# CAD-Llama: Leveraging Large Language Models for Computer-Aided Design Parametric 3D Model Generation

Jiahao Li    Weijian Ma    Xueyang Li    Yunzhong Lou    Guichun Zhou    Xiangdong Zhou*

School of Computer Science and Technology, Fudan University

{lijh23, mawj22, xueyangli21}@m.fudan.edu.cn    {yzlou20, gczhou19, xdzhou}@fudan.edu.cn
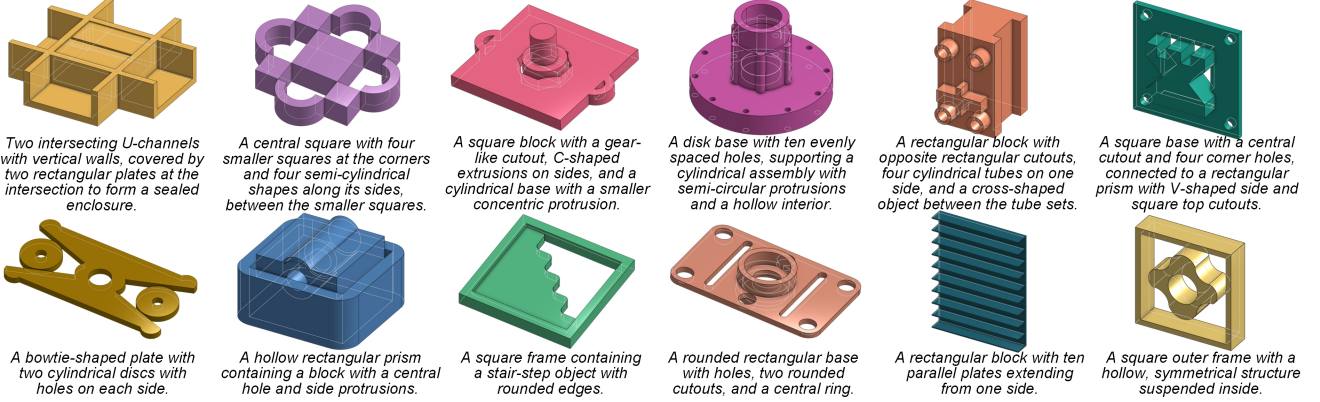
Figure 1. **A collection of generated CAD models with text prompts by using our method (CAD-Llama-INS).** Our approach enables the generation of more complex CAD models based on both abstract and detailed text prompts.

*The captions below the images:*
- Two intersecting U-channels with vertical walls, covered by two rectangular plates at the intersection to form a sealed enclosure.
- A central square with four smaller squares at the corners and four semi-cylindrical shapes along its sides, between the smaller squares.
- A square block with a gear-like cutout, C-shaped extrusions on sides, and a cylindrical base with a smaller concentric protrusion.
- A disk base with ten evenly spaced holes, supporting a cylindrical assembly with semi-circular protrusions and a hollow interior.
- A rectangular block with opposite rectangular cutouts, four cylindrical tubes on one side, and a cross-shaped object between the tube sets.
- A square base with a central cutout and four corner holes, connected to a rectangular prism with V-shaped side and square top cutouts.
- A bowtie-shaped plate with two cylindrical discs with holes on each side.
- A hollow rectangular prism containing a block with a central hole and side protrusions.
- A square frame containing a stair-step object with rounded edges.
- A rounded rectangular base with holes, two rounded cutouts, and a central ring.
- A rectangular block with ten parallel plates extending from one side.
- A square outer frame with a hollow, symmetrical structure suspended inside.

## Abstract

*Recently, Large Language Models (LLMs) have achieved significant success, prompting increased interest in expanding their generative capabilities beyond general text into domain-specific areas. This study investigates the generation of parametric sequences for computer-aided design (CAD) models using LLMs. This endeavor represents an initial step towards creating parametric 3D shapes with LLMs, as CAD model parameters directly correlate with shapes in three-dimensional space. Despite the formidable generative capacities of LLMs, this task remains challenging, as these models neither encounter parametric sequences during their pretraining phase nor possess direct awareness of 3D structures. To address this, we present CAD-Llama, a framework designed to enhance pretrained LLMs for generating parametric 3D CAD models. Specifically, we develop a hierarchical annotation pipeline and a code-like format to translate parametric 3D CAD command sequences into Structured Parametric CAD Code (SPCC), incorporating hierarchical semantic descriptions. Furthermore, we propose an adaptive pretraining approach utilizing SPCC, followed by an instruction tuning process aligned with CAD-specific guidelines. This methodology aims to equip LLMs with the spatial knowledge inherent in parametric sequences. Experimental results demonstrate that our framework significantly outperforms prior autoregressive methods and existing LLM baselines.*

## 1. Introduction

Computer-Aided Design (CAD) generative modeling has attracted increasing attention from research and industry communities. Large language models (LLMs) have recently demonstrated strong generative capabilities and impressive zero-shot performance across a broad range of downstream tasks [8, 32, 64, 67]. These models have also found widespread applications in the real world [6, 12, 29, 36]. However, the exploration of utilizing LLMs for generating parametric CAD construction sequences remains underexplored, thereby calling for further investigation on how to invoke the potential of LLM's learned priors onto the task of parametric CAD sequence generation and editing.

Leveraging LLMs for parametric CAD sequence generation is nontrivial. A substantial disparity exists between the original parameterized CAD sequences and the natural language familiar to LLMs, rendering the direct generation of parametric CAD sequences by LLMs a challenging task. Most previous works reconstruct parametric CAD sequences from various inputs, such as point clouds [34], text [20, 25, 58], B-rep models [55, 61], and partial CAD [62, 63], using encoder-decoder architectures trained solely

---

* Corresponding author.

1

on CAD dataset [20, 25, 58]. Some recent attempts demonstrate that LLMs can generate basic CAD construction sequences [3, 26, 59, 66] and have the potential to understand the semantics of symbolic graphic programs [43]. However, most of these methods suffer from weak generalization and lack the ability to generate complex CAD models, let alone generate CAD models based on complex text instructions.

We note that in order to effectively employing the generative capabilities of LLMs for CAD sequence generation necessitates a comprehensive understanding of both the characteristics of CAD data and the intrinsic capabilities of LLMs. The parametric CAD model, also referred to as CAD design history, consists of sequences of commands from CAD tools, yet it lacks semantic annotations pertaining to the design rationale and the geometry or shape of the respective CAD model. Consequently, without textual descriptions, it is challenging for LLMs to grasp the semantic implications of parametric CAD models. This limitation accounts for the fact that, in prior research, LLMs have typically only generated relatively simple CAD models. Conversely, LLMs excel in code generation owing to the extensive repository of code data accompanied by text comments and functional descriptions present in the training datasets.

Leveraging insights from the CAD modeling process and the remarkable language generation capabilities of LLMs, we propose CAD-Llama, an extensive framework that adapts open-source LLMs for the generation of CAD command sequences. For data acquisition, we introduce a novel hierarchical data annotation pipeline for CAD design history data, which is represented in the form of Python-like code, called **S**tructured **P**arametric **C**AD **C**ode (SPCC). During the annotation process, a visual language model (VLM) is utilized to annotate both the three-dimensional geometry and the two-dimensional sketch of each component with detailed textual descriptions. Subsequently, the comprehensive semantics and the interrelationships among components are captioned to yield the top-layer thorough descriptions. Regarding training methodologies, an adaptive pretraining paradigm, in conjunction with instruction tuning techniques for varied downstream tasks, is proposed to impart CAD modeling capabilities to the LLM and to adapt it for diverse downstream applications.

We conduct a series of experiments to evaluate our approach, covering both unconditional and conditional generation tasks. The results indicate that our method outperforms recent state-of-the-art parametric CAD generation models, as well as open-source models such as GPT-4 and LLaMA3, across various CAD-related tasks. We show that using rich and structured text descriptions of 3D shape and geometry to fine-tune LLMs leads to the emergence of the ability to generate professional parametric 3D CAD models under complex text instructions, as shown in Figure 1, which has not been explored or reported in previous studies.

In summary, our contributions include the following.
1. We present CAD-Llama, a novel unified framework to leverage LLMs' generative priors for parametric 3D CAD modeling based on text instructions.
2. We introduce a hierarchical annotation pipeline for 3D CAD models that captures both structured information and detailed textual descriptions of 3D shapes and geometry.
3. We propose an adaptive pretraining paradigm combined with instruction tuning on a multitask instructional dataset to align LLMs with CAD sequence modeling across a series of tasks.
4. Experimental results demonstrate that CAD-Llama can generate more accurate and complex parametric CAD models and achieve good performance in a series of downstream tasks.

## 2. Related Work

**Representation Learning of CAD models.** Building representations for understanding CAD models has become a long-sought problem throughout the vision history. Early research focused on utilizing shape-understanding methods to classify and segment CAD models in form of point clouds [41, 42], meshes [15, 49], voxels [30, 37, 47] and SDFs [7, 40]. However, methods in the shape domains fail to capture the exact shape parameters, leading to a difficulty in editing and reusing the created shapes. On the other hand, along with the emergence of large-scale parametric CAD datasets [21, 56, 58], language models have been adopted to model the parametric designs of CAD models. [33] also built a multimodal representation for CAD models based on point clouds and construction sequences. The sequence modeling ability of language models has opened up possibilities of generating precise parametric construction sequences. However, the granularity of control over parameters still remains a problem.

**Crossmodal CAD Generation.** Translating parametric CAD models from given conditions has been a problem of active research. Research in earlier times focused on precise translation from geometric shapes like point clouds or meshes into parametric sequences via heuristic primitive fitting methods like RANSAC [13, 48] or Hough Transform [4, 11]. Some follow-up works attempted to broaden the scope of input modalities. They are Point cloud[34, 58], partial CAD input [63], target B-reps [55, 61], voxel grids[22, 24], point clouds with[52] or without sequence guidance [24, 46], etc. However, all these works require detailed semantics of the target models, limiting their applications to the domains of concept design. Concurrent works on CAD model generation from text descriptions include Text2CAD [20] and CAD Translator [25], both of which employ encoder-decoder architectures to translate the text descriptions of CAD shapes into parametric CAD
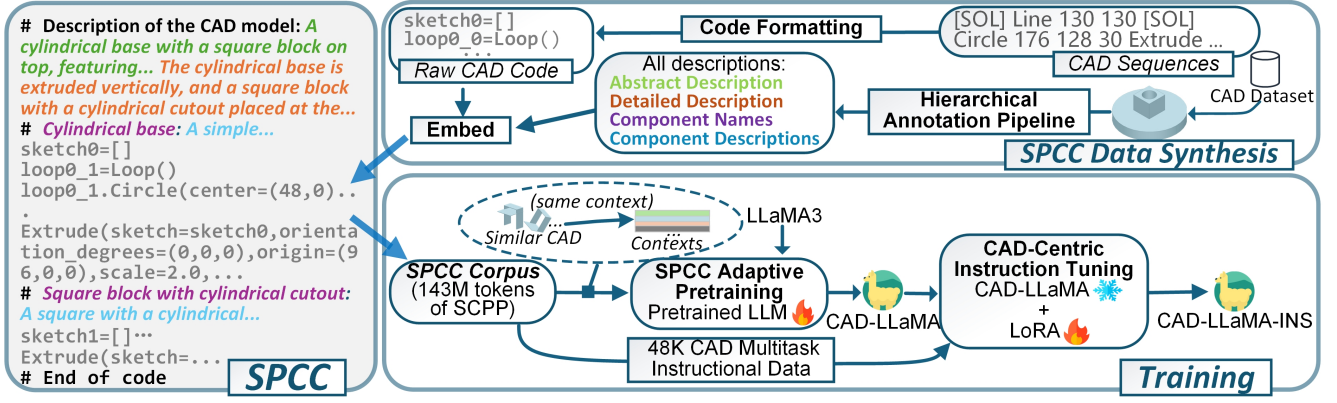
Figure 2. **Overview of the proposed framework CAD-Llama.** The framework consists of two parts: (1) SPCC data synthesis, which employs the hierarchical annotation pipeline to convert CAD sequences into SPCC representations, and (2) the pretraining and instruction tuning process, where the resulting SPCC corpus is leveraged to enhance model performance.

sequences. However, the limited capacity of the encoder-decoder architecture constrains its generalizability on out-of-distribution examples.

**Large Language Models and Computer-Aided Design.**
LLMs have demonstrated growing potential in many applications, ranging from mathematical problem solving and theorem proving assistance [8, 32, 64, 67] to aiding biological discovery [6, 12, 29, 36]. Applying LLMs for abstract content understanding and generation is also a popular direction of research. A recent investigation [43] shows that LLMs can understand symbolic graphic programs like SVG and CAD models via finetuning on VQA tasks. A few attempts tried to investigate the generation ability of LLMs on parametric CAD models. CAD-LLM [59] empirically investigates CAD generation on 2D domains. LLM4CAD [26] utilizes VLMs to perform zero-shot CAD generation tasks. CADTalk [65] generates semantic labels for CAD parts. OpenECAD [66] attempts to finetune a VLM with the assistance of CAD kernels like PythonOCC. Query2CAD [3] proposes a natural language translator into CAD code with an image-captioner in the loop. CAD-MLLM [60] and CAD-GPT [54] both leverage Multimodal Large Language Models (MLLMs) that generate CAD command sequences, with CAD-MLLM supporting diverse inputs like text, images, and point clouds, and CAD-GPT enhancing spatial reasoning for precise synthesis from single-view images or text. However, few previous work succeeded in leveraging LLM's strong generative prior on text to CAD construction sequence generation.

## 3. Method

In this section, we first propose the hierarchical annotation pipeline and the SPCC dataset synthesis for LLMs fine-tuning data preparation. Then, we propose a pretraining method to equip LLMs with CAD model generation capabilities, and an instruction tuning method that further leverage the LLM's ability to handle CAD-related downstream tasks. The framework of CAD-Llama is illustrated in Fig-

ure 2. Details are provided in the following subsections.

| Notation | Definition |
|----------|------------|
| $\mathcal{D}_j$ | The $i$-th CAD in dataset |
| $\mathcal{D}_{j \setminus \mathcal{P}_j^k}$ | CAD of $\mathcal{D}_j$ after removing $k$-th component |
| $C(\mathcal{D}_j)$ | CAD code representation of $\mathcal{D}_j$ |
| $\tilde{\mathcal{D}}_j$ | SPCC representation of $\mathcal{D}_j$ |
| $\mathcal{A}_j, \mathcal{T}_j$ | Overall abstract and detailed descriptions of $\mathcal{D}_j$ |
| $\mathcal{S}_j^i, \mathcal{I}_j^i$ | Name and description of the $i$-th component of $\mathcal{D}_j$ |

Table 1. Summary of key notations.

### 3.1. Notation

Denote the data set of the parametric CAD model as $\mathcal{D} = \{\mathcal{D}_1, \mathcal{D}_2, \ldots, \mathcal{D}_N\}$, where $N$ is the number of CAD models. Assume that $j$-th CAD model $\mathcal{D}_j$ contains $m$ components, represented as $\mathcal{D}_j = \{\mathcal{P}_j^1, \mathcal{P}_j^2, \ldots, \mathcal{P}_j^m\}$, where $\mathcal{P}_j^i$ refers to the parametric CAD sequence of the $i$-th component. In Table 1, we provide a brief description of the key notation. A detailed introduction to these notation is presented in the following two subsections.

### 3.2. Hierarchical Annotation Pipeline

A crucial step in fine-tuning or training a domain-specific LLM is constructing a domain dataset that bridges the gap between plain language which LLMs understand well, and domain-specific data. For parametric CAD model generation, this involves annotating 3D CAD models with text descriptions. Prior work has utilized VLMs to generate simple text labels or brief descriptions for training datasets. However, we believe that more detailed, structured textual descriptions of 3D shapes are essential for effective LLM fine-tuning, an aspect underexplored in previous studies.

Using VLMs for comprehensive CAD model annotations presents challenges, as a single prompt often fails to capture both fine-grained geometric properties and compositional relationships. To address this, we propose a two-
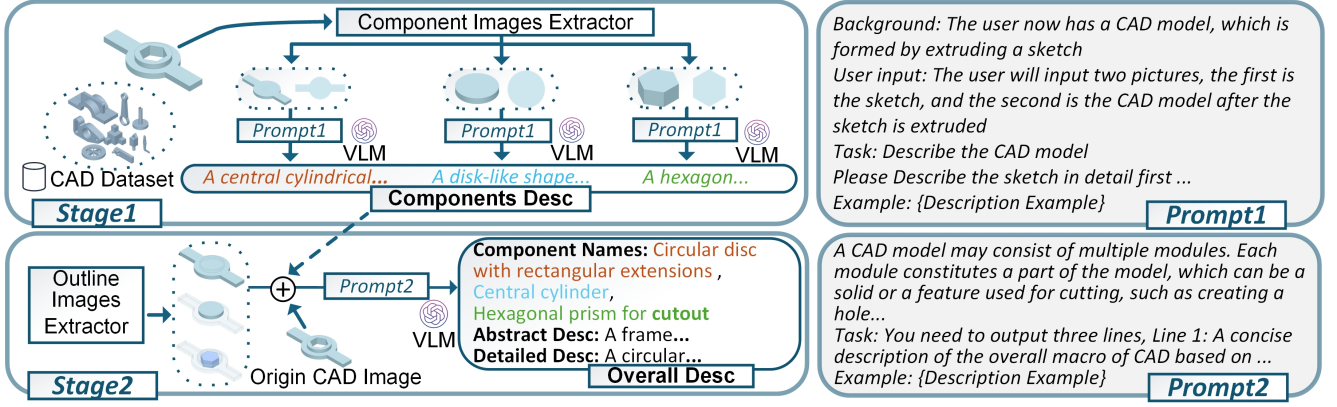
Figure 3. **Hierarchical Annotation Pipeline**. The figure illustrates our two-stage annotation process for CAD models. In the first stage, detailed descriptions of individual components are generated. In the second stage, a global description is produced, which includes both an abstract overview and detailed descriptions that capture the spatial relationships between components.

stage hierarchical annotation approach, as illustrated in Figure 3.

**Component Description Annotation** The first stage focuses on the detailed description of individual components. Formally, for the $i$-th component $\mathcal{P}_j^i$ of $\mathcal{D}_j$, we first produce the 3D image $I_j^i$ (obtained by projecting the 3D model into the 2D image) and the 2D sketch image $\hat{I}_j^i$ of $\mathcal{P}_j^i$ (obtained by rendering the corresponding 2D sketch commands). We then feed these images into VLMs (GPT-4o [2] used in our experiment), generating a detailed description $\mathcal{I}_j^i$ of the $i$-th component based on a pre-designed prompt:

$$\mathcal{I}_j^i = \text{VLM}(\text{prompt}_1, I_j^i, \hat{I}_j^i), \tag{1}$$

where $\text{prompt}_1$ is the pre-designed prompt used in stage one. By applying the above process to each component, we obtain detailed descriptions of all components $\mathcal{I}_j = \{\mathcal{I}_j^1, \mathcal{I}_j^2, \ldots, \mathcal{I}_j^m\}$. Additionally, we also include additional parameter information in the prompt, such as the extrusion direction and extrusion length. Taking component 1 in Figure 3 as an example, the generated description is: "*A central cylindrical disk, with two symmetrically positioned rectangular blocks extending from opposite sides of the disk's circumference, forming a shape that resembles a circular center with bar-like extensions, extruded upwards with an extrusion length of 12 units*".

**Overall Description Annotation** The second stage focuses on global descriptions, which include an abstract overview as well as a detailed description that explicitly addresses the spatial relationships and assembly process of the components. Additionally, since the global and local descriptions are obtained in different stages, there may be some semantic discontinuity. To bridge this gap, we let the VLM (GPT-4o) generate a short name for each component to link global and local descriptions. Specifically, for $m$ components in $\mathcal{D}_j$, we first generate its outline images $\dot{I}_j = \{\dot{I}_j^1, \dot{I}_j^2, \ldots, \dot{I}_j^m\}$ by enhancing the visibility of the target component and increasing the transparency of other components to clearly emphasize its specific location within

the CAD model. Components used for *Cutting* are rendered in blue. We then input these outline images $\dot{I}_j$, the original CAD image $I_j$, the descriptions for each module obtained in the first stage $\mathcal{I}_j$, and the $\text{prompt}_2$ used in the second stage into the VLM (GPT-4o) to generate the desired descriptions:

$$\mathcal{A}_j, \mathcal{T}_j, \mathcal{S}_j = \text{VLM}(\text{prompt}_2, I_j, \mathcal{I}_j, \dot{I}_j), \tag{2}$$

where $\mathcal{A}_j$ and $\mathcal{T}_j$ are the overall abstract and detailed descriptions, respectively, and $\mathcal{S}_j = \{\mathcal{S}_j^1, \mathcal{S}_j^2, \ldots, \mathcal{S}_j^m\}$ represents the short names for each component. For CAD models with a single component, the first-stage description serves as the final description.

These local and global hierarchical descriptions can be seamlessly integrated with the CAD data, which is designed similarly with a hierarchical structure, as detailed in Section 3.3. To enhance the stability and adaptability of VLM outputs to varying CAD model complexities, we classify CAD sequences into five complexity levels based on their length, providing 50 high-quality examples per level. All prompts employ a two-shot approach [5], selecting two examples from the corresponding level based on the complexity of the CAD model. This strategy reduce hallucinations [18, 28, 51] and improve the overall output quality.

### 3.3. SPCC Data Synthesis

Inspired by the considerable capabilities of LLMs in code generation [14, 35, 57], as well as some studies [53, 66] have attempted to convert various data types into a unified code format, we first convert parametric CAD sequences into a unified code format, as illustrated in the left part of Figure 2. Similarly to [66], we represent each sketch as a list of loops (e.g., `sketch_i.append(loop1)`), where each loop can call methods like *Line*, *Arc*, or *Circle* to draw. (e.g.,`loop1.Arc(endpoint=(87,-8),degrees=134,counterclockwise=True)`) Finally, the extrusion is performed by referencing the corresponding sketch to complete the operation. For the continuous parameters of the coordinates, we use the original 8-bit quantized

parameters from [58], where the starting point is set to (128, 128). To provide a more intuitive representation of scale information, we recenter the starting point to (0, 0). For angular parameters, we use discrete angle values within the range of 0 to 360 degrees. For more details, please refer to the supplementary materials.

**SPCC Corpus** Let $C()$ denote the CAD code formatting process and $F()$ represent our proposed annotation pipeline. For the CAD model $\mathcal{D}_j$, we obtain the parametric CAD code representation $C(\mathcal{D}_j) = \{C(\mathcal{P}_j^1), C(\mathcal{P}_j^2), \ldots, C(\mathcal{P}_j^m)\}$ and all necessary annotations $F(\mathcal{D}_j) = \{\mathcal{I}_j, \mathcal{A}_j, \mathcal{T}_j, \mathcal{S}_j\}$, where $\mathcal{I}_j = \{\mathcal{I}_j^1, \mathcal{I}_j^2, \ldots, \mathcal{I}_j^m\}$. Next, we integrate annotations by embedding them into specific segments of the CAD code, creating the SPCC. Specifically, for the $i$-th component $\mathcal{P}_j^i$ of $\mathcal{D}_j$, we combine its corresponding code and annotations to get its final representation: $\tilde{\mathcal{P}}_j^i = \{\text{concat}\{\mathcal{S}_j^i, \mathcal{I}_j^i\}, C(\mathcal{P}_j^i)\}$, where concat represents the concatenate operation. This process produces each component's final representation. We then add global descriptions as a prefix to obtain the final SPCC representation of $\mathcal{D}_j$, denoted as $\tilde{\mathcal{D}}_j = \{\mathcal{A}_j, \mathcal{T}_j, \tilde{\mathcal{P}}_j^1, \tilde{\mathcal{P}}_j^2, \ldots, \tilde{\mathcal{P}}_j^m\}$, resulting in the corpus $\tilde{\mathcal{D}} = \{\tilde{\mathcal{D}}_1, \tilde{\mathcal{D}}_2, \ldots, \tilde{\mathcal{D}}_N\}$ for training. Additionally, to enable LLMs to generate diverse CAD models from both detailed and abstract descriptions, we include data that contain only abstract descriptions, denoted as $\dot{\mathcal{D}}_j = \{\mathcal{A}_j, \tilde{\mathcal{P}}_j^1, \tilde{\mathcal{P}}_j^2, \ldots, \tilde{\mathcal{P}}_j^m\}$, in the final training corpus, represented as $\dot{\mathcal{D}} = \{\dot{\mathcal{D}}_1, \dot{\mathcal{D}}_2, \ldots, \dot{\mathcal{D}}_N\}$. For models with only one component, such as $\mathcal{D}_k$, we have $\tilde{\mathcal{D}}_k = \dot{\mathcal{D}}_k = \{\mathcal{I}_k^1, C(\mathcal{P}_k^1)\}$. Thus, the final SPCC corpus is $\mathcal{D}_{\text{SPCC}} = \{\tilde{\mathcal{D}}, \dot{\mathcal{D}}\}$.

**Multitask Instructional Dataset** To adapt CAD-Llama for downstream tasks, we construct a suite of CAD-centric instructional datasets, including *text-to-CAD*, *completion*, *caption* (CAD description generation), *addition*, and *deletion*. Table 2 presents detailed information about each task, including task descriptions, inputs, and outputs. Figure 6 also provides an example of CAD-related tasks, demonstrating how this series aids designers in continuously optimizing the model, from initial construction to iterative refinement. For *completion*, we use the initial portion (approximately 30% to 50% in our experiments) of $\tilde{\mathcal{D}}_j$ as input. For CAD editing tasks (*addition* and *deletion*), not all operations are logically valid. For example, in a CAD model consisting of a solid component and a cutting component, deleting the solid component while retaining only the cutting component is illogical. To effectively construct CAD editing instruction data using SPCC, we employ GPT-4o to identify the removable component $k$ within $\mathcal{D}_j$, explicitly justify the logical validity of its deletion, and generate corresponding deletion and inverse-addition instructions. We then remove module $k$ from $\mathcal{D}_j$, obtaining the edited CAD

| Task Name | Definition | Input | Output |
|---|---|---|---|
| *Text-to-CAD* | Given a description, generate SPCC. | $\mathcal{A}_j, \mathcal{T}_j$ | $\tilde{\mathcal{D}}_j$ |
| *Caption* | Given the CAD code, generate SPCC, which incorporates hierarchical descriptions. | $C(\mathcal{D}_j)$ | $\tilde{\mathcal{D}}_j$ |
| *Completion* | Given the partial SPCC, complete the missing part. | Partial($\tilde{\mathcal{D}}_j$) | $\tilde{\mathcal{D}}_j$ |
| *Addition* | Given the CAD code and an instruction, add a specific component to the CAD model. | $C(\mathcal{D}_{j\setminus\mathcal{P}_j^k})$ | $C(\mathcal{D}_j)$ |
| *Addition** | Given the SPCC and an instruction, add a specific component to the CAD model. | $\tilde{\mathcal{D}}_{j\setminus P_k}$ | $\tilde{\mathcal{D}}_j$ |
| *Deletion* | Given the CAD code and an instruction to delete a specific component, remove the component. | $C(\mathcal{D}_j)$ | $C(\mathcal{D}_{j\setminus\mathcal{P}_j^k})$ |
| *Deletion** | Given the SPCC and an instruction to delete a specific component, remove the component. | $\tilde{\mathcal{D}}_j$ | $\tilde{\mathcal{D}}_{j\setminus P_k}$ |

Table 2. The overview of CAD-related tasks.

model $\mathcal{D}_{j\setminus\mathcal{P}_j^k}$. Using both $\mathcal{D}_j$ and $\mathcal{D}_{j\setminus\mathcal{P}_j^k}$ along with the instructions, we construct the dataset for *addition* and *deletion*.

For the *addition* and *deletion* tasks, both input and output CAD representations are provided as CAD code, which lacks hierarchical descriptions. To demonstrate that SPCC enhances CAD editing performance and that CAD-Llama effectively understands the inherent structure of CAD Code, we designed two variant tasks: *deletion** and *addition**. During training, both input and output CAD models are ground-truth SPCC. During testing, we first use CAD-Llama-INS (instruction-tuned version of CAD-Llama) to caption the input CAD code, and the resulting SPCC serves as the final input CAD model. Taking *deletion** as an example, inputs and outputs at different stages are as follows:

(Train)　　　　　Input: $\tilde{\mathcal{D}}_j \rightarrow$ Output: $\tilde{\mathcal{D}}_{j\setminus P_k}$

(Test)　Input: $C(\mathcal{D}_j) \xrightarrow[\text{CAD-Llama-INS}]{\textbf{Caption}} \tilde{\mathcal{D}}_j \rightarrow$ Output: $\tilde{\mathcal{D}}_{j\setminus P_k}$

### 3.4. Training

**SPCC-Adaptive Pretraining** We select LLaMA3-8B [10] as our foundational LLM and conduct SPCC-adaptive pretraining on this LLM using the SPCC corpus. The traditional pretraining method creates input contexts by randomly concatenating pretraining data. However, in the same context, the preceding documents do not offer any assistance in predicting the content of the following document. Some CAD models have only minor differences, such as a change in the position of a component. To enable LLMs to capture these differences between similar CAD models for more efficient learning, similar to [50], we group similar CAD models together, so that each input context contains similar CAD models. Specifically, we use a pretrained CLIP [44] model to map each CAD model $\mathcal{D}_j \in \mathcal{D}$ to an embedding $\mathbf{E}(I_j)$ based on its image $I_j$. Then, we calculate the similarity between pairs of CAD models using cosine similarity:

$$s(\mathcal{D}_i, \mathcal{D}_j) = \cos(\mathbf{E}(\mathcal{D}_i), \mathbf{E}(\mathcal{D}_j)). \qquad (3)$$

Subsequently, we construct a CAD document graph based on the similarities and build input contexts for pretraining

by traversing this graph. After deriving the final input contexts through the aforementioned methods, SPCC-adaptive pretraining optimizes a standard autoregressive language modeling objective, which maximizes the conditional probabilities of each token given its preceding tokens as context. Formally, given an input context represented by tokens $\mathcal{X} = \{x_0, x_1, \ldots, x_{n-1}, x_n\}$, CAD-Llama applies this objective by maximizing the following log-likelihood:

$$\mathcal{L}(\mathcal{X}) = \sum_{i=1}^{n} \log P(x_i | x_{i-1}, x_{i-2}, \ldots, x_0; \Phi), \quad (4)$$

where $n$ is the context window size, $x_n$ is the special token `<|end_of_text|>`, and $\Phi$ denotes the model parameters. After pretraining, the model is equipped with essential capabilities for generating and understanding SPCC, and we name this model CAD-Llama.

**CAD-centric Instruction Tuning** Given the CAD-centric multitask instructional dataset $D = \{(X_i, Y_i)\}_{i=1}^{N}$, where $X_i$ represents the input along with the corresponding instruction description, and $Y_i$ represents the corresponding output, we fine-tune CAD-Llama on this dataset, employing LoRA [16] for parameter-efficient tuning, with the objective of maximizing the following log-likelihood:

$$\mathcal{L}(D) = \sum_{i=1}^{N} \log P(Y_i \mid X_i; \Theta), \quad (5)$$

where $\Theta$ is the trainable parameters of CAD-Llama. After this process, we obtained CAD-Llama-INS. The experiments in the following section demonstrate that the CAD-related instruction tuning process enhances the model's performance on a series of downstream tasks.

## 4. Experiments

In this section, we present the details of the experiments and the experimental results to evaluate the performance of our proposed method.

### 4.1. Experimental Setups

**Datasets** In our experiment we adopt DeepCAD[58] dataset, which contains approximately 178K parametric CAD models. We observed that many simple CAD models (e.g., cubes) in the dataset may introduce repetitive patterns, potentially degrading performance [19, 23]. We removed most of this data and applied a similar de-duplication method from [62, 63], leaving approximately 100K CAD models for training. The training data is processed using the method described in Section 3.3 to obtain the pretraining corpora for CAD-Llama. During the instruction tuning phase, we sampled 12K entries from each task in the training set to construct an instruction dataset, resulting in 48K entries.

**Metrics** For unconditional generation, we used metrics from [34, 58, 62, 63], which include: (1) Coverage (*COV*)

measures how well the generative model covers the real data distribution. (2) Minimum Matching Distance (*MMD*) calculates the minimum distance between generated samples and real samples. (3) Jensen-Shannon Divergence (*JSD*) quantifies the similarity between the distributions. (4) The success ratio ($S_R$) assesses the proportion of valid generated CAD models. (5) The *Novel* score quantifies the proportion of generated CAD sequences that do not appear in the training set.

For *text-to-CAD* task, the metrics include: (1) the accuracy of CAD model reconstructions $ACC_T$ [25], consists of command accuracy $ACC_{cmd}$, parameter accuracy $ACC_{param}$ [58], and success ratio ($S_R$), with these metrics combined to compute an overall accuracy: $ACC_T = \frac{1}{2}\left(\frac{ACC_{cmd}+ACC_{param}}{2} + S_R\right)$ (2) Median Chamfer Distance (*MCD*). (3) *MMD* and (4) *JSD*.

For CAD captioning, we use *BLEU* [39], *Rouge* [27] to measure the closeness of generated captions to reference captions. For the deletion task, we use Exact Match (*EM*) to evaluate whether the generated CAD model matches the ground truth. For the addition task, we use $ACC_{cmd}$ and $ACC_{param}$ to evaluate the accuracy of the added modules.

**Implementation details** We select LLaMA3-8B-HF [10] as our backbone. The learning rate is set to 2e-5 with the AdamW optimizer [31], and a linear learning rate warm-up is applied. The size of the context window is 2048 during SPCC-adaptive pretraining and 4096 during instruction tuning. To improve training efficiency, we use Deep-Speed [45], Flash-Attention [9]. Furthermore, we perform full fine-tuning during pretraining and use LoRA [16] for parameter-efficient training in instruction tuning, using a rank of 256 and a *lora_$\alpha$* value of 128.

**Baselines** We compare our method with a series of baseline methods. For unconditional generation, this includes parametric CAD generation models DeepCAD [58], Skex-Gen [62] and HNC-CAD [63]; For CAD-related downstream tasks, our baseline models include the open-source LLMs LLaMA3-8B [10] and Mistral-7B [17], as well as the proprietary models GPT-4 [2] and GPT-3.5 [38]. For the text-to-CAD task, our baselines also include CAD-Translator [25] and Text2CAD [20], both of which are based on the text-to-CAD transformer architecture.

### 4.2. Unconditional Generation

We use the pretrained CAD-Llama for unconditional generation, prompted by the common prefix in SPCC format: "Description of the CAD model". Each method generates 9,000 samples, with 2,000 points sampled for each one, which are compared to 3,000 randomly selected samples from the test set. Table 3 presents the main results on unconditional generation. For *COV*, CAD-Llama achieves results comparable to HNC-CAD, indicating that after pretraining on the SCPP corpus, CAD-Llama has developed
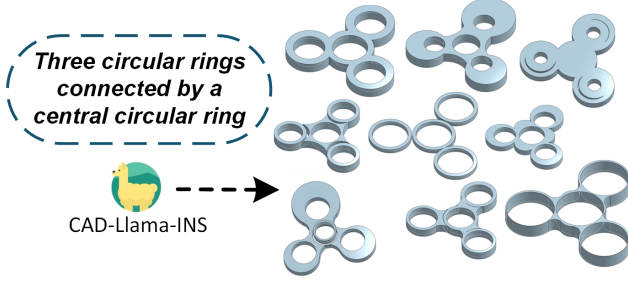
Figure 4. Qualitative example demonstrating CAD-Llama-INS generates CAD models that both conform to the description and exhibit diversity based on an abstract text description.

| Method | COV ↑ | MMD ↓ | JSD ↓ | $S_R$ ↑ | Novel ↑ |
|--------|-------|-------|-------|---------|---------|
| DeepCAD [58] | 66.68 | 1.19 | 2.59 | 61.84 | 91.7 |
| SkexGen [62] | 77.42 | 1.07 | 0.93 | 72.26 | **99.1** |
| HNC-CAD [63] | **80.46** | 0.98 | 0.74 | 79.11 | 93.9 |
| CAD-Llama (Ours) | 79.93 | **0.96** | **0.66** | **99.90** | 97.1 |

Table 3. Results on unconditional generation. We present the test set *Coverage* (**COV**) of generated CAD sequences, *Minimum Matching Distance* (**MMD**), *Jensen-Shannon Divergence* (**JSD**), *Success Ratio* (**$S_R$**) and **Novel** score. **COV**, **$S_R$** and **Novel** score are multiplied by 100%. **JSD** and **MMD** are multiplied by $10^2$. ↑: the higher the better, ↓: the lower the better.

the capability to generate diverse CAD models. In *MMD* and *JSD*, CAD-Llama demonstrates superior performance with scores of 0.96 and 0.66, indicating a narrower distribution over the target space. For $S_R$, CAD-Llama achieves the highest value of 99.90, indicating highly stable results, surpassing the other three transformer-based methods, which exhibit significantly lower $S_R$ values. Figure 4 qualitatively illustrates that, given an abstract description, CAD-Llama-INS has the ability to generate CAD models that are both consistent with the text and diverse in nature, providing wide range of options and offering inspiration. Additionally, Figure 5 shows the unconditional generation results of CAD-Llama, demonstrating the model's ability to generate CAD models of varying complexity and diversity.

## 4.3. Main Results on Text-to-CAD Task

In the *text-to-CAD* task, CAD-Llama-INS demonstrated superior performance compared to the transformer-based baseline methods, as well as GPT, LLaMA3, and others. As shown in Table 4, CAD-Llama-INS surpassed CAD-Translator and Text2CAD in accuracy by approximately 14% and significantly outperformed other LLMs. This demonstrates the efficacy of our approach in leveraging LLMs to produce CAD models that more accurately reflect textual descriptions. Furthermore, our method demonstrated substantial improvements over the baselines in metrics such as *MCD*, *MMD*, and *JSD*, indicating a closer geometric alignment with ground truth. These results underscore the limitations of transformer-based method, which,
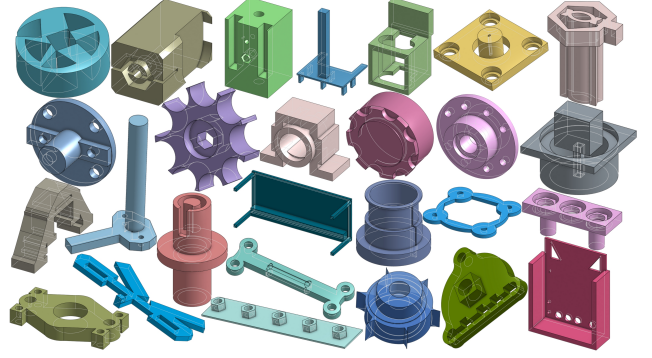


Figure 5. The unconditional generation results of CAD-Llama, demonstrate a wide range of complexity and diverse outputs.

despite their ability to predict corresponding commands, often struggle with accurately predicting the parameters essential for the precision of parameterized CAD models.

| Method | $ACC_T$ ↑ | MCD ↓ | MMD ↓ | JSD ↓ |
|--------|-----------|-------|-------|-------|
| GPT-3.5 [38] | 20.39 | 32.56 | 3.27 | 13.60 |
| GPT-4 [2] | 20.03 | 25.62 | 3.33 | 18.09 |
| LLaMA3 [10] | 17.26 | 17.33 | 4.10 | 12.36 |
| Mistral [17] | 13.12 | 32.79 | 4.71 | 18.42 |
| Text2CAD [20] | 69.91 | 20.64 | 3.02 | 9.98 |
| CAD-Translator [25] | 70.36 | 21.29 | 2.94 | 10.92 |
| CAD-Llama-INS (Ours) | **84.72** | **10.53** | **1.54** | **3.59** |

Table 4. Results on *text-to-CAD* task. The metric $ACC_T$ is multiplied by 100%. **MCD**, **MMD**, and **JSD** are multiplied by $10^2$.

## 4.4. Main Results on CAD-related Downstream Tasks

We evaluate CAD-Llama-INS and baselines on CAD-related tasks, with baselines in a two-shot setting. The main results are presented in Table 5. CAD-Llama-INS achieved an average score of 63.58%, surpassing GPT-4 by 15.7% and outperforming LLaMA3 and Mistral by approximately 30%. For all tested CAD-related tasks, CAD-Llama-INS outperformed almost all baseline LLMs. This indicates that fine-tuning on SPCC corpus significantly enhances the understanding and generation capabilities of LLMs for CAD.

For the *deletion** and *addition** tasks, CAD-Llama-INS significantly improved performance across all methods. Following structured annotation, GPT-4 leveraged its natural language reasoning capabilities to accurately identify modules for deletion, outperforming CAD-Llama-INS in the delete task. However, it struggles with the *addition** task, which requires generating CAD parameters. The experimental results indicate that SPCC offers a clear logical structure and semantic clarity, which improves the understanding and generation of LLMs. This also shows that
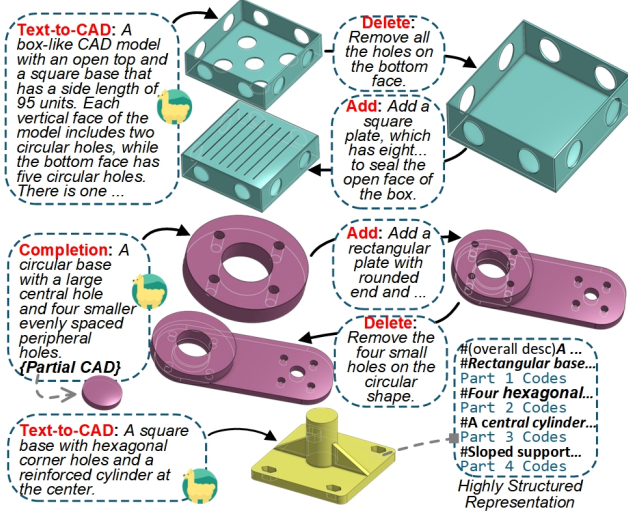
Figure 6. Examples of CAD-related tasks by using CAD-Llama-INS: the highly structured results with explicit annotations, with this series of tasks, aid designers in continuously optimizing the model, from initial construction to iterative refinement. For more detailed examples, please refer to the supplementary materials.

| Tasks | Metric | GPT-4 | GPT-3.5 | LLaMA3 | Mistral | CAD-Llama-INS |
|---|---|---|---|---|---|---|
| *Caption* | BLEU@1 | 27.34 | 25.03 | 22.41 | 21.26 | **43.44** |
| | BLEU@4 | 3.39 | 2.44 | 2.87 | 2.36 | **13.88** |
| | Rouge | 30.23 | 28.82 | 24.07 | 25.27 | **47.32** |
| *Completion* | $ACC_{cmd}$ | 51.18 | 42.98 | 22.68 | 27.30 | **73.87** |
| | $ACC_{param}$ | 38.96 | 36.29 | 35.90 | 28.56 | **57.14** |
| *Addition* | $ACC_{cmd}$ | 65.01 | 43.12 | 27.13 | 31.13 | **79.41** |
| | $ACC_{param}$ | 40.67 | 35.10 | 22.39 | 25.39 | **63.09** |
| *Addition*[*] | $ACC_{cmd}$ | 69.96 (+4.95) | 67.30 (+24.18) | 27.90 (+0.77) | 42.85 (+11.72) | **84.89** (+5.48) |
| | $ACC_{param}$ | 43.42 (+2.75) | 41.06 (+5.96) | 36.84 (+14.45) | 34.95 (+9.56) | **64.87** (+1.78) |
| *Deletion* | EM | 66.20 | 31.80 | 34.75 | 30.92 | **81.93** |
| *Deletion*[*] | EM | **90.41** (+24.21) | 53.03 (+21.23) | 43.69 (+8.94) | 39.81 (+8.89) | 89.55 (+7.62) |
| Average | / | 47.88 | 36.99 | 27.33 | 28.16 | **63.58** |

Table 5. Results (%) on multiple CAD-related tasks. *Deletion*[*] and *addition*[*] indicate the results of first using CAD-Llama-INS to generate SPCC, followed by Delete and Add edits. More experimental results can be found in the supplementary materials.

CAD-Llama-INS, after pretraining on the SPCC corpus, effectively interprets intrinsic structural information. Figure 6 presents two examples of *text-to-CAD* and a range of CAD-related tasks using CAD-Llama-INS.

### 4.5. Ablation Studies

Training data is crucial for the pretraining and fine-tuning of LLMs in this domain. We evaluate the impact of different representations of parametric CAD model training data on the *text-to-CAD* task. The evaluation methods are categorized based on whether the CAD data is represented in code format or as its original command sequence, and

| Method | $ACC_{cmd}\uparrow$ | $ACC_{param}\uparrow$ | $S_R\uparrow$ | $MCD\downarrow$ | $MMD\downarrow$ | $JSD\downarrow$ |
|---|---|---|---|---|---|---|
| *SDCS* | 39.17 | 25.56 | 18.14 | 18.03 | 2.29 | 6.30 |
| *SDCC* | 42.62 | 27.13 | 21.46 | 17.19 | 2.37 | 6.13 |
| *SPCS* | 73.13 | 47.32 | 98.71 | 14.08 | 1.64 | 3.79 |
| *SPCC* | **80.41** | **59.09** | **99.30** | **10.53** | **1.54** | **3.59** |

Table 6. Evaluation of different CAD representation methods in the Text-to-CAD task. **SDCS** uses a single textual description as a prefix along with CAD command sequences, while **SDCC** uses CAD code with a single description. **SPCS** incorporates hierarchical descriptions with CAD command sequences, and **SPCC** combines hierarchical descriptions with CAD code.

whether hierarchical or single descriptions of the 3D shape and geometry are used: (1) Single Description with CAD Sequences (*SDCS*) uses CAD command sequences with a single-prefix description that encompasses both general details and components information. (2) Single Description with CAD Code (*SDCC*) uses CAD code with the single-prefix description. (3) Structured Parametric CAD Sequences (*SPCS*) uses CAD command sequences with hierarchical descriptions. (4) Structured Parametric CAD Code (*SPCC*) uses CAD code with hierarchical descriptions. For more details, please refer to the supplementary materials.

The experimental results in Table 6 show that the SPCC method outperforms all other methods in the metrics, followed by the SPCS method. In contrast, the SDCS and SDCC methods, underperformed by approximately 30-40% in $ACC_{cmd}$ and $ACC_{param}$. These findings highlight the significant advantage of using hierarchical descriptions in improving LLMs' ability to comprehend and generate CAD models, resulting in more accurate CAD generation. Additionally, representing CAD sequences in code format further enhances performance. The structured CAD representation approach, which incorporates hierarchical descriptions, shows a significant high value $S_R$, indicating a substantial increase in the stability of CAD generation. In contrast, single-description methods are notably less effective in generating valid CAD models.

## 5. Conclusion

We introduce a novel paradigm that leverages the generative prior of LLMs into generating parametric CAD sequences. A hierarchical annotation pipeline is proposed to infuse textual descriptions of visual semantics and 3D shape at different levels of each CAD model via VLMs. A supervised fine-tuning paradigm is proposed to grant LLMs of general understanding and generation ability on parametric CAD models. An instruction tuning paradigm is proposed to fit into different downstream tasks of CAD model editing and operations. Experimental results show the superiority of our methods over traditional autoregressive methods as well as prevailing LLM baselines. In the future, with larger parameters and richer corpus, we believe that more exciting

results of LLMs for CAD will appear.

# References

[1] Opencascade. https://www.opencascade.com/. Accessed: 20-Oct-2023. 2

[2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 4, 6, 7

[3] Akshay Badagabettu, Sai Sravan Yarlagadda, and Amir Barati Farimani. Query2CAD: Generating CAD models using natural language queries. *arXiv preprint arXiv:2406.00144*, 2024. 2, 3

[4] Dorit Borrmann, Jan Elseberg, Kai Lingemann, and Andreas Nüchter. The 3d hough transform for plane detection in point clouds: A review and a new accumulator design. *3D Research*, 2(2):1–13, 2011. 2

[5] Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020. 4

[6] Sébastien Bubeck, Venkat Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. Large language models in medicine. *Nature Medicine*, 29:1936–1944, 2023. 1, 3

[7] Rohan Chabra, Jan Eric Lenssen, Eddy Ilg, Tanner Schmidt, Julian Straub, Steven Lovegrove, and Richard Newcombe. Deep local shapes: Learning local sdf priors for detailed 3d reconstruction. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 608–625, 2020. 2

[8] Katherine M. Collins, Albert Q. Jiang, Simon Frieder, Lionel Wong, Miri Zilka, Umang Bhatt, Thomas Lukasiewicz, Yuhuai Wu, Joshua B. Tenenbaum, William Hart, Timothy Gowers, Wenda Li, Adrian Weller, and Mateja Jamnik. Evaluating language models for mathematics through interactions. *Proceedings of the National Academy of Sciences*, 120(24):e2318124121, 2023. 1, 3

[9] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*, 35:16344–16359, 2022. 6

[10] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. 5, 6, 7

[11] Richard O Duda and Peter E Hart. Use of the hough transformation to detect lines and curves in pictures. *Communications of the ACM*, 15(1):11–15, 1972. 2

[12] Noelia Ferruz and Birte Höcker. Controllable protein design with language models. *Nature Machine Intelligence*, 4(6):521–532, 2022. 1, 3

[13] Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, 1981. 2

[14] Shuzheng Gao, Xin-Cheng Wen, Cuiyun Gao, Wenxuan Wang, Hongyu Zhang, and Michael R Lyu. What makes good in-context demonstrations for code intelligence tasks with llms? In *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 761–773. IEEE, 2023. 4

[15] Zhongpai Gao. Learning continuous mesh representation with spherical implicit surface. *arXiv preprint arXiv:2301.04695*, 2023. 2

[16] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 6

[17] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023. 6, 7

[18] Chaoya Jiang, Haiyang Xu, Mengfan Dong, Jiaxing Chen, Wei Ye, Ming Yan, Qinghao Ye, Ji Zhang, Fei Huang, and Shikun Zhang. Hallucination augmented contrastive learning for multimodal large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27036–27046, 2024. 4

[19] Nikhil Kandpal, Eric Wallace, and Colin Raffel. Deduplicating training data mitigates privacy risks in language models. In *International Conference on Machine Learning*, pages 10697–10707. PMLR, 2022. 6

[20] Mohammad Sadil Khan, Sankalp Sinha, Talha Uddin Sheikh, Didier Stricker, Sk Aziz Ali, and Muhammad Zeshan Afzal. Text2cad: Generating sequential cad models from beginner-to-expert level text prompts. *arXiv preprint arXiv:2409.17106*, 2024. 1, 2, 6, 7

[21] 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 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9601–9611, 2019. 2

[22] Joseph George Lambourne, Karl Willis, Pradeep Kumar Jayaraman, Longfei Zhang, Aditya Sanghi, and Kamal Rahimi Malekshan. Reconstructing editable prismatic cad from rounded voxel models. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–9, 2022. 2

[23] Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. Deduplicating training data makes language models better. *arXiv preprint arXiv:2107.06499*, 2021. 6

[24] Pu Li, Jianwei Guo, Xiaopeng Zhang, and Dong-Ming Yan. Secad-net: Self-supervised cad reconstruction by learning sketch-extrude operations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16816–16826, 2023. 2

[25] Xueyang Li, Yu Song, Yunzhong Lou, and Xiangdong Zhou. CAD Translator: An effective drive for text to 3d parametric computer-aided design generative modeling. In *Proceedings of the 32nd ACM International Conference on Multimedia (MM 2024)*, Melbourne, Australia, 2024. 1, 2, 6, 7

[26] Xingang Li, Yuewan Sun, and Zhenghui Sha. LLM4CAD: Multi-modal large language models for 3d computer-aided design generation. In *Proceedings of the ASME 2024 International Design Engineering Technical Conferences & Computers and Information in Engineering Conference (IDETC/CIE 2024)*, Washington, DC, USA, 2024. 2, 3

[27] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004. 6

[28] Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Aligning large multi-modal model with robust instruction tuning. *arXiv preprint arXiv:2306.14565*, 2023. 4

[29] Shengchao Liu, Jiongxiao Wang, Yijin Yang, Chengpeng Wang, Ling Liu, Hongyu Guo, and Chaowei Xiao. Conversational drug editing using retrieval and domain feedback. In *Proceedings of the Twelfth International Conference on Learning Representations (ICLR)*, 2024. 1, 3

[30] Zhijian Liu, Haotian Tang, Yujun Lin, and Song Han. Point-voxel CNN for efficient 3d deep learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 963–973, 2019. 2

[31] I Loshchilov. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 6

[32] Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, and Dongmei Zhang. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. *arXiv preprint arXiv:2308.09583*, 2023. 1, 3

[33] Weijian Ma, Minyang Xu, Xueyang Li, and Xiangdong Zhou. MultiCAD: Contrastive representation learning for multi-modal 3d computer-aided design models. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management (CIKM 2023)*, pages 1766–1776, 2023. 2

[34] Weijian Ma, Shuaiqi Chen, Yunzhong Lou, Xueyang Li, and Xiangdong Zhou. Draw step by step: Reconstructing CAD construction sequences from point clouds via multi-modal diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 27154–27163, 2024. 1, 2, 6

[35] Yingwei Ma, Yue Liu, Yue Yu, Yuanliang Zhang, Yu Jiang, Changjian Wang, and Shanshan Li. At which training stage does code data help llms reasoning? *arXiv preprint arXiv:2309.16298*, 2023. 4

[36] Ali Madani, Bryan Krause, Eric R. Greene, Sandeep Subramanian, Benjamin P. Mohr, James M. Holton, Jose L. Olmos Jr, Ce Xiong, Zhongkai Sun, Richard Socher, James S. Fraser, and Nikhil Naik. Large language models generate functional protein sequences across diverse families. *Nature Biotechnology*, 41:25–33, 2023. 1, 3

[37] Daniel Maturana and Sebastian Scherer. VoxNet: A 3d convolutional neural network for real-time object recognition. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 922–928, 2015. 2

[38] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022. 6, 7

[39] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002. 6

[40] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 165–174, 2019. 2

[41] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017. 2

[42] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017. 2

[43] Zeju Qiu, Weiyang Liu, Haiwen Feng, Zhen Liu, Tim Z. Xiao, Katherine M. Collins, Joshua B. Tenenbaum, Adrian Weller, Michael J. Black, and Bernhard Schölkopf. Can large language models understand symbolic graphics programs? *arXiv preprint arXiv:2408.08313*, 2024. 2, 3

[44] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 5, 3

[45] Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3505–3506, 2020. 6

[46] Daxuan Ren, Jianmin Zheng, Jianfei Cai, Jiatong Li, and Junzhe Zhang. Extrudenet: Unsupervised inverse sketch-and-extrude for shape parsing. In *European Conference on Computer Vision*, pages 482–498. Springer, 2022. 2

[47] Gernot Riegler, Ali Osman Ulusoy, and Andreas Geiger. OctNet: Learning deep 3d representations at high resolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3577–3586, 2017. 2

[48] Ruwen Schnabel, Roland Wahl, and Reinhard Klein. Efficient ransac for point-cloud shape detection. In *Computer*

*graphics forum*, pages 214–226. Wiley Online Library, 2007. 2

[49] Tianchang Shen, Zhaoshuo Li, Marc Law, Matan Atzmon, Sanja Fidler, James Lucas, Jun Gao, and Nicholas Sharp. Spacemesh: A continuous representation for learning manifold surface meshes. In *SIGGRAPH Asia 2024 Conference Papers (SA Conference Papers '24)*, page 11, New York, NY, USA, 2024. ACM. 2

[50] Weijia Shi, Sewon Min, Maria Lomeli, Chunting Zhou, Margaret Li, Gergely Szilvasy, Rich James, Xi Victoria Lin, Noah A Smith, Luke Zettlemoyer, et al. In-context pretraining: Language modeling beyond document boundaries. *arXiv preprint arXiv:2310.10638*, 2023. 5, 3

[51] Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. Aligning large multimodal models with factually augmented rlhf. *arXiv preprint arXiv:2309.14525*, 2023. 4

[52] Mikaela Angelina Uy, Yen-Yu Chang, Minhyuk Sung, Purvi Goel, Joseph G Lambourne, Tolga Birdal, and Leonidas J Guibas. Point2cyl: Reverse engineering 3d objects from point clouds to extrusion cylinders. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11850–11860, 2022. 2

[53] Jianing Wang, Junda Wu, Yupeng Hou, Yao Liu, Ming Gao, and Julian McAuley. Instructgraph: Boosting large language models via graph-centric instruction tuning and preference alignment. *arXiv preprint arXiv:2402.08785*, 2024. 4

[54] Siyu Wang, Cailian Chen, Xinyi Le, Qimin Xu, Lei Xu, Yanzhou Zhang, and Jie Yang. Cad-gpt: Synthesising cad construction sequence with spatial reasoning-enhanced multimodal llms. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 7880–7888, 2025. 3

[55] Karl DD Willis, Pradeep Kumar Jayaraman, Joseph G Lambourne, Hang Chu, and Yewen Pu. Engineering sketch generation for computer-aided design. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2105–2114, 2021. 1, 2

[56] Karl DD 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 construction from human design sequences. *ACM Transactions on Graphics (TOG)*, 40(4): 1–24, 2021. 2, 3

[57] Man-Fai Wong, Shangxin Guo, Ching-Nam Hang, Siu-Wai Ho, and Chee-Wei Tan. Natural language generation and understanding of big code for ai-assisted programming: A review. *Entropy*, 25(6):888, 2023. 4

[58] 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. 1, 2, 5, 6, 7, 3

[59] Sifan Wu, Amir Khasahmadi, Mor Katz, Pradeep Kumar Jayaraman, Yewen Pu, Karl Willis, and Bang Liu. CAD-LLM: Large language model for CAD generation. In *Proceedings of the Neural Information Processing Systems (NeurIPS) 2023*, 2023. 2, 3

[60] Jingwei Xu, Zibo Zhao, Chenyu Wang, Wen Liu, Yi Ma, and Shenghua Gao. Cad-mllm: Unifying multimodality-conditioned cad generation with mllm. *arXiv preprint arXiv:2411.04954*, 2024. 3

[61] Xianghao Xu, Wenzhe Peng, Chin-Yi Cheng, Karl DD Willis, and Daniel Ritchie. Inferring cad modeling sequences using zone graphs. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6062–6070, 2021. 1, 2

[62] Xiang Xu, Karl DD Willis, Joseph G Lambourne, Chin-Yi Cheng, Pradeep Kumar Jayaraman, and Yasutaka Furukawa. Skexgen: Autoregressive generation of cad construction sequences with disentangled codebooks. In *International Conference on Machine Learning*, pages 24698–24724. PMLR, 2022. 1, 6, 7

[63] Xiang Xu, Pradeep Kumar Jayaraman, Joseph G Lambourne, Karl DD Willis, and Yasutaka Furukawa. Hierarchical neural coding for controllable cad model generation. *arXiv preprint arXiv:2307.00149*, 2023. 1, 2, 6, 7

[64] Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*, 2023. 1, 3

[65] Haocheng Yuan, Jing Xu, Hao Pan, Adrien Bousseau, Niloy J Mitra, and Changjian Li. Cadtalk: An algorithm and benchmark for semantic commenting of cad programs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3753–3762, 2024. 3

[66] Zhe Yuan, Jianqi Shi, and Yanhong Huang. OpenECAD: An efficient visual language model for editable 3d-cad design. *Computers & Graphics*, 124:104048, 2024. 2, 3, 4

[67] Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. MAmmoTH: Building math generalist models through hybrid instruction tuning. *arXiv preprint arXiv:2309.05653*, 2023. 1, 3

# CAD-Llama: Leveraging Large Language Models for Computer-Aided Design Parametric 3D Model Generation

## Supplementary Material

## A. Overview

In the supplementary material, we put forward some details about the data selection and method design. Cost analysis as well as extra experiment results are also put forward. The remaining parts are organized as follows.

- First, we provide a cost analysis on both GPU resource and GPT-4o tokens.
- Then we illustrate the format of our CAD code used throughout the pretraining and instruction tuning stage.
- After that we introduce the hierarchical annotation pipeline in detail with respect to CAD components, image extractor and two-stage prompting strategy.
- Finally we provide extra experiment results both quantitatively and qualitatively.

## B. Training Cost and GPT-4o Token Cost

Both SPCC-adaptive pretraining and instruction tuning stages are conducted on 4 A100 GPUs. Table 7 summarizes the computational costs and token consumption for these stages. For generating finetuning data, during the SPCC-adaptive pretraining stage, altogether 70 million tokens are required to comprehend the image and generate prompts hierarchically. During the instruction tuning stage, 6 million tokens are used to generate instruction data. For the consumption of GPUs, SPCC-adaptive pretraining requires 38 A100-GPU hours and processes GPT-4o 70M tokens, while instruction tuning requires 12 A100-GPU hours and processes GPT-4o 6M tokens. This demonstrates that our model achieves efficient training with limited computational resources. Notably, during the Instruction Tuning phase, the model adapts effectively to various downstream tasks using only a small amount of data and training time.

| Stage | A100-GPU Hours | Tokens (GPT-4o) |
|---|---|---|
| SPCC-Adaptive Pretraining | 38 | 70M |
| Instruction Tuning | 12 | 6M |

Table 7. Training costs and token consumption during the two training stages. Tokens are used for prompt generation in each stage.

## C. Details of CAD Code Formatting

We follow the annotations of DeepCAD [58] dataset and denote the components of the CAD command sequence. The complete set of command parameters is defined as $p_i = [x, y, \alpha, f, r, \theta, \phi, \gamma, p_x, p_y, p_z, s, e_1, e_2, b, u]$. We normalize and quantize these parameters as follows: (1) For discrete coordinate parameters, including the sketch plane origin $(p_x, p_y, p_z)$, extrusion distances $(e_1, e_2)$, curve endpoint coordinates $(x, y)$, and the circle radius $r$, we quantize all continuous parameters into 256 levels, represent them with 8-bit integers, and recenter the origin from $(128, 128)$ to $(0, 0)$ for a more intuitive representation of scale. (2) For angular parameters, including the sketch orientation angles $(\theta, \phi, \gamma)$ within the range $[-\pi, \pi]$ and the arc's sweep angle $\alpha$ within $[0, 2\pi]$, we use discrete values within the ranges $[-180, 180]$ and $[0, 360]$ degrees, respectively. (3) The sketch profile scale $s$ is constrained within the range $[0, 2]$, while the boolean operation type $b$ can take one of the following values: *new body*, *join*, *cut*, or *intersect*. The extrusion type $u$ denotes one of three configurations: *one-sided*, *symmetric*, or *two-sided*. These parameters are utilized in their original forms. (4) The arc's counterclockwise flag $f$ is a binary indicator, which we represent as either True or False.

For converting the annotation of CAD construction sequence into a LLM-friendly format, we further extract the hierarchy of CAD construction sequences and organize them into python-like pseudocode. In particular, the SOL and EOS commands are abstracted as an object `Loop()` and an ending comment `# End of code`, respectively. Other commands, such as *Arc*, *Line*, *Circle*, and *Extrude*, are represented as function calls with corresponding parameters as inputs of the function. Detailed examples are illustrated in Figure 8 and 9.

## D. Details of Hierarchical Annotation Pipeline

### D.1. Definition of CAD Component

In our definition, a CAD model consists of one or more components. Typically, a single sketch-extrude pair is treated as an individual component. However, when multiple identical sketch-extrude pairs occur consecutively in a CAD sequence, such as 10 cylinders uniformly distributed in a circular arrangement, describing each pair individually leads to redundancy and poses challenges for vision-language models (VLMs) in accurately capturing such repetitive structures.

To address this, when identical sketch-extrude pairs occur consecutively and their count exceeds a specified threshold (set to 3 in our experiments), we collectively define them as a single component. Otherwise, each sketch-
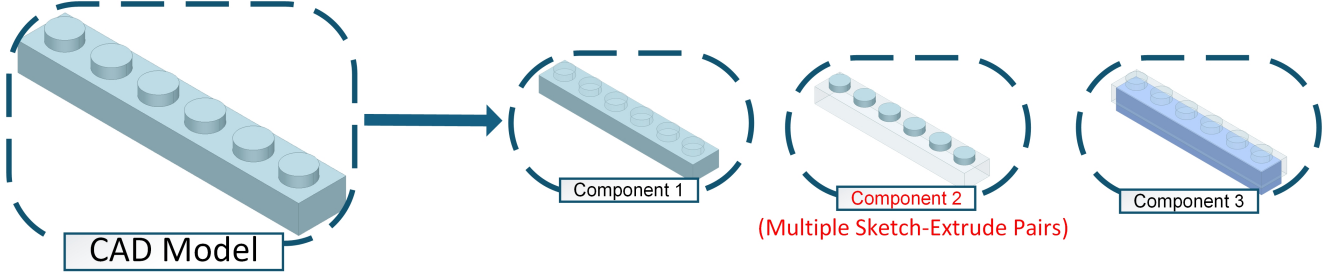
Figure 7. Illustration of defining a single component from consecutive equivalent sketch-extrude pairs based on specified criteria. Note that a large proportion of bottom cube in the final CAD model has been cut out, as is shown in component 3.

extrude pair is treated as an individual component. The equivalence of two sketch-extrude pairs is determined based on the following criteria: all commands and parameters must match, except for the sketch plane origin parameters $(p_x, p_y, p_z)$. As illustrated in Figure 7, the second component comprises multiple sketch-extrude pairs.

## D.2. Annotation Image Generation Pipeline

The hierarchical annotation pipeline contains two stages. Different images are fed into the VLM in different stages. We propose two kinds of image extractors which extracts different features of the CAD model, namely **Components Images Extractor** and **Outlines Images Extractor**, as shown in Figure 3 in the main paper. Both of them are python scripts rendering CAD construction sequences using PythonOCC [1] (Python version of OCCT) while focusing on different aspects of a single model. Taking the *i*-th component of the *j*-th CAD model $\mathcal{D}_j$ as an example, in the first stage, we use the **Components Images Extractor** and obtain the component image $I_j^i$ and its corresponding 2D sketch image $\hat{I}_j^i$. Specifically, $I_j^i$ is rendered from the default viewpoint by extracting the component's CAD command, while $\hat{I}_j^i$ is obtained by rendering the corresponding 2D sketch commands. In the second stage, we use the **Outlines Images Extractor** and obtain the outline image $\dot{I}_j^i$, achieved by increasing the transparency of other components (set to 0.85 in our experiments) while keeping the target component's transparency unchanged during rendering; if the target component is used for *cutting*, it is rendered in blue, as illustrated by component 3 in Figure 7.

## D.3. Two Stage Prompting Methods.

In this part we provide the detailed prompt used in Section D.2. In particular, two prompts are adopted where **prompt1** is for obtaining descriptions of individual components and **prompt2** is for acquiring both overall descriptions and component names. In **prompt1**, to enable GPT-4o to generate more detailed descriptions, we provide additional information that includes extrusion direction, extrusion length, and

quantity information. The extrusion direction is included only when the CAD model is extruded in a specific direction, such as up, down, left, right, front, or back. We observed that over 95% of the extrusion directions in the DeepCAD [58] dataset fall within these categories. Quantity information is added only when a component contains multiple sketch-extrude pairs (see Section D.1), which helps mitigate the hallucination phenomenon in VLM. The specific content of the two prompts is as follows:

   **Prompt1**: *Background: The user now has a CAD model, which is formed by extruding a sketch. User input: The user will input two pictures, the first is the sketch, and the second is the CAD model after the sketch is extruded. Task: Describe the CAD model. Please describe the sketch in detail first, include the additional information in the description and output the final description result as a single line. Additional information: {Extrusion direction and length information, Number information} Examples: {Two Description Example}*

   **Prompt2**: *A CAD model may consist of multiple modules. Each module constitutes a part of the model, which can be a solid or a feature used for cutting, such as creating a hole. The user has a CAD consisting of {num_parts} modules. The user will input {num_parts+1} pictures, the first image is the original CAD model, followed by {num_parts} images where each module is rendered with enhanced highlighting. These modules collectively form the original CAD model. Modules used for cutting are highlighted in blue. The subsequent description explains each of the four modules individually, following the order presented in the module images: {Component Descriptions} Task: You need to output three lines, Line 1: A concise description of the overall macro of CAD based on first image. Line 2: A detailed description that includes the specific characteristics of each of the {num_parts} modules mentioned above, as well as the process by which they are assembled based on all provided images and component descriptions. Line 3: Short names for {num_parts} modules. Example: {Two Description Example}*

# E. Details of Experiments

## E.1. Prompts Used in Baseline Methods

We provide the detailed prompt used in the baseline methods (GPT-4, GPT-3.5, LLaMA3, and Mistral) in Figure 11, where {task_definition} specifies the task instructions.

## E.2. Ablation Study Details in Main Experiment

In the ablation study of the main experiment, we explored the impact of using different CAD representations for pretraining on the Text-to-CAD task. The evaluation methods are categorized based on whether the CAD data is in code format or raw sequences, and whether hierarchical or single descriptions are used. The single description $\mathcal{SD}$ of the CAD model $\mathcal{D}_j$ is defined as:

$$\mathcal{SD} = \text{concat}\{\mathcal{A}_j, \mathcal{T}_j, \text{"Parts description:"}, D_j\}$$

where $D_j = \text{concat}\{\{\mathcal{S}_j^1, \mathcal{I}_j^1\}, \{\mathcal{S}_j^2, \mathcal{I}_j^2\}, \ldots, \{\mathcal{S}_j^k, \mathcal{I}_j^k\}\}$ represent the full description of all $k$ components of $\mathcal{D}_j$.

| Tasks | Metric | *w/o* ICP | *with* ICP |
|:---:|:---:|:---:|:---:|
| **Text-to-CAD** | $\text{ACC}_{\text{cmd}}$ | 79.89 | **80.41** |
| | $\text{ACC}_{\text{param}}$ | 59.04 | **59.09** |
| **Add** | $\text{ACC}_{\text{cmd}}$ | 77.73 | **79.41** |
| | $\text{ACC}_{\text{param}}$ | 62.16 | **63.09** |
| **Delete** | EM | 80.91 | **81.93** |

Table 8. Performance comparison of CAD-related tasks with and without In-Context Pretraining (ICP). The results show that ICP improves performance across all tasks

## E.3. Ablation Studies on Pretraining Method

This section presents a simple ablation study on several CAD-related tasks to validate the effectiveness of In-Context Pretraining (ICP) [50] in enhancing CAD-Llama-INS performance on downstream tasks. ICP is a method that groups related documents within the same input context, encouraging LLMs to read and reason across document boundaries. Similar to [50], we used a pretrained CLIP [44] model to encode CAD images and group similar CADs for pretraining based on their cosine similarity.

As shown in Table 8, ICP enhances the performance of downstream editing tasks, such as add and delete, by enabling LLMs to better capture the distinctions between different CAD structures during pretraining. This improved understanding allows the model to more effectively handle precise modifications required in these tasks. Additionally, ICL contributes to a marginal improvement in the Text-to-CAD task.

| Dataset | $\text{ACC}_T\uparrow$ | MCD↓ | MMD↓ | JSD↓ |
|:---:|:---:|:---:|:---:|:---:|
| DeepCAD | 84.72 | 10.53 | 1.54 | 3.59 |
| Fusion360 | 78.35 | 23.06 | 1.98 | 6.01 |

Table 9. Our model, CAD-Llama-INS, trained exclusively on the DeepCAD dataset, demonstrates strong generalization capabilities on the Fusion360 dataset in the Text-to-CAD task.

## E.4. Cross Dataset Generalization

To further evaluate the generalization ability of CAD-Llama-INS, we conducted experiments on the test set of the Fusion 360 [56] dataset for the Text-to-CAD task. Similar to DeepCAD [58], the Fusion 360 dataset also contains CAD construction sequences. We employed the hierarchical annotation pipeline to generate descriptions for the Fusion 360 dataset. These descriptions are used to prompt CAD-Llama-INS, which was pre-trained and fine-tuned exclusively on the DeepCAD dataset, to produce corresponding CAD models. The experimental results, as shown in Table 9, demonstrate that CAD-Llama-INS achieves strong generalization performance, achieving comparable or superior results on the Fusion 360 dataset despite being trained solely on DeepCAD. This highlights the effectiveness of our approach in adapting to new datasets. A qualitative analysis is also conducted, as illustrated in Figure 12. Based on textual prompts, CAD-Llama-INS demonstrates the capability to generate CAD models that closely align with the ground truth.

## E.5. Qualitative results

To comprehensively evaluate the performance of our approach, we provide qualitative results across multiple tasks. Specifically, qualitative results for text-to-CAD generation are illustrated in Figure 13 to Figure 23. Results for captioning tasks are presented in Figure 10, while results for unconditional generation are shown in Figure 24. Additionally, results for multi-task evaluation, encompassing the process from initial construction to iterative refinement, are shown in Figures 19 to 23.

## E.6. Examples of Failure Cases

Our experimental results also show some limitations of our method, in some cases there are parameter generation errors and mismatching between the input text instruction and the generated CAD command sequences, Figure 25 illustrates some failure cases of Text-to-CAD generation.

# Description of the model: A cylindrical disc with two rectangular cutouts on its top surface is combined with a thick-walled cylindrical ring with a hollow center. The thick-walled cylindrical ring is placed on top of the cylindrical disc, aligning the hollow center of the ring with the center of the disc, with two rectangular cutouts on the top surface of the disc.

# Cylinder base: A cylindrical base is extruded upwards to a length of 59 units.

```
sketch0 = []
loop0_0 = Loop()
loop0_0.Circle(center=(48, 0),radius=48)
sketch0.add(loop0_0)
Extrude(sketch=sketch0, orientation_degrees=(0,
0, 0), origin=(-96, 0, 0), scale=2.0,
extrude_distance=(59, 0), boolean_type='New',
extent_type='One-sided')
```

# Cylindrical extension with hole: A cylindrical shape with a central circular hole is extruded upwards to a length of 35 units.

```
sketch1 = []
loop1_0 = Loop()
loop1_0 .Circle(center=(48, 0),radius=48)
sketch1.add(loop1_0)
loop1_1 = Loop()
loop1_1.Circle(center=(48, 0),radius=34)
sketch1.add(loop1_1)
Extrude(sketch=sketch1, orientation_degrees=(0,
0, 0), origin=(-32, 0, 61), scale=0.68,
extrude_distance=(35, 0), boolean_type='Joining',
extent_type='One-sided')
```

# Rectangular cutouts: Two rectangular cutouts, positioned symmetrically on the surface of the cylindrical base and extending partially through its thickness, are extruded upwards to a length of 60 units.

```
sketch2 = []
loop2_0= Loop()
loop2_0.Line(endpoint=(55, -42))
loop2_0.Line(endpoint=(55, -37))
loop2_0.Line(endpoint=(40, -37))
loop2_0.Line(endpoint=(40, -42))
sketch2.add(loop2_0)
loop2_1 = Loop()
loop2_1.Line(endpoint=(55, 37))
loop2_1.Line(endpoint=(55, 42))
loop2_1.Line(endpoint=(40, 42))
loop2_1.Line(endpoint=(40, 37))
sketch2.add(loop2_1)
Extrude(sketch=sketch2, orientation_degrees=(0,
0, 0), origin=(-96, 0, 0), scale=2.0,
extrude_distance=(60, 0),
boolean_type='Cutting', extent_type='One-
sided')
```

# End of code



■ Abstract Description
■ Detailed Description
■ Component Names
■ Component Descriptions

# Description of the model: A hexagonal prism with a central cylindrical hole running through its length, combined with two concentric cylindrical tubes arranged horizontally. The hexagonal prism with a central cylindrical hole is extruded backward longitudinally. Then, two concentric cylindrical tubes with central circular holes are extruded backward longitudinally and arranged horizontally to fit into the hexagonal prism.

# Hexagonal prism with cylindrical hole: A hexagonal prism with a central cylindrical hole running through its length is extruded backward to a length of 96 units.

```
sketch0 = []
loop0_0 = Loop()
loop0_0.Line(endpoint=(18, -45))
loop0_0.Line(endpoint=(66, -50))
loop0_0.Line(endpoint=(95, -14))
loop0_0.Line(endpoint=(77, 31))
loop0_0.Line(endpoint=(29, 36))
loop0_0.Line(endpoint=(0, 0))
sketch0.add(loop0_0)
loop0_1 = Loop()
loop0_1.Circle(center=(48, -7), radius=24)
sketch0.add(loop0_1)
Extrude(sketch=sketch0, orientation_degrees=(90, -90, 90),
origin=(-51, 0, 8), scale=1.06, extrude_distance=(96, 0),
boolean_type='New', extent_type='One-sided')
```

# Concentric cylindrical tubes: Two concentric cylindrical tubes with central circular holes are arranged horizontally and extruded backward to a length of 26 units.

```
sketch1 = []
loop1_0 = Loop()
loop1_0.Circle(center=(48, 0), radius=48)
sketch1.add(loop1_0)
loop1_1 = Loop()
loop1_1.Circle(center=(48, 0), radius=22)
sketch1.add(loop1_1)
Extrude(sketch=sketch1, orientation_degrees=(90, -90,
90), origin=(-62, 0, 0), scale=1.29, extrude_distance=(51,
0), boolean_type='Joining', extent_type='One-sided')
sketch2 = []
loop2_0 = Loop()
loop2_0.Circle(center=(48, 0), radius=48)
sketch2.add(loop2_0)
loop2_1 = Loop()
loop2_1.Circle(center=(48, 0), radius=30)
sketch2.add(loop2_1)
Extrude(sketch=sketch2, orientation_degrees=(90, -90,
90), origin=(-27, 0, 0), scale=0.57, extrude_distance=(26,
0), boolean_type='Joining', extent_type='One-sided')
```
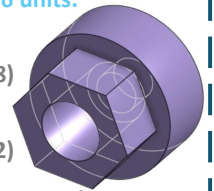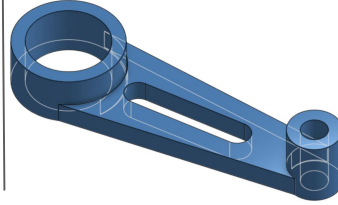
# End of code



Figure 8. Examples of SPCC data representation, generated by our CAD-Llama-INS.

# Description of the model: A tapered rectangular prism featuring an elongated oval slot through the center, with semi-circular cutouts on opposite sides, which connected to two cylindrical shapes with central circular holes. The tapered rectangular plate is positioned horizontally, and the two cylindrical shapes of different sizes are placed vertically on the semicircular cutouts at each end of the plate, aligning their central circular holes with the cutouts.
# Rectangular plate with cutouts and slot: A rectangular plate featuring two large semicircular cutouts on opposite sides and a central elongated oval slot is extruded upwards to a length of 7 units.

```
sketch0 = []
loop = Loop()
loop.Line(endpoint=(95, 10))
loop.Arc(endpoint=(94, 29), degrees=180,
counterclockwise=False)
loop.Line(endpoint=(0, 39))
loop.Arc(endpoint=(0, 0), degrees=168,
counterclockwise=False)
sketch0.add(loop)
loop = Loop()
loop.Line(endpoint=(64, 14))
loop.Arc(endpoint=(64, 24), degrees=180,
counterclockwise=True)
loop.Line(endpoint=(25, 24))
loop.Arc(endpoint=(25, 14), degrees=180,
counterclockwise=True)
sketch0.add(loop)
Extrude(sketch=sketch0, orientation_degrees=(0, 0, 0),
origin=(3, -17, 0), scale=0.89, extrude_distance=(7, 0),
boolean_type='New', extent_type='One-sided')
```

# Large cylindrical Shape: A cylindrical shape with a central circular hole, resembling a thick washer, is extruded upwards to a length of 17 units.

```
sketch1 = []
loop = Loop()
loop.Circle(center=(48, 0),radius=48)
sketch1.add(loop)
loop = Loop()
loop.Circle(center=(48, 0),radius=36)
sketch1.add(loop)
Extrude(sketch=sketch1, orientation_degrees=(0, 0, 0),
origin=(-17, 0, 0), scale=0.36, extrude_distance=(17, 0),
boolean_type='Joining', extent_type='One-sided')
```

# Small Cylindrical Shape: A cylindrical shape, featuring a central circular hole on one end and resembling a thick washer, is extruded upwards to a length of 22 units.

```
sketch2 = []
loop = Loop()
loop.Circle(center=(48, 0),radius=48)
sketch2.add(loop)
loop = Loop()
loop.Circle(center=(48, 0),radius=24)
sketch2.add(loop)
Extrude(sketch=sketch2, orientation_degrees=(0, 0, 0),
origin=(79, 0, 0), scale=0.18, extrude_distance=(22, 0),
boolean_type='Joining', extent_type='One-sided')
# End of code
```

- 🟩 Abstract Description
- 🟧 Detailed Description
- 🟪 Component Names
- 🟦 Component Descriptions

# Description of the model: A central cylinder is surrounded by six evenly spaced cylindrical pillars arranged in a circular pattern. The central cylinder is extruded vertically, and the six evenly spaced cylindrical pillars are also extruded vertically and arranged in a circular pattern around the central cylinder.
# Central Cylinder: A cylindrical shape is extruded upwards to a length of 72 units.

```
sketch0 = []
loop0_0 = Loop()
loop0_0.Circle(center=(48, 0), radius=48)
sketch0.add(loop0_0)
Extrude(sketch=sketch0, orientation_degrees=(0, 0, 0), origin=(-60, 0, 0),
scale=1.25, extrude_distance=(72, 0), boolean_type='New',
extent_type='One-sided')
```

# Circular array of cylindrical pillars: Six evenly spaced cylindrical pillars, arranged in a circular pattern, extruded upwards to a length of 96 units.

```
sketch1 = []
loop1_0 = Loop()
loop1_0.Circle(center=(48, 0), radius=48)
sketch1.add(loop1_0)
Extrude(sketch=sketch1, orientation_degrees=(0, 0, 0), origin=(-10, 60,
0), scale=0.2, extrude_distance=(96, 0), boolean_type='New',
extent_type='One-sided')
sketch2 = []
loop2_0 = Loop()
loop2_0.Circle(center=(48, 0), radius=47)
sketch2.add(loop2_0)
Extrude(sketch=sketch2, orientation_degrees=(0, 0, 0), origin=(-61, 30,
0), scale=0.2, extrude_distance=(96, 0), boolean_type='Joining',
extent_type='One-sided')
sketch3 = []
loop3_0 = Loop()
loop3_0.Circle(center=(48, 0), radius=47)
sketch3.add(loop3_0)
Extrude(sketch=sketch3, orientation_degrees=(0, 0, 0), origin=(-61,
-30, 0), scale=0.2, extrude_distance=(96, 0), boolean_type='Joining',
extent_type='One-sided')
sketch4 = []
loop4_0 = Loop()
loop4_0.Circle(center=(48, 0), radius=48)
sketch4.add(loop4_0)
Extrude(sketch=sketch4, orientation_degrees=(0, 0, 0), origin=(-10,
-60, 0), scale=0.2, extrude_distance=(96, 0), boolean_type='Joining',
extent_type='One-sided')
sketch5 = []
loop5_0 = Loop()
loop5_0.Circle(center=(48, 0), radius=47)
sketch5.add(loop5_0)
Extrude(sketch=sketch5, orientation_degrees=(0, 0, 0), origin=(42, -
30, 0), scale=0.2, extrude_distance=(96, 0), boolean_type='Joining',
extent_type='One-sided')
sketch6 = []
loop6_0 = Loop()
loop6_0.Circle(center=(48, 0), radius=47)
sketch6.add(loop6_0)
Extrude(sketch=sketch6, orientation_degrees=(0, 0, 0), origin=(42,
30, 0), scale=0.2, extrude_distance=(96, 0), boolean_type='Joining',
extent_type='One-sided')
# End of code
```
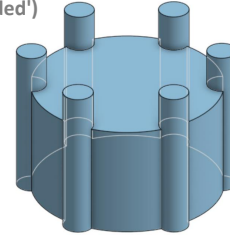
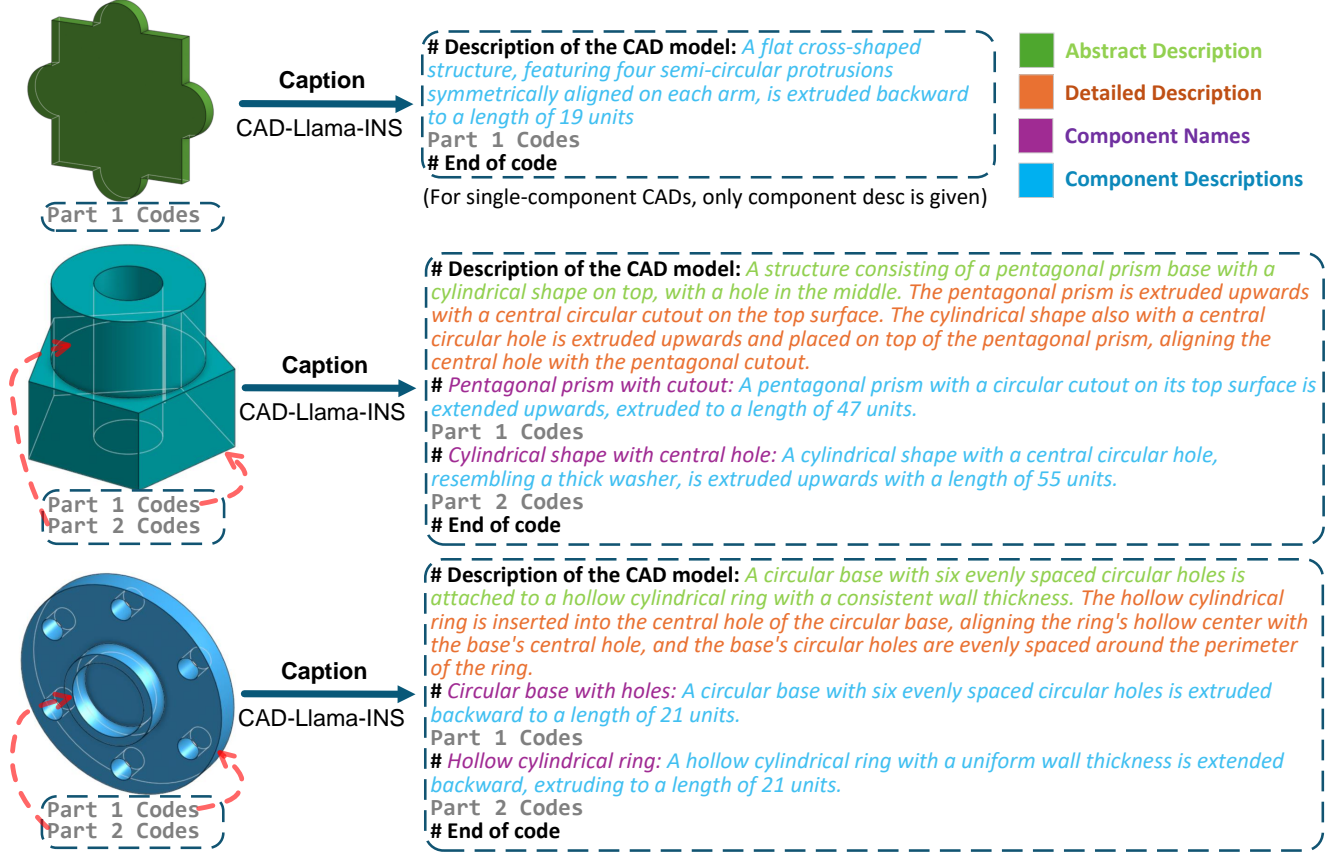Figure 9. Examples of SPCC data representation, generated by our CAD-Llama-INS.

Figure 10. Examples of results from the Caption task, demonstrating the capabilities of CAD-Llama-INS in understanding the internal structure of raw CAD code and its geometric shapes.

**Prompt for Baseline Methods (2-shot)**

*<background information>*
*Computer-Aided Design (CAD) is a technology used in engineering and manufacturing to create precise 2D and 3D models. In CAD modeling, common geometric operations include drawing lines (Line), arcs (Arc), and circles (Circle) to form sketches that define the shape of the object.*
*These sketches are usually composed of multiple loops, with each loop made up of various geometric entities like lines, arcs, and circles to form a closed surface. Once the sketch is created, an extrusion operation can be applied to extend the 2D profile into 3D space, thus forming a solid body. The extrusion process involves setting parameters such as the orientation of the sketch plane and the distance of extrusion.*
*In CAD systems, the creation of 3D objects usually begins with 2D sketches consisting of geometric entities like lines, arcs, and circles. Each of these entities is defined by specific parameters. The default starting point of each loop is (0,0).*
*`Line` command is determined by its `endpoint` coordinates.*
*`Arc` command requires its `endpoint`, a sweep angle (`degrees`), and a direction flag (`counterclockwise`) to indicate whether it is drawn clockwise or counterclockwise.*
*`Circle` command is defined by its `center` coordinates and `radius`. Once the sketch is complete.*
*`Extrude` command is used to transform the 2D sketch into a 3D object. The extrusion process involves parameters such as the orientation of the sketch plane, defined by three angles (`orientation_degrees`), and the sketch plane's position in 3D space (`origin`). The `scale` parameter controls the size of the sketch profile, while `extrude_distance` defines how far the sketch is extended in both directions. The `boolean_type` specifies how the new geometry interacts with the existing model, with options like 'New' (creating a new solid), 'Joining' (merging with the existing body), 'Cutting' (removing material), and 'Intersection' (keeping only the overlapping portion). The `extent_type` controls the extrusion direction, with options such as 'One-sided' (extending in one direction), 'Symmetric' (extending equally in both directions), or 'Two-sided' (with different distances on each side).*
*Note: All coordinate values are quantized to integers between 0 and 255, so your output coordinate value should be less than or equal to 255.*
*</background information>*
*<task>{task_definition}</task>*
*<example><input>{e1_input}</input><output>{e1_output}</output></example>*
*<example><input>{e2_input}</input><output>{e2_output}</output></example>*

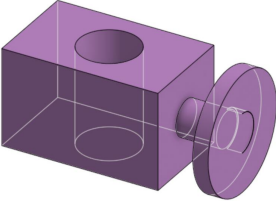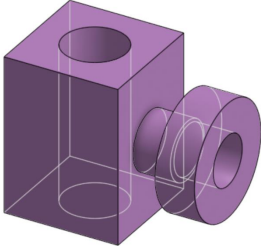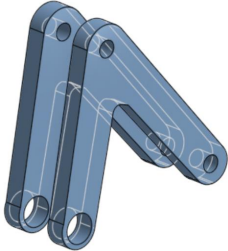Figure 11. Detailed prompt used in the baseline methods (GPT-4, GPT-3.5, LLaMA and Mistrial).
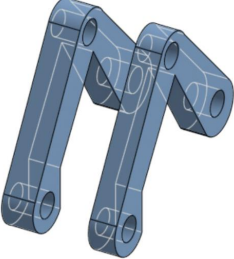
Figure 12. Comparison results of Text-to-CAD task on the Fusion 360 dataset.
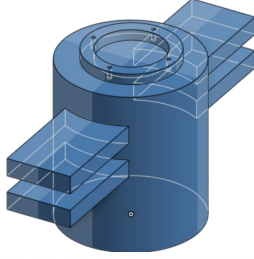
**Text-to-CAD**

*A horizontally oriented rectangular block with symmetrical notches on opposite sides, forming angular cutouts. A central circular hole runs through the length of the block, with additional cylindrical cutouts intersecting it. On either end, there are concentric cylindrical features, including a hollow cylindrical tube with a slanted edge. These cylindrical elements are mirrored on both sides of the block. The angular cutouts are formed by triangular prisms intersecting the block at a 45-degree angle, creating a stepped appearance along its sides*
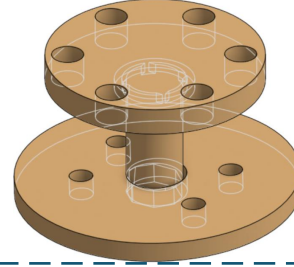
**Text-to-CAD**

*The model features a circular base with eight evenly spaced trapezoidal teeth extending radially outward from its perimeter. At the center of the base is a circular disk, which contains a hexagonal hole passing through it.*
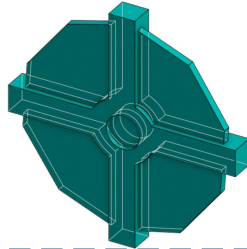
**Text-to-CAD**

*The CAD creation integrates simplified complexity. A cylinder forms the core structure, providing foundational stability and rotational symmetry. Mounted laterally, a rectangular solid block with a horizontal groove runs along its length, aiding in the precise linear arrangement of associated parts. Finally, a circular ring with four smaller, equally spaced holes is securely affixed to the top of the cylinder, providing multiple anchoring points for securing or interfacing with other components.*
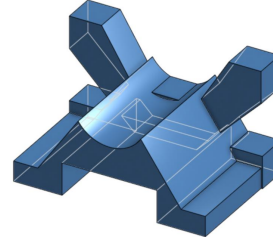
**Text-to-CAD**

*A circular disc with six holes is mounted on a cylindrical base, with a central octagonal support. The circular disc with six equally spaced holes is placed on top. Below it, a long main column serves as support, with a polygonal structure located beneath the column. A spherical ring with five holes is then placed around the base of the column.*
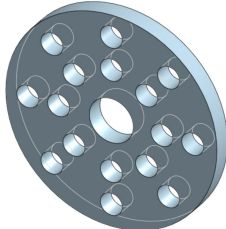
**Text-to-CAD**

*A dynamic cross-sectional framework embedded with a dodecagonal pattern. The cross-shaped structure, combined with a dodecagonal prism, aligns its central circular hole with one of the prism's faces, effectively creating an intricate framework with elongated arms supporting the dodecagonal edges.*

**Text-to-CAD**

*The base features a large rectangular hollow at its center on the bottom side. On top of the base, there are two symmetrical rectangular protrusions extending outward on either side. At the center of the base's top surface is a semi-circular concave feature that connects the structure above to the base. Rising from the base are two symmetrical sets of inclined planes, joined at the center by a smooth saddle-shaped curved surface. Four prismatic arms extend diagonally upward from the base, symmetrically positioned on either side.*

**Text-to-CAD**

*Incorporating a circular plate with a large central hole and two sets of eight smaller circular holes arranged in a ring pattern, one set located closer to the center and the other set positioned near the outer edge, symmetrically distributed around the central hole.*

**Text-to-CAD**

*The foundation is a circular plate with a central solid hub and four equally spaced rectangular cutouts extending radially towards the outer edge, creating a cross-like pattern. Below this, a cylindrical tube formed by extruding a circular ring with a central hole is mounted around the hub.*

Figure 13. Supplementary results of the Text-to-CAD task generated by CAD-Llama-INS based on text prompts.

**Text-to-CAD**
*A gear-like shape with eight rounded teeth with a cylindrical ring-like shape that has four equally spaced circular holes around its perimeter. The gear-like shape with eight evenly spaced, rounded teeth is extruded backward and assembled with the ring, forming a complex shape with multiple holes, which is also extruded backward.*

**Text-to-CAD**
*A rounded rectangular prism with semi-cylindrical ends, featuring a series of four evenly spaced grooves running longitudinally along each flat side.*

**Text-to-CAD**
*A circular base with a slightly recessed central area. At the center, there is a cylindrical hole surrounded by a rounded square recess, where the edges of the square are smoothly curved.*

**Text-to-CAD**
*A circular ring with a symmetrical arrangement of protruding segments around its outer edge. The ring has a hollow center and features six evenly spaced rectangular extensions projecting outward. Each extension has a flat top surface and sharp angular transitions to the base ring.*

**Text-to-CAD**
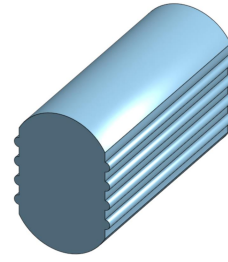*A central cylindrical body with a hollow cylindrical hole running along its axis. Extending outward symmetrically from the central body are four rectangular prismatic fins, evenly spaced at 90-degree intervals around the cylinder.*

**Text-to-CAD**
*a rounded rectangular block with four arrow-shaped features extending symmetrically inward toward the center. These arrow-like structures are recessed into the block, creating sharp-edged cutouts that converge centrally in a cross-like configuration.*

**Text-to-CAD**
*A central hexagonal plate with six evenly distributed circular loops extending outward from each side of the hexagon.*

**Text-to-CAD**
*The large circular ring with twelve evenly spaced smaller holes is centered on a rectangular base with four additional holes placed near each corner, ensuring stability and formal coherence in the combined layout.*

Figure 14. Supplementary results of the Text-to-CAD task generated by CAD-Llama-INS based on text prompts.

9

**Text-to-CAD**
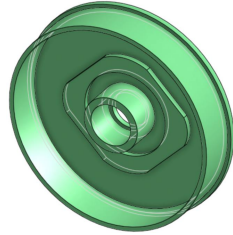*A cross-shaped plate with a cylindrical block with a star-like central hole. The cross-shaped plate, which has beveled ends on each arm, supports the cylindrical block at its center, where the star-like cutouts interlock with the central cutout of the cross-plate. .and a flower-patterned cylindrical top forms a dynamic centerpiece.*

**Text-to-CAD**
*The circular plate with a large central hole and four equally spaced rectangular cutouts is merged with a cylindrical shape having a central circular hole and four concave indentations along the outer edge, aligning the large central holes for cohesion while the indentations complement the rectangular cutouts.*

**Text-to-CAD**
*A cross-shaped structure with a large central rectangular base, featuring semi-circular and triangular cutouts, complimented by a cylindrical tube through its center. The cross shape is derived from two rectangular prisms intersecting perpendicularly, projecting from the large base with appropriate cutouts, while the cylindrical tube is integrated by centering it through the entire intersection.ng or interfacing with other components.*

**Text-to-CAD**
*A circular frame with an inner cross that divides the circle into four equal quadrants. At the center of the frame is a hexagonal prism with six equal sides and uniform height. On top of the hexagonal prism, a smaller cylinder is extruded vertically, creating a simple yet structured design.*

**Text-to-CAD**
*A shape featuring a circular ring with tabs and an circular solid with a cross pattern of holes. Attach the circular solid at the center of the circular ring so that the holes are aligned with the tabs, allowing for versatile mounting.*

**Text-to-CAD**
*An framework consists of a square frames and a square plate with a grid pattern. This structure initiates with a square frame having rounded corners. Above this configuration, a square plate with a 4x3 grid of circular cutouts is layered, aligning its perimeter with the framework's boundaries.*

**Text-to-CAD**
*A cylindrical structure with a hollow center, structured with four evenly spaced rectangular notches and cylindrical pillars in a square pattern around the center.*

**Text-to-CAD**
*A multi-tiered open frame with vertical support rods. The central structure is a rectangular frame with a hollow center. Beneath it are four identical rectangular prisms, serving as supports at the four corners, enhancing the frame's rigidity. Four cylindrical pillars are placed beneath the four rectangular prisms, increasing the elevation and structural integrity.*

Figure 15. Supplementary results of the Text-to-CAD task generated by CAD-Llama-INS based on text prompts.

**Text-to-CAD**

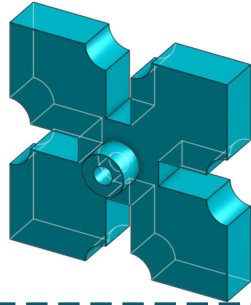*A rectangular block with a large central rectangular cutout and two small rectangular notches at the bottom edge, enhanced by a square plate with twelve equally sized circular holes arranged in a circular pattern fixed onto the left side of the block.*

**Text-to-CAD**

*An intricate structure combining a flat, symmetrical cross-shaped plate with rounded arms and a toroidal ring with a hollow circular center. The cross-shaped plate houses the toroidal ring in its center, with the circular holes near the ends of the arms providing attachment points.*

**Text-to-CAD**

*A circular base with a central cylindrical extrusion, intersected by a circular ring and four perpendicular bars, each bar having rectangular blocks attached to their ends, placed on top of the cylindrical extrusion.*

**Text-to-CAD**

*A circular disc with a central hole and four evenly spaced rectangular cutouts near the perimeter on top of the cylindrical protrusion, aligning the central hole with the cylinder.*

**Text-to-CAD**

*Two concentric circular rings connected by evenly spaced rectangular supports. The outer ring is larger and serves as the primary boundary, while the inner ring is smaller and located centrally.*

**Text-to-CAD**

*Two cylindrical discs, each featuring a central hexagonal hole. The discs are connected by a tapered rectangular plate, forming an elongated shape. One disc is larger in diameter and positioned at one end of the plate, while the smaller disc is located at the opposite end.*

**Text-to-CAD**

*Two hollow hexagonal frames are horizontally staggered and stacked together. The frames are connected along a shared edge, creating an overlapping configuration. Each frame has a uniform rectangular cross-section, and their hollow interiors remain open.*

**Text-to-CAD**

*A central hollow rounded square prism with four flat faces. The hollow center is divided into two equal parts by a thin rectangular plate. Four rectangular plates extend outward from the faces of the prism in an X-like configuration*

Figure 16. Supplementary results of the Text-to-CAD task generated by CAD-Llama-INS based on text prompts.

**Text-to-CAD**
*A regular octagonal plate featuring a large central circular hole, eight smaller circular holes arranged symmetrically around the center, and a raised octagonal wall along the outer edge.*

**Text-to-CAD**
*A circular disc with two central holes aligned vertically and evenly spaced square notches around the perimeter.*

**Text-to-CAD**
*A central cylindrical core with six radially arranged rectangular plates, each extending outward symmetrically from the core. The plates are evenly spaced around the core and feature a consistent thickness, creating a star-like arrangement when viewed from above.*

**Text-to-CAD**
*A cross-shaped profile with integrated semi-cylindrical arch bases. The cross-shaped profile features a central circular hole with four extending rectangular cutouts, each linked by triangular channels. On each end of the cross's extensions, flat rectangular bases merge seamlessly with a semi-cylindrical arch, providing structural support and aesthetic contrast.*

**Text-to-CAD**
*A cylindrical body with a central hexagonal hole, connected to a rectangular prism at top with a concave cylindrical cut, and a T-shaped structure composed of three rectangular prisms attached to the side opposite the cylindrical base to act as a sturdy support for the handle-like structure.*

**Text-to-CAD**
*A wrench-like assembly with a circular base featuring a hexagonal cutout, a cylindrical disc with a pentagonal hole, and a rectangular bar with concave ends and a central diamond-shaped cutout with rounded corners. The assembly starts with the circular base with a hexagonal cutout at one end. Next, the cylindrical disc with a pentagonal hole through its center is positioned at another end. Finally, the rectangular bar with concave ends and a central diamond-shaped cutout with rounded corners connects the circular base and the cylindrical disc.*

**Text-to-CAD**
*A rectangular structure with four circular holes on both sides and a semicircular cutout in the middle. The structure is formed by first placing two identical rectangular prisms vertically and longitudinally. Then, two rectangular blocks with two circular holes each are placed horizontally on top of the prisms. A rectangular prism with equal rectangular faces is placed in the center. Finally, a rectangular prism with a semicircular profile is extruded backward to create the semicircular cutout in the middle.*

**Text-to-CAD**
*Integrated onto the structure is a square frame with diagonal cross braces, positioned vertically at the opposite end of the initial rectangular cutouts, enhancing structural rigidity. Furthermore, a rectangular plate with four equally spaced circular holes is mounted on the top surface of the extruded horizontal blocks.*

Figure 17. Supplementary results of the Text-to-CAD task generated by CAD-Llama-INS based on text prompts.

**Text-to-CAD**
*A U-shaped support structure with a cross-shaped top element, providing stability with aesthetic symmetry. Begin with the U-shaped frame and position the cross-shaped block at the top center of the structure, aligning its semi-circular bottom end with the notches in the elongated arms of the frame.*

**Text-to-CAD**
*A circular disc with twelve evenly spaced circular holes near the perimeter, featuring a central cylindrical protrusion with a concentric circular cutout and a small central hole.*

**Text-to-CAD**
*The rectangular base is extruded backward. A cubical prism is cut out from the front surface. Two triangular prisms are cut out from the sides at an angle.*

**Text-to-CAD**
*Four rectangular prisms are organized in a square pattern, forming the corners of the structure. A hollow rectangular beam connects the adjacent prisms, creating a grid-like framework with open square compartments in the center. The arrangement is symmetric, with the beams evenly spaced to form a consistent 3x3 grid pattern*

**Text-to-CAD**
*A cylindrical structure with concentric circular grooves recessed into its top face. The grooves are evenly spaced, creating a series of nested rings that decrease in diameter toward the center.*

**Text-to-CAD**
*A rectangular shelf unit with four evenly spaced shelves is attached perpendicularly to a vertical panel on one side. The rectangular shelf unit consists of two horizontal panels. This is attached perpendicularly to four rectangular prisms of equal size, arranged in a linear sequence with uniform spacing between them, extruded backward and arranged vertically.*

**Text-to-CAD**
*A circular plate with a hexagonal hole at its center, which also contains a smaller circular hole within it. Around the perimeter of the plate, there are multiple evenly spaced cylindrical holes arranged in a circular pattern.*

**Text-to-CAD**
*A U-shaped frame with vertical supports and horizontal extensions. Attach the U-shaped structure horizontally, then add the two elongated rectangular prisms vertically as supports.*

Figure 18. Supplementary results of the Text-to-CAD task generated by CAD-Llama-INS based on text prompts.

**Text-to-CAD**: *A cylindrical tube*

**Add**: *please add a circular cutout on the top surface of the cylinder.*

**Add**: *please add a U -shape rectangular enclosure surrounding the cylindrical object.*

**Text-to-CAD**: *A cylindrical disc with a large central hole and four smaller equally spaced circular holes around it.*

**Add**: *please create two circular cutouts on the disk and add small cylindrical features in the center of each cutout.*

**Delete**: *please remove the central large circular hole.*

**Text-to-CAD**: *An octagonal prism with equal-length sides and a uniform height.*

**Add**: *please add a cross-shaped protrusion with a central circular hole on top of the octagonal base.*

**Add**: *please add a tube feature above the circular hole at the center of the cross-shaped protrusion.*

**Text-to-CAD**: *A circular ring with six evenly spaced circular holes around its perimeter.*

**Add**: *please add a hollow cylindrical protrusion in the center of the ring and embed it into the large hole of the base.*

**Add**: *please insert cylindrical rods that fit into the circular holes along the perimeter of the ring.*

**Text-to-CAD**: *A simple thin tube.*

**Add**: *please add an outer hollow cylindrical layer surrounding the inner thin tube.*

**Add**: *please add a hexagonal outer layer around the middle section of the existing structure.*

Figure 19. Supplementary working examples of Text-to-CAD, Delete, Add tasks using CAD-Llama-INS.

**Text-to-CAD**: *Four elongated rectangular bars arranged in a cross pattern with equal spacing between them.*

**Add**: *please add a flat plate to the center of the framework, filling the middle section.*

**Add**: *please create four circular holes on the flat plate, located at the corners of the central section.*

**Text-to-CAD**: *A hexagonal prism with uniform thickness.*

**Add**: *please add four vertical cylindrical extrusions around the hexagonal prism, arranged in a rectangular pattern.*

**Delete**: *please add two flat panels connecting adjacent vertical cylindrical extrusions, and add one parallel panel in the center between the previous two panels.*

**Text-to-CAD**: *A thin rectangular block*

**Add**: *please add a rectangular frame around the edges of the rectangular block.*

**Add**: *please add eight identical hollow cylindrical rings arranged in two rows of four on the surface of the framed plate.*

**Text-to-CAD**: *A square frame with a central square cutout.*

**Add**: *please add four vertical supports at the corners of the rectangular frame to form a table-like structure.*

**Add**: *please create four circular holes evenly distributed on the top plane and four triangular cutouts at each corner of the frame.*

**Text-to-CAD**: *A rectangular block.*

**Add**: *please add symmetrical grooves on the left and right sides of the rectangular plate.*

**Add**: *please add two hexagonal protrusions on the left and right grooves each containing a circular hole.*

Figure 20. Supplementary working examples of Text-to-CAD, Delete, Add tasks using CAD-Llama-INS.

**Text-to-CAD:** *A flat rectangular plate with a cylindrical extrusion in the center. There are six evenly distributed small holes near the edges of the plate.*

**Delete:** *please remove all holes of the rectangular plate.*

**Add:** *please create four evenly distributed square cutouts and one central circular cutout on the cylinder.*

**Text-to-CAD:** *A rectangular prism with two curved concave cuts along its top surface*

**Add:** *please create a through circular cutouts on the rectangular block.*

**Add:** *please add two rectangular blocks with concave curved edges on one face to the top and bottom of the rectangular structure.*

**Text-to-CAD:** *A thick circular ring with a hollow center.*

**Add:** *please add a circular flat surface at the bottom of the ring to create a closed base structure.*

**Add:** *please embed four evenly distributed rectangular protrusions along the edge of the ring, perpendicular to its side surface.*

**Text-to-CAD:** *A thin octagonal plate.*

**Add:** *Please add a central circular cutout and six evenly distributed cylindrical cutouts near the edges of the octagonal plate.*

**Add:** *Please add a cylindrical extrusion beneath the octagonal plate, with a hexagonal cutout running through the center of the cylinder.*

**Text-to-CAD:** *A symmetric, gear-like structure with six evenly spaced protrusions extending outward from a central core. Each protrusion has a concave curve on its outer face, creating a rounded, recessed shape between adjacent protrusions.*

**Add:** *please add a circular hole through the center of the gear-like structure.*

**Add:** *please add a cylindrical object beneath the central hole, featuring a hexagonal cutout through its center.*

Figure 21. Supplementary working examples of Text-to-CAD, Delete, Add tasks using CAD-Llama-INS.

**Text-to-CAD**: *A flat cross-shaped structure with rectangular arms extending outward in four perpendicular directions. Each arm has a circular hole near its outer edge, while the center of the cross also features a circular hole.*

**Delete**: *please remove all circular holes.*

**Text-to-CAD**: *A series of five concentric squares, each progressively smaller and stacked vertically.*

**Add**: *please add two more square layers in the existing stack.*

**Text-to-CAD**: *A long horizontal rectangular base with rounded edges on both ends. Perpendicular to the base, two thin vertical rods extend upward*

**Add**: *Please add three horizontal bars connecting the vertical rods to form a grid structure, ensuring equal spacing between the horizontal bars.*

**Text-to-CAD**: *A flat rectangular plate with a cylindrical extrusion in the center. There are six evenly distributed small holes near the edges of the plate.*

**Add**: *Please create five equally spaced rectangular through-holes along length of the rectangular.*

**Text-to-CAD**: *A U-shaped block with two parallel arms and a connecting base, featuring two horizontal cylindrical holes and two vertical cylindrical cutouts on the bottom.*

**Add**: *Please add a horizontal cylinder to fit into two horizontal cylindrical holes.*

**Text-to-CAD**: *A cylindrical structure with a wide, flat base and a smaller concentric cylinder on top, ensuring structural stability.*

**Add**: *Please add four identical hollow cylindrical rings evenly distributed around the central cylindrical structure.*

Figure 22. Supplementary working examples of Text-to-CAD, Delete, Add tasks using CAD-Llama-INS.

**Text-to-CAD**: *A cylindrical container is composed of a cylindrical base with a hollow center. The hollow cylindrical ring is placed on top of the base.*

**Add**: *Please add four evenly spaced cylindrical protrusions to the outer wall of the circular shell.*

**Text-to-CAD**: *A right-extruded rectangular prism.*

**Add**: *Please add four cylindrical rings of varying diameters and thicknesses above the right-extruded rectangular prism to enhance the existing structure.*

**Text-to-CAD**: *A cylindrical solid with a circular cross-section.*

**Add**: *Please add a series of evenly spaced circular holes along the surface of the cylindrical tube, arranged in a single line.*

**Text-to-CAD**: *A simple rectangular block.*

**Add**: *Please add a grid of sixteen circular holes on top face of the bock.*

**Text-to-CAD**: *The primary structure is a rectangular base with a series of rectangular notches along the top edge and a single rectangular cutout along the bottom edge. To the bottom edge is affixed a rectangular bar with two large circular holes near each end and two smaller circular holes positioned closer to the center.*

**Delete**: *Please remove all circular holes from the rectangular bar.*

**Text-to-CAD**: *Two flat rectangular plates connected perpendicularly along one edge to form an "L" shape. Each plate features symmetrical rectangular cutouts along its free edges, creating slots or notches.*

**Add**: *Please add multiple evenly spaced circular holes on both vertical plates, ensuring alignment and uniform spacing for structural consistency.*
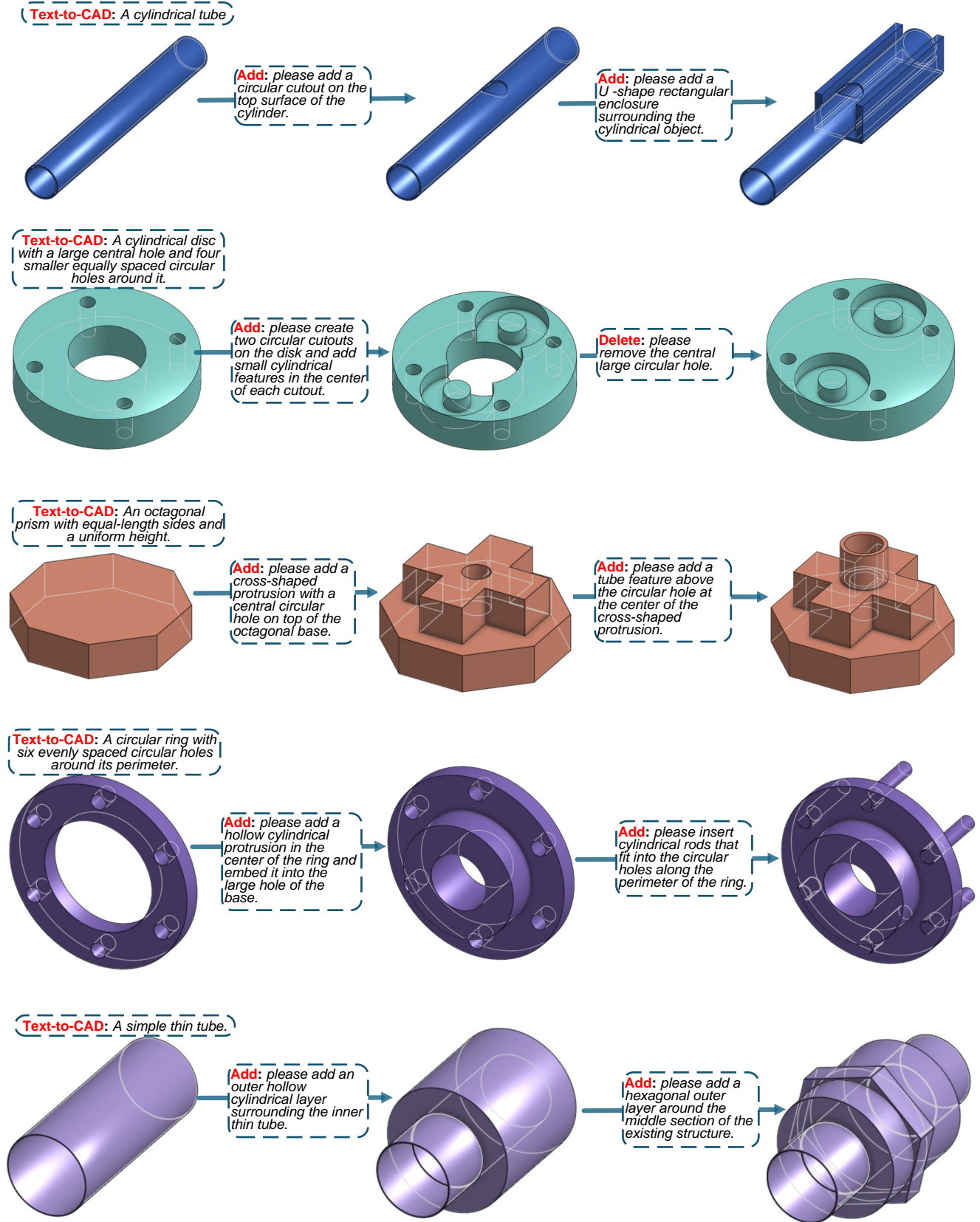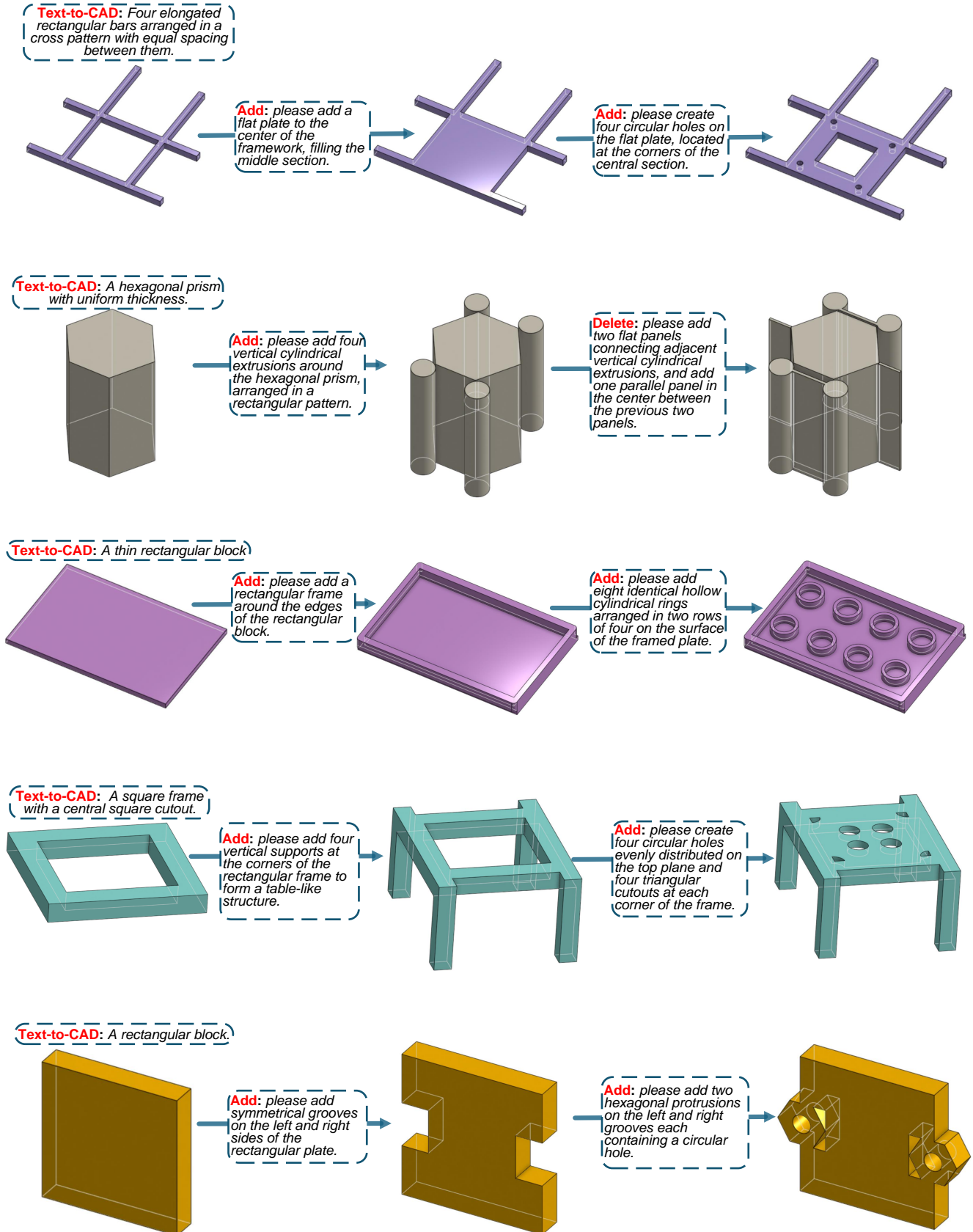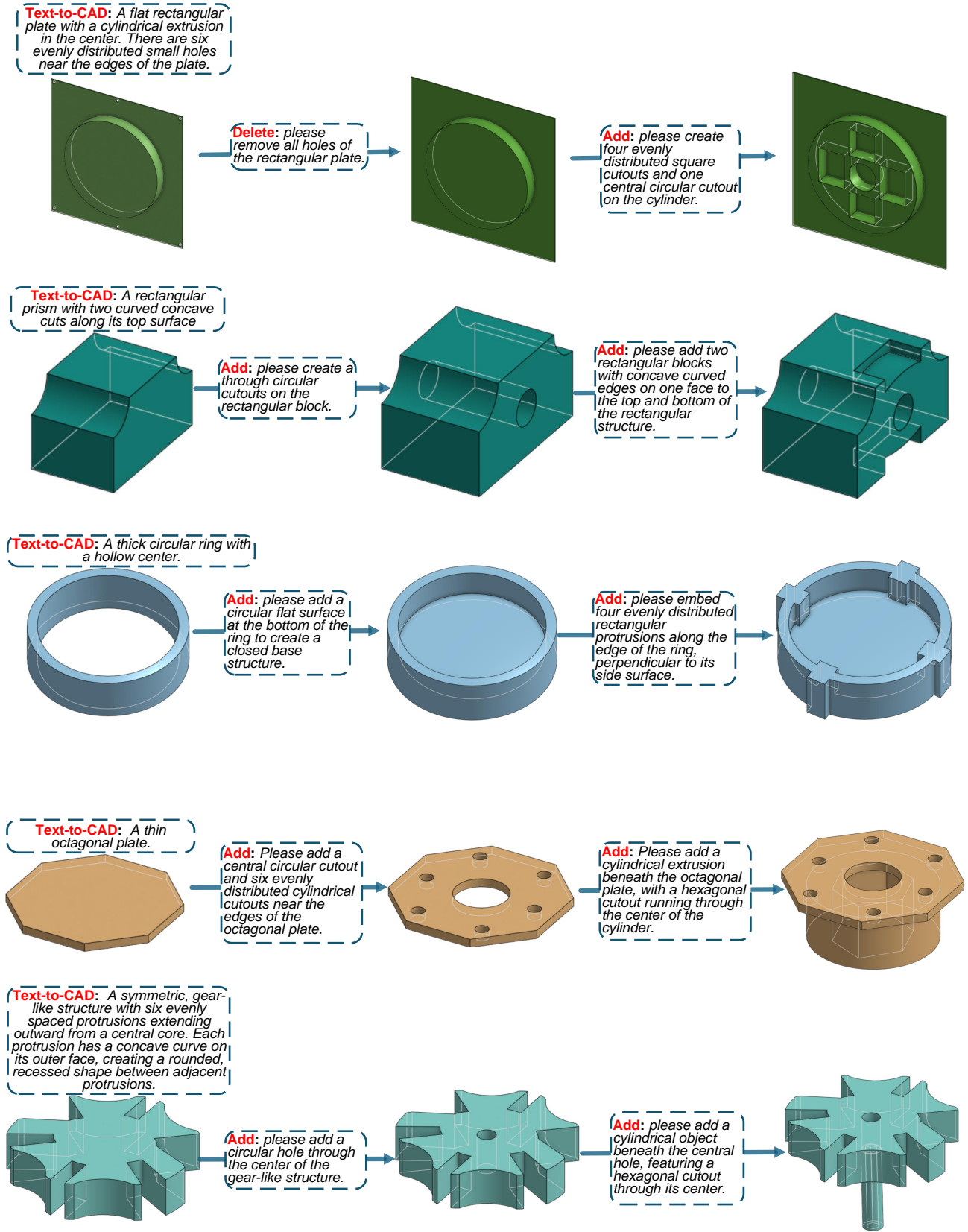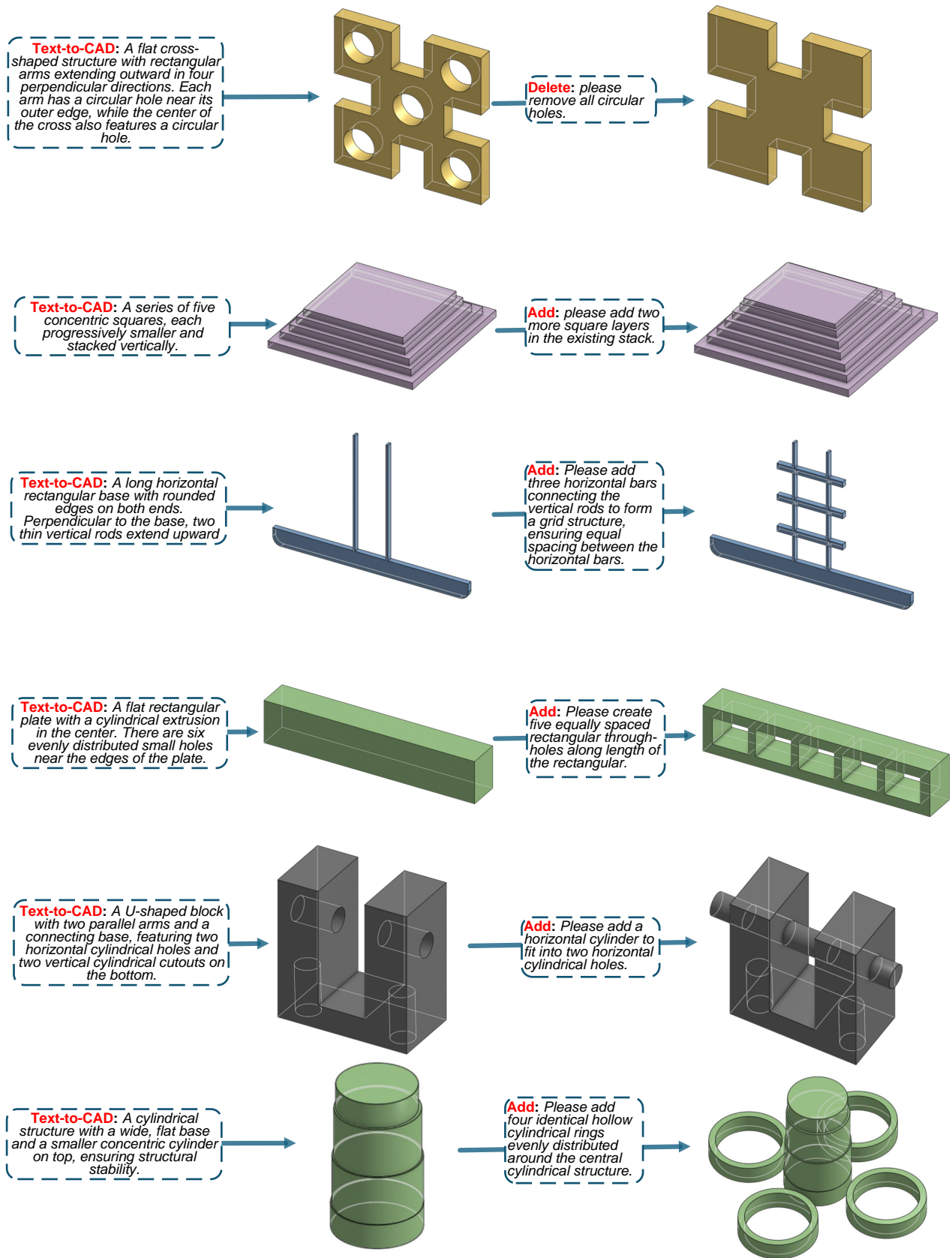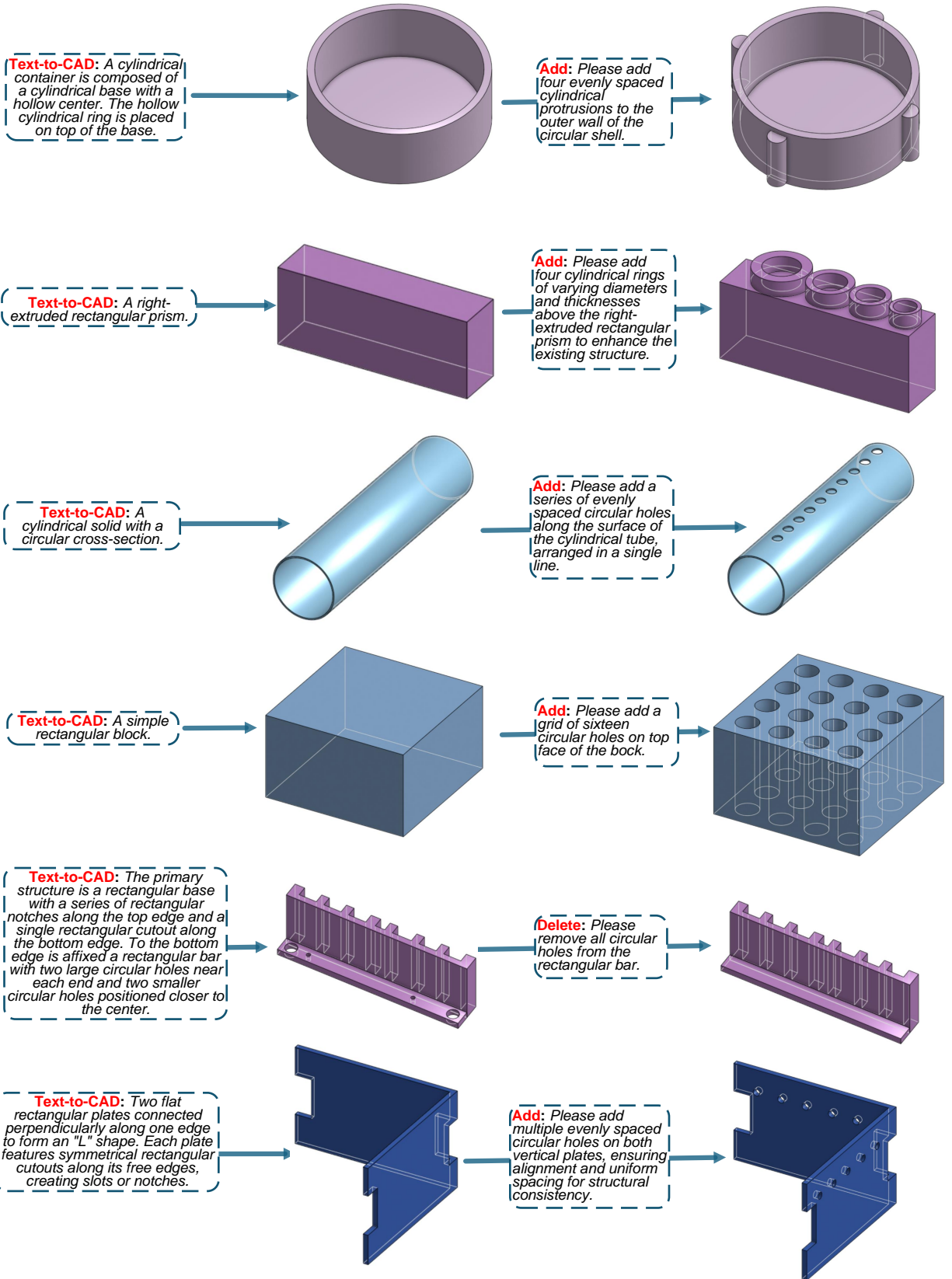
Figure 23. Supplementary working examples of Text-to-CAD, Delete, Add tasks using CAD-Llama-INS.

Figure 24. Supplementary results of unconditional generation produced by CAD-Llama.

**Input Text**

*A central hollow cylindrical body with two rectangular blocks extending symmetrically on opposite sides. Each rectangular block features a circular hole.*

*A central cylindrical structure with layered circular features on top and supported by four vertical cylindrical legs arranged symmetrically around the base.*

*The shape is a flat circular disk with a raised cylindrical feature in the center and **four holes** around it.*

*A hollow cube with **all six faces open**, defined by thin edges forming the cube's frame.*

**Generated CAD Models**

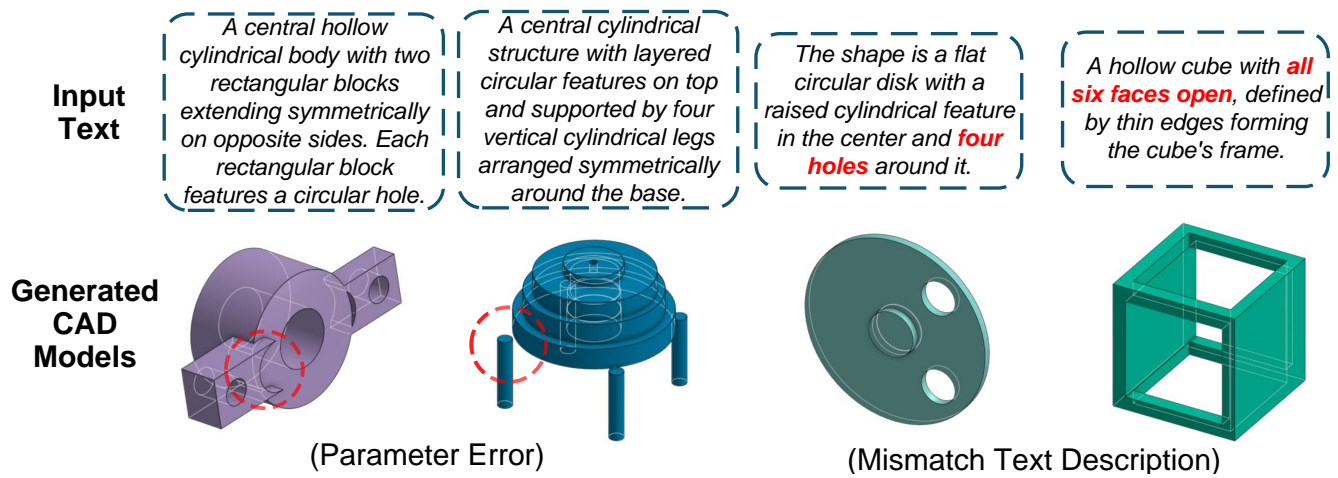(Parameter Error)          (Mismatch Text Description)

Figure 25. Failure cases for CAD-Llama-INS. We illustrate two types of errors: inaccuracies in parameter settings and misalignment with the text prompts.