

# REGULARIZED DIFFUSION MODELING FOR CAD REPRESENTATION GENERATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Computer-Aided Design (CAD) has significant practical value in various industrial applications. However, achieving high-quality and diverse shape generation, as well as flexible conditional control, remains a challenge in the field of CAD model generation. To address these issues, we propose CADiffusion, a diffusion-based generative model with a hierarchical latent representation tailored to the complexities of CAD design processes. To enhance the performance and reliability of the model in generating accurate CAD models, we have developed a specialized decoder with regularization strategies that navigate through the noise space of the diffusion model, smoothing the results. This approach not only improves the diversity and quality of the generated CAD models but also enhances their practical applicability, marking a significant advancement in the integration of generative models and automated CAD systems.

## 1 INTRODUCTION

In the field of 3D computer vision, exploring the generation of 3D shapes has emerged as a prominent issue in recent times. Various research studies have been conducted that encompass explicit representations such as point clouds Yang et al. (2019); Cai et al. (2020); Mo et al. (2019), polygon meshes Groueix et al. (2018); Wang et al. (2018); Nash et al. (2020), voxel grids Liao et al. (2018); Li et al. (2017), and implicit representations Park et al. (2019); Mescheder et al. (2019); Chen et al. (2020) such as neural radiance fields (NeRFs) and signed distance functions (SDFs). Despite the significant progress achieved by these methodologies, the effectiveness of generating 3D models for practical applications remains unsatisfactory, mainly due to constraints related to data availability and modeling capabilities of the models. In contrast to other forms of 3D data, Computer-Aided Design (CAD) data holds immense practical value and finds applications across various industrial domains, ranging from automotive and aerospace to manufacturing and architectural design. CAD software serves as the cornerstone for creating 3D shapes in these domains, facilitating intricate design processes, and streamlining manufacturing workflows. Therefore, exploring the generation of CAD models is highly significant as it has the potential to innovate many existing production processes. Our work primarily explores representations and the corresponding generative model structures that are more suitable for CAD generation.

Despite some existing research on the generation of CAD data, challenges persist to enhance the quality and diversity of the generated shapes. Current methodologies face challenges in generating highly diverse CAD shapes, primarily due to limitations in supporting the intricacies of CAD spatial modeling. DeepCAD Wu et al. (2021) generates CAD commands directly without modeling a space suitable for CAD generation. Moreover, the exploration of generative models remains limited to those such as GANs and autoregressive transformers Xu et al. (2022; 2023). Beyond the issues of quality and diversity, existing methods also struggle to ensure the consistency and realism of the generated results with the input conditions in conditional CAD generation.

To address these challenges, we introduce a hierarchical implicit space and, on top of it, we propose CADiffusion, a diffusion-based generative model with a latent representation structured in a tree logic. Most modern CAD design tools employ a “Sketch and Extrude” style workflow, where designers first draw loops of 2D curves as outer and inner boundaries to create 2D profiles, then extrude the 2D profiles to 3D shapes, and finally add or subtract 3D shapes to build complex CAD models. Therefore, a hierarchical representation perfectly aligns with the inherent logic of CAD itself. This

054 hierarchical structure also offers effective design control at different levels of the hierarchy. When  
055 modeling the latent space for CAD, we also need to consider the external, visible compositional logic  
056 of CAD. In our approach, the CAD data is divided into three distinct hierarchical levels, with each  
057 level employing a VAE to obtain latent representations. These representations are then organized into  
058 a tree logic latent structure.

059 To learn the probability distribution in the proposed latent space, we leverage diffusion models, which  
060 have recently achieved significant success in various 2D generation tasks. We find that diffusion  
061 modeling method also exhibits high-quality and highly diverse generation capabilities in CAD data.  
062 After using the diffusion model to fit the CAD data, to convert the CAD latent generated by the  
063 diffusion model into CAD models accurately and efficiently, we designed a corresponding decoder  
064 along with regularization terms. This involves navigating through the diffusion model’s noise space  
065 and smoothing the outcomes of the sampled noise. This process ensures that the decoder can translate  
066 any sample from the Gaussian noise space into a reasonable CAD model. The regularization approach  
067 not only enhances the decoder’s ability to handle variations but also contributes to the overall stability  
068 and reliability of CAD model generation in an automated setting. By integrating this regularized  
069 training methodology, we can bridge the gap between generating realistic data and the latent diffusion  
070 modeling, ensuring that the enhancements in generative model technology translate effectively into  
071 practical improvements in CAD systems. We evaluate the effectiveness of CADiffusion on benchmark  
072 datasets and show that it outperforms baseline approaches in a variety of metrics.

073 Therefore, our contributions can be summarized as follows: **1)** We are the first to explore the use of  
074 diffusion models for CAD generation, and have designed corresponding models and representations.  
075 **2)** We have introduced a new regularization strategy specifically for CAD latent diffusion models,  
076 enabling the decoder to produce more reasonable and higher quality results. **3)** Our CAD generation  
077 model achieves state-of-the-art performance, surpassing previous methods.

## 078 2 RELATED WORK

### 079 2.1 3D GENERATIVE MODELS

082 In recent years, significant attention has been directed towards the development of generative models  
083 for 3D shapes. Many existing approaches generate 3D shapes discretely, employing representations  
084 such as voxelized shapes Liao et al. (2018); Li et al. (2017), point clouds Yang et al. (2019); Cai  
085 et al. (2020); Mo et al. (2019), polygon meshes Groueix et al. (2018); Wang et al. (2018); Nash et al.  
086 (2020), and implicit signed distance fields Park et al. (2019); Mescheder et al. (2019); Chen et al.  
087 (2020). Despite their prevalence, these models often produce shapes with noise, limited geometric  
088 sharpness, and lack direct user editability. To address these limitations, recent research has focused  
089 on neural network architectures that generate 3D shapes through sequences of geometric operations.  
090 CSGNet Sharma et al. (2018), for instance, infers Constructive Solid Geometry (CSG) operations  
091 from voxelized shape inputs, while UCSG-Net Kania et al. (2020) enhances the inference process  
092 without relying on ground truth CSG trees. In addition, some approaches use domain-specific  
093 languages (DSLs) Mo et al. (2019); Jones et al. (2020) to synthesize 3D shapes. For instance,  
094 ShapeAssembly by Jones *et al.* Jones et al. (2020). introduces a DSL that constructs 3D shapes  
095 using hierarchical and symmetrical cuboid proxies, which can be generated through a variational  
096 autoencoder.

### 097 2.2 CAD GENERATION

099 Early research focused on direct CAD modeling without using any supervision from CAD modeling  
100 sequences. A common theme is to construct parametric curves Wang et al. (2020) and surfaces Sharma  
101 et al. (2020) with fixed Smirnov et al. (2021) or arbitrary topology for sketches Willis et al. (2021c)  
102 and solid models Wang et al. (2022); Guo et al. (2022); Jayaraman et al. (2022). In recent years,  
103 the availability of large-scale parametric CAD datasets has allowed learning-based methods to take  
104 advantage of data from CAD modeling sequences Willis et al. (2021b); Wu et al. (2021); Xu et al.  
105 (2022) and sketch constraints Seff et al. (2020) to generate engineering sketches and solid models.  
106 The resulting sequences can be processed using a solid modeling kernel to acquire editable parametric  
107 CAD files containing 2D engineering sketches Willis et al. (2021c); Para et al. (2021); Ganin et al.  
(2021); Seff et al. (2022) or 3D CAD shapes Wu et al. (2021); Xu et al. (2022). Furthermore, the

generation process may be influenced by the target B-rep Willis et al. (2021b); Xu et al. (2021), sketches Li et al. (2020); Seff et al. (2022), images Ganin et al. (2021), voxel grids Lambourne et al. (2022), or point clouds Uy et al. (2021), occasionally with sequence guidance Ren et al. (2022); Li et al. (2023).

More recently, there have been some advancements in the field of CAD model generation. DeepCAD Wu et al. (2021) directly generates CAD commands without modeling the data representations first. CAD models are graphically and geometrically complex, and generating commands can lead to overly simplified results. Consequently, the outcomes from DeepCAD are generally quite simple. Subsequent efforts, such as those by SkexGen Xu et al. (2022) and HNC Xu et al. (2023), have employed autoregressive transformer models. Their use of discrete formats excessively compresses the representations, failing to effectively capture the intrinsic logic of CAD data. The representations utilized by these methods are either too redundant or semantically sparse, which impairs the generative model’s performance in fitting them. Compared to other 3D data, CAD possesses parametric characteristics, and suitable representations and models for it are still under exploration.

### 2.3 DIFFUSION MODELS

Diffusion Probabilistic Models (DPMs) Sohl-Dickstein et al. (2015); Ho et al. (2020), commonly referred to as diffusion models, have emerged as a robust class of generative models. Unlike previous leading generative models such as the Generative Adversarial Network Goodfellow et al. (2020), Variational Autoencoder (VAE) Kingma & Welling (2014), and flow-based generative models Rezende & Mohamed (2015), diffusion models exhibit notable advantages in terms of training stability and generative diversity Croitoru et al. (2023). They have shown promising performance in image Ho et al. (2020); Dhariwal & Nichol (2021); Nichol et al. (2022); Rombach et al. (2022) and speech Chen et al. (2021); Kong et al. (2021) synthesis. In particular, diffusion model-based approaches have shown remarkable results in text-to-image synthesis Ramesh et al. (2022); Rombach et al. (2022); Saharia et al. (2022). In the realm of 3D computer vision, several studies have embraced diffusion models for generative 3D modeling Luo & Hu (2021); Zhou et al. (2021); Zeng et al. (2022). For example, PVD Zhou et al. (2021) used diffusion models to create 3D shapes using a point-voxel 3D representation. Luo *et al.* Luo & Hu (2021) considered points in point clouds as particles within a thermodynamic system with a heat bath. LION Zeng et al. (2022) introduced a VAE framework with hierarchical diffusion models in latent space. Similar attempts Chou et al. (2023); Cheng et al. (2023) have also applied diffusion models to the generation of SDFs. However, no exploration of diffusion models has been made on CAD data. We are the first to attempt using diffusion models to generate CAD data and have achieved very promising results.

## 3 METHOD

### 3.1 PRELIMINARIES: HIERARCHICAL CAD REPRESENTATION

CAD (Computer-Aided Design) models are inherently hierarchical because of the nature of the objects they represent. This hierarchical structure is essential to accurately represent and manipulate engineering designs, mechanical components, and architectural plans. Thus, in our approach, we employ a hierarchical representation for CAD that builds on the foundations laid by SkexGen Xu et al. (2022) and HNC Xu et al. (2023), which themselves are extensions of the pioneering work of TurtleGen Willis et al. (2021a) and DeepCAD Wu et al. (2021).

**CAD Representation:** Similar to HNC Xu et al. (2023), we conceptualize a CAD model as a tree, where it is organized into three levels: Solid, Profile, and Loop. At the lowest level, a “loop” represents the basic connected curve in the model. It is composed of a set of lines, arcs, and circles. Each such primitive is defined by two, three, or four xy-coordinates,  $L = \{(x_1, y_1), (x_2, y_2), \langle \text{SEP} \rangle, (x_3, y_3), \dots\}$ . Moving up the hierarchy, a “profile” defines a closed area on a sketch plane. It is constructed from a group of 2D bounding boxes. Each bounding box encompasses multiple loop elements that form part of the sketch,  $P = \{(x_i, y_i, w_i, h_i)\}_{i=1}^{N_i^{\text{loop}}}$ .  $(x_i, y_i)$  is the bottom-left corner of the bounding box.  $(w_i, h_i)$  is the width and height. Finally, at the top level, a “solid” represents a set of extruded profiles.

These extruded profiles are combined to form the entire 3D model,  $S = \{(x_j, y_j, z_j, w_j, h_j, d_j)\}_{j=1}^{N_j^{\text{profile}}}$ . The solid is described by a set of 3D bounding box parameters, providing a comprehensive representation of the volumetric aspects of the model.  $j$  is the index of the  $N_j^{\text{profile}}$  extruded profiles within a model.  $(x_j, y_j, z_j)$  is the bottom-left corner of the bounding box and  $(w_j, h_j, d_j)$  is its dimension.

**Hierarchical Latent Representation:** We use an adaptation of vector quantized VAE (VQ-VAE) van den Oord et al. (2017) consisting of a transformer encoder  $E$  and a decoder  $D$  to analyze and compress the CAD dataset. The dataset comprising sketch and extrude CAD models organized in a (S)olid-(P)rofile-(L)oop tree structure. This approach learns the distinct patterns inherent in the models by employing three discrete codebooks. Although the formats vary at each level, we have uniformly set the latent code length to 256 dimensions for ease of subsequent processing. Thanks to the powerful generative capabilities of our diffusion process, unlike HNC, we only need to use the simplest VQ-VAE without residual connections for modeling. Similar to hierarchical CAD representation, hierarchical codes are represented as a series of feature vectors, where each feature vector indicates a code or a separator token. We organize latent codes of three different levels into a tree structure, which is represented as follows: [S, ⟨SEP⟩, P, L, L, ⟨SEP⟩, P, L, L, L, L, ⟨END⟩], where uppercase letters represent latent codes of the corresponding levels. The boundary command ⟨SEP⟩ indicates a new grouping of profile and loop codes, ⟨END⟩ indicates the end of the data. We pad zeros after the ⟨END⟩ to unify the length of different CAD models and then form two-dimensional tensors.

### 3.2 DIFFUSION MODEL FOR CAD

**Diffusion models (DMs):** DMs learn a specific distribution by iteratively denoising a Gaussian variable through a fixed-length Markov chain, denoted as  $T$ . Specifically, given a data sample  $x_0$  drawn from the distribution  $q(x_0)$ , two distinct processes are defined: a forward process  $q(x_{0:T})$ , which progressively transforms a data sample into Gaussian noise, and a reverse process (generation process)  $p_\theta(x_{0:T})$ , which gradually denoises the Gaussian noise back into the real data.

$$q(x_{0:T}) = q(x_0) \prod_{t=1}^T q(x_t | x_{t-1}), \quad p_\theta(x_{0:T}) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1} | x_t), \quad (1)$$

both  $q(x_t | x_{t-1})$  and  $p_\theta(x_{t-1} | x_t)$  represent Gaussian transition probabilities formulated as

$$q(x_t | x_{t-1}) = \mathcal{N}\left(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t \mathbf{I}\right), \quad p_\theta(x_{t-1} | x_t) = \mathcal{N}\left(x_{t-1}; \mu_\theta(x_t, t), \beta_t \mathbf{I}\right). \quad (2)$$

The mean variable  $\mu_\theta(x_t, t)$  for the reverse transition  $p_\theta(x_{t-1} | x_t)$  can be represented as:

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right), \quad (3)$$

where  $\alpha_t = 1 - \beta_t$ ,  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$ , and  $\beta_t$  gradually decreases to 0 as  $t$  approaches 0. During the training stage of DMs, the evidence lower bound (ELBO) is maximized, eventually yielding the loss

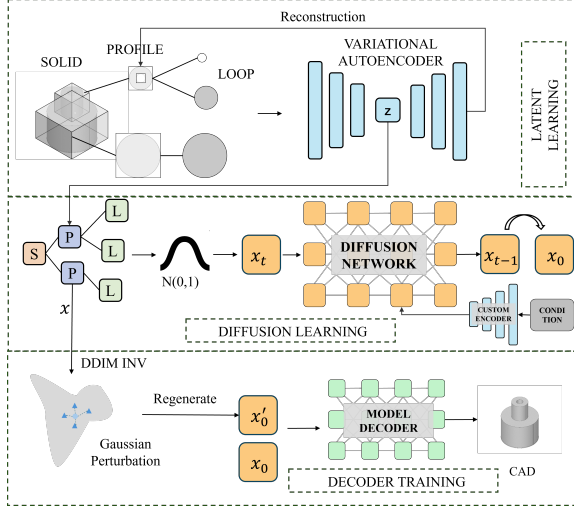


Figure 1: **Overview of CADiffusion.** We employ three distinct VQ-VAE models to perform data compression (**top**). Building upon this, we represent CAD data as corresponding tree-structured latent representations, which serve as input to our diffusion model and can be guided by various inputs (**middle**). To obtain a decoder suitable for CAD decoding, we designed specialized regularization strategies to ensure that samples from the Gaussian space generate reasonable CAD models (**bottom**).

function:

$$\mathcal{L}_{DM} = \mathbb{E}_{\mathbf{x}, t, \epsilon \sim \mathcal{N}(0,1)} \left[ \|\epsilon - \epsilon_{\theta}(\mathbf{x}_t, t)\|^2 \right] \quad (4)$$

In the process,  $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$ ,  $\epsilon$  represents a noise variable and  $t$  is uniformly sampled from the set  $\{1, \dots, T\}$ . The key component of denoising diffusion models is the neural network-based score estimator  $\epsilon_{\theta}(\mathbf{x}_t, t)$ , which serves as a time-step-conditioned denoising model.

**Diffusion Model for CAD:** Based on the description above, we organize tree logic latent representation into a 2-D tensor  $z$  corresponding to CAD. Since we have already captured the organizational logic of CAD in the tree-latent space, we aim to use the diffusion model to fit the data on a holistic level and generate coherent CAD models. Specifically, we obtain  $z_t, t \in \{1, \dots, T\}$  from a sample  $z_t$  by incrementally introducing Gaussian noise with a predetermined variance schedule. Subsequently, we employ a transformer-based time-conditional denoising model  $\epsilon_{\theta}$ . To train the denoising model, we utilize the simplified objective introduced by Ho *et al.* Ho et al. (2020) :

$$L_{\text{simple}}(\theta) := \mathbb{E}_{z, \epsilon \sim \mathcal{N}(0,1), t} \left[ \|\epsilon - \epsilon_{\theta}(z_t, t)\|^2 \right]. \quad (5)$$

During the inference phase, we generate  $\hat{z}_0$  by progressively removing noise from a variable sampled from the standard normal distribution  $\mathcal{N}(0, 1)$ .

### 3.3 CAD DECODER REGULARIZATION TERM

The unique challenges presented by our latent diffusion model, which generates hierarchical, tree logic latent representations, are presented below. Parsing generated latents into different hierarchies is prone to errors and can be overly cumbersome. Instead, the latents generated encapsulate the logic and components of CAD, enabling direct decoding into a CAD model. This approach is indeed more efficient; however, it is crucial to ensure its accuracy as well. Although the three different levels of latents is logically assembled together, decoding it into a realistic CAD model is more challenging than reconstructing a single level of CAD components separately. We found that training this decoder solely with the latents of training dataset is insufficient and leads to some unrealistic decoded results. To enhance the stability and fidelity of our CAD generation process, a novel regularization technique is developed. This technique involves perturbing the latent space to simulate variations that the diffusion model might generate, thus training the decoder to be resilient to these variations and ensuring smoother transitions between different CAD models. The regularization process consists of several steps, detailed below:

**Inverse Mapping to Noise Space:** Initially, the latent representation from dataset is mapped back to the noise space using the DDIM inversion method Song et al. (2021). The DDIM inversion process systematically reintroduces noise into a clean latent representation to reach a noised state that can then be diffused to regenerate the original latent, effectively serving as a way to explore variations in the generated CAD models. Starting from a latent at the initial time step, the DDIM inversion aims to compute a corresponding noised latent after  $T$  steps. The inversion process is governed by the following equation:

$$\hat{z}_t = \sqrt{\alpha_t} \frac{\hat{z}_{t-1} - \sqrt{1 - \alpha_{t-1}} \epsilon_{\theta}}{\sqrt{\alpha_{t-1}}} + \sqrt{1 - \alpha_t} \epsilon_{\theta}, \quad (6)$$

where  $\alpha_t$  is a pre-determined noise schedule parameter, and  $\epsilon_{\theta}$  is the neural network predicting the noise component.

**Gaussian Perturbation:** Once the DDIM inversion maps the clean latent  $z_0$  to a noised latent  $\hat{z}_T$ , we apply Gaussian perturbation as part of our regularization strategy. The perturbed noise latent  $\hat{z}'_T$  is then given by:

$$\hat{z}'_T = (1 - \sigma) \hat{z}_T + \sigma \mathcal{N}(0, I), \quad (7)$$

where  $\sigma$  is a scaling factor and  $\mathcal{N}(0, I)$  represents isotropic Gaussian noise. The perturbed noise vector is then used to regenerate a new latent vector  $\hat{z}'_0 = DDIM(\hat{z}'_T)$  through the forward diffusion process.

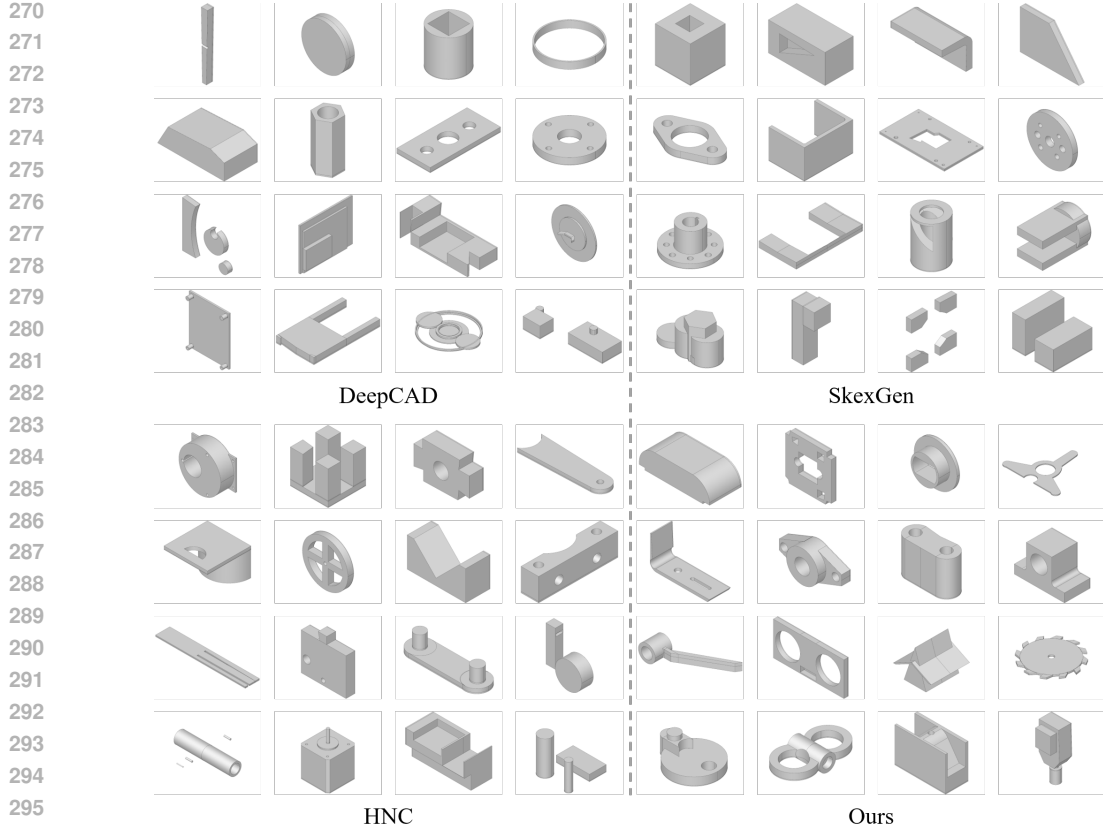


Figure 2: **Unconditional generation results of four different methods.** DeepCAD and SkexGen frequently resort to assembling simple components, resulting in CAD models that lack the rationality. While HNC’s outcomes display significant improvement, they occasionally contain artifacts with small components. In contrast, our method produces high quality results with well structure.

**Distance Minimization:** The perturbation does not deviate far enough from the original latent, we can minimize the distance between the decoded results of the perturbed latents and the original CAD:

$$\min_D \|D(\hat{z}'_0) - CAD\|. \tag{8}$$

This regularization term aims to train the decoder to produce smooth and consistent CAD models, reducing artifacts, and ensuring that the models are robust to variations in the latent input. The latent representations  $z_0$  of the original data are also used to train this decoder. This enhanced approach not only addresses the complexity of translating hierarchical latent structures into functional CAD designs but also significantly improves the adaptability and quality of the generated models.

### 3.4 CONDITIONAL GENERATION

The ability to randomly sample shapes offers limited scope for interaction, underscoring the importance of learning a conditional distribution for user applications. It is crucial to accommodate multiple forms of conditional inputs to address diverse scenarios effectively. Using the flexible conditional mechanism facilitated by the diffusion model, we integrate multiple conditional input modalities using task-specific encoders  $E_\phi$  and a cross-attention module. To enhance flexibility in controlling the distribution, we adopt classifier-free guidance for conditional generation. For training such conditional model, the objective function is formulated as follows:

$$L(\theta, \{\phi_i\}) := \mathbb{E}_{z, c, \epsilon, t} \left[ \|\epsilon - \epsilon_\theta(z_t, t, D \circ E_{\phi_i}(c_i))\|^2 \right] \tag{9}$$

The task-specific encoder  $E_{\phi_i}(c_i)$  is employed for the  $i^{\text{th}}$  modality, while  $D$  represents a dropout operation facilitating classifier-free guidance. In this work, we mainly explore two conditions with many practical applications: using point clouds and initial user input.

## 4 EXPERIMENTS

### 4.1 IMPLEMENTATION DETAILS

**Dataset:** Using the extensive DeepCAD dataset Wu et al. (2021), we acquire ground truth sketch-and-extrude models, comprising 178,238 instances. These models are divided into a training set (90%), a validation set (5%), and a test set (5%). To ensure the integrity of the data set, we implement methods similar to previous studies Willis et al. (2021c); Xu et al. (2022) to detect and eliminate duplicate models from the training set. In addition to removing duplicate models, we extract hierarchical properties for loops, profiles, and solids, and subsequently remove duplicate properties at each level. Furthermore, for training purposes, CAD models are included only if they meet specific criteria: a maximum of 5 solids, 20 loops per profile, 60 curves per loop, and a maximum of 200 commands in the sketch-and-extrude sequence. Following the duplicate removal and filtering processes, the training dataset comprises 102,114 solids, 60,584 profiles, and 150,158 loops for codebook learning. Additionally, 124,451 sketch-and-extrude sequences are retained for CAD model generation training. For CAD engineering drawings, we adopt the approach described in SkexGen Xu et al. (2022) and extract sketches from DeepCAD. A total of 99,650 sketches are utilized for training purposes after duplicate removal.

**Other Details:** The model is trained on a Nvidia RTX A100 GPU with a batch size of 256. Each VQ-VAE model and model decoder are trained for 250 epochs. For the randomly generated and conditionally generated diffusion models, we train them for 350 epochs and 500 epochs, respectively. We use the AdamW optimizer Loshchilov & Hutter (2018) with a learning rate of 0.001 after a linear warm-up for the first 2000 steps. The VQ-VAE network consists of 4 layers. For the diffusion model, there are six blocks, each comprising a self-attention layer and a fully-connected layer. If it is a conditional generation, each block also includes a cross-attention layer. During the generation process, we use DDIM for sampling, with a sampling step of 100 steps. For the corresponding scale of the CFG, we set it to 3.  $\sigma$  for perturbation is set it to 0.1. More details can be found in the appendix.

**Evaluation Metrics:** We use five established metrics to quantitatively assess random generation. Three metrics are based on point clouds sampled on the model surfaces. Two metrics scrutinizing generated tokens originating from sketch and extrude construction sequences. For point-cloud evaluation, 2,000 points are sampled from each generated and ground-truth dataset, facilitating a comparative analysis of the two sets. Descriptions of metrics are referred to in the appendix.

### 4.2 UNCONDITIONAL GENERATION

For the unconditional generation, we conduct comparisons with three CAD generation works, DeepCAD, SkexGen and HNC. The results of the other three methods were obtained using publicly available code, and we used their default settings. Each method produces 10,000 CAD models, which are then compared with a randomly selected subset of 2,500 ground truth models from the test set. We compute all metrics three times and take the average for comparison.

**Quantitative Evaluation:** As shown in Table 1, our method performs better than previous methods in all three metrics corresponding to point clouds, especially MMD and JSD, which are significantly better than the baseline methods, demonstrating notable improvements in both quality and diversity. The uniqueness score of our method is similar to the previous two methods and significantly better than DeepCAD. Although SkexGen’s novelty score is similar to ours, it fails to generate highly complex CAD models (see

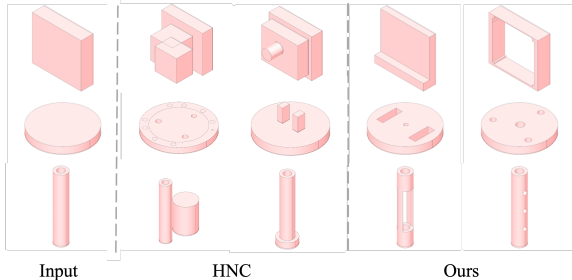


Figure 3: **Conditional generation results of CAD from initial user input.**

Table 1: Quantitative evaluations on the CAD generation task based on the Coverage (COV) percentage, Minimum Matching Distance (MMD), Jensen-Shannon Divergence (JSD), the percentage of Unique, Novel scores and Realism.

Method	COV % $\uparrow$	MMD $\downarrow$	JSD $\downarrow$	Novel % $\uparrow$	Unique % $\uparrow$	Realism % $\uparrow$
DeepCAD	79.98	1.21	3.34	90.5	86.4	36.4
SkexGen	83.58	1.11	0.91	99.2	99.8	42.3
HNC	86.62	1.03	0.74	94.1	99.7	44.1
Ours w/o reg	89.03	0.12	0.13	99.8	99.7	43.1
Ours	90.08	0.10	0.13	99.8	99.6	51.3

Table 2: Comparison with DeepCAD and Draw Step by Step, mean and median Chamfer Distance (CD) results. By employing more advanced conditional generation models, our method obtains more accurate reconstruction results.

Model	Mean CD $\downarrow (\times 10^3)$	Median CD $\downarrow (\times 10^3)$
DeepCAD	43.18	9.836
Draw Step by Step	39.16	7.821
Ours	<b>32.17</b>	<b>6.304</b>

Figure 2), as reported in previous methods. The comparison of these results demonstrates that our distribution fitting is quite effective. However, these metrics are intended to indicate how closely the model’s output matches the real distribution. They do not adequately measure the realism of the results. To better demonstrate the effectiveness of our approach from a quantitative perspective, we introduce a Human Evaluation similar to that in HNC Xu et al. (2023) to measure the realism of the generated complex results. For specific practices, please refer to the appendix. From this realistic comparison result in Table 1, it can further be seen that our method is capable of learning to generate complex and realistic models.

**Qualitative Evaluation:** As shown in Figure 2, the results of DeepCAD do not exhibit significant issues when generating a simple CAD. However, when tasked with generating complex structures, it often resorts to assembling original components and struggles to create CADs that resemble real-world examples in a rational manner. SkexGen may perform slightly better than DeepCAD, but it still encounters similar challenges. HNC’s results show considerable improvement, yet artifacts with small components are occasionally present. In comparison, our approach yields results with more well-structured characteristics, closely resembling real mechanical parts.

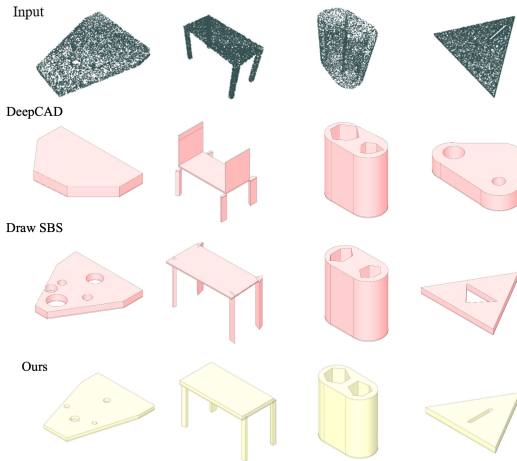


Figure 4: **Conditional generation results of CAD from point clouds.**

### 4.3 CONDITIONAL GENERATION

#### Autocompletion from User Input: We

consider generating a detailed model given an initial model, and we can also use this type of input as control conditions for the diffusion model to generate corresponding potential CAD models for automatic completion purposes. During the training process, we utilize random initial inputs as conditions for the conditional encoder, which serves as input for the diffusion model to reconstruct the corresponding complete CAD models. The encoder encodes the extruded profile parameters and shares the same structure as the encoder employed in training the codebook with VQ-VAE. Figure 3 shows the corresponding CAD autocomplete results with rich details from initial extruded profiles. Each row contains multiple generated results, each corresponding to different noise samples. It can be observed that the autocomplete results are generally reasonable and of high quality, which can assist



designers in CAD design. HNC has also implemented a similar functionality, and we have compared it. Due to the lack of a more suitable comparative method, we also employed human evaluation to test Realism. The results for the HNC Realism were 48.2%, while ours were 52.5%.

**Point to CAD:** 3D reverse engineering entails inferring a CAD model from a 3D scan, a process that demands the expertise of designers and often consumes considerable time. Our method can control the generation process through conditional input from point clouds to obtain CAD models that closely resemble the point clouds, thus achieving a relatively rapid and accurate reverse engineering process to some extent. The encoder used here is a pre-trained ULIP Xue et al. (2023) PointNet Qi et al. (2017). Next, we assess the proposed method for CAD generation based on the point cloud condition. As shown in Figure 4, given a point cloud, our method can obtain a CAD model that is generally similar in overall structure. DeepCAD has also implemented a similar functionality, and we have compared it with our method and another method called Draw Step by Step Ma et al. (2024). For the reconstructed CAD results, we assess them quantitatively against ground truth CAD models using mean and median Chamfer Distances (CD). The quantitative comparison results are shown in Table 2.

#### 4.4 ABLATION STUDY

To demonstrate the significance of the regularization term we introduced, we conducted a comparative analysis between models with and without this regularization. As illustrated in Figure 5, the decoder enhanced with regularization is capable of rationalizing outputs that were previously deemed unrealistic. Initially, some of these unsatisfactory results were attributed to noise and a lack of inherent symmetrical logic in the decoded outputs. By integrating a smoothness-inducing regularization, we observed that the decoder could produce CAD models that more closely align with the intrinsic logic of mechanical parts.

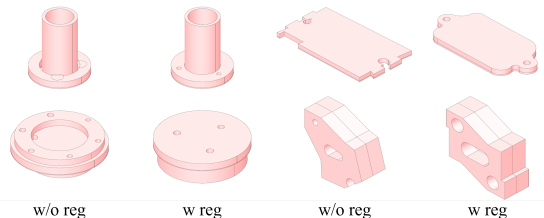


Figure 5: **Ablation study.** It can be observed that the results on the right, which incorporate regularization, are more aesthetically pleasing and logical, with less noise compared to those without regularization.

For more visual results, please refer to the appendix. As shown in the last two rows of Table 1, after adding the regularization terms, the realism of the results generated by our method has significantly improved. Our regularization improves other metrics as well, but not as significantly, because, as with HNC, the Realism metric measures the complex results with three or more extrusions, as only such complex results have evaluative value. Other metrics measure the average across all generated results, and our regularization does not significantly improve simple CAD results, which have limited potential for improvement. This leads to a relatively minor improvement in other metrics when averaged. No single metric can fully evaluate the quality of generated results, generating complex but unrealistic results can lead to high Novelty and Uniqueness scores, without significantly affecting other metrics that measure the diversity of generated results.

## 5 CONCLUSIONS AND FUTURE WORK

In conclusion, we have introduced CADiffusion, a novel diffusion-based generative model tailored for Computer-Aided Design (CAD) data generation. Our approach addresses the persistent challenges in producing diverse and high-quality CAD shapes by seamlessly integrating diffusion models and Vector Quantized Variational Autoencoders (VQVAE) to obtain latent representations. Through extensive experimentation, we have demonstrated the effectiveness of CADiffusion, achieving state-of-the-art performance on benchmark datasets. Additionally, we introduced a regularization method specifically tailored for training the decoder. This method employs perturbations in the Gaussian space to smooth the decoder’s outputs, thereby enabling it to produce more reasonable results. We believe that CADiffusion opens up new possibilities for advancing 3D shape generation in practical CAD modeling and design applications.

**limitation.** The current results of point to CAD do not fully match the input yet. In the future, we hope to explore better reverse engineering methods using diffusion priors.

## REFERENCES

- 486  
487  
488 Ruojin Cai, Guandao Yang, Hadar Averbuch-Elor, Zekun Hao, Serge Belongie, Noah Snaveley, and  
489 Bharath Hariharan. Learning gradient fields for shape generation. In *Proceedings of the European  
490 Conference on Computer Vision (ECCV)*, 2020.
- 491 Nanxin Chen, Yu Zhang, Heiga Zen, Ron J. Weiss, Mohammad Norouzi, and William Chan. Waveg-  
492 rad: Estimating gradients for waveform generation. In *ICLR*, 2021.
- 493  
494 Zhiqin Chen, Andrea Tagliasacchi, and Hao Zhang. Bsp-net: Generating compact meshes via binary  
495 space partitioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
496 Recognition*, pp. 45–54, 2020.
- 497 Yen-Chi Cheng, Hsin-Ying Lee, Sergey Tulyakov, Alexander G Schwing, and Liang-Yan Gui.  
498 Sdfusion: Multimodal 3d shape completion, reconstruction, and generation. In *Proceedings of the  
499 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4456–4465, 2023.
- 500  
501 Gene Chou, Yuval Bahat, and Felix Heide. Diffusion-sdf: Conditional generative modeling of signed  
502 distance functions. In *Proceedings of the IEEE/CVF international conference on computer vision*,  
503 pp. 2262–2272, 2023.
- 504 Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. Diffusion models in  
505 vision: A survey. *IEEE Trans. Pattern Anal. Mach. Intell.*, 45(9):10850–10869, 2023. doi: 10.1109/  
506 TPAMI.2023.3261988. URL <https://doi.org/10.1109/TPAMI.2023.3261988>.
- 507  
508 Prafulla Dhariwal and Alexander Quinn Nichol. Diffusion models beat gans on image synthesis. In  
509 Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman  
510 Vaughan (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on  
511 Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp.  
512 8780–8794, 2021. URL [https://proceedings.neurips.cc/paper/2021/hash/  
513 49ad23d1ec9fa4bd8d77d02681df5cfa-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/49ad23d1ec9fa4bd8d77d02681df5cfa-Abstract.html).
- 514 Yaroslav Ganin, Sergey Bartunov, Yujia Li, Ethan Keller, and Stefano Saliceti. Computer-aided  
515 design as language. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang,  
516 and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems 34:  
517 Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December  
518 6-14, 2021, virtual*, pp. 5885–5897, 2021. URL [https://proceedings.neurips.cc/  
519 paper/2021/hash/2e92962c0b6996add9517e4242ea9bdc-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/2e92962c0b6996add9517e4242ea9bdc-Abstract.html).
- 520 Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,  
521 Aaron C. Courville, and Yoshua Bengio. Generative adversarial networks. *Commun. ACM*, 63(11):  
522 139–144, 2020.
- 523  
524 Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. A  
525 papier-mâché approach to learning 3d surface generation. In *CVPR*, pp. 216–224, 2018.
- 526 Haoxiang Guo, Shilin Liu, Hao Pan, Yang Liu, Xin Tong, and Baining Guo. Complexgen: Cad  
527 reconstruction by b-rep chain complex generation. *ACM Trans. Graph. (SIGGRAPH)*, 41(4),  
528 July 2022. doi: 10.1145/3528223.3530078. URL [https://doi.org/10.1145/3528223.  
529 3530078](https://doi.org/10.1145/3528223.3530078).
- 530  
531 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In  
532 Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-  
533 Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Con-  
534 ference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12,  
535 2020, virtual*, 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/  
536 4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html).
- 537 Pradeep Kumar Jayaraman, Joseph G. Lambourne, Nishkrit Desai, Karl D. D. Willis, Aditya Sanghi,  
538 and Nigel J. W. Morris. Solidgen: An autoregressive model for direct b-rep synthesis. *arXiv  
539 Preprint*, 2022. doi: 10.48550/ARXIV.2203.13944. URL [https://arxiv.org/abs/  
2203.13944](https://arxiv.org/abs/2203.13944).

- 540 R. Kenny Jones, Theresa Barton, Xianghao Xu, Kai Wang, Ellen Jiang, Paul Guerrero, Niloy J.  
541 Mitra, and Daniel Ritchie. Shapeassembly: Learning to generate programs for 3d shape structure  
542 synthesis. *ACM Transactions on Graphics (TOG), Siggraph Asia 2020*, 39(6):Article 234, 2020.  
543
- 544 Kacper Kania, Maciej Zieba, and Tomasz Kajdanowicz. UcsG-net—unsupervised discovering of  
545 constructive solid geometry tree. *arXiv preprint arXiv:2006.09102*, 2020.  
546
- 547 Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *ICLR*, 2014.  
548
- 549 Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile  
550 diffusion model for audio synthesis. In *ICLR*, 2021.  
551
- 552 Joseph G. Lambourne, Karl D.D. Willis, Pradeep Kumar Jayaraman, Longfei Zhang, Aditya Sanghi,  
553 and Kamal Rahimi Malekshan. Reconstructing editable prismatic cad from rounded voxel models.  
554 In *SIGGRAPH Asia*, December 2022.
- 555 Changjian Li, Hao Pan, Adrien Bousseau, and Niloy J Mitra. Sketch2cad: Sequential cad modeling  
556 by sketching in context. *ACM TOG*, 39(6):1–14, 2020.  
557
- 558 Jun Li, Kai Xu, Siddhartha Chaudhuri, Ersin Yumer, Hao Zhang, and Leonidas Guibas. Grass:  
559 Generative recursive autoencoders for shape structures. *ACM Transactions on Graphics (Proc. of*  
560 *SIGGRAPH 2017)*, 36(4):to appear, 2017.
- 561 Pu Li, Jianwei Guo, Xiaopeng Zhang, and Dongming Yan. Secad-net: Self-supervised cad recon-  
562 struction by learning sketch-extrude operations. In *Proceedings of the IEEE/CVF Conference on*  
563 *Computer Vision and Pattern Recognition*, pp. 16816–16826, 2023.  
564
- 565 Yiyi Liao, Simon Donne, and Andreas Geiger. Deep marching cubes: Learning explicit surface repre-  
566 sentations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,  
567 pp. 2916–2925, 2018.  
568
- 569 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer-*  
570 *ence on Learning Representations*, 2018.  
571
- 572 Shitong Luo and Wei Hu. Diffusion probabilistic models for 3D point cloud generation. In *CVPR*,  
573 pp. 2837–2845, 2021.
- 574 Weijian Ma, Shuaiqi Chen, Yunzhong Lou, Xueyang Li, and Xiangdong Zhou. Draw step by  
575 step: Reconstructing cad construction sequences from point clouds via multimodal diffusion.  
576 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
577 27154–27163, 2024.  
578
- 579 Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger.  
580 Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the*  
581 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4460–4470, 2019.  
582
- 583 Kaichun Mo, Paul Guerrero, Li Yi, Hao Su, Peter Wonka, Niloy Mitra, and Leonidas Guibas.  
584 StructureNet: Hierarchical graph networks for 3d shape generation. *ACM Transactions on Graphics*  
585 *(TOG), Siggraph Asia 2019*, 38(6):Article 242, 2019.
- 586 Charlie Nash, Yaroslav Ganin, S. M. Ali Eslami, and Peter Battaglia. PolyGen: An autoregressive  
587 generative model of 3D meshes. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the*  
588 *37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine*  
589 *Learning Research*, pp. 7220–7229. PMLR, 13–18 Jul 2020. URL [http://proceedings.](http://proceedings.mlr.press/v119/nash20a.html)  
590 [mlr.press/v119/nash20a.html](http://proceedings.mlr.press/v119/nash20a.html).  
591
- 592 Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob  
593 McGrew, Ilya Sutskever, and Mark Chen. GLIDE: towards photorealistic image generation and  
editing with text-guided diffusion models. In *ICML*, pp. 16784–16804, 2022.

- 594 Wamiq Reyaz Para, Shariq Farooq Bhat, Paul Guerrero, Tom Kelly, Niloy J. Mitra, Leonidas J.  
595 Guibas, and Peter Wonka. Sketchgen: Generating constrained CAD sketches. In Marc’Aurelio  
596 Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan  
597 (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neu-  
598 ral Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp.  
599 5077–5088, 2021. URL [https://proceedings.neurips.cc/paper/2021/hash/  
600 28891cb4ab421830acc36b1f5fd6c91e-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/28891cb4ab421830acc36b1f5fd6c91e-Abstract.html).
- 601 Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deep-  
602 sdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the  
603 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 165–174, 2019.
- 604 Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets  
605 for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision  
606 and pattern recognition*, pp. 652–660, 2017.
- 607 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-  
608 conditional image generation with CLIP latents. *CoRR*, abs/2204.06125, 2022. doi: 10.48550/  
609 ARXIV.2204.06125. URL <https://doi.org/10.48550/arXiv.2204.06125>.
- 610 Daxuan Ren, Jianmin Zheng, Jianfei Cai, Jiatong Li, and Junzhe Zhang. Extrudenet: Unsupervised  
611 inverse sketch-and-extrude for shape parsing. In *ECCV, 2022*.
- 612 Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. In  
613 *ICML*, pp. 1530–1538, 2015.
- 614 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
615 resolution image synthesis with latent diffusion models. In *CVPR*, pp. 10674–10685, 2022.
- 616 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L. Denton,  
617 Seyed Kamyar Seyed Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim  
618 Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-  
619 image diffusion models with deep language understanding. In Sanmi Koyejo, S. Mo-  
620 hamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural  
621 Information Processing Systems 35: Annual Conference on Neural Information Process-  
622 ing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9,  
623 2022, 2022*. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/  
624 ec795aeadae0b7d230fa35cbaf04c041-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/ec795aeadae0b7d230fa35cbaf04c041-Abstract-Conference.html).
- 625 Ari Seff, Yaniv Ovadia, Wenda Zhou, and Ryan P. Adams. SketchGraphs: A large-scale dataset  
626 for modeling relational geometry in computer-aided design. In *ICML 2020 Workshop on Object-  
627 Oriented Learning*, 2020.
- 628 Ari Seff, Wenda Zhou, Nick Richardson, and Ryan P. Adams. Vitruvion: A generative model  
629 of parametric CAD sketches. In *The Tenth International Conference on Learning Represen-  
630 tations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL <https://openreview.net/forum?id=Ow1C7s3UcY>.
- 631 Gopal Sharma, Rishabh Goyal, Difan Liu, Evangelos Kalogerakis, and Subhransu Maji. Csgnet:  
632 Neural shape parser for constructive solid geometry. In *Proceedings of the IEEE Conference on  
633 Computer Vision and Pattern Recognition*, pp. 5515–5523, 2018.
- 634 Gopal Sharma, Difan Liu, Subhransu Maji, Evangelos Kalogerakis, Siddhartha Chaudhuri, and  
635 Radomír Měch. Parsenet: A parametric surface fitting network for 3d point clouds. In *ECCV*, pp.  
636 261–276. Springer, 2020.
- 637 Dmitriy Smirnov, Mikhail Bessmeltsev, and Justin Solomon. Learning manifold patch-based repre-  
638 sentations of man-made shapes. In *ICLR*, 2021.
- 639 Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised  
640 learning using nonequilibrium thermodynamics. In *ICML*, pp. 2256–2265, 2015.
- 641 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*,  
642 2021.

- 648 Mikaela Angelina Uy, Yen-yu Chang, Minhyuk Sung, Purvi Goel, Joseph Lambourne, Tolga Birdal,  
649 and Leonidas J. Guibas. Point2cyl: Reverse engineering 3d objects from point clouds to extrusion  
650 cylinders. *CoRR*, abs/2112.09329, 2021. URL <https://arxiv.org/abs/2112.09329>.
- 651 Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learn-  
652 ing. In *Proceedings of the 31st International Conference on Neural Information Processing*  
653 *Systems, NIPS'17*, pp. 6309–6318, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN  
654 9781510860964.
- 655 Kehan Wang, Jia Zheng, and Zihan Zhou. Neural face identification in a 2d wireframe projection of  
656 a manifold object. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition*  
657 *(CVPR)*, pp. 1612–1621, 2022. doi: 10.1109/CVPR52688.2022.00167.
- 658 Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2mesh:  
659 Generating 3d mesh models from single rgb images. In *Proceedings of the European Conference*  
660 *on Computer Vision (ECCV)*, pp. 52–67, 2018.
- 661 Xiaogang Wang, Yuelang Xu, Kai Xu, Andrea Tagliasacchi, Bin Zhou, Ali Mahdavi-Amiri, and Hao  
662 Zhang. Pie-net: Parametric inference of point cloud edges. In *Advances in Neural Information*  
663 *Processing Systems*, volume 33, pp. 20167–20178. Curran Associates, Inc., 2020.
- 664 Karl D. D. Willis, Pradeep Kumar Jayaraman, Joseph G. Lambourne, Hang Chu, and Yewen Pu. Engi-  
665 neering sketch generation for computer-aided design. In *IEEE Conference on Computer Vision and*  
666 *Pattern Recognition Workshops, CVPR Workshops 2021, virtual, June 19-25, 2021*, pp. 2105–2114.  
667 Computer Vision Foundation / IEEE, 2021a. doi: 10.1109/CVPRW53098.2021.00239. URL  
668 [https://openaccess.thecvf.com/content/CVPR2021W/SketchDL/html/](https://openaccess.thecvf.com/content/CVPR2021W/SketchDL/html/Willis_Engineering_Sketch_Generation_for_Computer-Aided_Design_CVPRW_2021_paper.html)  
669 [Willis\\_Engineering\\_Sketch\\_Generation\\_for\\_Computer-Aided\\_Design\\_](https://openaccess.thecvf.com/content/CVPR2021W/SketchDL/html/Willis_Engineering_Sketch_Generation_for_Computer-Aided_Design_CVPRW_2021_paper.html)  
670 [CVPRW\\_2021\\_paper.html](https://openaccess.thecvf.com/content/CVPR2021W/SketchDL/html/Willis_Engineering_Sketch_Generation_for_Computer-Aided_Design_CVPRW_2021_paper.html).
- 671 Karl D. D. Willis, Yewen Pu, Jieliang Luo, Hang Chu, Tao Du, Joseph G. Lambourne, Armando Solar-  
672 Lezama, and Wojciech Matusik. Fusion 360 gallery: A dataset and environment for programmatic  
673 cad construction from human design sequences. *ACM TOG*, 40(4), 2021b.
- 674 Karl DD Willis, Pradeep Kumar Jayaraman, Joseph G Lambourne, Hang Chu, and Yewen Pu.  
675 Engineering sketch generation for computer-aided design. In *CVPRW*, pp. 2105–2114, 2021c.
- 676 Rundi Wu, Chang Xiao, and Changxi Zheng. Deepcad: A deep generative network for computer-aided  
677 design models. In *ICCV*, pp. 6772–6782, October 2021.
- 678 Xiang Xu, Karl DD Willis, Joseph G Lambourne, Chin-Yi Cheng, Pradeep Kumar Jayaraman, and  
679 Yasutaka Furukawa. Skexgen: Autoregressive generation of cad construction sequences with  
680 disentangled codebooks. In *International Conference on Machine Learning*, pp. 24698–24724.  
681 PMLR, 2022.
- 682 Xiang Xu, Pradeep Kumar Jayaraman, Joseph G Lambourne, Karl DD Willis, and Yasutaka Furukawa.  
683 Hierarchical neural coding for controllable cad model generation. *arXiv preprint arXiv:2307.00149*,  
684 2023.
- 685 Xianghao Xu, Wenzhe Peng, Chin-Yi Cheng, Karl D.D. Willis, and Daniel Ritchie. Inferring cad  
686 modeling sequences using zone graphs. In *CVPR*, pp. 6062–6070, June 2021.
- 687 Le Xue, Mingfei Gao, Chen Xing, Roberto Martín-Martín, Jiajun Wu, Caiming Xiong, Ran Xu,  
688 Juan Carlos Nieves, and Silvio Savarese. Ulip: Learning a unified representation of language,  
689 images, and point clouds for 3d understanding. In *Proceedings of the IEEE/CVF Conference on*  
690 *Computer Vision and Pattern Recognition*, pp. 1179–1189, 2023.
- 691 Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge Belongie, and Bharath Hariharan.  
692 Pointflow: 3d point cloud generation with continuous normalizing flows. In *Proceedings of the*  
693 *IEEE/CVF International Conference on Computer Vision*, pp. 4541–4550, 2019.
- 694 Xiaohui Zeng, Arash Vahdat, Francis Williams, Zan Gojcic, Or Litany, Sanja Fi-  
695 dler, and Karsten Kreis. LION: latent point diffusion models for 3d shape genera-  
696 tion. 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/](http://papers.nips.cc/paper_files/paper/2022/hash/40e56dabe12095a5fc44a6e4c3835948-Abstract-Conference.html)  
697 [40e56dabe12095a5fc44a6e4c3835948-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/40e56dabe12095a5fc44a6e4c3835948-Abstract-Conference.html).

702 Linqi Zhou, Yilun Du, and Jiajun Wu. 3D shape generation and completion through point-voxel  
703 diffusion. In *ICCV*, pp. 5806–5815, 2021.

## 706 A EVALUATION METRICS

708 Due to the unique characteristics of CAD data modalities, we have provided detailed explanations of  
709 the five different metrics mentioned in the main paper here, to facilitate reader understanding.

- 711 • Coverage (COV) denotes the proportion of ground-truth models that contain at least one  
712 matched generated sample, where matching is determined based on the Chamfer distance  
713 (CD) or Earth Mover’s distance (EMD). COV serves as a measure of the diversity of  
714 generated shapes, revealing potential mode collapse if only a few ground-truth models are  
715 matched, resulting in low coverage scores.
- 716 • Minimum Matching Distance (MMD) calculates the average minimum matching distance  
717 between the ground truth and generated sets.
- 718 • Jensen-Shannon divergence (JSD) assesses the similarity between two probability distribu-  
719 tions, reflecting the degree to which ground truth and generated point clouds occupy similar  
720 locations. Utilizing voxelization, occupancy distributions are computed to derive the JSD  
721 score.
- 722 • Novelty indicates the percentage of generated CAD sequences absent in the training set,  
723 while uniqueness signifies the percentage of generated data that appear once within the  
724 generated set.

## 726 B IMPLEMENTATION DETAILS

728 The tree structure corresponding to the CAD is obtained directly by parsing according to the CAD  
729 logic, implemented by a segment of code. This part is the same as in HNC. The CAD latent used  
730 for diffusion learning is in the form of tensors with a shape of  $32 \times 256$ , the input and output  
731 shapes are the same. Therefore, all our model architectures use 1D transformer structures, which  
732 mainly include Self-Attention, Fully Connected Layers, and Cross-Attention layers. In the Self-  
733 Attention mechanism, the basic QKV form corresponds to different matrix transformations followed  
734 by attention calculations. In the Cross-Attention layer, the input conditions are also transformed  
735 into an  $n \times 256$  format by a specific encoder and then used as the KV components in the attention  
736 computation.

## 738 C HUMAN EVALUATION

740 Our approach to Human Evaluation is similar to HNC and SolidGen. For each method that need to  
741 be evaluated, we randomly select models with three or more extrusions from their results generated  
742 without any conditions. For each model created by a generation method, we randomly choose a  
743 real model from the dataset and display the rendered images of these two models side by side.  
744 We randomly select 100 results from each method, and the image pairs are presented to college  
745 participants, who are asked to assess which one appears more “realistic.”

## 747 D MORE EXPERIMENTAL RESULTS

749 Here, we include visualizations of more results. Randomly generated results are shown in Figure 6,  
750 autocompletion results are shown in Figure 7, and Point2CAD results are shown in Figure 8. For  
751 more ablation study results, please refer to Figure 9.

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809



Figure 6: **Random generation results of CADiffusion.** The CADiffusion model can generate diverse and realistic CAD models.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863



Figure 7: **Conditional generation results of CAD from initial user input.** The leftmost column represents initial conditional inputs. The autocomplete results demonstrate a consistent level of quality and reasonableness, providing valuable assistance to CAD designers in their design processes.



864  
 865  
 866  
 867  
 868  
 869  
 870  
 871  
 872  
 873  
 874  
 875  
 876  
 877  
 878  
 879  
 880  
 881  
 882  
 883  
 884  
 885  
 886  
 887  
 888  
 889  
 890  
 891  
 892  
 893  
 894  
 895  
 896  
 897  
 898  
 899  
 900  
 901  
 902  
 903  
 904  
 905  
 906  
 907  
 908  
 909  
 910  
 911  
 912  
 913  
 914  
 915  
 916  
 917

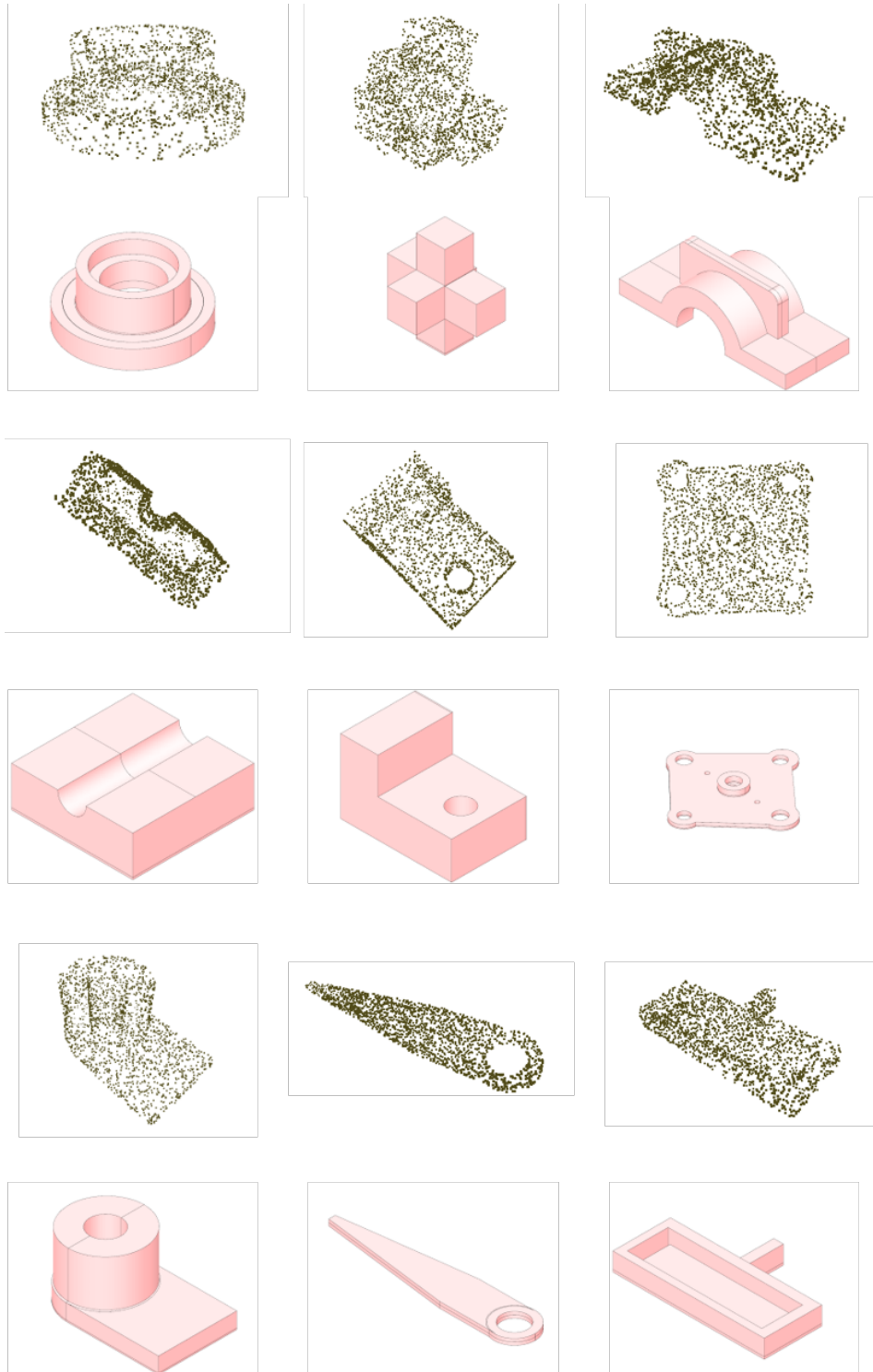


Figure 8: **Conditional generation results of CAD from point clouds.** The rows from top to bottom consist of input point clouds and their corresponding generated CAD models. Based on a provided point cloud, our approach is capable of deriving a CAD model that exhibits a broad similarity in its overall structure.

918  
 919  
 920  
 921  
 922  
 923  
 924  
 925  
 926  
 927  
 928  
 929  
 930  
 931  
 932  
 933  
 934  
 935  
 936  
 937  
 938  
 939  
 940  
 941  
 942  
 943  
 944  
 945  
 946  
 947  
 948  
 949  
 950  
 951  
 952  
 953  
 954  
 955  
 956  
 957  
 958  
 959  
 960  
 961  
 962  
 963  
 964  
 965  
 966  
 967  
 968  
 969  
 970  
 971

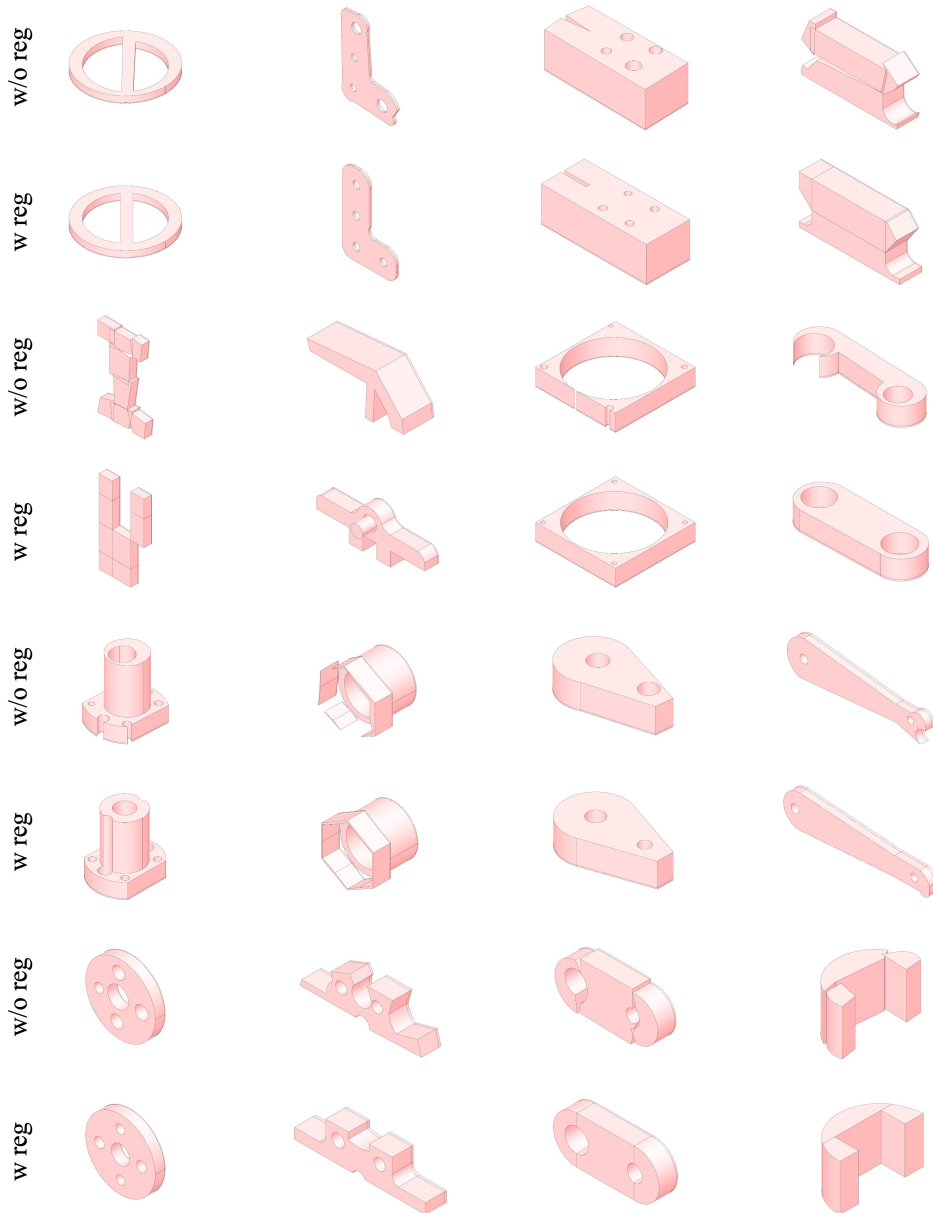


Figure 9: **Ablation studies.** It can be observed that the results which incorporate regularization, are more aesthetically pleasing and logical, with less noise compared to those without regularization.