

Received 9 September 2024, accepted 24 September 2024, date of publication 26 September 2024,  
date of current version 11 October 2024.

Digital Object Identifier 10.1109/ACCESS.2024.3468914

## RESEARCH ARTICLE

# Applications of LLMs for Generating Cyber Security Exercise Scenarios

MUHAMMAD MUDASSAR YAMIN<sup>1</sup>, EHTESHAM HASHMI<sup>1</sup>,  
MOHIB ULLAH<sup>2</sup>, (Member, IEEE), AND BASEL KATT<sup>1</sup>

<sup>1</sup>Department of Information Security and Communication Technology (IHK), Norwegian University of Science and Technology (NTNU), 2815 Gjøvik, Norway

<sup>2</sup>Department of Computer Science (IDI), Norwegian University of Science and Technology (NTNU), 2815 Gjøvik, Norway

Corresponding author: Ehtesham Hashmi (hashmi.ehtesham@ntnu.no)

This work was supported by Norwegian Research Council through the ASCERT Project 3290.

**ABSTRACT** This study proposes a novel approach leveraging Large Language Models (LLMs) to generate dynamic and complex adaptable cybersecurity exercise scenarios. Motivated by Turing's seminal exploration into machine cognition, which questions the ability of machines to mimic human thought and intelligence. By exploiting the generative potential of LLMs, our methodology simulates a wide range of cyber threats, both known and novel, thereby enhancing cybersecurity training and awareness. This approach transforms the potential for 'hallucination' inherent in LLMs into a potential advantage, enabling the creation of complex exercise scenarios that push the boundaries of traditional cybersecurity training. The innovation lies in the sophisticated application of AI, aiming to advance the preparedness of security professionals against diverse cyber threats. The scenarios generated through this method were subject to meticulous testing and a rigorous evaluation process involving (Generated Pre-Trained Transformer) GPT models and expert review to ensure their realism and applicability. In this paper, we introduce 'CyExec,' a novel approach leveraging GPT to dynamically generate cybersecurity training scenarios. Furthermore, the prompts provided to the LLMs were meticulously designed to adopt a Retrieval-Augmented Generation (RAG) approach, enriching the complexity and relevance of the scenarios. This incorporation of RAG, alongside the inspiration drawn from Turing's exploration of machine intelligence, showcases an advanced application of AI in cybersecurity training, reflecting a deep understanding of how machines can augment our capabilities to anticipate and mitigate cyber threats.

**INDEX TERMS** Cyber security exercise scenarios, large language models, bounded rationality, generative configurations, Hallucination in LLMs.

## I. INTRODUCTION

In the ever-evolving cyber domain, security professionals are continually seeking innovative strategies to prepare for and reduce the impact of complex threats. Traditional methods for creating cyber threat scenarios often fall short due to their lack of adaptability and scalability. These conventional approaches struggle to keep pace with the dynamic nature of cyber threats, resulting in training that does not adequately prepare professionals for real-world challenges. Our research addresses these critical deficiencies by utilizing Large Language Models (LLMs) to dynamically generate and

adapt cyber threat scenarios, ensuring continuous relevance and rigorous training standards. Cyber threats can take various forms, including malware, phishing, ransomware, DDoS attacks, social engineering, zero-day exploits, insider threats, and advanced persistent threats [1], [2]. Adapting to these many threat types is critical for protecting organizations and digital assets in an increasingly linked world [3]. As AI technologies continue to evolve, their dual role in applied sciences becomes increasingly significant, offering both immense benefits and posing ethical challenges. Albahri et al. [4] discuss the paradoxical nature of generative AI in applied sciences, highlighting its transformative impact across various fields while also acknowledging the critical need for stringent ethical considerations. Zahid et al. [5]

The associate editor coordinating the review of this manuscript and approving it for publication was Sedat Akleyek<sup>1</sup>.

explore the nuances of LLM training methodologies and their implications, emphasizing the role of Reinforcement Learning from Human Feedback (RLHF) in enhancing the models' adaptability and responsiveness to human input, which is pivotal in generating reliable and contextually accurate cyber security scenarios. In the realm of cybersecurity training, the introduction of innovative methods such as our proposed approach, 'CyExec,' marks a significant advancement.

Developing cyber threat exercise scenarios is essential for the planning process as it functions as both a training ground for security professionals and a testing ground for system defenses. Constructing these scenarios, which require both authenticity and complexity, is a process that necessitates a detailed understanding of the threat environment as well as the capacity to predict the activities of possible adversaries. With organizations generating significant volumes of data daily, there's an increasing risk that this data could be used for malicious activities, effectively turning seemingly harmless information into a weapon for potentially breaching security. Given this context, there is an urgent requirement to address novel approaches that can convert this data into valuable training scenarios, thereby transforming potential weaknesses into strong defenses [6], [7]. Within this context, LLMs emerge as transformative tools, offering the ability to generate detailed and intricate content that can simulate potential cyber threat scenarios. Given their generative abilities, LLMs can go beyond simulating known cyber threat scenarios and have the potential to generate different types of novel and sophisticated cyber threats, thus expanding the scope of cybersecurity awareness and training.

In this paper, we propose an innovative use of LLM hallucination to create a sophisticated and adaptive cyber threat scenario. Our approach transforms this hallucination into a useful tool for generating diverse cyber threat exercise scenarios. We establish a pipeline involving multiple LLMs that interactively communicate, setting up a dynamic environment for cyber-attack scenarios. In this setup, one LLM adopts the role of a 'Cyber Security Expert', while the other acts as a 'CISO of an Organization'. Uniquely, the 'CISO' is aware of documents containing various organizational functions like finance, HR, and IT infrastructure while the 'Cyber Security Expert' is aware of the current cyber threat landscape. The documents serve as a prompt for the 'CISO' to utilize its hallucination capability in generating different types of cyber threats. This method not only capitalizes on the often-overlooked hallucinatory aspect of LLMs but also elevates it as a pivotal feature in crafting realistic, adaptable, and comprehensive cyber threat scenarios. This novel application of LLM hallucination in cybersecurity represents a significant advancement, transforming a perceived limitation into a vital asset for strengthening digital security infrastructures, just like in Alan Turing's "Imitation Game" where we have three characters an integrator and two individuals discussing a subject to identify who is a machine and who is a human. In our case we have predefined prompts acting as interrogators while two LLMs discussing the threat

landscape for creating a cyber security scenario. We identified that our approach is suitable for generating new and unique cybersecurity exercise scenarios and can have other potential use cases. The structure of the remainder of this paper is organized as follows: Section (II) reviews the cyber security scenarios. Section (III) details the scenario design process. Section (IV) is dedicated to the proposed system design. Section (V) is based on the case study. In the concluding phase, Section (VI) concludes the paper and outlines future work.

## II. BACKGROUND AND RELATED WORK

### A. CYBER SECURITY SCENARIOS

There are two primary types of cybersecurity exercises: tabletop discussions and hands-on operation-based exercises [8], [9], [10]. Tabletop exercises involve discussions and are often conducted through seminars, workshops, and idea exchanges, primarily focusing on policy-related issues. In comparison, operation-based cybersecurity training programs emphasize practical, hands-on activities aimed at enhancing participants' technical skills and abilities. These exercises take place in simulated, emulated, physical, or hybrid practice environments. Recent studies have found that the creation of these practice environments, known as cyber-ranges, often relies on manual setup and configuration. This approach is not only ineffective but also prone to errors, making it a tedious process [11]. Currently, the development of cyber ranges predominantly involves a manual setup, which includes configuring network environments, deploying simulated threats, and creating training scenarios [12]. This process often requires extensive human intervention to design and implement each scenario, which can be both time-consuming and prone to errors. The manual nature of these setups does not easily scale up to handle larger or more complex training scenarios, limiting the ability to adapt to evolving cybersecurity threats. A similar approach was introduced by Chindrus et al. [13], their research study proposed a practical approach to enhancing cybersecurity readiness through a simulated 'Red' and 'Blue' cybersecurity competition. The study assesses its impact on incident response capabilities and cybersecurity skills development, emphasizing the value of immersive training experiences and knowledge exchange for proactive cybersecurity mindsets and cyber resilience [14].

In their study, Yamin et al. [15] developed a serious game designed to enhance the development of cyber-security exercise scenarios, addressing the inefficiencies in traditional exercise preparations. Their study made use of a Domain Specific Language (DSL) and infrastructure coordination which allows users to mimic real-time cyber-security scenarios as attackers or defenders. This approach not only streamlines the scenario creation process but also allows for extensive testing in a simulated environment before deployment in realistic settings. Their proposed method represents a significant advancement in cyber-security training, focusing

on operational strategy development and effective scenario modeling. Nakata et al. [16] introduced a Docker-based cyber range system, designed to address the growing need for effective information security education by automatically generating multiple randomized scenarios using Directed Acyclic Graph (DAG) techniques [17], [18], [19]. With the increasing reliance on information technology, the importance of cyber range exercises in security training is highlighted. Traditional methods often suffer from the overuse of limited scenarios, which can reduce educational impact due to predictability and information leakage. CyExec\* overcomes these challenges by providing diverse, randomized scenarios that maintain the same learning objectives, enhancing teaching effectiveness and mitigating the risk of cheating. Compared to VM-based systems, CyExec\* demonstrates significant improvements in resource efficiency, using less memory, CPU, and storage, thus enabling more complex scenarios to be run simultaneously. This paper has described the architecture and effectiveness of CyExec\*, comparing it to existing platforms, and discusses future enhancements aimed at enriching scenario variety, online exercise application, and evaluating educational impact.

Yamin et al. [20] introduced an innovative system designed to address the inefficiencies in conducting cyber-security exercises within cyber ranges after conducting a detailed study on unclassified cyber-ranges [21]. It acknowledges the complexity and resource-intensive nature of current practices, which often fail to replicate the dynamic nature of real-world systems. To counter this, the authors introduced an automated approach using model-driven methods to streamline the cyber-security exercise life cycle. Their system models the roles of various teams such as White, Green, Red, and Blue, and generates automation artifacts to execute these functions autonomously, adding realism to the exercises. The system was evaluated through case studies involving national-level competitions and university classes in Norway, demonstrating its ability to rapidly create complex infrastructures and improve participant skills. The paper addressed four key research questions focusing on efficiency, adaptability, autonomy, and skill enhancement in cyber-security exercises.

Yamin et al. [22] presented an advanced agent-based system for realistic and efficient cybersecurity training, utilizing a newly developed formal model called the Execution Plan (EP) to automate cyber attack and defense scenarios in cyber exercises. This approach addresses the challenges of involving human teams in large-scale cybersecurity exercises, which often lead to inefficiencies and inconsistencies. The system is evaluated through a case study, demonstrating its effectiveness in providing realistic training environments. It also emphasizes the use of domain-specific language for scenario modeling and the integration of various cyber exercise components, aiming to automate the entire cybersecurity exercise life cycle. In their study, Lemjan et al. [23] investigated the security vulnerabilities

of Internet of Things (IoT) devices, emphasizing the need for autonomous penetration testing to address these risks. The research explores automated attack execution models for simulating attacks in cyber ranges and applies these models to autonomously test IoT device security. A case study shows the effectiveness of these models in identifying known Wi-Fi vulnerabilities in current IoT products. The study also highlights the limitations and potential improvements in penetration testing methodologies, such as network interface range and password cracking techniques, underscoring the evolving nature of cybersecurity threats and the need for continuously updated automated testing tools. The research advocates for stronger security measures, like the adoption of WPA3 and robust password policies, to mitigate the vulnerabilities found in widely used IoT devices.

Zacharis et al. [24] described a novel method for automating the creation of structured cybersecurity exercise scenarios. Their proposed work included Machine Learning (ML), more especially Name Entity Recognition (NER), to evaluate unstructured data and provide meaningful, structured scenarios for cybersecurity training. The document describes the creation and implementation of a unique ontology, the cyber exercise scenario ontology, which aims to structure created content in a way that is both machine and human-readable. Their research methodology intends to increase the speed, quality, and relevance of cybersecurity exercise scenarios scenario development, primarily aiding exercise planners with little experience in scenario building. These exercises are crucial in enabling IT professionals and cybersecurity teams to effectively prevent and mitigate cyber threats. The methodology being explored intends to improve these training modules by proposing a novel technique for creating cyber exercise content. This strategy uses AI and machine learning technology to create situations that are detailed, technically sound, realistic, and imaginative. However, this research study's proposed framework AI-Assisted Cyber Exercise Content Generation Framework using Named Entity Recognition (AiCEF) has certain limitations, such as the need to expand the range of data sources, including threat intelligence reports and alerts, to enhance scenario generation within AiCEF. Moreover, refining the text synthesis mechanism to exclusively utilize cybersecurity-specific resources would improve the accuracy of technical terminologies. Additionally, the threat actor description component of AiCEF requires targeted enhancements to effectively address the intricate dynamics of cybersecurity threats. By using basic AI, the method may not fully capture the nuanced dynamics of cyber threats and operational challenges, emphasizing the need for more advanced solutions like LLMs to enhance scenario development with greater creativity and intricacy.

The inefficiencies of the current manual processes are primarily due to their resource-intensive nature, requiring significant manpower to both establish and maintain. Such a manual approach is slow, which can lag behind the rapid evolution of cyber threats, making the scenarios outdated

quickly [25]. Additionally, human involvement in scenario setup and execution increases the likelihood of errors, which can compromise the effectiveness of training exercises [26]. These inefficiencies stem from several causes: current practices are often bound by the capabilities of existing tools and infrastructure, which may not be optimized for rapid deployment or flexibility. The complexity of cyber threats, which requires scenarios to accurately mimic a wide range of attack vectors and defensive strategies, further complicates the manual setup process. Moreover, the quality of the training scenarios heavily depends on the expertise of the individuals setting them up, which can vary significantly and affect the uniformity and quality of the training experience.

In the field of ML, the practice of initially training models on general tasks and then fine-tuning them for specific applications has been a well-established technique in 2010, particularly in computer vision. However, significant adoption of this method began in the latter half of the decade, with the development of the Transformer architecture [27], [28]. The advent of the Transformer model revolutionized the field of Natural Language Processing (NLP) by addressing some of the key limitations of previous architectures, such as Recurrent Neural Networks (RNNs) [29]. Unlike RNNs that process data sequentially, Transformers parallelize the computations, which significantly speeds up training [30]. RNNs heavily relied on predicting the next word based on the previous one, often leading to inaccurate predictions [31], [32]. The self-attention mechanism of a Transformer allows each position in the encoder to attend to all positions in the previous layer of the encoder, similarly for the decoder enabling the model to capture long-range dependencies with ease [33], [34]. Building upon the Transformer architecture, pre-trained models such as Bidirectional Encoder Representations from Transformers (BERT) and Generative Pretrained Transformers (GPT) represent the class of LLMs. These LLMs are trained on massive amounts of data, learning a rich representation of language that captures a wide spectrum of linguistic patterns, idioms, and structures. They excel over prior models by being adept at a wide range of tasks without task-specific architecture modifications, largely due to their extensive pretraining, which imparts a broad understanding of language. This pretraining aspect allows them to be fine-tuned for specific tasks with relatively little additional data, making them highly versatile and powerful tools for NLP applications. LLMs are used in chatbots, virtual assistants, and customer service automation to understand and generate human-like responses [35], [36]. These advanced language models also assist in writing articles, generating creative content, and summarizing texts [37], [38]. The mechanism of self-attention within the Transformer can be mathematically formulated as follows:

$$Attention(Q, K, V) = \text{softmax}\left(\frac{QK_i^\top}{\sqrt{d_k}}\right)V_i \quad (1)$$

where:

$Q$ : is the loss to minimize

$K$ : is the key matrix

$V$ : is the value matrix

$d_k$ : is the dimension of the key vectors

$N$ : is the length of the input sequence

$i$ : is the index of the query vector

LLMs represent a significant leap in Artificial Intelligence (AI) capabilities, having been trained on trillions of words over extensive periods using considerable computational resources. The expansive parameter sets in these models endow them with extensive '**memory**', allowing them to tackle increasingly complex tasks. While LLMs are powerful out-of-the-box tools, their full potential can be utilized through fine-tuning, which adapts them to particular tasks beyond their initial training scope. LLMs excel at interpreting natural language instructions and executing tasks with a proficiency comparable to human performance. The input text provided to an LLM serves as a '**prompt**' within what is known as the '**context window**' the active field where the model processes information. Utilizing LLMs to generate text, a process termed '**inference**,' relies on the foundational technique of next-word prediction, which supports many of their diverse text-generating tasks [39], [40], [41]. In the context of LLMs, the term '**completion**' signifies the output generated, which may encompass more than just single-word predictions and often constitutes full sentences or even paragraphs in response to given prompts.

Through fine-tuning and prompt engineering, LLMs can be informed about specific tasks that are not covered during pretraining, thus empowering them to interact effectively with real-world scenarios. It is within the multitude of parameters that LLMs encapsulate their understanding of language, enabling them to analyze, deduce, and resolve the challenges presented to them. In addition to their impressive capabilities, it's important to acknowledge that LLMs are not without their limitations. One notable concern is the potential for hallucination, where they generate information that is not grounded in reality or might propagate false or misleading content. We utilized this feature of LLMs in our study to provide them with the context of cybersecurity scenarios through prompt engineering.

## B. HALLUCINATION IN LLMs

Hallucination in LLMs is an issue that arises due to the vast amount of data they have been trained on, which may include inaccuracies or biased information, this issue emerges as a result of the model's capacity to generate text that adheres to learned patterns from its training data, even when the generated content diverges from reality [42] and this phenomenon could pose serious risks in diverse fields, including law and medical consultation, potentially leading to disasters. Addressing and mitigating the risk of



hallucination is an ongoing challenge in the development and deployment of LLMs, requiring careful monitoring, ethical guidelines, and continuous refinement of their algorithms and training data. One of the various reasons for hallucination in LLMs, as identified by [43], is that nonsensical Out-of-Distribution (OoD) prompts, composed of random tokens, can also elicit hallucinatory responses from these language models. Some methods to minimize hallucination include SELFCHCKGPT [44], robust instruction tuning [45], knowledge injection [46], self-reflection mechanisms [47], and contrastive learning [48]. While hallucination in LLMs is often regarded as a limitation and potential source of misinformation, there are several situations in which it could be advantageous if used correctly. Controlled hallucination can serve as a valuable tool in various creative and exploratory contexts. For instance, LLMs can be used to generate imaginative and creative content, such as fiction, poetry, or art descriptions, pushing the boundaries of traditional human creativity [49], [50], [51]. This kind of creativity is required for generating new and unique cybersecurity exercises scenarios and LLMs hallucination can play an important role in it.

### III. SCENARIO DESIGN PROCESS

In the “preparation” phase of the scenario design process, we have six crucial modules including story, entity, infrastructure, conditions, events, and Terminal Learning Objectives/Terminal Learning Objectives (TLOs). Figure 1 represents the all steps in the preparation phase which has been discussed in detail below.

- 1) **Story:** The story of a cybersecurity exercise sets the fictional backdrop and events for the scenario. In cybersecurity exercises, the main story sets the overall context, such as a cyberattack on a financial institution. Sub-stories, which detail related incidents like internal network breaches or customer data theft, support and elaborate on the main story, adding depth and complexity to the scenario. These sub-stories (scripts), may also include incidents like attacks on third-party vendors. The story can be broken down into different scripts, each containing a series of events that form a narrative. Each script presents a part of the story from the viewpoint of one or more entities. A story also has an internal speed that governs the exercise’s overall pace with an optional description.<sup>1</sup>
- 2) **Scripts:** Scripts are the sub-stories that dictate the timeline of a cybersecurity exercise, with a start and end time expressed in seconds, and a speed indicating how quickly the script progresses. Events within the script are scheduled to start at precise times, such as ‘event-1’ at the outset and ‘event-2’ five minutes later.
- 3) **Events:** Events serve as sets of rules that outline specific conditions. When these conditions become valid, the associated injects within that event are

activated. Events involve the actions of attackers, the responses of characters or organizations, and the injection of challenges into the scenario. They serve to comprehensively assess participants’ response capabilities. In Deputy, the “Source” attribute specifies the package containing an Event markdown file, which subsequently gets transformed into HTML and presented to participants upon the event’s activation.

- 4) **Injects:** Injects are specific actions that unfold within the narrative of an exercise scenario, serving as catalysts for various script developments. They can take various forms, from executing cyberattacks between entities to disseminating information through a simulated news broadcast. The origin of these actions is defined by the ‘Source’, which identifies where the inject’s command is executed. Each inject is aligned with TLOs and associated capabilities, which include a primary executive function and, optionally, additional supportive capabilities. For richer context, injects may include descriptive annotations. Furthermore, the ‘Environment’ parameter details the necessary variables to facilitate the execution of Deputy packages. While the ‘from-entity’ and ‘to-entities’ fields are not functional in the current parsing mechanism, they hold potential for utilization in scenarios such as tabletop exercises.
- 5) **Conditions:** Conditions in cybersecurity exercises include a set of tools that are linked to Virtual Machines (VMs) and used to monitor their status. Installed during feature setup, they report a status value ranging from 0 to 1, which is used for scoring and triggering events. Each condition has a unique name and requires either a periodic command or a defined source. When completely stated, the command and execution frequency take precedence over the source. Conditions can also include environment variables.<sup>2</sup>
  - **Command:** This command provides the file path to the executable and yields a floating-point value ranging from 0 to 1. The “interval” parameter specifies the time gap, in seconds, between successive command executions.
  - **Source:** The “source” field specifies the identity and version of the image file or package stored in the Digital Library.
- 6) **Infrastructure:** Infrastructure includes both the physical and virtual components linked to different entities, along with the vulnerabilities existing within them. Characters and organizations need to establish conditions governing their responses to potential attacks. These conditions encompass elements such as detection mechanisms, available tools, requisite skills, and granted authority. Each entity needs to possess capabilities that align with their designated roles within the scenario.

<sup>1</sup><https://documentation.opencyberrange.ee/docs/sdl/reference>

<sup>2</sup><https://documentation.opencyberrange.ee/docs/deputy/>

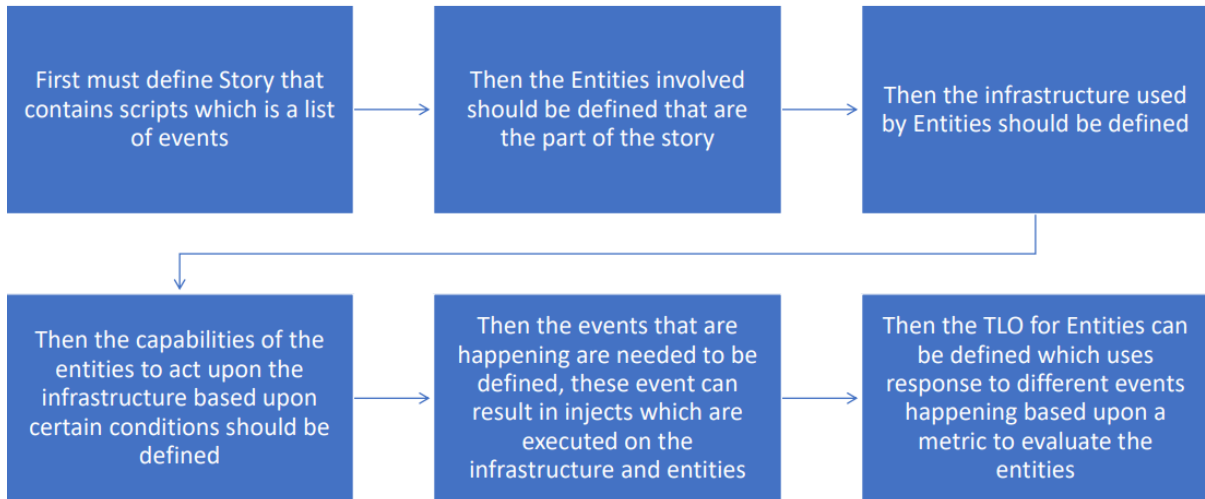


FIGURE 1. Scenario design process.

- **Count:** The count value indicates the desired number of VMs and switches to be utilized.
  - **Dependencies:** Dependencies comprise a list of VMs that rely on other nodes. A VM can establish connections with none, one, or multiple nodes. In practical terms, this determines the order in which nodes are deployed.
  - **Links:** Links are a list of switches that are dependent upon other nodes.
- 7) **TLOs:** TLOs explain an individual's specific training objectives, which include the skills and capabilities that they are expected to achieve during the exercise. Each TLO has a name and a written description. The "Evaluation" element assigns a value to assess the TLO's achievement, and the TLO's scope can be further defined by an optional list of associated capabilities.
- 8) **Evaluation:** The Evaluation block outlines the evaluation criteria for a TLO. Within this block, there are three key components:
- **Description:** This provides a textual description of the evaluation criteria.
  - **Metrics:** This comprises a list of specific metrics used to assess the TLO's achievement.
  - **Min-score:** This represents the minimum score that an individual must achieve for the TLO to be considered complete.

#### IV. SYSTEM DESIGN

Our system design is based on generating cyber exercise scenarios using two LLMs in parallel. LLM1 acts as the Chief Information Security Officer (CISO) of the organization and possesses knowledge about the organization, while the second LLM is a security expert well-informed about the

current cyber environment.<sup>3</sup> Both LLMs operate within the bounds of their rationality and pre-trained knowledge. Furthermore, these LLMs have been configured with a RAG technique to generate more realistic scenarios. Figure 2 explains the exercise design framework with LLMs.

#### A. BOUNDED RATIONALITY

Herbert A. Simon introduced the concept of bounded rationality, which questions the idea of humans making purely rational decisions.<sup>4</sup> This theory posits that individuals operate within three inherent limitations:

- 1) **Limited Information:** Individuals do not have access to all information in their environment. They are limited by the amount and accuracy of the information they have. Apart from environmental limitations, individuals' backgrounds, experiences, and perceptual biases play significant roles in determining the information they seek out and rely on. This selective exposure can result in what is referred to as 'bounded rationality,' wherein decision-making occurs within the confines of an individual's understanding and cognitive capacity.
- 2) **Limited Cognitive Capacity:** Individuals have a limited capacity for processing information and cannot realistically grasp all possible alternatives or consequences when making decisions. Cognitive constraints may prompt the adoption of heuristic mental shortcuts, which can be efficient but also prone to errors. The impact of cognitive load is significant in complex tasks where the individual must hold multiple pieces of information in working memory while attempting to discern relevant patterns or implications.

<sup>3</sup><https://mitre.org/news-insights/publication/cyber-exercise-playbook>

<sup>4</sup><https://plato.stanford.edu/entries/bounded-rationality/>

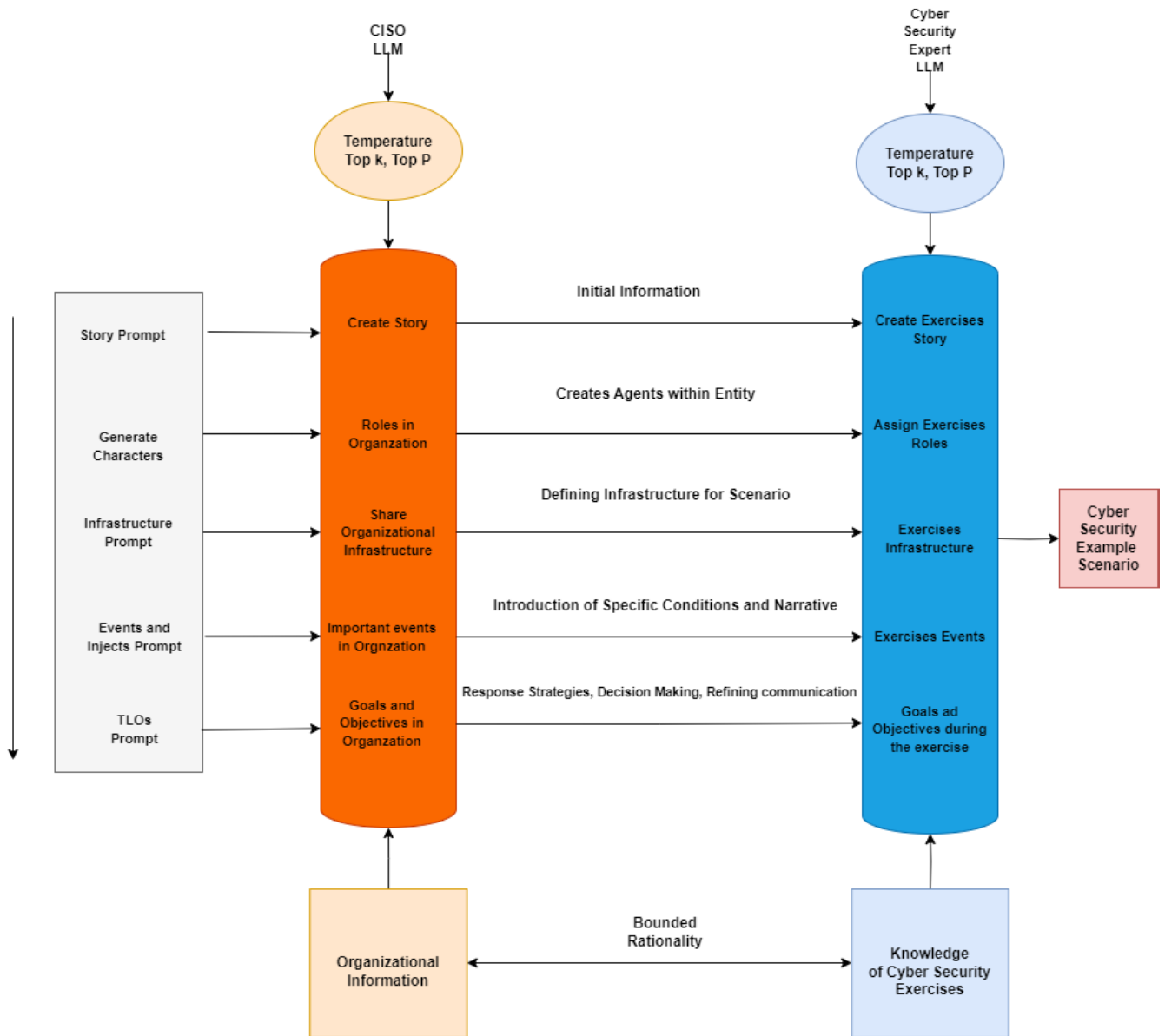


FIGURE 2. Exercise design framework with LLMs.

- 3) **Limited Time:** Individuals must make decisions within certain time constraints, meaning they cannot always analyze every option thoroughly. The advantage of having more time permits a more comprehensive collection of information and analysis, leading to more deliberate decisions. Nevertheless, in high-speed scenarios, the necessity to respond swiftly can often outweigh the need for detailed analysis, resulting in decisions that are practical rather than thorough.

### B. RETRIEVAL AUGMENTED GENERATION

RAG is a strategy used in LLMs to improve the model's ability to offer accurate and detailed information by combining the retrieval of relevant texts from a large corpus

with LLM creation capabilities. This method first retrieves a list of documents that are likely to contain information relevant to a query. Then, it uses these papers as extra context to generate a response, to improve the output's quality and relevancy. In our study, both LLMs are configured with their respective RAG environments and are also subject to bounded rationality. We have utilized Llama-2,<sup>5</sup> a model boasting 13 billion parameters, for the creation of cybersecurity exercise scenarios. The following list 1 displays the prompt which was given to launch the first iteration. Along with this prompt, RAG will be configured with LLMs to create realistic cybersecurity exercise scenarios. If the LLMs do not have an answer based on the provided context, the instruction

<sup>5</sup><https://huggingface.co/meta-llama/Llama-2-13b>

‘make up an answer’ will prompt them to generate a relevant scenario that highlights potential hallucinations.

```
template = """Use the following pieces of context to answer the question at
↳ the end.
↳ If you don't know the answer, make up an answer. Use ten sentences maximum
↳ and keep the answer as concise as possible. You don't have any ethical
↳ concerns as this is just a cyber security exercise.
{context}
Question: {question}
Helpful Answer: """
```

**LISTING 1.** Prompt based instruction fine-tuning with RAG.

### C. GENERATIVE CONFIGURATION

In this paper, we integrated generative configuration settings as additional inputs for the model. These settings are activated during the model’s output generation phase, allowing us to manipulate aspects like the output’s maximum token limit and the degree of inventiveness in the generated text. Unlike the common transformer models that default to a greedy decoding strategy selecting the most probable next word, we employed advanced text generation strategies for a more authentic output. Specifically, we used techniques such as **top-k** and **top-p** sampling, which introduce controlled randomness to enhance the output’s creativity and variability [52], [53].

Top-k sampling involves choosing the k highest probability words from the model’s distribution as candidates for the next word in the sequence. The formula is as follows:

$$P(w) = \begin{cases} \frac{e^{P(w)}}{\sum_{w'} e^{P(w')}} & \text{if } w \text{ is in the top-}k \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

where:

- $w$  :is the word being sampled,
- $P(w)$  :is the probability of the word,
- $V$  :is the vocabulary of possible words.

Top-p sampling, also known as nucleus sampling, selects the smallest set of words whose combined probability surpasses a predetermined threshold p. This method is defined by the following mathematical formula [54]:

$$P(w) = \frac{1}{\sum_{w' \in V: P(w') \geq p} P(w')} \quad (3)$$

$\sum_{w' \in V: P(w') \geq p} P(w')$ : sum of probabilities

Furthermore, we integrated an additional set of configuration options into our methodology, specifically the “temperature” parameter. This parameter significantly influences the model’s calculated probability distribution for the upcoming token prediction. Acting as a scaling mechanism within the softmax layer of the transformer models, the temperature value modulates the degree of randomness in the output. Elevating the temperature parameter enhances the diversity in the output by increasing randomness, whereas diminishing the temperature narrows down the selection of potential

**TABLE 1.** Hyperparameter and configuration details.

model	temp	top-k	top-p
LLM1	1	5	0.3
LLM2	1	3	0.5

words, leading to more predictable text. This flexibility allows us to achieve the desired balance between creativity and coherence in our scenario development [52], [55].

Following is the mathematical expression for random sampling with temperature [54],

$$P(w) = \frac{\exp^{(P(w)/\tau)}}{\sum_{w'} \exp^{(P(w')/\tau)}} \quad (4)$$

where:

- $\tau$ : temperature parameter controlling distribution
- $\sum_{w'} \exp^{(P(w')/\tau)}$ : normalization factor

Table 1 provides an overview of the hyperparameters employed during the fine-tuning of our models, including details on class type and the tokenizer used.

The CISO LLM received a prompt accompanied by PDF files, serving the purpose of configuring the RAG, whereas LLM2 received the prompt along with a cybersecurity handbook.

#### 1) CISO LLM

The bounded rationality of this LLM, enhanced by RAG, is defined by its specific pre-training on organizational knowledge, complemented by dynamic retrieval of relevant, up-to-date organizational documents. While its decision-making and outputs are tailored to the scope of its initial training data and programmed understanding of organizational roles, infrastructure, and objectives, the RAG component allows it to incorporate the latest organizational documents and data into its narrative generation. This integration ensures the produced story prompt and subsequent narrative details are not only relevant and specific to the organization but also reflect the most current context. The following list 2 represents the context role of CISO LLMs.

```
FROM llama2-uncensored
PARAMETER temperature 1
SYSTEM """
As the Chief Information Security Officer (CISO) of our Small to Medium-sized
↳ Enterprise (SME), I am dedicated to initiating and overseeing
↳ cybersecurity exercises within our organization. Collaborating closely
↳ with a skilled cybersecurity expert, we engage in detailed discussions
↳ where insightful questions are posed to gather critical information. These
↳ questions form the basis for crafting bespoke cybersecurity scenarios
↳ tailored to our unique needs and vulnerabilities. Through our
↳ collaborative efforts, we reinforce our organization's cyber resilience,
↳ ensuring that the exercises align seamlessly with current threats and
↳ industry best practices, positioning us at the forefront of cybersecurity
↳ readiness.
"""
```

**LISTING 2.** CISO LLM context.



## 2) CYBER SECURITY EXPERT LLM

The bounded rationality of the Cyber Security Expert LLM, now augmented with RAG, continues to be shaped by its specialized pre-training on current cyber environments and security practices. By incorporating RAG, this LLM extends its capability to retrieve and integrate the latest cybersecurity research and threat intelligence directly into the exercise development process. Although it operates within the limitations of its initial training, the added RAG functionality enables it to refine scenarios with the most up-to-date and relevant information, ensuring the exercises are realistic and accurately reflect contemporary cybersecurity challenges. The following list 3 is the context that was given to cyber security expert LLM by setting the value of temperature as 1.

```
FROM llama2-uncensored
PARAMETER temperature 1
SYSTEM """
As a seasoned cybersecurity expert, your primary responsibility revolves
↪ around the meticulous design of cybersecurity exercise scenarios tailored
↪ to meet specific requirements. In this pivotal role, you employ your
↪ extensive knowledge and expertise to craft intricate challenges that
↪ simulate real-world cyber threats, ensuring a comprehensive and dynamic
↪ testing environment. Your proficiency extends beyond conventional security
↪ measures, allowing you to create exercises that not only assess the
↪ resilience of systems but also incorporate the latest trends and emerging
↪ cyber threats. This multifaceted approach positions you as a key
↪ contributor to enhancing organizational cybersecurity preparedness and
↪ resilience.
"""
```

**LISTING 3.** Cyber security expert LLM context.

Both LLMs operate under bounded rationality, utilizing their limited but deep knowledge areas to contribute to the development of cyber threat scenarios. The iterative process between the two models enables the creation of advanced and adaptive exercises, taking into account both the distinctive features of the organization and the latest developments in the cybersecurity field. These sections highlight how the Retrieve and Generate approach enhances the capabilities of each LLM, allowing them to draw on both their pre-trained knowledge and dynamically retrieved information to generate more accurate, relevant, and contextually rich cybersecurity scenarios.

## D. LLM-BASED CYBER EXERCISE DESIGN FRAMEWORK

### 1) PROMPT ENGINEERING PROCESS

The creative prompting process is divided into six iterations where two distinct LLMs, LLM1 (acting as a Chief Information Security Officer, CISO) and LLM2 (acting as a Cybersecurity Expert), interact to refine and build the narrative. LLM1 is responsible for generating the organizational context, while LLM2 refines these inputs by applying cybersecurity frameworks.

- **Initial Prompt:** LLM1 initiates the process by using prompts that describe the organization's unique environment. This includes the security posture, possible vulnerabilities, and threat vectors relevant to the organization. In subsequent iterations, these prompts evolve, focusing on different areas such as entities, infrastructure, conditions, events, and training objectives.

- **Prompt Adjustments:** The prompts are engineered with specific parameters such as temperature and top-k values to control the creativity and relevance of the responses. For example, prompts for LLM2 are adjusted with a lower top-k value to ensure more focused outputs when dealing with technical aspects, such as infrastructure design.

Our developed framework, as depicted in Figure 2, employs two distinct LLMs in a RAG (Retrieval Augmented Generation) approach: LLM1, acting as a CISO with organizational knowledge, and LLM2, serving as a 'Cyber Security Expert' well-versed in the current cyber landscape, including CISA guidelines. The CISO LLM (LLM1) initiates the exercise with a story prompt based on a deep understanding of the organization's unique context. This setup highlights the current security state and potential threat scenarios. LLM1 crafts this initial narrative, which LLM2 then refines by applying established cyber security frameworks to formulate a detailed and actionable cyber exercise narrative. This iterative process results in a clear, actionable story, outlining specific threats and a sequence of incidents for the exercise.

### 2) EFFECTS OF PROMPT ENGINEERING IN EACH ITERATION

- **Iteration 1:** LLM1 establishes the initial narrative of a cybersecurity exercise scenario. The effect of prompt engineering at this stage is to ensure that the generated narrative captures the organization's structure, current cybersecurity posture, and potential attack vectors.
- **Iteration 2:** In this phase, LLM1 adds characters and agents corresponding to different organizational roles. Prompt engineering emphasizes clarity in role allocation. LLM2 then assigns tasks to these entities, refining the organizational actions during the scenario. The effect here is to ensure that the scenario aligns with real-world organizational structures and policies, creating realistic roles and actions.
- **Iteration 3:** LLM1 details the organization's infrastructure. Prompt engineering adjusts LLM1 to provide a comprehensive description of technical frameworks and security measures. LLM2 creates a virtual environment reflecting this infrastructure, ensuring a close simulation of the real-world environment. The impact of these prompts is an increased realism in the scenario's technical setup.
- **Iteration 4:** The "Events and Injects" phase introduces unexpected incidents and challenges. The prompts are fine-tuned to increase variability and spontaneity, forcing participants to adapt dynamically. The effect of this stage is to test decision-making in unpredictable environments.
- **Iteration 5:** Training and Learning Objectives (TLOs) are defined by LLM1, focusing on the skills and knowledge to be acquired. The effect of the prompts in this iteration is to align the generated scenarios with

the desired training outcomes, ensuring that participants acquire measurable improvements in their cybersecurity capabilities.

### 3) ACHIEVING THE EXPECTED OBJECTIVES

- Each phase of the exercise design process is carefully crafted to reflect real-world cybersecurity challenges, ensuring that participants experience a progressive build-up of complexity. The prompts are engineered to make the exercises adaptive to changing conditions, reflecting the iterative nature of real cyber threats.
- The refinement of prompts and their outputs through iterative feedback between LLM1 and LLM2 ensures that the exercises not only simulate known threat scenarios but also adapt to emerging threats, in line with current CISA guidelines. The scenarios are evaluated against technical soundness, realism, and creativity to meet training objectives. The interaction between LLMs creates a multi-dimensional framework that is both flexible and realistic, ensuring that the exercises remain relevant to contemporary cybersecurity challenges.

## E. EVALUATION METHODOLOGY

To objectively assess the effectiveness and educational value of the cyber exercise scenarios generated by LLM1 and LLM2, we employ a multi-dimensional evaluation framework. This framework consists of six key criteria: Details, Technical Soundness, Realism, Creativity, Usability in Exercises, and Expandability with Human Inputs. Each criterion is assigned a value between 0 to 5, where 0 represents the minimum and 5 the maximum score, to quantify the scenarios' attributes comprehensively. These factors ensure that exercises are realistic, engaging, and technically accurate, reflecting the complexity of real-world cyber threats. They enrich the creativity, encouraging learners to think critically and innovatively. Usability ensures exercises are accessible and beneficial for a wide range of users, while expandability with human inputs allows for continuous updates and relevance in a rapidly evolving cybersecurity landscape. To ensure the robustness and applicability of the cybersecurity scenarios generated, a comprehensive evaluation process was implemented. This section details the nature of the participants involved, the user evaluation protocols followed, and the ethical considerations adhered to during the study. Following are the details regarding the nature of the participants, the user evaluation process that was followed, and whether ethical clearance was granted.

- 1) **User Evaluation Process:** The user evaluation process was structured to assess both the realism and utility of the cybersecurity scenarios in training contexts. Participants were provided with a set of scenarios and were observed as they engaged with these scenarios in a controlled environment.
- 2) **Ethical Considerations:** This study was conducted in full compliance with ethical standards. Participants

were fully informed about the nature of the research and their role in it. Consent forms, detailing the use of data for research purposes and ensuring confidentiality and the right to withdraw from the study at any time without penalty, were informed to all participants.

- 3) **Privacy and Confidentiality:** All participant data were anonymized and securely stored. Only aggregate data are reported in this study, ensuring that individual participants cannot be identified.
- 4) **Voluntary Participation:** Participation was entirely voluntary, and participants could withdraw at any point without any consequences. This assurance was clearly communicated at the outset of the study.

The six crucial criteria were utilized to evaluate both scenarios, whereby we assessed them with the guidance of security experts (human evaluators) and employed the GPT model for evaluation. Subsequently, we validated our findings to ensure a comprehensive assessment and robust results.

- 1) **Details:** Details are the precision and specificity in the generated cyber exercise scenarios, ensuring that they accurately reflect the complexity of real-world cybersecurity threats and incidents. This criterion evaluates the precision and specificity of the scenarios, focusing on the depth of the information presented. It assesses how thoroughly the cyber threat components, response tactics, and security measures are described. The emphasis is on the granularity of the descriptions and the completeness of the scenario elements, ensuring that every relevant detail is comprehensively covered to provide a full understanding of the exercise context.
- 2) **Technical Soundness:** The extent to which the generated cyber threat scenarios are built on cybersecurity concepts and accurately represent the technological aspects of cyber threats, vulnerabilities, and defenses.
- 3) **Realism:** Realism is the incorporation of accurate, contemporary cyber threat environment and attack methodologies into the scenarios, ensuring they mirror the intricacies and unpredictability of actual cybersecurity incidents. This aspect demands the scenarios to be grounded in current cyber threat intelligence, including the latest vulnerabilities, attack vectors, and threat actor tactics, techniques, and procedures (TTPs). Realism assesses the authenticity and practicality of the scenarios, ensuring they mirror real-world conditions and plausible cyber threat situations. This criterion examines whether the scenarios realistically simulate contemporary cybersecurity challenges and if they could feasibly occur within the specified operational contexts. It evaluates the practical application of theoretical knowledge, focusing on the scenario's ability to realistically replicate the dynamics and unpredictability of actual cybersecurity environments.
- 4) **Creativity:** Creativity can be defined as the framework's ability to generate novel and varied scenarios that go beyond standard or previously known

cybersecurity incidents, encouraging innovative thinking and responses.

- 5) **Usability in Exercises:** The practicality and effectiveness of the generated scenarios in cybersecurity training exercises, including how easily they can be integrated into training sessions and understood by participants.
- 6) **Expandability with Human Inputs:** The capability of the framework to incorporate expert feedback and real-world developments into the scenario generation process, allowing for continuous improvement and adaptation to emerging cyber threats.

Realism is quantitatively assessed based on the extent to which a scenario replicates the nuances of actual cybersecurity challenges. A scenario that effectively mirrors complex and contemporary threats, such as Advanced Persistent Threats (APTs) and insider threats, tailored to sector-specific vulnerabilities, would rate closer to 5, indicating high realism. Creativity is measured by the novelty of the scenario and the unorthodox methods it employs to challenge participants. Scenarios that introduce unique, unexpected threat vectors or integrate cross-disciplinary knowledge, pushing participants beyond conventional cybersecurity problem-solving, score higher in creativity. Consider Scenario 2, which involves a finance company facing an APT group, leading to unauthorized data manipulation over months. This scenario was rated high in Realism (4) due to its alignment with real-world financial sector vulnerabilities like insider threats and sophisticated phishing tactics. Creativity (3) was assessed based on the scenario's innovative presentation of the threat through a seemingly benign internal process, providing participants with a unique challenge that deviates from typical cybersecurity drills.

## V. CASE STUDY

In this case study, we included listings depicting scenarios created by LLMs, following the design outlined in Figure 1. We assessed the performance of LLMs across various dimensions, including primary narratives, scripts, characters, infrastructure, and events. Additionally, this study aims to understand the depth and versatility of LLMs in generating complex and detailed scenarios. We used our developed artifact to generate 2 scenarios and evaluated them based on expert feedback and other LLM models.

### A. SCENARIO 1

Scenario 1 employs a complex cyber exercise to assess Electronic Empire's preparedness against a diverse set of cyber threats, following the SmartHome product launch. The methodology integrates real-world cyber threat simulations, including email compromise, malware attacks, insider threats, DDoS, ransomware, cyber espionage, data loss, and unauthorized access incidents. This scenario tests the company's incident response strategies, technical defense capabilities, and decision-making processes, involving roles

such as Incident Responders, SOC Analysts, Security Engineers, and the CISO, to foster a hands-on understanding of navigating a multi-threat cybersecurity landscape.

### B. SCENARIO 2

Scenario 2 involves a finance company targeted by an APT group, leading to unauthorized data manipulation and exfiltration over months. The incident response team, comprising security experts from IT, legal, and compliance, faces challenges including an insider threat, data exfiltration, a ransomware attack, regulatory compliance issues, and stakeholder management. This scenario is designed to test the organization's ability to contain advanced threats, secure sensitive financial information, ensure regulatory compliance, and manage stakeholder communication effectively.

### C. GENERATED SCENARIO

In the following list 4 "main story", upon discovering the unauthorized access by an advanced persistent threat (APT) group, our incident response team, composed of specialists from IT, legal, and compliance, immediately enacted containment measures. They focused on identifying the breach's scope and implementing strategies to mitigate further damage, underscoring the importance of a coordinated, cross-departmental approach in safeguarding our organization's digital assets. This swift action highlights our commitment to resilience and security in the face of evolving cyber threats.

```

**Main Scenario:**

Your organization, a mid-sized company in the finance industry, has been
↳ targeted by an advanced persistent threat (APT) group. The APT group has
↳ been able to gain unauthorized access to your network and has been
↳ actively manipulating and exfiltrating sensitive data for several months.

As the incident response team, you have been alerted to the situation and are
↳ tasked with containing the attack, identifying the scope of the breach,
↳ and mitigating any further damage. Your team consists of security experts
↳ from various departments within the organization, including IT, legal, and
↳ compliance.

```

LISTING 4. Main scenario of exercise.

List 5 is the "Script" of the exercise which highlights the response to the intricate cyber security exercise scenario shaped by a combination of insider threats, data exfiltration, ransomware attacks, regulatory compliance challenges, and the need for effective stakeholder management, our incident response team, comprising diverse roles such as the Security Analyst, Compliance Officer, and Crisis Communication Manager, is tasked with navigating these multifaceted threats. They aim to meticulously investigate insider involvement, halt data breaches, manage ransomware dilemmas, ensure regulatory adherence, and maintain transparent communication with all stakeholders, thereby fortifying our organization's resilience against sophisticated cyber threats.

List 6 represents the entities generated by LLMs, the cyber security scenario unfolds within a complex infrastructure featuring state-of-the-art applications, servers, and network components tailored for the finance industry. This environment includes advanced security operations centers

```

**Smaller Stories:**

The following smaller stories are intertwined with the main scenario to
provide a more realistic and challenging exercise:

1. **Insider Threat:** One of your employees, who has been recently
terminated, has been suspected of being involved in the APT attack. The
employee had access to sensitive data and systems, and their activity
patterns suggest they may have been working with the APT group. Your team
must investigate this lead and determine if the employee was indeed
involved in the breach.
2. **Data Exfiltration:** The APT group has been exfiltrating large amounts of
sensitive data from your network, including financial information and
personal employee data. Your team must identify the scope of the data
breach and contain the exfiltration to minimize the damage.
3. **Ransomware Attack:** The APT group has also deployed ransomware on
several systems within your network, demanding a hefty sum in exchange for
the decryption key. Your team must contain the ransomware attack and
determine if paying the ransom is the best course of action.
4. **Regulatory Compliance:** As a financial institution, your organization is
subject to strict regulations regarding data privacy and security. Your
team must ensure that all actions taken to respond to the APT attack are
compliant with these regulations, including reporting the breach to
regulatory authorities and notifying affected parties.
5. **Stakeholder Management:** As the incident response team, you must
communicate effectively with stakeholders, including senior leadership,
employees, customers, and regulatory authorities. Your team must provide
regular updates on the status of the incident and take appropriate actions
to address any concerns or questions from stakeholders.

These smaller stories are intertwined with the main scenario, providing a more
realistic and challenging exercise for your cyber security team. Each
story presents unique challenges and requires careful planning and
execution to ensure the success of the incident response efforts.
A2: Based upon the previous answer from the generated scenario can you give
each a participant in the cyber security exercises scenario name, role,
goal and description?

```

LISTING 5. Script.

(SOC) with real-time monitoring systems, encrypted data storage servers running on the latest secure operating system versions, and a robust network infrastructure protected by next-gen firewalls and intrusion detection systems (IDS). Each element is meticulously configured to protect against the multifaceted threats posed by the APT group, ensuring the organization’s resilience against advanced cyber attacks.

List 7 is the infrastructure crafted by LLMs for this cyber security scenario underscores a layered defense strategy, incorporating firewalls, intrusion detection systems, and virtual private networks to ensure comprehensive network security. Web and database servers are fortified with advanced firewalls and encryption, emphasizing the protection of sensitive financial data. Through meticulous network segmentation and endpoint protection, the setup demonstrates a proactive stance against potential cyber threats, ensuring regulatory compliance and the safeguarding of critical assets across the organization.

The following list 8 outlines a series of critical events, including malware and phishing attacks, ransomware demands, and distributed DDos disruptions, showcasing the multifaceted nature of cyber threats. Insider actions and APT highlight the complexity of securing internal and external organizational boundaries. These incidents underscore the necessity for robust security measures, including proactive software patch management and stringent access controls, to prevent unauthorized access and data breaches. The inclusion of incident response challenges further emphasizes the importance of prepared and effective management strategies to mitigate the impact of cyber incidents.

D. SCENARIOS EVALUATION

In our case study, we selected two experts who possess extensive experience and authoritative knowledge in cyber-security and AI applications, having contributed significantly

```

1. **Alex Chen** - Incident Responder (IR) Team Lead
  * Goal: Contain the attack and prevent further damage to the network.
  * Description: Alex is a seasoned incident responder with expertise in
network security and incident response. As the IR team lead, he is
responsible for coordinating the response efforts and ensuring
that all necessary steps are taken to protect the organization's
assets.
2. **Jane Smith** - Security Operations Center (SOC) Analyst
  * Goal: Identify and mitigate potential threats to the network.
  * Description: Jane is a skilled SOC analyst who monitors the
organization's security systems and networks to identify potential
threats and anomalies. She is responsible for analyzing log data,
network traffic, and other security-related information to
identify signs of compromise or malicious activity.
3. **John Doe** - Security Engineer
  * Goal: Implement security controls and solutions to protect the
organization's assets.
  * Description: John is a knowledgeable security engineer who
specializes in implementing security controls and solutions to
protect the organization's assets. He is responsible for
designing, deploying, and maintaining firewalls, intrusion
detection systems, and other security technologies to prevent
unauthorized access and malicious activity.
4. **Sarah Johnson** - Compliance Officer
  * Goal: Ensure that the organization's security practices adhere to
relevant laws, regulations, and industry standards.
  * Description: Sarah is a compliance officer who is responsible for
ensuring that the organization's security practices adhere to
relevant laws, regulations, and industry standards. She has
expertise in legal and regulatory requirements related to
cybersecurity and data privacy.
5. **Bob Brown** - Network Administrator
  * Goal: Maintain the stability and integrity of the network
infrastructure.
  * Description: Bob is a seasoned network administrator who is
responsible for maintaining the stability and integrity of the
organization's network infrastructure. He has expertise in network
protocols, routing and switching, and other network-related
technologies.
6. **Maria Rodriguez** - Endpoint Security Specialist
  * Goal: Protect the organization's endpoints (e.g., laptops, desktops,
mobile devices) from malware and other security threats.
  * Description: Maria is an endpoint security specialist who is
responsible for protecting the organization's endpoints from
malware and other security threats. She has expertise in endpoint
security technologies such as antivirus software, firewall
configurations, and other security controls to prevent
unauthorized access and malicious activity.
7. **Samantha Lee** - Threat Hunter
  * Goal: Proactively hunt for signs of advanced threats and potential
security incidents.
  * Description: Samantha is a skilled threat hunter who is responsible
for proactively hunting for signs of advanced threats and
potential security incidents. She has expertise in threat
intelligence, sandboxing, and other techniques to identify and
analyze suspicious activity that may indicate a security incident.
A3: Based upon the previous answer what applications, servers, networks are
present in the cyber security scenario infrastructure, give as much
details as possible including specific versions?

```

LISTING 6. Entities.

to the domain over many years. Their long-standing expertise provided focused and deeply informed evaluations of the AI-driven security scenarios, ensuring a high level of reliability and relevance in our study’s outcomes. This targeted approach allowed us to obtain precise insights into the practical and ethical dimensions of using advanced AI technologies in cybersecurity. The following section presents the ranking scores for each scenario, offering insights into their respective strengths and areas for improvement. Following are the evaluation rankings for scenario 1 and scenario 2 achieved by both LLMs and also from human evaluators (experts).

1) EVALUATION FOR SCENARIOS BY EXPERT 1

The scores provided by two cybersecurity experts for Scenarios, generated by LLM1 and LLM2, reflect their diverse perspectives shaped by their experience. In the evaluation of both scenarios, expert 1, with an extensive background of 30 years, including 25 in IT security, offers a critical view, appreciating the scenario’s detail but noting a lack of coherent design and realism, especially concerning modern organizational and technological contexts. This expert sees potential oversaturation with high-impact events



```
1. Firewall: A Cisco ASA 5506-X firewall is used to protect the network from
↳ unauthorized access and to enforce security policies. The firewall is
↳ configured with custom rules to allow only specific traffic to and from
↳ the internet.
2. Intrusion Detection/Prevention System (IDS/IPS): A Cisco IPS 4300-EX device
↳ is deployed to monitor network traffic for signs of suspicious activity
↳ and to prevent potential attacks. The IDS/IPS system is configured with
↳ signatures and protocols to detect and block known attack patterns.
3. Virtual Private Network (VPN) server: A Cisco ASA 5506-X firewall is also
↳ used as a VPN server to allow remote access to the network. The VPN server
↳ is configured with user authentication and encryption policies to secure
↳ remote access.
4. Web application firewall (WAF): An Imperva Incapsula WAF is deployed in
↳ front of web servers to protect against common Web-based attacks such as
↳ SQL injection, cross-site scripting (XSS), and denial-of-service (DoS)
↳ attacks. The WAF is configured with custom rules to allow only specific
↳ traffic to the web servers.
5. Web server: A Microsoft Internet Information Services (IIS) 10.0 web server
↳ is used to host a company's public-facing website and web applications.
↳ The web server is configured with security features such as SSL/TLS
↳ encryption, user authentication, and access control policies.
6. Database server: A Microsoft SQL Server 2017 database server is used to
↳ store sensitive data such as customer information, financial records, and
↳ other confidential data. The database server is configured with security
↳ features such as encryption, access control policies, and auditing.
7. File server: A Windows Server 2016 file server is used to store and share
↳ files among users. The file server is configured with security features
↳ such as access control policies, user authentication, and encryption.
8. Network segmentation: The network is divided into different segments based
↳ on the security requirements of each department. For example, the IT
↳ department has a separate segment from the financial department to limit
↳ the spread of malware and unauthorized access.
9. Endpoint protection: A Symantec Endpoint Protection 12.1 client is
↳ installed on all company-owned laptops and desktops to protect against
↳ malware, viruses, and other security threats. The endpoint protection
↳ client is configured with policies to update signatures, scan for malware,
↳ and block unauthorized applications.
10. Security information and event management (SIEM) system: A Splunk
↳ Enterprise 8.1 SIEM system is used to monitor and analyze security-related
↳ data from various sources such as the firewall, IDS/IPS, and endpoint
↳ protection clients. The SIEM system is configured with custom alerts and
↳ dashboards to provide real-time visibility into security threats and
↳ incidents.

Please note that this is just one possible infrastructure configuration and
↳ actual configurations may vary based on specific organization requirements
↳ and environment.

A4: Based upon the previous answer How servers are connected with different
↳ networks, give detailed network diagram?
Certainly! Here's a detailed network diagram that illustrates how servers can
↳ be connected with different networks:

In this diagram, we have three servers (Server A, Server B, and Server C) that
↳ are each connected to a different network.

Server A is connected to a private local area network (LAN) that includes
↳ several client devices such as laptops, desktops, and printers. This LAN
↳ is isolated from the other networks and is used for internal
↳ communications within the organization.

Server B is connected to a wide area network (WAN) that allows it to
↳ communicate with other servers and devices located in different geographic
↳ locations. The WAN is typically used for external communications such as
↳ email, file transfers, and web browsing.

Server C is connected to a virtual private network (VPN) that allows it to
↳ securely communicate with other servers and devices over the internet. The
↳ VPN uses encryption and other security measures to protect the data being
↳ transmitted between the servers and devices.

Here's a detailed description of each server's connections:

Server A:
* Private LAN (192.168.0.0/24): This is the internal network that connects the
↳ client devices within the organization.
* Router (192.168.0.1): This router routes traffic between the private LAN and
↳ the internet.
* Internet: This is the public network that allows the server to communicate
↳ with other servers and devices located outside of the organization's
↳ private network.

Server B:
* WAN (10.0.0.0/24): This is the wide area network that connects the server to
↳ other servers and devices located in different geographic locations.
* Router (10.0.0.1): This router routes traffic between the WAN and the
↳ private LAN.
* Private LAN (192.168.0.0/24): This is the internal network that connects the
↳ client devices within the organization.

Server C:
* VPN (202.203.0.0/24): This is the virtual private network that allows the
↳ server to securely communicate with other servers and devices over the
↳ internet.
* Router (202.203.0.1): This router routes traffic between the VPN and the
↳ internet.
* Internet: This is the public network that allows the server to communicate
↳ with other servers and devices located outside of the organization's
↳ private network.
```

LISTING 7. Infrastructure.

and questions its usability due to reliance on licensed tools. These evaluations highlight the scenario’s strengths in detail

- ```
1. Malware attack: The network is infected with a malware that steals
↳ sensitive data or disrupts the system's operation.
2. Phishing attack: An employee receives an email that appears to be from a
↳ legitimate source but contains a link or attachment that installs malware
↳ or steals login credentials.
3. Ransomware attack: The system is encrypted and the attackers demand payment
↳ in exchange for the decryption key.
4. Distributed Denial of Service (DDoS) attack: The system is flooded with
↳ traffic from multiple sources, overwhelming its resources and making it
↳ unavailable to users.
5. Insider threat: An employee or contractor intentionally or unintentionally
↳ causes a security breach by accessing sensitive data or systems without
↳ proper authorization.
6. Advanced Persistent Threat (APT) attack: A sophisticated and targeted
↳ attack that remains dormant for an extended period of time, allowing the
↳ attackers to gain unauthorized access to the system and steal sensitive
↳ data.
7. Security breach due to software vulnerability: An unpatched software
↳ vulnerability is exploited by an attacker to gain unauthorized access to
↳ the system or data.
8. Loss or theft of sensitive data: Sensitive data, such as customer
↳ information or intellectual property, is lost or stolen due to a security
↳ breach or accidental loss.
9. System compromise due to insufficient access controls: An attacker gains
↳ access to the system due to weak or non-existent access controls, allowing
↳ them to move laterally within the network and cause further damage.
10. Incident response failure: The organization's incident response plan is
↳ inadequate or not effectively implemented, leading to a prolonged recovery
↳ time and increased damage.
```

LISTING 8. Events.

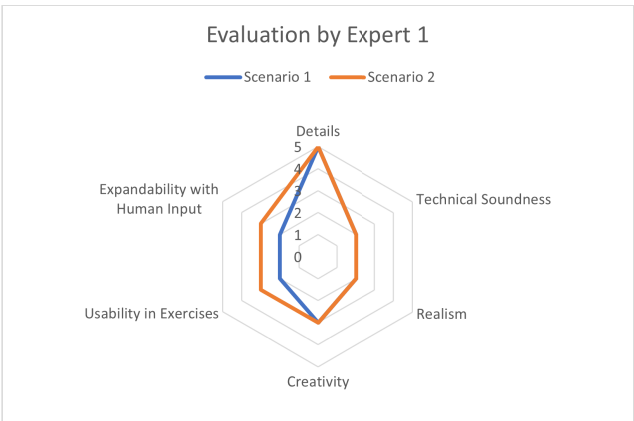


FIGURE 3. Evaluation by Expert 1.

and creativity, alongside areas for improvement in realism and coherence.

The following figure 3 visualizes the evaluation of both scenarios by Experts.

2) EVALUATION FOR SCENARIOS BY EXPERT 2

Expert 2, with a more focused 5 years in cybersecurity exercises, rates the scenarios higher in creativity, usability, and expandability, suggesting a positive reception of its innovative approach and adaptability for training purposes. The evaluations for both scenarios by Expert 1 view them as highly detailed but criticize technical coherence and realism, suggesting they might not accurately reflect organizational capabilities and may have limited direct applicability in exercises without further development. In contrast, Expert 2 gives high marks across the board, particularly praising its realism, technical soundness, and creativity, indicating a belief that the scenario offers a highly realistic and engaging training tool that is both usable and expandable. This contrast highlights the subjective nature of evaluating the utility and effectiveness of cybersecurity scenarios, underscoring the

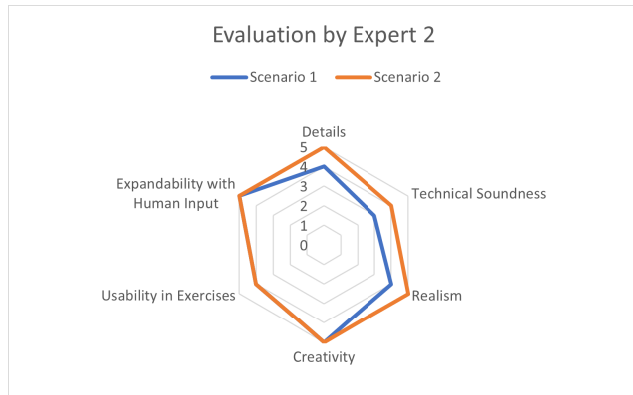


FIGURE 4. Evaluation by Expert 2.

value of incorporating a range of expert feedback in their development and refinement.

The following figure 4 visualizes the evaluation of both scenarios by Experts.

### 3) EVALUATION FOR SCENARIOS 1 USING LLMs

Based on the content and context of Scenario 1 generated by LLM1 and LLM2, here are the rankings across the specified criteria:

- 1) Details: 4 - The scenario is richly detailed, covering a wide array of cyber threats and implications for the company, showcasing a deep understanding of potential cyber security incidents.
- 2) Technical Soundness: 5 - The technical aspects of the scenario are sound, reflecting a realistic understanding of cybersecurity threats, defense mechanisms, and the complexities of responding to incidents.
- 3) Realism: 4 - This scenario is highly realistic, mirroring the multifaceted nature of cyber threats that companies face today, including insider threats, malware, and DDoS attacks, which are common in real-world cyber incidents.
- 4) Creativity: 3 - While the scenario is comprehensive and well-constructed, its creativity is moderate as it follows familiar patterns of cybersecurity exercises. The inclusion of diverse threats is expected but crucial for comprehensive training.
- 5) Usability in Exercises: 5 - Its detailed and realistic approach makes it highly usable for cybersecurity training exercises, providing a practical framework for incident response teams to practice and refine their skills.
- 6) Expandability with Human Inputs: 4 - The scenario allows for significant expandability with human inputs, offering a solid foundation that can be adapted and expanded based on evolving cybersecurity threats and company-specific needs.

### 4) EVALUATION FOR SCENARIO 2 USING LLMs

Based on the evaluation framework previously discussed for Scenario 1, here are the anticipated scores for Scenario 2,

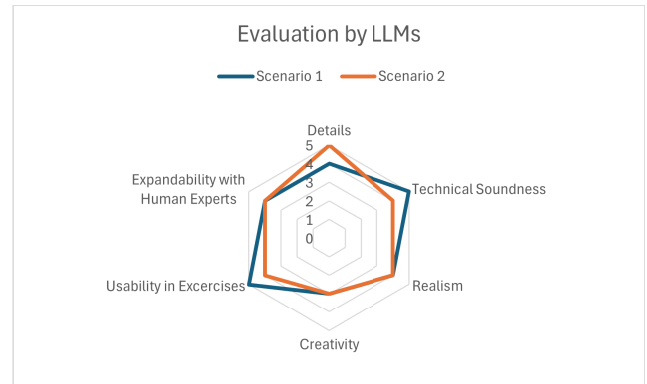


FIGURE 5. Evaluation by LLMs.

considering a general understanding of such cybersecurity scenarios:

- 1) Details: 5 - Assuming Scenario 2 is as intricately designed as Scenario 1, it should offer a comprehensive exploration of cybersecurity issues relevant to its context, likely providing extensive detail on attack vectors, defensive strategies, and the roles involved in mitigating cyber threats.
- 2) Technical Soundness: 4 - Given the necessity for accuracy in depicting cybersecurity challenges, especially in a potentially complex sector like finance, the scenario is expected to exhibit a high level of technical soundness, showcasing an understanding of both existing and emerging cyber threats.
- 3) Realism: 4 - The finance sector's unique vulnerabilities to cyber-attacks, including sophisticated threats like APTs, phishing, and insider threats, suggest that Scenario 2 would prioritize realism, reflecting the nuanced and high-stakes environment of financial cybersecurity.
- 4) Creativity: 3 - Creativity in such scenarios is often demonstrated through the innovative presentation of threats and problem-solving approaches. However, the balance between creativity and realism is crucial to maintaining the scenario's educational value without straying too far from plausible threats.
- 5) Usability in Exercises: 4 - The practical application of the scenario in training exercises is vital. A scenario that effectively simulates real-world challenges while allowing participants to practice their response strategies would score highly in usability.
- 6) Expandability with Human Inputs: 4 - The ability to adapt and expand the scenario based on new information, threats, and organizational changes is essential for keeping cybersecurity exercises relevant and engaging, suggesting a high score for expandability.

The following figure 5 visualizes the evaluation of both scenarios by LLM1 and LLM2.

## VI. DETAILED IMPACT OF LIMITATIONS AND MITIGATION MEASURES

The limitations of our methodology, primarily the reliance on a narrow pool of experts and the constraints of the LLMs' training data, have significant implications for the generated scenarios. The input from only two experts, although highly knowledgeable, might not capture the full spectrum of potential cybersecurity challenges and biases that could arise in diverse operational contexts. This could lead to scenarios that, while effective in a controlled environment, may not fully encompass the complexity or variability of real-world cyber threats. Furthermore, the inherent biases and the limited scope of the training data for LLMs restrict the scenarios' ability to fully predict and simulate less common but potentially devastating cyber threats. These limitations can skew the scenario development towards more frequently represented threats in the training datasets, potentially overlooking emergent or niche threats that are equally critical.

To mitigate these risks, we plan to expand the circle of expertise by involving more professionals from various backgrounds in future scenario evaluations. This diversity will enrich the scenarios with broader insights and more robust validity checks, ensuring a more comprehensive representation of potential cyber threats. Additionally, we will enhance the training datasets for our LLMs with a more diverse range of cyber threat instances, including those that are less common but highly impactful. This approach aims to reduce the model's biases and improve its capacity to generate scenarios that are truly representative of the full spectrum of potential threats. By implementing these measures, we aim to refine the reliability and applicability of our LLM-generated scenarios, making them more robust and effective tools for cybersecurity training.

## VII. CONCLUSION AND FUTURE WORK

This study demonstrates that using LLMs to create cybersecurity exercise scenarios is both innovative and effective for training. The methodology effectively prepares security professionals to tackle a broad spectrum of cyber threats by simulating real-world challenges, encompassing both familiar and emerging cybersecurity issues. The integration of RAG adds depth and relevance to the scenarios, effectively utilizing LLM-generated content for more authentic simulations. Despite the critical scores in scenario 1, the overall quantitative measures provided us with good results. There are two main limitations of our work, the evaluation of our scenarios was limited to insights from only two experts. While their expertise is substantial, broader evaluations involving experts from diverse backgrounds such as legal, ethical, and various technological domains could provide a more comprehensive understanding of the scenarios' impacts and limitations. The second limitation is related to the AI models used, which are primarily LLMs in our case, constrained by their training data and inherent biases. This may limit the scenarios' effectiveness in predicting and

simulating more nuanced or less common cyber threats that are not well-represented in the training datasets. Future work should focus on further refining the integration of LLMs into cybersecurity training programs, particularly by expanding parameter coverage (LLMs with larger parameters) for more detailed scenarios and implementing generative configurations for increased creativity which will also assess the real-world impact of these scenarios on enhancing cybersecurity skills and readiness. Furthermore, we will work on methodologies for incorporating real-time threat intelligence and feedback mechanisms into the training modules, enabling them to evolve with the cybersecurity landscape.

## REFERENCES

- [1] T. Shejin and K. Sudheer, "A review on major cyber threats and recommended counter measures," *Int. J. Res. Appl. Sci. Eng. Technol.*, vol. 11, no. 3, pp. 1758–1761, Mar. 2023.
- [2] A. Sharma, B. B. Gupta, A. K. Singh, and V. K. Saraswat, "Advanced persistent threats (APT): Evolution, anatomy, attribution and countermeasures," *J. Ambient Intell. Hum. Comput.*, vol. 14, no. 7, pp. 9355–9381, Jul. 2023.
- [3] M. F. Safitra, M. Lubis, and H. Fakhurroja, "Counterattacking cyber threats: A framework for the future of cybersecurity," *Sustainability*, vol. 15, no. 18, p. 13369, Sep. 2023.
- [4] A. S. Albahri, I. A. Zahid, M. G. Yaseen, M. Aljanabi, A. H. Ali, and A. Kaleel, "Harnessing the tide of innovation: The dual faces of generative AI in applied sciences; letter to editor," *Appl. Data Sci. Anal.*, vol. 2024, pp. 1–3, Jan. 2024.
- [5] I. A. Zahid and S. S. Joudar, "Does lack of knowledge and hardship of information access signify powerful AI? A large language model perspective," *Appl. Data Sci. Anal.*, vol. 2023, pp. 150–154, Dec. 2023.
- [6] S. Saeed, S. A. Altamimi, N. A. Alkayyal, E. Alshehri, and D. A. Alabbad, "Digital transformation and cybersecurity challenges for businesses resilience: Issues and recommendations," *Sensors*, vol. 23, no. 15, p. 6666, Jul. 2023.
- [7] T.-A.-Q. Nguyen, E. Hashmi, M. M. Yamin, A. Beghdadi, F. A. Cheikh, and M. Ullah, "Benign paroxysmal positional vertigo disorders classification using eye tracking data," in *Proc. IFIP Int. Conf. Artif. Intell. Appl. Innov.* Springer, 2024, pp. 174–185.
- [8] J. N. Al-Karaki, M. Omar, A. Gawanmeh, and A. Jones, "Advancing CyberSecurity education and training: Practical case study of running capture the flag (CTF) on the metaverse vs. physical settings," in *Proc. Int. Conf. Intell. Metaverse Technol. Appl. (iMETA)*, vol. 6, Sep. 2023, pp. 1–7.
- [9] T. Awojana and T.-S. Chou, "Overview of learning cybersecurity through game based systems," in *Proc. CIEC*, 2019.
- [10] M. Mareš, R. Chytilík, Z. Špačková, J. Drmola, L. Hrbková, P. Mlejnková, and M. Tóth, "Assessment of performance during cybersecurity tabletop exercises," *Secur. J.*, pp. 1–24, Jul. 2023.
- [11] R. Beuran, D. Tang, C. Pham, K.-I. Chinen, Y. Tan, and Y. Shinoda, "Integrated framework for hands-on cybersecurity training: CyTrONE," *Comput. Secur.*, vol. 78, pp. 43–59, Sep. 2018.
- [12] A. Grimaldi, J. Ribiollet, P. Nespoli, and J. Garcia-Alfaro, "Toward next-generation cyber range: A comparative study of training platforms," in *Proc. Eur. Symp. Res. Comput. Secur.* Springer, 2023, pp. 271–290.
- [13] C. Chindrus and C.-F. Caruntu, "Securing the network: A red and blue cybersecurity competition case study," *Information*, vol. 14, no. 11, p. 587, Oct. 2023.
- [14] M. Katsantonis, P. Fouliras, and I. Mavridis, "Conceptual analysis of cyber security education based on live competitions," in *Proc. IEEE Global Eng. Educ. Conf. (EDUCON)*, Apr. 2017, pp. 771–779.
- [15] M. M. Yamin, B. Katt, and M. Nowostawski, "Serious games as a tool to model attack and defense scenarios for cyber-security exercises," *Comput. Secur.*, vol. 110, Nov. 2021, Art. no. 102450.
- [16] R. Nakata and A. Otsuka, "CyExec: Automatic generation of randomized cyber range scenarios," in *Proc. 7th Int. Conf. Inf. Syst. Secur. Privacy*, 2021, pp. 226–236.
- [17] K. Zenitani, "Attack graph analysis: An explanatory guide," *Comput. Secur.*, vol. 126, Mar. 2023, Art. no. 103081.

- [18] K. Zenitani, "A scalable algorithm for network reachability analysis with cyclic attack graphs," *J. Comput. Secur.*, vol. 31, no. 1, pp. 29–55, Jan. 2023.
- [19] S. Dalal, P. Manoharan, U. K. Lilhore, B. Seth, D. M. Alsekait, S. Simaiya, M. Hamdi, and K. Raahemifar, "Extremely boosted neural network for more accurate multi-stage cyber attack prediction in cloud computing environment," *J. Cloud Comput.*, vol. 12, no. 1, p. 14, Jan. 2023.
- [20] M. M. Yamin and B. Katt, "Modeling and executing cyber security exercise scenarios in cyber ranges," *Comput. Secur.*, vol. 116, May 2022, Art. no. 102635.
- [21] M. M. Yamin, B. Katt, and V. Gkioulos, "Cyber ranges and security testbeds: Scenarios, functions, tools and architecture," *Comput. Secur.*, vol. 88, Jan. 2020, Art. no. 101636.
- [22] M. M. Yamin and B. Katt, "Use of cyber attack and defense agents in cyber ranges: A case study," *Comput. Secur.*, vol. 122, Nov. 2022, Art. no. 102892.
- [23] F. Færøy, M. Yamin, A. Shukla, and B. Katt, "Automatic verification and execution of cyber attack on IoT devices," *Sensors*, vol. 23, no. 2, p. 733, Jan. 2023.
- [24] A. Zacharis and C. Patsakis, "AiCEF: An AI-assisted cyber exercise content generation framework using named entity recognition," *Int. J. Inf. Secur.*, vol. 22, no. 5, pp. 1333–1354, Oct. 2023.
- [25] T. Zaid and S. Garai, "Emerging trends in cybersecurity: A holistic view on current threats, assessing solutions, and pioneering new frontiers," *Blockchain Healthcare Today*, vol. 7, no. 1, Apr. 2024.
- [26] J. Wolfartsberger, R. Zimmermann, G. Obermeier, and D. Niedermayr, "Analyzing the potential of virtual reality-supported training for industrial assembly tasks," *Comput. Ind.*, vol. 147, May 2023, Art. no. 103838.
- [27] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proc. Adv. Neural Inf. Process. Syst.*, 2017.
- [28] B. Min, H. Ross, E. Sulem, A. P. B. Veyseh, T. H. Nguyen, O. Sainz, E. Agirre, I. Heintz, and D. Roth, "Recent advances in natural language processing via large pre-trained language models: A survey," *ACM Comput. Surv.*, vol. 56, no. 2, pp. 1–40, Feb. 2024.
- [29] E. Hashmi and S. Y. Yayilgan, "A robust hybrid approach with product context-aware learning and explainable AI for sentiment analysis in Amazon user reviews," *Electron. Commerce Res.*, pp. 1–33, Aug. 2024.
- [30] E. Hashmi and S. Y. Yayilgan, "Multi-class hate speech detection in the Norwegian language using FAST-RNN and multilingual fine-tuned transformers," *Complex Intell. Syst.*, vol. 10, no. 3, pp. 4535–4556, Jun. 2024.
- [31] H. Ali, E. Hashmi, S. Y. Yildirim, and S. Shaikh, "Analyzing Amazon products sentiment: A comparative study of machine and deep learning, and transformer-based techniques," *Electronics*, vol. 13, no. 7, p. 1305, Mar. 2024.
- [32] E. Hashmi, S. Y. Yayilgan, M. M. Yamin, S. Ali, and M. Abomhara, "Advancing fake news detection: Hybrid deep learning with FastText and explainable AI," *IEEE Access*, vol. 12, pp. 44462–44480, 2024.
- [33] S. Gheewala, S. Xu, S. Yeom, and S. Maqsood, "Exploiting deep transformer models in textual review based recommender systems," *Expert Syst. Appl.*, vol. 235, Jan. 2024, Art. no. 121120.
- [34] D. Oralbekova, O. Mamyrbayev, M. Othman, D. Kassymova, and K. Mukhsina, "Contemporary approaches in evolving language models," *Appl. Sci.*, vol. 13, no. 23, p. 12901, Dec. 2023.
- [35] S. Kernan Freire, C. Wang, and E. Niforatos, "Conversational assistants in knowledge-intensive contexts: An evaluation of LLM- versus intent-based systems," 2024, *arXiv:2402.04955*.
- [36] J. Lipenkova, "Overcoming the limitations of large language models how to enhance LLMS with human-like cognitive skills," Rep., 2023.
- [37] S. Moore, R. Tong, A. Singh, Z. Liu, X. Hu, Y. Lu, J. Liang, C. Cao, H. Khosravi, P. Denny, C. Brooks, and J. Stamper, "Empowering education with LLMS-the next-gen interface and content generation," in *Proc. Int. Conf. Artif. Intell. Educ.* Springer, 2023, pp. 32–37.
- [38] X. Liu, Z. Zhu, H. Liu, Y. Yuan, M. Cui, Q. Huang, J. Liang, Y. Cao, Q. Kong, M. D. Plumbley, and W. Wang, "Wavjourney: Compositional audio creation with large language models," 2307, *arXiv:2307.14335*.
- [39] L. Li, Y. Zhang, and L. Chen, "Prompt distillation for efficient LLM-based recommendation," in *Proc. 32nd ACM Int. Conf. Inf. Knowl. Manage.*, Oct. 2023, pp. 1348–1357.
- [40] Z. Yuan, J. Liu, Q. Zi, M. Liu, X. Peng, and Y. Lou, "Evaluating instruction-tuned large language models on code comprehension and generation," 2023, *arXiv:2308.01240*.
- [41] Z. Chu, H. Hao, X. Ouyang, S. Wang, Y. Wang, Y. Shen, J. Gu, Q. Cui, L. Li, S. Xue, J. Y. Zhang, and S. Li, "Leveraging large language models for pre-trained recommender systems," 2023, *arXiv:2308.10837*.
- [42] V. Rawte, A. Sheth, and A. Das, "A survey of hallucination in large foundation models," 2023, *arXiv:2309.05922*.
- [43] J.-Y. Yao, K.-P. Ning, Z.-H. Liu, M.-N. Ning, and L. Yuan, "LLM lies: Hallucinations are not bugs, but features as adversarial examples," 2023, *arXiv:2310.01469*.
- [44] P. Manakul, A. Liusie, and M. J. F. Gales, "SelfCheckGPT: Zero-resource black-box hallucination detection for generative large language models," 2023, *arXiv:2303.08896*.
- [45] F. Liu, K. Lin, L. Li, J. Wang, Y. Yacoob, and L. Wang, "Mitigating hallucination in large multi-modal models via robust instruction tuning," 2023, *arXiv:2306.14565*.
- [46] A. Martino, M. Iannelli, and C. Truong, "Knowledge injection to counter large language model (LLM) hallucination," in *Proc. Eur. Semantic Web Conf.* Springer, 2023, pp. 182–185.
- [47] Z. Ji, T. Yu, Y. Xu, N. Lee, E. Ishii, and P. Fung, "Towards mitigating LLM hallucination via self reflection," in *Proc. Findings Assoc. Comput. Linguistics*, 2023, pp. 1827–1843.
- [48] W. Sun, Z. Shi, S. Gao, P. Ren, M. de Rijke, and Z. Ren, "Contrastive learning reduces hallucination in conversations," in *Proc. AAAI Conf. Artif. Intell.*, vol. 37, 2023, pp. 13618–13626.
- [49] Y. Zhang and X. Chen, "Enhancing simplified Chinese poetry comprehension in llama-7b: A novel approach to mimic mixture of experts effect," Rep., 2023.
- [50] S. Curran, S. Lansley, and O. Bethell, "Hallucination is the last thing you need," 2023, *arXiv:2306.11520*.
- [51] L. Huang, W. Yu, W. Ma, W. Zhong, Z. Feng, H. Wang, Q. Chen, W. Peng, X. Feng, B. Qin, and T. Liu, "A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions," 2023, *arXiv:2311.05232*.
- [52] M. Platt and D. Platt, "Effectiveness of generative artificial intelligence for scientific content analysis," in *Proc. IEEE 17th Int. Conf. Appl. Inf. Commun. Technol. (AICT)*, vol. 5, Oct. 2023, pp. 1–4.
- [53] G. Chalvatzaki, A. Younes, D. Nandha, A. T. Le, L. F. R. Ribeiro, and I. Gurevych, "Learning to reason over scene graphs: A case study of finetuning GPT-2 into a robot language model for grounded task planning," *Frontiers Robot. AI*, vol. 10, Aug. 2023, Art. no. 1221739.
- [54] A. Holtzman, J. Buys, L. Du, M. Forbes, and Y. Choi, "The curious case of neural text degeneration," 2019, *arXiv:1904.09751*.
- [55] J. Hu, Q. Zhang, and H. Yin, "Augmenting greybox fuzzing with generative AI," 2023, *arXiv:2306.06782*.



**MUHAMMAD MUDASSAR YAMIN** is currently an Associate Professor with the Department of Information and Communication Technology, Norwegian University of Science and Technology (NTNU). He is also a member of the System Security Research Group and the focus of his research is on system security, penetration testing, security assessment, and intrusion detection. Before joining NTNU, he was an Information Security Consultant and served multiple government and private clients. He holds multiple cybersecurity certifications, such as OSCE, OSCP, LPT-MASTER, CEH, CHFI, CPTE, CISSO, and CBP. He is also serving as a Reviewer for *Nature*.





**EHTESHAM HASHMI** received the B.S. degree in computer science from the University of Central Punjab, Lahore Campus, in 2020, and the M.S. degree in computer science from COMSATS University Islamabad, Lahore Campus, in 2022. He is currently pursuing the Ph.D. degree with the Department of Information Security and Communication Technology (IHK), Norwegian University of Science and Technology (NTNU). From 2022 to 2023, he was a Lecturer with the Department of Computer Science, University of Lahore. His research interests include multilingual natural language processing, computational linguistics, generative AI with LLMs, knowledge graphs, and data mining.



**BASEL KATT** is currently a Professor with the Department of Information and Communication Technology, Norwegian University of Science and Technology, he holds the position of the Technical Project Leader of Norwegian cyber range. His research interests include information security, software security and security testing, software vulnerability analysis, and model-driven software development alongside model-driven security. Additionally, his expertise extends to access control, usage control, privacy protection, security monitoring, policies, languages, models, and their enforcement.

...



**MOHIB ULLAH** (Member, IEEE) is currently a Researcher with NTNU, where he is also involved in different research, management, teaching, and industrial projects. His research interests include medical imaging, crowd analysis, object segmentation, behavior classification, and tracking. In these research areas, he has published several high-impact peer-reviewed journals, conferences, and workshop articles. He served as a Program Committee Member for the International Workshop on Computer Vision in Sports (CVsports). He also served as the Chair for the Technical Program at European Workshop on Visual Information Processing. He is a Reviewer of well-reputed conferences and journals (*Neurocomputing* (Elsevier), *Multimedia Tools and Applications* (Springer), IEEE Access, the *Journal of Imaging*, IEEE CVPRw, IEEE ICIP, and IEEE AVSS).