

Towards Explainable Network Intrusion Detection using Large Language Models

Paul R. B. Houssel

p.houssel@uq.edu.au

The University of Queensland

Siamak Layeghy

siamak.layeghy@uq.net.au

The University of Queensland

Priyanka Singh

priyanka.singh@uq.edu.au

The University of Queensland

Marius Portmann

marius@ieee.org

The University of Queensland

ABSTRACT

Large Language Models (LLMs) have revolutionised natural language processing tasks, particularly as chat agents. However, their applicability to threat detection problems remains unclear. This paper examines the feasibility of employing LLMs as a Network Intrusion Detection System (NIDS), despite their high computational requirements, primarily for the sake of explainability. Furthermore, considerable resources have been invested in developing LLMs, and they may offer utility for NIDS. Current state-of-the-art NIDS rely on artificial benchmarking datasets, resulting in skewed performance when applied to real-world networking environments. Therefore, we compare the *GPT-4* and *LLama3* models against traditional architectures and transformer-based models to assess their ability to detect malicious NetFlows without depending on artificially skewed datasets, but solely on their vast pre-trained acquired knowledge. Our results reveal that, although LLMs struggle with precise attack detection, they hold significant potential for a path towards explainable NIDS. Our preliminary exploration shows that LLMs are unfit for the detection of Malicious NetFlows. Most promisingly, however, these exhibit significant potential as complementary agents in NIDS, particularly in providing explanations and aiding in threat response when integrated with Retrieval Augmented Generation (RAG) and function calling capabilities.

1 INTRODUCTION

Large Language Models (LLMs) have recently emerged as a transformative force in natural language processing, showcasing proficiency in handling unstructured data. With the ability to generate text and comprehend context, LLMs have found widespread application in various domains, ranging from conversational chat agents to code generation and translation tasks. Despite these advances, the application of LLMs to cybersecurity, particularly Network Intrusion Detection Systems (NIDS), remains largely unexplored.

NIDS are critical for monitoring and analysing network traffic to detect malicious activities and security breaches. Existing NIDS employ a combination of signature-based and anomaly-based detection techniques. Signature-based methods rely on predefined patterns of known attacks, while anomaly-based methods identify deviations from established norms in network behaviour. Although deep learning-based NIDS have demonstrated near-perfect performance on benchmark datasets, their efficacy in real-world environments is often questioned due to the prevalence of synthetic and misrepresentative datasets used in academic research [1]. Additionally, these models lack explainability [2], making it difficult for security professionals to understand and trust their predictions and respond to the detected threat. It thus remains important to investigate if LLMs are a potential solution to these.

Several studies have used transformers and LLMs to directly detect threats or as complementary solutions for NIDS. However, no previous work has thoroughly investigated LLMs' adaptivity to the domain of NetFlows. Previous studies replace LLM's sequence-to-sequence layers with a classification head, thus losing the potential for using the same model to provide explanations. Using LLMs and their natural language output, both detection by classification and explainability can be achieved. It remains to be seen if performance is degraded when using these text prediction models as threat classifiers, a use case very different from their original design purpose. LLMs have shown the ability to analyse long-term dependencies and be adaptive to the applied domains [3]. Furthermore, their pre-training on large amounts of natural language data might incorporate the knowledge required to detect malicious NetFlows without relying on skewed artificial datasets. This paper explores the capability of LLMs to detect network attacks based on network flows while contrasting their benefits and limitations with traditional Machine Learning (ML) solutions.

We empirically assess the pre-trained OpenAI's *GPT-4* and Meta's *LLama3* models, in the context of zero-shot learning

and fine-tuning on standardised NetFlow datasets. Our initial experiments investigate the performance capabilities of LLMs to detect threats, a trait which traditional more lightweight ML solutions have already been proven to possess with near perfection [4]. We will evaluate our model variants on the *NF-UNSW-NB15-v2* and *NF-CSE-CIC-IDS2018-v2* datasets [5]. Finally, we investigate the capacity to explain the detection of threats and their computational cost of inference to fuel a discussion on the advantages and drawbacks of using LLMs for NIDS versus more lightweight solutions incorporating LLMs as a complementary solution to existing NIDS. Our findings aim to provide a first exploratory evaluation of LLMs in the context of NIDS, offering insights into potential future research directions and use cases. We do not aim for a comprehensive benchmark of models and parameters of LLMs for detecting malicious NetFlows.

2 RELATED WORK

2.1 Malicious NetFlow Detection

In contrast to NIDS, Host-based Intrusion Detection Systems (HIDS) typically deal with textual data such as system calls, logs, or application and memory traces. Natural Language Processing (NLP) techniques are employed to examine the semantic interactions and sequences of system calls and contextual data, enabling the effective identification of suspicious activities. This facilitates precise, real-time intrusion detection by analysing text data streams. Many solutions relying on NLP have been developed even before the advent of LLMs [6]. Natural language is less present in the NIDS domain, but surprisingly more literature exists.

For detection based on network traffic, the application of LLMs has started to gain traction, although the solutions often modify the standard LLM architecture. Specifically, these models typically strip the sequence-to-sequence head and replace it with a classification head to obtain a deterministic output to evaluate the model. The initial textual output head, which predicts the next textual token, is replaced by a classification head vector which has a dimension equivalent to the number of prediction classes, e.g two for a binary classification. Lira et al. [7] investigated the capabilities of LLMs such as *GPT-3.5*, *GPT-4*, and *ADA* to detect Distributed Denial of Service (DDoS) attacks on the *CIC-IDS2017* [8] and *Urban IoT* [9] datasets. Their findings indicate that LLMs, when fine-tuned or employed in few-shot learning scenarios, could detect DDoS attacks with significant accuracy proving their adaptivity trait. Manocchio et al. [10] explored various encoding and classification methods on different transformer architectures, including *GPT*'s and *BERT*'s to enhance network analysis. Ferrag et al. [11] used *BERT* for threat detection in Internet of Things (IoT)

networks, outperforming existing ML solutions by distinguishing between 14 attack classes and benign traffic using *PCAP* files. Alkhatib et al. [12] explored the usage of *BERT* for detecting threats on Controller Area Networks (CAN), showcasing the model's versatility across different cybersecurity applications. Among these, Chen et al. [13] is the only study which employed LLMs, which are pre-trained on natural language and possess a sequence-to-sequence head producing sequences of natural text tokens, to detect DDoS attacks by fine-tuning *BERT*.

2.2 Explainable NIDS

Studies aiming to render NIDS explainable typically employ existing explainable AI methods. These methods assign an importance factor to each input feature, indicating their influence on the model's output label [2]. For instance, Mallampati et al. [14], Senevirathna et al. [15], and Zebin et al. [16] utilise the SHapley Additive exPlanation (SHAP) method across various network domains. This approach expresses explainability as the weight of each feature and its impact on the output class, enabling the development of algorithms to respond to threats in real time. Wei et al. [17] propose a framework that leverages this feature's importance to design rules for active intrusion responses. Their XNIDS framework identifies features critical to the prediction class and generates actionable defence rules, which aid in understanding DL-NIDS behaviours and troubleshooting detection errors.

LLMs could potentially provide natural language explanations for human operators, not only limiting the explanations to features which statistically appear anomalous to the NIDS. This will allow to facilitate threat response and a better understanding of their NIDS functionality or potential malfunctioning. In this study, we further investigate the capability of LLMs to detect malicious NetFlows while explaining their reasoning. These could potentially be employed as a complementary solution, aiming to enhance the explainability of NIDS.

3 METHODOLOGY

This section outlines the methodology employed to investigate the potential of LLMs for NIDS. Our approach is based on the open-source *LLama3-8B-Instruct* and *GPT-4* model, designed to better handle instruction-based tasks compared to the foundational model. In contrast with foundational model's difficulties in answering in the required format, as it was not trained to understand instructions but simply to predict sequences of words. We focus on the network traffic under the format of flows, which are currently the most prominent and efficient source of network telemetry for threat

detection due to their scalability [1]. They summarise network traffic which shares destination and source port and IP addresses by aggregating packets into flows, making them ideal for detecting patterns indicative of malicious activities. Sarhan et al. [5] proposes a set of standardised NetFlow datasets, with a common set of overlapping features, to address the limited ability to evaluate ML-based NIDS across different network environments and attack scenarios. Among these, we solely consider the datasets outside of the IoT domain: *NF-UNSW-NB15-v2* and *NF-CSE-CIC-IDS2018-v2*.

3.1 Domain Specific Performance

To measure the performance on specific network environments, we built two baseline models using the *LLama3-8B-Instruct* and *GPT-4-0613* model by applying the zero-shot learning paradigm. Additionally, we will investigate two fine-tuning methods only on the *LLama3* model due to the inherited cost associated with fine-tuning using the *OpenAI* API. For this series of experiments, all datasets are split into training and testing sets with a 95% to 5% ratio while stratifying the attack type. All evaluations are performed following a 10-fold cross-validation method. We only consider a small testing set due to the sheer computational cost of LLM inference. All prediction results from our experiments are evaluated using the macro average of Precision and Recall for both the benign and malicious classes. The macro average is the unweighted sum of the individual performance scores for each class, representing each class equally. We provide the pre-trained LLM with an instruction to detect network attacks based on NetFlows, requesting it to output either "1" or "0" to indicate whether the provided flow is malicious or benign. After ingesting the instruction prompt, the model will take as an input a tabular NetFlow entry encoded as text using key-value pairs separated by commas to represent the feature name and value pairs. For instance, a network flow originally represented as a row within the CSV table of the dataset is transformed into text as follows: "*L4_DST_PORT: 80*". To ensure consistency and reduce variability in outputs, we set the model's temperature parameter to 10^{-1} . The temperature parameter, ranging from 0 to 1, influences the probabilities generated by the softmax output function. A low-temperature setting makes the model's outputs more deterministic by prioritising the most likely tokens, important for binary classification tasks [18]. Together with prompt engineering, it allows us to constrain the response of the model to a binary output without any exceptions.

Zero-shot Learning. Zero-shot learning consists of a pre-trained model tasked with making predictions on yet explicitly unseen data. In the context of a pre-trained LLM,

zero-shot learning leverages the knowledge and understanding embedded in the model from its training on a vast amount of natural language. When applied to zero-shot learning, the LLM can make informed predictions to identify malicious network traffic by understanding the provided instructions and applying its broad contextual knowledge to classify the traffic as either malicious or benign, even though it hasn't been specifically trained on labelled examples of network attacks.

Fine-tuning. Fine-tuning LLMs for specific tasks is a crucial step to enhance their performance and adapt them to domain-specific requirements. In our investigation, we focus on two state-of-the-art fine-tuning methods: Odds Ratio Preference Optimisation (ORPO) [19] and Kahneman-Tversky Optimisation (KTO) [20]. These methods have shown significant promise and represent the current state of the art [18]. ORPO [19], requires preference data in the format of prompts paired with rejected and accepted answers. For the task of NIDS binary classification, we can adapt these methods by treating the true label as the accepted prompt and the inverse of the label as the rejected prompt. In contrast, KTO [20] does not require paired preference data and instead relies on a series of prompts and corresponding responses along with a series of binary labels indicating wherever the provided response is relevant. In our case, we use the true label as the accepted response. Half of the prompts are randomly sampled with relevant responses using the label and the other half are samples as irrelevant by using the inverse of the label.

3.2 Explainability

To investigate and grasp the LLM's ability to explain the reasons for predicting a network flow as benign or malicious, we manually analyse the models' responses to find patterns and inconsistencies. More specifically, we analysed the explanations provided by our baseline model for the cases of True Positives (TP), False Positives (FP), True Negatives (TN), and False Negatives (FN) samples, as these will offer variability in answering. This analysis is not a quantitative evaluation of the quality of explanations but serves as an initial exploration to assess the importance of future research into making LLMs explainable for NIDS.

3.3 Computational Complexity

To assess the computational cost, we compared our baseline with the most common architectures with the default parameters, used in related work and industrial application of NIDS. All experiments were conducted on a machine equipped with an NVIDIA RTX3090 GPU and a 3.50 GHz CPU. The measured time for sequential inference corresponds to the average time taken to predict a single NetFlow

sample. This time was measured by calculating the average time to predict a single sample when processing the entire testing set as a batch, computing the average over 10 runs.

4 EMPIRICAL RESULTS

4.1 Domain Specific Performance

Zero-shot Learning. The results of these empirical experiments under the zero-shot evaluation (Table 1 and 2) reveal several important insights. Unsurprisingly, the performance of the pre-trained *LLama3* model is worse than random selection, which would have been, considering the macro average, 50% for all metrics, while *GPT-4* is slightly better. This outcome is expected because pre-trained LLMs are primarily designed for natural language tasks and not for the specific task of network classification. In contrast, classical ML models demonstrate near-perfect performance. For instance, on the *NF-UNSW-NB15-v2* dataset, a Random Forest and an LSTM neural network respectively achieve a weighted F1-score of 92.17% and 92.82% [21] while a *DANN* model achieves 97.81% [22].

Table 1: Domain-specific zero-shot performance of the model on the *NF-UNSW-NB15-v2* dataset.

Model	Macro Avg. Perf. Metrics↑	
	Precision	Recall
<i>LLama3-8b-Instruct</i>	48.56%	40.56%
<i>GPT-4</i>	50.85%	53.47%

Table 2: Domain specific zero-shot performance on the *NF-CSE-CIC-IDS2018-v2* dataset.

Model	Macro Avg. Perf. Metrics↑	
	Precision	Recall
<i>LLama3-8b-Instruct</i>	49.66%	49.73%
<i>GPT-4</i>	50.25%	51.02%

Fine-tuning. In comparison to the zero-shot learning results, both fine-tuning methods showed slight improvements but still failed to significantly outperform random selection or approach the performance of related works (Table. 3 and 4). Overall, *KTO* performed slightly better, indicating that the choice of fine-tuning method is an important parameter for future research. We see that fine-tuning can enhance performance with increased training, but the improvements observed were minimal. Fine-tuning tends to adjust the model’s responses to fit a desired style but

is insufficient to significantly alter the weights for the binary classification of network attacks. Pre-trained natural language models, which are trained on extensive natural language corpora, may not be capable of recognising the subtle patterns that indicate malicious NetFlows. They are more adept at generating appropriate responses to natural language queries without any further deep understanding.

Fine-tuning LLMs for detecting malicious NetFlows is unnecessary; instead, transformer-based models can be trained directly on entire datasets. Manocchio et al., through their proposed *FlowTransformer* framework [10], have shown that transformer architectures perform well on binary NetFlow classification tasks when trained on complete datasets. This performance is achieved without the need pre-training on natural language datasets. For instance, their model, using GPT’s architecture achieves an 89.98% F1-score.

Table 3: Domain specific performance of the *LLama3-8b-Instruct* model on the *NF-UNSW-NB15-v2* dataset.

Learning Paradigm	# Training Samples	Macro Avg. Perf. Metrics↑
		Precision
		Recall
zero-shot	0	48.56%
	1,000	52.69%
<i>KTO</i> Finetuned	10,000	53.98%
	50,000	55.18%
	1,000	48.88%
<i>ORPO</i> Finetuned	10,000	49.25%
	50,000	52.00%
		40.56%
		50.22%
		51.67%
		40.76%
		42.40%
		43.50%

Table 4: Domain specific performance of the *LLama3-8b-Instruct* model on the *NF-CSE-CIC-IDS2018-v2* dataset.

Learning Paradigm	# Training Samples	Macro Avg. Perf. Metrics↑
		Precision
		Recall
zero-shot	0	49.66%
	1,000	50.02%
<i>KTO</i> Finetuned	10,000	52.22%
	50,000	55.02%
	1,000	50.43%
<i>ORPO</i> Finetuned	10,000	50.89%
	50,000	53.60%
		49.73%
		50.09%
		55.13%
		58.79%
		50.95%
		51.01%
		51.58%

4.2 Explainability

The analysis of LLama3 generated explanations for various prediction cases reveal unique benefits and notable limitations to enrich future NIDS applications with explainability, a unique feature proper to LLMs. These sample prompts demonstrate that LLMs can provide detailed and contextually relevant explanations for their classifications of network flows but with notable limitations.

When asked about a TN sample (Fig. 1), the model correctly identifies the network flow as a DNS query. However, it inaccurately claims the source IP address is from China, whereas it is actually from Japan, and misidentifies the destination address as being in the USA, although it belongs to the University of New South Wales (UNSW), the dataset’s creator. Despite these errors, the LLM correctly identifies the byte size range of all packets in the flow as typical for a DNS query. These inaccuracies highlight the LLM’s tendency to hallucinate or generate incorrect information, particularly regarding geographical locations.

On a TP case, the LLM mistakenly states that protocol 139 corresponds to the *NetBIOS* protocol, whereas it corresponds to the *Host Identity Protocol*. Its confusion is evident as NetBIOS uses port 139. Most importantly, it operates at Layer 5, while the *protocol* feature identifies those of Layer 4. Additionally, the LLM incorrectly assumes that the use of ports with the value 0 indicates malicious activity. While unusual, ports set to 0 can occur in protocols not using port numbers, like *ICMP*. These errors show that the LLM has gaps in its understanding of protocol specifics and port usage, leading to incorrect conclusions about malicious behaviour.

In the FP scenario, unusual patterns, such as atypical ports and large file transfers are identified leading to the conclusion that these imply maliciousness. While the reasoning is plausible, it cannot correlate these features to specific attack signatures accurately.

For the FN sample, the model’s explanation is mostly consistent. However, it incorrectly identifies the main Layer 7 protocol user and mistakenly states that Layer 7 protocol number 7 corresponds to HTTP when it indicates the Internet Printing Protocol (IPP). The model classifies the network flow as benign due to small packet lengths and a low number of retransmitted packets, correctly identifying the traffic as typical HTTP traffic. While the model can accurately explain benign traffic by recognising common patterns, it fails to identify the malicious aspects of the sample.

- Protocol: The protocol used is UDP (17), which is commonly used for DNS queries. This suggests that the flow may be related to a legitimate DNS request.
- Destination port: The destination port is 53, which is the standard port for DNS traffic.
- Source and destination addresses: The source address is an IP address in China, while the destination address is an IP address in the United States. This could indicate that the flow is related to a legitimate DNS query from a Chinese user to a US-based DNS server.
- Packet sizes and lengths: The packet sizes are relatively small (73-89 bytes), which is typical for DNS.
- Flow duration: The flow duration is 0 milliseconds, indicating that this was likely a single, brief request-response exchange.

Figure 1: Explanation provided for a True Negative NetFlow.

Overall, from these manual experiments, we learned that LLM’s explanations are convincing and correctly retrieve information from the NetFlow data, being consistent with the presented features. However, they often fail to reason logically and augment the information with factual accuracy. The LLM tends to puzzle feature values together to identify traffic types (e.g., *HTTP* traffic, *DNS* queries) effectively, which can be useful for network operators. However, it struggles to correlate multiple features to identify attack signatures accurately, instead treating individual suspicious features as independent clues of malicious activity. One significant problem is their tendency to generate nonsensical or unfaithful content, known as “hallucination” [18]. These models lack a true comprehension of facts or logical reasoning, relying on learned patterns of word succession within context. This is even more evident in further prompts. When questioned about the units of measurement used in network flows, the LLM failed to provide a coherent answer. When provided with a malformed NetFlow, the LLM was unable to determine its soundness. Additionally, the LLM failed to explain a SYN flood attack, a straightforward task for signature-based detectors like *Cisco Snort*.

4.3 Computational Complexity

The feasibility of using LLMs for NIDS not only relies on detection performance but as importantly on inference time. Given that NIDS must handle hundreds, if not thousands, of NetFlows per second, these models must keep up with the network traffic and be able to report incidents promptly. Delays in inference can result in missed or late alerts, compromising network security.

Table 5: Complexity of benchmarked model compared to traditional NIDS architectures.

Method	Inference Time (μs)	#Parameters
<i>LLama3-8B-Instruct</i>	14,000.01	8,030,000,000
Random Forest	2.03	867,614 (nodes)
Decision Tree [23]	1.53	—
LSTM NN	25.02	56,555
DANN	28.40	67,208

Our comparison in Table. 5, highlights the significant difference in computational complexity among the models. The *LLama3* model exhibited significantly higher inference times due to its complex architecture and a large number of parameters. As *GPT-4* is closed-source and runs on the server side, we were unable to obtain statistics about its inference times and parameter count, which are in a similar range to those of *LLama3*. In contrast, traditional ML models like the Random Forest and the Decision Tree demonstrated swift inference times, making them more suitable for real-time NIDS applications. For instance, DL models can still offer a reasonable trade-off compromise, with costs potentially being manageable. However, considering the near 7000-fold slower inference time of LLMs compared to light-weight models, it is not feasible to implement LLMs in a real-time NIDS.

5 DISCUSSION

Current fine-tuning approaches have not significantly improved performance, highlighting the need for more effective methods tailored to classification tasks. Parameters such as temperature, data encoding and prompt formulation can also affect performance and need to be studied comprehensively for their impact on classification. Furthermore, as demonstrated by the cross-domain performance limitations of existing NIDS models [21], there is a need for more robust and adaptable architectures. Comparing LLMs, known for their adaptability, with other ML architectures in terms of their performance on unseen network behaviours could offer insights for future development in this field.

To improve explainability, the integration of a Retrieval-Augmented Generation (RAG) system [24], linked with a Cyber Threat Intelligence (CTI) source, could significantly enhance the reliability of LLMs. RAG can ground LLMs in accurate, up-to-date information, potentially reducing or eliminating hallucinations. Expanding RAG beyond CTI to include specifications of observed systems, such as endpoints, operating systems, SBOM, and firewall rules, could enhance both detection rates and response effectiveness. On the other hand, to provide explainable threat response, deployed LLMs should be equipped with access to function

calls. Function calling, also known as tool use or API calling, enables LLMs to interact with external systems [25]. By providing the LLM with a predefined set of functions or tools along with usage instructions, it can intelligently select and invoke appropriate functions to address a given threat scenario. For instance, it could propose specific functions with parameters to temporarily modify firewall rules. To evaluate the progress of future research in this area, it is essential to develop a methodology for quantitatively assessing the explainability of LLMs to reason on NIDS alerts. Currently, datasets necessary to fuel this methodology do not exist, highlighting future research.

While LLMs show promise as assistants for providing explanations, they require further research and development before they can be considered reliable in this capacity. Future work shall also focus on LLMs' understanding of network traffic, as it is crucial for their potential role as a virtual system administrator, complementing NIDS. This research would build upon the work of Donadel et al. [26], evaluating their understanding of network topology.

6 CONCLUSION

Our investigation into the application of LLMs for NIDS has revealed notable limitations in their performance as primary solutions. If used as a solution to detect threats based on NetFlows, future research shall rather focus on transformer-based architectures without relying on pre-trained LLMs. LLMs currently struggle to effectively detect domain-specific malicious NetFlows and have high inference time complexity, making them impractical as standalone NIDS solutions. Given these findings, we propose that LLMs should be harnessed as complementary solutions to existing state-of-the-art NIDS. Specifically, their strengths in providing explainability make them valuable when integrated with traditional ML-based systems, particularly for generating detailed and contextually relevant explanations when alerts are raised. This approach leverages the potential of LLMs to enhance the interpretability and transparency of NIDS, thus aiding in threat response and decision-making processes. While challenges such as hallucination and logical reasoning limitations persist, the promising results in explainability underscore the importance of further exploring LLMs to finally render NIDS explainable.

REFERENCES

- [1] S. Layeghy, M. Gallagher, and M. Portmann, "Benchmarking the Benchmark – Analysis of Synthetic NIDS Datasets," Apr. 2021.
- [2] S. Neupane, J. Ables, W. Anderson, S. Mittal, S. Rahimi, I. Banicescu, and M. Seale, "Explainable Intrusion Detection Systems (X-IDS): A Survey of Current Methods, Challenges, and Opportunities," *IEEE Access*, vol. 10, pp. 112392–112415, 2022.
- [3] H. Xu, S. Wang, N. Li, K. Wang, Y. Zhao, K. Chen, T. Yu, Y. Liu, and H. Wang, "Large Language Models for Cyber Security: A Systematic Literature Review," *arXiv:2405.04760*, 2024.
- [4] M. Sarhan, S. Layeghy, N. Moustafa, and M. Portmann, "NetFlow Datasets for Machine Learning-Based Network Intrusion Detection Systems," in *Big Data Technologies and Applications* (Z. Deze, H. Huang, R. Hou, S. Rho, and N. Chilamkurti, eds.), (Cham), pp. 117–135, Springer International Publishing, 2021.
- [5] M. Sarhan, S. Layeghy, and M. Portmann, "Towards a Standard Feature Set for Network Intrusion Detection System Datasets," *Mobile Networks and Applications*, vol. 27, pp. 357–370, Feb. 2022.
- [6] Z. T. Sworna, Z. Mousavi, and M. A. Babar, "NLP Methods in Host-based Intrusion Detection Systems: A Systematic Review and Future Directions," Nov. 2022. *arXiv:2201.08066* [cs].
- [7] O. G. Lira, A. Marroquin, and M. A. To, "Harnessing the Advanced Capabilities of LLM for Adaptive Intrusion Detection Systems," in *Advanced Information Networking and Applications* (L. Barolli, ed.), (Cham), pp. 453–464, Springer Nature Switzerland, 2024.
- [8] I. Sharafaldin, A. Habibi Lashkari, and A. A. Ghorbani, "Toward Generating a New Intrusion Detection Dataset and Intrusion Traffic Characterization," in *Proceedings of the 4th International Conference on Information Systems Security and Privacy*, (Funchal, Madeira, Portugal), pp. 108–116, SCITEPRESS - Science and Technology Publications, 2018.
- [9] A. Hekmati, N. Jethwa, E. Grippo, and B. Krishnamachari, "Correlation-Aware Neural Networks for DDoS Attack Detection In IoT Systems," Feb. 2023. *arXiv:2302.07982* [cs].
- [10] L. D. Manocchio, S. Layeghy, W. W. Lo, G. K. Kulatilleke, M. Sarhan, and M. Portmann, "FlowTransformer: A transformer framework for flow-based network intrusion detection systems," *Expert Systems with Applications*, vol. 241, p. 122564, May 2024.
- [11] M. A. Ferrag, M. Ndhlovu, N. Tihanyi, L. C. Cordeiro, M. Debbah, T. Lestable, and N. S. Thandi, "Revolutionizing Cyber Threat Detection with Large Language Models: A privacy-preserving BERT-based Lightweight Model for IoT/IoT Devices," Feb. 2024. *arXiv:2306.14263* [cs].
- [12] N. Alkhateib, M. Mushtaq, H. Ghauch, and J.-L. Danger, "CAN-BERT do it? Controller Area Network Intrusion Detection System based on BERT Language Model," in *2022 IEEE/ACS 19th International Conference on Computer Systems and Applications (AICCSA)*, (Abu Dhabi, United Arab Emirates), pp. 1–8, IEEE, Dec. 2022.
- [13] M. Guastalla, Y. Li, A. Hekmati, and B. Krishnamachari, "Application of Large Language Models to DDoS Attack Detection," in *Security and Privacy in Cyber-Physical Systems and Smart Vehicles*, vol. 552, pp. 83–99, Cham: Springer Nature Switzerland, 2024. Series Title: Lecture Notes of the Institute for Computer Sciences, Social Informatics and Telecommunications Engineering.
- [14] S. B. Mallampati and H. Seetha, "Enhancing Intrusion Detection with Explainable AI: A Transparent Approach to Network Security," *Cybernetics and Information Technologies*, vol. 24, pp. 98–117, 2024.
- [15] T. Senevirathna, B. Siniarski, M. Liyanage, and S. Wang, "Deceiving Post-Hoc Explainable AI (XAI) Methods in Network Intrusion Detection," in *2024 IEEE 21st Consumer Communications & Networking Conference (CCNC)*, pp. 107–112, Jan. 2024. ISSN: 2331-9860.
- [16] T. Zebin, S. Rezvy, and Y. Luo, "An Explainable AI-Based Intrusion Detection System for DNS Over HTTPS (DoH) Attacks," *IEEE Transactions on Information Forensics and Security*, vol. 17, pp. 2339–2349, 2022. Conference Name: IEEE Transactions on Information Forensics and Security.
- [17] F. Wei, H. Li, Z. Zhao, and H. Hu, "XNIDS: Explaining Deep Learning-based Network Intrusion Detection Systems for Active Intrusion Responses," in *Proceedings of the 32nd USENIX Security Symposium*, (Anaheim, CA, USA), 2023.
- [18] S. Minaee, T. Mikolov, N. Nikzad, M. Chenaghlu, R. Socher, X. Amatriain, and J. Gao, "Large Language Models: A Survey," Feb. 2024. *arXiv:2402.06196* [cs].
- [19] J. Hong, N. Lee, and J. Thorne, "ORPO: Monolithic Preference Optimization without Reference Model," 2024. *arXiv:2403.07691* [cs].
- [20] K. Ethayarajh, W. Xu, N. Muennighoff, D. Jurafsky, and D. Kiela, "KTO: Model Alignment as Prospect Theoretic Optimization," *arXiv:2402.01306*, 2024.
- [21] S. Layeghy and M. Portmann, "Explainable Cross-domain Evaluation of ML-based Network Intrusion Detection Systems," *Computers and Electrical Engineering*, vol. 108, p. 108692, May 2023.
- [22] S. Layeghy, M. Baktashmotlagh, and M. Portmann, "DI-NIDS: Domain invariant network intrusion detection system," *Knowledge-Based Systems*, vol. 273, p. 110626, Aug. 2023.
- [23] L. D. Manocchio, S. Layeghy, and M. Portmann, "Network Intrusion Detection System in a Light Bulb," in *2022 32nd International Telecommunication Networks and Applications Conference (ITNAC)*, pp. 1–8, IEEE, Nov. 2022. ISSN: 2474-154X.
- [24] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel, S. Riedel, and D. Kiela, "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks," in *Advances in Neural Information Processing Systems*, vol. 33, (Online), pp. 9459–9474, Curran Associates, Inc., 2020.
- [25] Y. Liang, C. Wu, T. Song, W. Wu, Y. Xia, Y. Liu, Y. Ou, S. Lu, L. Ji, S. Mao, Y. Wang, L. Shou, M. Gong, and N. Duan, "TaskMatrix.AI: Completing Tasks by Connecting Foundation Models with Millions of APIs," Mar. 2023. *arXiv:2303.16434* [cs].
- [26] D. Donadel, F. Marchiori, L. Pajola, and M. Conti, "Can LLMs Understand Computer Networks? Towards a Virtual System Administrator," Apr. 2024. *arXiv:2404.12689* [cs].