

Multi-Stage Retrieval for Operational Technology Cybersecurity Compliance Using Large Language Models: A Railway Casestudy

R. Bolton, M. Sheikhfathollahi, S. Parkinson, D. Basher, and H. Parkinson

Abstract—Operational Technology Cybersecurity (OTCS) continues to be a dominant challenge for critical infrastructure such as railways. As these systems become increasingly vulnerable to malicious attacks due to digitalization, effective documentation and compliance processes are essential to protect these safety-critical systems. This paper proposes a novel system that leverages Large Language Models (LLMs) and multi-stage retrieval to enhance the compliance verification process against standards like IEC 62443 and the rail-specific IEC 63452. We first evaluate a Baseline Compliance Architecture (BCA) for answering OTCS compliance queries, then develop an extended approach called Parallel Compliance Architecture (PCA) that incorporates additional context from regulatory standards. Through empirical evaluation comparing OpenAI-gpt-4o and Claude-3.5-haiku models in these architectures, we demonstrate that the PCA significantly improves both correctness and reasoning quality in compliance verification. Our research establishes metrics for response correctness, logical reasoning, and hallucination detection, highlighting the strengths and limitations of using LLMs for compliance verification in railway cybersecurity. The results suggest that retrieval-augmented approaches can significantly improve the efficiency and accuracy of compliance assessments, particularly valuable in an industry facing a shortage of cybersecurity expertise.

Impact Statement—This research addresses a critical gap in railway cybersecurity by introducing an AI-assisted compliance verification system. As critical infrastructure becomes increasingly digitized, railways face growing cybersecurity threats that could compromise safety and operations. Manual compliance verification is time-consuming, resource-intensive, and requires specialized expertise that is increasingly scarce in the industry. Our multi-stage retrieval system using Large Language Models demonstrates significant improvements in compliance assessment quality, with our Parallel Compliance Architecture (PCA) achieving 0.1 points higher correctness scores and 0.55 points higher reasoning scores than the baseline approach as seen in Table II. These improvements could substantially reduce verification time.

Index Terms—Operational Technology, Cybersecurity, Large

Submitted on 17/04/2025

R. Bolton is with Digital Transit Limited, 3M Buckley Innovation Centre, Huddersfield, HD1 3BD, West Yorkshire, UK (e-mail: regan.bolton@digitaltransit.co.uk).

M. Sheikhfathollahi is with the Department of Computer Science at University of Huddersfield, University of Huddersfield, Huddersfield, HD1 3DH, West Yorkshire, UK (e-mail: mohammadreza.sheikhfathollahi@hud.ac.uk).

S. Parkinson is with the Department of Computer Science at University of Huddersfield, University of Huddersfield, Huddersfield, HD1 3DH, West Yorkshire, UK (e-mail: s.parkinson@hud.ac.uk).

D. Basher is with Digital Transit Limited, 3M Buckley Innovation Centre, Huddersfield, HD1 3BD, West Yorkshire, UK (e-mail: dan.basher@digitaltransit.co.uk).

H. Parkinson is with Digital Transit Limited, 3M Buckley Innovation Centre, Huddersfield, HD1 3BD, West Yorkshire, UK (e-mail: hjparkinson@digitaltransit.co.uk).

Language Models, Retrieval-Augmented Generation

I. INTRODUCTION

CYBERSECURITY in Operational Technology (OT) remains a significant challenge, especially in critical infrastructure such as railways, where digital systems are used to monitor and control operations [1]. As the rail industry increasingly embraces digitalisation [2], the attack surface has expanded, leading to a rise in cybersecurity threats [3]. Without proper management of OT cybersecurity (OTCS), safety-critical rail systems can be compromised, leading to operational failures and safety hazards [4]. Until recently, rail systems operated largely within restricted environments, isolated from wide-area network communication. However, recent technological advances, combined with an increasing demand for operational data, have led to a rapid expansion in connectivity. As a result, many safety critical systems are now exposed to networked environments for which they were not originally designed. These systems, often developed decades ago with lifespans of 30 to 40 years, did not incorporate cybersecurity considerations, exposing them to new attack vectors that were outside of the scope of their initial design and deployment. This challenge is not unique to the rail sector; however, rail has been shown to be behind other industries in cybersecurity, such as aviation, making it especially vulnerable to attack [5]. Furthermore, due to the age of some systems, the experience required to adequately assess the deployed technology may no longer be available. Therefore, there is a need to assist individuals in the absence of specialist knowledge in performing compliance exercises, which involve reviewing documentation against standards to ensure it is secure.

Effective cybersecurity documentation, development, and review processes are essential parts of the compliance process and are required to proactively protect these systems from malicious attacks. The importance of cybersecurity has resulted in dedicated international standards to ensure that the rail industry is adequately protected. Failure to adhere to such cybersecurity standards can result in significant vulnerabilities that can be exploited by an adversary. Therefore, adherence to the OT standards of IEC 62443 [6] and the anticipated rail-specific IEC 63452 [7] is crucial. These standards outline the necessary processes to achieve compliance, ensuring that organisational cybersecurity processes are established from the outset and that effective cybersecurity activities are carried out throughout the development lifecycle of any railway asset.

Compliance with these standards is a considerable challenge due to their complexity and specialised subject knowledge, making it increasingly difficult to meet regulatory requirements [8]. The process of analysing the compliance of cybersecurity documentation with regulations and standards is typically performed by someone who has extensive knowledge of the domain and familiarity with the associated standards. In general, cybersecurity assessment methods are based on testing, examination, or interview [9]. Examination can be difficult due to the inherent ambiguity in the OTCS documentation. This complexity adds to the time and cost involved in conducting a thorough compliance assessment. These issues contribute to the motivation for the presented research.

The challenges of maintaining compliance are exacerbated as any change to a train’s digital systems will have an impact on the OTCS. This impact will be assessed in terms of both the new vulnerabilities and risks introduced by the new technology and the impact of the change on existing legacy systems. Therefore, it is necessary to thoroughly assess compliance to maintain a high level of security. In addition to the need to maintain compliance with the introduction of new systems, there is also the need to periodically perform an assessment to account for the evolution of cybersecurity standards.

There is clearly a need to assist in compliance exercises; otherwise, it is likely that an organisation might fall behind with their compliance due to the knowledge and time-intensive process. Researchers have previously worked on developing intelligent compliance-checking solutions. However; the work is in different domains, such as construction [10] and software engineering [11]. Many solutions use forms of Artificial Intelligence, most prominently natural language processing (NLP) [12]. These works provide promising results; however, they often lack the flexibility to handle a diverse domain. Large Language Models (LLMs) have recently demonstrated promise in other knowledge-intensive and time-intensive tasks. To the best of the authors’ knowledge, the use of LLMs to improve automated OTCS compliance has not yet been investigated.

LLMs excel at understanding large volumes of text and can generate human-like text based on patterns learnt from their vast amount of training data [13]. A recent example demonstrates their ability in multimode industrial diagnostics [14], student performance assessment [15], and safety case generation [16]. This capability enables them to grasp context, semantics, and nuances in human language effectively. With the introduction of increasingly powerful LLMs such as GPT-4 [17] and Meta’s open-source LLama-3 model [18], these models are becoming increasingly accurate and capable of demonstrating human-like reasoning. During the past year, advances in LLMs have prompted a shift in research focus from generative capabilities to reasoning abilities [19].

In this paper, a system is proposed that utilises LLMs and other advanced techniques to accelerate the compliance verification process. This has resulted in the following novel contributions:

- Development and evaluation of a multi-stage LLM retrieval system designed for compliance verification in the OTCS domain;
- A comparison of correctness and reasoning between the

Baseline Compliance Architecture (BCA) and the two Parallel Compliance Architecture (PCA) retrieval system architectures; and

- Established metrics for correct responses, logical reasoning, and hallucination for both architectures, highlighting their strengths and weaknesses.

The remainder of this paper is structured as follows. Section II describes studies related to the research presented in this paper. Section III presents the methodology, explaining the retrieval system, and discussing the two architectures developed and tested in this paper. Section IV describes the experimental setup. Section V presents the results. Section VI provides an analysis of the observations made from the dataset. Furthermore, this section includes a discussion and the limitations of the study, and Section VII concludes the paper with a conjecture and discusses avenues for future work.

II. RELATED WORK

LLMs are being used in multiple sectors for tasks that involve the analysis of textual artefacts in different ways. In many cases, researchers employ an ensemble of LLM techniques in many domains to enhance performance and accuracy [20], [21]. However, these techniques are not yet fully exploited in the OTCS domain.

One new technique in such research is retrieval augmented generation (RAG) [22] which combines the strengths of traditional information retrieval systems (such as databases) with the capabilities of generative large language models. By combining this additional knowledge with its own language skills, AI can write text that is more accurate, up-to-date, and relevant to the individual’s specific needs [23]. This is powered by a retrieval engine that works by efficiently retrieving informational nodes from an external database, using a retriever, and then incorporating this context into the context window of the LLM.

Research has already addressed compliance in the architecture, engineering and construction industry (AEC) [24]. One such approach focused on evaluating different prompt engineering techniques [25]. Their compliance goal was much simpler than that of this project. More specifically, the authors used pairs of fire safety regulations and building design specifications. Their conclusions determined that the design of prompt engineering largely determined the performance of the LLM. In this research, detailed prompts are constructed to instruct the LLM. Despite their small data set and less complex domain, LLMs showed promise in classifying compliance and non-compliance. The disadvantage of using long context is that it may not perform well with larger specifications. Regulations are likely to target key areas of the design specifications and, as a result, the context may become diluted. This dilution can make it challenging for the system to focus on the most relevant information, potentially reducing the accuracy and effectiveness of compliance verification.

An example, in the medical field, involves an investigation of how compliance can be assessed against reporting guidelines in clinical trials [26]. They manually extracted text-question pairs as their method. The findings showed that the

LLM demonstrated an acceptable classification accuracy of greater than 95% in its compliance evaluation. They attribute this success to the fine-tuning of their models as increased performance is observed when this occurs. Although fine-tuning the model is not performed, a lightweight solution can be applied using a retriever to add further context and enhance the system’s ability to provide correct answers. One conclusion from their paper is that it is likely that an analysis of all the papers (standards) will be required to further improve performance. The presented solution to this issue involves using an additional retriever on all of the user documentation. This enables the control of the size, amount, and overlap of nodes, which means that only specific parts of the documentation are retrieved.

Another study focuses on how assistive technology can be analysed for compliance with product specification standards [27]. Their approach is unique in the way that they use a retriever to trace product specifications to the relevant standards, then in a second stage analyse compliance using an LLM with manually inserted rules for compliance. The benefit of this approach is that it is more specific than other research, enabling the system to focus on the most relevant information. This multi-stage approach to identifying compliance is similar to the presented architecture, except that the system uses RAG to generate a contextualised response, automatically inferring the compliance rules.

Due to the different methods such as prompt engineering, retrieval augmented generation, and fine-tuning, it is likely that the research will eventually involve an ensemble of these methods. This is supported by the fact that complex domains, such as OTCS compliance, will have a wide variety of different types of documentation and complex rules that must be followed. At this preliminary stage, fine-tuning will not be performed as it has a large overhead and can be performed at the end of experimentation. However, other methods will be used. One downside to every technique is that, in one way or another, the solution is not easily generalisable or involves a manual separation or analysis. This research is uniquely applied to rail OTCS; however, the automation of the proposed system means that it can be deployed to any domain, with minor modification.

III. METHODOLOGY

In order to develop a system to use LLMs to assist with compliance verification, the following multi-stage methodology is presented. Several parameters are used in various presented prompt templates, these are explained in Table I. In this section, the concepts of RAG retrieval, OTCS standards, BCA, and PCA are defined and discussed.

A. RAG retrieval

To answer questions about documentation, there are several options. One approach would be to use long context; however, in this paper’s use case, which involves a massive volume of data, this would drastically increase input token usage and reduce precision. LLMs are known to ‘forget’ information when processing large context windows, often

TABLE I: A table of parameters used in prompt templates and their meaning

Parameter	Explanation
<i>query_str</i>	The input question for the <i>input component</i> .
<i>user_docs_str</i>	A string of retrieved-context nodes from the <i>document retriever</i> .
<i>context_str</i>	A string of retrieved-context nodes from the <i>context retriever</i> .

prioritising details at the beginning and end of token sequences while under-representing or overlooking information in the middle [28]. A more efficient solution is to add context to the context window during inference. Instead of manually adding document chunks, a retriever can be leveraged to automatically retrieve relevant document chunks at inference time. RAG retrieval has already proven to be very powerful in improving generation quality by adding additional context [29]. However, it is proposed that RAG can also be used to retrieve relevant information from lengthy OTCS documentation and answer targeted questions about documents.

In the presented RAG model, the retriever selects the most relevant K documents, and the generator uses these documents to create a probability distribution to generate the next output token. Specifically, the input query is denoted as x , and the retrieved document as y , which is used as an additional context to generate the target as z , the RAG retriever model is composed of two core components.

- 1) Retriever $p_\eta(y|x)$, This component, parametrised by η , retrieves a set of relevant documents based on the input query x . This provides a distribution over textual resources and returns the top K documents that are most likely to contain useful information to answer the query.
- 2) Generator $p_\theta(z_i|x, y, z_{1:i-1})$, The generator, parametrised by θ , predicts the next token z_i in the target sequence z . Based on this components, the probability of generating the target z given the query input x and the retrieved document y can be described as:

$$p_{\text{RAG}}(z|x) = \prod_i^N \sum_{y \in \text{top-K}(p(\cdot|x))} p_\eta(y|x) p_\theta(z_i|x, y_i, z_{1:i-1}) \quad (1)$$

The retrieval component $p_\eta(y|x)$ in the presented system uses a hybrid approach, which combines multiple types of retrieval to improve the retrieval process and a vector store containing embeddings from embed-english-v3.0 embedding model is used [30]. This retriever uses the hybrid query mode, enabling it to combine both vector-based similarity (cosine similarity) and keyword-based similarity (BM25) for more precise retrieval results. The retriever first retrieves the top- $K = 10$ documents using a hybrid retrieval approach, These top-10 documents are then ranked using a Cohere re-ranking model, and the top- $K = 2$ documents are selected for further processing. This two-step process ensures a more meaningful alignment with the input query x , leveraging both hybrid

retrieval and advanced re-ranking to provide highly relevant results for the generator.

$$S(x, y) = \alpha \cdot \text{Cosine Similarity}(x, y) + (1 - \alpha) \cdot \text{BM25 Similarity}(x, y) \quad (2)$$

where x is the input query and y is the retrieved document. α ranging from 0 to 1, controls the balance between cosine similarity and BM25. The retriever is parametrised with an alpha value of $\alpha = 0.75$, which controls the balance between the vector search and the textual search. This value $\alpha = 0.75$ has been generally tested and is considered the best practice in optimising retrieval performance.

The retrieved documents y are then used as context by the generator to produce the target sequence z .

To the best of the authors' knowledge, this is a novel approach to combine RAG and OTCS for this use case. Using these methods, the aim is to improve the efficiency of cybersecurity assessments by allowing OTCS assessors to easily be able to query extensive complex OTCS documentation. These questions will be largely based on the latest rail OTCS standard IEC6345 [7]. Considering that relevant information can be spread across multiple documents, it is expected that this system will significantly enhance the quality of the results by effectively retrieving precise relevant information based on the input query.

B. Baseline compliance architecture (BCA)

The presented methodology leverages the ability of a large language model (LLM) to interpret OTCS compliance within the rail domain. When using an LLM to address OTCS-compliance with long context document processing, the responses tend to lack precision and detail. However, by targeting queries with RAG, it is expected to significantly improve the quality of the analysis. Henceforth, the RAG retriever used for this purpose shall be referred to as the *document retriever*. Specifically, the retriever acts on a vector index containing the chunked internal OTCS case study.

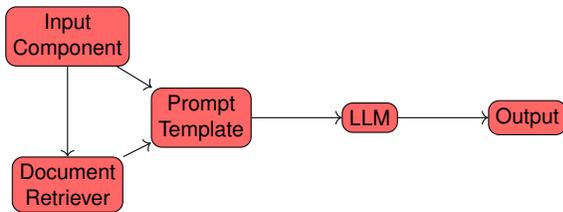


Fig. 1: Flowchart of the basic RAG system architecture.

To use and test this system, the architecture as depicted in Figure 1 is created. The prompt template in Figure ?? merges the *query_str* with the retrieved nodes (*document_retriever*) and adds a system prompt in Figure ?? for the LLM.¹

A detailed system prompt is used to increase the likelihood that the LLM will follow instructions accurately, avoid potential confusion, and generate contextually relevant high-quality

¹Note: Throughout the paper, the new lines in the code block have been modified for readability.

responses. This structured approach improves the likelihood that the model will remain focused on the task and provide accurate answers.

C. Parallel compliance architecture (PCA)

To improve the models understanding of the OTCS domain, one solution could be to perform fine tuning of a model; however, this is time consuming, expensive, and does not integrate easily with retrieval augmented systems. Instead, we define another retriever that uses exactly the same methods as the document retriever, except on a different vector index. This process is named the *context retriever*. The context retriever returns OTCS standards and regulation information. This allows the system to improve the reasoning process and understanding of queries, using additional regulatory context. In this work, the full IEC 62443 (1-1 to 4-2) and IEC 63452 are used as material.

The parallel architecture is depicted in Figure 4. The main difference is the addition of the context retriever and the prompt engineering techniques used to make the LLM understand the additional data.

In summary, the system prompt in Figure ?? instructs the LLM to respond to a question, referencing the user documentation. Furthermore, the *query_str* is rewritten to include the *user_docs_str* and *context_str* as seen in Figure ??.

IV. EXPERIMENT

This section describes the testing setup used in this research and justifies the used evaluation methodology.

A. Implementation details

Throughout the tests, GPT-4o [31] and Claude-3.5-Haiku [32] are used to generate responses,[18] with the Llama-3.1-405B-Instruct model used as an evaluator. All models have a max_tokens setting of 2048. All models are accessed using an API, provided by Openrouter [33]. Similarly, the Cohere embed-english-v3.0 model is used for the embeddings [30].

LlamaIndex [34] for Python [35] is used to connect custom data sources to the LLMs. The nodes are created with a chunk size of 1024 tokens and a chunk overlap of 20 tokens in both vector indexes described in Section III. The retrieval mechanism described in Section III-A is used for both retrievers, context, and documentation. LlamaParse, a genAI-native parsing platform, is used to preprocess OTCS documentation, including user documents and standards. Built for LLM use cases, LlamaParse ensures high-quality data through advanced features like table extraction and multi-file-type support, optimising the retrieval process [36]. This is essential for complex OTCS documentation which can be very large and include tables for things like requirements.

B. Dataset description

To test the approach, queries from a dataset of 44 questions covering various aspects of OTCS compliance are used, primarily using the model GPT-4o to analyse and respond to the questions (unless specified). The questions were designed

```

You are an AI assistant specialized in reviewing documentation.
Your primary task is to perform an expert analysis on the User Documentation.
Do NOT use any prior knowledge.
Your analysis should be detailed and based directly on evidence from the
User Documentation.

```

Fig. 2: System prompt for the BCA

```

You will be provided with some documentation.

===== User Documentation =====
{user_docs_str}
=====

Based solely on the User Documentation, please answer the
following Question.

Question: {query_str}
Important Guidelines:
- Do NOT use any prior knowledge or external information.
Your response must be in the following format:
- First provide step-by-step reasoning on how to answer the Question
- Then provide a summary of how you reached your answer.

```

Fig. 3: User prompt for the BCA

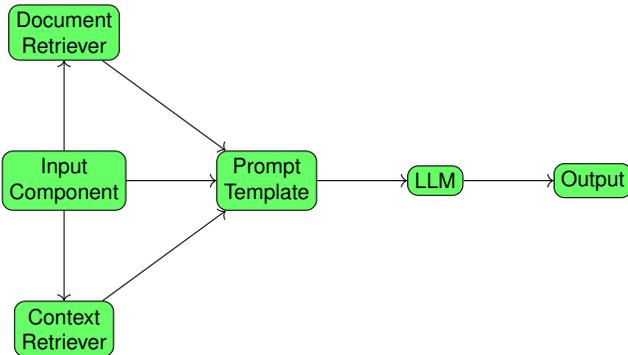


Fig. 4: Flowchart of the parallel system architecture.

with varying levels of difficulty and diversity to challenge the system. Primarily, the questions were taken from the new IEC63452 standard [7]. The expected answers included both compliant, non-compliant and partially compliant scenarios.

C. Evaluation methodology

The evaluation method throughout the experimentation consisted of a mixture of human- and LLM-based evaluations. Given the complexity of the OTCS domain, it is necessary to involve a human expert in OTCS. Recent work in legal analysis [37] and security operation centre analysis [38] both use expert evaluation. Building on this, in the presented approach, the expert must decide whether the answer is satisfactory by responding with 'Correct,' 'Not Correct,' or 'Partially Correct.' The expert should also provide a rationale for their decision.

Additionally, the clarity of the reasoning in the answers is assessed by responding with 'Strongest,' 'Weakest,' or 'Moderate', evaluated comparatively between the three tested architectures, to ensure a comprehensive evaluation of both accuracy and coherence between them. The human evaluation methodology used aligns with the research's aim; specifically, the aim is to increase the confidence in compliance verification by improving reasoning and correctness as much as possible.

For LLM-based evaluation, a technique called LLM-as-a-judge [39] is used, which prompts an LLM to assess another LLM's response based on predefined criteria. In this work, the focus is on evaluating correctness and reasoning ability while also analysing the impact of retrieved context on response quality. The evaluation was carried out using the Arize-Phoenix tool [40]. Recent research has demonstrated that LLMs tend to exhibit bias such as self-recognition and self-preference, when evaluating their own responses [41]. Furthermore, evidence suggests that larger models outperform smaller ones in evaluation tasks [42], prompting us to use Llama-3.1-405B-Instruct as the evaluation model [18].

In the presented evaluation process, the focus is on hallucination detection and identifying instances where the model generated factually incorrect or hallucinated information. This is motivated by the research aim of increasing the confidence of LLM systems. Other LLM-based evaluation was considered; however, it required expert subject knowledge or was too subjective to be reliably accurate.

You are an AI assistant specialized in reviewing documentation based on the provided User Documentation.

Your primary task is to perform an expert analysis on the **User Documentation** using the provided **Contextual Information** to enhance your analysis where appropriate and necessary. **Do NOT** use any prior knowledge or perform your analysis directly on the **Contextual Information**; it is provided **ONLY** to help you understand the **Question** and enhance your reasoning capabilities.

Your analysis should be detailed and based directly on evidence from the **User Documentation**.

Fig. 5: System prompt for the PCA

You will be provided with some documentation and supporting context:

```
===== User Documentation =====
{user_docs_str}
=====
```

```
----- Contextual Information -----
{context_str}
-----
```

Based **solely** on the **User Documentation** and by enhancing your analysis utilising the **Contextual Information** please answer the following question.

Question: {query_str}

Important Guidelines:

- **Do NOT** use any prior knowledge or external information.
- **Do NOT** perform an analysis of the **Contextual Information** in your answer.

Your response **must** be in the following format:

- First Provide step-by-step reasoning on how to answer the **Question**, potentially making use of the **Contextual Information** to refine your steps.
- Then provide a summary of how you reached your answer.

Fig. 6: User prompt for the PCA

V. RESULTS

A. BCA results

After generating responses from the dataset, the hallucination evaluation for the BCA is presented in Figure 7. Of the 44 questions tested, only one (2.3%) exhibited a hallucinated response. This is a positive outcome, emphasising the effectiveness of the retrieval system in providing highly relevant and accurate information to guide the model’s responses, thereby minimising the likelihood of errors commonly associated with LLMs.

In addition, a human evaluation was conducted to further assess the quality of responses. The results of this evaluation,

shown in Figure 8, provide insight into how cybersecurity experts classified the correctness of the model. This evaluation complements the automatic evaluation by offering a detailed, expert-driven perspective on the model’s performance.

As shown in Figure 8-left, out of the 44 questions in the dataset, 19 were classified as correct, 17 as partially correct, and 6 as incorrect. This distribution suggests that the system can accurately answer just under half of the questions through user documentation retrieval alone, with a significant portion requiring further refinement and a smaller fraction being entirely incorrect. The “partially correct” responses often contained elements of accuracy, but fell short in terms of

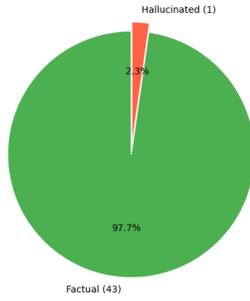


Fig. 7: Evaluation of BCA Hallucination by LLM: Factual vs. Hallucinated Answers, Using the Llama-3.1-405B-Instruct model.

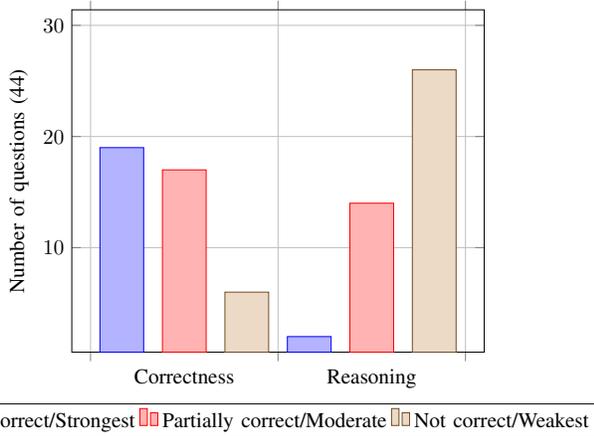


Fig. 8: Results of the human evaluation for BCA: correctness (left) and reasoning (right)

completeness or precision in addressing the specific nuances of the questions. For example, some answers provided general information relevant to OTCS, but did not capture the exact technical details or contextual specifics demanded by the question. Conversely, the “incorrect” responses were either factually inaccurate or entirely unrelated to the query, highlighting areas where the model’s comprehension or retrieval mechanisms require significant improvement.

Figure 8-right, shows the reasoning evaluation, where, of the same 44 questions, only 2 exhibited the strongest reasoning, 14 were moderate and 26 were weakest, highlighting a challenge in constructing well-reasoned, logically coherent responses despite retrieving relevant information. Although the system excels at avoiding hallucinations and providing relevant OTCS data, it requires additional context to improve accuracy, completeness, and reasoning quality.

B. PCA experiment

Although the BCA performs well in avoiding hallucinations and providing relevant information, it still requires more context on OTCS to improve correctness and reasoning. Based on this, the performance of the PCA is evaluated using the same 44 questions as the BCA. To determine whether the choice of LLM model significantly affects performance, a control experiment was conducted using Claude 3.5 Haiku

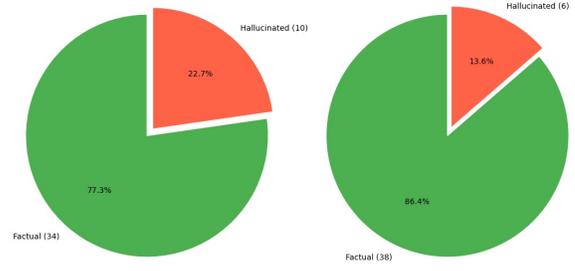


Fig. 9: Evaluation of PCA Hallucination by LLM: Factual vs. Hallucinated Answers, Using the Llama-3.1-405B-Instruct Model for evaluation and GPT-4o as the response model (left) and Claude-3.5-Haiku as the response model (right).

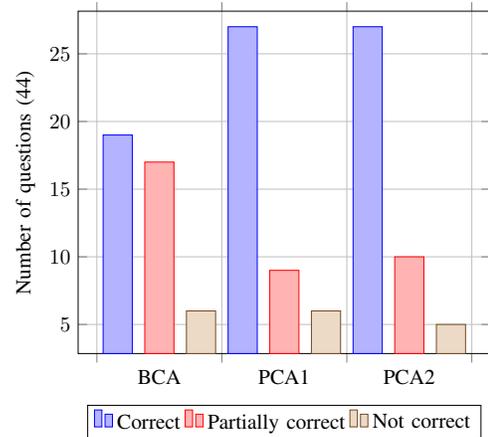


Fig. 10: Results of the human evaluation on correctness for BCA (left) using GPT-4o, PCA1 (middle) using GPT-4o, and PCA2 (right) using Claude-3.5-Haiku.

as an alternative provider to GPT-4o.

The results of the hallucination evaluation are illustrated in Figure 9, which provides a visual representation of the performance of the model in terms of hallucination. The figure demonstrates that the hallucination rate remains relatively consistent across both models, suggesting that the effectiveness of the retrieval system plays a crucial role in maintaining response quality. This consistency can be attributed to the robust performance of the retrieval system, which ensures that high-quality, relevant information is retrieved for processing by the models, as the quality of retrieved data plays a critical role in shaping the output. The increase in “hallucinated” answers from the BCA can be explained by the additional context that is used in the PCA and the necessity of using this context directly in the answer. It is important to note that the context often won’t appear directly in the response and as a result the entry is marked as “hallucinated” when in fact it has aided the understanding of the query for the LLM without direct reference in the response. That being said it is still useful to compare the hallucination between the two PCA despite this metric limitation.

Human evaluation was also performed to gain deeper insight

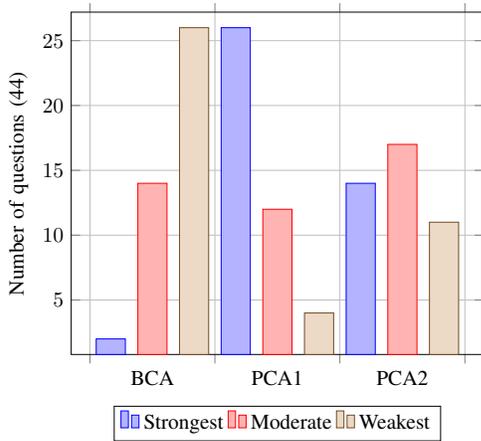


Fig. 11: Results of the human evaluation on reasoning for BCA (left) using GPT-4o, PCA1 (middle) using GPT-4o, and PCA2 (right) using Claude-3.5-Haiku.

into the quality of the models’ responses. The findings of this evaluation, carried out by a cybersecurity expert, are shown in Figure 10. The responses of PCA1 and PCA2 are nearly identical, confirming the effectiveness of the retrieval mechanism. Both architectures correctly answered 27 of the 44 questions, demonstrating similar performance. PCA1 provided partially correct answers for 9 questions, while PCA2 did so for 10 questions. Furthermore, PCA1 and PCA2 incorrectly answered 6 and 5 questions, respectively.

Compared to BCA, the PCA demonstrated a significant increase in the number of completely correct responses and a decrease in the number of partially correct responses. This indicates a shift toward more accurate and definitive answers. A key reason for this improvement is that the PCA uses both retrievers, allowing the system to synthesise responses more effectively. Unlike the BCA, which relies primarily on direct retrieval from user documents, the PCA incorporates context to reach a deeper understanding, resulting in more coherent and well-reasoned responses. Interestingly, we can observe that the number of incorrect responses are relatively similar between all the models; this demonstrates that when document retrieval fails, no amount of additional context can help the model give a more correct response. This shows that document retrieval is likely a bottleneck of the system.

As shown in Figure 11, the BCA performs significantly worse than PCA1 and PCA2 in terms of reasoning. The reasoning in PCA1, which uses GPT-4o, demonstrates the strongest performance. Among the 44 questions, PCA1 had the strongest reasoning in 26 cases, while PCA2 did so in 14. PCA1 provided moderate reasoning for 12 questions, compared to 17 for PCA2. Furthermore, PCA1 exhibited the weakest reasoning in 4 instances, whereas PCA2 did so in 11. These results demonstrate that GPT-4o is better fine-tuned for the task, in terms of providing good reasoning compared to that of Claude 3.5-Haiku. This is intuitive as “smarter” models such as 4o can likely produce more convincing reasoning.

TABLE II: LLM and Human Evaluation Results

Arch.	Model	Hall	Correctness	Reasoning
BCA	GPT-4o	0.023	0.65	0.21
PCA1	GPT-4o	0.227	0.75	0.76
PCA2	Claude-3.5	0.136	0.77	0.53

Note: The hallucinated percentages were calculated based on LLM evaluations using Llama3.1-405b as the evaluator. Correctness scores (0-1) are calculated by assigning 1 for correct, 0.5 for partially correct, and 0 for not correct responses, then averaging. Reasoning scores (0-1) are calculated by assigning 1 for strongest, 0.5 for modest, and 0 for weakest reasoning, then averaging. The headings Arch. and Hall. refer to architecture and hallucinated answers, respectively.

VI. ANALYSIS & OBSERVATIONS

Table II shows our evaluation of compliance architectures (BCA and PCAs). The data show that PCA variants, in general, provide more correct responses. This increase suggests that the PCA method works better for OTCS compliance assessment in an area where there are few automated solutions.

The hallucination rates between architectures need careful reading, as different retrieval methods affect these measurements, PCA and BCA cannot be compared in this metric. Our analysis found several key differences between the tested systems.

Both PCA experiments used more accurate technical terms compared to the more general language seen with BCA. The context information often enhanced the final answers, especially when compliance criteria appeared in the provided context. This contributes to the higher correctness scores in PCA variants.

Comparison of PCA1 and PCA2 showed different behaviours. PCA2 (Claude-3.5) used very brief reasoning steps, sometimes resulting in less detailed reasoning. Yet, this brevity sometimes helped by reducing over-thinking on simple questions and leaving less room for hallucination. The choice of LLM clearly affects both the reasoning process and the final answers.

PCA1 had some drawbacks, mainly focussing too much on specific query words and tending to be stricter and more negative. This version was less willing to make connections without clear evidence, a problem that could be fixed by testing more specific queries.

A repeated issue in PCA systems was sometimes mixing up context information with document content. Though rare, these cases usually included direct mention of the context material before deciding it was not relevant to the question. This is an odd behaviour given the direct instructions in the prompt to avoid this confusion. This highlights the need for better prompt refinement to define clearer use of the different types nodes, perhaps by providing an example. This would ultimately reduce the ambiguity in the LLM’s understanding of the prompt.

The key takeaway from our results and analysis is that the retrieval of the correct document chunks plays the largest part in proper reasoning and general correctness. This is particularly notable in longer queries that require a lot of information from the document retriever. This also encourages

hallucination in the response. By improving the retriever’s ability to filter irrelevant documents or by incorporating domain-specific embeddings, this could significantly reduce errors and improve system reliability. Since retrieval is essential to obtain accurate and relevant information, optimising this component is essential to improve overall performance.

A. Limitations

The case study used in the compliance verification analysis only partially represents what would be expected in a comprehensive suite of OTCS documentation. Specifically, it includes only an initial risk assessment, security requirements, and a definition of the initial zoning of a SuC. Furthermore, the remainder of the case study merely refers to other complementary documents which are presumed to exist. This limitation is due in part to the scarcity of detailed OTCS documentation within the rail industry and the reluctance to share sensitive cybersecurity data [43]. In addition, the document used in the presented research is significantly smaller in size than typical real-world compliance documents, where the method excels in managing large datasets. Despite these challenges, the case study serves as a valuable starting point, highlighting key OTCS issues and allowing a suitable analysis that can be inferred from the available information. As a result of the chosen case study, compliance-related questions tended to focus on the existing sections that are present in the case study.

Another limitation is the method used to analyse LLM responses. Since compliance verification is not simply a matter of checking whether something exists or not, the system must provide sufficient reasoning to extract actionable insights from OTCS documentation. This means that the involvement of a human expert in the loop is essential, although the analysis may be subjective and opinion-based. Despite the existence of LLM-based correctness metrics, it is deemed that they are inappropriate for this task, as the evaluation requires domain-specific knowledge for a detailed performance analysis, beyond just detecting hallucinations. It is suspected that automated systems cannot yet replicate the nuanced understanding and expertise needed to assess OTCS compliance verification systems.

Although automated evaluation methods are much faster and more cost effective, they are likely to be less accurate than human evaluation of LLM responses, particularly when it comes to identifying irrelevant information. LLM-based metrics may struggle to detect when a chunk of text, although factually correct, is not relevant to answering the question at hand. This is especially true in the OTCS domain, where nuanced understanding is crucial.

VII. CONCLUSIONS

This study evaluated the performance of two compliance assessment architectures, the BCA and the PCA, in the Operational Technology Cybersecurity (OTCS) domain, as described in Section III-B and Section III-C, by implementing a multi-stage retrieval system powered by Retrieval-Augmented Generation (RAG). The presented findings show that the retrieval system is a bottleneck in this process and plays a crucial role

in ensuring accurate responses from LLMs, as described in Section V-B. Although both architectures demonstrate promising results, they exhibit different strengths and weaknesses in handling compliance-related queries. The presented evaluation methodology combined LLM-based assessment (LLM-as-a-Judge) with expert human evaluations to ensure a comprehensive analysis.

The study highlights the need to optimise retrieval processes, improve prompt engineering, and explore model fine-tuning to improve the accuracy of compliance verification. In addition, integrating more comprehensive cybersecurity documentation and case studies will further validate the effectiveness of the presented approach. As LLMs continue to evolve, their potential to automate and improve cybersecurity compliance verification remains promising, paving the way for more intelligent, scalable, and reliable assessment systems.

Despite these findings, there is significant room for further experimentation and development. Given that the LLM in the presented approach is not fine-tuned and relies solely on its pre-trained knowledge, the reasoning capabilities are impressive, which approach near-human expert-like judgment using RAG. The proposed solution remains lightweight due to the lack of fine-tuning and is highly adaptable for various domains by simply augmenting the retrieved data. This flexibility enhances its practicality and scalability, making it a versatile tool for compliance verification and other complex analytical tasks beyond OTCS.

In terms of future work, optimising retrieval parameters such as the chunk size and retrieval amount based on the complexity of the query could improve performance. Furthermore, experimenting with agentic RAG to decide when sufficient document/context chunks have been retrieved to answer the query would be an interesting extension to the research.

REFERENCES

- [1] R. Kour, A. Patwardhan, A. Thaduri, and R. Karim, “A review on cybersecurity in railways,” *Proceedings of the Institution of Mechanical Engineers, Part F: Journal of Rail and Rapid Transit*, vol. 237, no. 1, pp. 3–20, 2023.
- [2] V. Jägare, R. Karim, P. Söderholm, P.-O. Larsson-Kräik, and U. Juntti, “Change management in digitalised operation and maintenance of railway,” in *International Heavy Haul Association (IHHA) STS 2019, 10-14th June 2019, Narvik, Norway.*, 2019, pp. 904–911.
- [3] M. Rekek, C. Gransart, and M. Berbineau, “Cyber-physical security risk assessment for train control and monitoring systems,” in *2018 IEEE Conference on Communications and Network Security (CNS)*. IEEE, 2018, pp. 1–9.
- [4] G. Sabaliauskaite and A. P. Mathur, “Aligning cyber-physical system safety and security,” in *Complex Systems Design & Management Asia: Designing Smart Cities: Proceedings of the First Asia-Pacific Conference on Complex Systems Design & Management, CSD&M Asia 2014*. Springer, 2015, pp. 41–53.
- [5] R. Kour, M. Aljumaili, R. Karim, and P. Tretten, “emaintenance in railways: Issues and challenges in cybersecurity,” *Proceedings of the Institution of Mechanical Engineers, Part F: Journal of Rail and Rapid Transit*, vol. 233, no. 10, pp. 1012–1022, 2019.
- [6] “Iec62443 suite of standards,” 2024. [Online]. Available: <https://www.isa.org/standards-and-publications/isa-standards/isa-iec-62443-series-of-standards>
- [7] BSI Group, “BS EN IEC 63452 Ed.1.0 Railway applications - Cybersecurity,” 2024. [Online]. Available: <https://standardsdevelopment.bsigroup.com/projects/2022-01003/section>
- [8] A. T. Tunggal, “Best practices for cybersecurity compliance monitoring,” 2024, updated April 21, 2024. [Online]. Available: <https://www.upguard.com/blog/compliance-monitoring>

- [9] J. M. Stewart, E. Tittel, and M. Chapple, *CISSP: Certified information systems security professional study guide*. John Wiley & Sons, 2011.
- [10] R. Amor and J. Dimyadi, "The promise of automated compliance checking," *Developments in the built environment*, vol. 5, p. 100039, 2021.
- [11] S. Malsane, J. Matthews, S. Lockley, P. E. Love, and D. Greenwood, "Development of an object model for automated compliance checking," *Automation in construction*, vol. 49, pp. 51–58, 2015.
- [12] J. Zhang and N. M. El-Gohary, "Semantic nlp-based information extraction from construction regulatory documents for automated compliance checking," *Journal of computing in civil engineering*, vol. 30, no. 2, p. 04015014, 2016.
- [13] S. Hore, "An introduction to large language models (llms)," 2024. [Online]. Available: <https://www.analyticsvidhya.com/blog/2023/03/an-introduction-to-large-language-models-llms/>
- [14] S. Jose, K. T. Nguyen, K. Medjaher, R. Zemouri, M. Lévesque, and A. Tahan, "Advancing multimodal diagnostics: Integrating industrial textual data and domain knowledge with large language models," *Expert Systems with Applications*, vol. 255, p. 124603, 2024.
- [15] C. Oh, M. Park, S. Lim, and K. Song, "Language model-guided student performance prediction with multimodal auxiliary information," *Expert Systems with Applications*, vol. 250, p. 123960, 2024.
- [16] M. Sivakumar, A. B. Belle, J. Shan, and K. K. Shahandashti, "Prompting gpt-4 to support automatic safety case generation," *Expert Systems with Applications*, vol. 255, p. 124653, 2024.
- [17] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altschmidt, S. Altman, S. Anadkat *et al.*, "Gpt-4 technical report," *arXiv preprint arXiv:2303.08774*, 2023.
- [18] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan *et al.*, "The llama 3 herd of models," *arXiv preprint arXiv:2407.21783*, 2024.
- [19] Y. Zhang, S. Mao, T. Ge, X. Wang, A. de Wynter, Y. Xia, W. Wu, T. Song, M. Lan, and F. Wei, "Llm as a mastermind: A survey of strategic reasoning with large language models," *arXiv preprint arXiv:2404.01230*, 2024.
- [20] B. Perak, S. Beliga, and A. Meštrović, "Incorporating dialect understanding into llm using rag and prompt engineering techniques for causal commonsense reasoning," in *Proceedings of the Eleventh Workshop on NLP for Similar Languages, Varieties, and Dialects (VarDial 2024)*, 2024, pp. 220–229.
- [21] F. Bianchini, M. Calamo, F. De Luzi, M. Macrì, and M. Mecella, "Enhancing complex linguistic tasks resolution through fine-tuning llms, rag and knowledge graphs (short paper)," in *International Conference on Advanced Information Systems Engineering*. Springer, 2024, pp. 147–155.
- [22] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel *et al.*, "Retrieval-augmented generation for knowledge-intensive nlp tasks," *Advances in Neural Information Processing Systems*, vol. 33, pp. 9459–9474, 2020.
- [23] G. Cloud, "Retrieval-augmented generation (rag)," 2024, accessed: August 7, 2024. [Online]. Available: <https://cloud.google.com/use-cases/retrieval-augmented-generation?hl=en>
- [24] J. Dimyadi and R. Amor, "Automated building code compliance checking—where is it at," *Proceedings of CIB WBC*, vol. 6, no. 1, 2013.
- [25] X. Liu, H. Li, and X. Zhu, "A gpt-based method of automated compliance checking through prompt engineering," 2023.
- [26] J. G. Wrightson, P. Blazey, K. M. Khan, and C. L. Ardern, "Gpt for rcts?: Using ai to measure adherence to reporting guidelines," *medRxiv*, pp. 2023–12, 2023.
- [27] C. Arora, J. Grundy, L. Puli, and N. Layton, "Towards standards-compliant assistive technology product specifications via llms," *arXiv preprint arXiv:2404.03122*, 2024.
- [28] T. Li, G. Zhang, Q. D. Do, X. Yue, and W. Chen, "Long-context llms struggle with long in-context learning," *arXiv preprint arXiv:2404.02060*, 2024.
- [29] A. Salemi and H. Zamani, "Evaluating retrieval quality in retrieval-augmented generation," in *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2024, pp. 2395–2400.
- [30] Cohere, "Cohere-embed-english-v3.0," <https://huggingface.co/Cohere/Cohere-embed-english-v3.0>, Cohere, 2024, accessed: 2025-02-26.
- [31] K. Wiggers, "Openai debuts gpt-4o'omni' model now powering chatgpt," *TechCrunch*. Retrieved May, vol. 16, p. 2024, 2024.
- [32] M. Rahman, S. Khatoonabadi, A. Abdellatif, and E. Shihab, "Automatic detection of llm-generated code: A case study of claude 3 haiku," *arXiv preprint arXiv:2409.01382*, 2024.
- [33] OpenRouter, "Openrouter: Api for accessing open-source and proprietary llms," <https://openrouter.ai/>, 2023, accessed: 2024-09-20.
- [34] J. Liu, "LlamaIndex," 11 2022. [Online]. Available: https://github.com/jerryliu/llama_index
- [35] G. Van Rossum and F. L. Drake Jr, *Python reference manual*. Centrum voor Wiskunde en Informatica Amsterdam, 1995.
- [36] J. Liu, "Introducing llamacloud and llamaparse - llamaindex - build knowledge assistants over your enterprise data," Feb 2024. [Online]. Available: <https://www.llamaindex.ai/blog/introducing-llamacloud-and-llamaparse-af8cedf9006b>
- [37] I. Cheong, K. Xia, K. K. Feng, Q. Z. Chen, and A. X. Zhang, "(a) i am not a lawyer, but...: Engaging legal experts towards responsible llm policies for legal advice," in *The 2024 ACM Conference on Fairness, Accountability, and Transparency*, 2024, pp. 2454–2469.
- [38] O. Oniagbi, A. Hakkala, and I. Hasanov, "Evaluation of llm agents for the soc tier 1 analyst triage process," 2024.
- [39] L. Zheng, W.-L. Chiang, Y. Sheng, S. Zhuang, Z. Wu, Y. Zhuang, Z. Lin, Z. Li, D. Li, E. Xing *et al.*, "Judging llm-as-a-judge with mt-bench and chatbot arena," *Advances in Neural Information Processing Systems*, vol. 36, pp. 46 595–46 623, 2023.
- [40] A. Al, "Phoenix: Open-source ml observability and performance debugging," <https://github.com/Arize-ai/phoenix>, 2023, accessed: 2024-09-20.
- [41] A. Panicssery, S. R. Bowman, and S. Feng, "Llm evaluators recognize and favor their own generations," *arXiv preprint arXiv:2404.13076*, 2024.
- [42] A. S. Thakur, K. Choudhary, V. S. Ramayapally, S. Vaidyanathan, and D. Hupkes, "Judging the judges: Evaluating alignment and vulnerabilities in llms-as-judges," *arXiv preprint arXiv:2406.12624*, 2024.
- [43] A. Patwardhan, A. Thaduri, and R. Karim, "Distributed ledger for cybersecurity: issues and challenges in railways," *Sustainability*, vol. 13, no. 18, p. 10176, 2021.