# Error Notebook-Guided, Training-Free Part Retrieval in 3D CAD Assemblies via Vision-Language Models

**Yunqing Liu, Nan Zhang & Zhiming Tan**[*]
Fujitsu R&D Center
Shanghai, China
{liuyunqing, zhangnan, zhmtan}@fujitsu.com

## Abstract

Effective specification-aware part retrieval within complex CAD assemblies is essential for automated design verification and downstream engineering tasks. However, directly using LLMs/VLMs to this task presents some challenges: the input sequences may exceed model token limits, and even after processing, performance remains unsatisfactory. Moreover, fine-tuning LLMs/VLMs requires significant computational resources, and for many high-performing general-use proprietary models (e.g., GPT or Gemini), fine-tuning access is not available. In this paper, we propose a novel part retrieval framework that requires no extra training, but using *Error Notebooks + RAG* for refined prompt engineering to help improve the existing general model's retrieval performance. The construction of *Error Notebooks* consists of two steps: (1) collecting historical erroneous CoTs and their incorrect answers, and (2) connecting these CoTs through reflective corrections until the correct solutions are obtained. As a result, the *Error Notebooks* serve as a repository of tasks along with their corrected CoTs and final answers. RAG is then employed to retrieve specification-relevant records from the *Error Notebooks* and incorporate them into the inference process. Another major contribution of our work is a human-in-the-loop CAD dataset, which is used to evaluate our method. In addition, the engineering value of our novel framework lies in its ability to effectively handle 3D models with lengthy, non-natural language metadata. Experiments with proprietary models, including GPT-4o and the Gemini series, show substantial gains, with GPT-4o (Omni) achieving up to a 23.4% absolute accuracy improvement on the human preference dataset. Moreover, ablation studies confirm that CoT reasoning provides benefits especially in challenging cases with higher part counts ($> 10$).

## 1 Introduction

Recent efforts demonstrate the promise of LLMs/VLMs in the engineering design and manufacturing domain. For instance, Alrashedy et al. (2025) applied LLMs to generate Computer-Aided Design (CAD) code from natural language descriptions, which can then be executed to render 3D objects. Such approaches show that general-purpose models can automate the CAD modeling process. Additionally, Wu et al. (2021) developed deep generative models to create 3D CAD structures directly (e.g., by modeling sequences of CAD operations), hinting at the potential of combining language and vision for CAD design tasks. Recent work has also shown that LLMs can assist in design ideation and automation, such as guiding parametric modeling, generating shape grammars, and integrating CAD workflows with natural language instructions ( Vardhan (2025); Li et al. (2025); Akhtar et al. (2025)). These studies further highlight the versatility of LLMs in supporting creative and engineering tasks within CAD environments. Despite this progress, a critical task remains challenging for LLMs/VLMs: specification-driven part retrieval within complex CAD assemblies. Each CAD assembly (often stored as a STEP file) can contain dozens of parts described by lengthy, non-natural language metadata. Retrieving specific parts that match a given design specification

---

[*]Corresponding author.

arXiv:2509.01350v1 [cs.AI] 1 Sep 2025

or relational description is essential for automated design verification and other downstream tasks, yet directly prompting LLMs or VLMs for this often yields poor results. A primary obstacle is the extreme sequence length of assembly data, which can exceed current model token limits. Even if the STEP data is processed, for example, into images, we found that off-the-shelf models still frequently misidentify parts because the task requires fine-grained reasoning about part relationships and attributes.

Fine-tuning a model on this task could improve performance, but it is sometimes impractical: many models (e.g., GPT or Gemini) are proprietary or lack fine-tuning access, and training a custom model would demand significant computational resources. **However, certain training techniques for LLMs and VLMs may serve as inspiration for enhancing the performance of methods that do not require training.** For example, in the mathematical domain, Pan et al. (2025) fine-tuned a model on a special dataset of erroneous reasoning chains paired with corrected solutions. This taught the model to reflect on and fix its own errors during generation. More broadly, research on reflection and self-correction in LLMs highlights several strategies that could inspire our training-free framework. One line of work leverages *external critics or verifier models* to provide feedback on intermediate reasoning steps, guiding the model away from incorrect trajectories ( An et al. (2023); Li et al. (2023); Tong et al. (2024); Shinn et al. (2023); Renze (2024)). Another line explores *intrinsic self-correction*, where models are fine-tuned on specially constructed datasets that pair erroneous reasoning trajectories with their corrections ( Weng et al. (2023); Yang et al. (2025); Zhang et al. (2024); Han et al. (2024); Yan et al. (2024)). To collect such data, prior studies often introduce errors by raising the decoding temperature or by sampling across multiple models, ensuring that the training set contains both flawed and corrected reasoning paths ( Xi et al. (2024)). These approaches enable models to revise their reasoning, and prevent error accumulation. Although our method does not involve weight updates, we draw inspiration from these techniques. In particular, the idea of coupling flawed reasoning with explicit reflection and correction motivates our *Error Notebook* design. Instead of using fine-tuning to encode these revision patterns into the model parameters, we operationalize them at inference time: by retrieving analogous past samples and their corrections, we provide the model with direct exemplars of reflection, thereby encouraging more reliable reasoning without any additional training cost.

In this work, we introduce a novel inference-phase strategy for vision-language part retrieval in 3D CAD assemblies. Rather than training or fine-tuning a new model, our approach enhances reasoning on-the-fly through retrospective error analysis and retrieval-augmented guidance. Central to our method is the *Error Notebook*, a mechanism that refines model reasoning at inference time by recording and organizing corrected reasoning trajectories. For each new assembly query, we retrieve analogous cases from the *Error Notebook* and provide them as few-shot exemplars to guide the model's chain-of-thought (CoT) using a retrieval-augmented generation (RAG) strategy. Our inference-time, error-driven augmentation leads to substantial performance gains on the part retrieval task. We evaluate several state-of-the-art VLMs (including GPT-4 variants and Google Gemini models) on our benchmark. In summary, our contributions are as follows:

(1) We propose a training-free reasoning framework that combines the *Error Notebook* and RAG for VLM inference. The *Error Notebook* is constructed by systematically collecting and correcting CoT errors. Our retrieval-augmented inference approach leverages the *Error Notebook* as a source of few-shot exemplars, significantly improving model performance without any fine-tuning.

(2) We reconstruct a new multimodal dataset of CAD assemblies with relational specifications, incorporating human annotations to capture human preferences.

(3) From the perspective of the engineering value, we design an effective two-stage VLM strategy that first generates part descriptions and then uses these descriptions for retrieval, thereby overcoming the challenge of processing extremely long STEP file inputs.

## 2 METHODOLOGY

### 2.1 DATASET CONSTRUCTION

Our study is based on the Fusion 360 Gallery Dataset (Willis et al., 2021b;a; Lambourne et al., 2021). Specifically, we utilize the *Assembly Dataset*, a subset of the Fusion 360 Gallery Dataset,
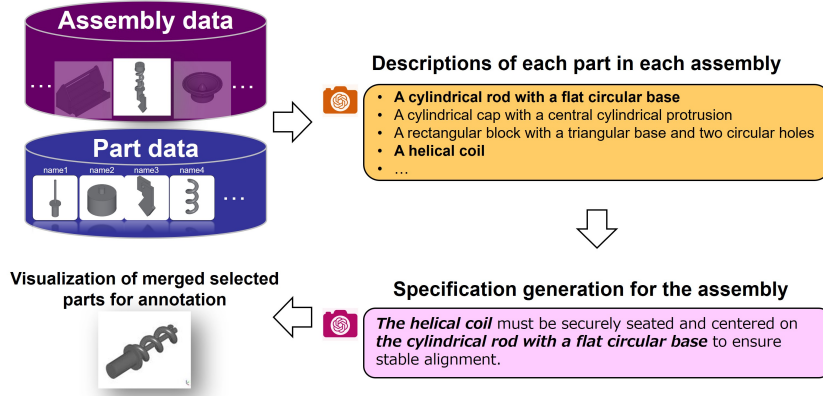
Figure 1: **Overview of the dataset construction pipeline.** For each assembly, a vision-language model is used to generate concise and discriminative natural language descriptions for every part. Subsequently, the model generates assembly-level specification sentences describing the required relationship or fit between selected parts. To support human annotation, the specified parts are merged and visualized as a single 3D model image.

which comprises multi-part CAD assemblies enriched with detailed information regarding joints, contact surfaces, and holes. For this work, we focus on the first archive (a1.0.0_00), which contains 752 assemblies. Each assembly project within this archive includes a single assembly and the corresponding part information (such as PNG images, STEP files, and additional metadata). The PNG files provide 2D image representations of the 3D models. STEP files (Standard for the Exchange of Product model data), as defined by ISO 10303, are neutral file formats that facilitate the exchange of 3D model data across different CAD software platforms, preserving geometry, structure, and other essential attributes. Figure 1 shows the overview of the dataset construction pipeline.

To begin, we catalog all part names and count the number of parts per assembly. Next, we utilize the GPT-4o (Omni) vision-language model to generate concise and descriptive noun phrases for each individual part. For each part, we provide both the overall assembly image and the part image as input, so that the model can generate the part description with full awareness of the assembly context. Each phrase is intended to succinctly describe the part's primary shape and distinguishing features, thereby allowing it to be differentiated from other parts within the same assembly. We provide several few-shot examples to guide the model toward generating higher-quality descriptions. Figure A.2 presents the prompt for this process.

Subsequently, we leverage the same model to further generate high-level specifications for the 3D assemblies. Each specification is focused on relationships between selected parts within the assembly. The process is as follows: First, the model reviews the assembly image and the corresponding list of part descriptions. It then selects two part descriptions that are most likely to exhibit a direct physical, spatial, or functional relationship (e.g., fit, mounting, alignment, or coupling). For each pair, the model generates a specification sentence that articulates the relationship, fit, or assembly condition between the two parts. The resulting set of filenames, $f_i$, is subsequently adopted as the ground truth for downstream part retrieval tasks. Figure A.3 presents the prompt for this process.

Finally, to facilitate the construction of a human preference database, which enables the evaluation of model performance from a human-centric perspective, we incorporate a human annotation stage. Specifically, we first employ Python's FreeCAD library to merge the selected parts into a single STEP file, and then use PythonOCC to render an image. Each annotation bundle includes the merged part image, the original assembly image, and the relevant specification sentence. Professional annotators review these materials and filter items according to the following procedure:

(1) Examine the assembly image to gain a comprehensive understanding of the overall structure.

(2) Items with overly similar part descriptions are discarded, as such cases can lead to ambiguity and multiple possible answers during part retrieval.

(3) Assemblies in which the overall structure is nearly indistinguishable from one or more of its constituent parts are also filtered out.

(4) Any other scenarios that may introduce ambiguity or permit multiple correct answers in part retrieval are excluded.

## 2.2 PART RETRIEVAL FRAMEWORK

Given a 3D assembly $\mathcal{A}$ consisting of $n$ parts $\{P_1, P_2, \ldots, P_n\}$, and a natural language assembly specification $S$, our goal is to retrieve the subset of parts $\mathcal{P}^* \subseteq \{P_1, \ldots, P_n\}$ that satisfy the specified relation described in $S$. The retrieval process is formulated as a two-stage VLM reasoning pipeline:

**Stage 1: Part Description Generation.** For each part $P_i$, we provide both the image of the complete assembly $\mathcal{I}_{\text{assembly}}$ and the image of the individual part $\mathcal{I}_{P_i}$ as input to a model $f_{\text{desc}}(\cdot)$. The model is prompted to generate a concise and discriminative noun phrase $d_i$ describing $P_i$ with explicit reference to the assembly context:

$$d_i = f_{\text{desc}}(\mathcal{I}_{\text{assembly}}, \mathcal{I}_{P_i}, \text{prompt}_{\text{desc}})$$

where $\text{prompt}_{\text{desc}}$ is a designed instruction that encourages the model to focus on salient geometric and functional features.

**Stage 2: Specification-Aware Part Retrieval via CoT Reasoning.** Given the assembly image $\mathcal{I}_{\text{assembly}}$, the mapping (JSON) from part filenames (IDs) to their descriptions $\mathcal{D} = \{filename_i : d_i\}_{i=1}^n$, and the specification $S$, we prompt the model $f_{\text{retr}}(\cdot)$ to identify the relevant parts:

$$\hat{\mathcal{P}}^* = f_{\text{retr}}(\mathcal{I}_{\text{assembly}}, \mathcal{D}, S, \text{prompt}_{\text{retr}})$$

where $\text{prompt}_{\text{retr}}$ requires the model to reason step-by-step (CoT) and produce both an interpretable rationale and the final answer in the form of a subset of part filenames.

## 2.3 ERROR NOTEBOOK CONSTRUCTION

To further improve model reasoning, we construct an *Error Notebook* that leverages the ability of VLMs to self-reflect and correct mistakes within their step-by-step reasoning process. Figure 2 shows this process.

Given, for each assembly, the assembly image $\mathcal{I}_{\text{assembly}}$, the mapping from part filenames to their descriptions $\mathcal{D}$, a specification $S$, the previous CoT reasoning $R^{\text{prev}}$, and the ground-truth filenames $\mathcal{P}^{*(\text{gt})}$. The goal is to generate a **corrected reasoning trajectory** $R^{\text{corr}}$ that leads to the correct solution in a human-like manner.

In theory, we formalize the step-by-step reasoning process as a trajectory $R = (s_1, s_2, \ldots, s_n, \hat{a})$, where $s_i$ are intermediate reasoning steps and $\hat{a}$ is the predicted answer. A suboptimal trajectory, $R^{\text{prev}}$, may contain both correct steps and erroneous steps, ultimately leading to an incorrect prediction. Models are expected to identify and revise the first erroneous step in $R^{\text{prev}}$. We define a *corrected reasoning trajectory* $R^{\text{corr}}$ as the concatenation of: 1) all steps up to the first error, 2) a natural language reflection that pinpoints and transitions from the error, and 3) the corrected reasoning steps that ultimately yield the ground-truth answer $\mathcal{P}^{*(\text{gt})}$. Formally, if $R^{\text{prev}} = (s_1^g, \ldots, s_k^g, s_1^b, \ldots, s_m^b, a^b)$, where $s_i^g$ are correct steps and $s_j^b$ are erroneous, we extract the subsequence ending at the first error, $R_{\text{sub}}^{\text{prev}} = (s_1^g, \ldots, s_k^g, s_1^b)$. The corrected trajectory is then constructed as:

$$R^{\text{corr}} = R_{\text{sub}}^{\text{prev}} \oplus \text{TR} \oplus R^g$$

where TR is a transition phrase, and $R^g$ is the correct trajectory from the correction point to the ground-truth answer $\mathcal{P}^{*(\text{gt})}$.

In our approach,

$$R^{\text{corr}} = f_{\text{corr}}\big(\mathcal{I}_{\text{assembly}}, \mathcal{D}, S, R^{\text{prev}}, \mathcal{P}^{*(\text{gt})}, \text{prompt}_{\text{corr}}\big).$$

The $\text{prompt}_{\text{corr}}$ instructs the model to:

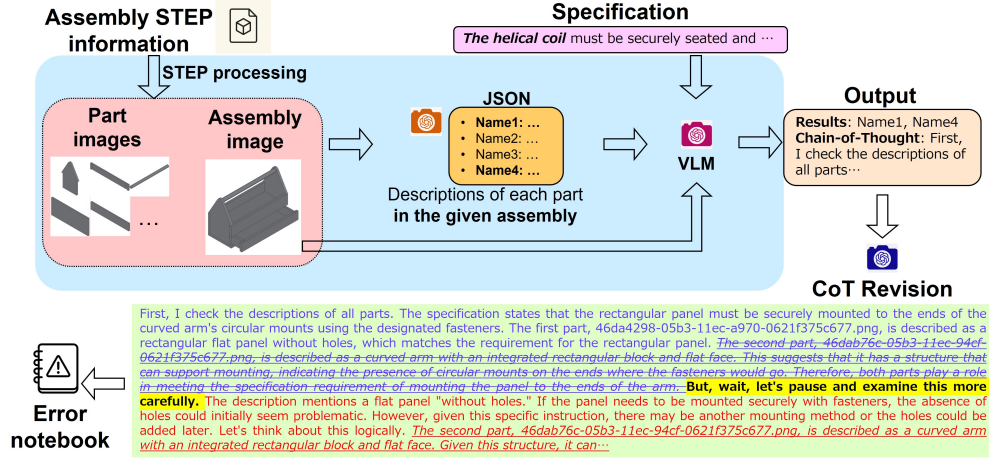Figure 2: **Error Notebook construction process. We define a corrected reasoning trajectory as the concatenation of: 1) all steps up to the first error, 2) a natural language reflection that pinpoints and transitions from the error, and 3) the corrected reasoning steps that ultimately yield the ground-truth answer.** (The overall pipeline processes to generate part-level descriptions using a VLM. Then the VLM infers relevant parts and produces an initial CoT to be revised.)

(1) Read and follow the previous reasoning $R^{\mathrm{prev}}$ step by step.

(2) Upon encountering the first logical or factual error, stop and explicitly articulate the transition, pointing out the mistake in a natural, self-reflective manner.

(3) From that point onward, independently correct the error, reasoning step by step until reaching $\mathcal{P}^{*(\mathrm{gt})}$.

(4) If no errors are detected, simply reproduce the previous correct reasoning and answer.

## 2.4 ERROR NOTEBOOK + RAG-BASED INFERENCE

In the inference stage, we adopt a RAG strategy that leverages examples from the *Error Notebook* as few-shot exemplars. Specifically, it retrieves the top-$n$ most relevant samples from the *Error Notebook* based on their similarity to the current assembly specification, using the corrected CoT trajectories from these entries to inform and guide the model's reasoning. Figure 3 shows the overview of the overall inference process based on VLMs.

Given an instance defined by the assembly image $\mathcal{I}_{\mathrm{assembly}}$, the mapping from part filenames to descriptions $\mathcal{D}$, and a specification $S$, the RAG-based inference proceeds as follows:

(1) **Sample Retrieval:** Let $\mathcal{E} = \{e_1, \ldots, e_M\}$ denote the set of entries in the *Error Notebook*, each comprising a specification $S_j$, part descriptions $\mathcal{D}_j$, and a corrected CoT trajectory $R_j^{\mathrm{corr}}$. For the current query, compute the similarity $\mathrm{sim}(S, S_j)$ between $S$ and each $S_j$ in $\mathcal{E}$. **To avoid data leakage, the current query instance $e_{\mathrm{cur}}$ is excluded from retrieval and will never appear among its own few-shot exemplars.** The top-$n$ most similar samples are selected:

$$\{e_{k_1}, \ldots, e_{k_n}\} = \arg \max_{e_j \in \mathcal{E} \setminus \{e_{\mathrm{cur}}\}} \mathrm{sim}(S, S_j)$$

where $e_{\mathrm{cur}}$ denotes the current query instance.

(2) **Few-Shot Prompt Construction:** For each retrieved sample $e_{k_i}$, construct a prompt block containing the assembly context, part descriptions, specification, the corrected CoT $R_{k_i}^{\mathrm{corr}}$, and the corresponding final answer. These prompt blocks are concatenated to serve as few-shot exemplars for the current query.

(3) **Main Query Prompt:** The final model input consists of (i) the few-shot exemplars constructed above and (ii) the current query context, which includes the assembly image $\mathcal{I}_{\mathrm{assembly}}$, part descrip-
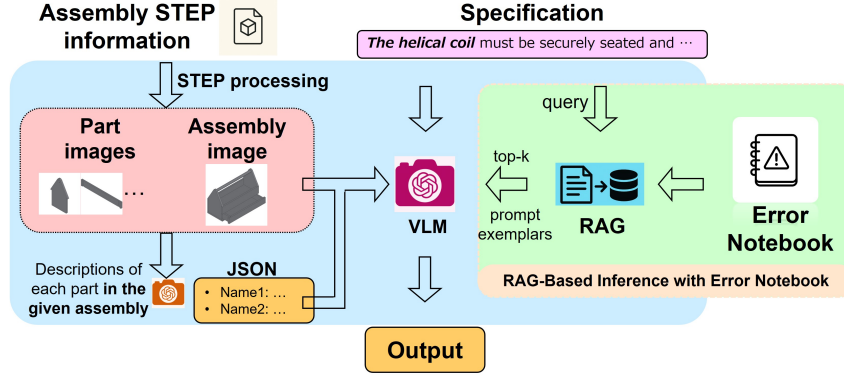
Figure 3: **Overview of the *Error Notebook* + RAG-based inference process.** Given the assembly specification, the system retrieves the most relevant examples from the Error Notebook according to the assembly specification, incorporates these as few-shot exemplars, and then performs step-by-step reasoning to generate the final answer.

tions $\mathcal{D}$, and specification $S$. The model is prompted to perform step-by-step reasoning, leveraging the retrieved exemplars as references.

Formally, let $F$ denote the few-shot prompt constructed from the top-$n$ retrieved entries. The model's output is given by:

$$R^{\text{RAG}} = f_{\text{rag}}\left(F, \mathcal{I}_{\text{assembly}}, \mathcal{D}, S, \text{prompt}_{\text{main}}\right)$$

where $R^{\text{RAG}}$ is the model's answer, and $\text{prompt}_{\text{main}}$ provides the instructions for the inference task.

## 3 EXPERIMENTS

### 3.1 IMPLEMENTATION DETAILS

Our pipeline interacts with VLMs (e.g., GPT-4o, Gemini) via API endpoints. For each inference call, images are encoded as base64 data URLs. We implement error handling with exponential backoff and up to 3 retries in the event of API errors. To process the dataset efficiently, all major computation steps are parallelized for asynchronously executing functions using multiple threads. Each assembly is processed as an independent unit. The generated part descriptions, which serve as intermediate outputs, are stored in JSON format. For fair comparison, both/all experiments on the same model/group employ identical description JSON files. Unless otherwise specified, the value of $k$ for RAG's top-$k$ retrieval is equal to the number of exemplars in the part retrieval stage, which defaults to 2.

### 3.2 MAIN RESULT

**(1) Our experimental results demonstrate that the proposed *Error Notebooks* with RAG framework enhances retrieval accuracy across all evaluated models and assembly complexities, as summarized in Table 1. The performance gains are particularly pronounced on the human preference dataset.** For example, GPT-4o (Omni) improves from 41.7% to 65.1% overall on the Human preference dataset, marking an absolute gain of 23.4%, while its performance on the self-generated dataset also rises from 28.5% to 48.3% (+19.8%). Similar trends are observed for other models: GPT-4o mini increases from 19.3% to 35.4% (+16.1%), Gemini 2.0 Flash Non-streaming from 44.2% to 56.8% (+12.6%), and Gemini 1.5 Pro Non-streaming from 43.0% to 46.7% (+3.7%). Another clear trend is that improvements are not limited to small assemblies: while the largest absolute gains often appear in cases with fewer parts (e.g., $< 10$ parts, GPT-4o Omni rises from 47.9% to 75.5%), consistent accuracy improvements are observed across all part-count intervals, including the more challenging $> 50$ parts group. These results highlight the effectiveness and generality of the proposed *Error Notebooks* + RAG strategy, which enhances inference across different proprietary (GPT, Gemini) models, without requiring additional training.

Table 1: Accuracy comparison of general models with and without Error Notebook–RAG integration on self-generated and human-preference CAD datasets. The best result is highlighted in **bold**. We divided the data from both datasets into 4 groups based on the number of parts in each assembly, reflecting the varying difficulty levels.

| Strategy | Self-generated dataset | | | | | Human preference dataset | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Overall | $< 10$ | $10 - 20$ | $20 - 50$ | $> 50$ | Overall | $< 10$ | $10 - 20$ | $20 - 50$ | $> 50$ |
| **GPT-4o (Omni)** | | | | | | | | | | |
| w/o E-Notebook | 28.5 | 40.7 | 22.4 | 15.3 | 5.0 | 41.7 | 47.9 | 32.4 | 26.5 | 0.0 |
| w/ E-Notebook | **48.3** | **66.8** | **35.9** | **29.7** | **16.3** | **65.1** | **75.5** | **42.6** | **41.2** | **21.4** |
| **GPT-4o mini** | | | | | | | | | | |
| w/o E-Notebook | 13.6 | 20.5 | 10.9 | 4.2 | 1.3 | 19.3 | 24.8 | 10.3 | 0.0 | 0.0 |
| w/ E-Notebook | **24.9** | **34.9** | **25.0** | **5.9** | **7.5** | **35.4** | **41.5** | **29.4** | **8.8** | **7.1** |
| **Gemini 2.5 Pro Non-streaming** | | | | | | | | | | |
| w/o E-Notebook | 36.5 | 55.1 | 25.6 | 14.4 | 6.2 | 54.0 | 65.2 | 35.3 | 20.6 | 0.0 |
| w/ E-Notebook | **42.2** | **60.9** | **30.8** | **21.2** | **11.3** | **59.5** | **69.5** | **42.6** | **29.4** | **14.3** |
| **Gemini 2.0 Flash Non-streaming** | | | | | | | | | | |
| w/o E-Notebook | 30.9 | 46.8 | 21.2 | 12.7 | 5.0 | 44.2 | 53.5 | 23.5 | 20.6 | 14.3 |
| w/ E-Notebook | **40.4** | **58.2** | **31.4** | **19.5** | **8.7** | **56.8** | **67.0** | **39.7** | **23.5** | **14.3** |
| **Gemini 1.5 Pro Non-streaming** | | | | | | | | | | |
| w/o E-Notebook | 29.9 | 44.3 | 21.2 | **13.6** | 6.2 | 43.0 | 52.1 | 23.5 | 17.6 | 14.3 |
| w/ E-Notebook | **32.4** | **49.3** | **22.4** | 11.9 | **6.2** | **46.7** | **57.1** | **25.0** | **17.6** | **14.3** |
| **Cloud Vision (Image) + Gemini 2.0 Flash Non-streaming** | | | | | | | | | | |
| w/o E-Notebook | 35.0 | 51.8 | 25.0 | 14.4 | 8.7 | 50.0 | 58.9 | 38.2 | 17.6 | 7.1 |
| w/ E-Notebook | **40.4** | **59.3** | **30.1** | **16.1** | **11.3** | **57.8** | **66.3** | **47.1** | **29.4** | **7.1** |

While Table 1 demonstrates the performance gap between models with and without Error Notebooks, Table 2 further shows that once Error Notebooks are incorporated, the **number** of exemplars retrieved by RAG has only a minor effect on final accuracy. For instance, on the self-generated dataset, the overall accuracy of the Non-CoT group varies only slightly between 49.4% (1 exemplar) and 52.7% (50 exemplars). A similar trend holds for the CoT group, where performance remains stable in the narrow range of 49.4% to 51.7%. Consistent patterns are observed on the human preference dataset. These results indicate that the key factor driving improvements is the presence of Error Notebooks themselves, and the effect of the specific number of exemplars sampled is negligible.

*Further discussion on Table 1.* Across the Gemini series, Gemini 2.5 Pro consistently achieves the best performance, regardless of whether the Error Notebook is used. On the self-generated dataset, GPT-4o (Omni) still holds an advantage, likely because this dataset was originally generated by GPT-4o itself. However, on the human-preference dataset, Gemini 2.5 Pro without the Error Notebook already surpasses GPT-4o (Omni) (54.0 vs 41.7). This suggests that Gemini 2.5 Pro aligns more closely with human preferences, highlighting the value of constructing the human-preference dataset in our evaluation.

**(2) The results in Figure 4 show that incorporating CoT reasoning from the *Error Notebook* is particularly valuable for challenging cases with higher part counts ($> 10$).** For assemblies with fewer parts ($< 10$), the Non-CoT group, where only final answers are given, often performs comparably or even slightly better, suggesting that in simple scenarios, direct access to the final correct solution is sufficient. By contrast, for complex assemblies with 10–50 parts, the CoT group consistently outperforms the Non-CoT group across nearly all exemplar sizes, confirming that step-by-step reasoning provides crucial guidance for harder queries. This trend is observed across all exemplar group sizes, with one notable exception: when using 50 exemplars, the CoT group shows a drop in accuracy. We attribute this to excessively long prompts caused by concatenating many CoTs, which may interfere with the model's judgment. A second important observation is that for
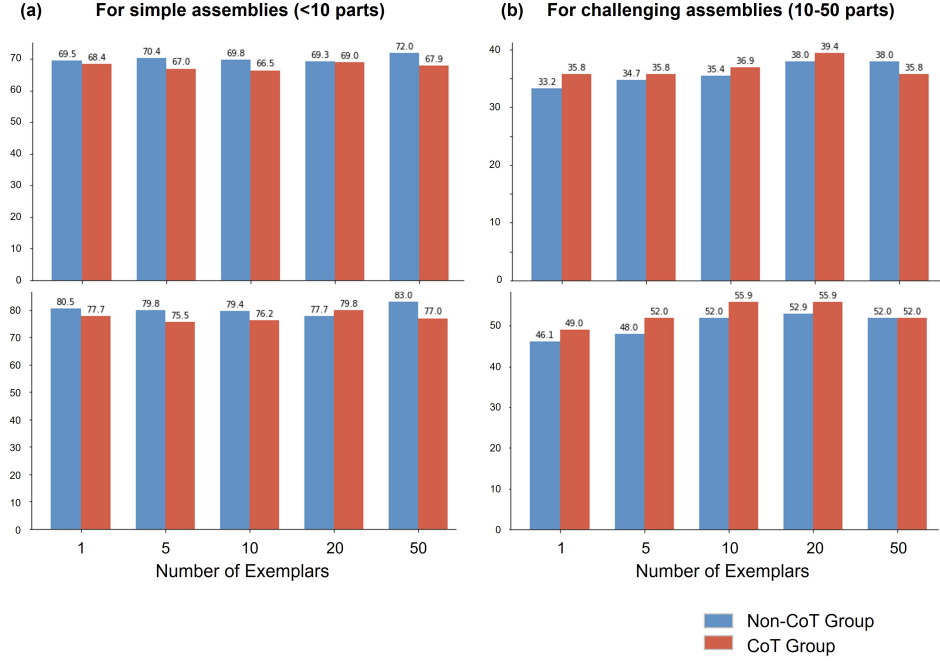
Figure 4: **Effect of CoT reasoning and exemplar number on retrieval accuracy across different assembly complexities and datasets.** Top row: results on the **self-generated dataset**; bottom row: results on the **human preference dataset**. (a) For simple assemblies ($< 10$ parts). (b) For more complex assemblies (10–50 parts). The $x$-axis indicates the **number of exemplars retrieved from the Error Notebook**, where each exemplar consists of either (i) the final corrected answer only (Non-CoT group) or (ii) the corrected CoT reasoning steps plus the final answer (CoT group).

Table 2: Ablation study on the number of exemplars retrieved from the *Error Notebook*. We also analyze the effect of excluding explicit CoT reasoning in each exemplar. *CoT Group* indicates that each retrieved exemplar includes explicit step-by-step reasoning, while *Non-CoT Group* omits such reasoning in the exemplars and includes ground truth only. The data from both datasets are divided into four groups based on the number of parts in each assembly, reflecting varying difficulty levels.

| Number of Exemplars | Self-generated dataset | | | | | Human preference dataset | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Overall | $< 10$ | $10 - 20$ | $20 - 50$ | $> 50$ | Overall | $< 10$ | $10 - 20$ | $20 - 50$ | $> 50$ |
| **Non-CoT Group** | | | | | | | | | | |
| 1 | 49.4 | 69.5 | 37.8 | 27.1 | 13.8 | 69.3 | 80.5 | 50.0 | 38.2 | 14.3 |
| 5 | 50.1 | 70.4 | 38.5 | 29.7 | 11.3 | 69.1 | 79.8 | 51.5 | 41.2 | 7.1 |
| 10 | 50.6 | 69.8 | 37.8 | 32.2 | 16.3 | 70.4 | 79.4 | 55.9 | 44.1 | 21.4 |
| 20 | 50.8 | 69.3 | 42.3 | 32.2 | 11.3 | 69.1 | 77.7 | 60.3 | 38.2 | 14.3 |
| 50 | 52.7 | 72.0 | 42.3 | 32.2 | 16.3 | 72.9 | 83.0 | 57.4 | 41.2 | 21.4 |
| **CoT Group** | | | | | | | | | | |
| 1 | 49.7 | 68.4 | 39.7 | 30.5 | 12.5 | 67.8 | 77.7 | 54.4 | 38.2 | 7.1 |
| 5 | 49.4 | 67.0 | 38.5 | 32.2 | 16.3 | 67.8 | 75.5 | 52.9 | 50.0 | 28.6 |
| 10 | 49.4 | 66.5 | 42.3 | 29.7 | 15.0 | 68.8 | 76.2 | 61.8 | 44.1 | 14.3 |
| 20 | 51.7 | 69.0 | 42.3 | 35.6 | 16.3 | 71.1 | 79.8 | 57.4 | 52.9 | 7.1 |
| 50 | 49.5 | 67.9 | 37.8 | 33.1 | 13.8 | 68.1 | 77.0 | 51.5 | 52.9 | 7.1 |

simple assemblies, increasing the number of exemplars has little effect, regardless of whether CoT is used. In contrast, for complex assemblies, accuracy steadily improves as the number of exemplars increases, up to around 20 exemplars.

**(3) We demonstrate the effectiveness of the proposed pipeline.** As shown in Figure 5, the proposed two-stage pipeline for part retrieval in 3D CAD assemblies achieves significantly higher accuracy compared to the image-only reasoning baseline. In the image-only setup, both the assembly image and individual part images are directly fed to the VLM in a single inference step, relying solely on visual input. In contrast, our proposed method first utilizes the VLM to generate concise part descriptions within the assembly context, and then performs part retrieval as a second reasoning step with the assistance of these textual descriptions. This design introduces an additional layer of interpretability and context-awareness, leading to consistent performance improvements across all part count groups. Quantitatively, the image-only baseline yields an overall accuracy of $15.0\%$ (107/715). The proposed pipeline achieves an overall accuracy



Figure 5: **Accuracy comparison between proposed pipeline and image-only reasoning.** Performance is shown for the proposed pipeline, which leverages part descriptions as intermediate references, versus the one that directly reasons over images.

of $33.6\%$ (240/715), with $51.2\%$ (185/361) for $< 10$ parts, $23.7\%$ (37/156) for 10–20 parts, $11.9\%$ (14/118) for 20–50 parts, and $5.0\%$ (4/80) for $> 50$ parts. These results demonstrate the effectiveness of incorporating part descriptions as intermediate representations.

## 4 CONCLUSION

In this work, we introduced a novel *Error Notebook*-guided, training-free part retrieval approach for complex 3D CAD assemblies. Our framework leverages retrospective error analysis and RAG to enhance VLM reasoning without additional training or fine-tuning. By systematically constructing *Error Notebooks* that capture and correct flawed reasoning trajectories, and by retrieving specification-similar exemplars at inference time, our method consistently improves accuracy across multiple proprietary VLMs, with gains of up to 23.4% absolute accuracy on human-preference benchmarks. Beyond retrieval accuracy, our study highlights two broader contributions: (i) a new human-in-the-loop CAD dataset that emphasizes specification-driven evaluation, and (ii) a two-stage VLM strategy that decomposes long assembly inputs into interpretable part descriptions, improving both scalability and interpretability.

Future work will explore open-source VLM integration, larger-scale datasets, and cross-domain applications of *Error Notebooks*, aiming to establish a more general paradigm for training-free reflective reasoning in multimodal AI.

## 5 THE USE OF LARGE LANGUAGE MODELS (LLMS)

We used LLMs as the experiment subject to study the improvement of our method on existing LLMs.

## 6 ETHICS STATEMENT

Our dataset construction process relies on professional human annotators, who were compensated fairly and provided clear annotation guidelines. Care was taken to exclude ambiguous or misleading cases to avoid introducing bias into the dataset. No personally identifiable information or sensitive data is involved. The proposed methods are intended for engineering and design applications, such as automated verification in CAD workflows, and do not pose foreseeable risks of misuse.

## 7 REPRODUCIBILITY STATEMENT

All code for dataset preprocessing, part description generation, Error Notebook construction, and inference experiments will be released . Detailed descriptions of dataset construction (including filtering and annotation protocols) are provided in Section 2.1. Experimental settings, including API interaction details, hyperparameters, and error-handling mechanisms, are documented in Section 3.1. Reproduction of our results only requires access to the Fusion 360 Gallery Assembly dataset and VLM APIs (e.g., GPT-4o, Gemini).

## REFERENCES

Naveed Akhtar et al. Large language models for computer-aided design: A survey. *arXiv preprint arXiv:2505.08137*, 2025.

Kamel Alrashedy, Pradyumna Tambwekar, Zulfiqar Zaidi, Megan Langwasser, Wei Xu, and Matthew Gombolay. Generating cad code with vision-language models for 3d designs. In *International conference on learning representations (ICLR)*, 2025.

Shengnan An, Zexiong Ma, Zeqi Lin, Nanning Zheng, Jian-Guang Lou, and Weizhu Chen. Learning from mistakes makes llm better reasoner, 2023.

Haixia Han, Jiaqing Liang, Jie Shi, Qianyu He, and Yanghua Xiao. Small language model can self-correct. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 18162–18170, 2024.

Joseph G. Lambourne, Karl D.D. Willis, Pradeep Kumar Jayaraman, Aditya Sanghi, Peter Meltzer, and Hooman Shayani. Brepnet: A topological message passing system for solid models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 12773–12782, June 2021.

Jiahao Li, Weijian Ma, Xueyang Li, Yunzhong Lou, Guichun Zhou, and Xiangdong Zhou. Cad-llama: Leveraging large language models for computer-aided design parametric 3d model generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 18563–18573, 2025.

Ming Li, Lichang Chen, Jiuhai Chen, Shwai He, and Tianyi Zhou. Reflection-tuning: Recycling data for better instruction-tuning. In *NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following*, 2023. URL https://openreview.net/forum?id=xaqoZZqkPU.

Zhuoshi Pan, Yu Li, Honglin Lin, Qizhi Pei, Zinan Tang, Wei Wu, Chenlin Ming, H. Vicky Zhao, Conghui He, and Lijun Wu. Lemma: Learning from errors for mathematical advancement in llms. *arXiv preprint arXiv:2503.17439*, 2025.

Matthew Renze. The effect of sampling temperature on problem solving in large language models. In *Findings of the Association for Computational Linguistics: EMNLP*, pp. 7346–7356, 2024.

Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning, 2023.

Yongqi Tong, Dawei Li, Sizhe Wang, Yujia Wang, Fei Teng, and Jingbo Shang. Can llms learn from previous mistakes? investigating llms' errors to boost reasoning. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 3065–3080, 2024.

Harsh Vardhan. Generative ai for cad automation: Leveraging large language models for 3d modelling. In *arXiv preprint arXiv:2508.00843*, 2025.

Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and Jun Zhao. Large language models are better reasoners with self-verification. In *Findings of the Association for Computational Linguistics: EMNLP*, pp. 2550–2575, 2023.

Karl D. D. Willis, Yewen Pu, Jieliang Luo, Hang Chu, Tao Du, Joseph G. Lambourne, Armando Solar-Lezama, and Wojciech Matusik. Fusion 360 gallery: A dataset and environment for programmatic cad construction from human design sequences. *ACM Transactions on Graphics (TOG)*, 40(4), 2021a.

Karl DD Willis, Pradeep Kumar Jayaraman, Hang Chu, Yunsheng Tian, Yifei Li, Daniele Grandi, Aditya Sanghi, Linh Tran, Joseph G Lambourne, Armando Solar-Lezama, and Wojciech Matusik. Joinable: Learning bottom-up assembly of parametric cad joints. *arXiv preprint arXiv:2111.12772*, 2021b.

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 (ICCV)*, pp. 6772–6782, October 2021.

Zhiheng Xi, Dingwen Yang, Jixuan Huang, Jiafu Tang, Guanyu Li, Yiwen Ding, Wei He, Boyang Hong, Shihan Do, Wenyu Zhan, et al. Enhancing llm reasoning via critique models with test-time and training-time supervision. *arXiv preprint arXiv:2411.16579*, 2024.

Yuchen Yan, Jin Jiang, Yang Liu, Yixin Cao, Xin Xu, Xunliang Cai, Jian Shao, et al. S3c-math: Spontaneous step-level self-correction makes large language models better mathematical reasoners. *arXiv preprint arXiv:2409.01524*, 2024.

Zhe Yang, Yichang Zhang, Yudong Wang, Ziyao Xu, Junyang Lin, and Zhifang Sui. Confidence v.s. critique: A decomposition of self-correction capability for llms. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (ACL)*, 2025.

Yunxiang Zhang, Muhammad Khalifa, Lajanugen Logeswaran, Jaekyeom Kim, Moontae Lee, Honglak Lee, and Lu Wang. Small language models need strong verifiers to self-correct reasoning. In *Findings of the Association for Computational Linguistics: ACL*, pp. 15637–15653, 2024.

# A APPENDIX

## A.1 ABBREVIATIONS

| | |
|---|---|
| VLM | Vision-Language Model |
| LLM | Large Language Model |
| CAD | Computer-Aided Design |
| STEP | Standard for the Exchange of Product Model Data (ISO 10303) |
| CoT | Chain-of-Thought |
| RAG | Retrieval-Augmented Generation |
| API | Application Programming Interface |
| GT | Ground Truth |
| GPT | Generative Pre-trained Transformer |

## A.2 FULL ENGINEERING PIPELINE ILLUSTRATION

To clarify the broader engineering context of our method and help better understand the meaning of *part retrieval* in practical CAD assembly analysis, we provide in Figure A.1 a complete overview of our proposed pipeline in an engineering setting. This illustration highlights the processing and visualization stages that do not require large model participation. Specifically, the left side depicts the *STEP processing* stage: an input assembly (in STEP format) is decomposed into its constituent parts using freecad, and subsequently rendered into 2D images using the pythonocc library. This generates intermediate representations (part-level STEP files and rendered images) that provide concrete references for the VLM-based retrieval process. On the right side, a textual specification is provided, and the VLMs enhanced with Error Notebook + RAG reasoning produce candidate part identifiers. These are then fused back into the assembly using freecad, and the resulting structure can be visualized with pythonocc.



Figure A.1: **Full engineering pipeline for specification-driven part retrieval.** The assembly STEP file is first decomposed into part-level STEP files using freecad, and both the assembly and part images are generated via pythonocc. Given a textual specification, VLMs enhanced with Error Notebook + RAG output candidate part identifiers, which are then fused back into the assembly with freecad for visualization.

## A.3 CASE STUDIES

Table A.1: Case studies of assembly-level part retrieval by GPT 4o (Omni) with Error Notebook. Each row shows the assembly image, the part count, the specification, and the retrieved results in image format.

| ID | Assembly Image | Part Count | Specification | Retrieval Results |
|---|---|---|---|---|
| 1 |  | 16 | The cylindrical protrusion on the vertical plate must align and securely fit into the curved channel of the rectangular housing. |  |
| 2 |  | 10 | The concave plate with a central circular hole on a short cylindrical base must be securely seated on the cylindrical base with radial grooves, ensuring proper alignment and fit. |  |
| 3 |  | 5 | The curved tapered arm with detailed thumb and fingers must fit snugly within the arm-shaped cavity of the curved block, ensuring full contact and proper alignment. |  |
| 4 |  | 10 | The semi-cylindrical block must fit securely onto the circular grid's central hub without obstructing the radial struts. |  |
| 5 |  | 8 | The cylindrical cap with integrated spout and loop handle must be securely screwed onto the threaded top collar of the cylindrical bottle body, ensuring a leak-proof seal. |  |
| 6 |  | 8 | The curved cylindrical shackle must be securely fitted into one of the lateral round holes on the cylindrical body. |  |

| ID | Assembly Image | Part Count | Specification | Retrieval Results |
|---|---|---|---|---|
| 7 |  | 10 | A flat rectangular plate with diagonal cutouts and rounded corners;A rectangular plate with a larger cut-out featuring a stylized raspberry design |  |
| 8 |  | 6 | The helical coil must be securely seated and centered on the cylindrical rod with a flat circular base to ensure stable alignment. |  |
| 9 |  | 5 | The threaded shaft of the knurled cylindrical knob must be securely fastened into the threaded hole of the curved lever arm to ensure proper functionality and alignment of the assembly. |  |
| 10 |  | 9 | The long, curved cylindrical tube must be snugly inserted into the perforated cylindrical opening of the elbow-shaped casing for a secure fit without gaps. |  |
| 11 |  | 4 | The cylindrical rod with a flat end must be fully inserted into the internal square socket of the cylindrical housing, ensuring secure attachment. |  |
| 12 |  | 4 | The hollow cylindrical cap must be securely fitted over the central circular protrusion of the curved base block, ensuring no gaps between the mating surfaces. |  |

A.4   SUPPLEMENTARY IMPLEMENTATION DETAILS

> You are an expert mechanical engineer. Given Image 1 (the assembly) and Image 2 (an individual part from the assembly), please generate a concise and descriptive noun phrase (not a full sentence). The phrase should briefly describe the part's main shape and any key features, in a way that clearly distinguishes it from the other parts in the assembly. Avoid generic names like 'part' or 'component'. Be specific about the shape and any holes, slots, or functional features. Your output should be a single noun phrase.
>
> ...................................................................................................
> For example:
> - A conical mount with a forked top;
> - A cylindrical pin;
> - Two plates with each having holes;
> - A flat round disk with three small holes;
> - A rectangular bracket with two mounting slots.

Figure A.2: **Prompt used to generate part-level descriptions in the dataset construction pipeline.**

> You are an expert mechanical engineer. Given an image of an assembled product (assembly) and a list of its part descriptions below:
> Part descriptions:
> {desc_list_str}
> ...................................................................................................
> Your task:
> 1. Review the assembly image and the list of part descriptions.
> 2. Choose any two part descriptions that are most likely to have a direct physical, spatial, or functional relationship in the assembly (such as fit, mounting, alignment, or coupling).
> 3. Generate one specification sentence (inspection/check item) that describes the required relationship, fit, or assembly condition between these two parts, as would appear in a manufacturing or assembly checklist.
> 4. Your specification should be clear, specific, and professional, mentioning both selected part descriptions explicitly.
> 5. Output only one specification sentence. Do not explain your reasoning.
> 6. Output format: The selected two part descriptions (exactly as shown above, separated by a semicolon), then a line break, then the specification sentence.
> ...................................................................................................
> For example, given descriptions like:
> 1. A cylindrical pin
> 2. A flat plate with holes
> Output:
> A cylindrical pin;A flat plate with holes
> The cylindrical pin must be fully inserted into one of the holes on the flat plate.

Figure A.3: **Prompt used to generate specification for each assembly in the dataset construction pipeline.**

You are an expert mechanical engineer with a sharp analytical mind. You are given the assembly image, the descriptions of all parts (each as 'filename: description'), the inspection specification, and a previous reasoning process (including its step-by-step thoughts and its Final Answer).

..........................................................................................

Your job:
1. Carefully read the previous reasoning step-by-step. Follow along and reproduce the steps until you encounter the first error or mistake.
2. Once you spot the first mistake, stop following the previous reasoning and use a natural transition phrase (such as: "But, wait, let's pause and examine this more carefully." or "Wait, something seems off. Let's pause and consider what we know so far.") to point out the error and correct it.
3. From that point on, continue the reasoning process in your own words, step-by-step, until you reach the correct answer (i.e., the filenames consistent with the correct ground-truth solution).
4. Do not mention "previous attempt" or "ground-truth solution" explicitly. Make your reasoning sound like a student discovering and correcting their own mistake in real time.
5. If the previous reasoning is already correct, simply reproduce the previous reasoning and the final answer as is.
6. End your output with a "Final Answer:" line followed by the filenames (from the keys above), separated by semicolons (;), with no extra words or punctuation.

Figure A.4: **Prompt used to revise CoTs.**

Now, for the following question, use the above reasoning as reference and answer step-by-step:
Assembly image:
[image attached]
Part descriptions:{desc_lines}
Specification:{spec}

..........................................................................................

Your task:
1. Think step by step (Chain-of-Thought) and explain how you identify the required part(s).
2. In the last line, write 'Final Answer:' followed by only the selected part filenames (from the keys above), separated by semicolons (;), with no extra words or punctuation.
Example output:
Chain-of-Thought:
First, I check the descriptions of all parts. Only part1.png and part2.png are described as cylindrical pins. Therefore, the required parts are part1.png and part2.png.
Final Answer:
part1.png;part2.png

Figure A.5: **Prompt used to generate the part retrieval results.**