

# Self-Supervised Representation Learning for CAD

Benjamin T. Jones\*  
University of Washington  
Seattle, WA 98195

benjones@cs.washington.edu

Michael Hu  
University of Washington  
Seattle, WA 98195

mkhu@cs.washington.edu

Vladimir G. Kim†  
Adobe Research  
Seattle, WA 98103

vokim@adobe.com

Adriana Schulz‡  
University of Washington  
Seattle, WA 98195

adriana@cs.washington.edu

## Abstract

*The design of man-made objects is dominated by computer aided design (CAD) tools. Assisting design with data-driven machine learning methods is hampered by lack of labeled data in CAD’s native format; the parametric boundary representation (B-Rep). Several data sets of mechanical parts in B-Rep format have recently been released for machine learning research. However, large scale databases are largely unlabeled, and labeled datasets are small. Additionally, task specific label sets are rare, and costly to annotate. This work proposes to leverage unlabeled CAD geometry on supervised learning tasks. We learn a novel, hybrid implicit/explicit surface representation for B-Rep geometry, and show that this pre-training significantly improves few-shot learning performance and also achieves state-of-the-art performance on several existing B-Rep benchmarks.*

## 1. Introduction

Almost every human-made object that exists today started its life as a model in a CAD system. As the prevalent method of creating 3D shapes, repositories of CAD models are extensive. Further, CAD models have a robust structure, including geometric and program representations that have the potential to expose design and manufacturing intent. Learning from CAD data can therefore enable a variety of applications in design automation and design- and fabrication-aware shape reconstruction and reverse engineering. An important challenge in learning from CAD is

that most of this data does not have labels that can be leveraged for inference tasks. Manually labeling B-Rep data is time consuming and expensive, and the specialized format requires CAD expertise, making it impractical for large collections.

In this work we ask: how can we leverage large databases of *unlabeled* CAD geometry in analysis and modeling tasks that typically require labels for learning?

Our work is driven by a simple, yet fundamental observation: the CAD data format was not developed to enable easy visualizing or straightforward geometric interpretation, it is a format designed to be compact, have infinite resolution, and allow easy editing. Indeed CAD interfaces consistently run sophisticated algorithms to convert the CAD representation into geometric formats for rendering. Driven by this observation our key insight is to leverage large collections of unlabeled CAD data to learn to geometrically interpret the CAD data format. We then leverage the networks trained over the geometric interpretation task in supervised learning tasks where only small labeled collections are available. In other words, we use geometry as a model of *self-supervision* and apply it to *few-shot-learning*.

Specifically, we learn to rasterize local CAD geometry using an encoder-decoder structure. The standard CAD format encodes geometry as parametric Boundary Representations (B-Rep). B-reps are graphs where the nodes are parametric geometry (surfaces, curves, and points) and edges denote the topological adjacency relationships between the geometry. Importantly, the parametric geometry associated with each node is *unbounded* and bounds are computed from the topological relationships: curves bounding surfaces and points bounding curves. As shown in Figure 2, the geometry of B-rep face is computed by *clipping* the surface primitive, where the clipping masked is constructed from

\*<http://www.bentodjones.com>

†<http://www.vovakim.com>

‡<https://homes.cs.washington.edu/~adriana>

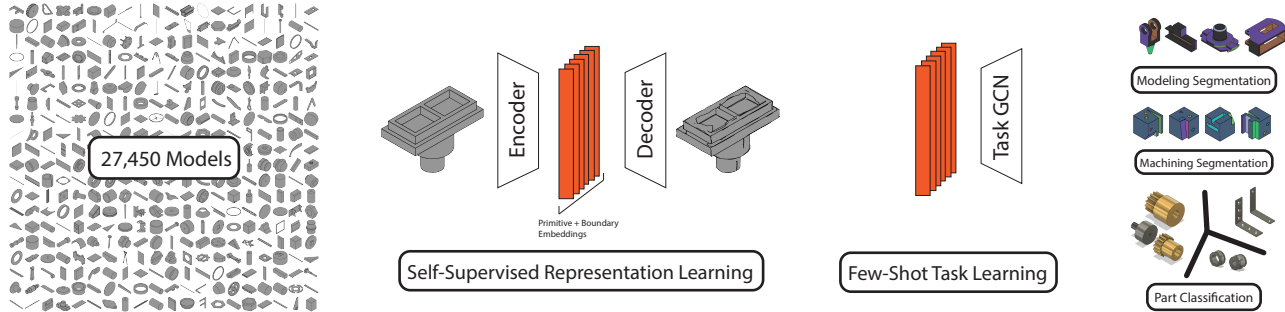


Figure 1. Overview of our technique. We train a geometric self-supervision task of a large, unlabeled dataset of CAD Boundary Representations (B-Rep) to learn geometrically relevant representations for each B-Rep face. These pre-trained representations are then used to train few-shot segmentation and classification learning tasks on labeled B-Rep datasets.

adjacent edges.

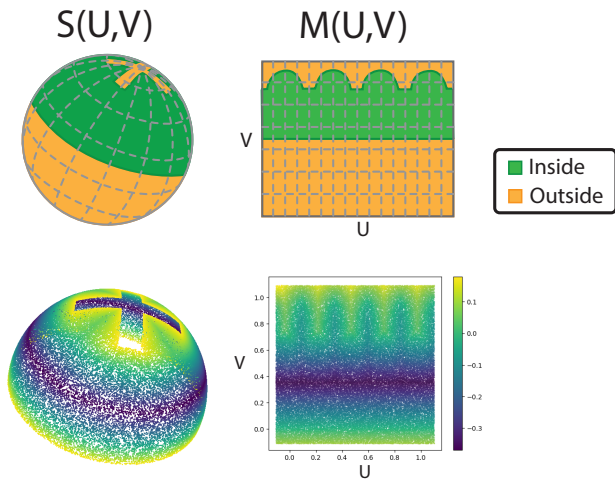


Figure 2. Face representation in B-Rep geometry. Above: schematic of a parameterized surface (sphere, left) with its unrolled clipping plane (right). Adjacent edges (dark green) implicitly define the clipping boundary. Below: sampled signed distance field (SDF) of the clipping plane in flattened (right) and embedded (left) form; this is the function we actually learn.

This means that B-reps are constructed piecewise by *explicitly* defined surfaces with *implicitly* defined boundaries. This is the observation that drives our proposed learning architecture: reconstructing faces by jointly decoding the explicit surface parameterization as well as the implicit surface boundary. Our proposed encoder uses message passing on the topological graph to capture the boundary information to encode B-Rep faces. To handle the graph heterogeneity (nodes comprised of faces, edges, and, vertices), we use a hierarchical message passing architecture inspired by the Structured B-Rep GCN [16]. Our decoder uses the learned embeddings as latent codes for two per-face neural function evaluators: one mapping from  $\mathbb{R}^2 \rightarrow \mathbb{R}^3$  that

encodes the face’s parametric surface, and one mapping  $\mathbb{R}^2 \rightarrow \mathbb{R}$  that encodes the face’s boundary as a signed distance field (SDF) *within* the parametric surface.

We apply our proposed model of B-rep self-supervision to learn specialized B-Rep tasks from very small sets of labeled data—10s to 100s of examples, vs 10k to 100k. This is done by using the embeddings learned on self-supervision as input features to supervised tasks. We evaluate our approach on three tasks and data-sets from prior work [3, 6, 21]. We validate our findings across varying training set sizes and show that our model consistently outperforms the prior supervised approaches with a significant improvement in performance on smaller training sets. This result shows that self-supervision on unlabeled CAD data can enable supervised learning on small labeled datasets paving the way to many exciting application is the domain—from modeling, to analysis, reconstruction and beyond. By using less data, our approach is also significantly faster to train, enabling applications that depend on training speed. Since our self-supervision task is in a differentiable CAD rasterizer, we use it to prototype a reverse engineering task, demonstrating how it can be used in the future for different applications that require gradient-based optimization over B-Reps.

## 2. Related Work

**Learning from CAD Collections** Recently, interest in learning from CAD data has increased due to advances in machine learning and the release of large CAD datasets, including collections of B-reps and program structure [20,36], CAD sketches [31], and CAD assemblies [16]. Prior work has leveraged such collections for many applications including segmentation [6], classification [3], assembly suggestions [16,35], and generative design of CAD sketches [10,28,31,32], B-Reps [13,14], and CAD programs [36,37]. There is, however, a fundamental gap between the capabilities shown in past work and real-world applications that can fundamentally impact CAD design,

such as auto-complete modeling interfaces or reconstruction of complex geometries: the lack for task-specific labels in large datasets. For example, while most prior work leveraged the Onshape Public dataset [20], this dataset contains mostly designs created by novice CAD users, not capturing the design process of CAD experts. The Fusion 360 dataset [36] is significantly smaller compared to the Onshape collection, and other public resources, such as GrabCAD [1], contain designs from multiple CAD systems that are mostly unlabeled. In this work, we make a fundamental step toward enabling novel applications on CAD learning by proposing a novel direction for leveraging unlabeled data in supervised learning of CAD geometry.

**Learning on B-Reps.** A CAD B-Rep is a specialized data structure that encodes solid geometry as a graph of parametric shapes and their topological relationships, which has the advantages of arbitrary spatial resolution and easy programmatic editability. B-reps can be exported to other common geometric representations, such as polygonal meshes, using geometric CAD kernels [2]. Several techniques for learning on B-reps use message passing networks. BRepNet [21] and UV-Net [15] create a reduced graph of B-Rep faces, while SB-GCN [16] propose a hierarchical structure over the 4 classes of topological entities (faces, loops, edges, and vertices). These methods, however, require a CAD kernel in the loop to generate the features from the B-Rep format and large labeled sets for learning. While we use a CAD kernel to train our encoder-decoder on unlabeled data, our task-specific networks can be trained on small sets of labeled data without the need of a CAD kernel at inference time.

**Neural Shape Representations.** Neural shape generation is an active research area with a large number of representations used by prior techniques. Some methods operate over a fixed discretization of the domain into points [4, 9], voxel grids [5, 24], or vertex coordinates of a mesh template [33]. Due to irregularity of geometric data, functional representation is often used, learning to represent shapes as continuous functions, such as surface atlases [12, 38] or signed distance fields defined over a volume [7, 25, 29]. In this work we chose to use functional representation of the output shape as a neural occupancy and 3D mapping function over UV domains. Functional representation is well-suited for heterogeneous geometry with varied levels of detail, since it does not require choosing a fixed sampling rate. Furthermore, our representation that combines implicit fields and atlas-like embeddings directly produces a surface and can be easily supervised with the ground truth B-Rep data.

**Few Shot Learning.** Performance of strongly-supervised methods is commonly hampered by lack of training data. Few-shot learning techniques often rely on learning rich

features in a self-supervised fashion, and then using a few examples to adapt these features to a new task [34]. One common self-supervision strategy is to withhold some data from the original input, and train a network to predict it. For example, one can remove color from images and train a network to colorize [22, 40], remove part of the image and train a network to complete it [30], randomly perturb orientation and predict the upright position [11], or randomly shuffle patches and predict their true ordering [8, 26]. Another commonly used tool is auto-encoders that encode input to a lower-dimensional space and then attempt to reconstruct it with a decoder [18, 19]. Our approach is closest to the auto-encoders, except our input and output are in two different representations: CAD B-Reps and surface rasterizations. This enables us to learn features related to the actual 3D geometry.

### 3. Geometric Self-Supervision

Our goal is to learn relevant features of CAD B-Reps that can be used on a number of different modeling and analysis tasks. We use an encoder-decoder architecture on B-Rep surfaces to learn a latent space of relevant features. Based on the insight that the learned features should include geometric understanding of CAD shapes, we train our encoder-decoder to learn to rasterize CAD models. We further choose to learn this embedding at the face level as opposed to a feature per part. This is driven by our application domain, where tasks typically require understanding local topological information. An overview of our architecture is shown in Figure 1.

**Decoder** As previously discussed, the format we selected for the decoder output is driven by the observation that B-Reps are compositions of explicitly represented surfaces with implicitly represented boundaries. For example, the geometry associated with faces are unbounded parametric surfaces  $S : \mathbb{R}^2 \rightarrow \mathbb{R}^3$  and the boundary of such surfaces is captured by the geometry of the neighboring edges. Edges in turn are defined by unbounded parametric curves bounded by neighboring vertices. Vertices are represented as points. In addition to the domain parameters, parametric geometry has additional, fixed parameters (like radius or cone angle) that we call *shape parameters*. A bounded surface is called a *clipped surface*, and the bounding function, which is defined *implicitly* by the bounding topology, is its *clipping mask*, illustrated in Figure 2.

For self-supervision, we need to choose an output geometric representation for surface patches with boundary. Bounded patches do not have a natural parameterization and so are not suitable to learn directly as parametric functions. They also are difficult to represent as neural implicit. [27] However, the clipping function itself defines a closed region *within* the parameterization of the supporting surface,

which is a function that lends itself well to implicit representation in 2D. We therefore choose to learn an implicit function of the clipping region as a 2D signed distance function. Since this representation relies on an explicit surface parameterization, we also learn the supporting surface parameterization. Crucially, this does not require parameterizing a boundary.

We choose to use a conditional neural field as our decoder, since this representation can capture both explicit and implicit geometry. Specifically, the explicit parametric surface is an  $\mathbb{R}^2 \rightarrow \mathbb{R}^3$  function, mapping  $(u, v)$  coordinates of a face’s supporting surface to the 3D position of that point  $(x, y, z)$ . The clipping mask encoded as an SDF over the parametric domain is a function  $\mathbb{R}^2 \rightarrow \mathbb{R}$ , that maps the same  $(u, v)$  coordinates to a value  $(d)$  measuring the signed distance to the boundary. We combine these two function to learn a single field that maps  $\mathbb{R}^2 \rightarrow \mathbb{R}^4$ . We parameterize this field as a 4-layer fully-connected ReLU network where the input coordinates and conditioning vector are concatenated to the output at every hidden layer (similar to DeepSDF [29]). The conditioning vector is the output of our face encoder described below. We call the evaluation of this field “rasterization” because raster sampling the field and filtering by  $d$  yields a 3D surface rasterization.

To normalize the field input range, and constrain the uv-space we must sample while rasterizing, we reparameterize the uv-space of each surface prior to training so that the clipping mask fits snugly within the unit square  $([0, 1]^2)$ . This way, our encoder-decoder is learning the explicit surface, the implicit boundary, and the support range of the implicit boundary mask.

**Encoder** Our encoder design is driven by the same observation that B-Reps are compositions of explicitly represented surfaces with implicitly represented boundaries. For the encoder, this means that we must capture adjacent topological entities. We propose to encode this using message passing on the topological graph.

As show in Figure 3, the B-rep topological graph has a hierarchical structure; high dimensional topological elements are only adjacent to the immediate lower dimensional entities: faces are adjacent to their bounding edges, which are adjacent to their end-point vertices. Driven by this observation, use as our face-encoder a hierarchical message passing network inspired by the Structured B-Rep GCN (SB-GCN) [16]: a graph convolutional network that leverages the hierarchical structure of B-reps to learn embeddings for each topological entity.

From SB-GCN, we adopt the idea of passing messages up from bounding entities to those they bound (vertices to edges and edges to faces). Since we are only attempting to learn local features, we stop message passing at this stage. Unlike SB-GCN, we only use 3 layers of topological enti-

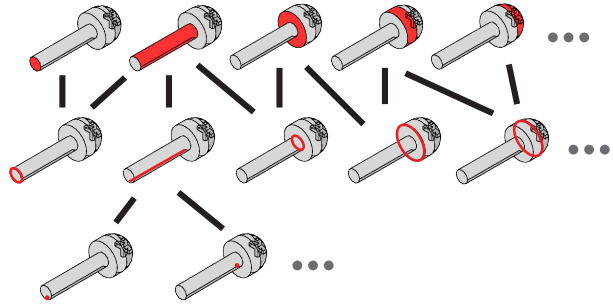


Figure 3. A boundary representation is made of parametric surface patches, bounded by parametric edges, bounded by vertices. The relationship between these entities forms a hierarchical graph.

ties – vertices, edges, and faces – avoiding an intermediate *loop* level that aggregates edges into closed curves labeled by whether they are interior or exterior face boundaries. To achieve a similar goal in our architecture we use multi-headed graph attention to handle large and variable numbers of B-Rep graph neighbors, and graph-edge features indicating the type of relationship between B-Rep topological entities (is a face to the left or right on a B-Rep edge, and is a vertex the start or end of a B-Rep edge). [39]

Finally, SB-GCN (and other B-Rep representation learning networks) use evaluated geometry information as input features in addition to the pure parametric geometry, such as surface bounding box and area. These features are computed by a CAD kernel: modeling software that understands how to construct and evaluate CAD geometry. Since we want to force our network to learn the evaluation function of a CAD kernel, as well as be fully differentiable back to the shape parameters, we only use the parametric geometry definitions as input features. We encode each entity with attached geometry (vertices, edges, and faces) as a one-hot encoding of the type of surface or curve concatenated with the continuous parameters of the function (for example, radius, center point, and axis for a cylinder), and a flag indicating orientation of the entity relative to the geometry (is the curve parameterization or surface normal reversed). Geometric functions with fewer parameters are zero-padded to create a fixed size input vector for each topological entity. In order to have a fixed input size, we limit our input to B-Reps that have geometry with a fixed number of parameters: planes, cylinders, cones, spheres, and tori for surfaces, and lines, circles, and ellipses for edges. We validated this choice by filtering the parts in the Automate data set [16], and found that 72.16% contained only these primitives.

**Training** We trained our encoder-decoder to minimize  $L^2$  losses at randomly sampled uv-points on each face’s supporting surface. We used a biased random sampling in the re-parameterized uv-square  $(u, v) \in [-0.1, 1.1]^2$  to ensure

positive and negative SDF samples were collected. CAD kernels support efficiently querying if a point in a surface parameterization is inside or outside the clipping function, so we approximate the SDF by sampling a large number of uv-points (5000) and matching each inside and outside point to its nearest neighbor in the opposite set, using a KD-Tree to accelerate these queries. We keep  $N = 500$  of these points on each face for training. To bias our sample towards the 0-level set, as is common when training implicit SDFs [29], we sort by  $|d|$  and task 40% of our sample to be the nearest points to the boundary, and randomly sample the rest.

We trained our geometric self-supervision over 23266 shapes from the training set of Fusion 360 Gallery Segmentation Dataset [36]. Optimization was performed using the Adam optimizer [17] with a learning rate of .0005. For our experiments, we used an embedding size of 64 for all message passing layers, 2 attention heads for the vertex to edge layer, 16 attention heads for the edge to face layer, and a hidden size of 1024 for all self-supervision decoder layers.

## 4. Few-Shot Learning

Since in typical task-specific applications it is challenging to find or construct large collections of labeled CAD data, we propose leveraging our rich latent space to enable supervised learning over very small collections. In addition to enabling training on very little data our approach also ensures that we do not need to use computationally expensive CAD kernel functionality for generating the features at inference time.

We propose two few-shot learning setups. The first are B-Rep segmentation tasks, which assign a task-specific label to each face of a B-Rep; for example, what type of machining technique can be used to construct that surface. The second are B-Rep classification tasks, which assign a label to an entire B-Rep; for example the mechanical function of that part (gear, bracket, etc.) We frame both of these as a multi-class classification problems.

**B-Rep Segmentation Network** Since our encodings are learned at the face level, they can be used directly as input for face level predictions. Because classification of a face is often dependent on its context within a part, we use a small graph convolutional network to capture this context; a 2-layer Residual MR-GCN. [23] We use pre-computed face embeddings from our geometric self-supervised learning as node features, and use face-face adjacency as edges; we construct this graph by removing vertex nodes from the B-Rep graph and contracting edge nodes, preserving multiple edges between faces if they exist. Output predictions for each face are then made with a fully connected network with two hidden layers.

**B-Rep Classification Network** While we do not have latent codes for entire B-Rep shapes, we take a similar approach to previous B-Rep learning architectures and pool learned features for each face. To do this we project each face’s embedding into a new vector for pooling using a fully connected layer, followed by a second projection of the pooled part features into the prediction output space.

## 5. Results

We validate the application of our proposed approach by applying our method to three few-shot learning tasks. We further evaluate our method by analyzing the rasterization results.

### 5.1. Few-Shot Construction-Based Segmentation

The first task we apply our method to is segmentation of B-Rep geometry by the modeling operation used to construct each face (extrusion, revolution, chamfer, etc.). This requires the model to understand how geometry is constructed in CAD software. For this task we use 27450 parts from the the Fusion 360 Gallery dataset, a collection of user-constructed parts annotated with one of 8 face construction operations on each face [21].

We pre-trained our geometric self-supervision network over this dataset, then our face-level prediction network using the face embeddings from the self-supervision. Some classification results are shown in Figure 4. Our method is able to achieve 65% accuracy after seeing just 10 training examples, and is 96% accurate when given all 23266 in the training set.

We compare our method against two network architectures for B-Rep learning; BRepNet [21] and UV-Net [15]. BRepNet defines a convolution operator relative to co-edges in a B-Rep structure and uses parametric face and edge types, as well as concavity and edge length features. UV-Net is a message passing network between adjacent faces, and uses grid sampled points of faces and edges as its features. We trained and tested each network on random subsets of the training data ranging from 10 examples to 23266. We repeated this 10 times at each training set size (all three methods saw the same training sets).

Figure 5 plots the average accuracy at each train size across these runs with a bootstrapped 95% confidence interval. Our method outperforms the baselines at all dataset sizes, and does so by a significant margin in the few-shot regime. It also has a smaller confidence interval, indicating that pre-training on our rasterization task makes classifications more robust to choice of training example<sup>1</sup>. In addition to these quantitative results, we show classification comparisons at the 100 example level in Figure 6. Our

<sup>1</sup>BRepNet’s very small confidence interval at small training set sizes is an artifact of it predominantly or entirely predicting one class.

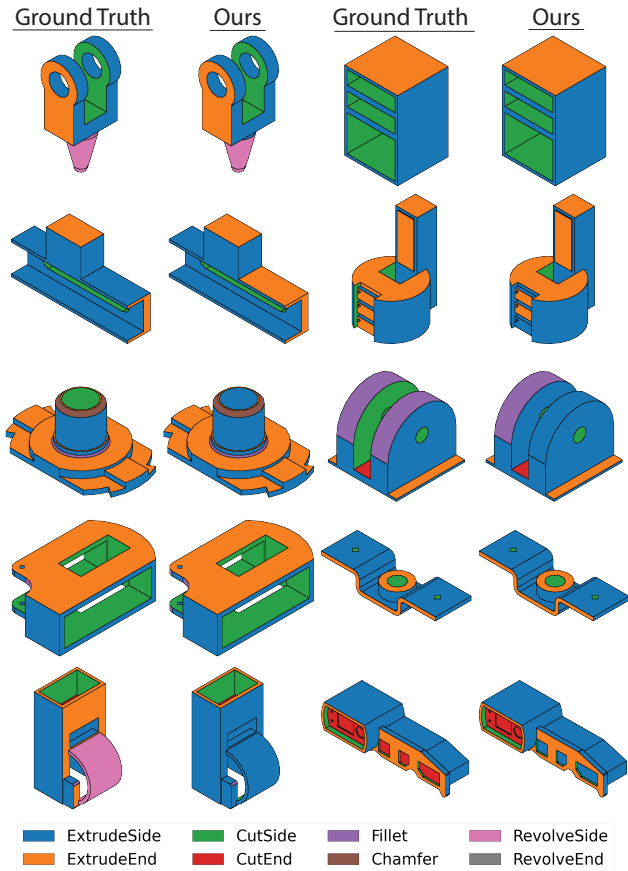


Figure 4. Face classification example results for the Fusion 360 Segmentation Task. The CAD operation which created each face is indicated by color. Left is our model’s prediction, right is the ground truth.

method is able to start generalizing with 10s to 100s of examples, whereas the baselines fail to generalize and mostly learn the most frequent label at this data scale.

We further note that training our method is an order of magnitude faster than either baseline (See Figure 7). This raises the possibility of segmentation tasks trained on-the-fly to enable predictors to be trained and deployed in the modeling process.

## 5.2. Few-Shot Manufacturing-Driven Segmentation

The second segmentation problem we considered is face segmentation that classifies how each face is manufactured. We evaluated this on the MFCAD dataset [6]. This is a synthetic dataset of 15488 CAD models where each face is labeled with one of 16 types of machining feature (chamfer, through hole, blind hole, etc.) The parts in this dataset are generated by applying machining operations to square stock, and so consist entirely of planar faces and straight

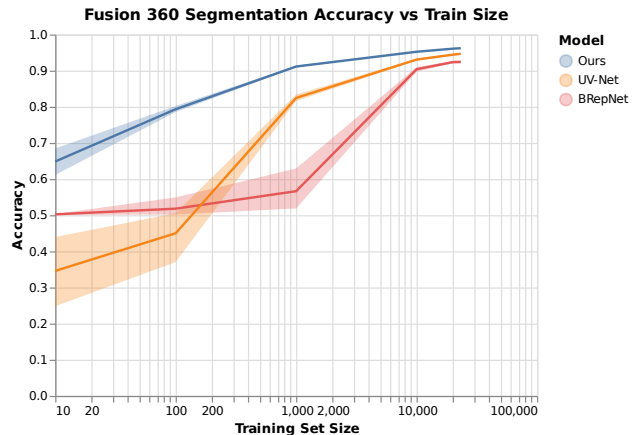


Figure 5. Comparison of our method versus UV-Net and BRepNet for Construction-Based Segmentation on the Fusion 360 Gallery Segmentation dataset. The mean accuracy across 10 training runs is plotted with a bootstrapped 95% confidence interval.

line edges. As compared to the Fusion 360 Gallery dataset, this task has twice as many classes to choose from, and will rely more on neighborhood and boundary information since all faces are planar. It also shows the ability of our approach to generalize pre-trained features to out-of-distribution data, since we use face codes pre-trained on Fusion 360 Gallery models (in this and all examples), which is more naturalistic in comparison.

We compare this task to the same two baselines, illustrated qualitatively in Figure 8 and quantitatively in Figure 9. We had to remove curvature features from BRepNet since they are universally zero on this dataset, and BRepNet requires all input features to have non-zero standard deviation. As before, our method outperforms the baselines at few-shot learning, and achieves comparable accuracy at large data sizes (all three methods are essentially perfect given sufficient training data). The stability of our prediction accuracy as measured by confidence interval is significantly better than both baselines.

## 5.3. Few-Shot Part Classification

In addition to segmentation, we also applied our method to part classification. For this, we used the FabWave [3] dataset, a hand-labeled subset of GrabCAD [1] which has been categorized into classes of mechanical parts (gears, brackets, washers, etc.) Many parts from within a category are parametric variations of each other. This task is important for understanding the function of mechanical parts in assemblies.

As a baseline comparison we only compare against UV-Net since BRepNet does not support classification. Again using face codes pre-trained on Fusion 360 Gallery B-Reps, we train our B-Rep Classification Network over the 26

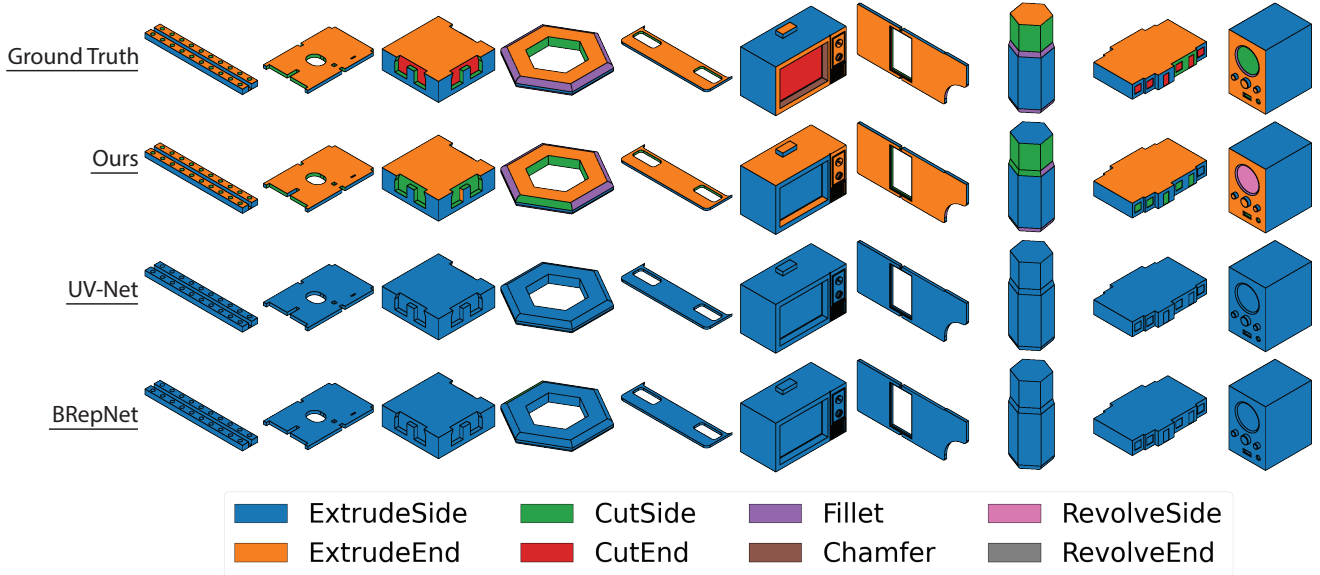


Figure 6. Few shot classification results from our method versus baseline models on the Fusion 360 gallery task, trained over just 100 models. Face color indicates the modeling operation used to create a face. Our model performs significantly better in the low data regime; it can differentiate operations, whereas the baselines are largely guessing the most common operation: ExtrudeSide.

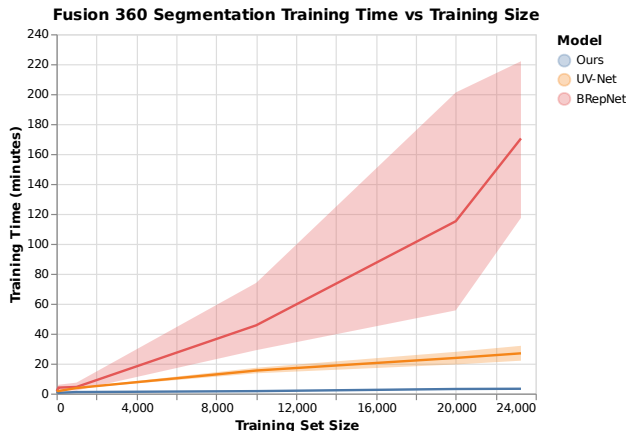


Figure 7. Training times on the Fusion 360 Segmentation task. All models were trained on an NVIDIA RTX 2080 Ti, and reported time is the time to minimum validation loss. Face codes pre-trained using a truncated SB-GCN.

classes which have at least 3 examples compatible with our network (for train, test, and validation), and use stratified sampling to create training subsets which contain at least one part from each category in each split. Since UV-Net cannot be run on B-Reps that do not contain edges, we restricted our test set to only models with edges, as well as removing 2 classes which UV-Net could not distinguish since they differ only in part orientation (our technique is able to distinguish these classes, so keeping them in the evaluation

would both improve our accuracy and decrease UV-Net’s). Figure 10 shows the results. Our method outperforms UV-Net across training sizes and even achieves 100% accuracy at higher training set sizes.

#### 5.4. Rasterization Evaluation

In Figure 11 we present a collection of part renderings to illustrate qualitatively our reconstruction and classification results. It shows a gallery of rasterization results on unseen test examples from the Fusion 360 Gallery segmentation dataset. We see that the explicit prediction of support surface coordinates works very well. While the implicit prediction of clipping planes works well in many cases, it sometimes misses inner loops and complex boundary geometry.

#### 5.5. Ablations

**Self-Supervision Ablations** In addition to the network described in Section 3, we also tried using a truncated SB-GCN (only using its upwards pass) as the encoder network, using the same shape parameter input features. We evaluate both explicit surface and SDF accuracy in Table 1, and see that our architecture significantly outperforms SB-GCN in both measures.

**Segmentation Ablations** We tried 3 types of face-level prediction network using our self-supervised face embeddings to determine which was best. The first two try to directly classify faces from the embedding directly, one us-

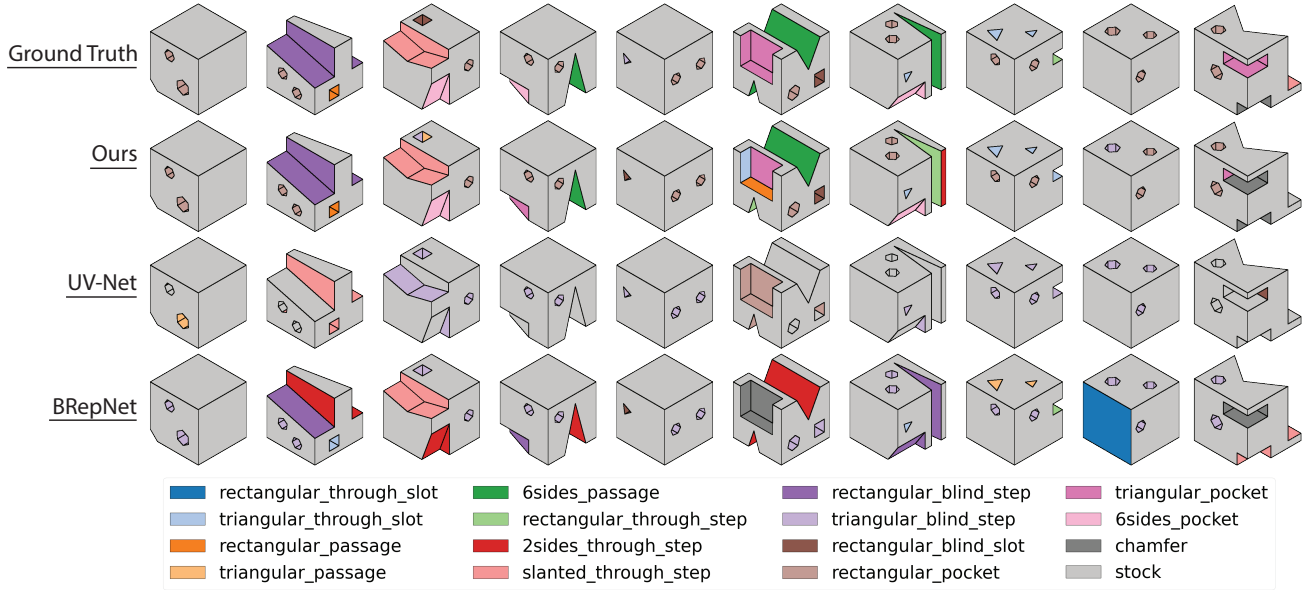


Figure 8. MFCAD few shot segmentation comparison at 100 train samples. Top row is the ground truth labeling, followed by Ours, UV-Net, and BRepNet predictions.

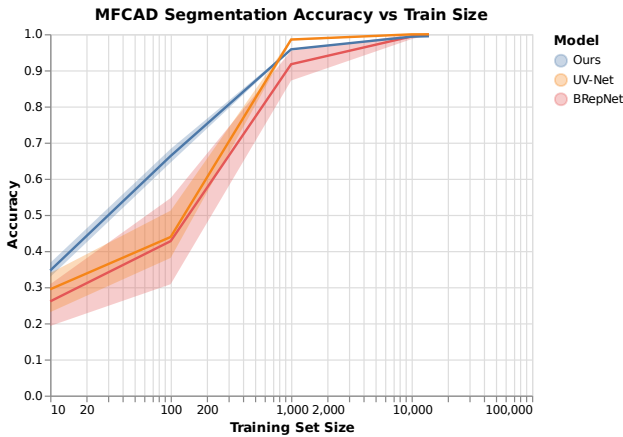


Figure 9. Comparison of our method versus UV-Net and BRepNet for Manufacturing-Based Segmentation on the MFCAD dataset. The mean accuracy across 10 training runs is plotted with a bootstrapped 95% confidence interval.

Model	XYZ Error	SDF Error
Ours	<b>0.0256</b>	<b>0.0147</b>
SB-GCN	0.035	0.0214

Table 1. Self-Supervision ablations. XYZ Error is the average pointwise distance between predicted and sampled surface position. SDF Error is the average absolute difference between predicted and actual SDF value.

ing a Linear Support Vector Machine (SVM), and the other

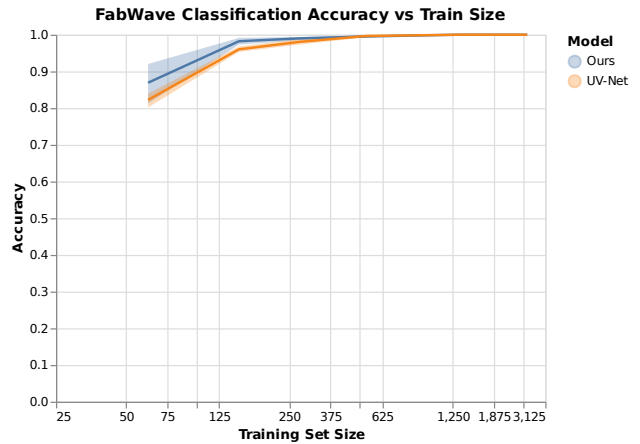


Figure 10. Comparison of our method versus UV-Net for Part Classification on the FabWave Dataset. The mean accuracy across 10 training runs is plotted with a bootstrapped 95% confidence interval.

using a Multi-Layer perceptron. These test linear and non-linear decision boundaries based on face-codes alone. The third is the message passing scheme described in Section 4, which tests if neighborhood context is necessary.

To evaluate how well each method worked across labeled dataset sizes, we trained each model repeatedly on all supported tasks at a variety of training set sizes using a similar scheme to our baseline comparisons. Table 2 records the average face segmentation accuracy across dataset sizes for both segmentation tasks. Results were consistent across



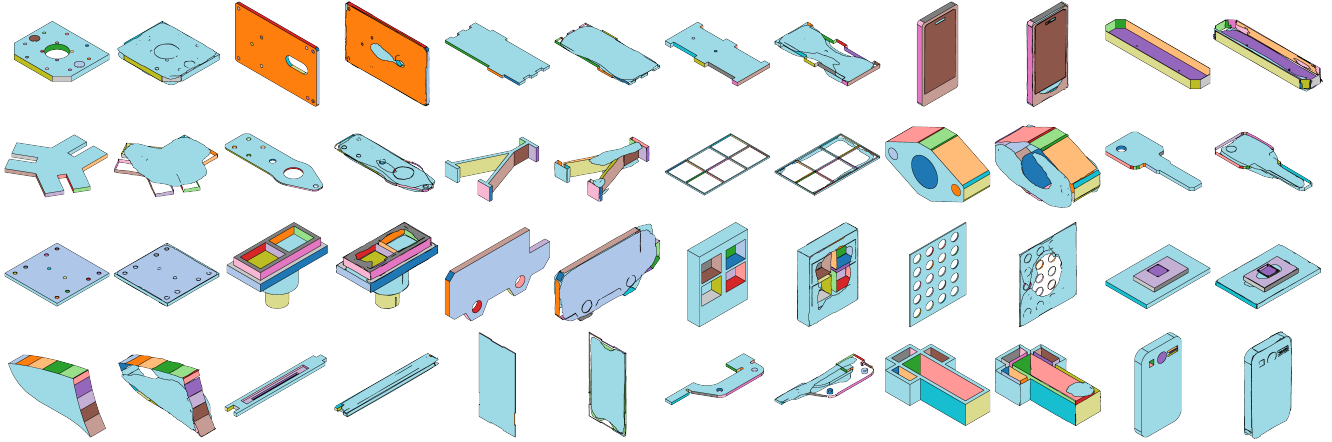


Figure 11. Shape reconstructions (right) from face embeddings on the Fusion 360 Segmentation test set, compared to ground truth (left). Each B-Rep face is given a unique color, which is consistent between ground truth and reconstruction. Reconstructions were created by sampling a  $100 \times 100$  (u,v) grid in the range  $[-0.1, 1.1]$  for each face to create a mesh for the supporting surface, then removing mesh vertices outside the predicted clipping plane SDF. Supporting surface reconstructions are highly accurate, as are the clipping masks for typical shapes with a single boundary. Complex interior and exterior boundaries of clipping planes are sometimes predicted inaccurately.

datasets and training sizes; a non-linear decision boundary is more accurate than a linear one, and neighborhood information improves accuracy, leading us to choose the message passing network as our method.

**Classification Ablations** We also tested adding additional message passing layers prior to pooling for the classification network. As Table 3 shows, this additional complexity did not yield real improvement on

## 5.6. Limitations and Future Work

Our work has three main limitations. The first is the restriction of B-Reps to those with fixed-size parameter sets for all attached geometry. This is a limitation of the choice of encoder network, and could be alleviated by using a sequence or tree encoder to compute fixed-size embeddings for the generic functional geometric expressions of each B-Rep topology.

The second limitation is that we only self-supervise on local information. This means that we rely on additional message passing layers in our task specific decoders to gather neighborhood features.

Finally, we only create encodings for faces, which limits the kinds of tasks we can learn. Adding a second decoder and loss term to rasterize edges could extend this work to edge-based tasks, however there are not currently edge-specific tasks in the literature to compare against, and we encountered frequently CAD kernel errors when sampling edges, so we did not take this step.

Improving self-supervision accuracy will be a fruitful direction for future work. Figure 12 shows that the accuracy of downstream learning tasks is correlated with the accuracy

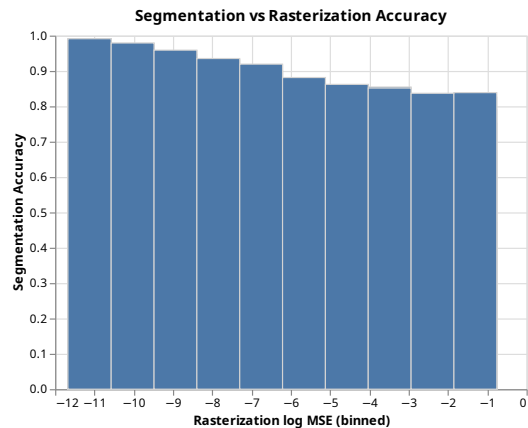


Figure 12. Segmentation accuracy compared to rasterization accuracy on the Fusion 360 segmentation task. X-axis is the log of the per-face MSE of our rasterization, binned into 10 groups, and the Y-axis is the fraction of faces segmented accurately within each bin. Data is aggregated across all segmentation training set sizes and seeds (the trend is similar across all sizes, with lower accuracies at lower training set sizes).

of our model’s rasterization. Thus improvements in rasterization performance should yield improvements in classification performance in both few-shot and large data regimes.

Improving rasterization performance could also unlock future applications in reverse engineering. Since operating our self-supervision network as a rasterizer creates, in effect, a differentiable CAD renderer, we can use it for gradient-based optimization of B-Rep *shape parameters*. To demonstrate the potential of such an application, we prototyped a shape matching application, where an input B-

Task / Model		Training Set Size					
Fusion 360 Segmentation	Accuracy@:	10	100	1000	10000	20000	23266
Self-Supervision + SVM		0.50	0.66	0.75	0.76	0.77	0.77
Self-Supervision + MLP		0.55	0.67	0.82	0.90	0.91	0.91
Self-Supervision + MP		<b>0.56</b>	<b>0.68</b>	<b>0.83</b>	<b>0.92</b>	<b>0.94</b>	<b>0.94</b>
MFCAD	Accuracy@:	10	100	1000	10000	13940	--
Self-Supervision + SVM		0.15	0.50	0.56	0.59	0.58	
Self-Supervision + MLP		0.35	0.55	0.83	0.90	0.91	
Self-Supervision + MP		<b>0.38</b>	<b>0.65</b>	<b>0.95</b>	<b>0.99</b>	<b>0.99</b>	

Table 2. Segmentation ablations. Reported face classification accuracies are the mean of 10 runs at each data set size with the train set subset at different random seeds (each model sees the same 10 random subsets). Model selected by best validation loss on a random 20% validation split, except for the SVM models. Bold indicates indicates the best accuracy at each train size for each task. Self-Supervision codes pre-trained using a truncated SB-GCN.

Task / Model		Training Set Fraction						
FabWave	Accuracy@:	1%	5%	10%	20%	50%	75%	100%
Self-Supervision + Pooling		<b>0.72</b>	<b>0.93</b>	0.97	<b>0.99</b>	<b>0.999</b>	.999	<b>1.00</b>
Self-Supervision + MP + Pooling		0.72	0.93	<b>0.97</b>	0.99	0.999	<b>1.00</b>	<b>1.00</b>

Table 3. Classification ablations. Reported accuracies are mean of 10 runs, similar to Table 2. Adding message passing prior to pooling does not confer an advantage, so we do not use it in our reported results. Self-Supervision codes were pre-trained using a truncated SB-GCN.

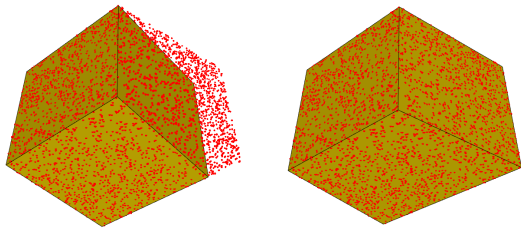


Figure 13. Shape matching with differentiable rasterization. The input cube (left) was optimized using SGD over its B-Rep shape parameters (right) to match a target point cloud (red).

Rep is optimized via stochastic gradient descent to match a target point cloud. Figure 13 shows the results of using this optimization to angle the face of a cube. We find that this optimization technique struggles on more complex shapes, which may overcome by improved rasterization performance.

## 6. Conclusion

In this work we propose to learn a spatial embedding of B-reps and apply it to few-shot learning on supervised tasks. We validate that this approach is effective compared to prior work on supervised learning. Our results show comparable results on large training sets and significantly better performance on smaller sets. By experimenting over three different tasks and data-sets, we posit that this method will

be widely applicable on a plethora of CAD applications. Being faster to train can enable many applications in this domains; particularly, user-guided annotation for customized predictions. Finally, our method enables a fully differential embedding of B-rep geometry, compared to prior work that required non-differentiable CAD kernels, paving the way to exciting future work on CAD optimization and reverse engineering.

## References

- [1] Grabcad. <https://grabcad.com/>. Accessed: 2022-05-19. 3, 6
- [2] Siemens. parasolid cad kernel. <https://www.plm.automation.siemens.com/global/en/products/plm-components/parasolid.html>. Accessed: 2022-05-19. 3
- [3] *Development of a Pilot Manufacturing Cyberinfrastructure With an Information Rich Mechanical CAD 3D Model Repository*, volume Volume 1: Additive Manufacturing; Manufacturing Equipment and Systems; Bio and Sustainable Manufacturing of *International Manufacturing Science and Engineering Conference*, 06 2019. V001T02A035. 2, 6
- [4] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas J Guibas. Learning representations and generative models for 3d point clouds. *arXiv preprint arXiv:1707.02392*, 2017. 3
- [5] André Brock, Theodore Lim, James M. Ritchie, and Nick Weston. Generative and discriminative voxel modeling

- with convolutional neural networks. *CoRR*, abs/1608.04236, 2016. 3
- [6] Weijuan Cao, Trevor Robinson, Yang Hua, Flavien Bousuge, Andrew R. Colligan, and Wanbin Pan. Graph representation of 3d cad models for machining feature recognition with deep learning. volume Volume 11A: 46th Design Automation Conference (DAC) of *International Design Engineering Technical Conferences and Computers and Information in Engineering Conference*, 08 2020. 2, 6
- [7] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. *IEEE Computer Vision and Pattern Recognition (CVPR)*, 2019. 3
- [8] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1422–1430, 2015. 3
- [9] Haoqiang Fan, Hao Su, and Leonidas Guibas. A point set generation network for 3d object reconstruction from a single image. *CVPR*, 2017. 3
- [10] Yaroslav Ganin, Sergey Bartunov, Yujia Li, Ethan Keller, and Stefano Saliceti. Computer-aided design as language. *Advances in Neural Information Processing Systems*, 34, 2021. 2
- [11] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. In *ICLR 2018*, 2018. 3
- [12] Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. Atlasnet: A papier-mache approach to learning 3d surface generation. *arXiv preprint arXiv:1802.05384*, 2018. 3
- [13] Hao-Xiang Guo, Shilin Liu, Hao Pan, Liu Yang, Xin Tong, and Baining Guo. Complexgen: Cad reconstruction by b-rep chain complex generation. *ACM Transactions on Graphics (TOG)*, 39(4):106:1–106:14, 2022. 2
- [14] Pradeep Kumar Jayaraman, Joseph G Lambourne, Nishkrit Desai, Karl DD Willis, Aditya Sanghi, and Nigel JW Morris. Solidgen: An autoregressive model for direct b-rep synthesis. *arXiv preprint arXiv:2203.13944*, 2022. 2
- [15] Pradeep Kumar Jayaraman, Aditya Sanghi, Joseph G Lambourne, Karl DD Willis, Thomas Davies, Hooman Shayani, and Nigel Morris. Uv-net: Learning from boundary representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11703–11712, 2021. 3, 5
- [16] Benjamin Jones, Dalton Hildreth, Duowen Chen, Ilya Baran, Vladimir G. Kim, and Adriana Schulz. Automate: A dataset and learning approach for automatic mating of cad assemblies. *ACM Transactions on Graphics*, 40(6), dec 2021. 2, 3, 4
- [17] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [18] Durk P Kingma, Shakir Mohamed, Danilo Jimenez Rezende, and Max Welling. Semi-supervised learning with deep generative models. In *Advances in neural information processing systems*, pages 3581–3589, 2014. 3
- [19] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 3
- [20] Sebastian Koch, Albert Matveev, Zhongshi Jiang, Francis Williams, Alexey Artemov, Evgeny Burnaev, Marc Alexa, Denis Zorin, and Daniele Panozzo. ABC: A Big CAD Model Dataset for Geometric Deep Learning. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9593–9603, Long Beach, CA, USA, June 2019. IEEE. 2, 3
- [21] Joseph G. Lambourne, Karl D. D. Willis, Pradeep Kumar Jayaraman, Aditya Sanghi, Peter Meltzer, and Hooman Shayani. BRepNet: A topological message passing system for solid models. *arXiv:2104.00706 [cs]*, Apr. 2021. arXiv: 2104.00706. 2, 3, 5
- [22] Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. Learning representations for automatic colorization. In *European Conference on Computer Vision*, pages 577–593. Springer, 2016. 3
- [23] Changjian Li, Hao Pan, Adrien Bousseau, and Niloy J Mitra. Sketch2cad: Sequential cad modeling by sketching in context. *ACM Transactions on Graphics (TOG)*, 39(6):1–14, 2020. 5
- [24] Jerry Liu, Fisher Yu, and Thomas Funkhouser. Interactive 3d modeling with a generative adversarial network. *International Conference on 3D Vision (3DV)*, 2017. 3
- [25] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4460–4470, 2019. 3
- [26] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European Conference on Computer Vision*, pages 69–84. Springer, 2016. 3
- [27] David Palmer, Dmitriy Smirnov, Stephanie Wang, Albert Chern, and Justin Solomon. DeepCurrents: Learning implicit representations of shapes with boundaries. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [28] Wamiq Para, Shariq Bhat, Paul Guerrero, Tom Kelly, Niloy Mitra, Leonidas J Guibas, and Peter Wonka. Sketchgen: Generating constrained cad sketches. *Advances in Neural Information Processing Systems*, 34, 2021. 2
- [29] Jeong Joon Park, Peter Florence, Julian Straub, Richard A. Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. *CVPR*, 2019. 3, 4, 5
- [30] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2536–2544, 2016. 3
- [31] Ari Seff, Yaniv Ovadia, Wenda Zhou, and Ryan P. Adams. SketchGraphs: A Large-Scale Dataset for Modeling Relational Geometry in Computer-Aided Design. *arXiv:2007.08506 [cs, stat]*, July 2020. arXiv: 2007.08506. 2
- [32] Ari Seff, Wenda Zhou, Nick Richardson, and Ryan P Adams. Vitruvion: A generative model of parametric cad sketches. *arXiv preprint arXiv:2109.14124*, 2021. 2

- [33] Qingyang Tan, Lin Gao, Yu-Kun Lai, and Shihong Xia. Variational autoencoders for deforming 3d mesh models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 3
- [34] Yaqing Wang, Quanming Yao, James T. Kwok, and Lionel M. Ni. Generalizing from a few examples: A survey on few-shot learning. *ACM Comput. Surv.*, 53(3), jun 2020. 3
- [35] Karl DD Willis, Pradeep Kumar Jayaraman, Hang Chu, Yunsheng Tian, Yifei Li, Daniele Grandi, Aditya Sanghi, Linh Tran, Joseph G Lambourne, Armando Solar-Lezama, et al. Joinable: Learning bottom-up assembly of parametric cad joints. *arXiv preprint arXiv:2111.12772*, 2021. 2
- [36] Karl D. D. Willis, Yewen Pu, Jieliang Luo, Hang Chu, Tao Du, Joseph G. Lambourne, Armando Solar-Lezama, and Wojciech Matusik. Fusion 360 Gallery: A Dataset and Environment for Programmatic CAD Reconstruction. *arXiv:2010.02392 [cs]*, Oct. 2020. arXiv: 2010.02392. 2, 3, 5
- [37] Rundi Wu, Chang Xiao, and Changxi Zheng. Deepcad: A deep generative network for computer-aided design models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6772–6782, 2021. 2
- [38] Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. Foldingnet: Point cloud auto-encoder via deep grid deformation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 3, 2018. 3
- [39] Jiaxuan You, Rex Ying, and Jure Leskovec. Design space for graph neural networks. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS’20, Red Hook, NY, USA, 2020. Curran Associates Inc. 4
- [40] Richard Zhang, Phillip Isola, and Alexei A Efros. Colorful image colorization. In *European conference on computer vision*, pages 649–666. Springer, 2016. 3