

Towards Proactive Text-to-CAD Generation

Anonymous ACL submission

Abstract

Computer-Aided Design (CAD) systems are indispensable in mechanical engineering and product development processes. Nowadays, text-to-CAD methods can significantly reduce the learning cost of complex CAD systems and has attracted increasing attention. However, such methods fail to achieve alignment among user expectations, textual descriptions and CAD models. To address this limitation, we propose a new paradigm, “Proactive Text-to-CAD Generation”, which first employs large language models to proactively elicit and formulate text enriched with comprehensive CAD design details, then generates CAD models from these refined descriptions. To support this paradigm, we construct the first actively interactive text-to-CAD dataset, Proactive-Text2CAD, which contains 4,590 high-quality dialogues. Moreover, building upon this dataset, we propose a novel agentic framework for this task, named “Proactive Agent”, which is driven by a hierarchical finite state machine accompanying with three carefully designed modules. Extensive evaluation and comprehensive analysis on the Proactive-Text2CAD dataset demonstrate the effectiveness of both our proposed paradigm and agentic framework, with our method achieving significant improvements in both textual detail refinement and final CAD model generation quality.

1 Introduction

Computer-Aided Design (CAD) systems serve as fundamental tools in mechanical engineering and product development, revolutionizing prototyping methodologies (Robertson and Allen, 1993). In traditional CAD software (e.g., Autodesk, FreeCAD, SolidWorks and Onshape), users create and modify geometric entities and constraints through graphical user interfaces (GUIs). However, this requires considerable expertise and proficiency, which can

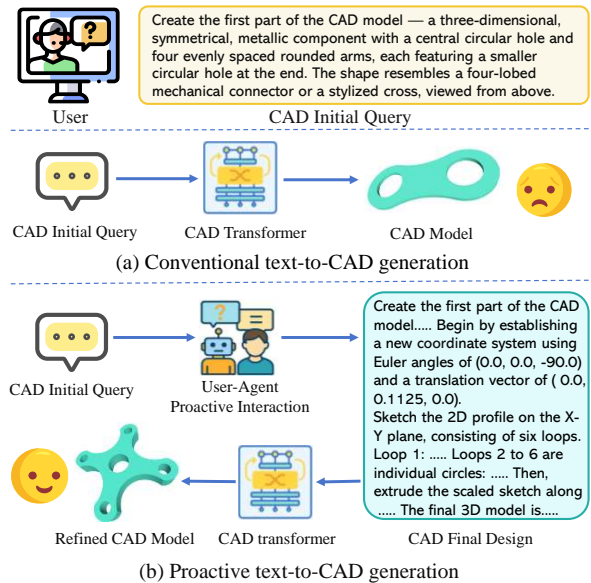


Figure 1: The conventional Text-to-CAD generation way versus our proposed proactive text-to-CAD generation way.

be challenging for non-specialists to master (Deng et al., 2024; Zhou and Camba, 2025).

To address this challenge, integrating natural language input with CAD systems through highly capable large language models (LLMs), which excel at interpreting and synthesizing structured data such as command sequences, streamlines parametric CAD generation. Recently, numerous studies have been devoted to this topic, such as Nelson et al. (2023); Khan et al. (2024b); Li et al. (2024b); Kapsalis (2024); Li et al. (2025); Zhang et al. (2025); Xie and Ju (2025). Such a text-to-CAD generation paradigm can definitely minimize the need for users to directly interact with complex GUIs, significantly lowering the barrier to CAD modeling (Zhou and Camba, 2025).

However, such a paradigm is not without its flaws, particularly in achieving alignment among user expectations, textual descriptions and CAD models. In detail, current text-to-CAD meth-

ods (Nelson et al., 2023; Khan et al., 2024b; Li et al., 2024b; Kapsalis, 2024; Badagabettu et al., 2024; Li et al., 2025; Zhang et al., 2025) can effectively handle the alignment between textual descriptions and CAD models, but they struggle to maintain satisfactory alignment when dealing with vague or abstract textual inputs. Some compelling experimental results (Khan et al., 2024b) also reveal that as textual descriptions become more abstract, the accuracy and precision of generated CAD models deteriorate significantly. This challenge is particularly pronounced because novices and non-specialists tend to provide high-level descriptions of model appearance rather than detailed parametric specifications. Moreover, even experienced CAD engineers find it difficult to produce text descriptions that are both sufficiently detailed and compliant with CAD generation principles without the aid of a visual interface. In short, merely relying on user-input text descriptions cannot bridge the gap between user intent and the final CAD model.

To fill this gap, we propose proactive text-to-CAD generation (shown in Figure 1). Such a way actively engages users in iterative questioning to seek missing CAD design details in the initial user provided textual description. By refining incomplete parametric specifications and enhancing the text’s descriptive quality through this dialogue-driven process, the framework ensures higher-quality input before final CAD generation. To achieve proactive text-to-CAD generation, we make the following efforts in this paper:

First, we build the first actively interactive text-to-CAD dataset, Proactive-Text2CAD. Our goal is to integrate the incomplete user-provided CAD text descriptions with a proactive information-seeking dialogue, enabling the agent to ask targeted questions actively when encountering missing or unclear CAD parameters. To achieve that, based on Text2CAD’s expert-level text annotation dataset and CAD dataset (Khan et al., 2024b), we further involve two key phases: (1) User initial query generation and (2) Proactive dialogues generation. Through such a dataset construction process, we obtain 4,590 high-quality dialogues, each containing an incomplete metadata entry and an active dialogue, totaling around 28,110 question-answer pairs. Building upon our dataset, we further establish a family of strong and representative baselines, which can be roughly categorized into three type (Deng et al., 2025): (1) Standard Prompting,

like Proactive Prompting (Deng et al., 2023); (2) CoT Prompting, like ProCoT (Deng et al., 2023) and PS+ Prompting (Wang et al., 2023), and (3) Multi-agent Prompting, like MACRS (Fang et al., 2024).

Second, we propose a novel agentic framework for this task, named “Proactive Agent”. In detail, through experiments, we find that current proactive methods, like ProCoT, PS+ Prompting and MACRS, underperform in refining textual descriptions of CAD models, especially exhibiting significant limitations in comprehensively identifying missing operational and parametric details in the CAD model, due to the complex CAD topology. For example, when the subtypes face and loop are missing in the sketch type of the CAD model, since loop is a subtype of face and a CAD model can contain multiple faces, current baselines fail to proactively ask the user for specific information about the face first. To address this issue, we involve three modules in the proposed agentic framework: strategic, tactical, and operational modules. The strategic module analyzes the current CAD design text at a macro perspective, captures missing details, and generates an agent workflow. The tactical and operational modules each contain independent modules to address specific functions, decoupling the complex workflow analysis and execution process. To interconnect these three modules and enable autonomous operation, we design a hierarchical finite state machine (HFSM) that drives the reasoning logic from workflow generation to execution through state transitions.

In summary, the contributions are as follows:

- We propose “Proactive Text-to-CAD Generation” paradigm that engages users in iterative questioning to seek missing CAD design details in the initial user provided textual description, then generates CAD models from the refined textual descriptions.
- We construct Proactive-Text2CAD, the first dataset for studying actively interactive text-to-CAD generation, which contains around 4,500 high-quality dialogues.
- We present “Proactive Agent”, an innovative agentic framework driven by a hierarchical finite state machine accompanying with workflow generation, task dispatching, questioning, and summarization functional modules.

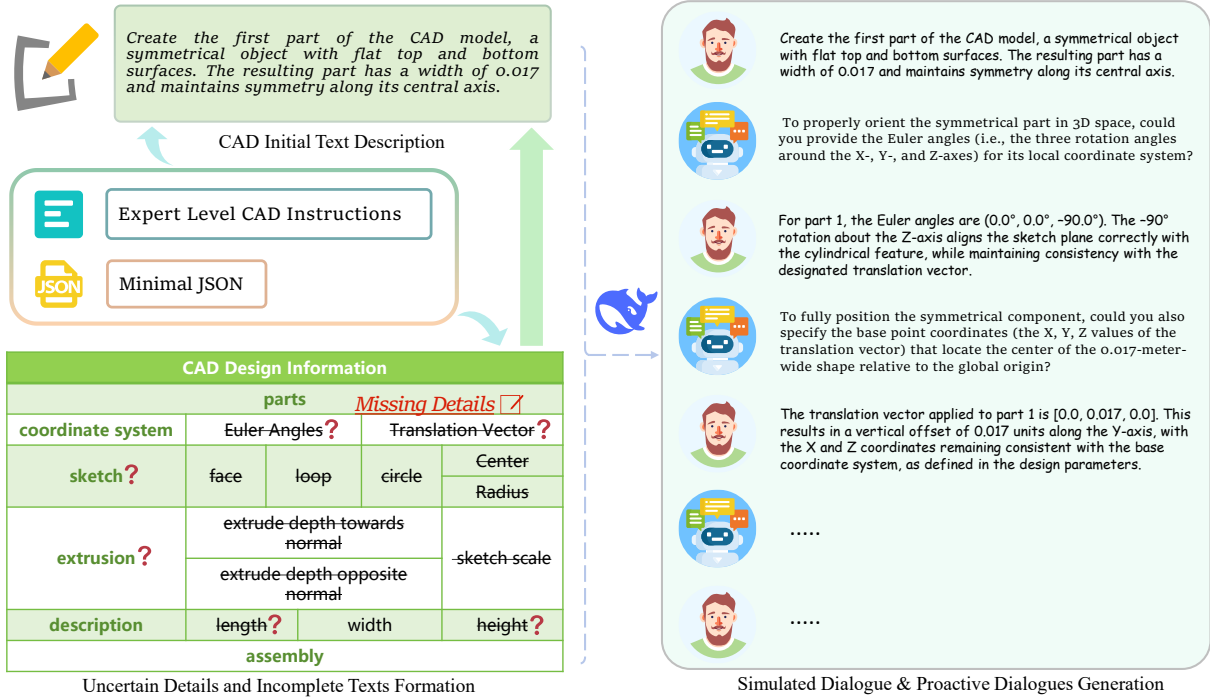


Figure 2: The whole process of Proactive-Text2CAD construction.

2 Proactive-Text2CAD Construction

2.1 Dataset Construction

Our objective is to combine incomplete user-provided initial CAD text descriptions with proactive information-seeking dialogue, enabling the agent to actively inquire when facing CAD parameter details that are missing or incomplete. As illustrated in Figure 2, the dataset construction process consists of two key phases: (1) User Initial Query Generation; (2) Proactive Dialogues Generation.

2.1.1 User Initial Query Generation

At this stage, our task is to generate user initial query with uncertain detail missing, which are divided into four specific steps:

First, through collecting the multi-level textual description dataset of CAD constructed by Text2CAD and the minimal metadata dataset (minimal json) (Wu et al., 2021a; Khan et al., 2024b), we generate an initial dataset containing precise geometric descriptions, text descriptions with relative measurements, and concise representations of shape attributes and their relational properties in CAD designs.

Next, based on the minimal metadata, we categorize the details of CAD textual descriptions into six main types: component quantity, sketch, extrusion, coordinate system, appearance description, and assembly method. Among these, the sketch,

extrusion, coordinate system, and appearance description types further contain multiple subtypes (e.g., subtypes of sketch: face, loop, line, etc.)¹.

Then, we omit certain CAD parameter type information to create missing types. We subsequently use Deepseek-R1 (DeepSeek-AI et al., 2025) to rewrite the text descriptions based on the missing types, generating user initial query that lack the missing category information, i.e., incomplete texts².

Finally, since there exists sequential relationships among CAD operation types and dependency relationships among parameter types, we employ topological sorting to arrange the questioning order of missing types according to the dependency graph³.

2.1.2 Proactive Dialogues Generation

In this stage, our task is to generate the whole dialogue, which are divided into four specific steps:

First, we utilize DeepSeek-R1 (DeepSeek-AI et al., 2025) to generate one question for each missing type to inquire about the details of the missing information. These questions serve as the questions posed by the system to the user.

Second, for each generated question, we retrieve specific information about the missing types from

¹Detailed type information is shown in Appendix A.2.1.

²The rewriting prompts are provided in Appendix A.3.1

³Specific information can be found in Appendix A.2.1

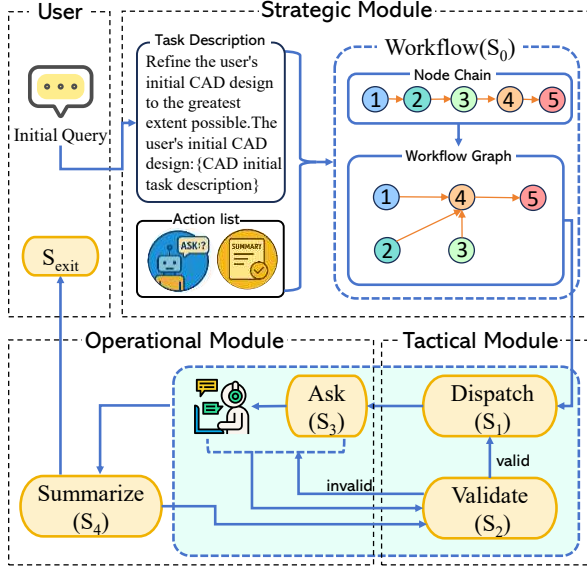


Figure 3: Our proposed agentic framework: “Proactive Agent”.

tically identify the missing operational and parametric details in the current CAD model. To address the aforementioned issue, we propose a novel multi-agent framework named “Proactive Agent” (shown in Figure 3), which consists of strategic, tactical, and operational modules. To interconnect these three modules and enable autonomous operation, we design a hierarchical finite state machine (HFSM) that drives the reasoning logic from workflow generation to execution through state transitions.

3.3.1 Strategic Module

The strategic module (S_0) analyzes the current CAD design text at a macro level, captures missing details, and generates an agent workflow. There are two key phase in this module.

Specifically, first, the strategic module generate a topological sequence \mathcal{C} of the workflow graph via a give LLM $_{\theta}$ ⁶, which can be denoted as

$$\mathcal{C}(V) \leftarrow \text{LLM}_{\theta}(d, t, \mathcal{A}). \quad (1)$$

In the above equation, d represents the task description of workflow generation, t represents the initial textual description of the user’s CAD modeling task, \mathcal{A} denotes the action set {“ask”, “summarize”}, where “ask” indicates proactively querying the user for CAD information, and “summarize” indicates extracting the CAD design scheme from the dialogue history, and the nodes

$V = \{v_1, v_2, \dots, v_n\}$ represent tasks to be finished, such as “ask for the coordinate system details” and “ask for the details of the sketch to be created”.

Then, based on the dependencies among the details of each CAD operation and parameter, the topological sequence \mathcal{C} then forms a directed acyclic graph \mathcal{G} according to topological sorting (Qiao et al., 2024), which can be denoted as:

$$\mathcal{G}(V, E) \leftarrow \mathcal{C}(V) \quad (2)$$

where $E = \{(v_i, v_j)\}$, where $1 \leq i \neq j \leq n$, represent the execution relationships between nodes (v_j must be executed after v_i).

3.3.2 Tactical Module

To achieve efficient workflow execution, the tactical module (S_{tact}) consists of two sub-modules: dispatch (S_1) and validation (S_2).

The dispatch sub-module is responsible for receiving and parsing the workflow graph \mathcal{G} . Via Kahn’s Algorithm (Kahn, 1962), the parsing function τ can dispatch nodes V to either the asking or summary sub-modules (introduced in §3.3.3, and denoted as S_3 and S_4 , respectively), which can be represented as:

$$\tau(V; \mathcal{G}(V, E)) \rightarrow \{S_3, S_4\}, \quad (3)$$

The validation sub-module evaluates the question q (produced by the asking sub-module) and summary p (generated by the summary sub-module) against the task instruction v_i via the given LLM $_{\theta}$ ⁷. Based on verification results, it triggers regeneration of q or p through corresponding sub-modules.

3.3.3 Operational Module

To complete the underlying task details at the micro-level and improve the execution efficiency of workflow G , the operational module (S_{oper}) consists of two independent sub-modules: asking (S_3) and summary (S_4).

The asking sub-module is responsible for proactively asking questions to the user. After receiving a task instruction v_i from the dispatch sub-module, the asking sub-module can leverage the give LLM $_{\theta}$ to generate question q ⁸, which can be represented as

$$q \leftarrow \text{LLM}_{\theta}(v_i) \quad (4)$$

⁷The prompt is shown in the Appendix A.3.2.

⁸The prompt is shown in the Appendix A.3.2

⁶The prompt is shown in the Appendix A.3.2.

Method	Clarif.Acc. \uparrow		ROUGE-L \uparrow		Rule-based Score \uparrow		BertScore \uparrow	
	Micro	Macro	Micro	Macro	Micro	Macro	Micro	Macro
Standard	0.10	0.09	0.05	0.04	0.07	0.06	0.05	0.04
CoT	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Proactive	4.12	3.72	2.33	2.02	7.53	6.84	2.74	2.38
ProCoT	67.85	71.22	19.55	18.75	35.76	31.10	20.50	19.74
PS+	44.67	37.27	14.75	13.84	12.27	9.93	14.90	13.50
MACRS	99.54	99.74	17.24	15.72	34.68	28.90	21.33	19.93
Ours	99.89	99.93	21.04	19.11	45.53	38.14	24.62	22.21

Table 1: Experimental results at the turn-level.

Method	Completion	Effectiveness
	Rate \uparrow	Rate \uparrow
Standard	0.04	0.07
CoT	0.00	0.00
Proactive	3.22	3.17
ProCoT	12.45	5.95
PS+	6.78	8.56
MACRS	8.27	2.68
Ours	31.67	17.34

Table 2: Experimental results at the dialogue-level.

The summary sub-module is responsible for generating the final CAD design solution based on the dialogue history \mathcal{H} . After the summary module receives the task instruction, it generates a CAD design solution in combination with the task instruction v_i , which can be represented as:

$$p \leftarrow \text{LLM}_\theta(v_i, \mathcal{H}) \quad (5)$$

where the dialogue history \mathcal{H} includes the user’s initial CAD text description, the questions posed by the asking sub-module, and the user’s responses.

3.3.4 Hierarchical Finite State Machine

To coordinate the integration and autonomous operation of the aforementioned modules, we leverage the Hierarchical Finite State Machine (HFSM) (Alur and Yannakakis, 2001). Such a machine can employ state-based reasoning to systematically govern the workflow from generation to execution. The HFSM is formally modeled as a quintuple $\{S, O, \mu, S_0, \{S_{\text{exit}}\}\}$:

- $S = S_0 \cup S_{\text{tact}} \cup S_{\text{oper}} \cup S_{\text{exit}}$ is a hierarchical state set, where S_0 denotes the strategic module and is the initial state in the machine, $S_{\text{tact}} = \{S_1, S_2\}$ is the tactical module, $S_{\text{oper}} = \{S_3, S_4\}$ represents the operational module, and while S_{exit} denotes the terminal state.
- O represents the set of all possible outputs from the aforementioned modules S .
- $\mu : S \times O \rightarrow S$ is a transition function that determines the next state in the reasoning process based on the current state and the execution result of the corresponding module⁹.

Figure 3 illustrates the whole working mechanism of the HFSM for the proposed proactive agent.

⁹We put the detail of the transition function in the Appendix A.4 due to the page limit.

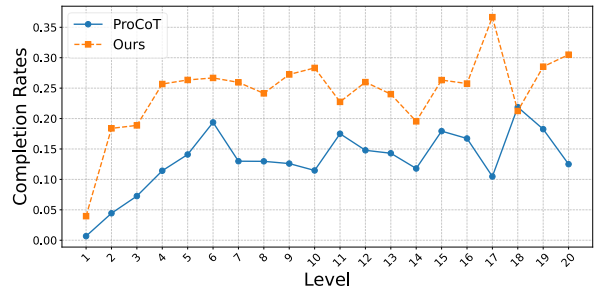


Figure 4: Completion rates for dialogues with varying numbers of missing details.

4 Experiment

4.1 Performance Comparisons

In this paper, we employ GPT-4o Mini as the primary model for our main experiments.

4.1.1 Results at the Turn-Level

Table 1 shows that our proactive agent outperforms all baseline methods across every metric, confirming its effectiveness in identifying missing CAD details and generating relevant proactive queries. CoT yields deficient results, failing to analyze omissions or initiate queries, which underscores its limitations for this task. In comparison, MACRS and ProCoT maintain solid performance in both detail analysis and query generation, with ProCoT showing particularly strong questioning capability in this evaluation.

4.1.2 Results at the Dialogue-Level

Table 2 presents the dialog-level evaluation results. Our proactive agent achieves the highest Completion Rate and Effective Rate, outperforming all baselines in comprehensively identifying missing CAD description details. Consistent with turn-level observations, CoT fails to detect omissions or initiate meaningful inquiries. While MACRS maintains

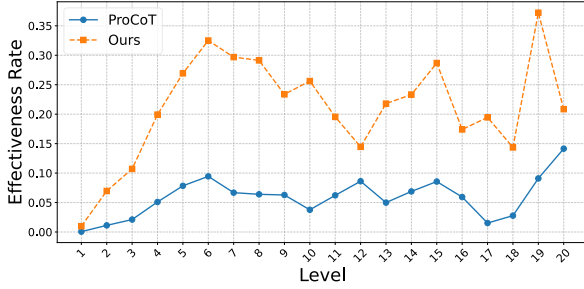


Figure 5: Effectiveness rates for dialogues with varying numbers of missing details.

Method	Median CD ↓	Mean CD ↓	Avg. Rank ↓
Initial Query	45.52	142.61	\
Standard	47.66	141.56	3.56
CoT	48.48	141.50	3.82
Proactive	36.42	129.54	2.78
ProCoT	28.00	116.13	2.20
PS+	34.73	131.80	2.61
MACRS	27.03	119.69	2.64
Ours	20.60	99.70	1.60

Table 3: Experimental results at the final CAD-level.

moderate completion performance, its effectiveness rate remains notably low.

We further categorize dialogues into 20 levels based on missing information types and operational complexity. As shown in Figure 4 and 5, our method consistently surpasses the strongest baseline (ProCoT) across all classification levels, demonstrating robust performance granularity.

4.1.3 Results at the CAD-Level

Table 3 presents CAD-level results, showing that our method and most baselines improve CD values over initial queries, confirming proactive interaction enhances alignment. Our agent achieves the best CD performance, verifying its effectiveness. The negative correlation between CD results and dialogue-level completion rates indicates that refined text prompts improve geometric alignment. Interestingly, MACRS attains a slightly better median CD than ProCoT, likely owing to its polarized refinement outcomes, yet ProCoT achieves a higher average CD, indicating more consistent quality.

4.1.4 Results of changing different base LLMs

We investigate base LLM impact using GPT-4o mini, Gemini-2.0-flash-lite, and DeepSeek-R1. As ProCoT is the strongest baseline, we compare it with our method on dialogue and CAD generation. Table 4 shows base LLM selection significantly af-

fects performance. DeepSeek-R1 achieves the best results, while our method consistently outperforms ProCoT across all metrics with every LLM.

4.1.5 Results of ablation studies

Table 6 presents the ablation study. Removing the validation module lowers completion rate, as it improves description completeness but introduces redundant questions that slightly hurt geometric precision. However, since LLMs lack proactive reasoning under information gaps (Zhou et al., 2025), this module remains essential. The Strategic Module’s action list standardizes the action space, substantially improving completion rate, effectiveness rate, and geometric precision (as shown by lower CD scores), confirming structured actions enhance proactive parameter acquisition.

4.2 Human Evaluation

We follow (Chen et al., 2025a) to evaluate user experience across four dimensions: intent understanding, interactivity, operational experience, and emotional feedback. Annotators of three expertise levels conduct paired comparisons (Zheng et al., 2023; Si et al., 2024) on 100 instances, generating target CAD models and selecting the preferred system across seven metrics and overall. Final decisions are made by majority vote, with a Fleiss’ Kappa of 0.752 indicating moderate agreement.

As shown in Table 5, our system outperforms both Text2CAD and the strongest baseline ProCoT in all evaluation dimensions. Users favor our system for enhancing novice onboarding, proficient-user productivity, and operational experience. Notably, the advantage over Text2CAD is significantly larger than that over ProCoT, indicating a clear user preference for interactively refining initial designs over providing complete descriptions in a single attempt.

4.3 Case Study

Figure 6 compares geometric fidelity across user-input CAD descriptions, the strongest baseline (ProCoT), our method, and ground-truth models. Our approach produces geometrically richer outputs with closer structural correspondence to the ground truth. In Case 1, it accurately reconstructs the stacked cylindrical-prismatic geometry, a well-defined structure not captured by the initial input or ProCoT. For Case 3’s open-top container with uniform walls and orthogonal corners, our method reconstructs all structural components except the

LLM	Method	Dialogue-level Evaluation		CAD Model Evaluation	
		Completion Rate (%) \uparrow	Effectiveness Rate (%) \uparrow	Median CD \downarrow	Mean CD \downarrow
GPT-4o mini	ProCoT	12.45	5.95	28.00	116.13
	Ours	31.67	17.34	20.60	99.70
Gemini-2.0-flash-lite	ProCoT	14.58	8.64	25.11	114.54
	Ours	24.76	9.54	20.91	106.24
DeepSeek-R1	ProCoT	14.31	10.22	30.55	115.44
	Ours	48.75	29.71	17.60	105.14

Table 4: The experimental results by changing base LLMs.

Framework	Status	Functional		Interactive			Emotional		Overall
		DIC	TaskEff	GE	Learnability	IC	ASA	IES	
Text2CAD vs. Ours	Text2CAD	3%	13%	5%	7%	9%	5%	9%	4%
	Tie	6%	8%	13%	22%	17%	23%	18%	20%
	Ours	91%	79%	82%	71%	74%	72%	73%	76%
ProCoT vs. Ours	ProCoT	24%	28%	27%	25%	23%	11%	19%	21%
	Tie	9%	10%	15%	11%	16%	19%	12%	7%
	Ours	67%	62%	58%	64%	61%	70%	69%	72%

Table 5: The experimental results of human evaluation.

	Dialogue-level Evaluation		CAD Model Evaluation	
	Completion Rate(%) \uparrow	Effectiveness Rate(%) \uparrow	Mean CD \downarrow	Median CD \downarrow
Ours (w/o Validation)	24.09	23.46	17.62	113.72
Ours (w/o Action list)	17.45	15.93	23.58	118.35
Ours	31.67	17.34	20.60	99.70

Table 6: Ablation study results on core modules.

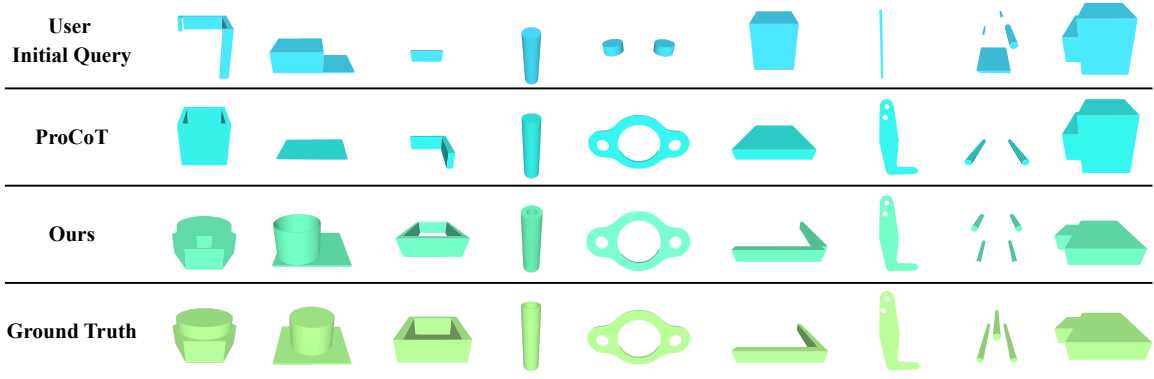


Figure 6: Case studies to illustrate the effectiveness of our method.

base plate, whereas ProCoT only adds one side panel to the initial input. These results demonstrate our framework’s improved precision in converting textual specifications to complex CAD geometries.

5 Conclusion

We propose a proactive text-to-CAD paradigm that identifies missing design details through iterative

user questioning. To achieve this goal, we build the first dataset for studying actively interactive text-to-CAD generation, Proactive-Text2CAD. Besides, we present “Proactive Agent”, an innovative agentic framework. Through extensive experiments, we demonstrate our proposed paradigm and agentic framework can achieve significant improvements in both textual detail refinement and final CAD model generation quality.

536 Limitation

537 Our method is currently a pipeline way with two
538 parts, where Part 1 is designed to proactively query
539 users for missing CAD design information and gener-
540 ate the final CAD design text, and while Part 2
541 is configured to employ a CAD Transformer for
542 converting final CAD design text into CAD mod-
543 els. In future work, we may explore an end-to-end
544 way from the user’s initial text to the CAD model
545 through active interaction and implement joint op-
546 timization of the user’s initial text and the CAD
547 model.

548 Besides, through experiments, we find that the
549 validation sub-module in our proposed agentic
550 framework can increase the proactivity of inquiry,
551 but also introduce the issue of generating invalid
552 questions, which can slightly reduce the efficiency
553 of our agentic framework. In future work, we will
554 explore new method to improve the validation sub-
555 module.

556 Declaration of LLM Usage

557 Large language models were used solely in a lim-
558 ited, assistive capacity during manuscript prepara-
559 tion, primarily for improving language quality
560 and clarity. All scientific content, analyses, in-
561 terpretations, and conclusions were developed by
562 the authors, who reviewed and approved the fi-
563 nal manuscript. The literature review and refer-
564 ence selection were conducted independently by
565 the authors using publicly available and verifiable
566 sources.

567 References

568 1962. Topological sorting of large networks. *Communi-*
569 *cations of the ACM*, 5(11):558–562.

570 Rajeev Alur and Mihalis Yannakakis. 2001. Model
571 checking of hierarchical state machines. *ACM Trans-*
572 *actions on Programming Languages and Systems*
573 *(TOPLAS)*, 23(3):273–303.

574 Autodesk. <https://www.autodesk.com/products/autocad/>. Online: accessed 2025-04-20.
575

576 Akshay Badagabettu, Sai Sravan Yarlagadda, and
577 Amir Barati Farimani. 2024. [Query2cad: Generating cad models using natural language queries](#). *Preprint*,
578 arXiv:2406.00144.
579

580 Jiaqi Chen, Yanzhe Zhang, Yutong Zhang, Yijia Shao,
581 and Diyi Yang. 2025a. Generative interfaces for lan-
582 guage models. *arXiv preprint arXiv:2508.19227*.

Valerie Chen, Alan Zhu, Sebastian Zhao, Hussein 583
Mozannar, David Sontag, and Ameet Talwalkar. 584
2025b. Need help? designing proactive ai assis- 585
tants for programming. In *Proceedings of the 2025* 586
CHI Conference on Human Factors in Computing 587
Systems, pages 1–18. 588

Yi Cheng, Wenge Liu, Jian Wang, Chak Tou Leong, 589
Yi Ouyang, Wenjie Li, Xian Wu, and Yefeng Zheng. 590
2024. Cooper: Coordinating specialized agents to- 591
wards a complex dialogue goal. In *Proceedings of* 592
the AAAI Conference on Artificial Intelligence, vol- 593
ume 38, pages 17853–17861. 594

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, 595
Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, 596
Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, 597
Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhi- 598
hong Shao, Zhuoshu Li, Ziyi Gao, and 181 others. 599
2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). *Preprint*,
600 arXiv:2501.12948. 601
602

Yang Deng, Lizi Liao, Liang Chen, Hongru Wang, 603
Wenqiang Lei, and Tat-Seng Chua. 2023. Prompt- 604
ing and evaluating large language models for proac- 605
tive dialogues: Clarification, target-guided, and non- 606
collaboration. *arXiv preprint arXiv:2305.13626*. 607

Yang Deng, Lizi Liao, Wenqiang Lei, Grace Hui Yang, 608
Wai Lam, and Tat-Seng Chua. 2025. Proactive con- 609
versational ai: A comprehensive survey of advance- 610
ments and opportunities. *ACM Transactions on In-* 611
formation Systems, 43(3):1–45. 612

Yuanzhe Deng, James Chen, and Alison Olechowski. 613
2024. What sets proficient and expert users apart? re- 614
sults of a computer-aided design experiment. *Journal* 615
of Mechanical Design, 146(1). 616

Haoqiang Fan, Hao Su, and Leonidas Guibas. 2016. 617
[A point set generation network for 3d object reconstruction from a single image](#). *Preprint*,
618 arXiv:1612.00603. 619
620

Jiabao Fang, Shen Gao, Pengjie Ren, Xiuying Chen, 621
Suzan Verberne, and Zhaochun Ren. 2024. A multi- 622
agent conversational recommender system. *arXiv* 623
preprint arXiv:2402.01135. 624

Joseph L Fleiss. 1971. Measuring nominal scale agree- 625
ment among many raters. *Psychological bulletin*, 626
76(5):378. 627

FreeCAD. <https://www.freecadweb.org/>. Online: 628
accessed 2025-04-20. 629

Jan Hartmann, Alistair Sutcliffe, and Antonella De An- 630
geli. 2008. Towards a theory of user judgment of 631
aesthetics and user interface quality. *ACM Trans-* 632
actions on Computer-Human Interaction (TOCHI), 633
15(4):1–30. 634

Timo Kapsalis. 2024. Cadgpt: Harnessing natural 635
language processing for 3d modelling to enhance 636
computer-aided design workflows. *arXiv preprint* 637
arXiv:2401.05476. 638

639	Mohammad Sadil Khan, Elona Dupont, Sk Aziz Ali, Kseniya Cherenkova, Anis Kacem, and Djamila Aouada. 2024a. Cad-signet: Cad language inference from point clouds using layer-wise sketch instance guided attention. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 4713–4722.	696
640		697
641		698
642		699
643		700
644		701
645		
646	Mohammad Sadil Khan, Sankalp Sinha, Talha Uddin Sheikh, Didier Stricker, Sk Aziz Ali, and Muhammad Zeshan Afzal. 2024b. Text2cad: Generating sequential cad models from beginner-to-expert level text prompts. <i>arXiv preprint arXiv:2409.17106</i> .	702
647		703
648		704
649		705
650		
651	Namyoun Kim, Kai Tzu iunn Ong, Yeonjun Hwang, Minseok Kang, Iiseo Jihn, Gayoung Kim, Minju Kim, and Jinyoung Yeo. 2025. Principles: Synthetic strategy memory for proactive dialogue agents. <i>Preprint</i> , arXiv:2509.17459.	706
652		707
653		708
654		709
655		
656	Chuang Li, Yang Deng, Hengchang Hu, Min-Yen Kan, and Haizhou Li. 2024a. Incorporating external knowledge and goal guidance for llm-based conversational recommender systems. <i>Preprint</i> , arXiv:2405.01868.	710
657		711
658		712
659		713
660		714
661	Jiahao Li, Weijian Ma, Xueyang Li, Yunzhong Lou, Guichun Zhou, and Xiangdong Zhou. 2025. Cad-llama: Leveraging large language models for computer-aided design parametric 3d model generation. <i>Preprint</i> , arXiv:2505.04481.	715
662		716
663		717
664		718
665		719
666		720
667	Xueyang Li, Yu Song, Yunzhong Lou, and Xiangdong Zhou. 2024b. Cad translator: An effective drive for text to 3d parametric computer-aided design generative modeling. In <i>Proceedings of the 32nd ACM International Conference on Multimedia</i> , pages 8461–8470.	721
668		722
669		723
670		724
671		725
672	Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	726
673		727
674		728
675		
676	Tianjian Liu, Hongzheng Zhao, Yuheng Liu, Xingbo Wang, and Zhenhui Peng. 2024. Compeer: A generative conversational agent for proactive peer support. In <i>Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology</i> , UIST '24, New York, NY, USA. Association for Computing Machinery.	729
677		730
678		731
679		732
680		733
681		734
682		
683	Xingyu Bruce Liu, Shitao Fang, Weiyan Shi, Chien-Sheng Wu, Takeo Igarashi, and Xiang'Anthony' Chen. 2025. Proactive conversational agents with inner thoughts. In <i>Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems</i> , pages 1–19.	735
684		736
685		737
686		738
687		739
688		
689	Yaxi Lu, Shenzhi Yang, Cheng Qian, Guirong Chen, Qinyu Luo, Yesai Wu, Huadong Wang, Xin Cong, Zhong Zhang, Yankai Lin, Weiwen Liu, Yasheng Wang, Zhiyuan Liu, Fangming Liu, and Maosong Sun. 2024a. Proactive agent: Shifting llm agents from reactive responses to active assistance. <i>Preprint</i> , arXiv:2410.12361.	740
690		741
691		742
692		743
693		744
694		745
695		
	Yaxi Lu, Shenzhi Yang, Cheng Qian, Guirong Chen, Qinyu Luo, Yesai Wu, Huadong Wang, Xin Cong, Zhong Zhang, Yankai Lin, and 1 others. 2024b. Proactive agent: Shifting llm agents from reactive responses to active assistance. <i>arXiv preprint arXiv:2410.12361</i> .	746
		747
		748
		749
		750
	Matt D Nelson, Brady L Goenner, and Bruce K Gale. 2023. Utilizing chatgpt to assist cad design for microfluidic devices. <i>Lab on a Chip</i> , 23(17):3778–3784.	746
		747
		748
		749
		750
	Jakob Nielsen and 1 others. 2012. Usability 101: Introduction to usability.	746
		707
	Onshape. https://www.onshape.com/en/ . Online: accessed 2025-04-20.	740
		709
	Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In <i>Proceedings of the 36th annual acm symposium on user interface software and technology</i> , pages 1–22.	710
		711
		712
		713
		714
		715
	Cyril Picard, Kristen M Edwards, Anna C Doris, Brandon Man, Giorgio Giannone, Md Ferdous Alam, and Faez Ahmed. 2025. From concept to manufacturing: Evaluating vision-language models for engineering design. <i>Artificial Intelligence Review</i> , 58(9):288.	716
		717
		718
		719
		720
	Shuofei Qiao, Runnan Fang, Zhisong Qiu, Xiaobin Wang, Ningyu Zhang, Yong Jiang, Pengjun Xie, Fei Huang, and Huajun Chen. 2024. Benchmarking agentic workflow generation. <i>arXiv preprint arXiv:2410.07869</i> .	721
		722
		723
		724
		725
	David Robertson and Thomas J Allen. 1993. Cad system use and engineering performance. <i>IEEE Transactions on Engineering Management</i> , 40(3):274–282.	726
		727
		728
	Omar Shaikh, Shardul Sapkota, Shan Rizvi, Eric Horvitz, Joon Sung Park, Diyi Yang, and Michael S Bernstein. 2025. Creating general user models from computer use. In <i>Proceedings of the 38th Annual ACM Symposium on User Interface Software and Technology</i> , pages 1–23.	729
		730
		731
		732
		733
		734
	Chenglei Si, Yanzhe Zhang, Ryan Li, Zhengyuan Yang, Ruibo Liu, and Diyi Yang. 2024. Design2code: Benchmarking multimodal code generation for automated front-end engineering. <i>arXiv preprint arXiv:2403.03163</i> .	735
		736
		737
		738
		739
	SolidWorks. https://www.solidworks.com/ . Online: accessed 2025-04-20.	740
		741
	Nan Sun, Bo Mao, Yongchang Li, Di Guo, and Huaping Liu. 2024. Assistantx: An llm-powered proactive assistant in collaborative human-populated environment. <i>arXiv preprint arXiv:2409.17655</i> .	742
		743
		744
		745
	Muhammad Usama, Mohammad Sadil Khan, Didier Stricker, and Muhammad Zeshan Afzal. 2025. Nurbgen: High-fidelity text-to-cad generation through llm-driven nurbs modeling. <i>arXiv preprint arXiv:2511.06194</i> .	746
		747
		748
		749
		750

849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899

A Appendix

A.1 Related Work

Text-to-CAD Generation Recent advances in natural language processing and large language models (LLMs) have enabled the direct generation of parametric CAD models from text, significantly streamlining design workflows and diminishing reliance on traditional inputs such as sketches. Current text-to-CAD methods primarily follow two core paradigms: the first is a sequence generation-based paradigm, which typically utilizes Transformer (Khan et al., 2024b; Li et al., 2024b) or LLMs (Kapsalis, 2024; Badagabettu et al., 2024; Xie and Ju, 2025; Zhang et al., 2025; Picard et al., 2025; Usama et al., 2025) as the backbone network to map textual descriptions to discrete sequences of CAD operations or modeling scripts (e.g., parametric commands like Sketch, Extrude, Fillet). Through training on high-quality text-parameter sequence pairs, the model can autoregressively generate logically coherent modeling steps to reconstruct 3D shapes. The second is a code/function-call paradigm, which focuses on leveraging the powerful code generation and spatial reasoning capabilities of LLMs to directly generate executable CAD script code from natural language instructions (Badagabettu et al., 2024; Xie and Ju, 2025), bypassing intermediate, task-specific command sequences to streamline the generation pipeline and enhance code reusability. However, these methods fundamentally suffer from the following limitations: the drawback of the single-pass generation mode, where they typically adopt a one-shot generation approach. When textual descriptions are incomplete, ambiguous, or vague, the models struggle to autonomously resolve such uncertainties, often resulting in a misalignment between the output and user intent; and inadequate utilization of multimodal signals, as traditional sequence generation primarily relies on parametric sequence signals while overlooking the inherent multimodal nature of CAD models (i.e., parametric sequences and rendered visual objects). Some cutting-edge work (e.g., (Wang et al., 2025)) has begun to explore incorporating visual feedback into the training loop, using visual scores to correct parametric sequences in order to enhance the visual quality of generated models and the accurate understanding of instructions. Our proposed proactive text-to-CAD paradigm aims to fundamentally address the limitations of the single-pass generation mode. Through

iterative refinement and an active interaction process, our method can progressively bridge the gap between descriptions and models, effectively improving the alignment of generation results with user expectations.

Proactive Agent Proactive agents have garnered significant research attention in recent years (Wang et al., 2024). From a behavioral perspective, current research advances primarily focus on enhancing agent proactivity along two dimensions: Proactive Conversation AI (Liu et al., 2024; Zhang et al., 2024a; Liu et al., 2025; Sun et al., 2024; Yang et al., 2025a; Wu et al., 2025), which concentrates on enabling agents to initiate, intervene in, and sustain dialogues, providing timely and relevant linguistic assistance or information based on user context and latent needs to optimize human-computer interaction experiences; and Proactive Generative/Action AI (Lu et al., 2024a; Zhao et al., 2025b,a; Yang et al., 2025b; Chen et al., 2025b; Park et al., 2023), which emphasizes the ability of agents to autonomously make decisions, plan, and execute tasks or generate content, allowing them to predict user intentions and independently take actions to achieve goals, thereby shifting from passive response to active service. This work is primarily concerned with the former, and in recent years, technological developments in proactive conversational interactive agents have evolved along two core trajectories. The first is a research direction centered on proactive single-agent systems (Deng et al., 2023; Wei et al., 2022; Deng et al., 2023; Wang et al., 2023; Kim et al., 2025; Lu et al., 2024b), which focus on enabling agents to exhibit self-driven, task-oriented behaviors in dynamic environments through internal state tracking, goal decomposition, and hierarchical policy planning. Such systems have made significant progress in autonomous decision-making, goal-driven dialogue, and achieving adaptability and proactivity in single-agent settings. The second trajectory is a research direction focused on proactive multi-agent systems (Fang et al., 2024; Zhang et al., 2024b,a; Cheng et al., 2024; Li et al., 2024a; Shaikh et al., 2025), emphasizing the resolution of collaborative mechanisms, including structured communication, consensus formation, and dynamic role assignment. This path aims to enable multiple agents to efficiently collaborate in open, multi-participant scenarios to achieve system-level objectives collectively, representing a key approach to building

951	proactivity in complex applications.	
952	A.2 Dataset Specifications	
953	A.2.1 User Initial Query Generation	
954	We first define the main types and subtypes of CAD	
955	operations and parameters, and remove a subset of	
956	the type information contained in the minimal meta-	
957	data. Subsequently, we perform a topological sort	
958	on the missing CAD types based on the dependen-	
959	cies among the absent type information.	
960	Main Types and Subtypes of CAD Operations	
961	and Parameters. We adopt the same represen-	
962	tation method proposed by Khan et al. (2024a),	
963	which uses a sketch-and-extrude format. Each 2D	
964	sketch consists of multiple faces, each face consists	
965	of multiple loops, and each loop either contains a	
966	line and an arc or a circle. Loops are always closed	
967	(i.e., the start and end coordinates are the same).	
968	The specific descriptions of the main types and sub-	
969	types in our dataset are presented in Table 12.	
970	Dependencies Between CAD Types The dependen-	
971	cies between CAD types are shown in Figure 7.	
972	After deleting missing types (including both main	
973	types and subtypes), we use topological sorting to	
974	order the missing types, laying the groundwork for	
975	the sequence of proactive dialogue generation.	
976	A.2.2 Proactive Dialogues Generation	
977	To facilitate the reproducibility of our dataset, we	
978	employ DeepSeek-R1 to generate proactive dia-	
979	logues with users. The system proactively asks	
980	users questions based on the missing type objects,	
981	and users retrieve corresponding information from	
982	the initial dataset to answer the system’s questions.	
983	The prompts for the system’s responses and user	
984	inquiries are shown in Table 13.	
985	The initial text of the user, namely the initial text	
986	description of the CAD design, is adapted from the	
987	text descriptions in the initial dataset. We also use	
988	DeepSeek-R1 for this adaptation. The prompts for	
989	the adaptation are shown in Table 14.	
990	A.2.3 Dataset Statistics	
991	We classify the samples according to the number	
992	of missing types in each sample. Since the sub-	
993	types within the main types vary across dataset	
994	samples (e.g., different sketches may contain dif-	
995	ferent numbers of loop subtypes), which makes it	
996	inconvenient to count, we uniformly consider the	
997	absence of one main type and one subtype as a	
998	single type of missing information, and include it	
999	in the count of missing types.	
	Our dataset contains approximately 4,590 dia-	1000
	logues, with 28,110 question-answer pairs. The	1001
	distribution of the number of missing types per	1002
	sample is naturally right-skewed and discrete, con-	1003
	centrated in the range of 4 to 9, with a long right	1004
	tail, consistent with the characteristics of a skewed	1005
	discrete count data distribution. The specific distri-	1006
	bution of the dataset is shown in Figure 8.	1007
	Level Distribution Due to the complex de-	1008
	pendency and nesting relationships between CAD	1009
	operation parameters, we observe that: (1) miss-	1010
	ing one main category often cascades to multiple	1011
	missing subcategories, and (2) while the number	1012
	of main categories remains fixed per CAD model,	1013
	subcategory counts exhibit significant inter-model	1014
	variation.	1015
	We assign equal missing probability to each type	1016
	during dataset construction and count the final num-	1017
	ber of missing types per sample. The natural distri-	1018
	bution of missing types in the dataset primarily falls	1019
	within the range of 1 to 20, as shown in Figure 8.	1020
	A.2.4 Inter-Rater Reliability (IRR)	1021
	Assessment Report	1022
	Methodology. We evaluated inter-rater agree-	1023
	ment using Fleiss’ Kappa (Fleiss, 1971) across	1024
	three dimensions of binary classification tasks, with	1025
	three independent raters per sample:	1026
	1. Initial Query Completeness.	1027
	Assesses whether user queries contain all type	1028
	information except explicitly missing types)	1029
	2. Question Coverage.	1030
	Verifies if questions address all missing cate-	1031
	gory details)	1032
	3. Response Specificity.	1033
	Evaluates if responses provide sufficiently de-	1034
	tailed information	1035
	Annotation Protocol. Each sample was indepen-	1036
	dently rated by 3 annotators using binary scoring	1037
	(1=meets requirements, 0=fails).	1038
	Calculation Process	1039
	1. Sample-level Agreement.	1040
	For each sample, compute agreement score	1041
	based on raters’ category distributions using	1042
	standard Fleiss’ formula.	1043
	2. Mean Observed Agreement (\bar{P}).	1044
	Average agreement scores across all samples,	1045
	reflecting actual consensus.	1046

1047 3. **Expected Random Agreement** (\bar{P}_e).
 1048 Calculate chance agreement probability based
 1049 on overall category frequencies.

4. **Kappa Coefficient.**

$$\kappa = \frac{\bar{P} - \bar{P}_e}{1 - \bar{P}_e}$$

1050 The results demonstrate strong inter-rater agree-
 1051 ment across all dimensions:

- 1052 1. Initial query completeness ($\kappa = 0.828$)
- 1053 2. Question coverage of missing details ($\kappa =$
 1054 0.786)
- 1055 3. Response specificity ($\kappa = 0.807$)

1056 These metrics demonstrate strong inter-
 1057 annotator agreement, validating both the
 1058 annotation protocol’s consistency and the resultant
 1059 dataset’s integrity.

1060 **A.3 Prompt List**

1061 **A.3.1 Dataset Construction**

1062 Prompts for generating dialogues between the sim-
 1063 ulated system and user are presented in Table 13.
 1064 Prompts for rewriting the text of the initial dataset
 1065 to generate CAD initial queries are presented in
 1066 Table 14.

1067 **A.3.2 Proposed Method: Proactive Agent**

1068 The prompts for generating workflows by the strate-
 1069 gic module are presented in Table 17. The prompts
 1070 for dispatching tasks by the dispatch sub-module,
 1071 for validating task completion by the validation
 1072 sub-module, for initiating questions by the asking
 1073 sub-module, and for generating summaries by the
 1074 summary sub-module are presented in Table 16.

1075 **A.3.3 Baseline**

1076 The prompts for the baseline methods, including
 1077 their detailed specifications, are presented in Ap-
 1078 pendix A.5.

1079 **A.3.4 Evaluation**

1080 The prompts for simulating user responses in the
 1081 dialogue-level experiments are presented in Table
 1082 15.

A.4 Details of Our Methods 1083

1084 The transition function μ of a Hierarchical Finite
 1085 State Machine (HFSM) is represented in Table 7.

1086 In this context, y represents the input to the cur-
 1087 rent state, and the current state transitions to the
 1088 next state through the transition function μ com-
 1089 bined with the input y and condition e , where e
 1090 is the output generated by the operational sub-
 1091 module.

1092 The overall execution logic of our hierarchical
 1093 finite state machine (HFSM) begins with the user’s
 1094 initial CAD design text description t , which first
 1095 enters the workflow sub-module of the strategic
 1096 module. At this point, it is in the initial state (S_0),
 1097 where it plans how to ask the user for missing infor-
 1098 mation in order to ultimately generate a complete
 1099 CAD design text, thereby generating the workflow
 1100 of the operational sub-module. Then, the gener-
 1101 ated workflow is input into the tactical module and
 1102 enters the dispatch sub-module for workflow pars-
 1103 ing. The workflow parsing generates a series of
 1104 subtasks, which the dispatch sub-module assigns to
 1105 the operational sub-module. At this point, it transi-
 1106 tions to state S_1 . Next, the asking sub-module or
 1107 summary sub-module of the operational module,
 1108 which receives the sub-tasks, executes the tasks,
 1109 transitioning to state S_3 or S_4 . The completion of
 1110 tasks by the asking and summary modules is ver-
 1111 ified by the validation sub-module of the tactical
 1112 module, which means transitioning to state S_2 for
 1113 validation. If the validation sub-module determines
 1114 that the task of the asking sub-module is complete,
 1115 it re-enters the dispatch module, transitioning to
 1116 the next task assignment state S_1 . Otherwise, it
 1117 returns to the asking sub-module to regenerate the
 1118 question q . The validation of the summary module
 1119 is similar to that of the asking module, except that
 1120 when the validation sub-module determines that the
 1121 task is complete, i.e., the final CAD design solution
 1122 has been generated, it transitions to the termination
 1123 state S_{exit} , completing the entire execution logic of
 1124 the HFSM.

A.5 Baseline 1125

A.5.1 Baseline Methods 1126

1127 Since no existing large language model (LLM)
 1128 methods currently address proactive querying or
 1129 interactive capabilities for CAD applications, we
 1130 establish six baseline approaches based on our
 1131 newly constructed Interactive-CAD dataset and
 1132 novel CAD generation paradigm. These baselines

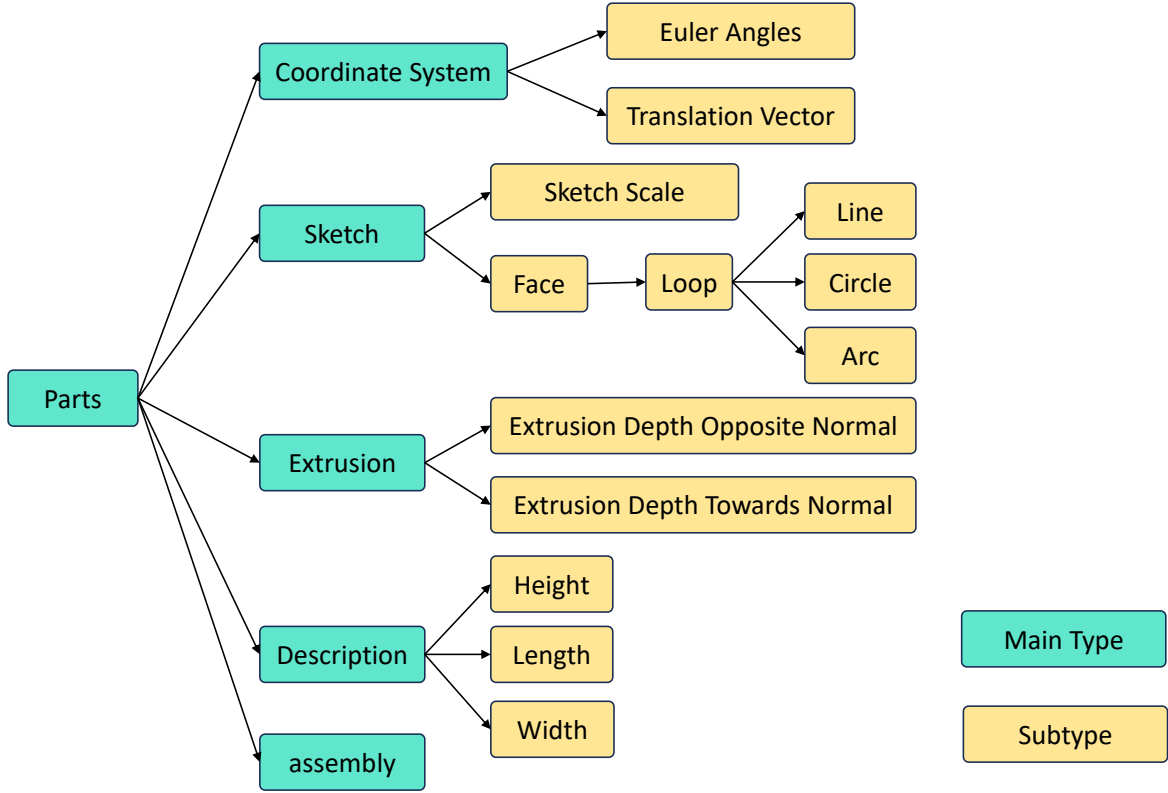


Figure 7: Dependencies between CAD types.

incorporate both commonly-used and state-of-the-art (SOTA) LLM-based proactive dialogue methodologies.

Standard Prompting (Deng et al., 2023): Given a task description q , the user’s initial CAD text description t , and dialogue history C , we instruct the LLM to perform CAD proactive querying and generate response r . The task descriptions and prompt templates are provided in Table 19 of Appendix A.5.2. This prompting scheme can be formally represented as:

$$p(r \mid q, t, C)$$

CoT (Wei et al., 2022): A chain-of-thought prompting approach that generates intermediate reasoning steps to derive the final response. In the task description, we require the system to simulate the next response based on the current dialogue history. The prompt template is provided in Table 20 of Appendix A.5.2.

Proactive Prompting (Deng et al., 2023): Proactive Prompting is designed to provide the LLM with alternative options for determining appropriate actions to take in responses, formally represented as:

$$p(a, r \mid q, t, C, A)$$

Given the task description q , the user’s initial CAD text description t , dialogue history C , and a set of possible dialogue actions A , this approach instructs the LLM to: (1) select the most suitable dialogue action $a \in A$, and (2) generate the corresponding response r . To adapt this to CAD proactive querying tasks, we define the dialogue actions as either “querying the user about missing CAD design details” or “summarizing the CAD design solution”. The prompt templates are provided in Table 21 of Appendix A.5.2.

ProCoT (Deng et al., 2023): The proactive chain-of-thought prompting scheme (ProCoT) involves dynamic reasoning and planning to analyze subsequent actions for achieving dialogue objectives, formally represented as:

$$p(c, a, r \mid q, t, C, A)$$

where c denotes the cognitive description of the decision-making process for subsequent actions, while q , t , C , A , and r maintain the same definitions as in **Proactive Prompting**. For CAD proactive querying tasks, we define c as the analysis of missing detail types in the current CAD design. The prompt templates are provided in Table 22 of Appendix A.5.2.

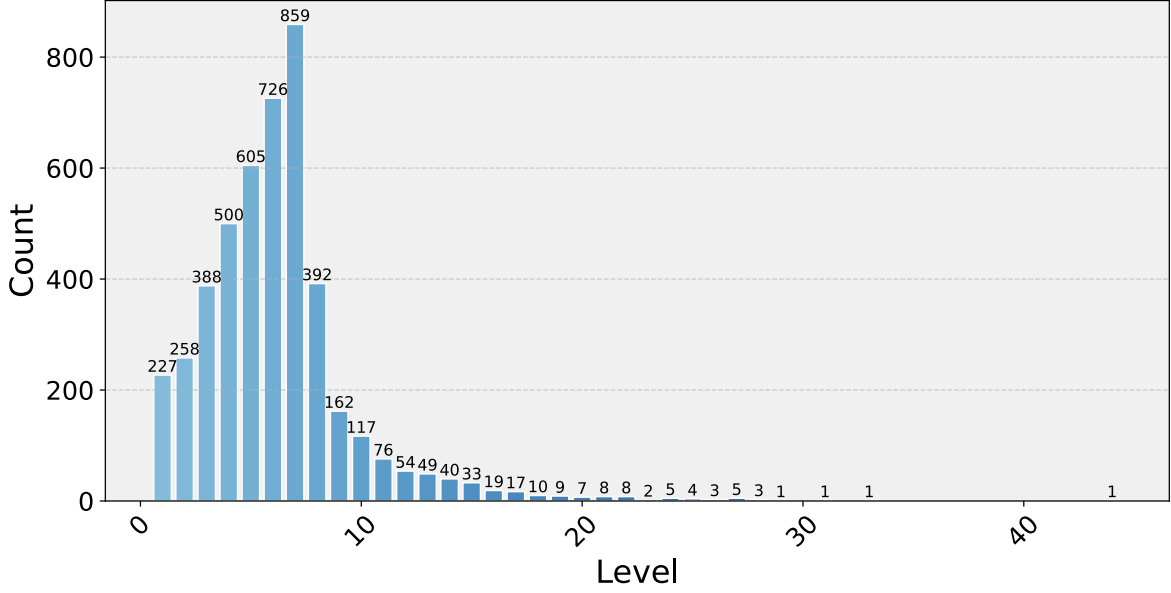


Figure 8: The statistical distribution of our dataset.

Current State	Condition/Input	Next State
S_{exit}	$y = t$	S_0
S_0	$y = \mathcal{G}(V, E)$	S_1
S_1	action = ask	S_3
S_1	action = summarize	S_4
S_3	$e = q$	S_2
S_4	$e = p$	S_2
S_2	action = ask \wedge task completion is valid	S_1
S_2	action = summarize \wedge task completion is valid	S_{exit}
S_2	action = ask \wedge task completion is invalid	S_3
S_2	action = summarize \wedge task completion is invalid	S_4
Other	All other cases	S_{exit}

Table 7: State transition function $\mu(S_i, y)$

PS+ Prompting (Wang et al., 2023): Plan-and-Solve Prompting (PS prompting) is a two-stage methodology that first generates both the reasoning process and potential answers through logical inference, then employs answer extraction prompts to derive the final solution. We adopt PS+ prompting with more detailed instructions, defining the task questions as: (1) “identifying missing details in the current CAD design” and (2) “formulating user queries about these missing details”, with the dialogue history incorporated into the questions to provide complete contextual information about the current CAD design. The reasoning prompt templates and answer extraction prompts are detailed in Table 23 of Appendix A.5.2.

MACRS (Fang et al., 2024): Multi-Agent Con-

versational Recommender System (MACRS) is a collaborative multi-agent framework that integrates four LLM-based agents to plan diverse dialogue behaviors, employing a feedback-aware reflection mechanism for agent adaptation. Specifically, MACRS incorporates four specialized agents designed to perform distinct dialogue functions: questioning, small talk, recommendation, and planning. For proactive querying in CAD design tasks, we define the user profile U_p in MACRS’s memory module as representing a CAD modeler, whose objective is to create a fully detailed CAD model. The system initializes the interaction using the user’s initial CAD text description as input. The questioning agent, small talk agent, and recommendation agent generate three candidate responses (R_{ask} ,

1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214

R_{chat} , and R_{rec}) based on dialogue history and the user profile. The planning agent (π_p) then performs multi-step reasoning to select the most appropriate response from these candidates, determining the final system response R_s :

$$R_s = \pi_p(I_{\text{plan}}, R_{\text{ask}}, R_{\text{rec}}, R_{\text{chat}}, D_h, A_h)$$

where I_{plan} represents the instruction for the planning agent π_p , D_h denotes the dialogue history, and A_h corresponds to the history of dialogue actions. The prompt templates for each agent can be found in Table 18 of Appendix A.5.2.

A.5.2 Prompts in Baseline Methods

The prompts for standard prompting are listed in Table 19, for CoT in Table 20, for Proactive prompting in Table 21, for ProCoT in Table 22, for PS+ prompting in Table 23, and for MACRS in Table 18.

A.6 Automated Evaluation

We design several evaluation metrics at three different levels: turn-level, dialogue-level and CAD-level.

A.6.1 Turn-Level Evaluation

Following Zhang et al. (2024c), we use **Clarification Accuracy (Clari. Acc.)** to measure the system’s ability to actively query incomplete CAD design texts, where a score of 1 is assigned if the response is a valid question, otherwise 0. Besides, we also involve **Rule-based Score** to assess whether responses target the correct missing attribute category, where a binary score (1/0) is determined by the presence of keywords relevant to the target attribute. Moreover, **ROUGE-L** (Lin, 2004) and **BERTScore** (Zhang et al., 2020) is used to quantify semantic similarity between responses and reference questions. Note that, we evaluate all turn-level metrics through both Micro and Macro averaging, where Micro-averaging computes the mean score per turn, while Macro-averaging calculates the mean score per dialogue.

A.6.2 Dialogue-Level Evaluation

Completion Rate. This metric measures the degree to which missing information in the CAD text description is completed during the dialogue, calculated as:

$$\text{Completion Rate} = \frac{N_c}{N_t} \quad (6)$$

where N_t refers to the count of initially absent CAD operations or parameter types in the user’s input, and N_c denotes the count of missing categories addressed (queried) during the dialogue.

Effectiveness Rate. This metric evaluates the proportion of effective questions in the dialogue, defined as:

$$\text{Effectiveness Rate} = \frac{N_e}{N_i} \quad (7)$$

where N_e denotes the number of effective questions that satisfy all of the following criteria: (a) targeting a missing category, (b) having not been previously asked, and (c) receiving a non-"unknown" response, and N_i refers to the total number of interaction turns.

A.6.3 CAD-Level Evaluation

To evaluate the quality of the final generated CAD, we leverage **Chamfer Distance (CD)** (Fan et al., 2016; Wu et al., 2021b), which can measure geometric similarity between generated 3D CAD models and ground truth, where lower values indicate higher geometric fidelity. In addition to automatic evaluation, we also involve manual evaluation to comprehensively assess the CAD model quality. **Average Rank** metric is used in manual evaluation, which can be denoted as

$$\text{Average Rank} = \frac{1}{N} \sum_{i=1}^N R_i \quad (8)$$

where N is the total number of CAD models generated by different methods and R_i is the rank of the i -th model (lower values denote better alignment)

A.7 Human Evaluation Framework

A.7.1 Evaluation Metrics

Functional Perception. We evaluate this dimension using the following two metrics:

- **Design Intent Consistency (DIC):** The degree to which the generated CAD model aligns with the target design intent, including dimensions, constraints, and structural features.
- **Task Efficiency (TaskEff):** The extent to which the system enables users to complete the design with minimal time and effort.

Interaction Perception. This dimension is assessed through three metrics:

- **Guidance Effectiveness (GE):** The effectiveness of the system’s proactive inquiries in clarifying ambiguous user intent.
- **Learnability:** The ease with which a new user can start using the system without prior training.
- **Information Clarity (IC):** The clarity and organization of parameters and feature-tree information.

Emotional Perception. The following two metrics are used to measure this aspect:

- **Aesthetic or Stylistic Appeal (ASA):** The extent to which the generated CAD models demonstrate aesthetic quality and stylistic coherence, including visual elegance, proportional harmony, and overall design attractiveness.
- **Interaction Experience Satisfaction (IES):** The overall satisfaction with the interaction process and responsiveness of the system.

A.7.2 Human Evaluation Questionnaire

We present our questionnaire in Table 10 and Table 11.

A.7.3 Annotator Demographics

All annotators possessed fundamental computer operation skills and were either students in computer-related fields or CAD design practitioners, each with prior data annotation experience. Participants were recruited as uncompensated volunteers through our academic collaboration network, and they contributed their expertise voluntarily. All participants were explicitly informed that the data would be used exclusively for academic research purposes. Given that the participants are students or professionals in related fields, their involvement is closely aligned with their own learning or professional interests. Therefore, this volunteer-based model, founded on shared interests and professional development, is considered appropriate and sufficient. The annotator group was evenly distributed in terms of CAD expertise, consisting of approximately 50% beginners (e.g., junior students) and 50% proficient users (e.g., senior students, graduate students, and practitioners). The two cohorts were identified through predefined screening questions, such as familiarity with CAD concepts and years of experience using relevant

software, to ensure the data encompasses both novice and expert perspectives.

A.7.4 Human Annotation Filtering

To maintain annotation reliability, a multi-stage filtering approach incorporating text evaluation, trap questions, and consistency checks was employed.

- **Text Evaluation:** We manually reviewed the initial design descriptions provided by annotators. If an annotator’s description significantly deviated from the intended model, they were identified as engaging in perfunctory behavior, and their questionnaire was deemed invalid.
- **Trap Questions:** Some questionnaires included embedded trap questions containing explicit instructions, such as “Select Option A for all items in this questionnaire.” Annotators who failed to follow these instructions were considered inattentive, and all their submissions were invalidated.
- **Consistency Check:** We manually compared each annotator’s selected options with their provided textual justifications. Responses showing inconsistency between the two were discarded. Additionally, annotators with low agreement compared to others underwent manual review. If their responses exhibited signs of random selection or perfunctory behavior, all questionnaires from that annotator were excluded.

The filtering procedure guarantees the retained annotations’ attentiveness and internal consistency, which enhances evaluation quality and ensures accurate, reliable results.

Method	Median CD ↓		Mean CD ↓	
	Text2CQ	Text2CAD	Text2CQ	Text2CAD
Initial Query	94.68		160.30	
ProCoT	66.07	28.00	144.89	116.13
Ours	64.27	20.60	123.50	99.70

Table 8: The experimental results by changing CAD generation module.

A.8 Discussion

To further validate the generality of our proposed paradigm across CAD generation architectures, we decoupled our framework by replacing the Text2CAD Transformer in the CAD generation module with a fine-tuned Qwen2.5-3B within the

Text-to-CadQuery framework (Xie and Ju, 2025), and tested it on our dataset. The experimental results are presented in Table 8. The findings demonstrate that the quality of CAD models generated through our proactive inquiry paradigm significantly surpasses that of models generated without proactive inquiry. Moreover, our method remains superior to the strongest baseline even after replacing the CAD generation module. Interestingly, the Text-to-CadQuery framework yields lower-quality CAD models compared to the Text2CAD Transformer, which contradicts the conventional view that code generation methods generally outperform traditional approaches (Xie and Ju, 2025). We hypothesize that this discrepancy arises because CadQuery relies on geometry kernels such as OpenCASCADE, which are highly sensitive to numerical precision. This sensitivity may lead to reduced geometric accuracy in CAD models generated from incomplete design text inputs. In summary, the experimental results confirm that our proposed paradigm effectively enhances the alignment between generated CAD models and user expectations across different CAD generation architectures, thereby validating the efficacy of our approach.

A.9 Software Environment

All experiments were implemented in Python 3.12.0 within a Conda environment. The version information of the major software components is presented in Table 9.

Software	Version
Python	3.12.0
Conda	24.1.2
PyTorch	2.2.1
Transformers	4.38.2
CUDA	12.4
PEFT	0.9.0
Flash Attention	2.5.8
Accelerate	0.27.2
Datasets	2.18.0
NumPy	1.26.4
OpenCV-Python	4.9.0.80
AV (PyAV)	12.0.0
scikit-learn	1.4.1.post1
pandas	2.2.1
tqdm	4.66.2
nltk	3.8.1
rouge-score	0.1.2
bert-score	0.3.13
PyYAML	6.0.1
loguru	0.7.2
requests	2.31.0

Table 9: Software Environment

Survey Questionnaire: Comparative Evaluation of Two CAD Generation Systems
<p>This questionnaire is designed to compare the performance of two CAD generation systems (System A and System B) under the same design task. Based on your experience with both systems, please evaluate them across several dimensions and provide your overall assessment.</p>
<p>Q1. Design Intent Consistency</p>
<p>Which system produced models more consistent with the intended design? [Better]: Model matches target, with correct dimensions, constraints, and features. [Weaker]: Model shows notable deviations or missing parameters. Choice: Option A: System A Option B: System B Option C: Tie / No significant difference Reason:</p>
<p>Q2. Task Efficiency</p>
<p>Which system allowed you to finish the task more efficiently? [Better]: Simple and efficient, minimal steps/time required. [Weaker]: Time-consuming or requires extra steps. Choice: Option A: System A Option B: System B Option C: Tie / No significant difference Reason:</p>
<p>Q3. Guidance Effectiveness</p>
<p>Which system can more accurately recognize users' vague intents? [Better]: Identifies ambiguous intent and provides actionable options. [Weaker]: Fails to clarify intent, with vague/unhelpful prompts. Choice: Option A: System A Option B: System B Option C: Tie / No significant difference Reason:</p>
<p>Q4. Learnability</p>
<p>Which system is easier for new users to learn? [Better]: Intuitive and straightforward, easy to master without prior training. [Weaker]: Complex and non-intuitive, requiring significant effort to learn. Choice: Option A: System A Option B: System B Option C: Tie / No significant difference Reason:</p>

Table 10: Human Evaluation Questionnaire (a)

<p>Q5. Information Clarity</p> <p>Which system provides clearer information about parameters and feature trees? [Better]: Parameters well-organized, feature tree easy to interpret/modify. [Weaker]: Parameters unclear, feature tree disorganized or confusing. Choice: Option A: System A Option B: System B Option C: Tie / No significant difference Reason:</p>
<p>Q6. Aesthetic or Stylistic Appeal</p> <p>Which system produces CAD models with stronger aesthetic or stylistic appeal? [Better]: Strong aesthetic quality, appealing proportions, coherent style, polished design. [Weaker]: Lacks refinement, awkward proportions, inconsistent style, rough design. Choice: Option A: System A Option B: System B Option C: Tie / No significant difference Reason:</p>
<p>Q7. Interaction Experience Satisfaction</p> <p>Which system provides a more satisfying interaction experience overall? [Better]: Smooth interaction, prompt response, overall pleasant. [Weaker]: Fragmented or slow, reducing satisfaction. Choice: Option A: System A Option B: System B Option C: Tie / No significant difference Reason:</p>
<p>Q8. Overall Preference</p> <p>Overall, which system do you prefer? Choice: Option A: System A Option B: System B Option C: Tie / No significant difference Overall Reason: Additional Suggestions/Comments:</p>

Table 11: Human Evaluation Questionnaire (b)

Main Category	Description
Parts	The number of parts included in the current CAD model.
Coordinate System	Defines the coordinate system of the part, which includes the following subtypes: Euler Angles: Three parameters (θ, ϕ, γ) that determine the orientation of the sketch plane. Translation Vector: Three parameters (τ_x, τ_y, τ_z) that describe the translation of the sketch plane.
Sketch	Defines the geometry of the 2D sketch. Face: Defines a face that contains a closed loop, which includes the following subtypes: Loop: Defines a loop composed of lines, arcs, or circles. Line: Contains start and end coordinates. Arc: Contains start, mid, and end coordinates. Circle: Contains center and top-most coordinates. Sketch scale: The scaling factor for the 2D sketch.
Extrusion	Defines the parameters for the extrusion operation, which includes the following subtypes: Extrusion depth towards normal: The extrusion depth in the direction of the normal of the sketch plane. Extrusion depth opposite normal: The extrusion depth in the direction opposite to the normal of the sketch plane.
Description	Provides a description of the CAD model, which includes the following subtypes: Length: The length of the part. Width: The width of the part. Height: The height of the part.
Assembly	The assembly relationships between parts.

Table 12: Main types and subtypes of CAD operations and parameters.

System Simulation
<p>You are an AI assistant helping a user design a CAD model. The user has provided some information, but the {category} is missing.</p> <p>Generate a natural, conversational question asking the user to provide the missing information. Make your question specific to the context of the CAD model being designed. Keep your question concise and focused only on the missing {category}. Please only generate the question without any additional content.</p>
User Simulation
<p>You are a user designing a CAD model. An AI assistant has asked you about missing {category} information.</p> <p>Referential Design: {information about the missing type in the initial dataset}</p> <p>Please generate a natural, conversational response where you are playing the role of a user answering a question. Use the parameters from the referential design provided above in your response. Keep your answer concise and focused only on the {category}. Please only generate the answer without any additional content.</p>

Table 13: Prompts for simulating system-User interaction.

User Initial Query Generation
<p>I need to update a CAD model description query to reflect recent modifications. Certain information has been removed from the original JSON data.</p> <p>Task: Please revise the query to include only the information that remains available in the current JSON, ensuring all non-null details and relevant parameters are retained.</p> <p>Original JSON: {initial metadata}</p> <p>Current JSON: {metadata with partially missing types}</p> <p>Original query: {Initial Text Containing Complete CAD Information}</p> <p>Requirements for the Updated Query:</p> <ol style="list-style-type: none"> 1. Include only the information still present in the current JSON. 2. Exclude any fields that are now null. 3. Maintain relevance to CAD model design. 4. Ensure clarity and conciseness. <p>Please only reply with the modified query without generating any additional text.</p> <p>Example:</p> <p>Original query: "Create the first part of the CAD model, a rectangular prism with a curved side. Begin by creating a new coordinate system with Euler angles of [0.0, 0.0, 0.0] and a translation vector of [0.0, 0.0, 0.0]. Next, create a 2D sketch on the X-Y plane of the coordinate system. The sketch consists of four lines. The first line has a start point at [0.0, 0.0] and an end point at [0.6, 0.0]. The second line has a start point at [0.6, 0.0] and an end point at [0.6, 0.375]. The third line has a start point at [0.6, 0.375] and an end point at [0.0, 0.375]. The fourth line has a start point at [0.0, 0.375] and an end point at [0.0, 0.0]. Scale the 2D sketch by a factor of 0.6. Then transform the 2D sketch into a 3D sketch with Euler angles of [0.0, 0.0, 0.0] and a translation vector of [0.0, 0.0, 0.0]. Finally, extrude the 3D model with an extrusion depth towards the normal of 0.075 and an opposite normal depth of 0.0. Scale the sketch by 0.6. The first part of the CAD model has the following dimensions: a length of 0.6 units, a width of 0.6 units, and a height of 0.075 units, with a curved side. The part is centered in the image."</p> <p>Modified query: "Create the initial part of the CAD model, which is a rectangular prism featuring a curved side. Start by defining a new coordinate system with Euler angles set to [0.0, 0.0, 0.0]. Then, generate a 2D sketch on the X-Y plane of this coordinate system. The sketch comprises four lines: The first line starts at [0.0, 0.0] and ends at [0.6, 0.0]. The second line extends from [0.6, 0.0] to [0.6, 0.375]. The third line extends from [0.6, 0.375] to [0.0, 0.375]. The fourth line connects [0.0, 0.375] back to [0.0, 0.0]. Apply a scaling factor of 0.6 to the 2D sketch. Next, convert it into a 3D sketch with Euler angles of [0.0, 0.0, 0.0]. Proceed by extruding the 3D model with an extrusion depth along the normal direction of 0.075, while keeping the opposite normal depth at 0.0. Scale the sketch again by a factor of 0.6. The resulting first part of the CAD model has the following dimensions: a width of 0.6 units and a height of 0.075 units, maintaining a curved side. The part is centrally positioned within the image."</p>

Table 14: Prompts for generating user CAD initial text.

User Simulation in the Dialogue-level Evaluation
<p>Assume you are a user who needs to respond to the system's question. The system's question is: "{system_message}" Your designed corresponding parameters are: "{parameters}" Please generate a response to the system's question based on your designed parameters. If the parameter is "unknown" or other non-parameter information, respond with "unknown".</p>

Table 15: Prompts for dialogue-Level evaluation of user simulation.

Dispatch Agent
<p>You are a dispatcher agent. Your job is to classify the task into one of two categories based on the tags in the preceding task <>: 'ask' or 'summarize'. Return only the category name without any additional text. Classify the following task into either 'ask' or 'summarize' category: {task}</p>
Validation Agent
<p>You are a validation agent. Your job is to determine if the output meets the requirements of the task. Return only 'yes' or 'no' without any additional text. Task: {task} Output: {output} Does this output fulfill the task requirements? Answer with 'yes' or 'no'.</p>
Asking Agent
<p>You are an assistant for CAD operations. Generate a clear and concise question based on the task description to ask the user. Generate a question based on task instruction: {task instruction}</p>
Summary Agent
<p>You are an assistant for CAD operations. Generate a comprehensive summary based on the conversation history and the task description. Generate a summary based on the task instruction and the conversation history. Task instruction: {task instruction} Conversation history: [...]</p>

Table 16: Prompts for asking, summary, validation, and dispatch.

Workflow Generation

You are a helpful and intelligent task planner, and your target is to decompose the assigned task into multiple subtasks for task completion and analyze the precedence relationships among subtasks.

At the beginning of your interactions, you will be given the task description and actions list you can take to finish the task, and you should decompose the given task into subtasks that can be accomplished using the provided actions. And then, you should analyze the precedence relationships among these subtasks, ensuring that each subtask is sequenced correctly relative to others. Based on the analysis, you should construct a workflow consisting of the identified subtasks to complete the task.

You should use "Node:

1. <subtask 1>
2. <subtask 2>" to denote subtasks, and use (x,y) to denote that <subtask x> is a predecessor of <subtask y>, (START,x) to indicate the beginning with <subtask x>, and (x,END) to signify the conclusion with <subtask x>. Remember that x, y are numbers.

Your response should use the following format:

Node:

1.<subtask 1>
2.<subtask 2>

...

Edges:(START,1) ... (n,END)

Now it's your turn.

Task: Refine the user's initial CAD design to the greatest extent possible.

The user's initial CAD design: [...]

The action list you can take: ['ask', 'summary']

'ask' means to inquire from the user to obtain specific CAD modeling information.

'summarize' refers to formulating a CAD design plan by integrating the CAD design information from the conversation history.

Remember that the format of the Node must be:

Node:

1: <action type> subtask 1
2: <action type> subtask 2

...

Edges:(START,1) ... (n,END)

Table 17: Prompts for generating workflows.

<p>Asking Agent</p> <p>CAD details design task: Based on the conversation history and CAD design priority rules, generate a question for the user regarding the highest-priority CAD design detail that has not yet been completed in the conversation history. The given conversation history is [...]</p> <p>You are a knowledgeable and enthusiastic CAD design details recommender chatbot.</p> <p>Your goal is to engage in friendly, casual conversation about CAD design details. Follow these guidelines:</p> <ul style="list-style-type: none"> - Don't say that you can't give recommendations directly. - As you are a chatbot, speak casually but not too informally. - Respond appropriately to the seeker's answers in line with your role. <p>You should elicit CAD design details by asking questions.</p> <p>If user asked any question at previous turn, You should answer the question.</p> <p>If there is nothing to respond to, please use the response "Alright!"</p> <p>Response should be equal or less than 15 words.</p>
<p>ChitChat Agent</p> <p>CAD details design task: Based on the conversation history and CAD design priority rules, generate a question for the user regarding the highest-priority CAD design detail that has not yet been completed in the conversation history. The given conversation history is [...]</p> <p>You are a knowledgeable and enthusiastic CAD design details recommender chatbot.</p> <p>Your goal is to engage in friendly, casual conversation about CAD design details. Follow these guidelines:</p> <ul style="list-style-type: none"> - Don't say that you can't give recommendations directly. - As you are a chatbot, speak casually but not too informally. - Respond appropriately to the seeker's answers in line with your role. <p>You should elicit CAD design details by asking questions.</p> <p>If user asked any question at previous turn, You should answer the question.</p> <p>If there is nothing to respond to, please use the response "Alright!"</p> <p>Response should be equal or less than 15 words.</p>
<p>Planning Agent</p> <p>CAD details design task: Based on the conversation history and CAD design priority rules, generate a question for the user regarding the highest-priority CAD design detail that has not yet been completed in the conversation history. The given conversation history is [...]</p> <p>You are a knowledgeable and enthusiastic planning agent decide which response to generate.</p> <p>Your goal is to engage in friendly, casual conversation about CAD design details. Follow these guidelines:</p> <ul style="list-style-type: none"> - Don't say that you can't give recommendations directly. - As you are a chatbot, speak casually but not too informally. - Respond appropriately to the seeker's answers in line with your role. <p>response from asking agent: {asking agent response}</p> <p>response from chit-chatting agent: {chit-chatting agent response}</p> <p>From the conversation history, determine whether user CAD design details is sufficient or not.</p> <p>Must choose one of the candidate responses based on three different dialogue acts. These three dialogue acts are: asking, or chit-chatting.</p> <p>If there is nothing to respond to, please use the response "Alright!"</p>

Table 18: Prompts for MACRS.

Standard Prompting

Your current task is to determine the user's intentions and satisfy their needs based on the provided conversation between the user and the system. The given conversation history is [...]. Please generate the response. If there is nothing to respond to, please use the response "Alright!"

Table 19: Prompts for standard prompting.

CoT

Your current task is to determine the user's intentions and satisfy their needs based on the provided conversation between the user and the system. The given conversation history is [...]. Let's think step by step. If there is nothing to respond to, please use the response "Alright!"

Table 20: Prompts for CoT.

Proactive Prompting

Your current task is to determine the user's intentions and satisfy their needs based on the provided conversation between the user and the system. The given conversation history is [...]. You may choose to either inquire about a missing detail in the current CAD design or, if the current CAD design is already complete, simply respond with "Alright!"

Table 21: Prompts for proactive prompting.

ProCoT

Your current task is to determine the user's intentions and satisfy their needs based on the provided conversation between the user and the system. The given conversation history is [...]. First, let's analyze step by step whether there are any missing details in the current CAD design. Then you may choose to either inquire about a missing detail in the current CAD design or, if the current CAD design is already complete, simply respond with "Alright!"

Table 22: Prompts for ProCoT.

PS+ Prompting

Your current task is to determine the user's intentions and satisfy their needs based on the provided conversation between the user and the system. The given conversation history is [...]. Let's first understand the user's intentions and devise a plan to satisfy their needs. Then, let's carry out the plan to satisfy their needs step by step.

Table 23: Prompts for PS+ prompting.