcal{V}$ in proportion to our estimates of the aggregate feature vector (that is, taking a weighted sum over the 40 values $W_{v,:}^{\mathcal{V}}$ in the approximation). The approximation at a voxel is computed by casting a set of rays into the scene and intersecting them with scene geometry; each ray provides a measured category c and depth d relative to the ray origin. The contribution of this sample to the aggregate feature vector is f_c integrated over x, y , thus approaching in the limit the likelihood of the depth distribution of an image with maximal field of view at the ray center. Before taking the aggregate sum, weights may be applied on a per-category basis to further bias the search space. In our implementation, we use weights set to rebalance the category frequencies in the 3D scenes, SUNCG, to the frequencies in the example set, NYUDv2.

Once voxels have been weighted in proportion to their estimated likelihood of selection, we generate candidate views in voxels with probability proportional to the voxel weight. To select a view direction, we do uniform sampling, except for the camera tilt. As a data-driven prior is also available for tilt in our chosen example dataset, NYUDv2, through accelerometer data, we exploit it by estimating the accelerometer distribution \mathcal{A} and selecting tilts in proportion to that distribution. This distribution is not a major requirement of the algorithm proposed here: a gaussian prior can likely perform similarly.

The final step to candidate generation is a view filter. Since there are many obviously poor samples drawn from the view distribution, candidates from each room are first filtered to those with the highest aggregate scores. This is useful because candidates can be generated in parallel, while set selection is sequential. Thus, far more candidates can be generated than evaluated, and a method to improve

the average candidate quality is useful for reducing runtime.

We run the candidate generation algorithm for all scenes in a synthetic training set, and union the outputs to form an overall candidate set for the view selection algorithm in the final step.

3.3. View set selection

The last step of our process is to select a subset of the candidate views that jointly reflect the example distribution.

Defining the goal in this task is tricky, because the true distribution of views is not known – only example images are provided. Assuming the example natural image dataset was generated according to some latent distribution, an ideal set would have high probability of having been generated by this latent distribution, without relying on replicating identical views to those in the example set. More specifically, because we wish to have the capability to select a set that is far larger than the given dataset, we must avoid the pitfall of selecting images that individually have high probability of coming from the distribution, but as a set are unlikely because either: 1) they are all near the mean of the distribution, or 2) they form clusters around one or more of the example images without bridging the support space in between. Our goal is to extend the distribution without replicating it.

The first of the above concerns alone justifies the principle utility of a set-based distribution matching algorithm: it is impossible for any algorithm to determine from a single image whether it was generated according to a particular latent distribution, even if the distribution is known. Therefore, in principle any approach to viewpoint selection that uses only an embedding and generation scheme cannot match an adversarial latent distribution; a set-aware approach is required.

The second of the above concerns precludes any method that selects or encodes image sets based on the distance to the nearest neighbor in the example distribution. In order to get around both of these challenges simultaneously, we formalize the problem of selecting a set S from the set of candidate images \mathcal{I} given their 40 dimensional likelihood estimates $W_{I,c}^{\mathcal{I}}$, with the following submodular objective (written for simplicity to assume that all images have distinct score vectors):

$$\begin{aligned} \max_{S \subseteq \mathcal{I}} \sum_{c \in \mathcal{C}} W_{I,c}^{\mathcal{I}} \mathbb{1}(\exists V \subset S, |V| \geq |S| - h_c \forall v \in V W_{v,c}^{\mathcal{I}} < W_{I,c}^{\mathcal{I}}) \\ \text{s.t. } |S| \leq k \end{aligned} \quad (1)$$

The integer values h_c represent the relative weights of the categories to the optimization, and should increase with the desired output count k ; the higher this proportion, the less a single high likelihood image can cover the output space. In our implementation, we set $h_c := \frac{k}{40}$ for all categories, as the task we compare to is semantic segmentation and mean IOU is the target metric. We next prove that

this optimization is tractable, even at the scale of millions of images, by exploiting the submodular maximization approaches described in Krause et al [18].

Theorem 1. Equation (1) is a nonnegative monotone submodular objective function.

Proof. First, note that the function being maximized, which we refer to as $F(S) : S \subseteq \mathcal{I} \rightarrow \mathcal{R}$, is a summation over the weights $W_{I,c}^{\mathcal{I}}$. As those are nonnegative, the function itself is also nonnegative. Next, consider any two image sets A and B , where $A \subset B \subseteq \mathcal{I}$. Take any image $I \in \mathcal{I} \setminus B$. We must show $\Delta(I|A) \geq \Delta(I|B)$. What can I contribute to $\Delta(I|A)$? For any $c \in \mathcal{C}$, if the indicator function for the pair I, c is positive, there are at most $h_c - 1$ images $v \in V \subset A$ such that $W_{v,c}^{\mathcal{I}} > W_{I,c}^{\mathcal{I}}$. Order the weights $\{W_{v,c}^{\mathcal{I}} : v \in A\}$ by decreasing value and name that ordering $W_{1,c}^{\mathcal{I},A}, W_{2,c}^{\mathcal{I},A}, \dots, W_{|A|,c}^{\mathcal{I},A}$. Similarly consider the same ordering for the corresponding category’s weights for B : $W_{1,c}^{\mathcal{I},B}, W_{2,c}^{\mathcal{I},B}, \dots, W_{|B|,c}^{\mathcal{I},B}$. Then $W_{h_c,c}^{\mathcal{I},A} \leq W_{h_c,c}^{\mathcal{I},B}$, because $A \subset B$. So $W_{I,c}^{\mathcal{I}} - W_{h_c,c}^{\mathcal{I},A} \geq W_{I,c}^{\mathcal{I}} - W_{h_c,c}^{\mathcal{I},B}$. But these are the contributions of category c to $\Delta(I|A)$ and $\Delta(I|B)$, respectively. Since this was done WLOG c , we have $\Delta(I|A) \geq \Delta(I|B)$.

The final step is to show monotonicity. Consider again the proof that $\Delta(I|A) \geq \Delta(I|B)$. Substitute $B := A \cup \{e\}$ for any $e \in \mathcal{I}$. Then $\forall c \in \mathcal{C}, i \in [|A|] W_{i,c}^{\mathcal{I},A} \leq W_{i,c}^{\mathcal{I},A \cup \{e\}}$, implying $F(A) \leq F(A \cup \{e\})$. \square

Since we have shown our objective is nonnegative, monotone, and submodular, there is a greedy $\frac{1}{2}$ -approximation algorithm due to Nemhauser et al [25]. The algorithm is adapted to this problem by computing the discrete derivative $\Delta(I|S)$ for a candidate image as the sum over categories c of the marginal gain of $W_{I,c}^{\mathcal{I}}$ over the least quantity in V still contributing to the optimization. We use a set of minheaps to maintain the sets quickly in practice, where each minheap always contains exactly h_c values. The integer weights h_c per category are selected by dividing the image budget k with image counts evenly amongst all categories, though unbalanced weights, such as those derived by rounding from the relative frequencies of the categories in NYUDv2 training, are also possible.

4. Experiments

In this section we describe a series of experiments to evaluate how well different viewpoint selection algorithms match the statistics of real-world data, and how much difference view selection strategies can make on training performance for semantic segmentation on RGBD data.

Datasets. For all our experiments, we use SUNCG [36] as the source of synthetic 3D scenes to render. We select a set of 171,496 rooms from the SUNCG scenes, by filtering

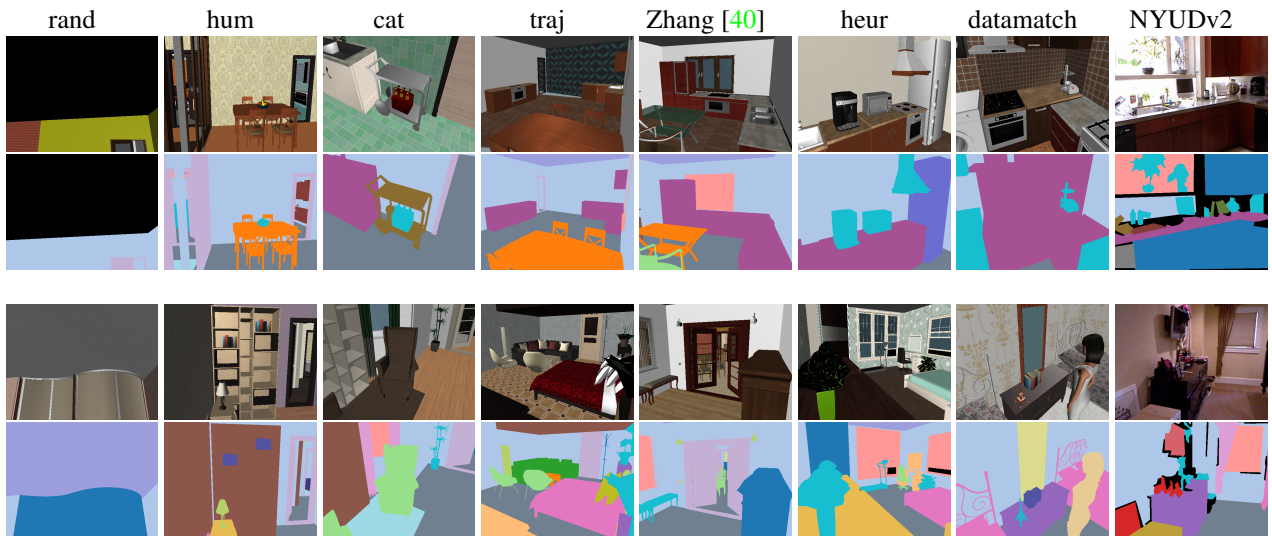


Figure 2: Each column shows viewpoints selected from SUNCG 3D scenes using a different selection algorithm (last column shows an example NYUDv2 image from the same room category). The first and third rows show color, and the second and fourth show NYU40 category semantic segmentations.

for rooms that have a floor area of 8 to 100 squared meters, with at least four walls and five objects present. We choose NYUDv2 [34] as our target data for our data matching algorithm and for evaluating semantic segmentation performance. This allows us to test the benefit of both synthetic depth and synthetic color information. We use the NYU40 category set as defined by Gupta et al. [11] for semantic segmentation and for computing all category-level statistics. The standard NYUDv2 split of 795 training images and 654 testing images is used for our algorithm’s data matching and the semantic segmentation evaluation.

4.1. Comparison algorithms

We compare images generated using variants of our algorithm against a variety of alternative methods for viewpoint selection including random cameras, single object closeup views, human-like navigation trajectories, and manual heuristics from prior work. For each algorithm, we generate a set of twenty thousand camera views (see Figure 2 for examples). We describe the different methods below.

Random cameras (rand). For each room, a single random view is computed. A camera position is randomly sampled within the room volume, rejecting samples inside the axis aligned bounding box of scene objects. The view direction is randomly sampled uniformly. This is similar to the approach described by Handa et al. [13] except that they select viewpoints with at least 3 objects.

Random human eye level (hum). In addition to the purely random view selection, we also consider a baseline algorithm where we restrict the camera positions to be at an average eye height of 1.55 meters, corresponding to holding a camera just below eye height.

Object closeups (cat). The category-based camera approach focuses on spending the image budget equally across the NYU40 categories. Four NYU40 categories not present in SUNCG, and the wall, ceiling, and floor categories are not allocated a budget. The remaining 33 categories are each allocated approximately 1/33rd of the 20,000 images. For each image of a given category, an object is chosen with replacement from the set of all objects of that category. For each chosen object, an image focusing on that object from eye height is selected. A set of candidate images in a ring in the eye height plane around the object are generated, and the image with the greatest fraction of pixels belonging to the object is selected.

Trajectories (traj). Trajectory cameras are selected from an unobstructed path in the room. A two dimensional positional grid is imposed on each room, and grid points are flagged as obstructed if they are inside or close to the surface of an object. For each such room, a point is selected. The furthest grid point from that point under an unobstructed distance metric is selected as one endpoint of the trajectory. The second endpoint is selected to be the furthest point from that point, again constraining the candidate set to paths not within the obstacles of the room. Images are densely sampled from the trajectory. For each image, the view direction is towards the centroid of object centroids, weighted by the number of faces of each object. Views are subsampled randomly to reach 20K images. At a high level, traj is an approach similar to the two-body trajectory algorithm used by McCormac et al. [23].

Zhang et al. [40]. We compare against the views generated by Zhang et al. [40]. Their method scores views based on object count and pixel coverage, and filters ray-traced images using a color and depth histogram similarity

method	depth	x	y	mean
rand	0.56	0.166	0.474	0.400
hum	0.404	0.169	0.244	0.272
cat	0.434	0.177	0.192	0.268
traj	0.346	0.107	0.145	0.199
Zhang et al. [40]	0.377	0.100	0.212	0.230
heur	0.349	0.104	0.183	0.212
datamatch	0.346	0.094	0.134	0.191

Table 1: EMD with thresholded ground distances for each axis of object occurrence, averaged over the NYU40 categories. Distances above are normalized to [0-1], where lower values are closer to NYUDv2.

to NYUDv2 images. The authors provided their generated view set which we subsampled randomly to match the total set size used for all other methods.

Hand-tuned heuristic (heur). We implement a hand-tuned heuristic that maximizes the number and pixel occupancy of object categories in the output views. First, views are densely sampled at eye height in a grid of unobstructed positions in the room, sampling several view directions per position, each of which is given a fixed downward tilt. Each sampled view is assigned a score $s = \sum_{i \in v_i} \mathbb{1}(p_i > m)(\log(p_i) - \log(m))$ where p_i is the number of pixels of object i in view v_i , and m is a minimum number of pixels per object. The top scoring images across all rooms are chosen as the training set. Intuitively, this method attempts to maximize the salient information in the frame, and is an attempt at improving Zhang et al.’s method.

Data matching (datamatch). We apply our data-driven view set selection algorithm to match RGBD frames from the NYUDv2 [34] training images, in terms of the distribution within-frame of the NYU40 categories. We rebalance the category frequencies in SUNCG to the frequencies of the NYU40 categories.

4.2. Evaluation with earth mover’s distance

We perform three types of experiments to investigate characteristics of our view set selection algorithm. The first is focused on quantitative evaluation of how well it generates distributions of images matching an example set. To address this question, we use an earth mover’s distance (EMD) to measure the difference between two RGBD image sets. Specifically, we use the thresholded distance metric as described in Pele et al. [28] as it is particularly suitable for image data. For each category in an image set, we generate three separate 1D histograms: the x , y , and depth occurrences of that category throughout the set. For each such histogram, the EMD to the corresponding histogram from the training & validation set of NYUDv2 is computed, and the results are averaged over all 40 NYUDv2 categories.

Table 1 shows the results. As we can see, the `datamatch` algorithm achieves the lowest EMD to NYUDv2 frames overall. The `traj` algorithm performs second best, most likely due to the fact that NYUDv2 was collected by following human-feasible trajectories through real rooms. In any case, this experiment confirms that our algorithm can produce distributions of images more similar to an example set than alternatives.

4.3. Semantic segmentation results

In our second experiment, we investigate the impact of our view set selection approach on the performance of pre-trained networks for semantic segmentation.

Procedure. We follow the training procedure of Zhang et al. [40] to train a fully convolutional network [20]. Initial weights are taken from FCN-VGG16, the weights of VGG 16-layer net [35] pretrained on ImageNet [6] converted to an FCN-32s network. We then pretrain this network using the images rendered from a given algorithm’s 20k selected views in SUNCG for five epochs, and finally we finetune on the NYUDv2 training set for 100K iterations. We run experiments using depth or color images. For depth, we use the HHA depth encoding of Gupta et al. [11]. For color images, we render an RGB image for each view using the OpenGL rendering toolkit provided by Song et al. [36]. We run all our experiments with the Caffe library [16] and the public FCN implementation [20], on Nvidia Titan X GPUs.

Results. Table 2 reports weighted and unweighted accuracy and IoU metrics for the NYUDv2 test set, as defined by Long et al. [20]. We see that good view selection makes a big difference for depth with `datamatch` views improving the baseline FCN HHA mean IoU by 5.2 points (20.6% relative improvement). A smaller impact is seen for synthetic color, where `datamatch` leads to a 0.6 point improvement (1.9% relative improvement). We hypothesize this is due to the large color disparity between the synthetic color images and the real world color. For depth, using the HHA encoding, `heur` and `datamatch` both outperform the camera selection algorithm used in [40]¹. The `rand` algorithm performs especially bad since the HHA encoding relies on a built-in assumption of a dominant front orientation of the camera. Moreover, `hum`, `traj`, and `rand` significantly reduce baseline performance in color indicating that pretraining with bad synthetic views can hurt instead of help.

Influence of submodular optimization for set selection.

To isolate the impact of the submodular optimization part of our algorithm we also compare against a version of `datamatch` that only selects the best k scoring images instead. The views generated by this ablated version of the algorithm result in reducing performance on HHA to 61.8

¹Reported HHA performance in Zhang et al. is significantly lower at 23.4% mean IoU.

method	HHA				RGB			
	Pixel Acc.	Mean Acc.	Mean IoU	F.W. IoU	Pixel Acc.	Mean Acc.	Mean IoU	F.W. IoU
FCN-32s [20]	58.2	36.1	25.2	41.9	61.8	44.7	31.6	46.0
rand	6.1	3.1	0.6	1.3	59.3	41.0	28.4	44.2
hum	60.0	39.8	28.1	44.3	61.4	42.2	30.0	45.7
cat	61.3	41.6	29.8	45.5	62.6	44.5	31.6	46.9
traj	60.7	40.9	28.6	44.7	61.2	42.5	30.2	45.1
Zhang et al. [40] ¹	61.3	41.9	29.8	45.5	62.7	43.7	31.8	46.3
heur	61.7	42.2	30.1	45.9	62.6	44.5	31.7	46.9
datamatch	61.2	42.7	30.4	45.1	62.4	45.1	32.1	46.8

Table 2: Semantic segmentation results on the NYUDv2 test set after pretraining on different synthetic image datasets.

pixel accuracy, 42.2 mean accuracy, 30.2 mean IoU, and 45.8 frequency weighted IoU.

4.4. Human view selection judgments

In a third experiment, we investigate the correlation between semantic segmentation performance and human view selection by running a study where we present people with a target set of NYUDv2 images in a given room type (e.g., bedroom) and ask them to select from a random set of five synthetic views sampled uniformly from the images of all seven algorithms considered in the previous experiments. Our goal in this experiment is to characterize whether human judgement of good views is as good for semantic segmentation as our proposed algorithm.

We recruited 261 Amazon Mechanical Turk workers to select any number of synthetic images that match the types of views of the target NYUDv2 set. Each worker selected images from a total of 30 random sets of five images. Overall, the workers saw 27,564 images and selected 10,046 as matching the NYUDv2 target set. To evaluate these results, we measure the fraction of time that an image from a given algorithm was selected as matching. The mean values are rand: 0.04, hum: 0.29, cat: 0.34, traj: 0.51, Zhang et al. [40]: 0.50, heur: 0.53, datamatch: 0.41. Interestingly, views produced by algorithms that provide greater EMD results and lesser semantic segmentation results than datamatch are selected more often by people. This phenomenon suggests that people might not be a good judge of how images fit into distributions, corroborating the need for automatic algorithms for this task.

We compare the human judgments of match to NYUDv2 against EMD as predictors of semantic segmentation performance. Figure 3 plots each algorithm as a point and shows the linear fit (along with 95% confidence contours), and the R^2 measure of goodness of fit. As expected, EMD is inversely correlated with segmentation performance, and human judgments are mostly positively correlated with performance. The higher R^2 numbers for EMD-to-segmentation fit indicate that EMD is a better predictor of performance than human judgment. The higher correlation of EMD

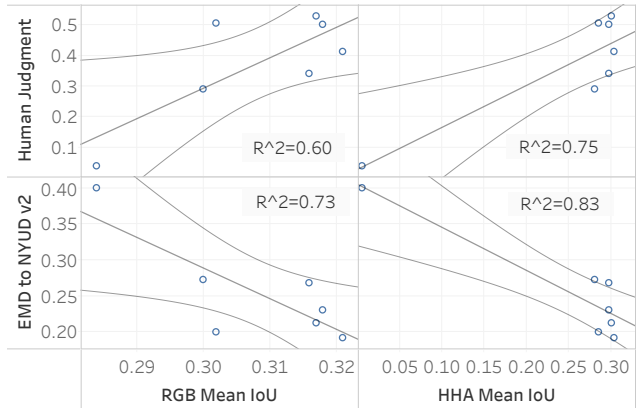


Figure 3: Semantic segmentation performance (left: RGB, right: HHA) plotted against both EMD distance of images to NYUDv2 and subjective human judgment of how close images are to the set of NYUDv2 images.

to semantic segmentation performance suggests that our choice of image representation is effective.

5. Conclusion

In summary, this paper has investigated the idea of using data-driven algorithms to select a set of views when creating synthetic datasets of rendered images for pre-training deep networks. It provides a formulation for representing views based on pdfs of object observations of specific semantic categories, generating candidate views with an efficient and practical search algorithm, and selecting a set of views matching a target distribution with submodular maximization. Results of experiments show that careful view select with this approach can improve the performance of semantic segmentation networks.

This is just the tip of the iceberg on how to select views for training deep networks for computer vision tasks. One way to improve our algorithm is to remove the independence assumption between categorical likelihoods (i.e., explicitly model object co-occurrence probabilities). We chose a formulation motivated by semantic segmentation, but oth-

ers may be better suited for different tasks. We establish that view set selection is an important problem in this domain and conjecture that much future work is warranted.

References

- [1] W. Bares. A photographic composition assistant for intelligent virtual 3D camera systems. In *International Symposium on Smart Graphics*, pages 172–183. Springer, 2006. 1, 2
- [2] W. Bares, S. McDermott, C. Boudreaux, and S. Thainimit. Virtual 3D camera composition from frame constraints. In *Proceedings of the eighth ACM international conference on Multimedia*, pages 177–186. ACM, 2000. 1, 2
- [3] P. S. Blaer and P. K. Allen. Data acquisition and view planning for 3-D modeling tasks. In *Intelligent Robots and Systems, 2007. IROS 2007. IEEE/RSJ International Conference on*, pages 417–422. IEEE, 2007. 2
- [4] S. Choi, Q.-Y. Zhou, and V. Koltun. Robust reconstruction of indoor scenes. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5556–5565. IEEE, 2015. 2
- [5] M. Christie, P. Olivier, and J.-M. Normand. Camera control in computer graphics. In *Computer Graphics Forum*, volume 27, pages 2197–2218. Wiley Online Library, 2008. 1, 2
- [6] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A large-scale hierarchical image database. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009. 7
- [7] A. Doumanoglou, R. Kouskouridas, S. Malassiotis, and T.-K. Kim. Recovering 6D object pose and predicting next-best-view in the crowd. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2
- [8] K. A. Ehinger and A. Oliva. Canonical views of scenes depend on the shape of the space. 2011. 2
- [9] S. Freitag, B. Weyers, A. Bönsch, and T. W. Kuhlen. Comparison and evaluation of viewpoint quality estimation algorithms for immersive virtual environments. In *ICAT-EGVE*, pages 53–60, 2015. 2
- [10] B. Gooch, E. Reinhard, C. Moulding, and P. Shirley. Artistic composition for image creation. In *Rendering Techniques 2001*, pages 83–88. Springer, 2001. 1, 2
- [11] S. Gupta, R. Girshick, P. Arbeláez, and J. Malik. Learning rich features from RGB-D images for object detection and segmentation. In *European Conference on Computer Vision*, pages 345–360. Springer, 2014. 6, 7
- [12] A. Handa, R. A. Newcombe, A. Angeli, and A. J. Davison. Real-time camera tracking: When is high frame-rate best? In *European Conference on Computer Vision*, pages 222–235. Springer, 2012. 2
- [13] A. Handa, V. Patraucean, V. Badrinarayanan, S. Stent, and R. Cipolla. SceneNet: Understanding real world indoor scenes with synthetic data. *arXiv preprint arXiv:1511.07041*, 2015. 1, 2, 6
- [14] A. Handa, T. Whelan, J. McDonald, and A. J. Davison. A benchmark for RGB-D visual odometry, 3D reconstruction and SLAM. In *2014 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1524–1531. IEEE, 2014. 2
- [15] J. Jankowski and M. Hachet. Advances in interaction with 3D environments. In *Computer Graphics Forum*, volume 34, pages 152–190. Wiley Online Library, 2015. 2
- [16] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the 22nd ACM international conference on Multimedia*, pages 675–678. ACM, 2014. 7
- [17] M. Krainin, B. Curless, and D. Fox. Autonomous generation of complete 3D object models using next best view manipulation planning. In *Robotics and Automation (ICRA), 2011 IEEE International Conference on*, pages 5031–5037. IEEE, 2011. 2
- [18] A. Krause and D. Golovin. Submodular function maximization. *Tractability: Practical Approaches to Hard Problems*, 3(19):8, 2012. 5
- [19] T. Liu, J. McCann, W. Li, and T. Funkhouser. Composition-aware scene optimization for product images. In *Computer Graphics Forum*, volume 34, pages 13–24. Wiley Online Library, 2015. 1, 2
- [20] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015. 2, 7, 8
- [21] A. Mavrincac and X. Chen. Modeling coverage in camera networks: A survey. *International journal of computer vision*, 101(1):205–226, 2013. 1, 2
- [22] N. Mayer, E. Ilg, P. Häusser, P. Fischer, D. Cremers, A. Dosovitskiy, and T. Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. *arXiv preprint arXiv:1512.02134*, 2015. 2
- [23] J. McCormac, A. Handa, S. Leutenegger, and A. J. Davison. SceneNet RGB-D: 5M photorealistic images of synthetic indoor trajectories with ground truth. *arXiv preprint arXiv:1612.05079*, 2016. 1, 6
- [24] E. Mezuman and Y. Weiss. Learning about canonical views from internet image collections. In *Advances in Neural Information Processing Systems*, pages 719–727, 2012. 2
- [25] G. L. Nemhauser, L. A. Wolsey, and M. L. Fisher. An analysis of approximations for maximizing submodular set functions. *Mathematical Programming*, 14(1):265–294, 1978. 5
- [26] P. Olivier, N. Halper, J. Pickering, and P. Luna. Visual composition as optimisation. In *AISB symposium on AI and creativity in entertainment and visual art*, pages 22–30, 1999. 1, 2
- [27] J. Papon and M. Schoeler. Semantic pose using deep networks trained on synthetic RGB-D. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 774–782, 2015. 2
- [28] O. Pele and M. Werman. Fast and robust earth mover’s distances. In *Computer vision, 2009 IEEE 12th international conference on*, pages 460–467. IEEE, 2009. 7
- [29] S. R. Richter, V. Vineet, S. Roth, and V. Koltun. Playing for data: Ground truth from computer games. *arXiv preprint arXiv:1608.02192*, 2016. 2

- [30] G. Ros, L. Sellart, J. Materzynska, D. Vazquez, and A. M. Lopez. The SYNTHIA dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3234–3243, 2016. 2
- [31] J. L. Schönberger, E. Zheng, J.-M. Frahm, and M. Pollefeys. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision*, pages 501–518. Springer, 2016. 1
- [32] A. Secord, J. Lu, A. Finkelstein, M. Singh, and A. Nealen. Perceptual models of viewpoint preference. *ACM Transactions on Graphics (TOG)*, 30(5):109, 2011. 2
- [33] A. Shafaei, J. J. Little, and M. Schmidt. Play and learn: Using video games to train computer vision models. *arXiv preprint arXiv:1608.01745*, 2016. 2
- [34] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus. Indoor segmentation and support inference from RGBD images. In *European Conference on Computer Vision*, pages 746–760. Springer, 2012. 1, 3, 6, 7
- [35] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2015. 7
- [36] S. Song, F. Yu, A. Zeng, A. X. Chang, M. Savva, and T. Funkhouser. Semantic scene completion from a single depth image. In *Proceedings of 30th IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 1, 2, 5, 7
- [37] S. Tulsiani and J. Malik. Viewpoints and keypoints. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1510–1519. IEEE, 2015. 2
- [38] P.-P. Vázquez, M. Feixas, M. Sbert, and W. Heidrich. Viewpoint selection using viewpoint entropy. In *VMV*, volume 1, pages 273–280, 2001. 2
- [39] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao. 3D ShapeNets: A deep representation for volumetric shapes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1912–1920, 2015. 2
- [40] Y. Zhang, S. Song, E. Yumer, M. Savva, J.-Y. Lee, H. Jin, and T. Funkhouser. Physically-based rendering for indoor scene understanding using convolutional neural networks. In *Proceedings of 30th IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 1, 6, 7, 8
- [41] E. Zheng, E. Dunn, V. Jovic, and J.-M. Frahm. Patchmatch based joint view selection and depthmap estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1510–1517, 2014. 1

Appendix

A. Example Viewpoint Sets

The viewpoint set selection problem we pose in our paper takes as input a 3D scene and provides as output a set of viewpoints in that scene that best match a prior data distribution. Figure 4 shows two example 3D scenes and output viewpoint sets selected using our `datamatch` algorithm.

B. Single Image Matching

The `datamatch` algorithm was designed to match the object distributions seen in sets of images, and to generate sets of viewpoints that exemplify similar distributions. However, it can be directly applied to single image inputs and restricted to generate one viewpoint that best matches the input, given a particular target 3D scene. Figure 5 shows some example viewpoints generated to match a particular view of a kitchen counter.

C. RGB Image Matching

A key strength of the `datamatch` algorithm we present is that it can be directly applied to match viewpoint distributions in arbitrary target datasets. Moreover, though the main paper evaluation targets the RGBD data in NYUDv2, using depth data is not a requirement of the algorithm. We demonstrate how our algorithm generalizes to a viewpoint distribution that differs drastically from the views in NYUDv2. To this end, we collected a small dataset of semantically annotated images from an RGB camera mounted on a Roomba robot and generated a viewpoint set for a living room SUNCG scene. Figure 6 shows the results. To replace the voxel weighting step, uniform voxel weights were used for all voxels in the 3D room volume; the sample count was not increased.

D. User Study Details

Figure 7 shows screenshots of the interface that we used to carry out the user study asking people to select generated views that match the NYUDv2 viewpoints. Each participant saw 30 randomly sampled sets of viewpoints, and matched them against a randomly sampled set of NYUDv2 images from the same room category.

E. Constants

The 40 weights w_c per category were chosen to approximately rebalance SUNCG frequencies to NYUDv2 frequencies, and are given in Table 3. The rebalancing was object occurrence in NYUDv2 divided by object occurrence in SUNCG. Weights marked N/A correspond to categories not present in SUNCG, and the wall, ceiling, and floor categories were set to a weight of 1. The candidate generation



Figure 4: Two example 3D scenes and their viewpoint sets selected by the `datamatch` algorithm trained on NYUDv2. Viewpoints are shown in the overview images as view frustum cones drawn in red and orange outlines. Corresponding images for a subset of the viewpoints in each are shown around the scene. For the first scene, the top 3 views for each room were selected in the set, whereas for the second scene, the top 20 were selected.

step used 1500 candidates per room, 200 raycasts per voxel, and filtering down to 20 candidates per room. These settings resulted in roughly comparable runtime to the `heur` algorithm.

When generating images from only a single image or

handful of images as is the case in this supplemental document, the constants and `datamatch` algorithm are adjusted to compensate for the increased noise in the semantic category pdfs. First, for each exemplifying image in the input, a univariate gaussian with standard deviation equal to

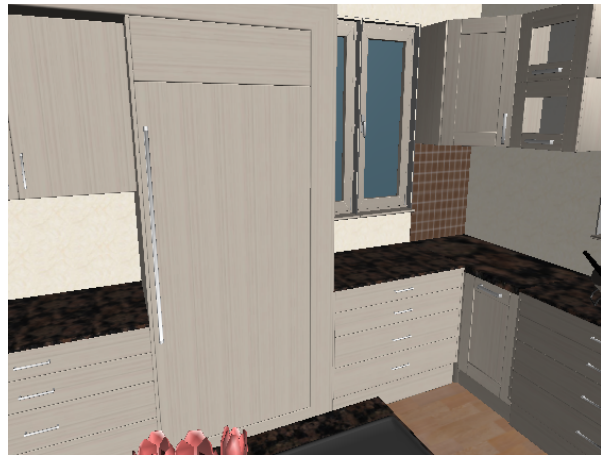


Figure 5: **Top left:** input image exemplifying desired framing. **Top right:** output viewpoint in same room, generated using the datamatch algorithm. **Middle and bottom:** output viewpoints in other rooms. Note that the objects seen and their framing in the view is similar to the input image, despite differences in the room layout.

10% of the histogram resolution is added to each semantic category's pdf along the depth axis, rather than a single point per pixel. Next, to compensate for the narrower dis-

tribution being matched, the candidate sample count is increased to 3500 to ensure good samples are still found. The candidate image vector likelihood values are not divided by

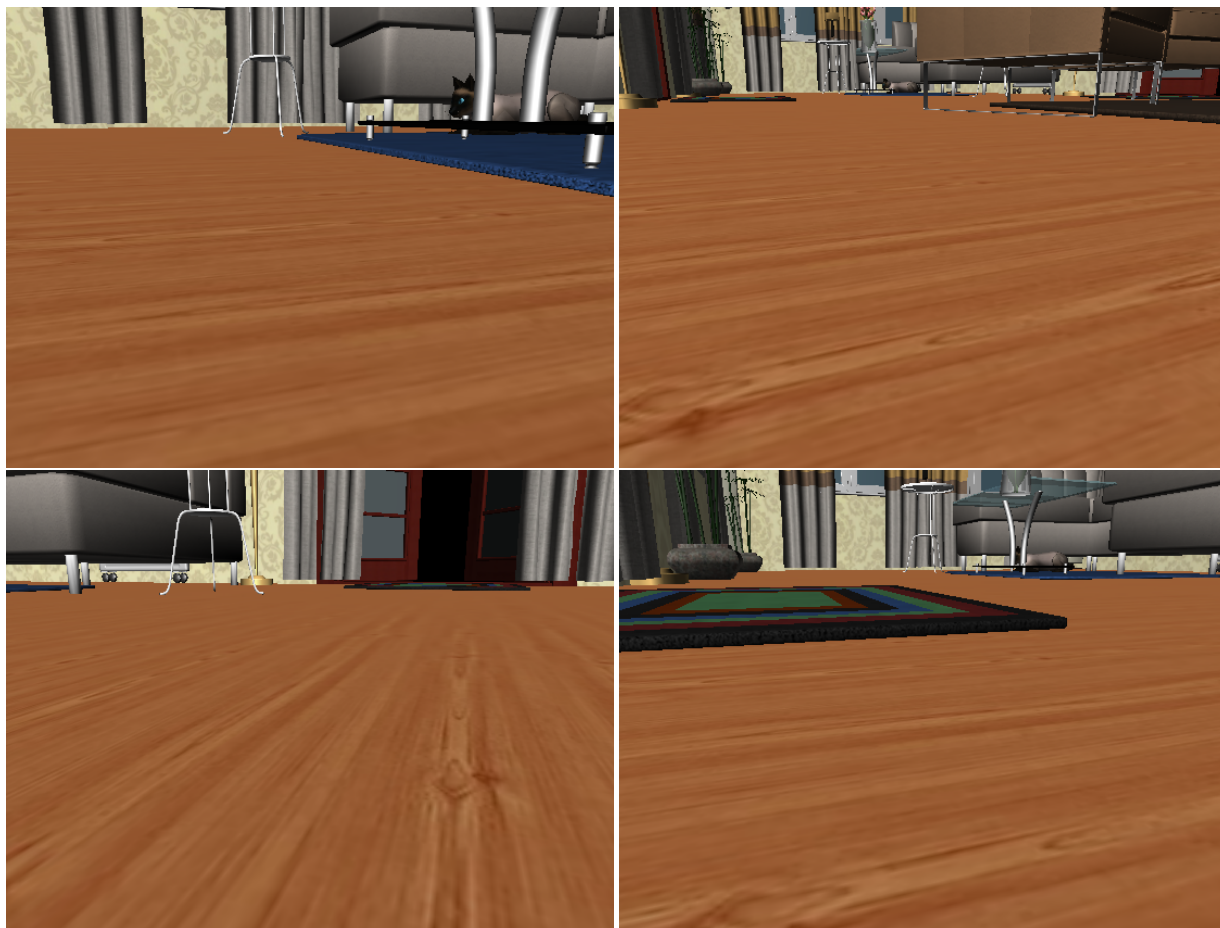


Figure 6: **Top:** RGB images collected using an iPad Air2 camera attached to a Roomba robot. We semantically annotated these images and used the `datamatch` algorithm to generate matching viewpoints in a 3D scene. **Bottom:** Generated viewpoints which exhibit the same view distribution as the input data. Note that depth data was not used, and the algorithm was not changed in any way to account for differences between Roomba viewpoints and human viewpoints.

the observation frequency; this helps stabilize filtering when the input image count is small. Finally, as the distribution is no longer intended to match NYUDv2, uniform category weights w_c are used. Note also that as the images generated for the figures here are created from a single room’s output rather than that of 20,000 rooms, the `SELECT` portion of the `datamatch` algorithm is superfluous.

F. Runtime

The runtime of the `datamatch` algorithm is comparable to that of other methods. When run on 20,000 selected scenes in SUNCG, the total candidate generation time was 898 hours, 37 minutes; this was run in parallel on a cluster with approximately 500 simultaneously allocated cores.

In this task, we ask you to choose images that are most similar to another set of images. Images are similar if the **VIEW** is similar (the position of the camera, and what objects are shown in the image). The different colors in the image represent different kinds of objects. Focus on the positioning of the view, and the types of objects shown in the image.

You will see a total of 30 sets of images. For each set, there will be 5 images in the center of the screen. Pick as many as you think are good matches to the smaller target images at the top of the screen. You can click on an image, or type the number of the image to select it (it will be outlined in blue when selected).

Each set should take about less than a minute to complete. The task should take about 5 minutes in total.

Start

1/30



Please select the images that are most similar in view to the above images



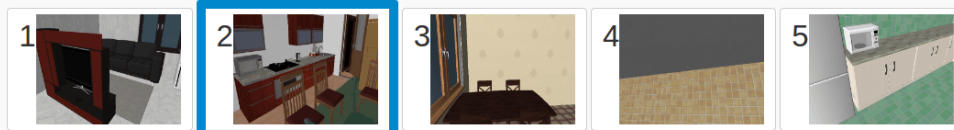
Click on 'Next' to continue

Next

2/30



Please select the images that are most similar in view to the above images



Click on 'Next' to continue

Next

8/30



Please select the images that are most similar in view to the above images



Click on 'Next' to continue

Next

Figure 7: Screenshots of the instructions and interface for the view set match selection user study we carried out. Participants were instructed to select generated viewpoints images in the bottom row that match the viewpoints of the top example images from NYUDv2.

wall	1.0000
floor	1.0000
cabinet	1.0770
bed	0.6267
chair	0.8133
sofa	0.7622
table	0.7552
door	0.1293
window	0.2678
bookshelf	6.3694
picture	2.3617
counter	0.2667
blinds	1.6238
desk	0.6391
shelves	0.4318
curtain	0.4273
dresser	0.5381
pillow	11.8598
mirror	0.7994
floor_mat	0.2805
clothes	38.6042
ceiling	1.0000
books	4.0946
refridgerator	0.5520
television	0.1626
paper	N/A
towel	N/A
shower_curtain	0.1942
box	N/A
whiteboard	1.6733
person	0.3813
night_stand	0.2648
toilet	0.1204
sink	0.4990
lamp	0.1518
bathtub	0.2869
bag	N/A
otherstructure	2.8287
otherfurniture	497.8754
otherprop	1.1364

Table 3: Category weights w_c for the NYU40 categories. These weights were chosen to rebalance the category occurrence frequencies of SUNCG to NYUDv2 images.

In addition, the set selection stage required 22 minutes, 11 seconds. By comparison, the total runtime for the `heur` method on the same dataset was 880 hours, 54 minutes. For both algorithms, the primary computational bottleneck is rendering candidate images in order to score them. Note that the times reported above do not include rendering the selected images to disk for training, as OpenGL renders are used and therefore this is an IO-bound task.