

Generating Fast and Slow: Scene Decomposition via Reconstruction

Mihir Prabhudesai¹ Anirudh Goyal² Deepak Pathak¹ Katerina Fragkiadaki¹
¹Carnegie Mellon University ²Mila, University of Montreal
 {mprabhud, dpathak, katef}@cs.cmu.edu anirudhgoyal9119@gmail.com

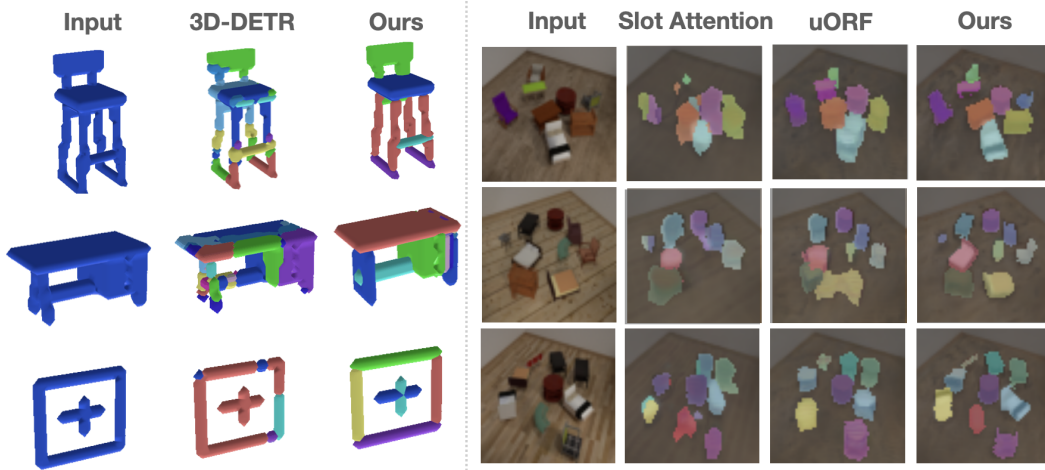


Figure 1: Point-cloud and image decompositions with **Generating Fast and Slow Networks (GFS-Nets)**. GFS-Nets parse completely novel scenes into familiar entities via slow inference, i.e., gradient descent on the reconstruction error of the scene example under consideration. *Left*: GFS-Nets outperform a state-of-the-art 3D-DETR detector by 50% in segmentation accuracy in out-of-distribution 3D point clouds, when trained on the same training data. *Right*: GFS-Nets outperform state-of-the-art unsupervised generative models of Slot Attention [42] and uORF [67] in RGB image decomposition.

Abstract

We consider the problem of segmenting scenes into constituent entities, i.e. underlying objects and their parts. Current supervised visual detectors though impressive within their training distribution, often fail to segment out-of-distribution scenes into their constituent entities. Recent slot-centric generative models break such dependence on supervision, by attempting to segment scenes into entities unsupervised, by reconstructing pixels. However, they have been restricted thus far to toy scenes as they suffer from a reconstruction-segmentation trade-off: as the entity bottleneck gets wider, reconstruction improves but then the segmentation collapses. We propose GFS-Nets (Generating Fast and Slow Networks) that alleviate this issue with two ingredients: i) curriculum training in the form of primitives, often missing from current generative models and, ii) test-time adaptation per scene through gradient descent on the reconstruction objective, what we call slow inference, missing from current feed-forward detectors. We show the proposed curriculum suffices to break the reconstruction-segmentation trade-off, and slow inference greatly improves segmentation in out-of-distribution scenes. We evaluate GFS-Nets in 3D and 2D scene segmentation benchmarks of PartNet, CLEVR, Room Diverse++, and show large ($\sim 50\%$) performance improvements against SOTA supervised feed-forward detectors and unsupervised object discovery methods.

1 Introduction

Scenes are composed of objects and objects are composed of parts, and this compositional organization of the perceptual representations is considered a critical component for the level of combinatorial generalization that humans are capable of, that extends far beyond their direct experiences [25]. Even when encountered with a truly novel image, humans try to parse it in terms of known entity components [26]. Such an entity-centric understanding or reasoning of scenes allows humans to generalize known concepts to *out-of-distribution* examples [25, 39, 44].

An expensive way to build this entity-centric understanding is to supervise the part parsing using human labels for training object and parts detection or segmentation in images and 3D point clouds [56, 69, 5, 71]. However, these part detectors are category specific and often fail outside the training distribution [29, 19, 2]. Consider the abstract shape shown in Figure 1 (last row on the left). We can intuitively figure out the meaningful entities that this shape could be broken into. Yet, a state-of-the-art transformer based 3D-DETR detector [5], when trained supervised for segmenting chairs, fails to decompose these abstract shapes into entities, even though these shapes contain familiar (chair) parts. An underlying issue is that these models do not receive any top-down feedback [3] from entities, i.e., they do not have any entity-centric reasoning rather, they are purely feed-forward in nature, i.e., the information is routed from input to output layers in a sequential manner. How do we enable top-down feedback between the scene and its constituent parts?

There has been interest recently in building models that segment scenes into entities in an unsupervised way by optimizing a reconstruction objective [42, 17, 24, 63, 57, 35, 16, 4, 23, 68, 55, 13, 22, 73, 34]. These methods differ in details but share the notion of having a fixed set of entities, also known as *slots* or *object files*. Each slot extracts information about a single entity during encoding, and it is “rendered” back to the input during decoding, thus, exhibiting top down feedback. Yet, it has been hard to scale the results of these models beyond toy datasets.

We argue this is because these models answer two questions at once: understanding *what* the entity is and *where* it is present in the input. This creates a chicken-and-egg problem because unless we have a representation for an entity, we can’t find where it is in the scene, and unless we localize the entity, we can’t learn a representation of what it is. Hence, a trivial solution often encountered is to assign one slot to the input and other slots to nothing, as noted in Engelcke *et al.* [15] and also confirmed by our experiments in Section 4.1.1 and 4.2.

We propose Generating Fast and Slow Networks (GFS-Nets), a slot-centric autoencoder that segments scenes while optimizing a reconstruction objective. Our model builds on top of Slot Attention[42], with two key insights: **First**, it uses curriculum to learn the *what* question (what entities are) before going to the *where* question (where they are) it does this by training the model to autoencode primitive entities using a single slot bottleneck. This collection of primitive entities can either be collected by design [60] or can be taken from a collection of part segmentations [49]. **Second**, inference in GFS-Nets is performed in two ways: i) *Fast inference* where a visual scene is simply encoded and decoded and no model weights are updated. ii) *Slow inference* where the weights of the model are adapted by gradient descent while optimizing the reconstruction objective. In this way, GFS-Nets adapts to a truly novel scene by *projecting* it onto a learned distribution of entities, while also slowly adapting the entity distribution. Slow inference in GFS-Nets successfully parses completely unfamiliar scenes into familiar entities, implementing a form of *search for explanations* of the scene at inference, as we show in Figure 3.

Our slow inference is related to test-time adaptation or optimization [72, 54, 59, 58, 74, 65]. However, test-time adaptation has not been used so far in entity-centric models to improve scene decomposition. In fact, it has been recently noted that reconstruction loss can make entity segmentation worse, to quote Engelcke *et al.* [15], “If the bottleneck is too narrow, segmentation and reconstruction degrade. If it is too wide, reconstruction is reasonable but the model collapses to a single component and no useful segmentation is learned.” GFS-Nets alleviate from this issue by coupling curriculum training with slow inference for entity-centric learning.

Project page: https://mihirp1998.github.io/project_pages/gfsnets/

We test GFS-Nets on the following datasets: PartNet [49], ShapeNet[7], CLEVR [30] and a difficult version of Room Diverse [67]. We evaluate GFS-Nets’s ability to parse out-of-distribution scenes and compare it against state of the art entity-centric generative models [42, 67], program synthesis models [60], 3D unsupervised part discovery models [21, 64] and state of the art supervised visual detectors [5, 43] trained with labeled data to segment entities. We show improvements across all baselines in our ability to segment novel scenes. Additionally, we ablate different design choices of GFS-Nets. We will make our code and datasets publicly available to the community.

2 Related Work

Entity-centric generative models for scene and structure decomposition Entity-centric models attempt to segment a scene into objects and parts while autoencoding images. MONet[4], GENESIS[16] and IODINE[23] perform multiple steps between their encoder and the decoder to output a set of independent latent vectors that describe objects in the image. This iterative encode-decode inference helps them in receiving top-down feedback. However, this also results in making the system computationally inefficient and inflexible. Slot Attention [42] on the other hand receives top-down feedback within a single encoding step. The idea of having a single-encode step makes the process of extracting out symbols modular and computationally efficient. However, all methods above fail on non-toy scenes. We attribute this failure to the *what* or *where* problem, a.k.a. symmetry problem, pointed out in the ‘Instance Slots’ Section in Greff *et al.*[25]: “bottom-up information is unable to break the symmetry amongst slots and thus ends up assigning the same content to each one.” In GFS-Nets, we propose to circumvent this problem using curriculum training.

Shape program synthesis and analysis-by-synthesis GFS-Nets is also related to works in analysis-by-synthesis [37], program synthesis for shape prediction [60, 14, 41], as well as earlier works on Computer Vision, such as Marr’s 3D sketch [45] which involves representing a scene in terms of generalized cylinders and their syntactic relations to each other. In place of data-driven Markov Chain Monte Carlo search of analysis-by-synthesis methods that require good initialization, our slow inference searches in the space of primitives by gradient descent. In contrast to program synthesis methods, it does not require a predefined domain-specific language (DSL) or program annotations for visual structures [41], rather, it discovers compositions over primitives via its slow inference.

Unsupervised 3D Part Discovery There are numerous methods that attempt the decomposition of complex 3D shapes into primitive parts without primitive supervision[33, 21, 51, 18, 12, 62, 11, 9]. Traditional primitives include cuboids [62, 50], superquadrics [53, 51], and convexes [11, 8]. [20] proposes a 3D representation that decomposes space into a structured set of implicit functions [21]. Neural Parts [52] represents arbitrarily complex genus-zero shapes and thus yields comparatively expressive parts. However the resulting parts of these methods[20, 52] are still not semantically meaningful and the decomposition is highly dependent on the number of parts initialized. The work of [66] does 3D part reconstruction directly from a 2D image input without access to any ground-truth 3D shapes for training. However, both [66] and [52] take as input the number of parts, and different decompositions are predicted with varying part numbers. There is no clear way to select the right number of parts. In our case, parts can be quite complex: pairs of parallel surfaces, quadruplets of legs, as we use implicit functions to represent them. Moreover GFS-Nets’s dynamic attention-based routing allows it to infer different number of parts for each input scene.

3 Generating Fast and Slow (GFS-Nets)

The goal of GFS-Nets is to decompose an unfamiliar scene into familiar entities. The model encodes the scene, which can be a 3D point cloud or an RGB image, into a set of slot vectors and decodes back the scene using a decoder shared across slots. The model uses the slot attention feature-to-slot mapping of Locatello *et al.*[42], where visual features are softly partitioned across slots through attention, and update the slot vectors.

GFS-Nets is trained with a curriculum. First, it is first trained to autoencode individual entities using a single slot in their bottleneck, as shown in the left of Figure 2 and described in Section 3.2. Next, GFS-Nets is trained to autoencode complex scenes using multiple slots in their bottleneck, as shown

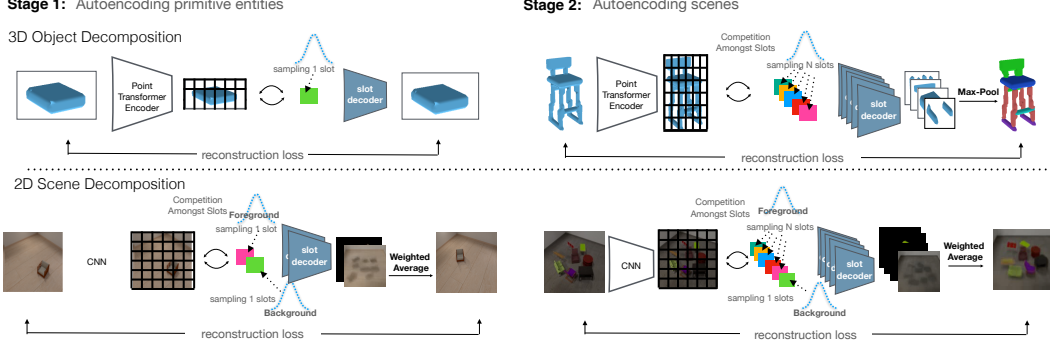


Figure 2: **Stages of training in GFS-Nets.** *Left:* Autoencoding primitive entities. Primitive entities are point clouds of individual parts in 3D, and single object RGB images in 2D. GFS-Nets use one and two slots respectively, for primitive entity autoencoding in 3D point clouds and RGB images. *Right:* Autoencoding scenes. Scenes can be 3D points clouds of whole object or multi-object RGB images. In this stage, GFS-Nets use more slots in their bottleneck. During inference, some slots get mapped to nothing in the input scene. In this way, GFS-Nets is able to infer decompositions with a varying number of entities, in an unsupervised manner.

on the right of Figure 2 and described in Section 3.3. We first describe slot attention of [42] in Section 3.1, as well as encoders and decoders for 3D point clouds [70, 40] and RGB images [47, 67] which our method builds upon.

3.1 Background

3.1.1 Routing from input to slots with iterative attention

Many approaches instantiate slots (a.k.a. query vectors) from 2D visual feature maps or 3D point feature clouds [5, 48]. Most works use standard cross attention operations [61, 22, 42] or iterative cross (features to slots) and self-attention (slot-to-slots) operations [5] to map a set of N input feature vectors to a set of K slot vectors. Competition amongst slots and iterative routing introduced in [22, 42] encourages the slots not to encode information in a redundant manner.

Given a visual scene encoded as a set of feature vectors $M \in \mathbb{R}^{N \times C}$ and K randomly initialized slots sampled from a multivariate Gaussian distribution with a diagonal covariance $s \sim \mathcal{N}(\mu, \text{Diag}(\sigma^2)) \in \mathbb{R}^{K \times D}$, where $\mu, \sigma \in \mathbb{R}^C$ are learnable parameters of the Gaussian, Slot Attention[42] computes an attention map a between the feature map M and the slots s :

$$a = \text{Softmax}(k(M) \cdot q(s)^T, \text{axis}=1) \in \mathbb{R}^{N \times K}. \quad (1)$$

k , q , and v are learnable linear transformations that map inputs and slots to a common dimension D . The softmax normalization over slots ensures competition amongst them to attend to a specific feature vector in M . We then generate the update vector for each slot given the computed attention and input feature maps:

$$\text{updates} = a^T v(M) \in \mathbb{R}^{K \times C}, \text{ where } a_{i,1} = \frac{a_{i,1}}{\sum_{i=0}^N a_{i,1}} \quad (2)$$

which we then use to do a gated update on each slot feature using a GRU[10]: $\text{slots} = \text{GRU}(\text{state} = \text{slots}, \text{input} = \text{updates})$. We iterate over the steps of calculating attention and updating the slots for 3 timesteps.

Since the slot vectors are sampled from a shared Gaussian distribution, slot attention[42] can instantiate a variable number of slots in different scenes, in contrast to DETR encoder [5] where a constant number of (deterministic) slot vectors are learned and used in each scene. For a more detailed description, we refer the reader to [42].

Slot attention is one of many ways of mapping inputs to slots. GFS-Nets is agnostic to the method of mapping visual features to slot vectors and we ablate different encoders in our experiment section 4.1.2.

3.1.2 Encoding and decoding 3D point clouds

GFS-Nets featurize a 3D point cloud using a 3D point transformer [70] which maps the 3D input points to a set of M feature vectors of C dimensions each. We set M to 128 and C to 64 in our experiments.

GFS-Nets decodes 3D point clouds from each slot using implicit functions [46]. Specifically, each decoder takes in as input the slot vector s_i and an (X, Y, Z) location and returns the corresponding occupancy score $o_{x,y,z} = \text{Dec}(s_i, (x, y, z))$, where Dec is a multi-block ResNet MLP similar to that of Lal et al. [40].

3.1.3 Encoding and decoding RGB images

GFS-Nets featurize an RGB image concatenated with its pixel coordinates using a convolutional neural network (CNN) and produce a feature map of $H \times W \times C$ where H, W are the spatial dimensions and C is the feature dimension. We reshape the output feature map to a set of vectors $M \in \mathbb{R}^{(H \times W) \times C}$. We set C as 64.

GFS-Nets decodes RGB pixels from each slot using conditional NeRFs [47] that predict the color c and volume density σ at each 3D location of the latent 3D volume. Thus given slot vector s_i , spatial location x and viewing direction d , the MLP based decoder network Dec predicts $(c, \sigma) = \text{Dec}(s_i, (x, d))$. The color c and density σ of the points are composed together across slots using density σ weighted averaging. Specifically, $\hat{\sigma} = \sum_{i=0}^K w_i \sigma_i$, $\hat{c} = \sum_{i=0}^K w_i c_i$ where $w_i = \sigma_i / \sum_{i=0}^K \sigma_i$. We then use differentiable volume rendering of the composed $\hat{\sigma}$ and \hat{c} using camera ray casting to generate the final RGB prediction similar to [67].

3.2 Autoencoding primitive entities

In this stage, we train GFS-Nets using a set of primitive entities. Primitive entities are generic or semantic parts of objects in 3D point clouds or single object multi-view images, as shown on the left of Figure 2. We experiment with different primitive sets in the experimental section.

For encoding 3D point cloud primitive entities, we sample **a single slot** in the entity bottleneck of GFS-Nets, that is, $s \in \mathbb{R}^{1 \times D}$. For encoding a 2D RGB image, we sample two slots, one from a foreground learnable Gaussian distribution and one from a background learnable Gaussian distribution, that is $s \in \mathbb{R}^{2 \times D}$. In a forward pass of the encoder, we then update slots s as described in Section 3.1.

3.3 Autoencoding scenes

In this stage, we initialize the parameters of the model with the parameters learnt in the previous section. The only difference is that instead of instantiating one or two slots we now instantiate K slot vectors, where

K is the maximum number of entities any example in the dataset could possibly have. For RGB input, we sample K slots from the foreground Gaussian and one slot from the background Gaussian since there can be many foreground objects but only a single background.

Finally we pass each of our slot vectors through the slot decoder $out_k = \text{Dec}(s_k)$, whose weights are shared across slots. For decoding RGB images, GFS-Nets uses separate decoders for foreground and background. For 3D point cloud input, out_k represents the probability of sampled points being occupied by a given slot s_k ; we aggregate them by max pooling across the slot dimension. We then use binary cross entropy loss between the ground truth occupancies and the predictions.

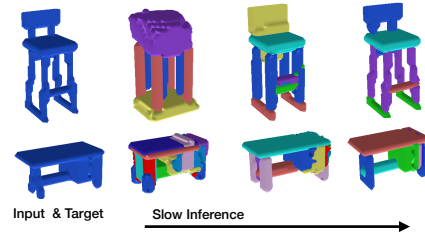


Figure 3: **Slow Inference in GFS-Nets:** First columns shows the input point cloud. In 2nd to 4th columns, we show the segmentation results throughout training iterations of the autoencoding process. Segmentation improves over time with reconstruction quality, despite that the autoencoding loss only penalizes the reconstruction error.

For RGB input, out_k represents the color and density of points being occupied by a given slot s_k ; we aggregate across slots as described in Section 3.1 and use pixel mean squared error loss.

Fast and Slow Inference Fast inference refers to a single forward pass through the trained network. Slow inference refers to test-time adaptation of all the learnable parameters of the model for autoencoding a single scene using the method described in Section 3.3. Our experiment show that GFS-Nets with slow inference always achieves much better scene decomposition than GFS-Nets with fast inference, especially in out-of-distribution scenes. In contrast, slow inference on entity-centric generative models without curricula, such as [42], results in improved reconstruction accuracy but worse scene decomposition performance, due to the reconstruction-segmentation trade off noted in [15].

4 Experiments

We test GFS-Nets in the tasks of segmenting object parts in 3D point clouds and segmenting objects and predicting views in 2D multiview images under varying amounts of label supervision. We compare GFS-Nets against previous state-of-the-art models in each task. We further quantify contribution of various design choices, namely, slow inference through reconstruction feedback, curriculum learning, and encoder architecture. We evaluate segmentation performance on scenes within the training distribution as well as out-of-distribution scenes.

In each setup, our experiments aim to answer the following questions:

- How does GFS-Nets compare against state-of-the-art 2D and 3D scene decomposition models ([43, 67, 5]) under varying amounts of label supervision?
- How much, if any, slow inference through reconstruction feedback improves segmentation accuracy in GFS-Nets and its variants?
- How much curriculum or supervision during training contributes to segmentation performance?

4.1 Segmenting parts in 3D object point clouds

In this task, the input is a complete 3D point cloud of an object, and the goal is to segment the 3D point cloud into different semantic part instances. We consider two setups with different segmentation supervision:

1. Segmenting without any part segmentation labels.
2. Segmenting with part segmentation labels from a related object category.

4.1.1 Segmenting 3D pointclouds without segmentation supervision

In this setup, we test our model and baselines in segmenting the test objects of the Chair and Table categories in the PartNet benchmark [49] without access to any ground-truth 3D segmentations at training time. For training, we provide our model and baselines access to **unsegmented** objects of the Chair category in the trainset of the PartNet benchmark. We evaluate the segmentation performance of our model and baselines with the Adjusted Random Index (ARI) [28] using the implementation of Kabra *et al.* [32], which calculates the similarity between two-point clusters while being invariant to the ordering of the cluster centers. Neither our model nor the baselines have access to ground-truth 3D segmentations during training; as a result, they may output 3D parts of coarser or finer resolution. PartNet contains three different levels of ground-truth segmentation labels with progressively finer segmentation granularity. We evaluate the ARI score across all three levels of PartNet and report the best of three for all the models.

Our model and two amongst three of our baselines, specifically (**PQ-Nets**[64] of Wu *et al.* and **Shape2Prog** [60] of Tian *et al.*) expect access to a set of 3D primitive parts during training. In this setup, we consider the primitive part dataset introduced by Shape2Prog [60]. The dataset consists of differently sized cubes, cuboids, and discs; we call it the generic primitive set since it is not related to any semantic object category but rather consists of generic 3D parts that remind of generalized

cylinders of Marr *et al.*[45]. Please refer to the supplementary file for visualizations of the generic primitive set.

Baselines: We compare our model against the following baselines: **(i) PQ-Nets** of Wu *et al.*[64] assume access to a set of primitive 3D parts, similar to our model, and learn a primitive part autoencoder. Whole object autoencoding uses a sequential model that encodes the 3D point cloud into a 1D latent vector and sequentially decodes parts using the part decoder. We use the publicly available code to train the model to autoencode the generic primitive set and the unlabelled 3D point clouds of the instances of the Chair category in PartNet. **(ii) Shape2Prog** of Tian *et al.*[60] is a shape program synthesis method that is trained supervised to predict shape programs from object 3D point clouds. The program represents the part category, location, and the symmetry relations among the parts (if any). **Shape2Prog** introduced two synthetically generated datasets that helped the model parse 3D pointclouds from ShapeNet [6] into shape programs without any supervision: i) the generic primitive set we discussed earlier which they use to train their part decoders, and ii) a synthetic whole shape dataset of chairs and tables generated programmatically alongside its respective ground-truth programs. Their model requires supervised pre-training on the dataset of synthetic whole shapes paired with programs. We therefore use their publicly available model weights trained on synthetic whole shapes to further train on PartNet Chairs. Note that no other baseline nor GFS-Nets assumes access to the synthetic whole shape dataset. **(iii) StructImplicit** of Genova *et al.*[21] encodes the object point cloud into a 1D latent and then uses an MLP to map it to a *fixed* number of part vectors, each represented with a 3D implicit function. We use the publicly available code to train the model on the unlabelled Chair PartNet dataset. We set the number of part vectors based on the maximum number of parts any example in the dataset can have.

Method	in-dist (Chairs)		out-of-dist (Tables)	
	Fast Infer.	Slow Infer.	Fast Infer.	Slow Infer.
Shape2Prog [60]	0.28	0.53	0.23	0.40
PQ-Nets [64]	0.20	0.31	0.17	0.21
StructImplicit [21]	0.25	0.23	0.17	0.17
GFS-Nets w/o curriculum	0.41	0.35	0.47	0.38
GFS-Nets-PrimitiveOnly	0.42	0.48	0.49	0.57
GFS-Nets	0.57	0.62	0.60	0.69

Table 1: **ARI Segmentation accuracy (higher the better)** in the test set of Chair (in-distribution) and Table category of PartNet(out-of-distribution). GFS-Nets significantly outperform all of the baselines.

We consider the following ablative versions of GFS-Nets: **(i) GFS-Nets w/o curriculum** is trained without the primitive learning stage. Instead, it is trained to directly autoencode the aggregated dataset of both the generic primitive set and 3D point clouds of training chairs (as described in Section 3.3). This version of our model coincides precisely with the previous work of Slot Attention of Locatello *et al.* [42]. **(ii) GFS-Nets-PrimitiveOnly** is trained only on the generic primitive set (as described in Section 3.2) and is not trained to autoencode and 3D point clouds of training chairs.

Any model that has a decoder that reconstructs the input can be additionally trained with slow inference with reconstruction feedback at each test scene independently. GFS-Nets, PQ-Nets [64], Shape2Prog [60], and StructImplicit [21] are all equipped with such decoders. We thus test both our model and the baselines in both modes of fast (regular) and slow inference; we use the notation [modelname]-Slow and [modelname]-Fast to denote the model with the corresponding inference type.

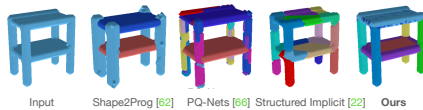


Figure 4: Unsupervised 3D segmentations of out-of-distribution scenes for GFS-Nets and baselines.

We show quantitative and qualitative results of our model and the baselines in Table 4 and Figure 4 respectively. From the Table 4 we draw the following conclusions:

Segregation into slot like entities using attention helps. GFS-Nets significantly outperform PQ-Nets [64] and StructImplicit [21] that map the input object 3D pointcloud into a 1D latent vec-

tor instead of maintaining segregated representations into multiple slot latent vectors. Moreover, GFS-Nets’s ability to dynamically route information per example via attention allows it to automatically estimate the number of parts required per input 3D pointcloud due to which, different instances can have a different number of active slots, while StructImplicit segregates all object instances into a constant number of parts, which is a hyperparameter of the model.

Curriculum training helps. GFS-Nets-Fast outperform by a large margin GFS-Nets w/o curriculum-Fast, the version of GFS-Nets that uses exactly the same training data and inference method, but is trained without curriculum.

Slow inference through reconstruction feedback helps in the presence of curriculum and hurts in the absence of it. GFS-Nets, GFS-Nets-PrimitiveOnly, PQ-Nets and Shape2Prog all have curriculum through a primitive learning stage, while GFS-Nets w/o curriculum and StructImplicit do not. GFS-Nets-Slow outperform GFS-Nets-Fast and GFS-Nets-PrimitiveOnly-Slow outperform GFS-Nets-PrimitiveOnly-Fast, PQ-Nets-Slow outperform PQ-Nets-Fast and ShapeProg-Slow outperform ShapeProg-Fast. In contrast, GFS-Nets w/o curriculum-Slow performs worse than GFS-Nets w/o curriculum-Fast and StructImplicit-Slow does worse than StructImplicit-Fast. Such trade-off between reconstruction and segmentation in generative model for scene decomposition has been pointed out by Engelcke *et al.* [15] and our results verify their findings.

Unsupervised autoencoding of complex (non primitive) 3D pointclouds helps. GFS-Nets outperforms GFS-Nets-PrimitiveOnly.

Please see the supplementary file for more experimental details. We visualize intermediate iterations during slow inference in our supplementary video.

4.1.2 Segmenting 3D pointclouds with weak segmentation supervision

In this setup, we test our model and baselines in segmenting the test objects of the Chair and Table categories in the PartNet benchmark [49], with access to ground-truth pointcloud segmentation labels of level-3 in the training 3D pointclouds of the Chair category in PartNet. We use the ARI segmentation metric to evaluate the level 3 ground-truth labels of PartNet for our model and baselines.

GFS-Nets assume access to a set of 3D primitive parts for primitive entity learning. In this setup, we consider each 3D part from each training chair 3D pointcloud example as a separate 3D primitive. We augment these primitives by random rotations and translations; we call it the semantic primitive set.

Method	in-dist (Chair)		out-dist (Table)	
	Fast Infer.	Slow Infer.	Fast Infer.	Slow Infer.
3D-DETR [5]	0.67	n/a	0.41	n/a
Learning2Group [43]	0.62	n/a	0.46	n/a
GFS-Nets w/o supervision	0.40	0.44	0.31	0.48
GFS-Nets w/o curriculum	0.61	0.66	0.48	0.60
GFS-Nets w/o Gaussian	0.65	0.62	0.38	0.58
GFS-Nets w/o SlotAttention	0.58	0.55	0.31	0.44
GFS-Nets	0.64	0.67	0.51	0.63

Table 2: **ARI Segmentation scores (higher the better)** on the test set of the Chair category (in-distribution) and on the test set of the Table category (out-of-distribution). All the models are trained using the training set of the Chair category in PartNet.

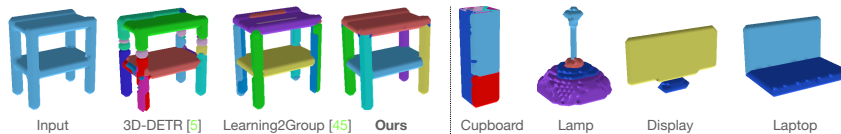


Figure 5: **Out-of-distribution 3D segmentation results** for GFS-Nets and baselines. *Left:* 3D segmentation results on out-of-distribution shapes for GFS-Nets and baselines. *Right:* 3D segmentation results for out-of-distribution categories of PartNet for GFS-Nets trained semi-supervised only on the Chair category.

Baselines We compare our model against the following baselines: (i) **Learning2Group** of Luo *et al.* [43] progressively groups points into segments by learning pairwise grouping decisions parameterized by features of the two point clusters. Given the intermediate grouping decisions are not supervised, and the non-differentiability of their grouping functions they use reinforcement learning gradients for training using supervision from the final segmentation. We used the publicly available code and trained the model in the training set of the Chair category. (ii) **3D-DETR** is the equivalent of 2D-DETR of Carion *et al.* [5] for 3D point-cloud segmentation. A point transformer [70] encodes the input 3D point cloud into a set of point features which are mapped to a learnable set of 16 parametric queries using 6 layers of transformer encoder (self-attention) and decoder (cross-attention) blocks. Then, each query decodes a 3D point segmentation map by computing multi-head attention with the encoded point features. The attention for each query is then upsampled using a point transformer decoder. The predicted 3D part segments are then matched against ground-truth segments using Hungarian matching [36], and error is computed for each pair of matched ground-truth part and paired prediction through binary cross-entropy loss.

We train GFS-Nets in a semi-supervised way combining an (unsupervised) autoencoding loss and a supervised part segmentation loss. Specifically alongside the autoencoding loss of the training 3D point clouds, we apply a part segmentation loss via Hungarian matching between the decoded 3D part shapes from the slots and the ground-truth part 3D segments in the training set, similar to our 3D-DETR baseline model. For this case we did not consider the models of PQ-Nets [64], Shape2Prog [60] or [21] as it is not clear on how to train them in a semi-supervised manner using segmentation supervision.

We consider the following ablative versions of GFS-Nets: (i) **GFS-Nets w/o supervision** is trained without the supervised part segmentation loss, so the model only receives gradients from autoencoding the primitives and whole chairs. (ii) **GFS-Nets w/o Gaussian** does not initialize the slot vectors by sampling from a learnable Gaussian but instead uses separate learnable query vectors for each slot, similar to DETR. (iii) **GFS-Nets w/o SlotAttention** does not use Slot Attention for mapping point features to slots. Instead it maps 3D point features to slots via iterative layers of cross (query to point) and self (query-to-query) attention layers on learnable query vectors similar to 3D-DETR and DETR [5]. Please note that slots and queries represent the same thing, but we use the terminology of DETR [5] in this case.

We report the fast and slow inference results for all ablative versions of our model. Our baselines in this case, 3D-DETR and Learning2Group [43] are feedforward in nature, they are not equipped with decoders, and thus cannot be evaluated with slow inference.

We show quantitative results in Table 2 and qualitative results in Figure 5. GFS-Nets significantly outperform the baselines, and GFS-Nets-Slow results in a significant boost in performance ($\sim 50\%$) over the feedforward inference in our model, GFS-Nets-Fast. We draw the following conclusions from Table 2:

Supervision can be used as a replacement to curriculum training. GFS-Nets only slightly outperforms GFS-Nets w/o curriculum in this setup.

Generic primitives generalize a lot better than category-specific primitives or supervision. GFS-Nets in Table 4 (that uses the generic primitive set) outperforms GFS-Nets and GFS-Nets-w/o supervision in Table 2 (that use the semantic primitive set) in OOD generalization on the Table category.

Competition amongst slots helps. GFS-Nets outperforms GFS-Nets w/o SlotAttention, thus showing competition amongst slot vectors during encoding helps generalization.

Slow inference through reconstruction feedback helps OOD generalization of GFS-Nets. Our baselines 3D-DETR and Learning2Group [43] are feed-forward in nature, they lack any form of reconstruction feedback, and thus cannot adapt as our model through such feedback.

4.2 RGB image segmentation and view prediction without supervision

In this task, the input to the model is a single RGB image that contains multiple objects and the goal is to segment the image into objects and predict images from alternative scene viewpoints. We test our model in the following two benchmarks: (i) **CLEVR** [31]: this is a commonly used benchmark for unsupervised image segmentation. The dataset contains scenes that consist of 5 to 7 spheres, cylinders and cubes, in various colors and materials. We use the multi-view version of this dataset introduced by uORF [67] as the training set.



Figure 6: **RGB image decomposition and view prediction** for GFS-Nets and baselines. *Left*: From 2nd to 5th column we show reconstructed RGB image superimposed with the predicted segmentation masks for our model and baselines. *Right*: From 2nd to 5th column we show RGB renderings across multiple novel viewpoints for GFS-Nets.

(ii) **Room Diverse++**: Room Diverse++ builds upon Room Diverse, which is a multi-view RGB dataset introduced by Yu *et al.* [67] where they placed multiple ShapeNet Chairs with different solid colors (Red, Blue etc) onto a textured background. We make it more challenging by applying real-world ShapeNet textures to the object meshes instead of solid colors. We also add more categories such as Beds, Benches, Tables, Cupboards, and Sofas to the Chair category. We sample about 6 to 10 objects in each scene.

Baselines We consider the following baselines: (i) **Slot attention** Locatello *et al.* [42] autoencodes single view RGB images, and is described in Section 3.1. We use the publicly available code and train the model on the training set for single image autoencoding.

(ii) **uORF** Yu *et al.* [67] is a recent work that extends Slot Attention of [42] to multiview RGB images.

During primitive training for our model, we assume access to posed RGB images of single object scenes as our primitive dataset, as shown in Figure 2 bottom row left. Our implementation builds upon [67] with the difference that we consider curriculum and slow inference at test time. We ablate slow and fast versions of our model.

Method	CLEVR	Room Diverse++
uORF [67]	0.86	0.65
Slot Attention [42]	0.04	0.28
GFS-Nets-Fast	0.86	0.74
GFS-Nets-Slow	0.92	0.81

Table 3: **ARI Segmentation accuracy (higher the better)** on CLEVR and RoomDiverse++ Dataset for RGB input. As can be seen GFS-Nets-Slow outperform all other baselines in segmentation accuracy.

We show quantitative results of GFS-Nets and the baselines in Table 3 and qualitative view prediction and segmentation results in Figure 6. GFS-Nets-slow outperforms the baselines significantly. This suggests that curriculum and slow inference are key to out-of-distribution generalization. Slot Attention [42] achieves very low accuracy on CLEVR, as it is unable to disentangle background from foreground, as also reported in Yu *et al.* [67]. We further visualize GFS-Nets’s ability to render novel viewpoints in the supplementary video.

Limitations and Future Work: GFS-Nets has limitations that offer several promising directions for future work. First, GFS-Nets does not model pairwise interactions between slots. Incorporating such interactions is a direct avenue for future work. Second, entity decomposition is inherently hierarchical: objects - parts - subparts while GFS-Nets currently models a flat non-hierarchical list of entities. Third, amortizing the slow inference in feedforward predictions could be another interesting extension of the current framework. Finally, although GFS-Nets is tested on far more realistic data than considered in prior works, scaling the model to more realistic datasets used by mainstream supervised visual detectors, is a direct avenue for future work.

5 Conclusion

We presented GFS-Nets, a model that segments scenes into entities while autoencoding the input scene through slot-based encoders and shared slot decoders. We showed that curriculum training in the form of primitive entities, single object scenes, or labelled segmentations is necessary to break the trade-off between segmentation and reconstruction observed in previous works [15]. This permits GFS-Nets to use slow inference through reconstruction feedback and drastically improve the decompositions of out-of-distribution scenes. We believe that a compositional approach to AI, in terms of grounded symbol-like representations [44, 38] is of fundamental importance for realizing human-level generalization [1], and we hope that GFS-Nets may contribute towards that goal.

References

- [1] Bahdanau, D., Murty, S., Noukhovitch, M., Nguyen, T.H., de Vries, H., Courville, A.: Systematic generalization: what is required and can it be learned? arXiv preprint arXiv:1811.12889 (2018)
- [2] Barbu, A., Mayo, D., Alverio, J., Luo, W., Wang, C., Gutfreund, D., Tenenbaum, J., Katz, B.: Objectnet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. *Advances in neural information processing systems* **32** (2019)
- [3] van Bergen, R.S., Kriegeskorte, N.: Going in circles is the way forward: the role of recurrence in visual inference. *Current Opinion in Neurobiology* **65**, 176–193 (Dec 2020). <https://doi.org/10.1016/j.conb.2020.11.009>, <http://dx.doi.org/10.1016/j.conb.2020.11.009>
- [4] Burgess, C.P., Matthey, L., Watters, N., Kabra, R., Higgins, I., Botvinick, M., Lerchner, A.: Monet: Unsupervised scene decomposition and representation. arXiv preprint arXiv:1901.11390 (2019)
- [5] Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: *European Conference on Computer Vision*. pp. 213–229. Springer (2020)
- [6] Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., Xiao, J., Yi, L., Yu, F.: ShapeNet: An Information-Rich 3D Model Repository. Tech. Rep. arXiv:1512.03012 [cs.GR], Stanford University — Princeton University — Toyota Technological Institute at Chicago (2015)
- [7] Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., et al.: Shapenet: An information-rich 3d model repository. arXiv preprint arXiv:1512.03012 (2015)
- [8] Chen, Z., Tagliasacchi, A., Zhang, H.: Bsp-net: Generating compact meshes via binary space partitioning. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 45–54 (2020)
- [9] Chen, Z., Yin, K., Fisher, M., Chaudhuri, S., Zhang, H.: Bae-net: branched autoencoder for shape co-segmentation. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 8490–8499 (2019)
- [10] Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., Bengio, Y.: Learning phrase representations using rnn encoder-decoder for statistical machine translation. arXiv preprint arXiv:1406.1078 (2014)
- [11] Deng, B., Genova, K., Yazdani, S., Bouaziz, S., Hinton, G., Tagliasacchi, A.: Cvxnet: Learnable convex decomposition. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 31–44 (2020)
- [12] Deprelle, T., Groueix, T., Fisher, M., Kim, V.G., Russell, B.C., Aubry, M.: Learning elementary structures for 3d shape generation and matching. arXiv preprint arXiv:1908.04725 (2019)
- [13] Du, Y., Smith, K., Ulman, T., Tenenbaum, J., Wu, J.: Unsupervised discovery of 3d physical objects from video. arXiv preprint arXiv:2007.12348 (2020)
- [14] Ellis, K., Wong, C., Nye, M., Sable-Meyer, M., Cary, L., Morales, L., Hewitt, L., Solar-Lezama, A., Tenenbaum, J.B.: Dreamcoder: Growing generalizable, interpretable knowledge with wake-sleep bayesian program learning. arXiv preprint arXiv:2006.08381 (2020)
- [15] Engelcke, M., Jones, O.P., Posner, I.: Reconstruction bottlenecks in object-centric generative models. arXiv preprint arXiv:2007.06245 (2020)
- [16] Engelcke, M., Kosiorek, A.R., Jones, O.P., Posner, I.: Genesis: Generative scene inference and sampling with object-centric latent representations. arXiv preprint arXiv:1907.13052 (2019)

- [17] Eslami, S., Heess, N., Weber, T., Tassa, Y., Szepesvari, D., Kavukcuoglu, K., Hinton, G.E.: Attend, infer, repeat: Fast scene understanding with generative models. arXiv preprint arXiv:1603.08575 (2016)
- [18] Gao, L., Yang, J., Wu, T., Yuan, Y.J., Fu, H., Lai, Y.K., Zhang, H.: Sdm-net: Deep generative network for structured deformable mesh. *ACM Transactions on Graphics (TOG)* **38**(6), 1–15 (2019)
- [19] Geirhos, R., Rubisch, P., Michaelis, C., Bethge, M., Wichmann, F.A., Brendel, W.: Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. arXiv preprint arXiv:1811.12231 (2018)
- [20] Genova, K., Cole, F., Sud, A., Sarna, A., Funkhouser, T.: Local deep implicit functions for 3d shape. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 4857–4866 (2020)
- [21] Genova, K., Cole, F., Vlasic, D., Sarna, A., Freeman, W.T., Funkhouser, T.: Learning shape templates with structured implicit functions. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 7154–7164 (2019)
- [22] Goyal, A., Lamb, A., Hoffmann, J., Sodhani, S., Levine, S., Bengio, Y., Schölkopf, B.: Recurrent independent mechanisms. arXiv preprint arXiv:1909.10893 (2019)
- [23] Greff, K., Kaufman, R.L., Kabra, R., Watters, N., Burgess, C., Zoran, D., Matthey, L., Botvinick, M., Lerchner, A.: Multi-object representation learning with iterative variational inference. In: *International Conference on Machine Learning*. pp. 2424–2433. PMLR (2019)
- [24] Greff, K., Rasmus, A., Berglund, M., Hao, T.H., Schmidhuber, J., Valpola, H.: Tagger: Deep unsupervised perceptual grouping. arXiv preprint arXiv:1606.06724 (2016)
- [25] Greff, K., van Steenkiste, S., Schmidhuber, J.: On the binding problem in artificial neural networks. *CoRR* **abs/2012.05208** (2020), <https://arxiv.org/abs/2012.05208>
- [26] Guerin, F.: Projection: A mechanism for human-like reasoning in artificial intelligence. *CoRR* **abs/2103.13512** (2021), <https://arxiv.org/abs/2103.13512>
- [27] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 770–778 (2016)
- [28] Hubert, L., Arabie, P.: Comparing partitions. *Journal of classification* **2**(1), 193–218 (1985)
- [29] Jo, J., Bengio, Y.: Measuring the tendency of cnns to learn surface statistical regularities. arXiv preprint arXiv:1711.11561 (2017)
- [30] Johnson, J., Hariharan, B., Van Der Maaten, L., Fei-Fei, L., Lawrence Zitnick, C., Girshick, R.: Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 2901–2910 (2017)
- [31] Johnson, J., Krishna, R., Stark, M., Li, L.J., Shamma, D., Bernstein, M., Fei-Fei, L.: Image retrieval using scene graphs. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2015)
- [32] Kabra, R., Burgess, C., Matthey, L., Kaufman, R.L., Greff, K., Reynolds, M., Lerchner, A.: Multi-object datasets. <https://github.com/deepmind/multi-object-datasets/> (2019)
- [33] Kato, H., Harada, T.: Learning view priors for single-view 3d reconstruction. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 9778–9787 (2019)
- [34] Kipf, T., Elsayed, G.F., Mahendran, A., Stone, A., Sabour, S., Heigold, G., Jonschkowski, R., Dosovitskiy, A., Greff, K.: Conditional object-centric learning from video. arXiv preprint arXiv:2111.12594 (2021)

- [35] Kosiorek, A., Kim, H., Teh, Y.W., Posner, I.: Sequential attend, infer, repeat: Generative modelling of moving objects. *Advances in Neural Information Processing Systems* **31**, 8606–8616 (2018)
- [36] Kuhn, H.W.: The hungarian method for the assignment problem. *Naval research logistics quarterly* **2**(1-2), 83–97 (1955)
- [37] Kulkarni, T.D., Kohli, P., Tenenbaum, J.B., Mansinghka, V.: Picture: A probabilistic programming language for scene perception. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 4390–4399 (2015)
- [38] Lake, B.M., Salakhutdinov, R., Tenenbaum, J.B.: Human-level concept learning through probabilistic program induction. *Science* **350**(6266), 1332–1338 (2015)
- [39] Lake, B.M., Ullman, T.D., Tenenbaum, J.B., Gershman, S.J.: Building machines that learn and think like people. *Behavioral and brain sciences* **40** (2017)
- [40] Lal, S., Prabhudesai, M., Mediratta, I., Harley, A.W., Fragkiadaki, K.: Coconets: Continuous contrastive 3d scene representations (2021)
- [41] Li, Y., Mao, J., Zhang, X., Freeman, W.T., Tenenbaum, J.B., Snavely, N., Wu, J.: Multi-plane program induction with 3d box priors. *arXiv preprint arXiv:2011.10007* (2020)
- [42] Locatello, F., Weissenborn, D., Unterthiner, T., Mahendran, A., Heigold, G., Uszkoreit, J., Dosovitskiy, A., Kipf, T.: Object-centric learning with slot attention (2020)
- [43] Luo, T., Mo, K., Huang, Z., Xu, J., Hu, S., Wang, L., Su, H.: Learning to group: A bottom-up framework for 3d part discovery in unseen categories. *arXiv preprint arXiv:2002.06478* (2020)
- [44] Marcus, G.F.: *The algebraic mind: Integrating connectionism and cognitive science*. MIT press (2003)
- [45] Marr, D.: *Vision: A Computational Investigation into the Human Representation and Processing of Visual Information*. W.H. Freeman and Co., New York, NY (1982)
- [46] Mescheder, L., Oechsle, M., Niemeyer, M., Nowozin, S., Geiger, A.: Occupancy networks: Learning 3d reconstruction in function space. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 4460–4470 (2019)
- [47] Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: *European conference on computer vision*. pp. 405–421. Springer (2020)
- [48] Misra, I., Girdhar, R., Joulin, A.: An end-to-end transformer model for 3d object detection. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 2906–2917 (2021)
- [49] Mo, K., Zhu, S., Chang, A.X., Yi, L., Tripathi, S., Guibas, L.J., Su, H.: Partnet: A large-scale benchmark for fine-grained and hierarchical part-level 3d object understanding. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 909–918 (2019)
- [50] Niu, C., Li, J., Xu, K.: Im2struct: Recovering 3d shape structure from a single rgb image. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 4521–4529 (2018)
- [51] Paschalidou, D., Gool, L.V., Geiger, A.: Learning unsupervised hierarchical part decomposition of 3d objects from a single rgb image. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 1060–1070 (2020)
- [52] Paschalidou, D., Katharopoulos, A., Geiger, A., Fidler, S.: Neural parts: Learning expressive 3d shape abstractions with invertible neural networks. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 3204–3215 (2021)

- [53] Paschalidou, D., Ulusoy, A.O., Geiger, A.: Superquadrics revisited: Learning 3d shape parsing beyond cuboids. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10344–10353 (2019)
- [54] Pathak, D., Agrawal, P., Efros, A.A., Darrell, T.: Curiosity-driven exploration by self-supervised prediction. In: International conference on machine learning. pp. 2778–2787. PMLR (2017)
- [55] Rahaman, N., Goyal, A., Gondal, M.W., Wuthrich, M., Bauer, S., Sharma, Y., Bengio, Y., Schölkopf, B.: S2rms: Spatially structured recurrent modules. arXiv preprint arXiv:2007.06533 (2020)
- [56] Ren, S., He, K., Girshick, R.B., Sun, J.: Faster R-CNN: towards real-time object detection with region proposal networks. CoRR **abs/1506.01497** (2015), <http://arxiv.org/abs/1506.01497>
- [57] Sabour, S., Frosst, N., Hinton, G.E.: Dynamic routing between capsules. arXiv preprint arXiv:1710.09829 (2017)
- [58] Sitzmann, V., Chan, E., Tucker, R., Snavely, N., Wetzstein, G.: Metasdf: Meta-learning signed distance functions. Advances in Neural Information Processing Systems **33**, 10136–10147 (2020)
- [59] Sun, Y., Wang, X., Liu, Z., Miller, J., Efros, A., Hardt, M.: Test-time training with self-supervision for generalization under distribution shifts. In: International Conference on Machine Learning. pp. 9229–9248. PMLR (2020)
- [60] Tian, Y., Luo, A., Sun, X., Ellis, K., Freeman, W.T., Tenenbaum, J.B., Wu, J.: Learning to infer and execute 3d shape programs. arXiv preprint arXiv:1901.02875 (2019)
- [61] Tsai, Y.H.H., Srivastava, N., Goh, H., Salakhutdinov, R.: Capsules with inverted dot-product attention routing. arXiv preprint arXiv:2002.04764 (2020)
- [62] Tulsiani, S., Su, H., Guibas, L.J., Efros, A.A., Malik, J.: Learning shape abstractions by assembling volumetric primitives. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2635–2643 (2017)
- [63] Van Steenkiste, S., Chang, M., Greff, K., Schmidhuber, J.: Relational neural expectation maximization: Unsupervised discovery of objects and their interactions. arXiv preprint arXiv:1802.10353 (2018)
- [64] Wu, R., Zhuang, Y., Xu, K., Zhang, H., Chen, B.: Pq-net: A generative part seq2seq network for 3d shapes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2020)
- [65] Xia, W., Zhang, Y., Yang, Y., Xue, J.H., Zhou, B., Yang, M.H.: Gan inversion: A survey. arXiv preprint arXiv:2101.05278 (2021)
- [66] Yao, C.H., Hung, W.C., Jampani, V., Yang, M.H.: Discovering 3d parts from image collections. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 12981–12990 (2021)
- [67] Yu, H.X., Guibas, L.J., Wu, J.: Unsupervised discovery of object radiance fields. arXiv preprint arXiv:2107.07905 (2021)
- [68] Zablotskaia, P., Dominici, E.A., Sigal, L., Lehmman, A.M.: Unsupervised video decomposition using spatio-temporal iterative inference. arXiv preprint arXiv:2006.14727 (2020)
- [69] Zhang, N., Donahue, J., Girshick, R., Darrell, T.: Part-based r-cnns for fine-grained category detection. In: European conference on computer vision. pp. 834–849. Springer (2014)
- [70] Zhao, H., Jiang, L., Jia, J., Torr, P.H., Koltun, V.: Point transformer. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 16259–16268 (2021)
- [71] Zhou, X., Girdhar, R., Joulin, A., Krähenbühl, P., Misra, I.: Detecting twenty-thousand classes using image-level supervision. arXiv preprint arXiv:2201.02605 (2022)

- [72] Zhu, J.Y., Krähenbühl, P., Shechtman, E., Efros, A.A.: Generative visual manipulation on the natural image manifold. In: European conference on computer vision. pp. 597–613. Springer (2016)
- [73] Zoran, D., Kabra, R., Lerchner, A., Rezende, D.J.: Parts: Unsupervised segmentation with slots, attention and independence maximization. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10439–10447 (2021)
- [74] Zuffi, S., Kanazawa, A., Berger-Wolf, T., Black, M.J.: Three-d safari: Learning to estimate zebra pose, shape, and texture from images" in the wild". In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5359–5368 (2019)

Appendix

The structure of this appendix is as follows: In Section 6 we cover the details on the datasets. In Section 7 we specify further implementation details. In Section 8 we provide additional qualitative and quantitative results for the experiments in Section 4 of our main paper.

6 Datasets

6.1 Point Cloud

For all the tasks we subsample the input point clouds to a standard size of 2048 points.

Generic primitive part dataset. We use the primitive dataset of [60] for learning the primitive distribution in Section 4.1.1. The dataset consists of 200K primitive instances sampled from the primitive templates that are visualized in Figure 7. Instances are sampled from the templates by changing their sizes and placing them in random locations, similar to [60].

Category-specific primitives from PartNet dataset. We separate out individual segments from the level-3 instance segmentation masks from the official train-split of Chair category of the PartNet dataset. We additionally do rotation and translation based augmentations on individual segments. We visualize some of these primitive examples in Figure 8. We use this dataset for learning the primitive distribution in Section 4.1.2. The dataset in total consists of about 100K primitive examples.

Synthetic whole shape dataset. This dataset was generated by Shape2Prog[60]. The segmentation labels from this dataset are used by them as a supervision signal for later generalizing to PartNet shapes. The dataset consists of about 120K synthetically generated Chairs and Tables, we visualize some of these synthetically generated tables and chairs in Figure 9. Note that no other baseline or GFS-Nets has access to this dataset.

PartNet dataset. We use the official level-3 train-test split of PartNet[49]. We use the train split of Chair category as our training set, we consider test split of Table category in PartNet our test categories. We use this as the train-test split in Section 4.1 of the main paper. We further compare our model on the other PartNet categories in supplementary section 8. We set the value of number of slots K as 16 for this dataset.

6.2 RGB

CLEVR dataset. We use the train-test split of CLEVR dataset opensourced by [67]. Their dataset is built on top of the original CLEVR [30] dataset, where they add additional multi-view rendering. Their train/test dataset consists of 5,6,7 objects from random geometric primitives(cylinder, cube and sphere) and random colors(red, blue, purple, gray, cyan, yellow, green, brown) placed in random locations in the scene. The multi-view setting consists of randomly sampling from 360 degree azimuth angles with a fixed elevation angle. Our train data consists of 1000 single object scenes and 1000 multi-object scenes. The only difference between their train-test data are the novel locations at which the objects would be placed in the scene. We set the value of number of slots K as 8. We visualize some of the primitive examples (multi-view single-object scenes) of the CLEVR dataset in Figure 10.

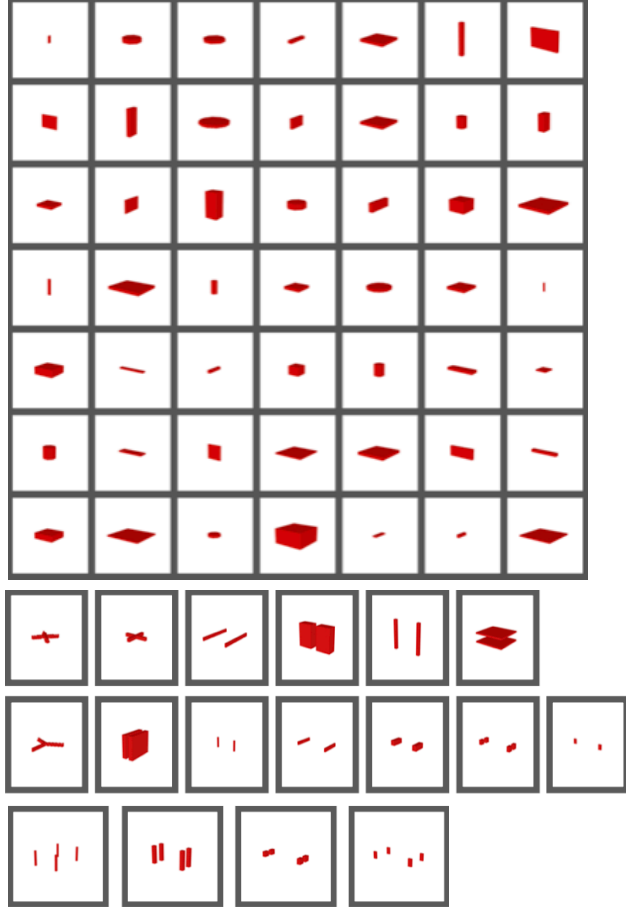


Figure 7: We visualize the generic primitive templates of [60], as you can see they mainly consists of Cubes, Cuboids and Discs.

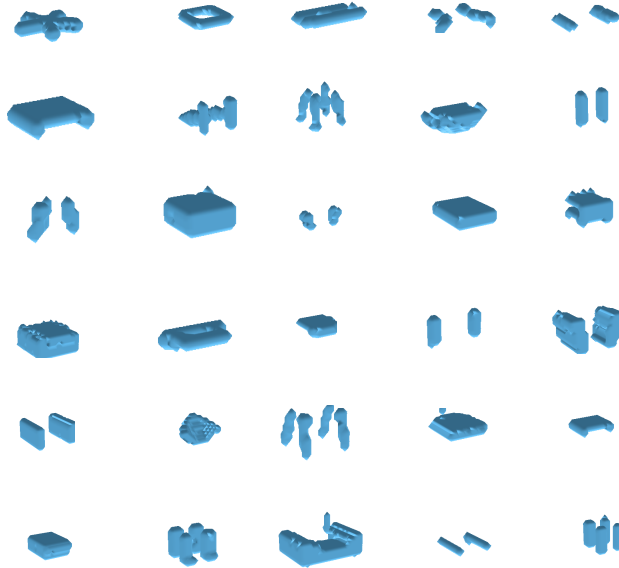


Figure 8: We visualize the category-specific primitives extracted from level-3 Chair category of PartNet.



Figure 9: We visualize the synthetic whole shape dataset of [60]. Shape2Prog supervises their model using the annotations from this dataset.

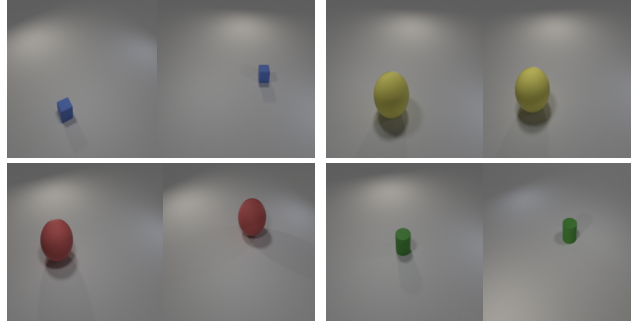


Figure 10: We visualize some examples of the primitives used in the CLEVR dataset. Primitives are represented as multi-view single object RGB scenes as shown in the figure.

RoomDiverse++ dataset. We build RoomDiverse++ on top of the RoomDiverse dataset of [67]. In which instead of using the solid CLEVR colors as textures of meshes we keep the original textures of the ShapeNet meshes, additionally along with Chairs we also add other ShapeNet categories such as Benches, Cabinets, Beds, Tables and Sofas in the dataset. In total we use 20 object meshes and 3 background textures to generate the dataset. We use the same multi-view rendering settings as CLEVR. Our train data consists of 1000 single object scenes and 2000 multi-object scenes. We use 6,7,8 objects placed at random locations for generating our multi-object train dataset. We use 500 of 7,8,9 object scenes as our test data. We set the value of number of slots K as 10. We visualize some of the primitive examples in the dataset in Figure 11.

7 Implementation details

7.1 Proposed Model

Code, training details and computational complexity. Proposed model is implemented in Python/PyTorch. We use a batch size of 8 for point cloud input and a batch size of 2 on RGB input. We set our learning rate as 10^{-4} . We use the Adam optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$. Our primitive learning model takes 24 hours (approximately 200k iterations) to converge. Similarly our



Figure 11: We visualize some examples of the primitives used in the Room Diverse++ dataset. Primitives are represented as multi-view single object RGB scenes as shown in the figure.

primitive composition learning model takes about 12hr (approximately 100k iterations) to converge. Our slow inference per example takes about 1 min (500 iterations). A forward pass through the proposed model takes about 0.15 secs. We use a single V100 GPU for training and inference.

Inputs. For RGB input we randomly select 2 views which include the RGB and ground-truth egomotion. We use a single RGB image for testing. For RGB, our input is a $64 \times 64 \times 3$ tensor, for point cloud our input is a 2048×3

Point Cloud Encoder. We adopt the point transformer[70] architecture as our encoder. Point transformer encoder is essentially layers of self attention blocks. Specifically a self attention block includes sampling of query points and updating them using their N most neighbouring points as key/value vectors. In the architecture we specifically apply 5 layers of self attention which look as follows: 2048-16-64, 2048-16-64, 512-16-64, 512-16-64, 128-16-64, 128-16-64. We use the notation of S - N - C , where S is the number of subsampled query points from the point cloud, N is the number of neighbouring points and C is the feature dimension. We thus get an output feature map of size 128×64 .

RGB Encoder. Our RGB encoder is similar to that of [67], which is basically a U-net. We concatenate the RGB values with the pixel coordinates that are normalized in the range of -1 to 1. We then pass the concatenated 2D image through a convolutional feature extractor. The architecture of the CNN is as follows: 3-2-64, 3-2-64, 3-2-64, 3-1-64, 3-1-64, 3-1-64. We use the notation of k - s - c , where k is the kernel size, s is the number of strides and c is the feature channel. Our final output size is $64 \times 64 \times 64$.

Point Cloud Decoder. We obtain point occupancies by querying the slot feature vector $slot_k$ at discrete locations (x, y, z) specifically $o_{x,y,z} = \text{Dec}(slot_k, (x, y, z))$. The architecture of Dec is similar to that of [40]. Given $slot_k$, which is one of the slot feature vector. We encode the coordinate (x, y, z) into a 64-D feature vector using a linear layer. We denote this vector as z . The inputs $slot_k$ and z are then processed as follows:

$$out_k = RB_i(RB_{i-1}(\dots RB_1(z + FC_1(slot_k)) \dots] + FC_{i-1}(slot_k)) + FC_i(slot_k)). \quad (3)$$

We set $i = 3$. FC_i is a linear layer that outputs a 64 dimensional vector. RN_i is a 2 layer ResNet MLP block [27]. The architecture of ResNet block is: ReLU, 64-64, ReLU, 64-64. Here, $i - o$ represents a linear layer, where i and o are the input and output dimension. Finally out_k is then passed through a ReLU activation function followed by a linear layer to generate a single value for occupancy.

RGB Decoder. Our decoder follows the architecture of [67]. The decoder is a 8 layer MLP, specifically ReLU, 64-64, ReLU, 64-64 64-4. The MLP takes in as input $slot_k$ feature and the location x, y, z and predicts the RGB and density values at that location.

Slot Extractor. Following is the pseudo-code of Slot Attention, which we use for learning primitive compositions.

Algorithm 1 Slot Attention module. Input is N vectors of dimension D_{inputs} which is mapped to a set of K slots of dimension D_{slots} . We sample the learnable multivariate gaussian to get the initial slot features, here the learnable parameters are: $\mu \in \mathbb{R}^{D_{\text{slots}}}$ and $\sigma \in \mathbb{R}^{D_{\text{slots}}}$. For our experiments we set $T = 3$.

```

1: Input:  $\text{inputs} \in \mathbb{R}^{N \times D_{\text{inputs}}}$ ,  $\text{slots} \sim \mathcal{N}(\mu, \text{diag}(\sigma)) \in \mathbb{R}^{K \times D_{\text{slots}}}$ 
2: Layer params:  $k, q, v$ : linear projections; GRU; MLP; LayerNorm(x3)
3:  $\text{inputs} = \text{LayerNorm}(\text{inputs})$ 
4: for  $t = 0 \dots T$ 
5:    $\text{slots\_prev} = \text{slots}$ 
6:    $\text{slots} = \text{LayerNorm}(\text{slots})$ 
7:    $\text{attn} = \text{Softmax}(\frac{1}{\sqrt{D}} k(\text{inputs}) \cdot q(\text{slots})^T, \text{axis}='slots')$ 
8:    $\text{updates} = \text{Avg}(\text{weights}=\text{attn} + \epsilon, \text{values}=v(\text{inputs}))$ 
9:    $\text{slots} = \text{GRU}(\text{state}=\text{slots\_prev}, \text{inputs}=\text{updates})$ 
10:   $\text{slots} += \text{MLP}(\text{LayerNorm}(\text{slots}))$ 
11: return  $\text{slots}$ 

```

For primitive learning stage on PointCloud we set the value of number of slots K as 1. Additionally we replace the $\text{Softmax}()$ activation in step 7 with a $\text{Sigmoid}()$ activation. For primitive learning stage on RGB scenes we set the value of number of slots K as 2. This is due to RGB images having a background component in them.

7.2 Baselines

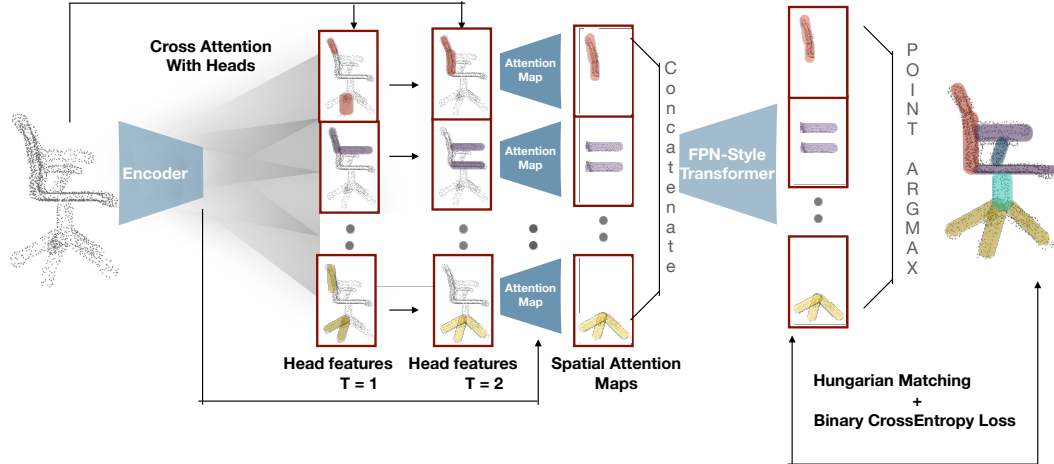


Figure 12: 3D-DETR model architecture. Our baseline model is very similar to 2D-DETR[5], where we replace their base 2D encoder and 2D FPN-style CNN with a 3D PointTransformer[70]. Given point clouds as input our encoder backbone featurizes the points into N feature vectors, we then do iterative self and cross attention using K learnable query heads, following the implementation of Carion *et al.*[5]. We then compute attention of the updated queries wrt to the encoded feature vectors. We then concatenate these attentions along the batch dimension and pass them through a FPN-style transformer that increases their resolution and outputs each query mask logits. We then do hungarian matching and binary cross entropy loss

3D-DETR. 3D-DETR is a SOTA segmentation architecture, we build this architecture on top of 2D-DETR[5], where we replace their 2D encoders with a 3D PointTransformer[70]. For fair comparison we make sure the number of parameters in 3D-DETR is the same as GFS-Nets.

The base encoder is a 5 layer architecture of 1024-16-64, 1024-16-64, 1024-16-64, 1024-16-64, 512-16-64, following the notation of S, N, C where S is the number of sampled query points, N is the number of selected neighbouring points and C is the feature dimension. We have 16 learnable

queries with six encoder-decoder layers of transformer attention, each layer consists of 8 heads. We then compute multi-head cross-attention between each query vector and the encoded 3d points. This gives us a map of size $M \times 512 \times 16$, where M is the number of heads in the multi-head attention.

Each attention map is individually upsampled through a PointTransformer decoder of the architecture 512-16-64, 1024-16-64, 1024-16-64, which gives us the final instance segmentation mask per query, similar to DETR[5]. We then use Hungarian matching to match the predicted masks against the ground truth masks and then apply binary cross entropy loss for each match. We aggregate the losses from each query and backpropagate. Figure 12 visualizes the architecture of 3D-DETR. Note that we do not follow the 2 stage training of [5], rather we train their model end-to-end for instance segmentation, similar to our model. We have found this trick to save compute and still not harm the end results.

Learning2Group. [43] We use their open-sourced architecture and code for comparison with our model. We train their model using our datasets from scratch.

Shape2Prog [60] We use their open-sourced architecture, code and pretrained checkpoints for comparison with our model. We change the value of number of blocks similar to the number of slots in our model for each dataset.

PQ-Nets. [64] We use their open-sourced architecture and code for comparison with our model. We train their model using our datasets from scratch. We change the value of number of slots in their model based on the maximum number of parts in the dataset.

Struct Implicit. [21] We use their open-sourced architecture and code for comparison with our model. We train their model using our datasets from scratch. We change the value of number of slots in their model based on the maximum number of parts in the dataset.

uORF. [67] We use their open-sourced architecture, code and pretrained checkpoints for comparison with our model on the CLEVR dataset. For RoomDiverse++ dataset, we train their model on our dataset using the hyperparameters they use for training their model on RoomDiverse dataset. We change the value of number of slots in their model similar to our model for each dataset.

Slot Attention. [42] We use their open-sourced architecture and code for comparison with our model on the CLEVR and RoomDiverse++ dataset. We train their model using our datasets from scratch. We change the value of number of slots in their model similar to our model for each dataset.

8 More Results

Segmenting 3D pointclouds without segmentation supervision. In this section we show additional qualitative and quantitative results for Section 4.1.1 of the main paper. We follow the experimental setup of Section 4.1.1, where we use the generic part dataset (visualized in Fig 7) for primitive learning. In Table 4, we further extend Table 1 of the main paper to include different levels in Chair and Table category. Here ($X/Y/Z$) refers to (level 1/level 2/level 3) scores respectively. We pick the best performing level in fast inference and specifically report it’s slow inference results. We further qualitatively compare our model against Shape2Prog (best performing baseline) in Figure 13

Method	in-dist (Chairs)		out-of-dist (Tables)	
	Fast Infer.	Slow Infer.	Fast Infer.	Slow Infer.
Shape2Prog [60]	0.28 /0.21/0.26	0.53	0.21/0.23/ 0.23	0.40
PQ-Nets [64]	0.20 /0.18/0.16	0.31	0.17 /0.14/0.16	0.21
StructImplicit [21]	0.18/0.22/ 0.25	0.23	0.14/0.15/ 0.17	0.17
GFS-Nets	0.51/0.48/ 0.57	0.62	0.51/0.55/ 0.60	0.69

Table 4: **ARI Segmentation accuracy (higher the better)** in the test set of Chair (in-distribution) and Table category of PartNet(out-of-distribution). GFS-Nets significantly outperform all of the baselines.

RGB image segmentation and view prediction without supervision In Figure 15, we show more qualitative results for Section 4.3 of the main paper, where we qualitatively compare our model for segmentation against uORF[67] and Slot Attention [42]. We see that slow inference prevents our model from missing to detect occluded objects in the scene, where as the baselines have a lot of False Positives i.e missed detections.

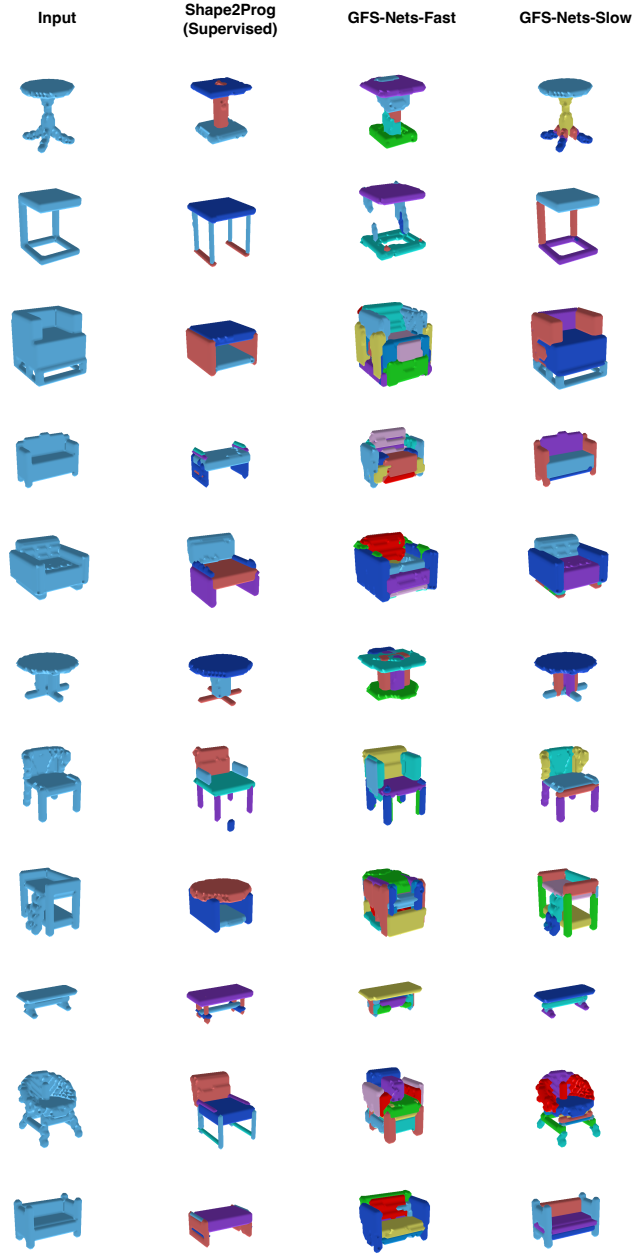


Figure 13: Additional shape decomposition results using generic primitives. (Section 4.1.1)

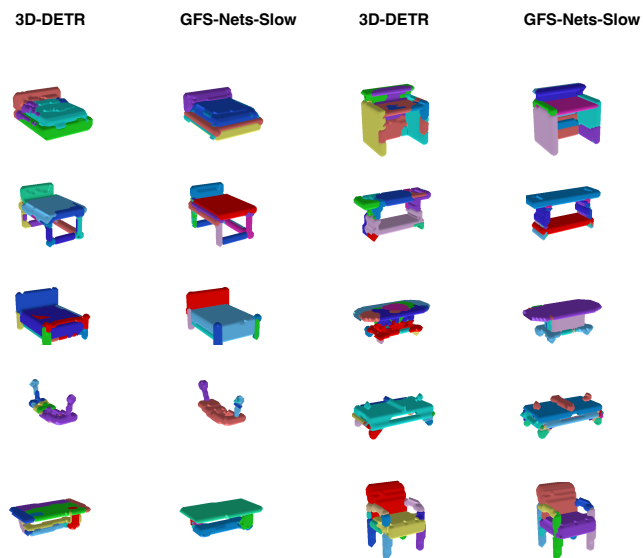


Figure 14: Additional shape decomposition results using category-specific primitives. (Section 4.1.2)



Figure 15: Additional RGB segmentation results on RoomDiverse++. (Section 4.3)