

RAG vs. Long Context: Examining Frontier Large Language Models for Environmental Review Document Comprehension

Anonymous ACL submission

Abstract

Large Language Models (LLMs) have been applied to many research problems across various domains. One of the applications of LLMs is providing question-answering systems that cater to users from different fields. The effectiveness of LLM-based question-answering systems has already been established at an acceptable level for users posing questions in popular and public domains such as trivia and literature. However, it has not often been established in niche domains that traditionally require specialized expertise. To this end, we construct the NEPAQuAD1.0 benchmark to evaluate the performance of three frontier LLMs – Claude Sonnet, Gemini, and GPT-4 – when answering questions originating from Environmental Impact Statements prepared by U.S. federal government agencies in accordance with the National Environmental Environmental Act (NEPA). We specifically measure the ability of LLMs to understand the nuances of legal, technical, and compliance-related information present in NEPA documents in different contextual scenarios. For example, we test the LLMs’ internal prior NEPA knowledge by providing questions without any context, as well as assess how LLMs synthesize the contextual information present in long NEPA documents to facilitate the question/answering task. We compare the performance of the long context LLMs and RAG powered models in handling different types of questions (e.g., problem-solving, divergent). Our results suggest that RAG powered models significantly outperform the long context models in the answer accuracy regardless of the choice of the frontier LLM. Our further analysis reveals that many models perform better answering closed questions than divergent and problem-solving questions.

1 Introduction

As Large Language Models (LLMs) become increasingly commonplace, researchers have dis-

covered that these models are useful for many tasks beyond text generation. Specifically, LLMs have shown potential utility in niche domains (like science) that would traditionally require specialized expertise, both in a pure text setting (Horawalavithana et al., 2022; Munikoti et al., 2024) and by incorporating data of various modalities (Dollar et al., 2022; Horawalavithana et al., 2023). Recent work has been done to evaluate these models (Acharya et al., 2023; Munikoti et al., 2023; Cai et al., 2024) and to assess their uncertainty (Wagle et al., 2024). Despite extensive research, constructing LLMs for answering domain-specific questions has proven challenging (Kasneci et al., 2023).

One such challenge for LLM-based question-answering systems occurs when systems are tasked with surfacing answers to questions from the content of long documents in specialized domains. Existing LLMs allow users to include a paragraph as context along with the content of the question; however, LLMs generally limit the size of that paragraph to a specific number of tokens. This restriction forces users to truncate or manually summarize the content of lengthy documents into short passages. Another approach users can take includes submitting only the question and relying on the models to find the correct document from a corpus and relevant content needed to answer the question. This strategy often works well for answering questions from well-known domains (e.g., sports or education); however, it is not as successful for less pervasive topics (Munikoti et al., 2023). Because LLMs are data-driven, they are not as apt to provide accurate answers for questions about less popular, more specialized domains such as Law (Kapoor et al., 2024) and Energy (Buster et al., 2024).

In this work, we focus on assessing the long-context LLMs in the environmental reviews conducted under the National Environmental Policy Act (NEPA)¹. NEPA is a U.S. environmental law

¹<https://www.epa.gov/nepa>

designed to protect the environment. An environmental impact statement (EIS) is required by Section 102(2)(C) of NEPA for any proposed major federal action significantly affecting the quality of the environment. An EIS is a detailed document that describes a proposed action, alternatives to the proposed action, and potential effects of the proposed action and alternatives on the environment. An EIS contains information about environmental permitting and policy decisions and considers a range of reasonable alternatives, analyzes the potential impacts resulting from the proposed action and alternatives, and demonstrates compliance with other applicable environmental laws and executive orders.

Along with the fact that EIS documents are usually lengthy (often several hundred pages) and are created by NEPA experts, another factor that can hinder the application of LLMs in this domain is that the development of an EIS document requires NEPA experts with various subject matter expertise to engage in preparation over multiple years, often citing older articles from as far back as the 1990s. For example, the Executive Order (EO) 12898, issued in 1994, is cited on page 60 of the EIS documents prepared for the First Responder Network Authority project². This could present significant challenges for current LLMs in helping NEPA users automatically retrieve answers from LLM-based question-answering systems. To our knowledge, there is no ground-truth benchmark built specifically for this domain to evaluate the quality of LLMs’ output for QA task when the questions pertain to EIS documents.

In this work, we leverage both long context and Retrieval-Augmented Generation (RAG) (Lewis et al., 2021) to develop LLM capability for question-answering over EIS documents (Figure 1). We select frontier LLMs for our experiments: Claude Sonnet (Team and Collaborators), Gemini (Team and Collaborators, 2024), and GPT-4 (OpenAI, 2024). To assess the efficacy of our proposed RAG model compared to other context-augmentation strategies, we also conduct rigorous experiments evaluating LLMs with various types of contexts for NEPA documents. To evaluate our approach, we establish a benchmark using ground truth triples of questions, answers, and corresponding contexts, generated through a semi-supervised

²<https://www.energy.gov/nepa/eis-0530-nationwide-public-safety-broadband-network-programmatic-environmental-impact>

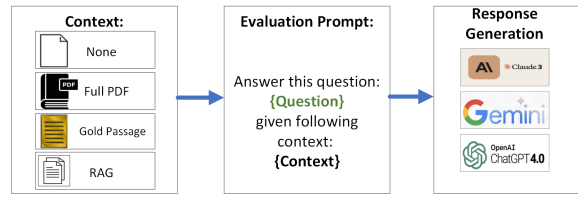


Figure 1: Illustration of varied EIS contexts used in the comparison.

method employing GPT-4. Overall, we make the following contributions:

1. Created the first-ever preliminary benchmark (NEPAQuAD1.0) to automatically evaluate the performance of LLMs in a question-answering task for EIS documents
2. Evaluated the capability of LLMs in question-answering tasks over long documents
3. Conducted rigorous comparisons of LLMs using zero-shot prompting versus context-driven prompting (i.e., passage, PDF, and RAG) to assess model performance.

The structure of this paper is outlined as follows: In Section 2, we describe the benchmark creation to assess the quality of our model in comparison to models derived from different contexts. The Section 3 section lays out our approach and the various contexts used for evaluation implementation, followed by a detailed analysis of our performance in Sections 4. Section 5 then discusses other works in literature that deal with long context and RAG for long documents. We finish with the conclusion and limitations of our work in Sections 6 and 7.

2 NEPAQuAD Benchmark

In this section, we describe the construction of a ground-truth benchmark to evaluate the quality of automated responses generated from LLMs. Due to the high costs associated with manually creating human-written questions and answers, and the inability to use ground-truth benchmarks from other domains, we adopt a semi-automatic approach to generate the NEPAQuAD1.0 (National Environmental Policy Act Question and Answering Dataset) benchmark. The general idea of our evaluation benchmark generation process is to extract meaningful passages from a set of EIS documents, then use GPT-4 to generate questions based on these passages. To ensure the quality of

the generated benchmark, two authors of this study, who are subject matter experts in NEPA, measure the quality of the ground-truth answers by comparing the provided proofs against the original context from which the questions were derived. Our generated ground-truth benchmark is a set of triples containing a question, an answer, and the proof (i.e., the text directly related to the answer, derived from the context from which the question originated). The process of benchmark generation is illustrated in Figure 2.

To evaluate performance of LLMs for the EIS question-answering task, we first select high-quality documents from the EIS document database and extract paragraphs as context to be used in the evaluation. Then, we identify the types of questions that we want to use to evaluate the models. Next, we use GPT-4 (OpenAI, 2024) to generate question-answer pairs based on the selected contexts by using carefully designed prompts. Finally, we use these generated questions to evaluate different LLMs with various contexts, with the generated answers acting as the ground-truth. We describe the process in detail below.

2.1 Gold Passage Selection

Document Selection NEPA experts select nine EIS documents from different government agencies that were most representative of various NEPA actions. These document exhibit great variations in content and focus depending on the authoring government agency, as each agency may interpret and implement the NEPA guidelines distinctively. For instance, the U.S. Forest Service might emphasize forest management and wildfire mitigation, while the Army Corps of Engineers could prioritize water resource development and infrastructure impacts. Table 3 shows the statistics about the selected documents (see Appendix). Each document has around 400 pages on average while the longest document contain more than 600K tokens.

Excerpt Selection For each of nine selected documents, we attempt to select excerpts that have important content of each document. Again, the default approach of randomly extracting excerpts poses the risk of resulting in *low-quality* excerpts, such as parts of appendix or images’ captions. To avoid this risk, NEPA experts manually select excerpts from the documents. They divided each document into three sections: beginning, middle, and end, and then selected two, six, and six ex-

cerpts from each of these sections respectively, for a total of 10 excerpts from each document. We use these excerpts as the ground-truth context, called gold passages, for question benchmark generation.

2.2 Question Type Selection

Once we identified the gold passages, NEPA experts select the type of questions to include in the benchmark. We started with a list of 15 types of questions³, and eventually narrowed it down to 10 types of questions after extensive discussions. These types are shown in Table 1. In addition to selecting the question types, the NEPA experts also created sample questions for each category for the EIS document domain. For a more detailed description of the question types, as well as example questions, please see Appendices A and B.

Question Type	#Questions	
Closed	789	49%
Comparison	64	4%
Convergent	109	7%
Divergent	121	8%
Evaluation	64	4%
Funnel	127	8%
Inference	101	6%
Problem-solving	11	1%
Process	108	7%
Recall	105	7%
Total	1599	100%

Table 1: Statistics on question types used in the NEPAQuAD1.0 benchmark

2.3 Prompt Design

The next step is to design a prompt that can instruct the question generation model to create high quality questions and answers. To ensure that the prompt can instruct generative model efficiently, we took advice from the NEPA experts to create the prompts. We also use the sample questions created by the NEPA experts to augment the original prompts and create an “enhanced” prompt. The template for the prompt and benchmark creation process is displayed in Figure 2. We restricted the output for each prompt in a CSV format with three fields: question, answer, and proof. The “proof” column stored the part of the gold passage that the model picked as evidence for the provided answer to the question.

³<https://tinyurl.com/3akej8ct>

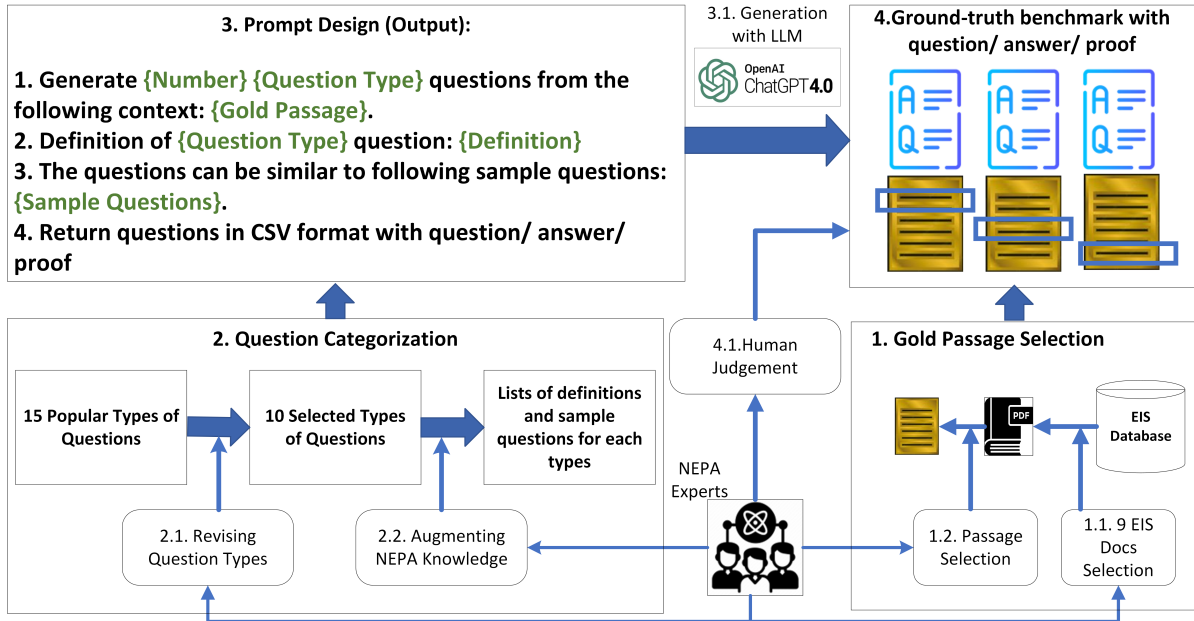


Figure 2: Steps of Ground-truth benchmark generation for evaluating LLMs over varied contexts for question-answering over EIS documents

2.4 Ground-truth Benchmark Generation

Automated Generation We selected GPT-4 as the generative model for this task, as GPT-4 has been used for generating questions and answers for documents in other domains such as agriculture (Balaguer et al., 2024). Specifically, we use GPT-4 Azure OpenAI service with default setting to execute the question generation prompts. For each document in our nine selected documents, we have 10 gold passages, results in 90 gold passages in total. For each gold passage, we generated 10 sets of questions for each of the 10 question types. We then filtered the generations for incorrect formats. Overall, we generated a benchmark of 1599 ground-truth triples of question-answer-proof over the nine selected EIS documents.

Quality Check In order to judge the generated benchmark’s correctness, we randomly select 100 sample questions for validation. The validation was done by the NEPA experts (two co-authors). Each of them independently went through the sample questions, checking both the correctness of question type (i.e. whether the generated question was the same type of question as requested) and the correctness of the answer. For each question, if either of the evaluators marked the question or answer as correct, we labelled them as correct. Overall, our generated benchmark achieved 77% of answer correctness and 82% of question type correctness.

3 Experimental Setup

In this work, we experimented with three frontier LLMs: Claude-3 Sonnet, Gemini, and GPT-4. For the context provided to the model, we had four possible variations: no context, PDF documents, silver passages (RAG setup), or gold passages. The combination of the models and context settings resulted in a total of 12 unique configurations (Figure 1). We explain these configurations in details in Section 3.1 and report the evaluation metrics in Section 3.2.

3.1 Context Variation

No Context: In this setting, we simply query the models with the questions with no context provided, the same as in other general domains. We do not provide any additional context about the origin of the questions, so the models were expected to answer questions from their existing knowledge. While this strategy can work well with popular domains such as sport or literature, we assume that NEPA domain may be challenging for the LLM models to get the accurate answer. In the other word, this setting can be said to be a test of the LLMs to answer out-of-general-domain questions. As such, this setting is usually expected to return low performance.

Full PDF as Context: In this scenario, in addition to the question, we also provide the model the

PDF (text) document from which the context to the question was extracted. Since we do not inform the model which part of the document to look at, the generated responses’ accuracy will heavily depend on the models’ ability to pick the correct context from the very large scale textual information provided. We expect this setting to yield performance better than no context.

RAG Context (Silver Passages): In RAG models, when a question is inputted for LLM generation, the corresponding context is extracted as a relevant passage from a given EIS document. We use the standard RAG setup where we encode both question and retrieved passages with BGE embedding model (Xiao et al., 2023). We use the cosine similarity score to assess the similarity between the question and the contexts. The number of top-ranked relevant passages extracted, referred to as top-K silver passages, is set at $k = 3$.

Gold Passage as Context: In this configuration, we include the actual context from which the question was generated in the prompt, alongside the question content. While the scenario where users manually identify the correct passage is rare in practice, we simulate this scenario to measure how well LLMs can perform if we were able to extract relevant passages with very high accuracy. We expect this setting to perform the best.

3.2 Evaluation Metrics

In order to evaluate the performance of the models in different configurations, We compare the answers from the generated responses of the model across these various configurations. Overlap based metrics such as BLEU score (Papineni et al., 2002), while used by many prior works, simply measure the syntactic similarity, and as such is not suitable to perform evaluations where semantics is more important. As such, for our work, we use the RAGAs score proposed by Es et al. (2023): the Answer Correctness (called Correctness in this paper).

The Answer Correctness score combined two aspects of factual correctness and semantic correctness for its calculation. While factual correctness captured the correctness at phrase/clause level of input answer, the semantic score is achieved by comparing the similarity between vectors of expected answers and predicted answers by using embedding models. GPT-4 is used in calculating the answer correctness that quantifies the factual overlap between the generated answer and the ground truth

answer (Es et al., 2023). We use the BGE (Xiao et al., 2023) as the embedding model for semantic correctness calculation. We set the weight of factual correctness as 0.75 and the weight of semantic correctness as 0.25 for measuring the Answer Correctness.

4 Performance Analysis

In this section, we describe the overall performance of the LLMs in the question/answering task evaluated with NEPAQuAD1.0 (as presented in Section 2). First, we compare the performance of three frontier LLMs: Claude-3 Sonnet, Gemini, and GPT-4 across various QA contexts (Section 4.1). Second, we compare the model performance across various question types (Section 4.2). In Section 4.3, we evaluate how the models performing to the questions generated from different parts of the document. Finally, we analyze the performance.

4.1 Evaluating Different QA Contexts

Table 2 reports the overall performance of the models across various QA contexts used in the evaluation. We observed that for the task with no context, the Gemini model produces the most accurate results by far. However, when PDF documents are provided as context, this trend is reversed, with GPT-4 surpassing Gemini in correctness. Despite that Gemini is able to handle very long contexts (1.5M tokens), it is surprising to see its performance drop when provided with PDF documents as additional contexts. This may be due to the model struggling to reason over the large amount of relevant and irrelevant content in the EIS document.

Overall, RAG models perform better in comparison to the models provided with PDF documents as additional contexts. In RAG setup, The Claude model outperforms both other models in term of correctness, although the scores across the Claude and GPT-4 models are much closer. There is notable increase in Gemini’s performance in the RAG setup when compared with the PDF contexts.

As expected, all models perform best on average when provided with the gold passage in comparison to other context variations. In this scenario, model needs to synthesize information that directly contains the answer to the question posed to the model. Notably, models perform comparably when they are provided with the RAG and gold passage contexts.

Context	Claude	Gemini	GPT-4
None	21.50%	50.16%	20.28%
Complete PDF	23.47%	46.62%	56.40%
Silver Passages	68.74%	57.06%	66.86%
Gold Passage	68.41%	61.81%	67.66%

Table 2: Evaluation on the answer correctness of LLMs over different configurations of context over EIS documents. Silver passages are selected by the RAG model.

4.2 Evaluating Different Question Types

We analyze the performance of LLMs over different types of questions depicted in Figure 3. When analyzing the results based on the type of questions, we see that all three models have superior performance on closed questions when provided with either silver or gold data as context, while GPT-4 is the only model that performs well on these questions even when provided with just the PDFs as context. For all other categories, both Claude and GPT-4 exhibit similar behavior pattern when provided with none or PDF context, although GPT4’s performance is notably better than Claude’s in almost every category, with this difference particularly noticeable with PDF context.

Interestingly, Gemini performs really well when provided with no context at all, and the performance decreases for all categories except the convergent and recall questions when provided with the PDFs as context. For a majority of the question types, even the silver or gold data is not able to get the performance to the level of no-context, with the closed and recall questions being the exceptions. Overall, RAG models and models with other contexts performed best in answering closed questions and worst in answering divergent and problem-solving questions.

4.3 Evaluating Positional Knowledge

We also analyzed the performance of various question types based on the portion of the document from which they were derived, as shown in Figure 4. Across the models, we observed a general trend where the earlier the source text in the document, the better the performance of the models. A notable exception to this are the problem-solving questions, which perform better when sourced from the latter parts of the document. This pattern holds true for all three models. Additionally, we noticed that divergent questions yield better results when derived from the middle of the document. Over-

all, all three models exhibit similar or comparable patterns of performance across different document sections and question types. These results suggest that performance of long-context models may vary not only by the position of relevant information, but also due to the type of the question and the amount of reasoning that the model has to perform.

4.4 Discussion

RAG for Question-Answering Over Domain-Specific Documents The findings of this study underscore the significant role of RAG models as crucial strategies for addressing domain-specific questions. These models have shown remarkable superiority in performance compared to zero-shot knowledge and using the full PDF as context. While evaluating LLMs in well-established fields such as mathematics or biology can be straightforward, using numerous human-written sets of ground-truth questions and answers (Team and Collaborators, 2024), evaluating domain-specific LLMs necessitates unsupervised or semi-supervised approaches to generate evaluation benchmarks. We recognize that while our approach satisfies the need for an automated method in this domain, it still faces challenges, particularly that the selected question types might not be representative of other research areas. Therefore, researchers in other domains should carefully consider the types of questions they want to generate for their studies.

Patterns of Output From this study, we draw two overall conclusions regarding the output patterns. First, surprisingly, we did not observe specific patterns of correctness in relation to document metadata such as token counts. This finding contradicts our initial assumption that documents with the lowest token counts would achieve the lowest accuracy and vice versa. We believe one reason for this may be that we selected only 90 passages as gold passages from the document, which might not be representative. Second, we noted that each LLM model tends to have distinctive response types. For instance, Gemini’s responses tend to be straightforward when no context is provided, often stating "I don’t have the context." In contrast, Claude and GPT-4 are more likely to attempt clarifications of input questions, such as predicting and providing the full content of abbreviations. We encourage researchers in other projects involving RAG to analyze patterns of output to enhance their work.

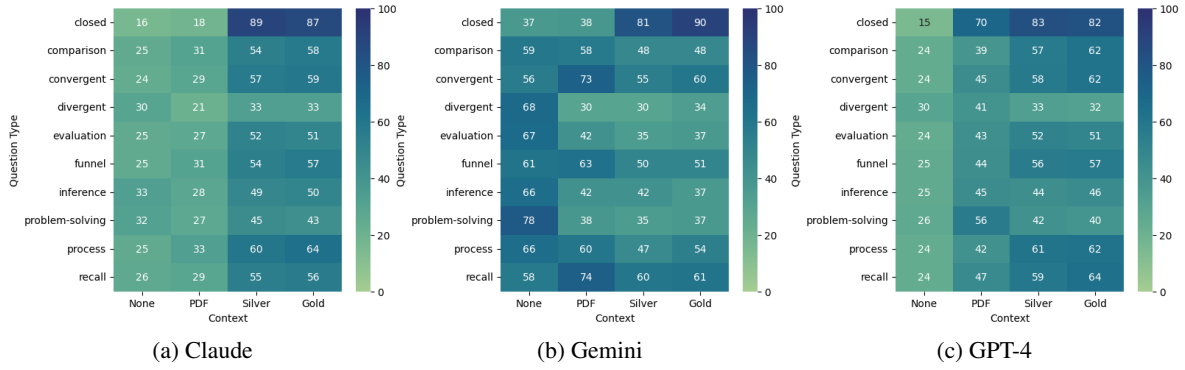


Figure 3: The evaluation results measured by the Answer Correctness scores of each LLM used with 4 scenarios of using context over each question types

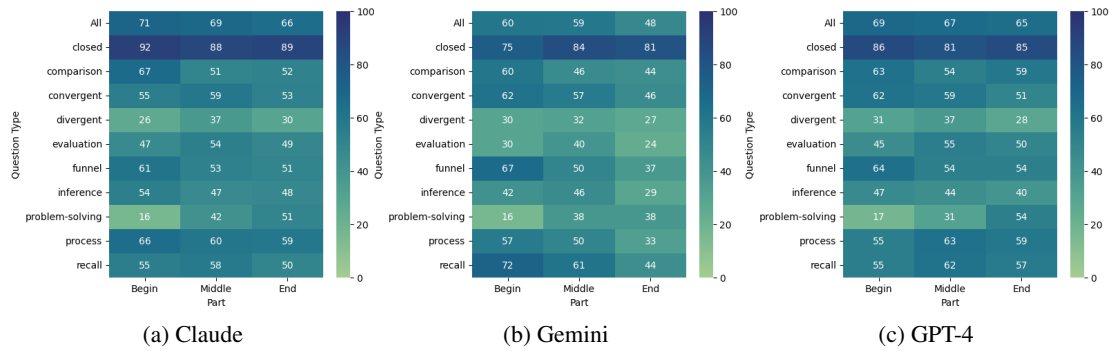


Figure 4: The evaluation results measured by the Answer Correctness scores of each LLM used over different parts with silver passages provided as context

Long Context Reasoning One of the primary objectives of this study is to assess how beneficial the long context models that can process 128K to 1.5M tokens context in answering questions from long EIS documents. We noticed that these models struggle to use long input contexts to answer more difficult questions that require multiple steps of reasoning (e.g, problem-solving). Given that we see the model performance varies over the positions and types of questions, we assume that effective question-type and -complexity aware reranking of retrieved documents may help to improve the performance (Jeong et al., 2024). For example, we can use of another LLM to adjust the order of retrieved chunks based on the problem-solving question type.

5 Related Work

Long Context Evaluation A popular technique to evaluate the long context LLMs is with a simple *needle in a haystack* analysis to test in-context retrieval ability (Chandrayan et al., 2024). Despite the claims made in these tests, it is shown current LLMs perform poorly in processing and

understanding long, context-rich sequences in rigorous scientific benchmarks (Li et al., 2024; Liu et al., 2024). For example, Li et al. (2024) constructed LongICLBench benchmark to assess a set of long-context LLMs in an extreme-label classification task as a long in-context learning task. They showed that long context understanding and reasoning is still a challenging task for the existing LLMs. Few studies showed that long-context LLMs are affected by the position of the relevant information in the input context (Liu et al., 2024; Ivgi et al., 2023). For instance, Ivgi et al. (2023) showed that encoder-decoder models have significantly higher performance when relevant information is placed at the start of the input context. In addition, Liu et al. (2024) showed that LLMs perform weakly when they must access relevant information in the middle of long contexts.

There are some other benchmarks proposed in multiple languages and domains for evaluating LLM’s long context understanding. Bai et al. (2023) proposed LongBench that covers six tasks, single-doc QA, multi-doc QA, summarization, few-shot learning, synthetic tasks, and code completion in English and Chinese languages. L-Eval Bench-

mark (An et al., 2023) contains 20 sub-tasks, 508 long documents, and over 2,000 human-labeled query-response pairs with diverse question styles, domains, and input length. Li et al. (2023) proposed LooGLE that includes around 6,000 questions across diverse domains and evaluated both commercial and open-sourced models. While they showed that commercial models outperform open-sourced models in short question-answering and cloze tasks they struggled in long dependency tasks. Furthermore, retrieval-based techniques showed significant advantages for answering short questions, whereas methods aimed at increasing the length of the context window had a minimal effect on the comprehension of longer contexts. ∞ Bench (Zhang et al., 2024) consists of 12 tasks with data length surpassing 100K tokens on average. They suggested that long context LLMs still require significant advancements to effectively process 100K+ context.

RAG for Long Documents RAG models offer a promising approach for enabling LLMs to search and extract relevant information from lengthy documents or extensive collections. The common strategy of splitting documents into smaller, more manageable chunks that fit within the LLM’s context window has its limitations, as highlighted by recent studies (Barnett et al., 2024). These limitations include the model’s failure to accurately extract answers even when they are present within the provided context, particularly when there is excessive noise or contradictory information. To address these challenges, researchers have proposed novel techniques. HippoRAG, a neurobiologically inspired long-term memory system designed for LLMs to handle long documents more effectively, aims to mitigate the limitations of current RAG models (Gutiérrez et al., 2024). Gao et al. (2023) provide a comprehensive survey of RAG methods, categorizing them into three paradigms: Naive RAG, Advanced RAG, and Modular RAG. The authors highlight the remarkable adaptability of Modular RAG, which allows for module substitution or reconfiguration to address specific challenges, surpassing the fixed structures of Naive and Advanced RAG. Modular RAG integrates new modules or adjusts interaction flow among existing ones, enhancing its applicability across different tasks. The survey also discusses the concept of adaptive retrieval in RAG, exemplified by methods like Flare (Jiang et al., 2023) and Self-RAG (Asai et al.,

2023). These approaches refine the RAG framework by enabling LLMs to actively determine the optimal moments and content for retrieval, enhancing the efficiency and relevance of the sourced information. Despite these advancements, Gao et al. (2023) emphasize that further research is needed to fully understand the intricacies of applying RAG to long documents and to develop more robust and reliable methods.

6 Conclusion

In this study, we conduct the initial investigation into the performance of LLMs within the domain of the National Environmental Policy Act and its associated documents. To facilitate this, we introduce NEPAQuAD, a question-and-answering benchmark designed to evaluate a model’s capability to understand the legal, technical, and compliance-related content found in NEPA documents. We assess three frontier LLMs designed for handling extensive contexts—Claude Sonnet, Gemini, and GPT-4—across various contextual settings. Our comprehensive analysis indicates that NEPA documents pose a significant challenge for LLMs, particularly in terms of understanding the complex semantics and effectively processing the lengthy documents. The findings reveal that models augmented with the RAG technique surpass those LLMs that are simply provided with the PDF content as long context. This suggests that incorporating more relevant knowledge retrieval processes can significantly enhance the performance of LLMs on complex document comprehension tasks like those found in the NEPA domain. In addition, we noticed that these LLMs struggle to use long input contexts to answer more difficult questions that require multiple steps of reasoning. For example, models performed best in answering closed questions and worst in answering divergent and problem-solving questions.

7 Limitations

Similar to other applications of LLMs, our proposed system for EIS long documents also has some limitations. We list those limitations as follows:

Restriction of token limitation on full PDF context. While we are able to use the Gemini model with token length as 1.5 million, we could only use 128K tokens per query for response generation with Claude and GPT-4. Thus, we need to

truncate the content of Full PDF to run these two LLM models. This might cause the performance drop on the context as Full PDF with questions from EIS documents. In future work, we should analyze more carefully about the impact of token truncation for Full PDF context.

Uncertainty of generated responses by LLMs. Due to budget constraints, we conducted only one phase of response generation across different configurations. This introduces a risk of uncertain outputs, meaning that LLMs might generate different responses each time, even with the same input, as demonstrated in another study (Wagle et al., 2024). In future work, we plan to run LLMs multiple times and analyze the effect of this uncertainty in response generation.

Challenges of human judgment. Currently, we leverage human evaluation as a preliminary proxy measure for qualitative analysis of benchmark. In future work, we plan to involve more NEPA experts in a more systematic manner to expand the dataset with human judgment results and to perform proper adjudication meetings between NEPA experts to reconcile conflicting results.

Bias in automated evaluation There might be a potential bias in the answer correctness evaluation process due to the use of GPT-4 to assess the outputs of various models. There is a concern that GPT-4 may inherently prefer the outputs generated by the same model over others in the factual correctness evaluation. This could lead to skewed evaluation results, where GPT-4’s outputs are rated more favorably, not necessarily because they are superior, but because of the inherent biases in the evaluation model (GPT-4).

To address the potential bias in the answer correctness evaluation process, we assess both factual and semantic correctness in the evaluation. For semantic correctness, we utilize the BGE (Xiao et al., 2023) as the embedding model and we calculate the semantic similarity between the model’s outputs and the ground-truth answers independently of GPT-4’s own evaluation mechanisms. By combining both factual and semantic correctness, we aim to accurately reflect the true performance of various models, including GPT-4.

8 Ethical Consideration

It has generally been the norm to assume that previously published work can be used as-is without having to consider the inherited ethical issues.

However, in present times, researchers should not “simply assume that [...] research will have a net positive impact on the world” (Hecht et al., 2021). We acknowledge that this applies not just to new work, but also when using existing work in the way that we have done.

While working on this project, care has been taken to ensure that any and all data was anonymized and no Personally Identifiable Information (PII) is present in the data used. We had domain experts in the team throughout the process, thereby ensuring they were aware of all the potential risks and benefits.

While we do not anticipate the novel work presented here to introduce new ethical concerns in and by themselves, we do recognize that there may also be pre-existing concerns and issues of the data, models, and methodologies we have used for this paper. In particular, it has been seen that LLMs, like the ones used in this work, exhibit a wide variety of bias – religious, gender, race, profession, and cultural – and frequently generate answers that are incorrect, misogynistic, antisemitic, and generally toxic (Abid et al., 2021; Buolamwini and Gebre, 2018; Liang et al., 2021; Nadeem et al., 2021; Welbl et al., 2021). However, when used within the parameters of our experiments detailed in this paper, we did not see such behaviour from any of the models. To our knowledge, when used as intended, our models do not pose additional ethical concerns than any other LLM.

References

- Abubakar Abid, Maheen Farooqi, and James Zou. 2021. Persistent anti-muslim bias in large language models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pages 298–306.
- Anurag Acharya, Sai Munikoti, Aaron Hellinger, Sara Smith, Sridevi Wagle, and Sameera Horawalavithana. 2023. Nuclearqa: A human-made benchmark for language models for the nuclear domain. *arXiv preprint arXiv:2310.10920*.
- Chenxin An, Shansan Gong, Ming Zhong, Mukai Li, Jun Zhang, Lingpeng Kong, and Xipeng Qiu. 2023. L-eval: Instituting standardized evaluation for long context language models. *arXiv preprint arXiv:2307.11088*.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *Preprint, arXiv:2310.11511*.

746	Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu,	Brent Hecht, Lauren Wilcox, Jeffrey P Bigham,	800
747	Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao	Johannes Schöning, Ehsan Hoque, Jason Ernst,	801
748	Liu, Aohan Zeng, Lei Hou, et al. 2023. Longbench:	Yonatan Bisk, Luigi De Russis, Lana Yarosh, Bushra	802
749	A bilingual, multitask benchmark for long context	Anjum, et al. 2021. It’s time to do something: Mit-	803
750	understanding. <i>arXiv preprint arXiv:2308.14508</i> .	igating the negative impacts of computing through	804
		a change to the peer review process. <i>arXiv preprint</i>	805
751	Angels Balaguer, Vinamra Benara, Renato Luiz de Fre-	<i>arXiv:2112.09544</i> .	806
752	itas Cunha, Roberto de M. Estevão Filho, Todd		
753	Hendry, Daniel Holstein, Jennifer Marsman, Nick	Sameera Horawalavithana, Ellyn Ayton, Shivam	807
754	Mecklenburg, Sara Malvar, Leonardo O. Nunes,	Sharma, Scott Howland, Megha Subramanian, Scott	808
755	Rafael Padilha, Morris Sharp, Bruno Silva, Swati	Vasquez, Robin Cosbey, Maria Glenski, and Svit-	809
756	Sharma, Vijay Aski, and Ranveer Chandra. 2024.	lana Volkova. 2022. Foundation models of scientific	810
757	Rag vs fine-tuning: Pipelines, tradeoffs, and a case	knowledge for chemistry: Opportunities, challenges	811
758	study on agriculture . <i>Preprint</i> , arXiv:2401.08406.	and lessons learned. In <i>Proceedings of BigScience</i>	812
		<i>Episode# 5–Workshop on Challenges & Perspec-</i>	813
759	Scott Barnett, Stefanus Kurniawan, Srikanth Thudumu,	<i>tives in Creating Large Language Models</i> , pages	814
760	Zach Brannelly, and Mohamed Abdelrazek. 2024.	160–172.	815
761	Seven failure points when engineering a retrieval		
762	augmented generation system. <i>arXiv preprint</i>	Sameera Horawalavithana, Sai Munikoti, Ian Stewart,	816
763	<i>arXiv:2401.05856</i> .	and Henry Kvinge. 2023. Scitune: Aligning large	817
		language models with scientific multimodal instruc-	818
764	Joy Buolamwini and Timnit Gebru. 2018. Gender	tions. <i>arXiv preprint arXiv:2307.01139</i> .	819
765	shades: Intersectional accuracy disparities in com-		
766	mercial gender classification. In <i>Conference on fair-</i>	Maor Ivgi, Uri Shaham, and Jonathan Berant. 2023. Ef-	820
767	<i>ness, accountability and transparency</i> , pages 77–91.	ficient long-text understanding with short-text mod-	821
768	PMLR.	els. <i>Transactions of the Association for Computa-</i>	822
		<i>tional Linguistics</i> , 11:284–299.	823
769	Grant Buster, Pavlo Pinchuk, Jacob Barrons, Ryan Mc-		
770	Keever, Aaron Levine, and Anthony Lopez. 2024.	Soyeong Jeong, Jinheon Baek, Sukmin Cho, Sung Ju	824
771	Supporting energy policy research with large lan-	Hwang, and Jong C Park. 2024. Adaptive-rag:	825
772	guage models. <i>arXiv preprint arXiv:2403.12924</i> .	Learning to adapt retrieval-augmented large lan-	826
		guage models through question complexity. <i>arXiv</i>	827
773	Hengxing Cai, Xiaochen Cai, Junhan Chang, Sihang	<i>preprint arXiv:2403.14403</i> .	828
774	Li, Lin Yao, Changxin Wang, Zhifeng Gao, Yongge		
775	Li, Mujie Lin, Shuwen Yang, et al. 2024. Sciassess:	Zhengbao Jiang, Frank F. Xu, Luyu Gao, Zhiqing Sun,	829
776	Benchmarking llm proficiency in scientific literature	Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie	830
777	analysis. <i>arXiv preprint arXiv:2403.01976</i> .	Callan, and Graham Neubig. 2023. Active retrieval	831
		augmented generation . <i>Preprint</i> , arXiv:2305.06983.	832
778	Kedar Chandrayan, Lance Martin, gkamradt, Lazaro		
779	Hurtado, arkadyark cohere, Ikko Eltociear	Sayash Kapoor, Peter Henderson, and Arvind	833
780	Ashimine, Pavel Král, and Prabha Arivalagan.	Narayanan. 2024. Promises and pitfalls of artificial	834
781	2024. gkamradt/LLMTestNeedleInAHaystack .	intelligence for legal applications. <i>arXiv preprint</i>	835
		<i>arXiv:2402.01656</i> .	836
782	Orion Walker Dollar, Sameera Horawalavithana, Scott		
783	Vasquez, W James Pfaendtner, and Svitlana Volkova.	Enkelejd Kasneci, Kathrin Sessler, Stefan Küche-	837
784	2022. Moljet: multimodal joint embedding trans-	mann, Maria Bannert, Daryna Dementieva, Frank	838
785	former for conditional de novo molecular design and	Fischer, Urs Gasser, Georg Groh, Stephan Günne-	839
786	multi-property optimization.	mann, Eyke Hüllermeier, Stephan Krusche, Gitta	840
		Kutyniok, Tilman Michaeli, Claudia Nerdel, Jür-	841
787	Shahul Es, Jithin James, Luis Espinosa-Anke, and	gen Pfeffer, Oleksandra Poquet, Michael Sailer,	842
788	Steven Schockaert. 2023. Ragas: Automated eval-	Albrecht Schmidt, Tina Seidel, Matthias Stadler,	843
789	uation of retrieval augmented generation . <i>Preprint</i> ,	Jochen Weller, Jochen Kuhn, and Gjergji Kasneci.	844
790	arXiv:2309.15217.	2023. Chatgpt for good? on opportunities and	845
		challenges of large language models for education .	846
791	Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia,	<i>Learning and Individual Differences</i> , 103:102274.	847
792	Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen		
793	Wang. 2023. Retrieval-augmented generation for	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	848
794	large language models: A survey. <i>arXiv preprint</i>	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	849
795	<i>arXiv:2312.10997</i> .	rich Küttler, Mike Lewis, Wen tau Yih, Tim Rock-	850
		täschel, Sebastian Riedel, and Douwe Kiela. 2021.	851
796	Bernal Jiménez Gutiérrez, Yiheng Shu, Yu Gu, Michi-	Retrieval-augmented generation for knowledge-	852
797	hiro Yasunaga, and Yu Su. 2024. Hipporag: Neu-	intensive nlp tasks . <i>Preprint</i> , arXiv:2005.11401.	853
798	robiologically inspired long-term memory for large		
799	language models. <i>arXiv preprint arXiv:2405.14831</i> .	Jiaqi Li, Mengmeng Wang, Zilong Zheng, and Muhan	854
		Zhang. 2023. Loogle: Can long-context language	855

856	models understand long contexts? <i>arXiv preprint</i>	Sridevi Wagle, Sai Munikoti, Anurag Acharya, Sara	909
857	<i>arXiv:2311.04939</i> .	Smith, and Sameera Horawalavithana. 2024. Em-	910
858	Tianle Li, Ge Zhang, Quy Duc Do, Xiang Yue,	pirical evaluation of uncertainty quantification in	911
859	and Wenhui Chen. 2024. Long-context llms strug-	retrieval-augmented language models for science.	912
860	gle with long in-context learning. <i>arXiv preprint</i>	In <i>Proceedings of the Workshop on Scientific Doc-</i>	913
861	<i>arXiv:2404.02060</i> .	ument Understanding (SDU), Vancouver, Canada.	914
862	Paul Pu Liang, Chiyu Wu, Louis-Philippe Morency,	Held in conjunction with the 38th AAAI Conference	915
863	and Ruslan Salakhutdinov. 2021. Towards under-	on Artificial Intelligence (AAAI 2024).	916
864	standing and mitigating social biases in language	Johannes Welbl, Amelia Glaese, Jonathan Uesato,	917
865	models. In <i>International Conference on Machine</i>	Sumanth Dathathri, John Mellor, Lisa Anne Hen-	918
866	<i>Learning</i> , pages 6565–6576. PMLR.	dricks, Kirsty Anderson, Pushmeet Kohli, Ben	919
867	Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paran-	Coppin, and Po-Sen Huang. 2021. Challenges	920
868	jape, Michele Bevilacqua, Fabio Petroni, and Percy	in detoxifying language models. <i>arXiv preprint</i>	921
869	Liang. 2024. Lost in the middle: How language	<i>arXiv:2109.07445</i> .	922
870	models use long contexts. <i>Transactions of the Asso-</i>	Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas	923
871	<i>ciation for Computational Linguistics</i> , 12:157–173.	Muennighoff. 2023. C-pack: Packaged resources	924
872	Sai Munikoti, Anurag Acharya, Sridevi Wagle, and	to advance general chinese embedding. <i>Preprint</i> ,	925
873	Sameera Horawalavithana. 2023. Evaluating the	arXiv:2309.07597.	926
874	effectiveness of retrieval-augmented large language	Xinrong Zhang, Yingfa Chen, Shengding Hu, Zi-	927
875	models in scientific document reasoning. <i>arXiv</i>	hang Xu, Junhao Chen, Moo Khai Hao, Xu Han,	928
876	<i>preprint arXiv:2311.04348</i> .	Zhen Leng Thai, Shuo Wang, Zhiyuan Liu, et al.	929
877	Sai Munikoti, Anurag Acharya, Sridevi Wagle,	2024. <i>infty</i> bench: Extending long context	930
878	and Sameera Horawalavithana. 2024. Atlantic:	evaluation beyond 100k tokens. <i>arXiv preprint</i>	931
879	Structure-aware retrieval-augmented language	<i>arXiv:2402.13718</i> .	932
880	model for interdisciplinary science. In <i>Proceedings</i>		
881	<i>of the Workshop on AI to Accelerate Science and</i>		
882	<i>Engineering (AI2ASE)</i> , Vancouver, Canada. Held		
883	in conjunction with the 38th AAAI Conference on		
884	Artificial Intelligence (AAAI 2024).		
885	Moin Nadeem, Anna Bethke, and Siva Reddy. 2021.		
886	<i>StereoSet: Measuring stereotypical bias in pre-</i>		
887	<i>trained language models</i> . In <i>Proceedings of the</i>		
888	<i>59th Annual Meeting of the Association for Compu-</i>		
889	<i>tational Linguistics and the 11th International Joint</i>		
890	<i>Conference on Natural Language Processing (Vol-</i>		
891	<i>ume 1: Long Papers)</i> , pages 5356–5371, Online. As-		
892	sociation for Computational Linguistics.		
893	OpenAI. 2024. <i>Gpt-4 technical report</i> . <i>Preprint</i> ,		
894	arXiv:2303.08774.		
895	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-		
896	Jing Zhu. 2002. Bleu: a method for automatic eval-		
897	uation of machine translation. In <i>Proceedings of the</i>		
898	<i>40th annual meeting of the Association for Compu-</i>		
899	<i>tational Linguistics</i> , pages 311–318.		
900	Anthropic Team and Collaborators.		
901	Article about claude 3 models.		
902	https://www-cdn.anthropic.com/		
903	de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/		
904	Model_Card_Claude_3.pdf . Accessed: 2024-05-		
905	06.		
906	Gemini Team and Collaborators. 2024. <i>Gemini: A fam-</i>		
907	<i>ily of highly capable multimodal models</i> . <i>Preprint</i> ,		
908	arXiv:2312.11805.		

933	A Question Definitions	B Sample Questions	979
934	NEPA experts reviewed and created the definitions	In this sections, we listed the sets of sample ques-	980
935	for each question types as following.	tions we used for each types of questions.	981
936	1. Closed questions: Closed questions have	B.1 Closed questions	982
937	two possible answers depending on how you	• Are there any federally recognized Tribes in a	983
938	phrase it: “yes” or “no” or “true” or “false.”	50-mile radius of [PROJECT]?	984
939	2. Comparison questions: Comparison ques-	• Are there any federally recognized species of	985
940	tions are higher-order questions that ask lis-	concern in a 50-mile radius of [PROJECT]?	986
941	teners to compare two things, such as objects,	• Did [AGENCY] approve the licensing action	987
942	people, ideas, stories or theories.	• Did the EIS consider [SUBJECT]?	988
943	3. Convergent questions: convergent questions	B.2 Comparison questions	989
944	are designed to try and help you find the so-	• Which Tribes were consulted in [PROJECT	990
945	lution to a problem, or a single response to a	1] and not [PROJECT 2] and vice-versa?	991
946	question.	• What are some differences between [STUDY	992
947	4. Divergent questions: Divergent questions	1] and [STUDY 2] that might account for dif-	993
948	have no right or wrong answers but rather	ferences in species count for [SPECIES]?	994
949	encourage open discussion. While they are	• Compare the considered alternatives in	995
950	similar to open questions, divergent questions	[PROJECT 1] with those in [PROJECT 2].	996
951	differ in that they invite the listener to share an	• Compare the outcomes of surveys from the	997
952	opinion, especially one that relates to future	new reactor EIS with the license renewal EIS	998
953	possibilities.	for [RPROJECT].	999
954	5. Evaluation questions: Evaluation questions,	B.3 Convergent questions	1000
955	sometimes referred to as key evaluation ques-	• Which other species of concern could logi-	1001
956	tions or KEQs, are high-level questions that	cally be in within the 50-mile radius around	1002
957	are used to guide an evaluation. Good evalu-	the [PROJECT]?	1003
958	ation questions will get to the heart of what	• How many similar projects could be built be-	1004
959	it is you want to know about your program,	fore the impact level for air quality was rated	1005
960	policy or service.	as high?	1006
961	6. Funnel questions: Funnel questions are al-	• If the area of effect for the proposed action	1007
962	ways a series of questions. Their sequence	were increased by 50%, what additional fed-	1008
963	mimics a funnel structure in that they start	eral species of concern would need to be ad-	1009
964	broadly with open questions, then segue to	ressed?	1010
965	closed questions.	B.4 Divergent questions	1011
966	7. Inference questions: Inference questions re-	• What considerations should the [AGENCY]	1012
967	quire learners to use inductive or deductive	addressed in the document but didn’t?	1013
968	reasoning to eliminate responses or critically	B.5 Evaluation questions	1014
969	assess a statement.	• Based on NEPA evaluations done in the vicin-	1015
970	8. Problem-solving questions. Problem-	ity of [PROJECT], does the conclusion of the	1016
971	solving questions present students with a sce-	Historical and Cultural resources section ap-	1017
972	nario or problem and require them to develop	propriately weigh the concerns of Tribal lead-	1018
973	a solution.	ers’?	1019
974	9. Process questions: A process question al-		
975	lows the speaker to evaluate the listener’s		
976	knowledge in more detail.		
977	10. Recall questions: A recall question asks the		
978	listener to recall a specific fact.		

1020	• Extrapolating using this and other NEPA evaluations, what is the long term outlook for [SPECIES] in the vicinity?	1064
1021		1065
1022		1066
1023	• How have [AGENCY’S] NEPA reviews trended over time and would this review have the same outcome 10 years ago or 10 years from now?	1067
1024		1068
1025		
1026		
1027	• In the license renewal EIS for [PROJECT], which impacts have changed from the initial EIS and why?	
1028		
1029		
1030	B.6 Funnel questions	
1031	• Which federally recognized Tribes are in a 50-mile radius of [PROJECT]? Which Tribes participated in this EIS? What were the concerns fo participating Tribes? What mitigations were made?	
1032		
1033		
1034		
1035		
1036	• Which federally recognized species of concern are in a 50-mile radius of [PROJECT]? What mitigations, if any, were made to project those species?	
1037		
1038		
1039		
1040	• Which alternatives were discussed? Which were considered? Why was [ALTERNATIVE] not considered?	
1041		
1042		
1043	• Which resource areas were discussed in the Affected Environment section of the document?	
1044		
1045		
1046	• What were the impacts of the proposed action on [SUBJECT]?	
1047		
1048	• Did [AGENCY] consider [X] when evaluating [SUBJECT]?	
1049		
1050	B.7 Inference questions	
1051	• If the federally recognized [TRIBE] has land in the vicinity of [PROJECT 1] like it does in the vicinity of [PROJECT 2], what concerns might [TRIBE] have with [PROJECT 1]?	
1052		
1053		
1054		
1055	• If the primary mitigation for [SPECIES] for [PROJECT TYPE] in the past has been [MITIGATION], what would you expect the mitigation to be for [PROJECT]?	
1056		
1057		
1058		
1059	• If [AGENCY 1] and [AGENCY 2] typically agree on impact levels and [AGENCY 1] found large impact in terrestrial ecology for an action in a nearby area, what would [AGENCY 2] find?	
1060		
1061		
1062		
1063		
	• If mitigations for air quality for [PROJECT 1] were effective and the same mitigations were applied to [PROJECT 2], what would we assume the outcome to be for [PROJECT 2]?	1064
		1065
		1066
		1067
		1068
	B.8 Problem-solving questions	1069
	• Given the following references, evaluate the effect of a new nuclear plant at [SITE] on cultural and historic resources in the vicinity.	1070
		1071
		1072
	• Given the location of the [PROJECT], create a list of aquatic species likely present in a 50-mile radius.	1073
		1074
		1075
	• Write an Abstract for [PROJECT]	1076
	• Given the list of reference in [SECTION] of [PROJECT 1] create a list of references applicable to [PROJECT 2]. Provide hyperlinks and ML numbers, if available.	1077
		1078
		1079
		1080
	B.9 Process questions	1081
	• How does this document define the NEPA process for consultation with Tribes?	1082
		1083
	• How does [AGENCY] define the area of effect for the proposed action?	1084
		1085
	B.10 Recall questions	1086
	• What references did [AGENCY] use in evaluating the effect of the applicant’s proposed action on [SPECIES]?	1087
		1088
		1089
	• Which resource areas indicated a moderate or large impact due to the proposed action?	1090
		1091
	C EIS Dataset	1092
	Table 3 reports the statistics of the EIS data that used to create the benchmark.	1093
		1094

Document Title	Agency	#Pages	#Tokens
Continental United States Interceptor Site	Missile Defense Agency, Department of Defense	74	41,742
Supplement Analysis of the Final Tank Closure and Waste Management for the Hanford Site, Richland, Washington, Offsite Secondary Waste Treatment and Disposal	Hanford Site Office, Department of Energy	63	43,167
Nationwide Public Safety Broadband Network Final Programmatic Environmental Impact Statement for the Southern United States	Department of Commerce	86	43,985
T-7A Recapitalization at Columbus Air Force Base, Mississippi	United States Department of the Air Force (DAF), Air Education and Training Command (AETC).	472	179,697
Oil and Gas Decommissioning Activities on the Pacific Outer Continental Shelf	The Bureau of Safety and Environmental Enforcement (BSEE) and Bureau of Ocean Energy Management (BOEM)	404	271,545
Final Environmental Impact Statement for the Land Management Plan Tonto National Forest	Department of Agriculture, Forest Service	472	325,641
Final Environmental Impact Statement for Nevada Gold Mines LLC's Goldrush Mine Project, Lander and Eureka Counties, NV	Bureau of Land Management, Interior.	454	413,083
Addressing Heat and Electrical Upgrades at Fort Wainwright, Alaska	Department of the Army, Department of Defense	618	514,003
Sea Port Oil Terminal Deepwater Port Project	The U.S. Coast Guard (USCG) and Maritime Administration (MARAD), Department of Transportation	890	613,214

Table 3: Statistics on the EIS documents used in the evaluation