

Learning to Infer Shape Programs Using Latent Execution Self Training

Homer Walke R. Kenny Jones Daniel Ritchie
Brown University

{homer_walke, russell_jones, daniel_ritchie}@brown.edu

Abstract

Inferring programs which generate 2D and 3D shapes is important for reverse engineering, enabling shape editing, and more. Supervised learning is hard to apply to this problem, as paired (program, shape) data rarely exists. Recent approaches use supervised pre-training with randomly-generated programs and then refine using self-supervised learning. But self-supervised learning either requires that the program execution process be differentiable or relies on reinforcement learning, which is unstable and slow to converge. In this paper, we present a new approach for learning to infer shape programs, which we call latent execution self training (LEST). As with recent prior work, LEST starts by training on randomly-generated (program, shape) pairs. As its name implies, it is based on the idea of self-training: running a model on unlabeled input shapes, treating the predicted programs as ground truth latent labels, and training again. Self-training is known to be susceptible to local minima. LEST circumvents this problem by leveraging the fact that predicted latent programs are executable: for a given shape $\mathbf{x}^* \in S^*$ and its predicted program $\mathbf{z} \in P$, we execute \mathbf{z} to obtain a shape $\mathbf{x} \in S$ and train on $(\mathbf{z} \in P, \mathbf{x} \in S)$ pairs, rather than $(\mathbf{z} \in P, \mathbf{x}^* \in S^*)$ pairs. Experiments show that the distribution of executed shapes S converges toward the distribution of real shapes S^* . We establish connections between LEST and algorithms for learning generative models, including variational Bayes, wake sleep, and expectation maximization. For constructive solid geometry and assembly-based modeling, LEST's inferred programs converge to high reconstruction accuracy significantly faster than those of reinforcement learning.

1. Introduction

Having access to a procedure which generates a visual datum reveals its underlying structure, facilitating high-level manipulation and editing by a person or autonomous agent. Thus, inferring such program from visual data is an important problem. Programs are particularly relevant for

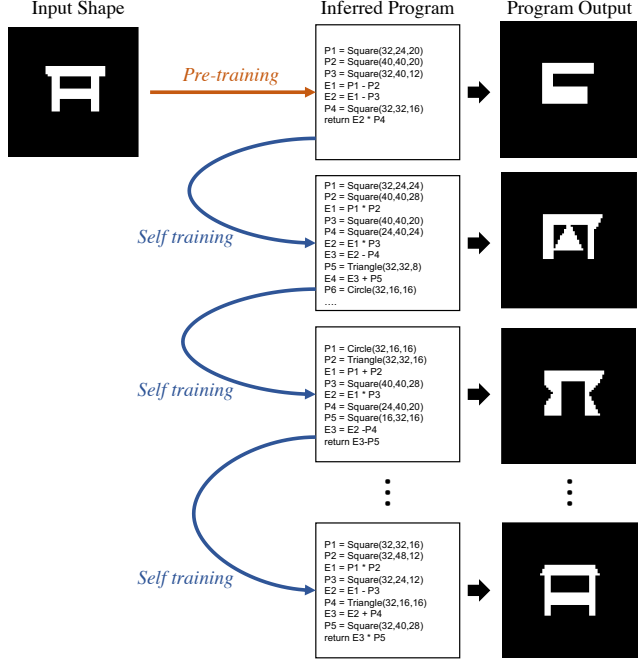


Figure 1. One can learn a model to infer shape programs from input shapes by training on randomly-generated programs. Our latent execution self training (LEST) approach extends this idea, re-training the inference model on its own inferred programs (and the shapes they produce) to further improve inference performance and reach convergence faster than reinforcement learning.

manipulating shapes, i.e. point sets defined by an indicator function. In \mathbb{R}^2 , inferring shape programs has applications in the design of diagrams, icons, and other 2D graphics. In \mathbb{R}^3 , it has applications in reverse engineering of CAD models, procedural modeling for 3D games, and 3D structure understanding for autonomous agents.

We formally define shape program inference as obtaining a latent program \mathbf{z} which generates a given observed shape \mathbf{x} , which we solve by constructing (e.g. learning) and sampling from a distribution $p(\mathbf{z}|\mathbf{x})$. This is a challenging problem: it is a structured prediction problem whose output is high-dimensional and features both discrete and continuous components (i.e. program control flow vs. program

parameters). It is also a one-to-many problem, as multiple latent programs can generate the same observed shape. Nevertheless, advances in deep neural networks have made it possible to learn models for $p(\mathbf{z}|\mathbf{x})$, provided that one has access to paired (\mathbf{x}, \mathbf{z}) data (i.e. a dataset of shapes and the programs which generate them). Unfortunately, while shape data is increasingly available in large quantities [1], these shapes do not typically come with their generating program. In fact, only some such shapes (those constructed in CAD modeling software) even have a ground-truth underlying program to begin with.

To circumvent the data availability problem, researchers have typically *synthesized* paired data by generating random programs and pairing them with the shapes they output, enabling supervised learning of $p(\mathbf{z}|\mathbf{x})$ models. However, the distribution of shapes S produced by these random programs P is typically quite different from that of the distribution of “real” shapes of interest S^* , which often causes the learned $p(\mathbf{z}|\mathbf{x})$ not to generalize well to real shapes $\mathbf{x}^* \in S^*$. To date, researchers have attacked this problem by fine-tuning $p(\mathbf{z}|\mathbf{x})$ with some form of self-supervised learning: optimizing the model’s parameters such that the outputs of its inferred programs are geometrically similar to their input shapes. Training such a model end-to-end requires the program executor be differentiable, though, which is rarely the case: shape programs typically make multiple discrete structural decisions. If one knows the functional form of the executor (i.e. has access to its source code), it may be possible to implement a differentiable relaxation of this function [12]. If not, one can try to learn a differentiable approximation of the executor’s behavior, though this approximation introduces errors [20]. The most generally-applicable fine-tuning method is reinforcement learning [19, 5], which treats the executor as a non-differentiable black box. However, RL is often unstable and slow to converge.

In this paper, we present a new approach to training shape program inference models $p(\mathbf{z}|\mathbf{x})$ that works with any black-box program executor yet converges better and faster than RL. Like prior work, our approach initializes $p(\mathbf{z}|\mathbf{x})$ by pre-training on randomly-generated data. At this point it diverges from prior work: instead of self-supervised learning, it follows a self-training procedure [18, 23]. In self-training, one uses a pre-trained model $p(\mathbf{z}|\mathbf{x})$ to infer latent \mathbf{z} ’s for input unlabeled \mathbf{x} ’s; these \mathbf{z} ’s then become “pseudo-labels” which are treated as ground truth for another round of supervised training. This approach has been shown to yield practical performance gains in a variety of domains but requires careful tuning, as training on too many incorrect pseudo-labels can cause learning to degrade. Our approach circumvents this problem by taking advantage of the fact that in shape program inference, an inferred \mathbf{z} is actually a program that can be executed to produce a shape \mathbf{x} that is consistent with the program; that is, \mathbf{z} is guaranteed

to be the “correct label” for \mathbf{x} . We call this procedure latent execution self training, or LEST.

We note that the idea of iteratively re-training a latent variable recognition model $p(\mathbf{z}|\mathbf{x})$ is closely related to methods for learning latent variable generative models $p(\mathbf{z}, \mathbf{x})$ which jointly train such a recognition model. These approaches include variational Bayes [14], wake sleep [10], and certain variants of expectation maximization [4]. Our approach takes advantage of the unique setting of shape program inference, in which the likelihood function $p(\mathbf{x}|\mathbf{z})$ is a known deterministic program executor, avoiding the need to learn a complex generative prior over programs $p(\mathbf{z})$.

In our experiments, we use LEST to learn to infer programs in two different shape domains: constructive solid geometry (CSG) and assembly-based modeling with ShapeAssembly, a domain-specific language for specifying the structure of manufactured 3D objects [11]. We show that pseudo-label programs inferred by the model $p(\mathbf{z}|\mathbf{x})$ produce shapes that gradually converge toward the target distribution (Figure 1). We also show that LEST converges to better shape reconstruction performance than RL in significantly less computation time.

In summary, our contributions are:

1. The latent execution self training (LEST) framework for learning to infer shape programs.
2. Implementation and evaluation of LEST for 2D CSG and ShapeAssembly program inference.

2. Background & Related Work

Shape program inference is a type of *visual program induction* problem [2]. Here, we will discuss prior work that has attacked the shape program inference problem, organized by the learning methodology used to construct $p(\mathbf{z}|\mathbf{x})$. While we focus on learning methodology, we note that several of these works also use classic search techniques (e.g. beam search), in which $p(\mathbf{z}|\mathbf{x})$ is used to guide the search.

Supervised Learning The most straightforward method for learning $p(\mathbf{z}|\mathbf{x})$ is supervised training with (\mathbf{x}, \mathbf{z}) pairs. Ellis et al. learn to infer 2D diagram programs using supervised training on randomly-generated synthetic programs [6]. Liu et al. learn to infer programs to describe simple scene images using the same approach [16]. Other recent works rely on supervised pre-training and then fine-tune by other means [19, 20, 5]. These all use randomly-generated programs, due to the difficulty of obtaining ground-truth programs for real-world shapes at scale. Shapes output by random programs are typically statistically different from real shapes, and this distribution shift leads to a performance hit when $p(\mathbf{z}|\mathbf{x})$ is applied on real shapes. LEST also starts with pre-training on random synthetic data, but it iteratively corrects for this distribution shift by self-training on its own predictions.

Method	Models Needed	Differentiable $p(\mathbf{x} \mathbf{z})$?	Notes
Supervised Learning	$p(\mathbf{z} \mathbf{x})$	No	Requires paired (\mathbf{x}, \mathbf{z}) data
Reinforcement Learning	$p(\mathbf{z} \mathbf{x})$	No	Unstable, slow convergence
Self-supervised Learning	$p(\mathbf{z} \mathbf{x})$	Yes	
Variational Bayes	$p(\mathbf{z} \mathbf{x}), p(\mathbf{z})$	Yes	
Wake-sleep, EM	$p(\mathbf{z} \mathbf{x}), p(\mathbf{z})$	No	
LEST	$p(\mathbf{z} \mathbf{x})$	No	

Table 1. Comparison of different methods for learning to infer shape programs \mathbf{z} from shapes \mathbf{x} , in terms of the models that must be trained, their differentiability requirements, and other considerations. In all cases, it is assumed that the recognition model $p(\mathbf{z}|\mathbf{x})$ is differentiable.

Reinforcement Learning The most generally-applicable method for fine-tuning a pre-trained $p(\mathbf{z}|\mathbf{x})$ is reinforcement learning: treating $p(\mathbf{z}|\mathbf{x})$ as a policy network and using policy gradient methods [21]. The geometric similarity of the inferred program’s output to its input is the reward function; the program executor $p(\mathbf{x}|\mathbf{z})$ can be treated as a (non-differentiable) black box. CSG-Net uses RL for fine-tuning [19], as does other recent work on inferring CSG programs from raw input geometry [5]. CSG-Net has been extended with a suite of improvement that allow it to converge without supervised pre-training [22]. The main problem with policy gradient RL is its instability due to high variance gradients; this typically necessitates a low learning rate and thus slow convergence. Like RL, LEST can treat the program executor as a black box, but (as we show experimentally) it is more stable and converges much faster.

Self-Supervised Learning If the functional form of the program executor $p(\mathbf{x}|\mathbf{z})$ is known and differentiable, then self-supervised learning is possible. Self-supervision works similarly to RL, except the gradient of the reward with respect to the parameters of $p(\mathbf{z}|\mathbf{x})$ can be computed directly, making policy gradient unnecessary. Shape programs are typically not fully differentiable, as they often involve discrete choices (e.g. which type of primitives to create). UC-SGNet uses a differentiable relaxation of constructive solid geometry to circumvent this issue [12]. Other work takes the more extreme approach of training a differentiable network to approximate the behavior of the program executor [20], which introduces errors. LEST does not require the program executor to be differentiable, yet it performs better than other approaches (e.g. RL) that share this property.

Generative Model Learning Shape program inference has also been explored in the context of learning a generative model $p(\mathbf{x}, \mathbf{z})$ of programs and the shapes they produce. The most popular approach for training such models at present is variational Bayes, in particular the variational autoencoder [14]. This method simultaneously trains a generative model $p(\mathbf{x}, \mathbf{z})$ and a recognition model $p(\mathbf{z}|\mathbf{x})$ by optimizing a lower bound on the marginal likelihood $p(\mathbf{x})$.

When the \mathbf{z} ’s are shape programs, the program executor is $p(\mathbf{x}|\mathbf{z})$, so training the generative model reduces to learning a prior over programs $p(\mathbf{z})$. To jointly train such models with gradient descent requires, as with self-supervised learning, that the executor $p(\mathbf{x}|\mathbf{z})$ be differentiable. When this is not possible, the wake sleep algorithm is a viable alternative [10]. This approach alternates training steps of the generative models and recognition model, training one using samples produced by the other. Recent work has used wake sleep for visual program induction [9, 7]. If one trains $p(\mathbf{x}, \mathbf{z})$ and $p(\mathbf{z}|\mathbf{x})$ to convergence before switching to training the other, instead of alternating training updates, this is equivalent to expectation maximization (under a view of EM as two alternating maximization procedures [17]). LEST also iteratively trains a $p(\mathbf{z}|\mathbf{x})$ to convergence, but it does not require also training a generative model $p(\mathbf{x}, \mathbf{z})$: it takes advantage of the fact that $p(\mathbf{x}|\mathbf{z})$ is a known deterministic function the (program executor) to avoid this.

Table 1 shows a summary of these different methods and their relevant properties.

3. Latent Execution Self Training (LEST)

LEST assumes four inputs, $(S^*, G, R_G, p(\mathbf{x}|\mathbf{z}))$. S^* is a set of shapes drawn from the distribution of shapes for which we would like to infer programs. G is a formal specification of the grammar of the language in which we would like to infer programs, R_G is a procedure for randomly sampling programs from this grammar, and $p(\mathbf{x}|\mathbf{z})$ executes programs in the language. The design of R_G is domain-dependent; Section 4 describes example random program generators for two shape domains.

Figure 2 shows a schematic overview of our latent execution self training approach, and Algorithm 1 provides pseudocode. The process begins by using R_G to generate a large set of random programs P . Each of these programs $\mathbf{z} \in P$ is executed to produce its corresponding output shape \mathbf{x} ; we use S to denote the set of such generated shapes. Pairs of (\mathbf{x}, \mathbf{z}) are then used to train an initial inference network $p(\mathbf{z}|\mathbf{x})$ using supervised learning. The architecture of $p(\mathbf{z}|\mathbf{x})$ is domain-dependent, though it typi-

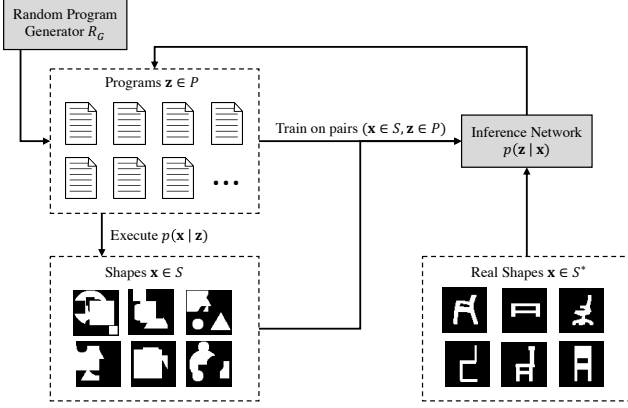


Figure 2. Visual illustration of latent execution self training. A random program generator R_G is used to create initial random programs P . These programs are executed to produce shapes S ; these shapes and their corresponding programs are used to train an inference network $p(z|x)$. The trained network is evaluated on shapes from a distribution of “real” shapes of interest S^* , and the resulting programs (and their output shapes) are used to begin another iteration of self-training.

Algorithm 1 Pseudocode for LEST

Input: $(S^*, G, R_G, p(z|x))$

Output: $p(z|x)$

```

 $P \leftarrow \{z \sim R_G\}$ 
 $S \leftarrow \{x \sim p(x|z) \mid z \in P\}$ 
 $\text{sim} \leftarrow 0$ 
repeat
   $\text{prev\_sim} \leftarrow \text{sim}$ 
   $p(z|x) \leftarrow \text{train}(\{x \in S, z \in P\})$ 
   $P \leftarrow \{z \sim p(z|x) \mid x \in S^*\}$ 
   $S \leftarrow \{x \sim p(x|z) \mid z \in P\}$ 
   $\text{sim} \leftarrow \text{similarity}(S, S^*)$ 
until  $\text{sim} < \text{prev\_sim}$ 

```

cally includes a recurrent language model for generating sequences of output program tokens. Once trained, $p(z|x)$ is invoked on the shapes x^* in the input set of “real” shapes S^* to produce a new set of programs P . At this point, the process repeats, with the new P serving the same role as their predecessors for the next round of training. This iteration terminates once the average similarity of program output shapes to their inputs stops increasing.

4. Experiments

We evaluate LEST in two shape program domains:

- **2D Constructive Solid Geometry:** Generating 2D shapes via the union, intersection, and difference of parametric primitive shapes.
- **ShapeAssembly:** Generating 3D shapes by creating cuboids and attaching them to one another [11].

In each domain, we compare LEST to multiple alternative approaches:

- **Supervised Pretraining (SP):** Using only supervised pretraining on randomly-generated synthetic data.
- **Reinforcement Learning (RL):** Training a randomly-initialized model with REINFORCE [21].
- **Supervised Pretraining + RL Fine-tuning (SP+RL):** Pretraining a model with randomly-generated data and then fine-tuning with REINFORCE.
- **Self-Training (ST):** Self-training without latent execution, i.e. training on $(x^* \in S^*, z \in P)$ pairs instead of $(x \in S, z \in P)$ pairs.

To evaluate the performance of each of these shape program inference methods, we consider the following metrics:

- **Reconstruction Accuracy:** How closely the output of a shape program matches the input shape from which it was inferred on a held out set of test shapes. The specific metric varies by domain.
- **Distributional Similarity:** How similar is the distribution of generated shapes S to the distribution of real shapes S^* ? The more similar, the more likely it is that the inferred programs P which generate S reflect meaningful shape structures. The specific metric varies by domain.

4.1. 2D Constructive Solid Geometry

We first evaluate LEST on constructive solid geometry (CSG) programs. In CSG, shapes are created by declaring parametric primitives (e.g. circles, boxes) and combining them with boolean operations (union, intersection, difference). CSG inference is non-trivial: as CSG uses non-additive operations (intersection, difference), inferring a CSG program does not just reduce to primitive detection. It is a good starting point for our evaluation because it is “well-conditioned”: small changes in a CSG program lead to small changes in its output shape.

Real Shapes S^* We use the CAD dataset from CSGNet [19]. The dataset consists of the front and side views of chairs, desks, and lamps downloaded from the Trimble 3D warehouse. Each image is rendered as a 64×64 binary mask. We split the dataset into $10K$ shapes for training, $3K$ shapes for validation, and $3K$ shapes for testing.

Random Program Generator R_G We use CSGNet’s random data: random binary CSG trees in which the leaves are randomly-sized squares, circles, and triangles, and the internal nodes are boolean operators. Primitive and operator types are uniformly distributed. Primitive locations are discretized on a 64×64 grid.

Network architecture for $p(z|x)$ We also use the same inference network architecture as CSGNet. This architecture uses a CNN to encode 64×64 binary mask shape images. The output of this encoder is used to initialize a GRU recurrent decoder for predicting program tokens.

Experiment details We initialize the network with the pretrained model from CSGNet, which was trained on the random data for 400 epochs (472 hours). We then use this pretrained network as the starting point for LEST, ST, and SP+RL. LEST and ST were run for 40-70 self-training rounds (until convergence). Each round was stopped when the model converged on validation set reconstruction performance. When inferring programs, we use beam search with beam width 10. We use the Adam optimizer with learning rate 0.001 [13]. For RL, we use the REINFORCE setup provided by CSGNet with chamfer distance reward. We use stochastic gradient descent with learning rate 0.01 and momentum 0.9. All timings were collected on a machine with a GeForce RTX 2080 Ti GPU and an Intel i7-7800X CPU.

Results Figure 3 top plots the distributional similarity between S and S^* over the course of LEST self-training, and Figure 3 bottom shows how the output of the inferred program for a few shapes $x^* \in S^*$ evolves over time. To measure distributional similarity, we use bidirectional average nearest neighbor distance between shapes in S and S^* (the average distance over all $x \in S$ to its nearest $x^* \in S^*$ plus the average distance over all $x^* \in S^*$ to its nearest $x \in S$). Supervised pretraining brings S considerably closer to S^* than the initial random shapes from R_G , and each successive round of LEST self-training improves upon this. The individual shape trajectories in Figure 3 bottom show how some of these improvements are achieved, e.g. replacing circular primitives with angular ones.

Figure 4 plots the reconstruction accuracy (Chamfer Distance) of LEST vs. other methods as a function of training time. The performance of supervised pretraining (SP), as well as RL, are reported as dotted reference lines; the time required for supervised pretraining is factored out for all other methods. Pure RL performs poorly; worse than supervised retraining alone. Starting from the SP baseline, all the methods improve significantly, though the self-training approaches converge much more quickly than RL and to a better optimum. Of note, standard self-training (ST) works well in this domain; this is likely due to the well-conditioned nature of CSG programs. Even more interestingly, ST eventually overtakes and outperforms LEST. Our investigations indicate that this is due to the domain gap between the S^* shapes and the S shapes produced by the CSG executor: when the CNN weights for ST and LEST are frozen after pretraining, both converge similarly (see supplemental for details). The best performance is achieved

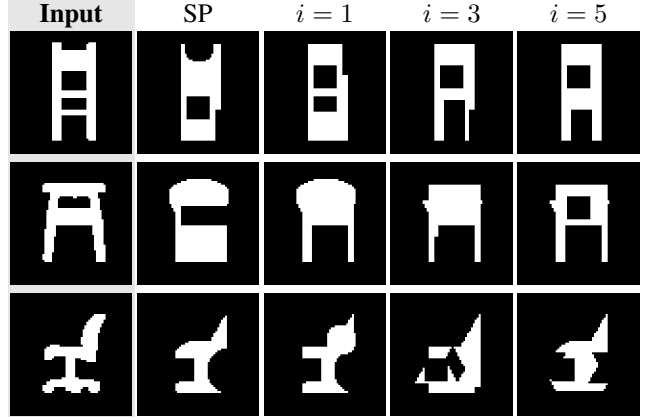
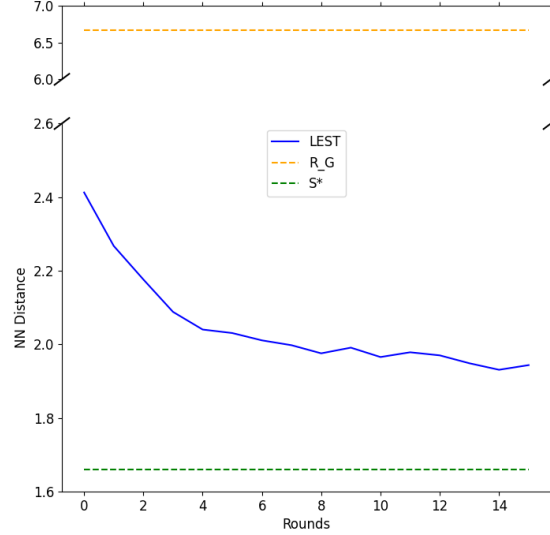


Figure 3. (Top) The distributional similarity of S to S^* with increasing LEST self training rounds, for 2D CSG; lower is better. (Bottom) How the output of the inferred program for specific shapes changes over LEST rounds.

by combining LEST and ST (LEST + ST), getting the benefits of LEST’s precise correspondence between observed shapes and latent programs *and* the benefits of ST’s ability to see shapes from the real distribution S^* . We implement this mode when training by randomly sampling for each program label if the input image should be the program execution or the ground truth shape.

Figure 5 shows some examples of reconstructions produced by LEST-inferred programs vs. RL-inferred ones. While RL’s programs tend to get the bulk structure of the shape correct, it misses finer details and often uses more verbose programs.

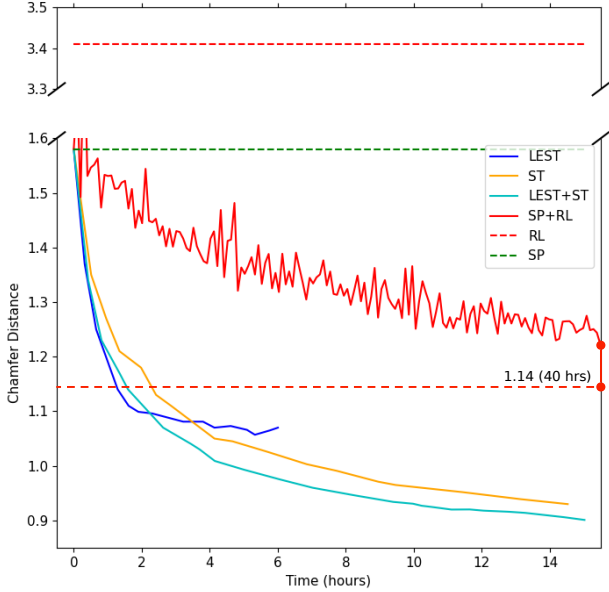


Figure 4. Test set reconstruction accuracy convergence rate for different learning methods on 2D CSG, measured by Chamfer Distance; so lower is better. LEST converges much faster compared with RL, but after several hours of training it reaches a local minima. This behavior can be avoided for 2D CSG by combining LEST and ST.

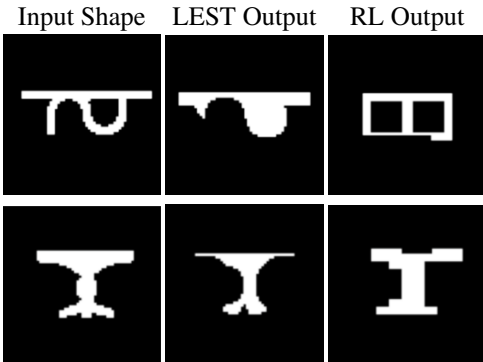


Figure 5. Qualitative comparison of how LEST-inferred CSG programs reconstruct shapes vs. RL-inferred programs.

4.2. ShapeAssembly

We next evaluate LEST’s ability to infer programs in the ShapeAssembly language [11]. ShapeAssembly is designed for specifying the hierarchical part structure of manufactured 3D objects. It creates objects by declaring cuboidal part geometries and then assembling those parts together via attachment and symmetry operators. Inferring ShapeAssembly programs is harder than inferring 2D CSG programs: the input shapes are 3D, and the semantics of ShapeAssembly’s program executor are considerably more complicated than the rules of CSG. Compared to CSG, ShapeAssembly is “ill-conditioned”: small changes to a pro-

gram can lead to large changes in its output shape (e.g. changing one parameter of a translational symmetry operator can create or delete significant pieces of geometry).

Real Shapes S^* Our target distribution of 3D shapes is derived from CAD models in ShapeNet [1]. We use 3,758 shapes from the *Chair* category, and represent each instance as a 64 x 64 x 64 voxel grid, using the voxelizations from [3]. For all experiments, we divide these shapes into train/validation/test splits (80/10/10).

Random Program Generator R_G ShapeAssembly programs can be created by sampling from the ShapeAssembly grammar. The grammar has operators that create primitives (Cuboid), attach primitives to one another (attach), or behave as function macros for higher order spatial relationships (reflect, translate, squeeze). Operators consume a mix of continuous and discrete parameters; for instance the starting dimensions of a primitive are continuous but the axis of a symmetry operation is a discrete choice. Additionally, ShapeAssembly programs are hierarchical, as primitives are allowed to expand into sub-programs with more detailed geometry.

We initially used an R_G that sampled unstructured ShapeAssembly programs, but as described in the supplemental material, our initial recognition network was unable to learn from this distribution. This is not surprising, as ShapeAssembly programs are “ill-conditioned”: small changes in program space can lead to large changes in geometry space. To overcome this, we added more structure to R_G . It first creates a root program by sampling a “skeleton” of vertically stacked primitives, and optionally extend this skeleton by connecting “limbs” to it. Then, primitives in the skeleton are allowed to decompose into sub-programs. Each sub-program contains primitives that attach to the (i) top of the bounding volume, (ii) the bottom of the bounding volume, or (iii) other primitives in the sub-program. We further constrain the quality of sampled programs via rejection sampling, keeping only programs where (i) each generated primitive occupies at least 8 voxels, (ii) 50% of each primitive’s occupied voxels are covered uniquely by that primitive, and (iii) a pretrained voxel-based classifier has over 50% confidence the resulting structure is a chair. We use this procedure to sample 75,000 ShapeAssembly programs, which we visualize in the supplemental material.

Network architecture for $p(z|x)$ Our ShapeAssembly program inference network is based on the decoder network from the original ShapeAssembly paper [11]. It is a hierarchical sequence model that uses a GRU cell to predict both program lines and the hierarchy structure. We use dropout for the network layers that re-encode program lines and

also add a small amount of Gaussian noise to the ground-truth continuous parameters in the input sequences, as these modifications improved reconstruction performance during supervised pretraining. The hidden state of the root level GRU cell is initialized with the output of a 3D CNN encoder that consumes the voxelized input shape. We reuse the 3D CNN architecture from CSGNet.

Experiment details We pretrained our supervised inference network for 140 epochs (24 hours) on the 75,000 randomly generated programs, stopping based on validation set reconstruction performance. We run 6 rounds of self-training, stopping each round once the model has converged on validation set reconstruction performance. Both of these stages were trained with a batch size of 64 and optimized using Adam [14] with a learning rate of 0.0001.

For RL, we copy the procedure of CSGNet whenever possible. We use the same Chamfer Distance based reward, although we compare point sets sampled from 3D shape surfaces (instead of 2D shape edges). Naively extending our decoder architecture that makes many continuous parameter predictions with REINFORCE training led to degenerate performance, as discussed in the supplemental. Thus, we design a specialized version of our program decoder where all continuous parameter predictions are replaced with discretized options: 11 bins for attachment parameters and 20 bins for cuboid parameters. During supervised pretraining, we replace corresponding L1 losses with Cross Entropy losses. This discretized inference network was also pretrained until convergence based on validation set reconstruction accuracy (39 epochs, 8 hours).

When inferring programs, we use a beam search procedure over the discrete program components with beam width of 10. All timings were collected on a machine with a GeForce RTX 2080 Ti GPU and an Intel i9-9900K CPU.

Results Figure 6 top plots the distributional similarity between S and S^* over the course of LEST self-training, and Figure 6 bottom shows how the output of the inferred program for a few shapes $\mathbf{x}^* \in S^*$ evolves over time. To measure distributional similarity, we calculate the Frechet Distance (FD) [8] between the test set of S^* and the inferred S in the feature space of a 3D CNN pretrained on a ShapeNet classification task. A lower FD implies that the two distributions are more similar. Through rounds of self-training, LEST’s inferred programs move away from the random data of R_G towards the target distribution S^* . This result is reinforced by qualitative analysis of the inferred shapes. Multiple rounds of LEST self-training leads to shapes that have better geometric and structural matches to the input.

Figure 7 plots the reconstruction accuracy of LEST vs. other methods as a function of training time. We evaluate reconstruction accuracy with the F1-score metric [15],

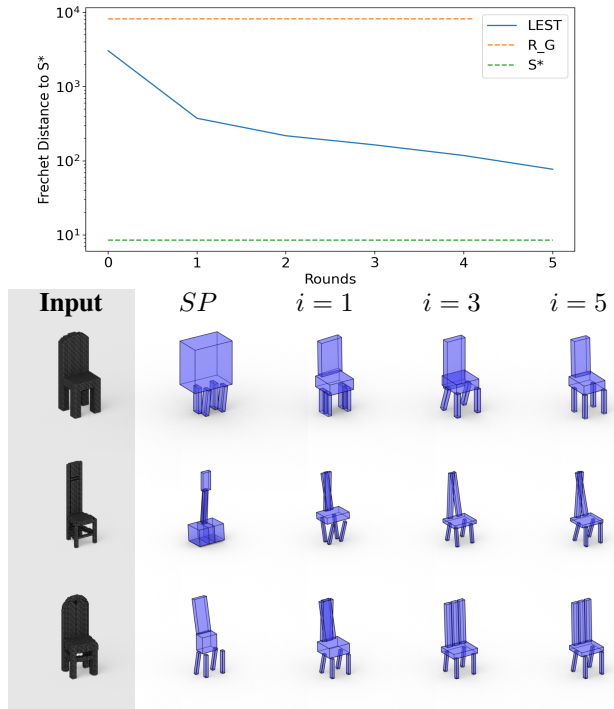


Figure 6. (Top) The distributional similarity of S to S^* with increasing LEST self training rounds, for ShapeAssembly; lower is better. (Bottom) How the output of the inferred program for specific shapes changes over LEST rounds.

by sampling a point cloud from the ground-truth ShapeNet mesh and the cuboid part proxies output by a ShapeAssembly program. Each point cloud is composed of 4096 points and we use a distance threshold of one voxel width. In the ShapeAssembly domain, LEST achieves the best F1-score and outperforms all other variants (notice that we plot 1 - F1-score to be consistent with our other plots). Once again, we see that LEST improves its performance with additional rounds of self-training, with each round taking about an hour to run. Notice that due to the difficulty of the ShapeAssembly problem space, ST and RL perform much worse than they did for 2D CSG, and, in fact, using RL or ST to fine-tune a SP model only hurts performance.

Figure 8 shows some examples of reconstructions produced by LEST-inferred programs vs. RL-inferred ones. Unlike 2D CSG, structural decisions made in ShapeAssembly programs can have non-local consequences. As a result, the RL models we train fail to converge to any intelligent behavior, even when initialized from a pretrained model, and their output programs are simplistic in nature. In comparison, the program inferences from LEST better reflect the structural and geometric elements of the input shapes.

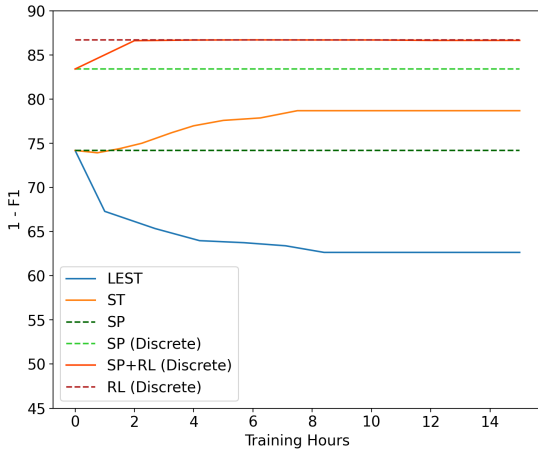


Figure 7. Test set reconstruction accuracy convergence rate for different learning methods on ShapeAssembly. We use 1 - F-Score as the evaluation metric, so lower is better. LEST is the only method that improves upon the performance of the pretrained supervised models. RL and ST perform poorly due to the high complexity of the domain.

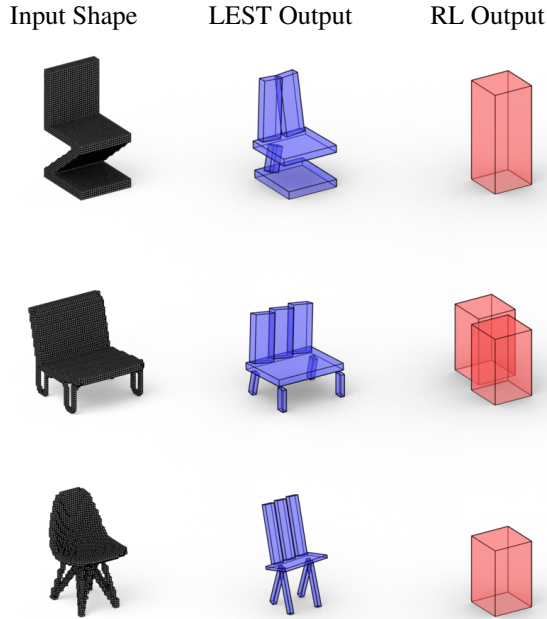


Figure 8. Qualitative comparison of how LEST ShapeAssembly programs reconstruct shapes vs. RL programs.

5. Conclusion & Future Work

In this paper, we presented Latent Execution Self Training (LEST), a new approach for unsupervised learning of shape program inference that works with black box program executors. LEST improves upon the idea behind self-training by leveraging the executable nature of latent vari-

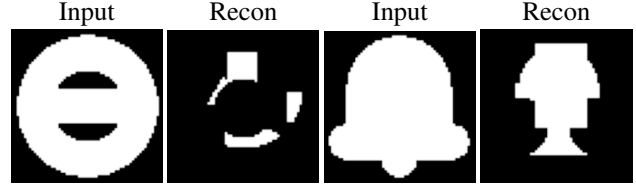


Figure 9. A 2D CSG inference network self-trained with LEST on CAD shapes (desk, chair, lamp) fails to generalize to icon shapes.

ables which are programs, thus producing pseudo-labels which are guaranteed to be “correct” for their associated inputs. The iterative re-training procedure that powers LEST is similar to variational methods for learning generative models, but LEST does not require learning a generative prior over programs. In experiments on inferring CSG programs and ShapeAssembly programs, LEST converges faster and to better results than other baselines, including reinforcement learning, which is the current de facto approach for black box program executors.

One limitation of LEST we have found is that learned inference networks do not generalize well beyond the distribution of real shapes S^* on which they were self-trained. Figure 9 shows examples of using the 2D CSG inference network trained on CAD shapes to infer CSG programs for icons from The Noun Project¹. This failure is not surprising: supervised learning (of which self-training is a variant) can typically be made to generalize to held-out samples from the *same* distribution as the training data, but not to held-out samples from a drastically *different* distribution.

The types of programs that LEST learns to infer are heavily influenced by the initial random programs from R_G . This dependence has two important consequences. First, as shown by our ShapeAssembly experiments, a “purely random” R_G does not always work to initialize the self-training process; complex shape grammars G may require an R_G that encodes some domain knowledge of typical program structure, which does take some effort for each new shape domain. Second, we note that while our work (and prior work in shape program inference) focuses on reconstruction quality, getting a good program structure matters just as much if the program is to be usable for e.g. editing and manipulation tasks. Currently, R_G is the only place where knowledge about what constitutes “good program structure” can be injected. Such knowledge must be expressed in procedural form, which may be more challenging to elicit from domain experts than declarative knowledge (i.e. “a good program has these properties” vs. “this is how you write a good program”). Finding efficient ways to elicit and inject such knowledge from people is an important future direction for all research into unsupervised and weakly-supervised shape program inference.

¹<https://thenounproject.com>

Acknowledgments

We would like to thank Stephen Bach for conversations about pseudo-labels and Michael Littman for the connection to expectation maximization. Alexander Simone created the Noun Project icon² we show as a failure case. This work was funded in part by NSF Award #1941808. Daniel Ritchie is an advisor to Geopipe, Inc. and owns equity in the company. Geopipe is a start-up that is developing 3D technology to build immersive virtual copies of the real world with applications in various fields, including games and architecture.

References

- [1] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An Information-Rich 3D Model Repository. *arXiv:1512.03012*, 2015. 2, 6
- [2] Siddhartha Chaudhuri, Daniel Ritchie, Jiajun Wu, Kai Xu, and Hao Zhang. Learning Generative Models of 3D Structures. *Computer Graphics Forum*, 2020. 2
- [3] Zhiqin Chen, Kangxue Yin, Matthew Fisher, Siddhartha Chaudhuri, and Hao Zhang. Bae-net: Branched autoencoder for shape co-segmentation. *Proceedings of International Conference on Computer Vision (ICCV)*, 2019. 6
- [4] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22, 1977. 2
- [5] Kevin Ellis, Maxwell Nye, Yewen Pu, Felix Sosa, Josh Tenenbaum, and Armando Solar-Lezama. Write, execute, assess: Program synthesis with a repl. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 2, 3
- [6] Kevin Ellis, Daniel Ritchie, Armando Solar-Lezama, and Josh Tenenbaum. Learning to Infer Graphics Programs from Hand-Drawn Images. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. 2
- [7] Kevin Ellis, Catherine Wong, Maxwell Nye, Mathias Sable-Meyer, Luc Cary, Lucas Morales, Luke Hewitt, Armando Solar-Lezama, and Joshua B. Tenenbaum. Dreamcoder: Growing generalizable, interpretable knowledge with wake-sleep bayesian program learning, 2020. 3
- [8] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NeurIPS*, 2017. 7
- [9] Luke B. Hewitt, Tuan Anh Le, and Joshua B. Tenenbaum. Learning to learn generative programs with memoised wake-sleep, 2020. 3
- [10] Geoffrey E Hinton, Peter Dayan, Brendan J Frey, and Radford M Neal. The “wake-sleep” algorithm for unsupervised neural networks. *Science*, 268(5214):1158–1161, 1995. 2, 3
- [11] R. Kenny Jones, Theresa Barton, Xianghao Xu, Kai Wang, Ellen Jiang, Paul Guerrero, Niloy J. Mitra, and Daniel Ritchie. Shapeassembly: Learning to generate programs for 3d shape structure synthesis. *ACM Transactions on Graphics (TOG), Siggraph Asia 2020*, 39(6):Article 234, 2020. 2, 4, 6
- [12] Kacper Kania, Maciej Zieba, and Tomasz Kajdanowicz. Ucs-g-net – unsupervised discovering of constructive solid geometry tree. In *arXiv*, 2020. 2, 3
- [13] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *ICLR 2015*, 2015. 5
- [14] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *International Conference on Learning Representations (ICLR)*, 2014. 2, 3, 7
- [15] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics*, 36(4), 2017. 7
- [16] Yunchao Liu, Zheng Wu, Daniel Ritchie, William T. Freeman, Joshua B. Tenenbaum, and Jiajun Wu. Learning to Describe Scenes with Programs. In *International Conference on Learning Representations (ICLR)*, 2019. 2
- [17] Radford M. Neal and Geoffrey E. Hinton. A new view of the em algorithm that justifies incremental and other variants. In *Learning in Graphical Models*, pages 355–368. Kluwer Academic Publishers, 1993. 3
- [18] H. Scudder. Probability of error of some adaptive pattern-recognition machines. *IEEE Trans. Inf. Theor.*, 11(3):363–371, Sept. 2006. 2
- [19] Gopal Sharma, Rishabh Goyal, Difan Liu, Evangelos Kalogerakis, and Subhansu Maji. CSGNet: Neural Shape Parser for Constructive Solid Geometry. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2, 3, 4
- [20] Yonglong Tian, Andrew Luo, Xingyuan Sun, Kevin Ellis, William T. Freeman, Joshua B. Tenenbaum, and Jiajun Wu. Learning to Infer and Execute 3D Shape Programs. In *International Conference on Learning Representations (ICLR)*, 2019. 2, 3
- [21] Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8, 1992. 3, 4
- [22] Chenghui Zhou, Chun-Liang Li, and Barnabas Poczos. Unsupervised program synthesis for images using tree-structured lstm, 2020. 3
- [23] Barret Zoph, Golnaz Ghiasi, Tsung-Yi Lin, Yin Cui, Hanxiao Liu, Ekin D. Cubuk, and Quoc V. Le. Rethinking pre-training and self-training, 2020. 2

²<https://thenounproject.com/AlexanderSimone/collection/simple-ui>

6. Supplementary Material

6.1. Additional 2D CSG Results

Generative Model Experiments As discussed in the main text, there are a number of methods related to LEST that train a generative model in addition to a recognition model, including variational Bayes, wake-sleep, and expectation maximization. As a proxy for these methods, we experimented with a variant of LEST that incorporates a generative model. In each round, a generative model, $p(\mathbf{x}, \mathbf{z})$, is trained to convergence on programs inferred in the previous round. Samples from the generative model are then used to train the recognition model, $p(\mathbf{z}|\mathbf{x})$, to convergence. Figure 10 shows how the reconstruction accuracy of this method, “LEST + Generative Model,” compares to LEST, ST, and RL. With the addition of the generative model, this new variant performs slightly worse than standard LEST, likely because samples from the generative model are farther from the distribution of target shapes than the executed programs, S , used by LEST. We also compared the quality of the generative model produced by this new variant with a generative model trained using the best inferred programs from LEST. We found the latter had a lower bidirectional average nearest neighbor distance on the test set (2.72 vs 2.56).

Domain Gap Experiments One reason we believe that ST eventually outperforms LEST in the domain of 2D CSG is because of the visual domain gap between the two methods. ST only sees images from S^* , LEST only sees images from S . To analyze this phenomenon experimentally, we compare LEST, ST and LEST+ST on reconstruction performance, while freezing the weights of the encoder CNN. As shown in Figure 11, freezing the weights of the CNN encoder hurts reconstruction performance for both ST and LEST+ST, but actually helps reconstruction performance for LEST. Looking at the best achieved test set Chamfer Distance in Table 2, the performance gap between ST and LEST disappears when the visual encoding remains static. This confirms our hypothesis that the domain gap between images from S^* and S explains why ST can outperform LEST when the CNN encoder is allowed to change.

Qualitative Comparisons Figures 13-15 show additional qualitative comparisons of LEST + ST, LEST, ST, and SP + RL.

6.2. Additional ShapeAssembly Results

Analysis of Structured and Unstructured R_G We experimented with different formulations of R_G for ShapeAssembly. Our structured formulation of R_G is described in Section 4.2 of the main text. Our unstructured formulation of R_G removes the notion of “skeleton” and

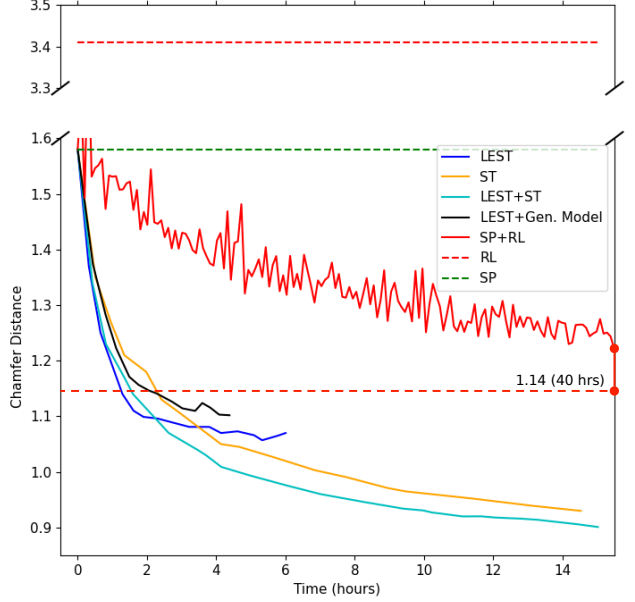


Figure 10. Test set reconstruction accuracy convergence rate for different learning methods on 2D CSG, measured by Chamfer Distance; so lower is better. Adding a generative model in between LEST rounds to generate the next set of training (program, image) pairs leads to slower convergence and a worse final reconstruction accuracy.

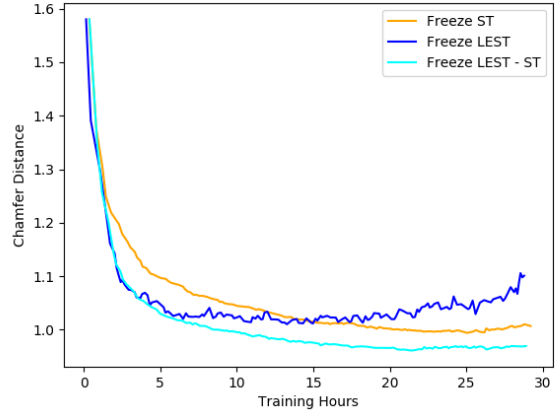


Figure 11. Test set reconstruction accuracy convergence rate for different learning methods on 2D CSG while keeping the weights of the CNN encoder frozen, measured by Chamfer Distance; so lower is better. ST and LEST converge to similar best reconstruction accuracies, while LEST+ST once again performs the best overall.

“limbs”, and instead randomly adds parts to the scene one at a time, by connecting a new cuboid to existing parts, and then optionally invoking symmetry commands. Sampled programs are rejected using the same criteria as the structured case (i.e. parts can’t overlap too much, must take up

Method	ST	LEST	LEST+ST
Default	0.930	1.041	0.901
Freeze CNN	0.994	1.010	0.961

Table 2. Best recorded test set Chamfer Distance on the 2D CSG domain for three learning methods: ST, LEST and LEST+ST. In the default row, the CNN encoder is allowed to update. In the freeze CNN row, the weights of the CNN encoder are kept static. The reconstruction gap between ST and LEST disappears when the visual encoding is kept constant. In both paradigms, LEST+ST achieves the best of both worlds.

at least 8 voxels, and a pre-trained model must think the resulting structure is “chair-like”). A qualitative comparison of how the structured and unstructured generations differ can be found in Figure 16.

In Figure 12, we quantitatively evaluate how using unstructured R_G compares with structured R_G measured by reconstruction performance. Of note, using unstructured R_G hurts the performance of the initial recognition model pretrained with supervised learning, as seen by comparing SP (Structured) with SP (Unstructured). With unstructured R_G , the starting F-Score is 12.9, while with structured R_G , the starting F-score is 25.8. Multiple rounds of LEST do improve the SP model trained on unstructured R_G , increasing the F-score from 12.9 to 18.04, but this is much below the best reconstruction performance of LEST when initialized with the SP model trained on structured R_G , which achieves an F-score of 37.3.

Reinforcement Learning with Continuous Predictions

In the main text, we compare ShapeAssembly LEST against a discrete RL variant as described in Section 4.2. We also trained an RL model using the default continuous ShapeAssembly model architecture. In order to train the continuous predictions with REINFORCE, we treat them all as Gaussians. The original prediction for each continuous value is interpreted as the mean of the distribution and we add an extra MLP head to the network for each predicted continuous value that we interpret as the standard deviation. As seen in Figure 12, this variant fails to learn anything meaningful, and performs much worse than the discrete RL counterpart.

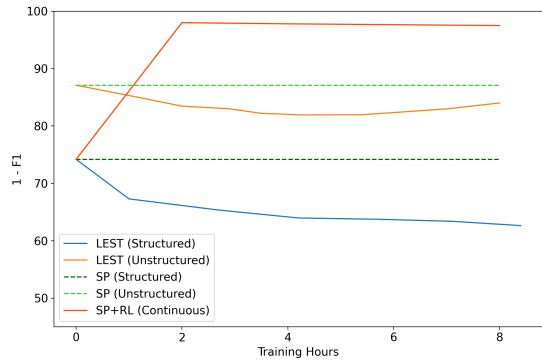


Figure 12. Test set reconstruction accuracy convergence rate for different learning methods on ShapeAssembly. We use 1 - F-Score as the evaluation metric, so lower is better. Using a structured R_G improves reconstruction performance compared with an unstructured R_G . The SP+RL variants fails to learn anything meaningful when it must predict continuous parameters.

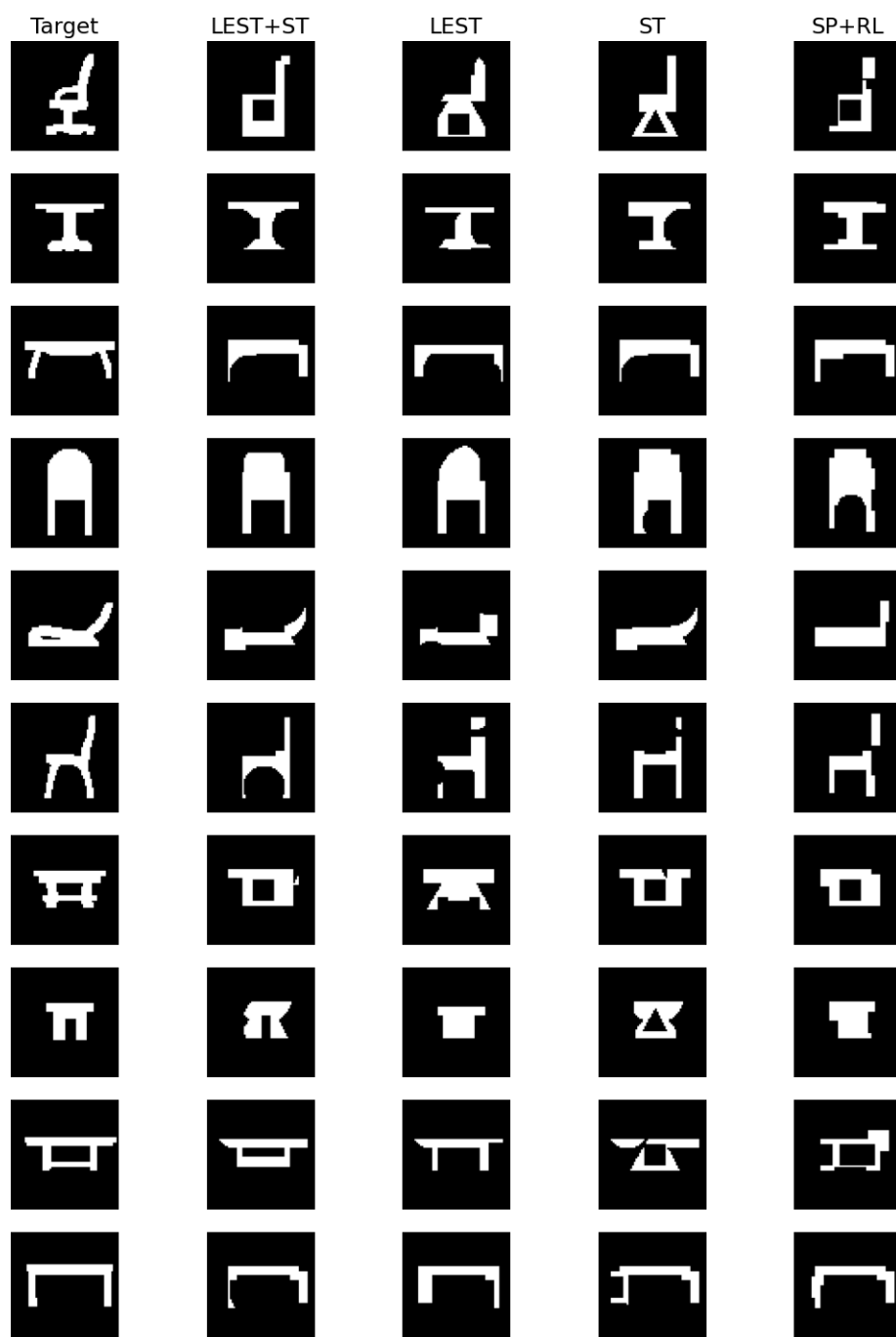


Figure 13. Qualitative comparison of CSG program reconstructions from LEST + ST, LEST, ST, and SP + RL.

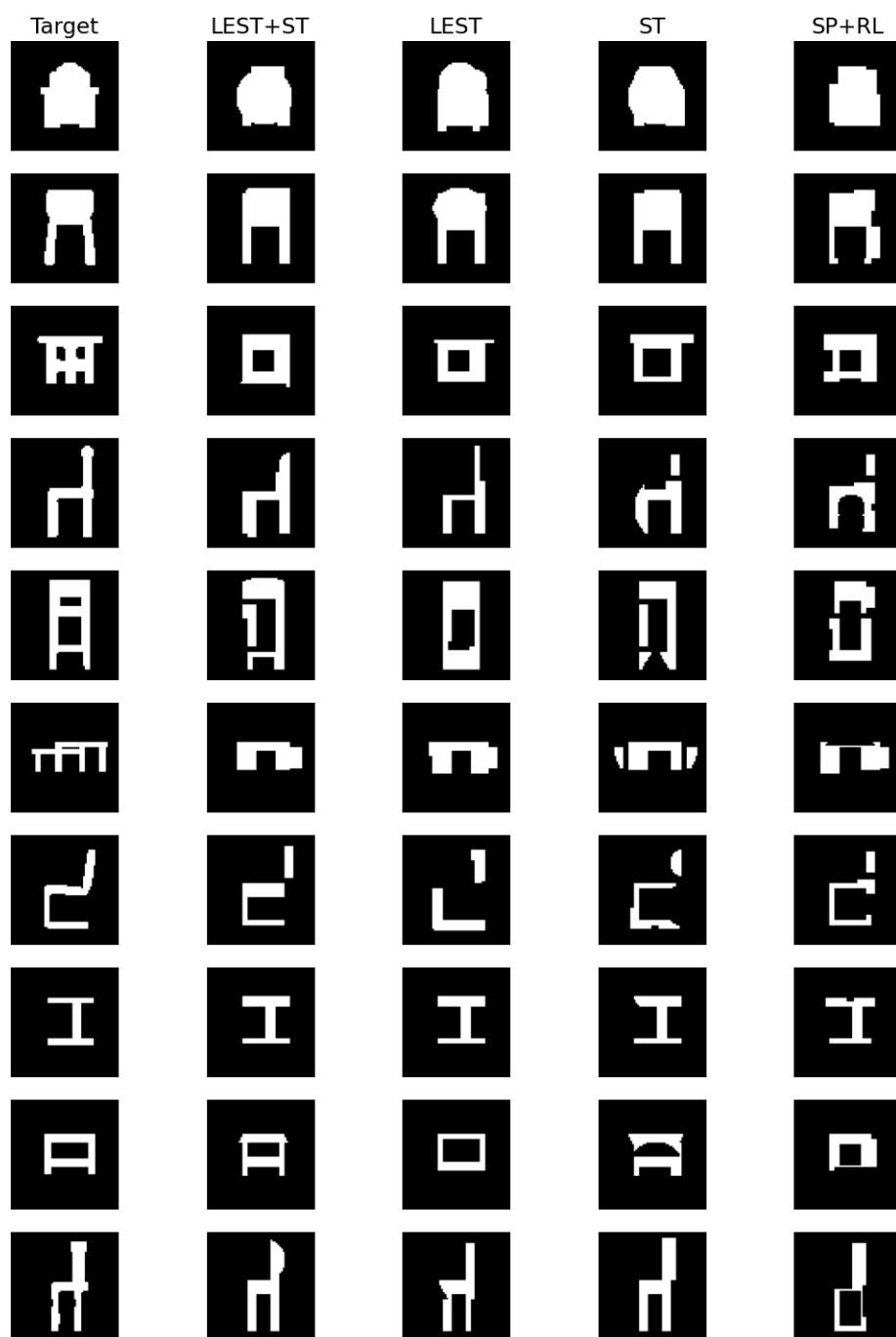


Figure 14. Qualitative comparison of CSG program reconstructions from LEST + ST, LEST, ST, and SP + RL.

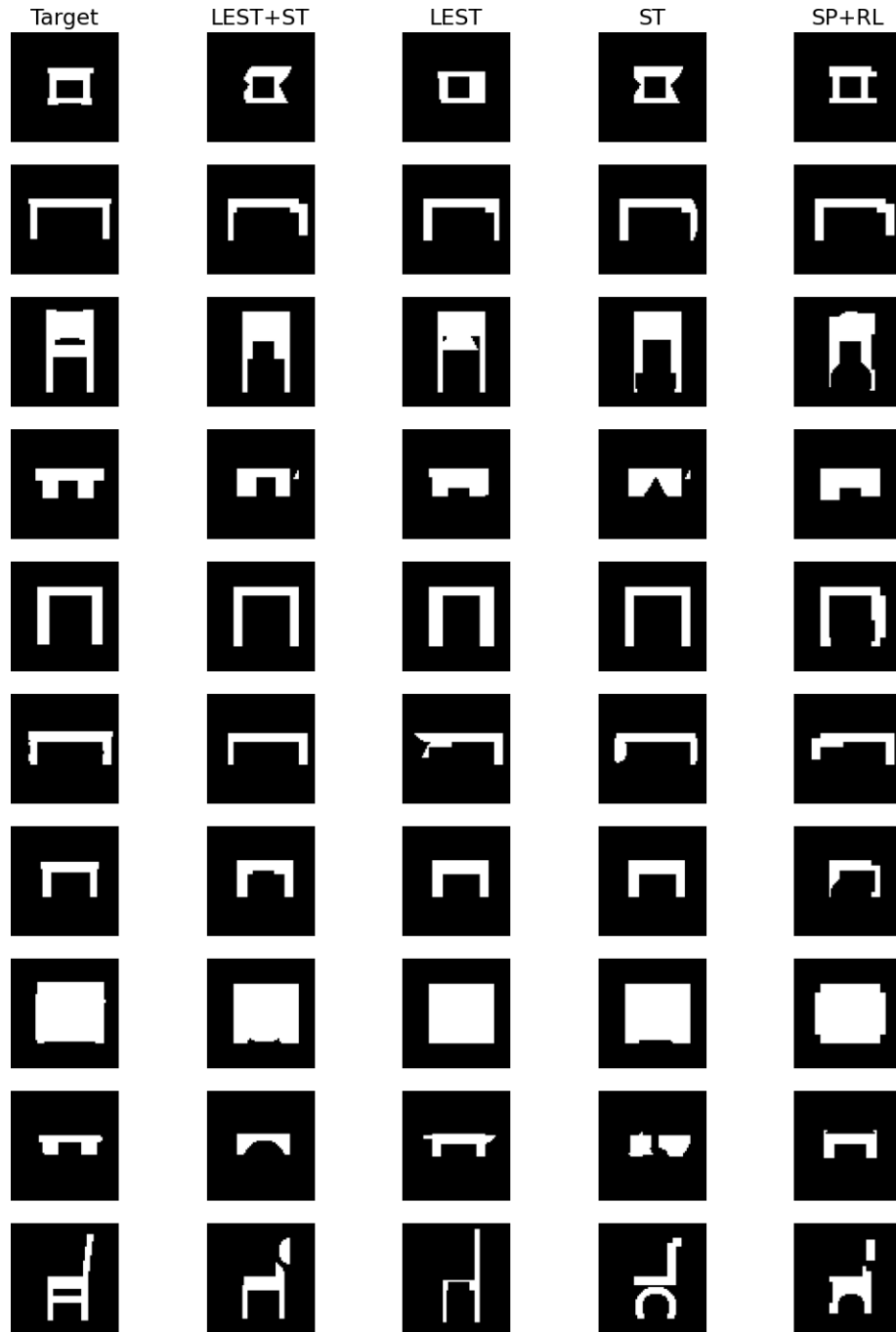


Figure 15. Qualitative comparison of CSG program reconstructions from LEST + ST, LEST, ST, and SP + RL.

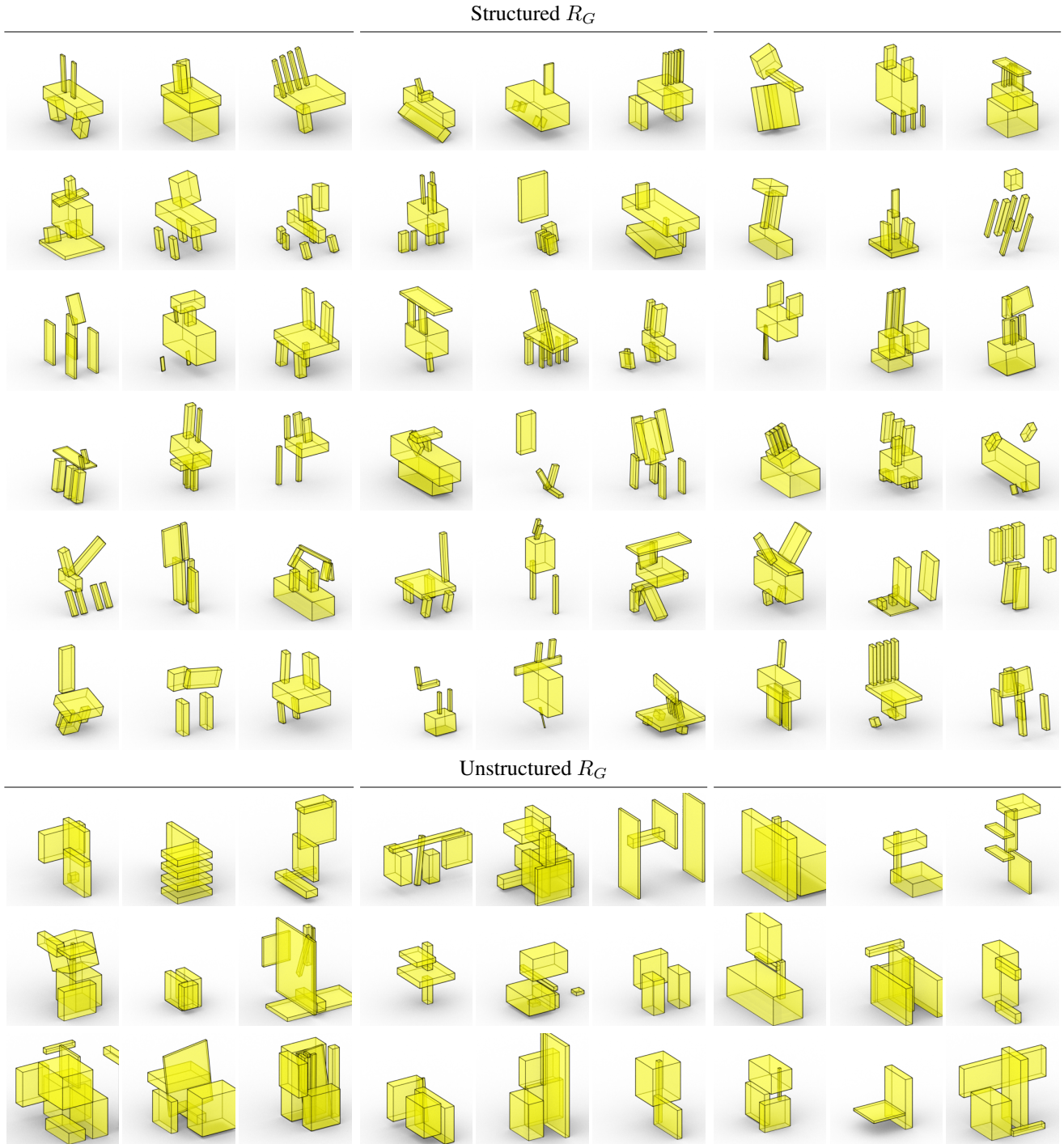


Figure 16. Qualitative samples of Structured R_G and Unstructured R_G for ShapeAssembly.