

DesignQA: A Multimodal Benchmark for Evaluating Large Language Models' Understanding of Engineering Documentation

Anna C. Doris^{1,*}, Daniele Grandi², Ryan Tomich³, Md Ferdous Alam¹, Hyunmin Cheong⁴, Faez Ahmed¹

¹Massachusetts Institute of Technology, Cambridge, MA

²Autodesk Research, San Francisco, CA

³MIT Motorsports, Cambridge, MA

⁴Autodesk Research, Toronto, ON, Canada

arXiv:2404.07917v1 [cs.AI] 11 Apr 2024

ABSTRACT

This research introduces *DesignQA*, a novel benchmark aimed at evaluating the proficiency of multimodal large language models (MLLMs) in comprehending and applying engineering requirements in technical documentation. Developed with a focus on real-world engineering challenges, *DesignQA* uniquely combines multimodal data—including textual design requirements, CAD images, and engineering drawings—derived from the Formula SAE student competition. Different from many existing MLLM benchmarks, *DesignQA* contains document-grounded visual questions where the input image and input document come from different sources. The benchmark features automatic evaluation metrics and is divided into segments—Rule Comprehension, Rule Compliance, and Rule Extraction—based on tasks that engineers perform when designing according to requirements. We evaluate state-of-the-art models like GPT4 and LLaVA against the benchmark, and our study uncovers the existing gaps in MLLMs' abilities to interpret complex engineering documentation. Key findings suggest that while MLLMs demonstrate potential in navigating technical documents, substantial limitations exist, particularly in accurately extracting and applying detailed requirements to engineering designs. This benchmark sets a foundation for future advancements in AI-supported engineering design processes. *DesignQA* is publicly available at: https://github.com/anniedoris/design_qa/.

NOMENCLATURE

GPT4 Refers specifically to *gpt-4-1106-vision-preview*
FSAE Formula SAE
LLaVA Refers specifically to *llava-1.5-13b*
LLM Large Language Model
MLLM Multimodal Large Language Model
RAG Retrieval-Augmented Generation
VQA Visual Question Answer (Benchmark)

*Corresponding author: adoris@mit.edu

1. INTRODUCTION

Large language models (LLMs), such as ChatGPT [1], are chat-bots that can engage in conversations based on user queries. Trained on data from much of the internet, LLMs are based on the Transformer architecture [2] and have learned to predict the next words based on an input sequence of text [3]. ChatGPT is the fastest adopted technology in history, with more than 100 million users two months after its release [4]. LLMs have garnered significant attention for their conversational abilities, and research studies have examined and quantified their abilities to answer questions on a range of topics, from medicine [5, 6] to education [7] to engineering [8, 9].

With the emergence of LLMs as conversational assistants, an important question is how they can help humans answer questions about engineering design problems. A key goal of design automation has been to have an AI helper that can make it easier and faster for human designers to create better products. Although Generative AI has made significant strides, this goal has been difficult to attain, since engineering design tasks necessitate synthesis of multimodal information across multiple sources. One such task, critical to engineering design, is designing products based on technical requirements, which list rules that consist of a metric and a value (e.g. maximum tire width can be no greater than five inches) [10]. Matching the complexity of many real-world designs, technical specifications can be lengthy and extremely detailed and often reference critical safety or regulatory specifications. Designing according to requirements necessitates that engineers or designers can interpret and synthesize multimodal data across sources (e.g. the requirements document, CAD, engineering drawings, documentation, standards, etc.).

Recently, models with multimodal capabilities [11–13], lengthy documents (long-text) processing capabilities [14, 15], and both multimodal *and* long-text capabilities [1] have been developed. These advances bring us closer to the reality of a multimodal AI assistant that could help automate engineering design according to requirement documents. As new MLLMs emerge, evaluating their capacity to fulfill these essential design-

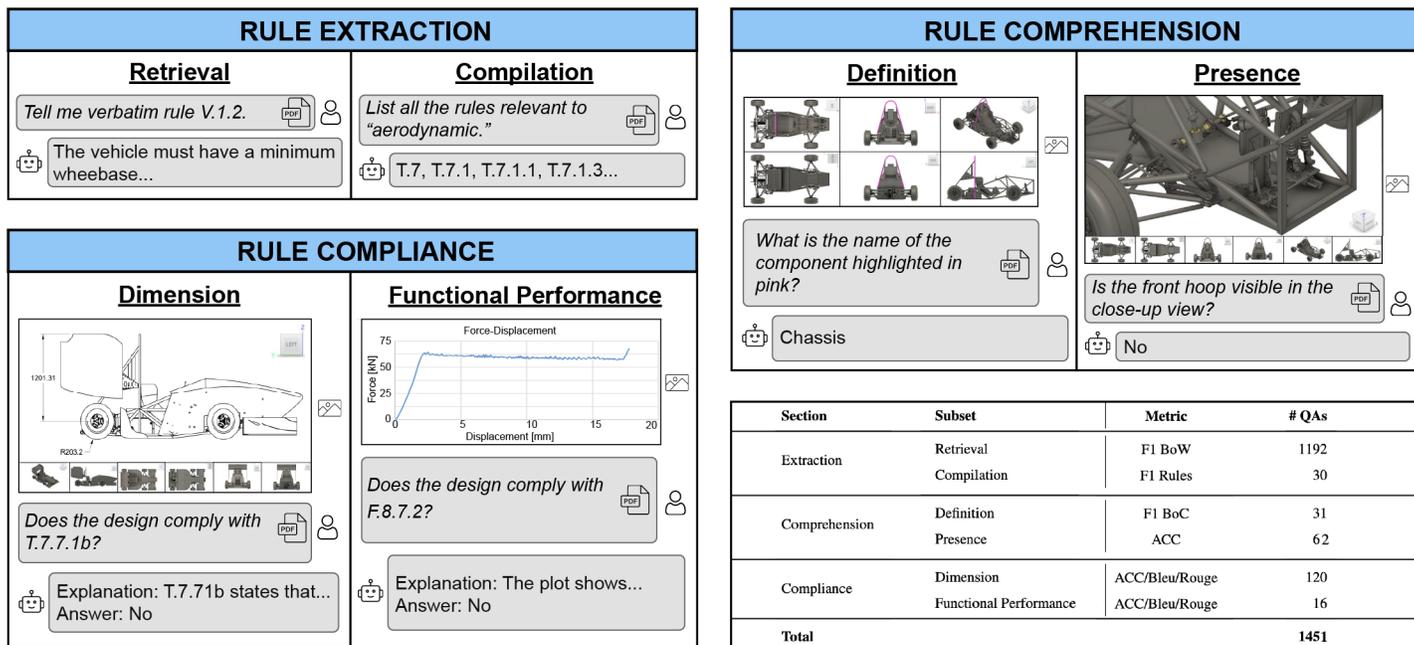


Figure 1: Overview of the three different segments (Rule Extraction, Rule Comprehension, and Rule Compliance) and six subsets in DesignQA. Prompts and images shown above are condensed versions of the actual prompts and images used. The bottom right table shows the metrics and the number of questions for each subset of the benchmark.

requirement-related tasks becomes imperative. This begs the question: How good are contemporary MLLMs at engineering design according to requirements? How can we measure tangible improvements in the efficacy of MLLMs at these tasks? Thus, we propose a novel benchmark aimed at assessing the proficiency of MLLMs in interpreting and adhering to the complex and multimodal demands of technical requirements in the design process.

We present DesignQA (Figure 1), the first zero-shot benchmark for technical requirements question-answering. The benchmark consists of 1451 questions and is based on the Formula SAE 2024 Rules and data (CAD, documentation, etc.) provided by the MIT Motorsports team. By developing this benchmark in conjunction with the MIT Motorsports team, we prioritized the generation of a dataset that is characteristic of real-world design requirement challenges. DesignQA also contains document-grounded reference-dependent visual question-answers (VQAs), one of a handful of benchmarks that tests models’ abilities to answer questions that require analysis across long-text documents *and* images. Notably, our benchmark assesses a model’s ability to synthesize information across an image and text from different sources, where the image was not seen by the model during its original training (pre-training).

In addition to developing the dataset, we used DesignQA to benchmark two state-of-the-art MLLMs, GPT4 and LLaVA, providing the FSAE rules to the models either via the context window or via a simple retrieval method. Of the models tested, we show that GPT4 (given the rules through its context window) performs the best on DesignQA. Based on observations of the performances of these models on the benchmark, we provide suggestions about how models might be modified for improved

results on DesignQA and design requirement questions generally. In summary, our contributions are:

1. **A Novel, Multifaceted Benchmark for MLLMs:** We introduce DesignQA, a benchmark that tests MLLMs’ understanding of design according to an engineering requirement document. DesignQA is unique in its need for models to analyze and integrate information from both visual and long-text inputs, emphasizing the complexity and multimodal nature of real-world engineering tasks.
2. **A Granular and Automatic Evaluation Framework:** We create DesignQA to be thorough and easy to use. The benchmark is divided into three segments - rule extraction, comprehension, and compliance - enabling a fine-grained investigation into a model’s strengths and weaknesses and enriching our understanding of AI in technical domains. Each subset of DesignQA has an automatic evaluation metric, permitting the quick evaluation of future MLLMs.
3. **High Quality, Real-World Question-Answer Pairs:** We develop a high-quality benchmark based on real-world data and problems. The question-answers in DesignQA are based on the FSAE competition rules and data provided by the MIT Motorsports team. Questions are designed and reviewed by members of the MIT Motorsports team, industry professionals, and engineering researchers.
4. **Evaluation of Contemporary MLLMs:** We evaluate MLLMs like GPT4 and LLaVA, unveiling the current limitations of AI and retrieval methods in handling multimodal

data and processing engineering requirements. The results of our evaluation both underscore the necessity for and expose possible avenues for improved models for engineering design.

2. RELATED WORK

In this section, we first provide an overview of existing work on AI for engineering design, showcasing that MLLMs have new potential to assist humans with design and design requirement problems. We then explore existing benchmarks for LLMs and MLLMs, highlighting the lack of benchmarks about engineering design and design requirements. We then categorize existing benchmarks by reference type. This sets the context for our contribution, which addresses the need for real-world, document-grounded benchmarks that bridge textual and visual information comprehensively.

2.1 AI for Engineering Design

Much of the prior work on AI for design has focused on single modalities [24], such as images or text. For text, several studies have investigated natural language processing (NLP) for technical engineering text. For example, [25] and [26] describe a technical language processing framework for circumnavigating the limitations of traditional NLP on unstructured engineering data. Expanding beyond text, [27] creates a deep learning architecture for technical document classification that factors in images as well as text. Despite significant advancements, many NLP and deep learning methods are specialized to a single domain and don't generalize well to other problems within engineering design.

LLMs offer more generalizable solutions to various problems within engineering design. [28] demonstrates the potential for LLMs (GPT-2 and GPT-3) to automate early-stage design concept generation. [8] explores how LLMs can assist engineers across an array of design and manufacturing tasks. While LLMs are useful for select engineering design tasks, many engineering tasks are highly multimodal (involving images, CAD, graphs, etc.). Therefore, recent advancements in MLLMs hold untapped potential for the automation of engineering design tasks. [9] investigates the potential of GPT-4 to automate engineering design tasks involving images, creating a dataset of over 1000 zero-shot queries. However, this dataset does not focus on engineering documentation. While a plethora of AI models like Google's Gemini, Meta's Llama family, and Anthropic's Claude models have emerged recently, their effectiveness is almost exclusively evaluated on non-engineering benchmarks. Critical for characterizing the abilities of MLLMs for engineering design tasks are benchmarks that can rigorously quantify their performances, which serves as inspiration for DesignQA.

2.2 LLM and MLLM Benchmarks

In this section, we explore existing benchmarks for LLMs and MLLMs based on domain and reference type. See Table 1 for a concise overview.

2.2.1 Benchmarks for Engineering Design and Design Requirements. Very few benchmarks exist for engineering design problems or design requirement-related tasks. Despite the

plethora of complex, multimodal QAs that could be generated from technically rich design requirement documents, very few datasets or benchmarks have been developed for this domain. PURE (Public REquirements) [29] is a dataset composed of 79 requirements documents scraped from the web. However, the dataset does not provide QA pairs and thus cannot be easily used for benchmarking. DesignQA harnesses the FSAE competition rules and MIT Motorsport design data to develop a benchmark of QAs pertaining to real-world design requirements.

2.2.2 LLM Reference-dependent Benchmarks. Text-based QA benchmarks that require a model to parse additional references to answer the posed question can be called "reference-dependent" benchmarks. Reference-dependent benchmarks differ from many classic reading comprehension benchmarks, like MCTest [16], which are "self-contained" and can be answered by short-text (approximately paragraph length) chunks accompanying the question. Since reference-dependent benchmarks require a model to locate the relevant information – usually across one or multiple long texts – and then apply that information to the posed question, they tend to be more complex questions. Following the distinction made by Dasigi *et al.* [19], reference-dependent benchmarks can further be divided into "open-domain" benchmarks and "document-grounded" benchmarks. An example of an open-domain question, taken from SQUAD, is: "Where do water droplets collide with ice crystals to form precipitation?" [17] Open-domain benchmarks, such as SQUAD [17] and WikiQA [18], test a model's ability to answer general-knowledge, factoid-type questions, the answer for which is usually contained in multiple sources in the model's pre-trained corpus.

In contrast, document-grounded benchmarks, like QASPER [19] and ZeroScrolls [14], test a model's ability to answer questions based on information provided in a specific long-text document. An example of a document-grounded question, taken from QASPER, is: "[In reference to a specific NLP paper] Which neural architecture do they use as a base for their attention conflict mechanisms?" [19] As noted by Dasigi *et al.*, document-grounded QAs tend to be more complex than open-domain QAs since they are anchored in user context and the answers are not widely available facts [19]. DesignQA, grounded in the FSAE rule document, fits within this document-grounded category. As a result, the questions posed in our benchmark are complex and rooted in user needs rather than common sense. Document-grounded questions are very characteristic of engineering design problems, as various types of documents – standards, manuals, documentation, etc. – often contain specific information that cannot be easily found on the internet.

2.2.3 MLLM Benchmarks. Multimodal datasets typically test an MLLM's capacity to analyze a non-text element with respect to question text. At the time of writing, most MLLMs can only accept images as non-text inputs, so multimodal QA benchmarks tend to consist of a visual (image) coupled with a question/answer pair. Visual question-answers (VQAs) can be categorized in the same way as text-based QA benchmarks. The vast majority of VQA benchmarks are self-contained. An example of a self-contained VQA question, taken from MME, is:

Table 1: Overview of select LLM and MLLM benchmarks, their domains, and reference-dependence. Our benchmark is unique in its focus on design requirements, and in that it contains multi-source document-grounded VQAs.

Benchmark	Domain	Reference Type				VQA
		Self-contained	Open-domain	Doc-Grounded Single Source*	Doc-Grounded Multi Source**	
McTest [16]	Narrative Children’s Stories	✓	-	-	-	✗
SQUAD [17]	Wikipedia	-	✓	-	-	✗
SQUAD [17]	Wikipedia	-	✓	-	-	✗
WikiQA [18]	Wikipedia	-	✓	-	-	✗
QASPER [19]	NLP Papers	-	-	✓	-	✗
ZeroScrolls [14]	Mixed: Wiki, Gov, etc.	-	-	✓	-	✗
MME [20]	COCO	✓	-	-	-	✓
MM-Bench [21]	Mixed: COCO, Llava, etc.	✓	-	-	-	✓
ScienceQA [22]	Open-source science materials	-	✓	-	-	✓
InfoSeek [23]	Wikipedia	-	-	✓	-	✓
DesignQA (Ours)	FSAE Rules Doc & Data	-	-	-	✓	✓

*Single source: image in question contained within the document; **Multi source: image in question not contained within the document

“[Pertaining to a photo showing doubles tennis partners] Are there two people in this image?” These are questions for which no context (other than the image) is needed to answer the question. The MME [20] and MMBench [21] benchmarks – highlighted by Chang and Wang et al.’s review paper [30] – are both self-contained VQAs. These benchmarks largely focus on basic tasks – primarily reasoning and perception – which are usually presented in a multiple choice format for ease of evaluation [20].

More challenging and less prevalent than self-contained VQAs, reference-dependent VQAs test a model’s ability to synthesize image analysis *with* additional knowledge, either from the open domain or from specific documents. ScienceQA [22], which contains elementary to high school-level science multiple choice questions, can be considered an open-domain VQA. The benchmark additionally encourages complete “train-of-thought” reasoning by providing a “lecture” – multiple sentences of general knowledge pertaining to the question – and “explanation” – reasoning for selecting a particular answer – for each VQA. InfoSeek [23] is the first document-grounded VQA, composed of questions about images in specific Wikipedia articles that can only be answered by consulting the corresponding article’s text.

However, there is still a significant need for document-grounded VQAs that are more characteristic of real-world tasks. InfoSeek has a direct match between visual and the document (i.e. the visual is contained within the document), while most visual questions asked by users would not fit this direct look-up framework (e.g. provided with an image of a broken machine and a manual for the machine, the exact image will not be contained within the manual). InfoSeek’s images and documents – which both come from Wikipedia – have also been seen by the model during pre-training; for many of the questions asked by users, either the document or the visual would not have been seen during pre-training. DesignQA is constructed to fill these gaps and is

more representative of these real-world scenarios.

3. DESIGNQA BENCHMARK

3.1 The Dataset

Our benchmark is developed in conjunction with the MIT Motorsports team, so all questions in the benchmark pertain to the 2024 Formula SAE Rules and are based on data provided by the team. The benchmark consists of 1451 question-answers (Figure 1) and is divided into three segments – Rule Extraction, Rule Comprehension, and Rule Compliance – each of which tests a different skill needed for designing according to a technical requirements document. Each segment of the dataset is further divided into two subsets, each of which corresponds with a specific task and has its own automatic evaluation metric. The QAs were created by a member of the MIT Motorsports team, a member of our team from industry (Autodesk), or a member of our team from academic research. All manually generated QAs – except those that are derivatives of other questions or Rule Compliance explanations – were reviewed by the two parties that didn’t write the question. To provide more context to the model, each of the QAs begins with the following preamble:

We are a student engineering team designing a vehicle for the FSAE competition. Attached is the FSAE rules document.

The rule document was provided to the model when asking each question. The following paragraphs describe in more detail how each segment of the benchmark was created.

3.1.1 Rule Extraction. A key - albeit usually simple - task for engineers is to locate a specific rule in a requirement document: answering a question about a rule is predicated on the ability to locate and extract the relevant rule. The Rule Extraction segment of the benchmark tests a model’s ability to extract

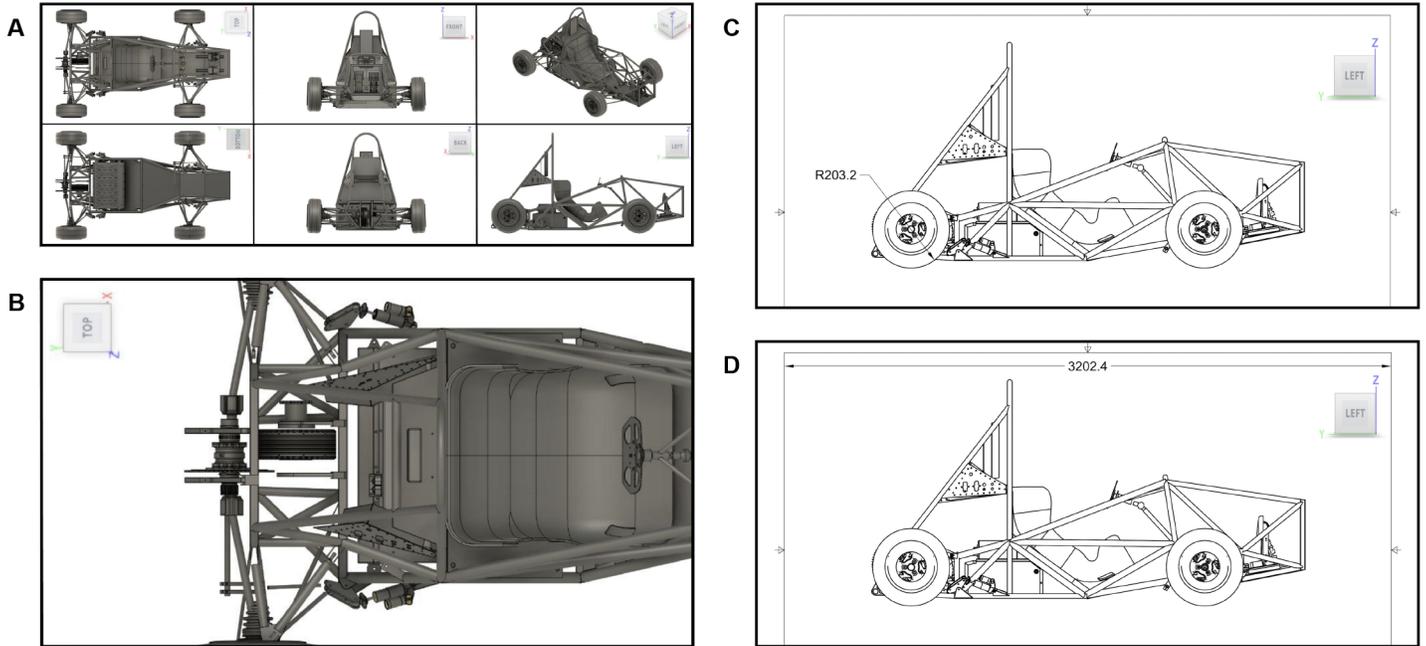


Figure 2: REPRESENTING 3D CAD MODELS IN 2D IMAGES. A) MULTI-VIEW CAD IMAGE. B) CLOSE-UP CAD IMAGE. C-D) ENGINEERING DRAWING IMAGES. C USES THE DIRECT DIMENSIONING METHOD AND D USES THE SCALE BAR DIMENSIONING METHOD.

rules from the 140-page FSAE rule document. This benchmark segment is further divided into two subsets: Retrieval QAs and Compilation QAs.

Retrieval questions test a model’s ability to extract specific information from a lengthy document. Given the large number of pages in the original rule document, retrieving the text of the rules word-for-word is a non-trivial task. While this retrieval task might become obsolete in the future as better models continue to be developed, it is critical for models to retrieve accurate information, as it is a necessary precursor for the other types of questions in this benchmark.

We programmatically create the Retrieval QAs by first extracting all the text from the PDF document, excluding the headers, footers, and page numbers. The rule document is well organized into numbered sections and subsections (which may or may not have titles) in the format ‘AA.##.##’. By using a combination of manually created scripts and regex patterns to identify the individual rules, we tabulate the set of rules and label the rule number, the rule title, and the rule text. Finally, we drop the rules that do not contain any text (while keeping the child rules) as well as sets of rules that pertain to other aspects of the race (e.g., Administrative Regulations, Document Requirements) and not to the design specifications that the vehicle must meet.

From the tabulated set of rules, we can then programmatically formulate the set of Retrieval QAs. To the preamble described in the previous section, we append the following:

What does rule $\{rule_number\}$ state exactly? Answer with only the text of the rule and no other words.

Where $\{rule_number\}$ is replaced by each of the selected rule numbers. This process results in 1192 Retrieval QAs.

Compilation questions assess a model’s ability to look for information spanning a long document. A common task that designers might perform when interacting with a design requirement document is the compilation of all rules relevant to a specific subject, such as the ‘suspension’ or ‘critical fasteners’. To create this set of QAs, we begin with a manually curated set of 30 common terms present in the rule document (nine of which include their synonyms, acronym, or plurals). This results in 30 questions with the following format:

Please list all rules relevant to $\{term\}$. Answer with only the rule numbers (i.e.: AA.1.1.1) separated by commas and no other words. The rules relevant to $\{term\}$ are:

Where $\{term\}$ is replaced by each of 30 common terms. To create the ground truth answers, we first compile the list of rules that include the term with a simple search through the tabulated rules, described in the previous section. We also include the children of each of the rules found, as well as any other rules that might be mentioned in both the parent and child rules.

3.1.2 CAD Representation. The Rule Comprehension and Rule Compliance segments of our benchmark ask questions about 3D CAD models of the designed vehicle. We develop QAs around

four different 3D CAD models provided by MIT Motorsports: the vehicle, the vehicle plus the aerodynamic package, the rear wheel package, and the powertrain. Before describing the details of these QAs, we devote this section to detailing how we provide 3D CAD model information to MLLMs. Since MLLMs cannot accept typical 3D CAD model formats (e.g. .stl, .step, etc.) at this time, we convert the CAD that we would like to show the model into 2D image forms, preserving as much 3D spatial information as possible. We represent 3D CAD models in 2D using three different kinds of images: 1) multi-view CAD images, 2) close-up CAD images, and 3) engineering drawing images (Figure 2).

Multi-view CAD images (Figure 2A) show six views of the CAD model: top, bottom, front, back, left, and isometric. Since the model cannot rotate a 3D CAD model in a CAD software GUI, these six views capture information about how the different views fit together to comprise the 3D model. Each view has a corresponding coordinate frame, so that it is clear to the viewer how each of the six views is related to the others.

Close-up CAD images (Figure 2B) show zoomed-in views of our CAD. The purpose of these images is to show finer detail in specific regions of the model. The close-up CAD images show a single view of the model with an orientation (and coordinate frame) matching one of the views in the corresponding multi-view CAD image.

Engineering drawing images (Figure 2C&D) display dimensional information about the 3D model. These images are created using engineering drawing software, so that the dimensions shown are highly accurate. They show a single view of the model with an orientation (and coordinate frame) matching one of the views in the corresponding multi-view CAD image. We used two different dimensioning systems to indicate the dimensions on these images. The first method was direct-dimensioning (as in Figure 2C), where dimensions relevant to a particular rule are explicitly indicated on the drawing. The second method was scale-bar-dimensioning (as in Figure 2D), where a scale bar is provided and from which a model could infer necessary dimensions. We used a mixture of these two dimensioning methods in our QAs, as we were interested in what effect the dimensioning method would have on model performance.

In the following sections, we describe how these three image types are employed in our QAs. Often, two of these image types are appended together to form an image that conveys more information. While we represent 3D CAD models using various 2D image types, this 3D model representation should be updated as MLLMs become more sophisticated and are able to parse inherent 3D model file formats.

3.1.3 Rule Comprehension. In order to understand how the rules relate to a design, engineers must first understand the terms presented in the rules and the names of the different components in the design. The Rule Comprehension segment of the benchmark evaluates a model’s ability to refer to elements of a 3D model according to the definitions and terminology presented in the rule document. This part of the benchmark is further divided into two subsets: Definition QAs and Presence QAs.

Definition questions test a model’s ability to identify the name of a highlighted component in a CAD model. From a list of 31 components, we created a multi-view CAD image where

the component-to-be-identified is highlighted in pink (Figure 1 and Figure 2A). Component synonyms were also collected (e.g. frame and chassis) for scoring purposes. Sometimes, it was necessary to hide some components in the CAD model so that the highlighted component could be better visualized. If components were hidden, it was noted in the prompt. Appended to the preamble is the following prompt, which resulted in the generation of 31 VQA pairs:

```
Also attached is an image showing six CAD views of our vehicle design. What is the name of the component(s) highlighted in pink?
{[If components hidden] Some parts of the design have been hidden so that the highlighted component(s) can better be visualized.}
Answer just with the name of the highlighted component(s) and nothing else.
```

We also tracked how the component-in-question was mentioned in the rule document: if it was mentioned explicitly in a “definition” section of the FSAE rule document (“definition component”), if it was not in a definition section but mentioned multiple times throughout the document (“multi-mention component”), or if it was not mentioned in the document at all (“no-mention component”). The intent behind this tracking was to understand whether the frequency and way in which a component’s name is mentioned in the rule document is correlated with the model’s ability to visually identify the component.

Presence questions assess a model’s ability to understand whether a particular component is present or not in a close-up CAD image. As such, these QAs are an easier variant of the Definition QAs, since they ask the model to provide a yes/no answer rather than the name of a component. Using the same list of 31 components from the Definition QAs, we generated two close-up CAD images (like that in Figure 2B), one which contained the component and another which did not. These 62 close-up CAD images were appended to the corresponding multi-view CAD image (like that in Figure 2A), which provided more 3D context for the close-up image. This resulted in 62 VQA pairs, each of which had the following prompt, where *{component_name}* was replaced with one of the 31 components:

```
Also attached is an image showing seven CAD views (each boxed in black) of our vehicle design. The top, big view shows a close-up view of the design. The six smaller views on the bottom of the image show different complete views of the CAD of the vehicle and are provided for context. Note that the close-up view orientation matches one of the six complete view orientations. The close-up view may also have some components hidden (with respect to the corresponding complete view) for visualization of specific components. Looking at the close-up view, is/are the {component_name} visible in the close-up view? Answer simply with yes or no.
```

3.1.4 Rule Compliance. Engineers frequently consult requirement documents to ensure that their designs comply with specific specifications. The Rule Compliance segment of the benchmark characterizes a model’s ability to check that a design conforms with a specific rule. This part of the benchmark is

further divided into two subsets: Dimension QAs and Functional Performance QAs, depending on the type of rule in question.

Dimension questions test a model’s ability to check that a design complies with a rule that stipulates dimensional constraints. From a list of 20 dimension rules, we generated three engineering drawing images for each rule:

1. A direct-dimensioned and rule-compliant image (as in Figure 2C).
2. A direct-dimensioned and rule-violating image. These were generated by editing the dimensions on the first image to explicitly violate the rule-in-question, or by modifying the CAD model so that the updated dimensions violated the rule.
3. A scale-bar-dimensioned and rule-compliant image (as in Figure 2D).

No scale-bar-dimensioned and rule-violating QAs were created. Since the CAD provided by MIT Motorsports is inherently rule-compliant, it is difficult to create negative examples when editing of direct-dimensions is not possible. Each of these engineering drawing images was appended to a corresponding multi-view CAD image (like that in Figure 2A) to provide context about the full model. This resulted in 60 VQAs with the following prompt, appended to the preamble:

Also attached is an image that shows an engineering drawing of the vehicle on the top accompanied by six CAD views of the vehicle on the bottom. The six CAD views each feature a different orientation of our design, so that 3D information about our design can be inferred. The CAD views are provided to contextualize the engineering drawing, which has the same orientation as one of the six CAD views. All units displayed in the engineering drawing have units of mm. Based on the engineering drawing, does our design comply with rule `{rule_number}` specified in the FSAE rule document?

{If direct-dimensioned:} Only use dimensions explicitly shown in the engineering drawing to answer the question. If a dimension is not explicitly shown, you can assume that it complies with the rules. }

{If scale-bar-dimensioned:} To answer the question, use the scale bar shown at the top of the engineering drawing to compute necessary dimensions in the drawing. }

First provide an explanation for your answer (begin it with ‘Explanation:’). Then provide just a yes/no answer (begin it with ‘Answer:’) that summarizes your response.

While these questions can be answered with a yes/no response, we also wanted to assess the model’s ability to explain *why* the design was or was not compliant, encouraging chain-of-thought reasoning [22]. For each direct-dimensioned question, we (or members of the MIT Motorsports team) wrote an explanation justifying the ground-truth yes/no answer. These explanations were not extensively reviewed, other than to ensure that they supported the corresponding ground truth yes/no answer. These human-written explanations can then be compared to generated model explanations.

From these 60 VQAs, we generated another set of 60 VQAs with additional context that would help with answering the ques-

tion. For the set of original 60 VQAs, we swapped out the multi-view CAD portion of the image with a different multi-view CAD image with components highlighted in pink that were relevant to the rule. We also added a line to the prompt explaining what the highlighted components were (e.g. “In the CAD views, the lower side impact structure is highlighted in pink”). In total, 120 Dimension VQAs were generated: 60 without additional context and 60 with additional context.

Functional Performance questions also test a model’s ability to check that a design complies with a rule given a relevant image. For this category however, either the rule, the image, or both is related to some functional performance of the design. Most questions involve a rule that imposes a constraint on some functional criteria of the design, and the image conveys the information required to check the rule. For example, there could be a restriction on the material choice for a part (hence the corresponding material strength) in the rule and the visualization of FEA results could indicate the maximum stress found in the part. When applicable, a pair of positive and negative examples is generated, where a variation is introduced to either the question or the image such that the first example would violate the rule while the second example would not. There is significant variation in the images and rules contained within this subset, and thus there is no standardized prompt. We encourage exploration of our code for more details. As these questions were more difficult to formulate due to limited availability of functional performance data, this subset has 16 VQAs. Like the Dimension questions, we generated a human-written explanation for each VQA.

3.2 Evaluation Metrics

Models tested on DesignQA can be evaluated completely automatically. For each of the six subsets of the benchmark, an appropriate automated evaluation metric was selected and implemented in our code so that a model’s (predicted) answer can readily be compared to the ground-truth answer. For several of the evaluation metrics, it can be difficult to intuit what score a prediction would receive relative to a ground-truth answer. As such, we’ve provided some examples in Figure 3 of model predictions (from our Model Evaluation section), their corresponding ground-truth answers, and the resulting scores. Each of the evaluation metrics is discussed in-depth below.

3.2.1 F1-score. Many subsets of our dataset are scored on a flavor of the F1-Score. The F1-Score is a popular metric for evaluating models’ performance on binary classification tasks, as the metric weighs both precision and recall. It is defined as: F1-Score: $F1\text{-Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$.

F1 Bag of Words: F1 applied to a Bag of Words (BoW) was used by [17–19] as an automatic metric for their benchmarks, which asked a model to pull verbatim phrases from a body of text to answer a question. As defined by [17], the metric first involves a cleaning step. The predicted answer (model response) and the ground-truth answer are converted to lower-case characters, extra white-space is removed, and punctuation and articles are taken out. They are then each “tokenized” into lists of words – predicted list (P) and ground-truth list (GT) – where F1-Score can be computed using Precision and Recall, where $Precision = \frac{P \cap GT}{\text{len}(P)}$ and $Recall = \frac{P \cap GT}{\text{len}(GT)}$.

Retrieval - F1 BoW			
	GPT4-AllRules	GPT4-RAG	LLaVA-RAG
Prediction	V.1.2 Wheelbase The vehicle must have a minimum wheelbase of 1525 mm	I'm sorry, but I cannot provide the text for rule S.3.7.1 as it is not included in the provided context from the Formula SAE rule document. If you can provide the specific text or section of the document where this rule is located, I would be able to help you further.	V.1.4 states that the vehicle must be open wheeled and open cockpit with four wheels that are not in a straight line.
	Score: 0.947	Score: 0.123	Score: 0.050
Ground Truth	Wheelbase The vehicle must have a minimum wheelbase of 1525 mm	A supplement to the Cost Report that reflects any changes or corrections made after the submission of the Cost Report may be submitted.	The distance to the ground below the Lower Side Impact Structure (F.6.4.5, F.7.5.1) at its lowest point should be 75 mm or less

Compilation - F1 Rules			
	GPT4-AllRules	GPT4-RAG	LLaVA-RAG
Prediction	T.3.2, T.3.3, T.4.3, IC.3.3, EV.4.6, EV.4.7, EV.7.2.2, EV.7.7.2, EV.9.3.2, EV.9.4.2, D.3.2.4, D.12.4.2, D.12.4.3, D.12.8.1, D.12.8.2, D.13.2.1, D.13.2.4.	EV.6.3, EV.6.3.1, EV.6.3.2, EV.6.3.3, EV.6.3.4, EV.6.3.5, EV.6.4, EV.6.4.1, EV.6.4.2, EV.6.4.3	EV.4.4, EV.4.4.1, EV.4.4.2, EV.4.4.3, EV.4.4.4, EV.4.4.5, EV.4.4.6, EV.4.4.7, EV.4.4.8
	Score: 0.242	Score: 0.343	Score: 0.0
Ground Truth	T.3.2, T.3.2.1, T.3.2.2, T.3.2.3, T.3.2.4, T.3.3.1, T.4.3.1, EV.4.7, EV.7.7.1, EV.7.7.4, EV.9.4.2, EV.4.7.1, EV.4.7.2, EV.9.3, EV.4.6, EV.4.5	F.11.1.3, T.1.8.4, T.4.2.9, T.4.3.3, T.8.3.2, IC.4.4.7, EV.4.1.3, EV.4.3.4, EV.6.3, EV.6.3.2, EV.6.3.3, EV.6.3.4, EV.6.3.5, EV.7.8.4, EV.6.3.1, EV.1.1, F.11.1.2, F.1.13, F.5.16, T.4.2.4, EV.7.8, EV.7.2.2, F.10.4, EV.6.1, EV.6.5	F.10.4.3, T.9.2.4, T.9.2.5, IC.9.3.2, EV.1.3, EV.6.6.3, EV.6.6.5, EV.7.7.2, EV.6.6.2, EV.7.2.2, EV.4.6

Definition - F1 BoC			
	GPT4-AllRules	GPT4-NoRules	LLaVA-NoRules
P	Roll hoop bracing	Suspension arms	steering
	Score: 0.609	Score: 0.455	Score: 0.235
GT	main hoop	uprights	roll hoop; roll hoops

Compliance Explanations - Bleu/Rouge			
	GPT4-AllRules	GPT4-RAG	LLaVA-RAG
Prediction	Rule T.7.5.a from the FSAE rule document states that in plan view, any part of any Aerodynamic Device must be no more than 700 mm forward of the fronts of the front tires. The engineering drawing provided shows a dimension of 705.28 mm from the front of the front tire to the most forward part of the nosecone, which is likely to be considered an Aerodynamic Device. Since this dimension exceeds the maximum allowed 700 mm, the design does not comply with rule T.7.5.a.	Unfortunately, the rule F.5.9.5 is not included in the provided excerpt from the FSAE rule document. Without the specific requirements of rule F.5.9.5, it is not possible to assess compliance based on the engineering drawing provided. To determine compliance, the exact wording and specifications of rule F.5.9.5 are necessary.	To determine if the design complies with rule F.3.2.1g, we need to compare the dimensions of the main hoop with the minimum size requirements given in the rule. In the engineering drawing, the main hoop is a steel tube with dimensions (height) x (outer diameter) = 300 mm x 35 mm. The rule F.3.2.1g states that the main hoop must be a single piece of uncut, continuous, closed section steel tubing meeting F.3.2.1.g. F.3.2.1.g does not provide any specific size requirements, but it refers to F.3.4.1. Since the main hoop in the image is a single piece of uncut, continuous, closed section steel tubing, we can conclude that it meets the requirements of F.3.2.1.g.
	Score: 0.320/0.468	Score: 0.070/0.215	Score: 0.097/0.268
Ground Truth	The rule states that no part of any aerodynamic device can be more than 700 mm forward of the fronts of the front tires. In the drawing, the front of the front wing is 705.28 mm (less than 700 mm) in front of the fronts of the front tires. Therefore, this design is rule compliant.	F.5.9 says: The included angle formed by the Main Hoop and the Main Hoop Braces must be 30° or more, and this image shows an angle of 28.76 degrees.	F.3.2.1g says the main hoop must meet size A requirements - 25mm outer diameter. The drawing shows the outer diameter of the main hoop is only 24.3mm

Figure 3: Responses from the different models evaluated across different subsets of the benchmark. We show the subsets that have evaluation metrics that can be harder to interpret, to provide references for various scores. The bolded portions of the predicted responses show what we interpreted to be correct.

Since our Retrieval QAs also ask the model to pull text from the rule document verbatim, we use this F1 Bag of Words metric to evaluate the Retrieval subset of the benchmark. We compute F1 BoW for each QA, and we report a macro-average across all questions.

F1 Rules: Similar to the Retrieval QAs, our Compilation questions ask a model to identify text (specifically rule numbers) contained within the rule document. We therefore use a very similar metric to the F1 Bag of Words used for the Retrieval QAs, except the lists P and GT are replaced by lists of rule numbers. We compute F1 Rules for each QA, and we report a macro-average

across all questions.

F1 Bag of Characters (BoC): Our Definition QAs ask a model to identify a component highlighted in a multi-view CAD image, using the rule document for reference. These QAs seemed like they should be scored similarly to the Retrieval QAs; however, since the model was now being asked for component names (several words) rather than complete rules (sentences), we did not want to penalize the model for small spelling errors or ending differences. For example, if the ground truth is “front hoop,” a predicted response of “front hooped” should be considered more correct than “front motor.” F1 Bag of Words would score “front

hooped" and "front motor" as equally correct. F1 Bag of Characters reflects the relative correctness of "front hooped" over "front motor." It is computed in the same way as F1 Bag of Words, except tokenization occurs on the character rather than word level. We compute the F1 BoC for each Definition QA across all synonyms, and we report the macro average of the highest score for each QA.

Accuracy: Several subsets of our benchmark – Presence, Dimension, and Functional Performance questions – ask the model to provide a yes/no answer. We score these using accuracy (ACC).

3.2.2 BLEU. The BLEU (Bilingual Evaluation Understudy) score was developed by Papineni et al. [31] for the automatic scoring of machine translations relative to human, reference translations. It has been employed in recent benchmarks [22, 32] and was used by [22] to evaluate the semantic similarity between an explanation provided by a generative model and an explanation written by a human. To compute BLEU, predicted and reference sentences are first broken up into n-grams, which are segments of n (a user-specified number) words. n-gram matches between the predicted and reference sentences are then found; once a predicted n-gram is matched to a reference n-gram, the n-gram is removed from the pool of reference n-grams in a process called "clipping." The number of matching n-grams is then divided by the number of n-grams in the predicted sentence, producing a modified precision score, p_n . The authors suggest to compute BLEU-4: p_n for $n = 1$ through $n = 4$, taking the geometric mean of the four p_n (since p_n decays exponentially with increasing n) [31]. However, since we only use one reference (one explanation), BLEU-4 scores were always near zero. As such, we report BLEU-2, which has non-zero scores but still preserves some of the information about adjacency of words. BLEU-2 (with max n-gram 2) can be computed as: $\log(\text{BLEU}) = \min(1 - \frac{r}{c}, 0) + \sum_{n=1}^2 \frac{1}{2} \log(p_n)$

The $\min()$ term is a "brevity penalty" and serves to penalize predictions (of length c) that are shorter than the ground truth (length r). For the explanation portions of the Rule Compliance questions, we report BLEU-2 to quantify the similarity between the model's generated explanation and the human-generated explanation.

3.2.3 ROUGE. The ROUGE (Recall-Oriented Understudy for Gisting Evaluation) metric was developed by Lin et al. [33] to measure the quality of a computer-generated text summary relative to a human, reference summary. The metric has been utilized in recent benchmarks [14, 22] to assess generative models' abilities to craft explanations and summarize text. ROGUE-L, which uses longest common sub-sequence (LCS), is useful for characterizing similarity in sentence-level word order. The LCS refers to the longest sequence of words that appear in the same order but not necessarily consecutively in the two sentences. ROGUE-L is computed as the F1 score of the LCS: $R_{LCS} = \frac{LCS(R,P)}{m}$, $P_{LCS} = \frac{LCS(R,P)}{n}$, and $ROUGE_L = \frac{2P_{LCS} \times R_{LCS}}{P_{LCS} + R_{LCS}}$,

where R is a reference sentence, P is a predicted sentence, m is the length of a reference sentence, n is the length of a predicted sentence, R_{LCS} is the recall, and P_{LCS} is the precision. For the explanation portions of the Rule Compliance questions, in addition to reporting BLEU-2, we also report ROGUE-L.

4. MODEL EVALUATION

We evaluate simple baselines and recent state-of-the-art MLLMs on our DesignQA benchmark to understand the state of current AI models in understanding engineering requirement documentation. The goal is to identify gaps in current AI's capabilities and encourage other researchers to build and train better AI models and frameworks for answering the questions we collected. Moreover, evaluating different MLLM models provides some insight into the relative difficulty of the questions, and the failure modes could serve as inspiration for better approaches to the benchmark.

4.1 Baselines and Models

Naive Baselines Similar to [14], we create basic baselines that rely on random selection so that it is easier to contextualize models' performances across the different subsets and metrics of the benchmark. For the Retrieval questions, we randomly choose a rule from the 1192 rules in our rule list. For the Compilation questions, we randomly pick 10 rules from our list of 1192. For the Definition questions, we randomly choose two consecutive words in the rule document. For the questions scored on accuracy (Presence, Dimension, and Functional Performance questions), we randomly select yes or no with 50% probability. Note that these baselines are evaluated to provide a sense of the lower threshold score (predicting randomly with no learning) for any machine learning model.

MLLM Models We consider two recent MLLMs in our evaluation: the closed-source model from OpenAI *gpt-4-1106-vision-preview* (GPT4) and the open-source model *llava-1.5-13b* (LLaVA) [1, 34]. GPT4 is chosen for evaluation because of its high performance on existing benchmarks [9, 14] while LLaVA is selected because of its promise (and the promise of its derivatives) as an open-source MLLM [35]. The extracted text from the FSAE rule document PDF is roughly 70,091 tokens in length. While GPT4 has a 128,000 token context window and can ingest the whole text in the prompt, LLaVA only has a 4,096 token context window and cannot. Thus, LLaVA requires the use of a Retrieval Augmented Generation (RAG) system to retrieve the appropriate context relevant to the question from the FSAE rule document. RAG is a natural language processing technique that enhances text generation by incorporating external knowledge, dynamically retrieving relevant information from a database or corpus to inform and improve the content being produced. It is especially helpful for models with small context windows, as it selectively pulls relevant information from an extensive document to assist in generating more informed and contextually accurate text outputs. While we need RAG to test the open-source LLaVA model, our goal is not to study different RAG techniques. To this end, we implement a very simple RAG system using LlamaIndex and using OpenAI's *text-embedding-3-large* to embed the information in the FSAE rule document [36]. We index the text from the rule document into 250-token chunks, with a 50-token overlap. From the embeddings, the cosine similarity between the question and each of the embedded chunks is then computed, and the top-15 (top-12 for Compliance QAs) most relevant pages are

Table 2: This table presents a detailed comparison of various MLLM models on our dataset, revealing GPT4’s superior performance of the models tested. The results underscore that despite the advancement of these models in understanding and applying complex technical specifications, there is significant room for improvement for AI models to understand engineering requirements.

Section	Evaluation Metric Subset (Metric)	Baseline		Model			
		Naive	GPT4-AllRules	GPT4-RAG	GPT4-NoRules	LLaVA-RAG	LLaVA-NoRules
Extraction	Retrieval (F1 BoW)	0.082	0.750	0.181	N/A	0.112	N/A
	Compilation (F1 rules)	0.137	0.298	0.362	N/A	0.281	N/A
Comprehension	Definition (F1 BoC)	0.358	0.470	N/A	0.420	N/A	0.393
	Presence (ACC)	0.5	0.629	0.532	N/A	0.484	N/A
Compliance	Dimension (ACC/Bleu/Rouge)	0.5/-/-	0.533/0.118/0.296	0.300/0.091/0.235	N/A	0.408/0.097/0.241	N/A
	Functional Performance (ACC/Bleu/Rouge)	0.5/-/-	0.563/0.167/0.342	0.563/0.121/0.306	N/A	0.536/0.163/0.321	N/A

then fed into the prompt as context, before posing the question in the benchmark.

While there are few images in the FSAE rule document, they were not considered for our evaluation, as they require further processing, and LLaVA was not trained on multiple image inputs. The few tables in the FSAE rule document will not receive any special treatment and will be fed into the models just as simple text from the PDF text-extracting script, together with the rest of the text on the page.

For all subsets of the benchmark, except the Definition questions, we evaluate the performance of GPT4 given the full rule document via its context window (GPT4-AllRules), GPT4 given the rule document via the LlamaIndex simple RAG (GPT4-RAG), and LLaVA given the rule document via the LlamaIndex simple RAG (LLaVA-RAG). While the Retrieval and Compilation questions in the dataset could be answered with text-only models, we pose the questions in this segment with a null image to the MLLMs. For the Definition questions, RAG always returns the same portion of the rule document since the questions across each QA are the same (i.e., “tell me the name of the highlighted component,” even though the images vary across QAs). Because of this, we do not test GPT4-RAG and LLaVA-RAG on the Definition questions, but we instead test GPT4 and LLaVA without any input of the rules (GPT4-NoRules, LLaVA-NoRules) in addition to GPT4-AllRules.

4.2 Results and Analysis

Table 2 shows all the baseline and model results. First, we discuss some overall findings and then we specifically delve into the results for each subset of the benchmark.

Overall Results For all subsets of the benchmark, models performed better than the naive baseline. While no single model performed the best across all subsets of the benchmark, GPT4-AllRules was the dominant model in all subsets except the Compilation questions, where GPT4-RAG was the best performer. This finding reflects that the simple LlamaIndex RAG model used was not very effective at providing relevant rule information to the model, while inputting the full rule text into the context window allowed GPT4-AllRules to access necessary information. For the Extraction and Comprehension segments of the benchmark, GPT4 (-RAG or -NoRules) performed better than the corresponding LLaVA model. The reverse was true for the

Compliance Segment: LLaVA-RAG was equivalent to or outperformed GPT4-RAG.

Rule Extraction: Retrieval Of the models tried, GPT4-AllRules does the best job at retrieving the requested rule verbatim (0.750 average BoW score). When the model’s answer diverged from the ground truth answer, we noticed that it was often because the model was reporting a nearby rule rather than the rule requested (e.g. V.3.2.5 instead of V.3.2.4). Sometimes GPT4-AllRule’s answer was different from the ground truth because it included all child rules in addition to the requested parent rule (e.g when asked for V.1, V.1 was reported along with V.1.1 and V.1.2). While not technically wrong, this result was not handled differently by our evaluation metric and may be a result we would want to handle specifically in the future.

Providing the rules to GPT4 via RAG instead of via context resulted in a much lower average BoW score (0.181 for GPT4-RAG). For a number of questions (like the GPT4-RAG Retrieval example shown in Figure 3) the model replied that it could not produce the text for the rule because it was not included in the rule text given to it, indicating that the simple LlamaIndex RAG failed to provide the relevant portion of the rule document to the model. When GPT4-RAG did provide an answer to the question, it had similar responses to that of GPT4-AllRules. LLaVA-RAG received the same portion of the rule document as GPT4-RAG, but its average BoW score was even lower (0.112). Unlike GPT4-RAG, LLaVA-RAG would hallucinate rules rather than indicate that the requested rule was not contained within the portion of the document it received via RAG. Furthermore, instead of returning rules verbatim, LLaVA-RAG would frequently offer an interpretation of the requested rule (as in the LLaVA-RAG Retrieval example in Figure 3).

Rule Extraction: Compilation GPT4-RAG performs the best at the Compilation questions (0.362 F1 rules). We noticed that GPT4-RAG returns many fewer rules per question on average (9.57 rules) than GPT4-AllRules (17.53 rules), as seen in Figure 3. This could explain the improved performance of GPT4-RAG relative to GPT4-AllRules: when given the full 140-page rule document, the model tends to over-predict relevant rules, resulting in low precision and reducing the average F1 score. When given the top-15 most relevant pages via RAG, the model has fewer rules that it can draw from and therefore does not include as many false positives in its responses. LLaVA-RAG has the lowest score of

the models tested (0.281). For five out of the 30 questions, its predicted rule list has a score of zero (no overlap with the ground truth rule list). In contrast, GPT4-RAG has a score of zero for two out of the 30 questions.

Table 3: Performance analysis of MLLM models on definition QA, highlighting the impact of component mention types within rule documents. GPT4-AllRules leads in direct definition mentions, while all models show varied effectiveness across multi-mention and no-mention scenarios, underlining the challenges in extracting definitions without explicit mentions.

F1 BoC	GPT4 -AllRules	GPT4 -NoRules	LLaVA -NoRules
Definition	0.74	0.56	0.50
Multi-mention	0.42	0.39	0.37
No-mention	0.35	0.35	0.37

Rule Comprehension: Definition As explained in Section 4.1, -NoRules versions of GPT4 and LLaVA were used for this subset of the benchmark rather than -RAG versions, since the identical question text across QAs resulted in RAG always incorrectly returning the same section of the rule text. Of the models tested, GPT4-AllRules has the highest score (0.470 average F1 BoC). However, GPT4-NoRules performs better than we expected (0.420 average F1 BoC), indicating that identification of some components in an FSAE vehicle can be performed with knowledge from the web. As explained in Section 3.1.3, we tracked how the component-in-question in these QAs was mentioned in the rule document (definition-component, multi-mention component, or no-mention component). As to be expected, the drop in F1 BoC score between GPT4-AllRules and GPT4-NoRules is due mostly to a drop in score for the definition-component QAs (Table 3), which without the definitions spelled-out in the rule document, become harder to answer.

Rule Comprehension: Presence GPT4-AllRules performed the best (0.629 accuracy) on the Presence questions of the models tested. GPT4-RAG has a lower average accuracy (0.532), likely because the model has access to limited (15 or fewer) rule document chunks that reference the component in question. Like the Dimension questions, LLaVA-RAG performs marginally worse than GPT4-RAG.

Rule Compliance: Dimension GPT4-AllRules is the best performing model on the Dimension questions, although its average accuracy (0.533) is marginally better than the naive baseline (random yes-no guessing). We noticed several recurring issues that caused the model to not answer questions correctly. First, the model sometimes struggled to reference the relevant rule (as seen in the Retrieval questions), and this meant that the model would not be set-up to answer the question correctly. For example, when asked three different questions pertaining to Rule F.5.7.5, GPT4-AllRules quotes the rule differently in each of its three explanations: one time it quotes the correct rule, another time it

mistakenly quotes the text of the preceding rule (F.5.7.4), and the third time it mistakenly quotes the text of a rule two rules ahead (F.5.7.7). Second, we notice that GPT4-AllRules sometimes has difficulty in extracting dimensions from the engineering drawing. For example, in two scale bar questions each which has a scale bar of 3202.4 mm, the model says that the scale bar is 1000 mm in one instance and 500 mm in another instance. Third, we note that the model faces difficulty in answering questions where the drawing contains two dimensions, and the correct answer can only be obtained by adding or subtracting the two together. For example, in two questions about the same rule, the model scrapes one of two dimensions from the provided engineering drawing and uses it to answer the compliance question, when the difference between the two dimensions should actually be used.

The drop in accuracy score between GPT4-AllRules and GPT4-RAG can largely be attributed to the fact that the RAG does not provide the model with the relevant portion of the rule document (as seen in the Retrieval questions). In many cases when this happens, GPT4-RAG refuses to provide a yes/no answer (as in the example in Figure 3), resulting in an accuracy score of zero for that question. LLaVA-RAG receives the same portions of the rules as GPT4-RAG, but it does not indicate when the portion of the rules doesn't contain the rule-in-question. As a result, LLaVA5-RAG provides yes/no answers to most questions (51/60) while GPT4-RAG provides yes/no answers to many fewer questions (21/60). Because providing a yes/no answer improves the chances of guessing correctly (while refusing to answer automatically results in a score of zero), LLaVA-RAG receives an inflated average accuracy score on the dimension questions compared to GPT4-RAG. We also observe that many of the explanations that the models provide are significantly longer than the reference explanations we have (see Figure 3). It is therefore difficult to capture the correctness of explanations through the Bleu and Rogue scores. The interpretability of these scores would be improved by obtaining more human reference explanations, so as to better capture the distribution of possible explanations that could be considered correct.

As explained in Section 3.1.2, one-third of the questions comprising the Dimension QAs used scale-bars in the engineering drawings, while the other two-thirds used direct dimensions. We were curious to see what impact these two different dimensioning systems had on model performance. As seen in Table 4, all models perform better on the direct-dimensioned drawings than on the scale-bar drawings. The models seemed to struggle with the scale-bar, sometimes indicating in their explanations that they were using it for "rough" or "estimated" dimensions rather than precise ones and in some cases explaining that their "image capabilities...do not include measuring dimensions." As explained in Section 3.1.4, half of the Dimension questions were given additional context that we believed would help in answering the question while the other half were not. Surprisingly, we did not see any obvious trends in model performance with versus without the additional context. GPT4-AllRules and LLaVA-RAG performed worse with additional context, while GPT4-RAG performed better. More investigation into the effect of additional context is needed.

Table 4: Effect of dimensioning system used in engineering drawings on model accuracy on Dimension QAs. These results highlight the fact that the MLLMs tested are better able to answer questions about engineering drawings that are direct-dimensioned rather than those that have a scale-bar.

ACC by Dimension System	GPT4-AllRules	GPT4-RAG	LLaVA-RAG
Direct	0.663	0.45	0.41
Multi-mention	0.275	0	0.400

Rule Compliance: Functional Performance All of the MLLMs tested received the same accuracy score on the Functional Performance questions, although they did not provide identical yes/no responses for each question. Similar to the Dimension questions, we notice that the models sometimes cite the wrong rule, setting up the rest of the question for failure. We also notice some instances of not extracting information correctly from the provided images (e.g. not identifying that an image of FEA has a numerical scale bar or not correctly identifying the peak force in a force-displacement graph). We also observe some instances where a model makes additional assumptions about the question that are not incorrect, but encourages it to choose the wrong yes/no answer. For example, in one question, GPT4-AllRules is presented with a graph showing gearbox temperature over time, which approaches but does not exceed 60 °C. The rule in question specifies that the maximum temperature should not exceed 60 °C; while data provided does not, the model explains that it could if tested for more time, so it is not rule compliant. This is a valid answer, but not what we intended. Additional context in the prompt could help improve the clarity of these questions. Moreover, this set of questions was limited to just 16 due to limited resources and the time and expertise required to obtain and write these questions. Future work will focus on expanding this subset to better characterize model capabilities on more functional performance-type questions.

Recommendations for Models for This Benchmark In the results presented thus far, we have demonstrated how two state-of-the-art models, with and without simple RAG, perform against our benchmark. Here, we present some observations on how one might modify these models to achieve improved performance on our benchmark.

First, experimentation with the RAG model would be worthwhile so that relevant portions of the rule document are better retrieved. While we generally saw better performance when the full rule document was included in the model’s context window instead of provided to the model via RAG (GPT4-AllRules vs. GPT4-RAG), including the full rule document in the context of every question is very expensive (around \$800 for all questions in the dataset, at the time of writing, when using the latest vision model through the OpenAI API). As such, effective RAG could help reduce the computational burden. As seen in the results presented for the Retrieval questions, the simple LlamaIndex RAG used in this paper often does not provide the relevant section of the rule document to the model. Since the benchmark contains

VQAs, we expect that models coupled with multimodal RAG – where the QA image could help in selecting the relevant portion of the rule document – could help improve scores on the benchmark. The RAG could also perhaps be improved by experimenting with different chunking methods. Nevertheless, as models become cheaper, more computationally efficient, and trained with larger context windows, RAG approaches might become less useful.

Second, improving the quality of the instructions provided to the MLLM models, by using techniques such as prompt-engineering, might result in improved performance on the benchmark. For example, extracting specific portions of the prompt to be fed to the RAG would likely improve the scores of models with RAG. Few-shot learning or including chain of thought reasoning [22] may also improve model scores. Third, we suspect that fine-tuned models would exhibit improved performance on the benchmark. Since many tasks in the benchmark require that the model can first extract a relevant rule and then answer a question, fine-tuning a model on the Retrieval QAs may result in improved performance across the entire benchmark. This capability could also be obtained by hybrid models that combine rule-based search and deep learning. Lastly, while we presented the results of two base models in this work, many other closed and open-source models could be evaluated on our benchmark and might perform better.

5. FUTURE WORK AND LIMITATIONS

Engineering design, inherently grounded in practical applications, necessitates a deep understanding of technical documents to ensure that designs not only meet but also adhere to stringent requirements and standards. One of the critical insights emerging from our benchmark evaluations is the realization that modern AI models are still in the nascent stages of truly understanding engineering documentation. This limitation highlights a significant gap in the path toward fully automated design processes. As highlighted in the prior section, there is therefore substantial future work to be done to improve MLLM models’ abilities to solve engineering design and design requirement problems.

While this benchmark presents a first step into formally evaluating a model’s ability to interpret and understand design requirement documentation, the size of our current benchmark is limited to just one requirement document, and the creation costs (requiring expert manual validation) pose limitations on the generalizability and scalability of the dataset. Moreover, the metrics used to evaluate models against the benchmark could be further improved, as Bleu and Rouge have known limitations and qualitatively do not perfectly capture the correctness of the generated answers.

6. CONCLUSION

This study introduces DesignQA, a novel MLLM benchmark with 1451 questions and answers based on data from MIT Motorsports and the FSAE competition rule document. The benchmark is designed to assess large language models’ abilities to answer questions about design according to technical requirements. The benchmark is divided into six subsets - Retrieval, Compilation, Definition, Presence, Dimension, and Functional Performance - representative of tasks performed by engineers

when designing according to technical specifications. Each subset of the benchmark has its own automatic evaluation metric, so that new MLLMs can be seamlessly tested. The questions in the benchmark are designed by humans - MIT Motorsports members, industry professionals, and researchers - to ensure a high-quality benchmark. Different from many existing MLLM benchmarks, DesignQA contains document-grounded VQAs, in which the input image and document come from differing sources, characteristic of many real-world scenarios.

Using DesignQA, we conducted a rigorous evaluation on variants of two state-of-the-art MLLMs: GPT4 and LLaVA. We uncovered significant limitations in their current abilities to accurately interpret complex technical documents, specifically in referencing relevant requirements and extracting numerical information from technical images. Our research highlights the need for advancements in AI models to enhance their comprehension of engineering requirements and documentation, suggesting directions for future efforts. Our work aims to bridge the gap in AI's capability to support engineering design processes more effectively, paving the way for sophisticated AI-assisted engineering solutions.

ACKNOWLEDGMENTS

The authors wish to express their gratitude to Henry Smith and the MIT Motorsports team for their critical contribution to dataset creation and verification, which significantly enhanced the quality of this research. This material is based upon work supported by the National Science Foundation Graduate Research Fellowship. Any opinion, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

REFERENCES

- [1] OpenAI. "GPT-4V(ision) System Card." 2023. URL <https://api.semanticscholar.org/CorpusID:263218031>.
- [2] Vaswani, Ashish, Shazeer, Noam, Parmar, Niki, Uszkoreit, Jakob, Jones, Llion, Gomez, Aidan N, Kaiser, Łukasz and Polosukhin, Illia. "Attention is all you need." *Advances in neural information processing systems* Vol. 30 (2017).
- [3] Bubeck, Sébastien, Chandrasekaran, Varun, Eldan, Ronen, Gehrke, Johannes, Horvitz, Eric, Kamar, Ece, Lee, Peter, Lee, Yin Tat, Li, Yuanzhi, Lundberg, Scott et al. "Sparks of artificial general intelligence: Early experiments with GPT-4. arXiv." *arXiv preprint arXiv:2303.12712* (2023).
- [4] Barbhuiya, Rejaul Karim. "Introduction to Artificial Intelligence: Current Developments, Concerns and Possibilities for Education." *Indian Journal of Educational Technology* Vol. 5 No. 2 (2023): p. 266.
- [5] Thirunavukarasu, Arun James, Ting, Darren Shu Jeng, Elangovan, Kabilan, Gutierrez, Laura, Tan, Ting Fang and Ting, Daniel Shu Wei. "Large language models in medicine." *Nature medicine* Vol. 29 No. 8 (2023): pp. 1930–1940.
- [6] Clusmann, Jan, Kolbinger, Fiona R, Muti, Hannah Sophie, Carrero, Zunamys I, Eckardt, Jan-Niklas, Laleh, Narmin Ghaffari, Löffler, Chiara Maria Lavinia, Schwarzkopf, Sophie-Caroline, Unger, Michaela, Veldhuizen, Gregory P et al. "The future landscape of large language models in medicine." *Communications Medicine* Vol. 3 No. 1 (2023): p. 141.
- [7] Kasneci, Enkelejda, Seßler, Kathrin, Küchemann, Stefan, Bannert, Maria, Dementieva, Daryna, Fischer, Frank, Gasser, Urs, Groh, Georg, Günnemann, Stephan, Hüllermeier, Eyke et al. "ChatGPT for good? On opportunities and challenges of large language models for education." *Learning and individual differences* Vol. 103 (2023): p. 102274.
- [8] Makatura, Liane, Foshey, Michael, Wang, Bohan, Hähnlein, Felix, Ma, Pingchuan, Deng, Bolei, Tjandrasuwita, Megan, Spielberg, Andrew, Owens, Crystal Elaine, Chen, Peter Yichen et al. "How Can Large Language Models Help Humans in Design and Manufacturing?" *arXiv preprint arXiv:2307.14377* (2023).
- [9] Picard, Cyril, Edwards, Kristen M, Doris, Anna C, Man, Brandon, Giannone, Giorgio, Alam, Md Ferdous and Ahmed, Faez. "From Concept to Manufacturing: Evaluating Vision-Language Models for Engineering Design." *arXiv preprint arXiv:2311.12668* (2023).
- [10] Ulrich, Karl T and Eppinger, Steven D. *Product design and development*. McGraw-hill (2016).
- [11] Zeng, Yan, Zhang, Hanbo, Zheng, Jiani, Xia, Jiangnan, Wei, Guoqiang, Wei, Yang, Zhang, Yuchen and Kong, Tao. "What Matters in Training a GPT4-Style Language Model with Multimodal Inputs?" *arXiv preprint arXiv:2307.02469* (2023).
- [12] Wang, Lei, Hu, Yi, He, Jiabang, Xu, Xing, Liu, Ning, Liu, Hui and Shen, Heng Tao. "T-SciQ: Teaching Multimodal Chain-of-Thought Reasoning via Large Language Model Signals for Science Question Answering." *arXiv preprint arXiv:2305.03453* (2023).
- [13] Liu, Haotian, Li, Chunyuan, Wu, Qingyang and Lee, Yong Jae. "Visual instruction tuning." *Advances in neural information processing systems* Vol. 36 (2024).
- [14] Shaham, Uri, Ivgi, Maor, Efrat, Avia, Berant, Jonathan and Levy, Omer. "ZeroSCROLLS: A Zero-Shot Benchmark for Long Text Understanding." *arXiv preprint arXiv:2305.14196* (2023).
- [15] Xiong, Wenhan, Liu, Jingyu, Molybog, Igor, Zhang, Hejia, Bhargava, Prajjwal, Hou, Rui, Martin, Louis, Rungta, Rashi, Sankararaman, Karthik Abinav, Oguz, Barlas et al. "Effective long-context scaling of foundation models." *arXiv preprint arXiv:2309.16039* (2023).
- [16] Richardson, Matthew, Burges, Christopher JC and Renshaw, Erin. "Mctest: A challenge dataset for the open-domain machine comprehension of text." *Proceedings of the 2013 conference on empirical methods in natural language processing*: pp. 193–203. 2013.
- [17] Rajpurkar, Pranav, Zhang, Jian, Lopyrev, Konstantin and Liang, Percy. "Squad: 100,000+ questions for machine comprehension of text." *arXiv preprint arXiv:1606.05250* (2016).
- [18] Yang, Yi, Yih, Wen-tau and Meek, Christopher. "Wikiqa: A challenge dataset for open-domain question answering."

- Proceedings of the 2015 conference on empirical methods in natural language processing*: pp. 2013–2018. 2015.
- [19] Dasigi, Pradeep, Lo, Kyle, Beltagy, Iz, Cohan, Arman, Smith, Noah A and Gardner, Matt. “A dataset of information-seeking questions and answers anchored in research papers.” *arXiv preprint arXiv:2105.03011* (2021).
- [20] Fu, Chaoyou, Chen, Peixian, Shen, Yunhang, Qin, Yulei, Zhang, Mengdan, Lin, Xu, Qiu, Zhenyu, Lin, Wei, Yang, Jinrui, Zheng, Xiawu et al. “MME: A Comprehensive Evaluation Benchmark for Multimodal Large Language Models.” *arXiv preprint arXiv:2306.13394* (2023).
- [21] Liu, Yuan, Duan, Haodong, Zhang, Yuanhan, Li, Bo, Zhang, Songyang, Zhao, Wangbo, Yuan, Yike, Wang, Jiaqi, He, Conghui, Liu, Ziwei et al. “Mmbench: Is your multi-modal model an all-around player?” *arXiv preprint arXiv:2307.06281* (2023).
- [22] Lu, Pan, Mishra, Swaroop, Xia, Tanglin, Qiu, Liang, Chang, Kai-Wei, Zhu, Song-Chun, Tafjord, Oyvind, Clark, Peter and Kalyan, Ashwin. “Learn to explain: Multimodal reasoning via thought chains for science question answering.” *Advances in Neural Information Processing Systems* Vol. 35 (2022): pp. 2507–2521.
- [23] Chen, Yang, Hu, Hexiang, Luan, Yi, Sun, Haitian, Changpinyo, Soravit, Ritter, Alan and Chang, Ming-Wei. “Can Pre-trained Vision and Language Models Answer Visual Information-Seeking Questions?” *arXiv preprint arXiv:2302.11713* (2023).
- [24] Song, Binyang, Zhou, Rui and Ahmed, Faez. “Multi-modal machine learning in engineering design: A review and future directions.” *Journal of Computing and Information Science in Engineering* Vol. 24 No. 1 (2024): p. 010801.
- [25] Dima, Alden, Lukens, Sarah, Hodkiewicz, Melinda, Sexton, Thurston and Brundage, Michael P. “Adapting natural language processing for technical text.” *Applied AI Letters* Vol. 2 No. 3 (2021): p. e33.
- [26] Brundage, Michael P, Sexton, Thurston, Hodkiewicz, Melinda, Dima, Alden and Lukens, Sarah. “Technical language processing: Unlocking maintenance knowledge.” *Manufacturing Letters* Vol. 27 (2021): pp. 42–46.
- [27] Jiang, Shuo, Hu, Jie, Magee, Christopher L and Luo, Jianxi. “Deep learning for technical document classification.” *IEEE Transactions on Engineering Management* (2022).
- [28] Zhu, Qihao and Luo, Jianxi. “Generative transformers for design concept generation.” *Journal of Computing and Information Science in Engineering* Vol. 23 No. 4 (2023): p. 041003.
- [29] Ferrari, Alessio, Spagnolo, Giorgio Oronzo and Gnesi, Stefania. “Pure: A dataset of public requirements documents.” *2017 IEEE 25th International Requirements Engineering Conference (RE)*: pp. 502–505. 2017. IEEE.
- [30] Chang, Yupeng, Wang, Xu, Wang, Jindong, Wu, Yuan, Zhu, Kaijie, Chen, Hao, Yang, Linyi, Yi, Xiaoyuan, Wang, Cunxiang, Wang, Yidong et al. “A survey on evaluation of large language models.” *arXiv preprint arXiv:2307.03109* (2023).
- [31] Papineni, Kishore, Roukos, Salim, Ward, Todd and Zhu, Wei-Jing. “Bleu: a method for automatic evaluation of machine translation.” *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*: pp. 311–318. 2002.
- [32] Blecher, Lukas, Cucurull, Guillem, Scialom, Thomas and Stojnic, Robert. “Nougat: Neural optical understanding for academic documents.” *arXiv preprint arXiv:2308.13418* (2023).
- [33] Lin, Chin-Yew. “Rouge: A package for automatic evaluation of summaries.” *Text summarization branches out*: pp. 74–81. 2004.
- [34] Liu, Haotian, Li, Chunyuan, Wu, Qingyang and Lee, Yong Jae. “Visual Instruction Tuning.” (2023).
- [35] Liu, Haotian, Li, Chunyuan, Li, Yuheng and Lee, Yong Jae. “Improved baselines with visual instruction tuning.” *arXiv preprint arXiv:2310.03744* (2023).
- [36] Liu, Jerry. “LlamaIndex.” (2022). DOI [10.5281/zenodo.1234](https://doi.org/10.5281/zenodo.1234). URL https://github.com/jerryjliu/llama_index.