

Large Language Models for Network Intrusion Detection Systems: Foundations, Implementations, and Future Directions

Shuo Yang, Xinran Zheng, Xincheng Zhang, Jinfeng Xu, Jinze Li,
Donglin Xie, Weicai Long and Edith C.H. Ngai*, *Senior Member, IEEE*

Abstract—Large Language Models (LLMs) have revolutionized various fields with their exceptional capabilities in understanding, processing, and generating human-like text. This paper investigates the potential of LLMs in advancing Network Intrusion Detection Systems (NIDS), analyzing current challenges, methodologies, and future opportunities. It begins by establishing a foundational understanding of NIDS and LLMs, exploring the enabling technologies that bridge the gap between intelligent and cognitive systems in AI-driven NIDS. While Intelligent NIDS leverage machine learning and deep learning to detect threats based on learned patterns, they often lack contextual awareness and explainability. In contrast, Cognitive NIDS integrate LLMs to process both structured and unstructured security data, enabling deeper contextual reasoning, explainable decision-making, and automated response for intrusion behaviors. Practical implementations are then detailed, highlighting LLMs as processors, detectors, and explainers within a comprehensive AI-driven NIDS pipeline. Furthermore, the concept of an LLM-centered Controller is proposed, emphasizing its potential to coordinate intrusion detection workflows, optimizing tool collaboration and system performance. Finally, this paper identifies critical challenges and opportunities, aiming to foster innovation in developing reliable, adaptive, and explainable NIDS. By presenting the transformative potential of LLMs, this paper seeks to inspire advancement in next-generation network security systems.

Index Terms—Large Language Models, Intrusion Detection Systems, Network Security, Generative Artificial Intelligence.

I. INTRODUCTION

The growing complexity and scale of modern network infrastructures have led to an exponential rise in cyber threats. Network Intrusion Detection Systems (NIDS) play a critical role in safeguarding these infrastructures by monitoring and analyzing network traffic for suspicious activities. However, traditional NIDS, which rely on predefined signatures or statistical methods, often struggle to detect sophisticated or novel attacks. The integration of Artificial Intelligence (AI) has significantly enhanced the capabilities of NIDS, enabling more intelligent detection mechanisms. Among these advancements, Large Language Models (LLMs) have emerged as promising tools due to their unparalleled abilities to understand, process, and generate human-like text.

LLMs, such as GPT-4¹, LLaMA², and DeepSeek³, have demonstrated remarkable success in diverse applications, ranging from natural language processing (NLP) to reasoning and decision-making tasks. Their ability to extract meaningful insights from vast datasets makes them well-suited for the complex and dynamic nature of intrusion detection. By leveraging LLMs, NIDS can transition from intelligent systems to cognitive systems capable of contextual reasoning, offering faster response and enhanced explainability.

Despite these promising capabilities, research and applications of NIDS are still in their early stages, and systematic exploration remains limited. This paper aims to bridge this gap by investigating the transformative role of LLMs in advancing NIDS. Specifically, we establish a foundational understanding of NIDS and LLMs, tracing their evolution and exploring their potential synergies. We then examine how LLMs can be integrated into key stages of the NIDS pipeline, including data, model, and response. Furthermore, we propose LLM-centered Controller, a novel architecture for orchestrating NIDS operations by coordinating various tools and components to enhance efficiency and effectiveness.

However, integrating LLMs into NIDS presents several challenges, including validity, complexity, and privacy. For instance, LLMs may generate false positives due to biases in training data, struggle with inconsistent threat assessments, and even produce misleading outputs that could hinder security responses. Additionally, their high computational demands would impact real-time processing capabilities, while the analysis of sensitive network traffic raises significant privacy risks. This paper identifies these challenges and outlines future research directions, focusing on enhanced multimodal integration, real-time analysis detection, privacy-preserving collaboration, and multi-agent systems. Through a comprehensive review, we aim to catalyze innovation and development in next-generation NIDS, paving the way for more reliable, adaptive, and explainable network security systems.

The rest of this paper is organized as follows. Section II provides a foundational background on NIDS and LLMs. Section III examines the transition of AI-driven NIDS from intelligent to cognitive systems, focusing on the general pipeline, existing gaps, and enabling technologies. Section IV discusses the practical implementation of NIDS with LLMs, including LLM-enhanced Processor, LLM-based Detector, LLM-driven Explainer, and LLM-centered Controller. Section V outlines

S. Yang, X. Zhang, J. Xu, J. Li and E. C.H. Ngai are with Department of Electrical and Electronic Engineering, The University of Hong Kong, Hong Kong ASR, China (shuo.yang@connect.hku.hk).

X. Zheng is with Department of Electronic Engineering, Tsinghua University, Beijing, China. D. Xie is with National Institute of Health Data Science, Peking University, Beijing, China. W. Long is with Data Science and Analytics Thrust, Hong Kong University of Science and Technology (Guangzhou), Guangzhou, China.

* Corresponding author: Edith C.H. Ngai (chngai@eee.hku.hk).

¹GPT-4: <https://openai.com/index/gpt-4>

²LLaMA: <https://github.com/meta-llama>

³DeepSeek: <https://www.deepseek.com/>

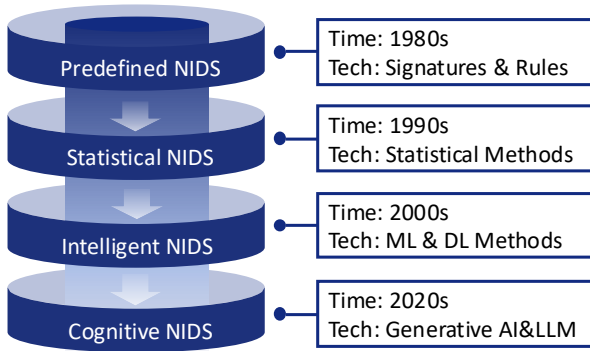


Fig. 1. The Roadmap of NIDS.

open challenges and future directions, and the conclusion is drawn in Section VI.

II. BACKGROUND OF NIDS AND LLMs

A. NIDS

NIDS are essential to modern cybersecurity, continuously evolving to counter increasingly sophisticated attacks. The roadmap in Fig. 1 outlines this progression, highlighting key technological advancements.

The earliest *Predefined NIDS* emerged in the late 1980s with systems like Haystack⁴ and SNORT⁵, which relied on static rule-based detection. While effective against known threats, they required frequent updates to address new attack patterns. To overcome this limitation, *Statistical NIDS* appeared in the mid-1990s, introducing anomaly-based detection. Systems such as SPADE⁶ established statistical baselines of normal network behavior, flagging deviations as potential intrusions. However, high false-positive rates often overwhelm analysts, reducing practicality in dynamic environments.

By the early 2000s, *Intelligent NIDS* leveraged Machine Learning (ML) and Deep Learning (DL) to adapt to evolving threats. Modern tools like Splunk⁷ and Suricata⁸ moved beyond rule-based detection, using clustering, supervised learning, and ensemble methods to improve accuracy. However, increasing model complexity led to concerns about interpretability. The advent of LLMs has driven the evolution toward *Cognitive NIDS*, capable of processing both structured and unstructured data, including logs, threat intelligence reports, and emails. By leveraging LLMs' reasoning capabilities, Cognitive NIDS enhance contextual awareness, provide actionable security insights, and enable more effective incident response.

B. LLMs

LLMs represent a major advancement of AI, particularly in NLP domain. These models are trained on vast amounts of data to understand and generate human-like text, excelling in tasks

such as summarization, translation, sentiment analysis. Their deep learning architecture, typically based on transformers, enables them to scale to billions of parameters, supporting their impressive performance across various applications.

LLMs have demonstrated broad applicability across industries, highlighting their versatility. In healthcare⁹, they assist in analyzing patient records, generating medical reports, aiding preliminary diagnoses, and supporting interactive virtual assistants. In finance¹⁰, LLMs help analyze complex financial data, generate insights, and improve decision-making processes. In customer service¹¹, they provide 24/7 support, handle inquiries, and enhance user experiences through personalized interactions. These applications illustrate the transformative potential of LLMs across multiple domains, reinforcing their role in driving advancements in cybersecurity and beyond.

III. AI-DRIVEN NIDS: FROM INTELLIGENT TO COGNITIVE

This section explores the evolution of NIDS from Intelligent to Cognitive, both of which fall under the broader category of AI-driven NIDS. The advancement of AI technologies has significantly enhanced the capabilities of NIDS, enabling them to address increasingly sophisticated security challenges.

A. General Pipeline of AI-driven NIDS

The effectiveness of AI-driven NIDS lies in their ability to analyze network data, distinguish between normal and anomalous behavior, and respond dynamically to threats. To understand the transition from Intelligent to Cognitive NIDS, we first examine the general pipeline of AI-driven NIDS. Fig. 2 illustrates five interconnected stages: data collection, data processing, intrusion detection, event analysis, and incident response, forming a backbone of effective AI-driven NIDS.

1) *Data Collection*: This stage aggregates raw data from packet captures, system logs, application logs, user activity records, and threat intelligence feeds using tools like Wireshark¹² and tcpdump¹³. Comprehensive data collection ensures full network visibility, enabling accurate analysis.

2) *Data Processing*: Raw data is cleansed, normalized, and structured to remove noise and redundancy. This step filters irrelevant traffic, extracts protocol-specific details (e.g., HTTP headers, DNS queries), and enhances data quality, reducing false positives and missed detections.

3) *Intrusion Detection*: AI-driven NIDS detect both known and emerging threats by continuously monitoring network activities. Intelligent NIDS rely on ML/DL models, while Cognitive NIDS leverage Generative AI and LLMs for deeper contextual understanding, improving detection of complex, multi-stage attacks.

⁴Smaha, Stephen E. "Haystack: An Intrusion Detection System." Fourth Aerospace Computer Security Applications Conference. Vol. 44. 1988.

⁵<https://www.snort.org/>

⁶<https://github.com/infosecdr/spade>

⁷<https://www.splunk.com/>

⁸<https://suricata.io/>

⁹Agent-Hospital: <https://www.taixex.cn/agent-hospital>

¹⁰FinGPT: <https://github.com/AI4Finance-Foundation/FinGPT>

¹¹ChatGPT-On-CS: <https://github.com/cs-lazy-tools/ChatGPT-On-CS>

¹²Wireshark: <https://www.wireshark.org/>

¹³tcpdump: <https://github.com/the-tcpdump-group/tcpdump>

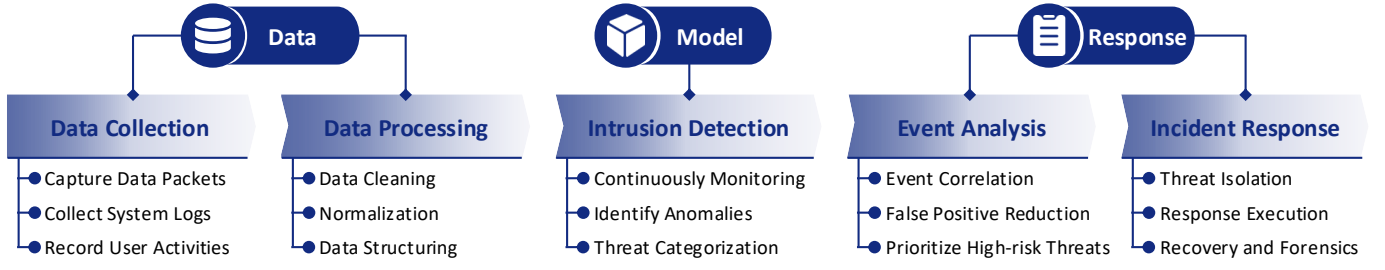


Fig. 2. The Pipeline of AI-driven NIDS.

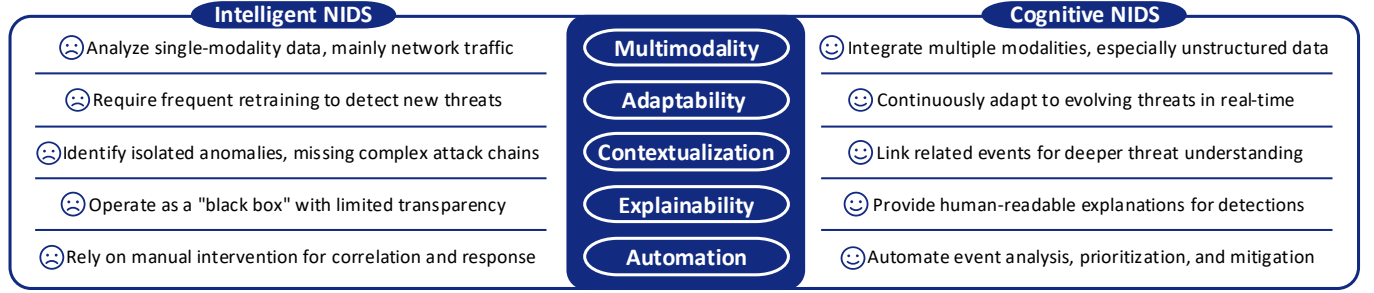


Fig. 3. Capability Comparison of Intelligent and Cognitive NIDS.

4) *Event Analysis*: Upon detecting an intrusion, NIDS correlate system logs, user activities, and threat intelligence to assess threat severity and minimize false positives. Integrating external sources like MITRE ATT&CK¹⁴ and Common Vulnerabilities and Exposures (CVE)¹⁵ helps identify zero-day exploits and Advanced Persistent Threat (APTs), improving situational awareness.

5) *Incident Response*: This stage mitigates threats through automated or manual actions such as isolating compromised systems, blocking malicious traffic, and generating alerts. It also includes forensic investigations and system recovery, ensuring resilience against future attacks.

B. The Gap Between Intelligent and Cognitive NIDS

While Intelligent NIDS utilize ML and DL methods for threat detection, they struggle with unstructured data processing, complex attack detection, explainability, and workflow optimization. Cognitive NIDS, powered by LLMs, overcome these limitations by enhancing contextual understanding, automation, and adaptive decision-making. A comparative overview of their capabilities is summarized in Fig. 3.

1) *Unstructured Data Handling*: Intelligent NIDS focus on structured data like IP headers and packet payloads, overlooking insights from unstructured sources such as logs, emails, and threat reports. Cognitive NIDS leverage multimodality, integrating both structured and unstructured data to enhance threat detection. With advanced NLP capabilities, they analyze logs for anomalies, extract Indicators of Compromise (IoCs), and detect phishing or social engineering threats across diverse data formats. This multimodal approach provides a more comprehensive and context-aware view of network security,

enabling richer threat intelligence and more accurate intrusion detection.

2) *Detection of Evolving and Multi-Stage Attacks*: Intelligent NIDS rely on pattern-based detection, often struggling to adapt to evolving threats and multi-stage attacks like APTs, where individual actions appear benign but collectively form a coordinated attack. Cognitive NIDS enhance adaptability by leveraging contextual reasoning and continuously learning from new attack patterns. They dynamically correlate seemingly unrelated events, such as suspicious logins, anomalous file downloads, and data exfiltration, enabling early detection of complex attack chains before significant damage occurs.

3) *Explainability and Reporting Mechanisms*: Intelligent NIDS, particularly those based on deep learning, often function as "black boxes", providing little transparency in decision-making. Cognitive NIDS improve explainability by generating natural language descriptions for detected threats, detailing not only what was flagged but also why. They produce human-readable incident reports tailored to different stakeholders, from technical analysts to business leaders, improving decision-making, trust, and response efficiency.

4) *Automation and Workflow Optimization*: Intelligent NIDS require manual intervention for event correlation, threat prioritization, and response execution, increasing analyst workload and response time. Cognitive NIDS leverage real-time decision-making and contextual automation, dynamically prioritizing threats and even executing mitigation actions, such as isolating compromised endpoints, blocking malicious IPs, and generating automated incident reports. By reducing manual effort, Cognitive NIDS enhance efficiency and accuracy in intrusion detection.

¹⁴MITRE ATT&CK: <https://attack.mitre.org/>

¹⁵CVE: <https://cve.mitre.org/>

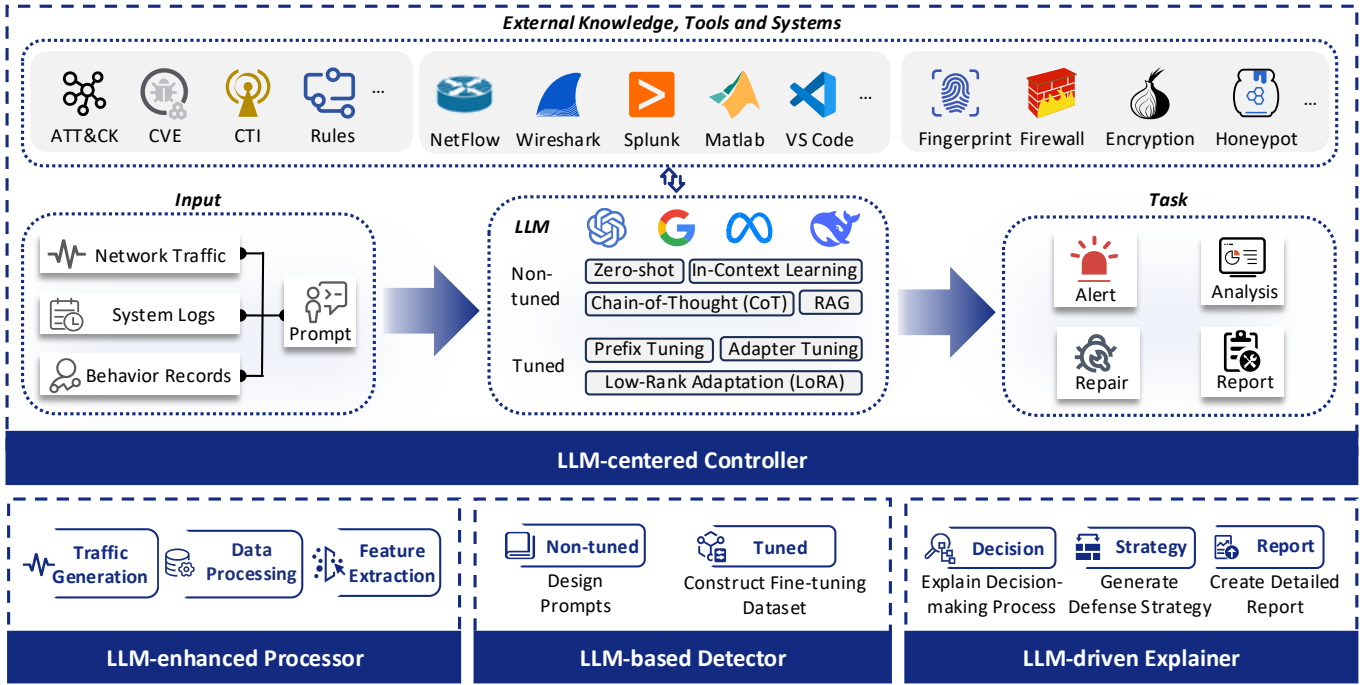


Fig. 4. The Overview of Cognitive NIDS.

C. Enabling Techniques for Cognitive NIDS

Integrating LLMs to transition from Intelligent to Cognitive NIDS presents unique challenges. Although LLMs have extensive knowledge of common tasks, they may lack domain-specific expertise, and their output can be inherently random. To overcome these challenges, various enabling techniques, categorized into non-tuned and tuned methods, can facilitate the effective integration of LLMs into NIDS.

1) Non-tuned Methods:

a) *Zero-shot Prompt*: Traditional NLP techniques require large labeled datasets for supervised training. In contrast, Zero-shot prompting uses pre-trained LLMs directly for tasks without needing task-specific training. However, this method faces challenges, such as the pre-trained model's insufficient knowledge of network intrusion detection and biases inherent in the training data.

b) *In-Context Learning (ICL)*: ICL enhances LLMs by using formatted prompts with task descriptions and examples, allowing the model to learn domain knowledge dynamically without modifying its parameters. This method improves the LLM's adaptability and task performance by effectively guiding the model in understanding specific requirements and applying domain-specific knowledge.

c) *Chain-of-Thought (CoT)*: CoT prompts are designed to guide LLMs through complex reasoning tasks by incorporating intermediate reasoning steps. Unlike ICL, which focuses on input-output pairs, CoT prompts help the LLM connect logical steps, enhancing accuracy and stability in solving problems that require abstract or symbolic reasoning.

d) *Retrieval-Augmented Generation (RAG)*: RAG combines information retrieval with LLMs' generative capabilities, allowing models to retrieve and integrate relevant information

from external sources to generate more accurate, context-aware responses. This approach bridges the knowledge gap between general models and up-to-date information, enabling more effective task performance.

2) *Tuned Methods*: Tuned methods are categorized into Full Fine-Tuning (FFT) and Parameter-Efficient Fine-Tuning (PEFT). FFT updates all model parameters, making it resource-intensive and computationally demanding. In contrast, PEFT modifies only a subset of parameters, providing a more efficient and scalable alternative, particularly suited for network intrusion detection. Below are key PEFT techniques.

a) *Prefix Tuning*: Prefix tuning introduces task-specific prefix vectors to the input sequence, optimizing the model for specific tasks without altering the original parameters. This approach is efficient, reducing computational overhead and enhancing flexibility and interpretability by adapting the model with minimal changes.

b) *Adapter Tuning*: Adapter tuning integrates lightweight modules into pre-trained models, enabling task adaptation without altering the model's core architecture. These modular adapters provide flexibility for experimentation and efficient task optimization.

c) *Low-Rank Adaptation (LoRA)*: LoRA introduces small, low-rank matrices into critical layers of a model, reducing the number of trainable parameters required for fine-tuning. This method maintains model performance while minimizing storage and computational costs, making it highly suitable for large models, including transformers.

The above techniques could equip LLMs with domain-specific capabilities for NIDS, enabling their seamless integration and advancing the development of Cognitive NIDS.

TABLE I
IMPLEMENTATION SUMMARY OF NIDS WITH LLMs

Ref.	LLM-enhanced Processor			LLM-based Detector		LLM-driven Explainer			LLM-centered Controller	LLMs	Dataset
	Traffic Generation	Data Processing	Feature Extraction	Non-tuned	Tuned	Decision	Strategy	Report			
[1]	●	○	○	○	○	○	○	○	○	GPT-2	ISXW2016 etc.
[2]	●	○	○	○	○	○	○	○	○	NetGPT	ISCTXor2016 etc.
[3]	●	○	○	○	○	○	○	○	○	GPT-3	ToN IoT Dataset
[4]	○	●	○	○	○	○	○	○	○	GPT-4	LogHub
[5]	○	●	○	○	○	○	○	○	○	GPT-3.5, GPT-4	N/A
[6]	○	○	●	○	○	○	○	○	○	GPT-4	Private
[7]	○	○	●	●	○	○	○	○	○	GPT-4, LLaMA2 etc.	DTL-IDS 5G
[8]	○	○	○	●	●	●	●	○	○	GPT-4, Mistral etc.	Private
[9]	○	○	○	●	●	●	○	○	○	GPT-4, LLaMA3	UNSW-NB15 etc.
[10]	○	○	○	○	●	○	○	○	○	Zephyr-7b- β	Private
[11]	○	○	○	○	○	●	○	○	○	GPT-4	NF-BoT
[12]	○	○	○	○	○	○	○	○	○	Mistral	CICIDS2017
[13]	○	○	○	○	○	○	●	○	○	GPT-3.5	N/A
[14]	○	○	○	○	○	●	○	○	○	GPT-3.5	KDD99 dataset
[15]	○	●	○	●	○	●	○	○	●	GPT-3.5, GPT-4o etc.	ACI-IoT23 etc.

IV. THE IMPLEMENTATION OF NIDS WITH LLMs

This section examines the role of LLMs in advancing Cognitive NIDS. With their extensive internal knowledge and advanced reasoning capabilities, LLMs can enhance or integrate with existing NIDS to improve efficiency and effectiveness at various stages of the pipeline discussed in § III-A. Fig. 4 provides an overview of Cognitive NIDS, which comprises four promising components: *LLM-enhanced Processor* (§ IV-A), *LLM-based Detector* (§ IV-B), *LLM-driven Explainer* (§ IV-C), and *LLM-centered Controller* (§ IV-D). The following sections provide a comprehensive analysis of these four components, with representative implementations summarized in Table I.

- LLM-enhanced Processor improves data processing efficiency and extracts informative features to enhance threat identification.
- LLM-based Detector leverages LLMs to detect complex attack patterns, significantly enhancing intrusion detection performance.
- LLM-driven Explainer enhances explainability by interpreting decision-making processes, generating defense strategies, and producing detailed reports.
- LLM-centered Controller optimizes workflow orchestration, facilitating seamless collaboration between various system components.

By leveraging LLMs, NIDS benefit from enhanced data processing, sophisticated detection capabilities, improved explainability, and optimized workflows. Cognitive NIDS provide stronger, more adaptive protection against the ever-evolving cybersecurity threat landscape.

A. LLM-enhanced Processor

The LLM-enhanced Processor leverages LLMs to improve the data processing capabilities of NIDS. In intrusion detection, the monitored data primarily consists of network traffic, user behavior logs, system event logs, application messages, and alerts from security devices. This data serves as the foundation for identifying potential security threats.

1) *Traffic Generation*: Network traffic data is fundamental for constructing effective NIDS, with model performance largely dependent on the quality and diversity of the training data. The scarcity of high-quality datasets in network intrusion detection presents significant challenges, especially in real-world applications. Existing datasets are often outdated or insufficient in both quality and quantity, as organizations are reluctant to share data due to privacy concerns or the protection of commercial secrets. As a result, traffic generation has emerged as a valuable approach, utilizing synthetic techniques to simulate network traffic for performance evaluation and benchmarking of NIDS. For instance, Meng *et al.* [1] developed NetGPT, which integrates multi-modal traffic modeling by encoding heterogeneous network traffic headers and payloads into a unified text input for traffic understanding and generation tasks. Qu *et al.* [2] introduced TrafficGPT, which directly generated PCAP files from token lists, achieving high quality and authenticity as demonstrated by metrics like JS divergence. Additionally, Kholgh *et al.* [3] created PAC-GPT, which fine-tuned GPT-3 to generate ICMP and DNS packets. These developments showcase the feasibility of using generative models such as LLMs to generate high-quality traffic datasets.

2) *Data Processing*: Transforming raw network data into a usable format often involves complex tasks such as data cleaning, normalization, and integration, particularly when dealing with heterogeneous data from various sources. LLMs, with their advanced semantic understanding and reasoning capabilities, can automate this process by identifying and correcting errors, missing values, and redundant information, thereby enhancing data quality and consistency. For example, Zhang *et al.* [4] proposed a log parsing framework based on information entropy sampling and CoT prompting, eliminating the need of manual rules and achieving high efficiency in log parsing for downstream tasks. Furthermore, LLMs can identify and correlate data from diverse sources, enabling a more comprehensive understanding of network attacks. Daniel *et al.* [5] demonstrated the use of ChatGPT for labeling NIDS rules with MITRE ATT&CK techniques, showcasing LLMs'

ability to process and integrate multi-source data effectively. In summary, the application of LLMs in data processing significantly reduces the need for manual intervention, enhances data accuracy, and lays a strong foundation for subsequent detection and response strategies.

3) *Feature Extraction*: LLMs excel at extracting meaningful features from data, which is essential for identifying potential security risks. For instance, Wang *et al.* [6] used ChatGPT to obtain embeddings from IoT traffic payloads, which were then used to train a deep learning model for IoT traffic anomaly detection. Their results highlighted the importance of the embedding layer in capturing contextual and lexical nuances, which significantly improved detection accuracy and reduced false alarms. Zhang *et al.* [7] applied LLMs for feature selection by ranking the importance of selected features at three levels: very important, kind of important, and not very important. This approach streamlined feature length and minimized the impact of irrelevant features on detection performance. These studies demonstrate that LLMs can significantly enhance feature extraction, improving the detection of anomalous behaviors and boosting the performance of intrusion detection models.

B. LLM-based Detector

While several studies have explored the integration of native LLMs with cybersecurity, only a few have proposed specific implementations for intrusion detection. Building on the techniques discussed in Section III-C, we introduce these advancements from both tuned and non-tuned perspectives.

1) *Non-tuned LLM-based Detector*: [9] have applied Zero-shot prompting to directly use LLMs for detection tasks. However, a key challenge is that pre-trained LLMs may lack sufficient knowledge of network attacks to perform complex detection tasks effectively. To address this limitation, Zhang *et al.* [7] employed a simple ICL approach by providing labeled examples to the model using GPT-4. Their method achieved over 90% accuracy on a simple dataset containing only five attack types, demonstrating that ICL can effectively improve the detection performance of LLMs. In another work, Bui *et al.* [8] enhanced frozen LLMs by supplementing them with useful information from ChromaDB¹⁶ and Langchain framework¹⁷, such as malicious payloads and threat intelligence related to suspected attacker IPs, achieving significant improvements over few-shot prompting. These studies suggest that leveraging the inherent capabilities of LLMs, along with supplemental information, can significantly enhance their performance for intrusion detection without updating the pre-trained model parameters.

2) *Tuned LLM-based Detector*: While pre-trained LLMs excel in a wide range of language-understanding tasks, they may not perform optimally on specific tasks such as intrusion detection. To address this, fine-tuning—further training on task-specific data—can be applied, enabling the model to understand and perform detection tasks more effectively. Compared to relying solely on prompting, tuned methods

offer a more direct and efficient optimization by adjusting model parameters specifically for the task at hand. Fine-tuned LLMs integrate more seamlessly with existing NIDS, reducing dependence on complex prompts and improving performance. For example, Houssel *et al.* [9] fine-tuned the LLaMA3 model using the NetFlow dataset, applying Odds Ratio Preference Optimization (ORPO) and Kahneman-Tversky Optimization (KTO) techniques. They proposed that LLMs should serve as complementary solutions to state-of-the-art NIDS, enhancing their capabilities. Fine-tuning also allows smaller models to achieve performance on par with larger ones, reducing computational costs and latency during inference. Rigaki *et al.* [10] used LoRA to fine-tune a 7-billion-parameter pre-trained LLM (Zephyr-7b- β), achieving performance comparable to that of more powerful models such as GPT-4. These studies show that fine-tuning not only improves model performance for network intrusion detection but also optimizes resource utilization, maintaining or even enhancing performance with lower computational overhead.

C. LLM-driven Explainer

The LLM-driven Explainer leverages LLMs to enhance decision-making and strategy formulation in cybersecurity. By analyzing current security incidents and historical data, it provides decision suggestions, response strategies, and automated reports that summarize detected events, actions, and justifications. This improves the explainability and transparency of NIDS. In this section, we explore how LLMs contribute to NIDS explainability across three key dimensions: decision, strategy, and report.

1) *Decision-level Explainability*: At the decision level, LLMs correlate multi-source data and explain why a specific threat is flagged by detailing the key features, patterns, or anomalies that influence the detection decision. By incorporating both traditional network metrics and semantic information from network communications, LLMs offer a more comprehensive analysis of network behavior. Through in-context learning, LLMs could enhance decision explainability by understanding event relationships and contextual information. Ziems *et al.* [11] demonstrated that LLM-generated explanations for decision tree-based NIDS models align closely with human assessments in terms of readability, quality, and contextual knowledge, aiding in a clearer understanding of decision boundaries. Additionally, Khediri *et al.* [12] integrated SHAP explanations with LLM-generated descriptions, providing feature importance and human-readable justifications for NIDS model predictions.

2) *Strategy-level Explainability*: At the strategy level, LLMs clarify how the system approaches various types of threats and multi-stage attacks. The LLM-driven Explainer generates response strategies such as isolating suspicious users, restricting access, or monitoring traffic. These LLM-generated recommendations enable security teams to take appropriate actions, improving strategy accuracy and response efficiency. Bui *et al.* [8] introduced a method to measure the contribution of each token to a predicted output label, allowing for concise incident summaries. This contextual information

¹⁶ChromaDB: <https://www.trychroma.com/>

¹⁷Langchain: <https://www.langchain.com/>

enables human analysts to interact with the model for further Q&A and assistance. Jüttner *et al.* [13] developed ChatIDS, an early-stage security alert explanation system powered by ChatGPT, demonstrating the potential of LLMs in providing intuitive interpretations of security alerts and actionable security measures.

3) *Report-level Explainability*: In cybersecurity incident response, LLMs can generate detailed analytical reports and visualizations, supporting security teams with comprehensive event descriptions, detection processes, impact assessments, and recommended mitigation measures. Automating report generation significantly reduces the burden of manual documentation while improving situational awareness. By analyzing historical data, LLMs can identify trends and patterns in security events, assisting teams in retrospective analyses and predicting future attack behaviors. This enables organizations to proactively prepare for emerging threats. Ali *et al.* [14] developed HuntGPT, a system that provides a dashboard summarizing attack types, key features influencing model decisions, and their impacts. Additionally, it allows users to interactively obtain real-time response recommendations and detailed investigative reports.

D. LLM-centered Controller

The LLM-centered Controller serves as the orchestrator of the intrusion detection workflow, ensuring seamless collaboration among various tools and components to maximize system efficiency and effectiveness. As illustrated in Fig. 4, the Controller leverages LLMs to manage processes across all stages, from data processing to incident response. It facilitates communication between modules, dynamically adjusts system configurations, and ensures smooth task execution. For instance, by integrating with traffic collection tools like Wireshark, the LLM can capture and analyze network traffic in real-time, extract key features, and automate traffic processing, optimizing data utilization and analysis.

Beyond data handling, the Controller enhances incident response by coordinating actions across multiple security tools, such as fingerprint recognition, encryption systems, and firewalls, providing a unified view with visual analytics that enhances interpretability. It can automate the isolation of compromised systems, enforce firewall rules for mitigation, or trigger in-depth forensic investigations on flagged endpoints. This integrated approach improves security posture assessment, reduces manual effort, shortens response times, and strengthens the organization's overall cybersecurity defenses. Additionally, it dynamically updates threat intelligence by retrieving information from internal and external knowledge bases such as MITRE ATT&CK, CVE, and Cyber Threat Intelligence (CTI¹⁸). This continuous adaptation enables the system to respond effectively to evolving security challenges. By serving as an intelligent unifying orchestrator, the LLM-centered Controller ensures that all NIDS components work harmoniously, adapting to dynamic threats and delivering real-time, robust protection against sophisticated attacks.

Early research has demonstrated the feasibility of this architecture. Li *et al.* [15] proposed IDS-AGENT, an LLM-based intrusion detection agent that employs an iterative reasoning-action pipeline. IDS-AGENT extracts traffic data, preprocesses it, calls various machine learning models for classification, retrieves knowledge from internal and external sources, and generates final detection reasoning. While IDS-AGENT achieves competitive performance compared to traditional NIDS, its focus remains on managing detection models rather than orchestrating the entire intrusion detection lifecycle or integrating external knowledge, tools, and systems comprehensively.

V. CHALLENGES AND FUTURE DIRECTIONS

This section discusses the challenges of integrating LLMs into NIDS and outlines key future directions to maximize their potential in intrusion detection.

A. Challenges

1) *Validity*: A primary challenge in utilizing LLMs for NIDS is the reliability and consistency of their outputs. Due to the inherent unpredictability of LLMs, threat assessments may be compromised by hallucinations¹⁹, where models generate inaccurate, misleading, or contextually irrelevant outputs. Additionally, biased training datasets can skew detection results, leading to false positives or blind spots in identifying specific attack patterns. Another critical issue is the security of LLMs themselves—adversarial attacks could manipulate or deceive the model, allowing malicious actors to evade detection.

2) *Complexity*: Deploying LLMs in real-time network environments poses significant computational and infrastructural challenges. Large-scale models require substantial processing power, memory, and storage, limiting their feasibility in resource-constrained environments such as edge networks and IoT ecosystems. Moreover, high latency in processing vast amounts of network traffic would delay detection and response time, potentially leaving systems vulnerable to fast-moving cyber threats.

3) *Privacy*: Integrating LLMs into NIDS raises privacy concerns, particularly regarding the handling of sensitive data within network traffic and logs. Since LLMs process extensive datasets, they may inadvertently expose private or confidential information, violating data protection regulations such as GDPR²⁰. Additionally, security risks associated with data leakage during training and inference must be addressed to prevent unauthorized access or misuse of sensitive information.

B. Future Directions

1) *Multimodal Integration for Robust Threat Detection*: One promising avenue for future research is the integration of multimodal data to enhance the performance and reliability of NIDS with LLMs. By combining multiple data sources,

¹⁹Huang, Lei, et al. "A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions." *ACM Transactions on Information Systems* 43.2 (2025): 1-55.

²⁰GDPR: <https://gdpr-info.eu/>

¹⁸CTI: <https://github.com/OpenCTI-Platform/opencti>

such as network traffic, system logs, behavioral analytics, and external threat intelligence, LLMs can provide a more contextualized and accurate understanding of threats. This integration enhances detection accuracy, reduces false positives, and ensures models are not overly dependent on single-modality inputs, mitigating biases and improving generalization.

2) *Optimizing Real-Time Intrusion Detection for Edge Networks*: For LLMs to be effectively applied in edge networks, advancements in real-time analysis and detection are crucial. Future research should focus on developing lightweight models (e.g. Small Language Model²¹) that require fewer computational resources and while maintaining high performance. Techniques such as model quantization, knowledge distillation, and efficient architecture can help reduce latency in threat detection, making real-time systems more feasible. Deploying Small Language Models at the edge (e.g. IoT gateways, routers, and mobile devices) will enable low-latency detection, improving response times and making Cognitive NIDS feasible in resource-constrained environments.

3) *Secure and Privacy-Preserving Threat Intelligence Collaboration*: To balance privacy and detection effectiveness, the NIDS with LLMs should integrate privacy-preserving techniques such as federated learning, homomorphic encryption, and secure multi-party computation. These approaches allow multiple organizations to collaborate on improving intrusion detection models without sharing sensitive data. Additionally, Zero-Trust Architectures (ZTA²²) and confidential computing can ensure that LLM inference is performed securely, reducing risks associated with data exposure during processing. By integrating LLMs with existing security tools, such as firewalls, SIEMs, and endpoint detection systems, Cognitive NIDS can provide more effective threat intelligence sharing and adaptive security measures while maintaining privacy compliance.

4) *Multi-Agent Systems for Collaborative Intrusion Detection*: A promising direction for Cognitive NIDS is the development of multi-agent systems (MAS²³), where autonomous LLM-agents collaborate to enhance threat detection, response, and decision-making. Unlike centralized models, MAS distribute security tasks among specialized agents, each analyzing different data sources such as network traffic, logs, and behavioral patterns. These agents share real-time intelligence, improving adaptive threat detection and cross-context correlation. By operating in parallel, MAS reduces computational overhead, enhances scalability, and enables faster incident response. Additionally, coordinated agents can autonomously mitigate threats, such as isolating compromised nodes or adjusting firewall rules. Future research should focus on optimizing MAS architectures to achieve greater accuracy, automation, and resilience against evolving cyber threats.

VI. CONCLUSION

This paper has examined the transformative role of Large Language Models (LLMs) in advancing Network Intrusion

Detection Systems (NIDS), highlighting their potential to transition the field from intelligent to cognitive systems. By exploring practical implementations, we showcase LLMs as processors, detectors, explainers, and controllers, with the LLM-centered Controller emphasizing the orchestration capabilities to streamline workflows and enhance system efficiency. While challenges such as validity, complexity, and privacy remain, proposed future research directions—including multimodal integration, real-time analysis, privacy-preserving collaboration, and multi-agent systems—provide a pathway to address these issues. We hope this work inspires further efforts to harness LLMs for reliable, adaptive, and explainable NIDS.

REFERENCES

- [1] X. Meng, C. Lin, Y. Wang, and Y. Zhang, "Netgpt: Generative pretrained transformer for network traffic," *arXiv preprint arXiv:2304.09513*, 2023.
- [2] J. Qu, X. Ma, and J. Li, "Trafficgpt: Breaking the token barrier for efficient long traffic analysis and generation," *arXiv preprint arXiv:2403.05822*, 2024.
- [3] D. K. Kholgh and P. Kostakos, "Pac-gpt: A novel approach to generating synthetic network traffic with gpt-3," *IEEE Access*, 2023.
- [4] W. Zhang, H. Guo, A. Le, J. Yang, J. Liu, Z. Li, T. Zheng, S. Xu, R. Zang, L. Zheng *et al.*, "Lemur: Log parsing with entropy sampling and chain-of-thought merging," *arXiv preprint arXiv:2402.18205*, 2024.
- [5] N. Daniel, F. K. Kaiser, A. Dzega, A. Elyashar, and R. Puzis, "Labeling nids rules with mitre att & c&k techniques using chatgpt," in *European Symposium on Research in Computer Security*. Springer, 2023, pp. 76–91.
- [6] T. Wang, Z. Zhao, and K. Wu, "Exploiting llm embeddings for content-based iot anomaly detection," in *2024 IEEE Pacific Rim Conference on Communications, Computers and Signal Processing (PACRIM)*. IEEE, 2024, pp. 1–6.
- [7] H. Zhang, A. B. Sediq, A. Afana, and M. Erol-Kantarci, "Large language models in wireless application design: In-context learning-enhanced automatic network intrusion detection," *arXiv preprint arXiv:2405.11002*, 2024.
- [8] M.-T. Bui, M. Boffa, R. V. Valentim, J. M. Navarro, F. Chen, X. Bao, Z. B. Houidi, and D. Rossi, "A systematic comparison of large language models performance for intrusion detection," *Proceedings of the ACM on Networking*, vol. 2, no. CoNEXT4, pp. 1–23, 2024.
- [9] P. R. Houssel, P. Singh, S. Layeghy, and M. Portmann, "Towards explainable network intrusion detection using large language models," *arXiv preprint arXiv:2408.04342*, 2024.
- [10] M. Rigaki, C. Catania, and S. Garcia, "Hackphyr: A local fine-tuned llm agent for network security environments," *arXiv preprint arXiv:2409.11276*, 2024.
- [11] N. Ziems, G. Liu, J. Flanagan, and M. Jiang, "Explaining tree model decisions in natural language for network intrusion detection," *arXiv preprint arXiv:2310.19658*, 2023.
- [12] A. Khediri, H. Slimi, A. Yahiaoui, M. Dourdour, H. Bendjenna, and C. E. Ghenai, "Enhancing machine learning model interpretability in intrusion detection systems through shap explanations and llm-generated descriptions," in *2024 6th International Conference on Pattern Analysis and Intelligent Systems (PAIS)*. IEEE, 2024, pp. 1–6.
- [13] V. Jüttner, M. Grimmer, and E. Buchmann, "Chatids: Explainable cybersecurity using generative ai," *arXiv preprint arXiv:2306.14504*, 2023.
- [14] T. Ali and P. Kostakos, "Huntgpt: Integrating machine learning-based anomaly detection and explainable ai with large language models (llms)," *arXiv preprint arXiv:2309.16021*, 2023.
- [15] Y. Li, Z. Xiang, N. D. Bastian, D. Song, and B. Li, "Ids-agent: An llm agent for explainable intrusion detection in iot networks," in *NeurIPS 2024 Workshop on Open-World Agents*.

²¹Zhang, Peiyuan, et al. "TinyLLaMA: An Open-Source Small Language Model." *arXiv preprint arXiv:2401.02385* (2024).

²²ZTA: <https://www.nist.gov/publications/zero-trust-architecture>

²³Talebirad, Yashar, and Amirhossein Nadiri. "Multi-Agent Collaboration: Harnessing the Power of Intelligent LLM Agents." *arXiv preprint arXiv:2306.03314* (2023).