IDS-AGENT: AN LLM AGENT FOR EXPLAINABLE IN TRUSION DETECTION IN IOT NETWORKS

Anonymous authors

Paper under double-blind review

ABSTRACT

Emerging threats to IoT networks have accelerated the development of intrusion detection systems (IDSs), characterized by a shift from traditional approaches based on attack signatures or anomaly detection to approaches based on machine learning (ML). However, current ML-based IDSs often lack result explanations and struggle to address zero-day attacks due to their fixed output label space. In this paper, we propose IDS-Agent, the first IDS based on an AI agent powered by large language models (LLMs). For each input network traffic and a detection request from the user, IDS-Agent predicts whether the traffic is benign or being attacked, with an explanation of the prediction results. The workflow of IDS-Agent involves iterative reasoning by a core LLM over the observation and action generation informed by the reasoning and retrieved knowledge. The action space of IDS-Agent includes data extraction and preprocessing, classification, knowledge, and memory retrieval, and results aggregation – these actions will be executed using abundant tools, mostly specialized for IDS. Furthermore, the IDS-Agent is equipped with a memory and knowledge base that retains information from current and previous sessions, along with IDS-related documents, enhancing its reasoning and action generation capabilities. The system prompts of IDS-Agent can be easily customized to adjust detection sensitivity or identify previously unknown types of attacks. In our experiments, we demonstrate the strong detection capabilities of IDS-Agent compared with ML-based IDSs and an IDS based on LLM with prompt engineering. IDS-Agent outperforms these SOTA baselines on the ACI-IoT and CIC-IoT benchmarks, with 0.97 and 0.75 detection F1 scores, respectively. IDS-Agent also achieves a recall of 0.61 in detecting zero-day attacks, outperforming previous approaches specially designed for this task.

034

037

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

029

031

032

1 INTRODUCTION

In recent years, the Internet of Things (IoT) has emerged as a transformative technology, increasingly adopted across a wide range of applications (Chataut et al., 2023). Alongside its rapid development, security concerns have arisen within IoT networks due to the typically large number of devices with potential trustworthiness issues (Alghofaili & Rassam, 2022). Therefore, the deployment of Intrusion Detection Systems (IDSs) has become essential, as they play a critical role in monitoring network traffic and identifying malicious activities (Khraisat et al., 2019).

044 Many traditional IDS employ signature-based methods, which rely on signatures of known attacks stored in databases (Ioulianou et al., 2018; Nagaraju et al., 2021; AlYousef & Abdelmajeed, 2019). 046 They suffer from high false negative rates when variations in attack methodologies do not exactly 047 match the existing signatures. Alternatively, IDS can also be designed as an anomaly detector to 048 identify distributional deviations from normal traffic (Khraisat et al., 2018; Jia et al., 2021). However, defining the normal behavior in a network can be challenging, especially in diverse and dynamic environments where normal activity can change over time. Machine learning (ML)-based IDS was 051 then proposed to address these problems by leveraging the representation power of ML models, such as deep neural networks (DNNs), to capture complex attack patterns from extensive training 052 data (Maseer et al., 2021). However, ML-based IDS still face limitations due to constraints in the model architecture and training data, which hinder their detection capabilities, especially when

addressing zero-day attacks (Ullah et al., 2019). Furthermore, the detection results from ML-based
IDS often lack clear explanations, which can diminish their credibility, particularly in safety-critical
IoT scenarios where understanding the rationale behind alerts is crucial (Catarinucci et al., 2015;
Ahmad & Zhang, 2021).

058 Recently, AI agents empowered by large language models (LLMs) have been created to handle 059 complicated tasks in various application domains (Yu et al., 2023; Abbasian et al., 2024; Shi et al., 060 2024; Li et al., 2024; Deng et al., 2023; Gur et al., 2024; Zheng et al., 2024). These agents are 061 characterized by their integrated capabilities to observe the environment, generate a reasonable 062 plan, and then take appropriate actions according to the user's requests. Typically, LLM agents are 063 equipped with a knowledge/memory base and a toolbox that includes tools such as local functions 064 and third-party APIs. They utilize one or more LLMs for knowledge/memory-based reasoning and planning, including the selection of the most effective tools to employ. The intelligence of LLM 065 agents in reasoning (as well as analysis and criticism (Shinn et al., 2023)) makes them both powerful 066 operators and effective intermediaries between task execution and human users. 067

068 In this paper, we propose IDS-Agent, the first LLM agent designed for intrusion detection, featuring 069 capabilities for explanation, customization, and adaptation to zero-day attacks. IDS-Agent takes 070 an input request for intrusion detection with a target raw data flow, and outputs detection results 071 with a detailed explanation. The agent adopts a reasoning-followed-by-action pipeline ((Yao et al., 2023)) with a specialized action space for intrusion detection. Specifically, knowledge-enabled 072 reasoning based on long-term memory from previous sessions is performed by the core LLM of the 073 agent to decide the optimal tools (and their settings/parameters) for data extraction, preprocessing, 074 classification, and results aggregation. 075

Compared with existing ML-based IDSs, IDS-Agent achieves a stronger detection performance and better interpretability by harnessing the power of multiple ML models and external knowledge in a comprehensive way. It aggregates the classification results from multiple ML models by prompting an LLM with the top-k label predictions and their confidence scores for each model. The prompt also includes external knowledge (e.g. regarding particular attack types) obtained by calling a search engine and additional instructions, for example, to customize detection sensitivity or to reveal new attack types. The LLM is instructed to produce structured outputs, including detection results and an explanation. Our main technical contributions are summarized below:

083 084

085

087

088

089

091

092

094

- We propose IDS-Agent, the first LLM-powered agent for intrusion detection, featured by its capabilities of results explanation, detection customization, and revealing of zero-day attacks.
- We propose a reasoning-followed-by-action pipeline for IDS-Agent with an action space and toolbox specialized for network traffic processing and intrusion detection results aggregation.
- We demonstrate the effectiveness of IDS-Agent on two IDS benchmarks, ACI-IoT'23 and CIC-IoT'23. IDS-Agent achieves higher detection accuracy compared with the latest LLM-based method, various ML models, and the majority voting classifier.
- Experiments show that when classifiers produce discrepancy predictions, IDS-Agent can utilize inherent and external knowledge to help the decision-making. IDS-Agent also demonstrates clearly better performance than existing approaches in detecting zero-day attacks. Moreover, we find that IDS-Agent effectively follows the sensitivity instructions without requiring expert intervention or additional tuning.
- 095 096

098

2 RELATED WORK

099 **Conventional IDS** An Intrusion Detection System (IDS) is designed to detect malicious activities 100 on computer systems, helping to ensure system security (Khraisat et al., 2019). IDSs are generally 101 classified into two main types: Signature-based Intrusion Detection Systems (SIDS) and Anomaly-102 based Intrusion Detection Systems (AIDS). SIDS relies on pattern-matching techniques to identify 103 known attacks (Ioulianou et al., 2018; Nagaraju et al., 2021; AlYousef & Abdelmajeed, 2019). 104 However, the rise in zero-day attacks has increasingly diminished the effectiveness of SIDS, as these 105 new attacks lack existing signatures (Ullah et al., 2019). In contrast, AIDS constructs a model of normal system behavior using machine learning, statistical, or knowledge-based techniques. Any 106 significant deviation from the constructed model is flagged as an anomaly, potentially indicating 107 an intrusion (Khraisat et al., 2018; Jia et al., 2021). Knowledge-based AIDS requires creating a

knowledge base that reflects the legitimate traffic profile by using techniques such as description
language (Studnia et al., 2018) and expert system (Kim et al., 2014). However, in a dynamically
changing environment, a time-consuming regular update on the knowledge base is needed. Compared
with these conventional IDS, our IDS-Agent can adopt a search engine to obtain up-to-date
knowledge and can handle zero-day attacks, as shown by our experiments.

113 ML-based IDS Many machine learning models, such as MLP (Bajaj & Arora, 2013), KNN (Li et al., 114 2014), Decision Tree (Guezzaz et al., 2021), SVM (Mohammadi et al., 2021), have been explored for 115 anomaly-based intrusion detection. For IoT intrusion detection, Verma & Ranga (2020) conducted a 116 comprehensive comparison of ensemble and individual classifiers, including Random Forest (RF), 117 AdaBoost (AB), and Gradient Boosted Machines (GBM). Roy et al. (2022) proposed a lightweight IDS model utilizing machine learning to detect cyber-attacks and anomalies in resource-constrained 118 IoT systems. Davis et al. (2024) advanced this line of research by applying a quantum-annealing 119 approach for feature selection in IoT intrusion detection. Compared with ML-based IDSs, our 120 IDS-Agent not only achieves better empirical performance but also provides a detailed explanation 121 of each intrusion detection result. 122

123 LLM-based IDS Large language models (LLMs), especially generative pre-trained commercial transformers, like GPTs, have recently demonstrated outstanding ability in information comprehension and 124 125 reasoning tasks. This has motivated some studies in applying LLMs to abnormal detection tasks, such as compiler optimization (Gu, 2023) and software vulnerability detection (Guo et al., 2024). Zhang 126 et al. (2024) is the first to use LLMs for IDS by employing a straightforward in-context learning 127 approach with GPT-4, which provides it with a few labeled examples. Their method achieved over 128 90% accuracy on a simple dataset containing only five types of attacks. However, in this paper, we 129 demonstrate that the performance of their method drops significantly when tested on more complex 130 and diverse datasets. Different from this LLM-based IDS, our IDS-Agent uses LLM for reasoning, 131 action planning, and tool selection, leading to a huge performance gain on diverse datasets. 132

133

134 135

136

3 Method

3.1 Overview

137 IDS-Agent is designed to produce a prediction result with an explanation for each user request for 138 IoT traffic inference, i.e., to determine if the traffic is benign or belongs to any particular attack type. 139 We also allow IDS-Agent to handle requests for customized detection sensitivities or to detect 140 new attack categories from the given IoT traffic flow. IDS-Agent is equipped with a) an abundant 141 toolbox containing special IDS tools such as ML models for classification and general tools such 142 as search engines to retrieve external knowledge and b) a memory and knowledge base storing the current session information, long-term memory from previous sessions, and supportive documents. 143 These tools, memory, and external knowledge will be integrated to guide the decision-making of 144 IDS-Agent in a structured manner, as detailed in the sequel. 145

146 147

3.2 PIPELINE OF IDS-AGENT

The pipeline of our IDS-Agent is inspired by the ReAct agent (Yao et al., 2023). The user request is fulfilled by executing a sequence of action steps $\{a_1, a_2, \dots\}$, where each action step is generated by a core LLM based on previous reasoning and observations. For any input user request with a specification of the traffic flow to be inferred, an initial observation o_0 is constructed by concatenating the user request with a system prompt, including a description for each available tool. This initial observation serves as the context for the agent to understand the task, facilitating subsequent reasoning and action generation.

¹⁵⁵ Specifically, IDS-Agent iterates over the following three steps after the construction of o_0 :

1) Reasoning: The core LLM generates a thought (in plain text) about the next action by $r_i = LLM(s_i)$, where $s_i = \{o_0, \{r_1, a_1, o_1\}, \dots, \{r_{i-1}, a_{i-1}, o_{i-1}\}\}$ is the short-term memory of the current session up to the (i - 1)-th iteration (with $s_1 = \{o_0\}$). Reasoning can optionally adopt long-term memory from previous sessions for in-context demonstration.

160 2) Action generation: The action is generated based on the reasoning/thought by $a_i = \text{LLM}(r_i, s_i)$.

161 Notably, each action we generate is a structured JSON file containing an *action name* and an *action input*. The action input consists of the name of the tool(s) to be used and the associated settings or



Figure 1: Overview of IDS-Agent. Top (framework of IDS-Agent): IDS-Agent adopts 191 a core LLM to generate thoughts (i.e. reasoning) and actions based on input traffic and user 192 prompt. It is equipped with a toolbox for action execution and a memory base for knowledge 193 retrieval. IDS-Agent iteratively conducts thought generation, action generation and execution, and 194 observation update. **Bottom** (an example of reasoning trace): in this example, several classifiers 195 are adopted by the agent, with same number of classifiers predicting the input traffic as 'Benign' 196 and 'MITM-ArpSpoofing'. Based on this observation, IDS-Agent decides to perform 'Knowledge Retrieval' and 'Memory Retrival', and finally aggregate these observations, which leads to correct 197 attack detection.

- 199
- 200

parameters of each tool. Such a structured generation of the action allows its efficient and accurateexecution using the specified tools.

3) Observation update: After executing the generated action a_i , we obtain a new observation o_i by converting the outputs of the tool(s) into plain text.

The iterations terminate when the observation is updated by a 'final answer' headline followed by a JSON file. This JSON file, which encapsulates the final prediction on the traffic data and related analysis and explanation, will be the output of IDS-Agent.

208 209

3.3 ACTION SPACE AND TOOL DESIGN

Our IDS-Agent is designed with a comprehensive action space, allowing it to handle various tasks in the pipeline of data processing and classification through iterative reasoning and execution. The action space includes the key actions described below.

Data Extraction: The goal is to accurately extract network traffic records x specified in the user request from the dataset for further analysis. We design the data extraction tool as a function that takes the given line number or the flow ID as the input and outputs a structured traffic sample. Preprocessing: The goal is to clean, normalize, and transform the extracted data into a standard format for classification. Our preprocessing tools are designed as functions for diverse data analysis operations, including feature scaling, data encoding, handling missing values, and selecting important features for classification.

220 **Classification:** This action applies machine learning models to the preprocessed data to obtain 221 classification results, i.e., to predict whether the traffic is benign or falling into any malicious category. 222 The inputs to the classification tool include the preprocessed traffic features and a classifier, while the 223 outputs include the top-k labels and their corresponding confidence scores. Note that our classification 224 tool is more than a classifier; it advances ML-based IDS by adapting to diverse model types and 225 incorporating more information from the classification results into the inference procedure. The types 226 of classifiers used include Random Forest, SVM, MLP, Decision Tree, KNN, etc. It is important to note that the classification toolbox is extensible. In real-world deployments, users can easily add 227 new classifiers to the toolbox without the need to fine-tune the LLM. This flexibility allows for easy 228 updates and adaptation to new attack patterns or changes in the network environment. Moreover, 229 the classifiers can also be open-sourced models trained by third parties (callable through APIs), or 230 models locally trained based on the data collected by the user. 231

Knowledge Retrieval: This action aims to obtain knowledge regarding the particular types of attacks
 predicted by the classifiers. The knowledge can be either external knowledge and retrieved by calling
 search engines such as Google and Wikipedia API, or stored knowledge base or predefined rules and
 retrieved by RAG. The retrieved knowledge will be used to guide the aggregation step.

Long Memory Retrieval: Long Memory carrying information from previous sessions (discussed in Sec. 3.4). We add an instruction in the system prompt in o_0 such that long-term memory retrieval will be activated (through reasoning).

239 Aggregation: This action aims to comprehensively integrate the results from multiple steps of 240 classification action (based on different classifiers) to generate a structured final inference decision. 241 The core of the aggregation tool is an LLM where the prompt is designed to include 1) the detailed 242 results from the classifiers, 2) the short-term memory, and 3) demonstrative inputs and outputs 243 aggregated by the LLM from previous sessions. Note that the short-term memory also includes the 244 external knowledge previously extracted and the system prompts in o_0 . This system prompt, as shown 245 in Fig. 3, includes a specification for the detection sensitivity; it can also include an instruction for revealing new attack categories. Compared to naive aggregation, such as majority voting, our method 246 incorporates more information to resolve any discrepancies between different model outputs in a 247 more comprehensive way. 248

249 250

3.4 MEMORY AND KNOWLEDGE BASE

The memory and knowledge base of IDS-Agent stores 1) the short-term memory, 2) the long-term memory, and 3) supportive documents for IDS.

Short-term Memory (STR). The STR, as described in Sec. 3.2, includes the historical reasoning
 trace, actions, and observations of the current session in a structured format, and is renewed after
 each observation update. The major goal of the STM here is to track the agent's iterative reasoning
 process and ensure consistency between steps in real-time.

258 Long-term Memory (LTM). The LTM consists of agent decisions and contextual information from 259 previous use cases (Wang et al., 2024; Zhong et al., 2024). Here, a structured LTM example is defined by $\phi = \{t, x, R, A, O, \hat{y}\}$, where t is a timestamp, x is the feature vector after preprocessing, R =260 $[r_1, \dots, r_n]$ is the reasoning trace, $A = [a_1, \dots, a_n]$ contains all the action steps, $O = [o_1, \dots, o_n]$ 261 are the observations, and \hat{y} is the final label prediction. The long-term memory base can be initialized 262 by running the agent on a validation dataset. During the inference, only sessions with correct agent 263 decisions will be stored in the long-term memory base. In a real-world intrusion detection scenario, 264 such correctness can be validated by human experts. 265

In our framework, LTM retrieval provides the agent with additional information while aggregating the results for individual classifiers. Here, we set the LTM retriever input as the current timestamp t and observations from previous data processing and classification actions, denoted by $\tilde{O} = [o_1, \ldots, o_{m-1}]$, where m is an arbitrary iteration where the LTM retrieval kicks in. The retriever obtains the top-k relevant final reasoning based on the weighted sum of timestamp distance and the cosine similarity between the embedding of \tilde{O} and the observation embeddings $O^{(j)}$ of previous LTM examples $\{\phi^{(1)}, \dots, \phi^{(L)}\}$. Specifically, we obtain the top-k solutions to

$$\operatorname{argmax}_{j}[\lambda_{1}r(t, t^{(j)}) + \lambda_{2}\operatorname{cosim}(\mathcal{E}(\tilde{O}), \mathcal{E}(O^{(j)}))],$$
(1)

where $\mathcal{E}(\cdot)$ is the encoder. $r(t, t^{(j)}) = 1 - (t - t^{(j)}) / \max_k(t - t^{(k)})$ is the recency of the memory. The equation ensures that both recent observations and content-wise similar observations are considered to address the evolving nature of intrusion data. Then, the observation o_m for the iteration m contains the input-prediction pairs $(x^{(j)}, \hat{y}^{(j)})$ of each retrieved $\phi^{(j)}$. In our experiments, we set k = 5 to retrieve the top 5 relevant structured LTM examples.

280 External Knowledge. In addition to the external knowledge obtained by calling search engines, such as Google and Wikipedia, IDS-Agent is also equipped with a vector database $\{\psi^{(1)}, \dots, \psi^{(K)}\}$ 281 containing related research papers and intrusion detection blogs (both parsed into chunks with fixed 282 token length). The retrieval from this knowledge base is similar to the retrieval of LTM. We obtain 283 the top-k solutions to $\operatorname{argmax}_i \operatorname{cosim}(\mathcal{E}(q), \mathcal{E}(\psi^{(j)}))$, where q is the query generated by the core 284 LLM (based on the reasoning) as the action input to action step a_m for an arbitrary iteration m 285 for knowledge retrieval. The retrieved document chunks are summarized (for compression) using 286 an LLM (may be the same as the core LLM or an independent LLM) and are used to update the 287 observation o_m . The definition, characteristics, and detection methods for various attacks recorded 288 in the retrieved chunks will facilitate IDS-Agent to better understand the potential risks while 289 aggregating the classification results for the ML models. For example, if an attack type can potentially 290 lead to catastrophic results, IDS-Agent will be more sensitive to it when any classifier makes such 291 a prediction.

292 293

295

296

273

4 EXPERIMENTS

4.1 EXPERIMENT SETTINGS

Dataset. This paper focuses on intrusion detection in the IoT environment, which presents more complexities and challenges than traditional networks. We consider the following two datasets commonly used in previous works.

ACI-IOT'23 (Bastian et al., 2023): This dataset contains both benign and malicious network traffic captured from a variety of IoT devices. The dataset includes simulations of several attack types, such as Reconnaissance (e.g., Host Discovery, OS Scan, Ping Sweep, and Port Scan), DoS (e.g., ICMP Flood, SYN Flood, UDP Flood, and Slowloris) and Brute Force (e.g., Dictionary Attacks).
We randomly select 10% of the data to train the ML-models for classification. For evaluation, we construct a test dataset from the remaining samples in ACI-IoT'23 by randomly selecting 200 benign samples and 20 samples per attack category.

2) CIC-IoT 2023 (Neto et al., 2023): This dataset simulates large-scale, real-time IoT environments and comprises 33 distinct types of attacks, which is even more difficult for intrusion detection than ACI-IoT'23. The dataset includes network traffic from a broad IoT topology with 105 devices. From the dataset, we selected the 24 most common attack types along with benign samples to create our training and testing datasets. The remaining 9 attack types were excluded from the training data and designated as unknown attacks for evaluating zero-day attack detection performance. Again, we use 10% of the data for training the machine learning classifiers. The test dataset is constructed with 100 benign samples and 10 samples for each attack type.

314

Evaluation Metrics. We are interested in the performance of IDS-Agent in both the binary classification of benign and malicious flows and the multi-classification that also requires recognizing the specific attack type when a flow is deemed malicious. For binary classification, we use *accuracy*, and *false alarm rate* (FAR) as the evaluation metrics. Accuracy is the ratio of correctly predicted samples to the total number of samples in the dataset, measuring the overall effectiveness of the IDS-Agent in detecting both benign and malicious flows. FAR measures the proportion of false positives (benign flow incorrectly classified as malicious).

For *multi-classification*, we treat each attack category (and the benign category) as a class. We use the per-class *precision*, *recall*, and *F1-score* as the evaluation metrics. Detailed results for the detection of each attack category are deferred to the appendix due to space limitations. In the main paper, we report the *macro-averaged* precision, recall, and F1-score across all classes as the overall performance
 of the IDS being evaluated. This macro-average is computed by averaging the metric over all classes
 with equal weights.

4.2 IMPLEMENTATION DETAILS

The core LLM of IDS-Agent. In our experiments, we consider three LLM choices: GPT-3.5 Turbo, GPT-40-mini, GPT-40.

332

328

329

333 **Tool Design.** 1) **Data preprocessing tool**. We implement Python functions that will be called 334 sequentially for data preprocessing. First, we remove from each flow the fields irrelevant to the traffic 335 features, including the label, time-stamp, and flow ID. Second, we encode the non-numeral fields 336 into numbers, including the connection type and protocol type. Third, we conduct feature selection 337 based on an F-test for linear dependency between features and labels. Finally, the extracted features are standardized. 2) Classification tool. The core of the classification tool is an ML model for 338 intrusion detection, which can be self-trained or an off-the-shelf model trained by a third party. Here, 339 we pretrained six (multi-)classifiers, including random forest (RF), K-Nearest Neighbors (KNN), 340 logistic regression (LR), decision tree (DT), multi-layer perceptrons (MLP), and support vector 341 classifier (SVC) using the training set from our benchmarks. The output of the classification tool is 342 the top-3 label predictions with their confidence scores. 3) Knowledge retrieval tool. We construct a 343 knowledge base for various IoT attacks by collecting 50 online blogs and 50 research papers. These 344 documents are then split into chunks of 1000 tokens with an overlap of 200 tokens. These chunks 345 were embedded using the OpenAI encoder and stored in a vector database powered by ChromaDB.

346 347 348

4.3 BASELINES

349 1) Machine learning methods. We compare IDS-Agent with the state-of-the-art ML-based IDS,
 350 which uses a quantum-annealing method for feature selection (Davis et al., 2024). In Appendix A.3,
 351 we also compare our IDS-Agent with six IDSs using the six ML models we pretrained, respectively.

352 2) LLM-based methods. We compared our IDS-Agent with the latest GPT-4-based intrusion 353 detection approach (Zhang et al., 2024), which leverages the model's reasoning capabilities and in-354 context demonstrations. The authors particularly demonstrate that providing GPT-4 with a few labeled 355 examples can improve the accuracy of intrusion detection. We further improve the performance of 356 this baseline by first clustering their in-context examples using a Gaussian Mixture Model (GMM) 357 and then selecting in-context examples from different clusters to cover as many attack cases as 358 possible. Compared to this baseline with a fixed set of demonstrations, IDS-Agent retrieves LTM for demonstration dynamically based on input similarity. 359

360 3) Ensemble learning methods (majority vote). We create a strong baseline by ensemble the results
 a61 from the six ML models we pretrained through majority voting.

363 4.4 EXPERIMENT RESULTS

365 Quantitative Analysis. The quantitative results on the ACI-IoT'23 dataset are summarized in Table 366 1. Our IDS-Agent achieves the best general performance when GPT-40 is used as the core LLM. 367 Comparable performance is achieved when GPT-4o-mini is used as the core LLM, which is over 368 60% more cost-effective than GPT-3.5 Turbo. IDS-Agent also shows a clearly better recall than the baselines. In particular, when detecting UDP flood attacks, IDS-Agent achieves a recall of 369 0.80 compared with 0.20 for the LLM baseline (GPT-40) and 0.55 for the majority voting baseline. 370 We also found that for some attack types such as DNS-Flood, Slowloris, and Dictionary attacks, 371 IDS-Agent demonstrates slightly lower precision than majority voting. This decrease in precision 372 is due to the IDS-Agent's heightened sensitivity to high-threat attacks, leading it to classify a 373 sample as an attack even with less than 50% voting from the classifiers involved. 374

It is worth noting that the vanilla GPT-4-based method (Zhang et al., 2024), despite enhancements
 through in-context learning, performs unsatisfactorily. Moreover, when there are more attack cate gories, more in-context demonstrations will be required (to effectively inform all attack categories)
 for their method, posing a significant challenge due to GPT-4's limited token input length and the

378 Table 1: Binary-classification and multi-classification performance of IDS-Agent compared with 379 baseline approaches on the ACI-IoT'23 and CIC-IoT'23 datasets. We compare our method with 380 vanilla GPT-4o-based method enhanced by in-context learning (Zhang et al., 2024), a Random Forest classifier that uses a quantum-annealing method for feature selection (Davis et al., 2024), and a strong 381 baseline of majority voting by six ML classifiers. 382

Dataset	Metric Types	Metrics	GPT-40	RF	Majority Vote	IDS-Agent (GPT-3.5)	IDS-Agent (GPT-4o- mini)	IDS-Agent (GPT-40)
	Binary-Class	Binary-Class Accuracy↑ FAR↓	0.721 0.497	0.890 0.060	0.960 0.020	0.954 0.050	0.963 0.041	0.965 0.030
ACI-IoT'23	B Multi-Class	Multi-Class Accuracy ↑ Macro Avg Precision ↑ Macro Avg Recall ↑ Macro Avg F-Score ↑	0.678 0.682 0.754 0.682	0.790 0.785 0.760 0.750	0.980 0.980 0.961 0.962	0.972 0.971 0.952 0.923	0.976 0.980 0.972 0.974	0.980 0.982 0.972 0.975
	Binary-Class	Binary-Class Accuracy↑ FAR↓	0.750 0.144	0.825 0.050	0.882 0.040	0.876 0.050	0.894 0.030	0.904 0.030
CIC-IoT'23	B Multi-Class	Multi-Class Accuracy ↑ Macro Avg Precision ↑ Macro Avg Recall ↑ Macro Avg F-Score ↑	0.610 0.580 0.450 0.510	0.751 0.755 0.692 0.680	0.771 0.760 0.700 0.699	0.762 0.759 0.694 0.700	0.788 0.790 0.723 0.722	0.795 0.800 0.733 0.750

399 400

increase in cost. Our method integrates ML models and utilizes RAG to retrieve the most related 401 knowledge from the memories, significantly reducing the token cost. 402

403 Table 1 also presents the evaluation results on the CIC-IoT'23 dataset. Powered by both reasoning 404 ability and tool calls, the IDS-Agent achieves higher accuracy than the LLM baseline and majority 405 voting method. Moreover, IDS-Agent achieves high detection accuracy on some very challenging 406 attacks, such as ArpSpoofing and Host Discovery. The detailed results for each attack type are 407 deferred to the Appendix A.3.

408 **Case study** While the majority voting baseline shows relatively strong performance, their clas-409 sification results usually lack interpretability (Yang et al., 2022). In contrast, we leverages the 410 reasoning capabilities of LLMs to enhance interpretability which may help to improve the detection 411 performance. One example is illustrated on the right of Figure 1 where the IDS-Agent concludes 412 that the MITM ARP-Spoofing classification is more likely, as it appears in three classifiers' top-3 predictions with significant confidence in general. Another case for the decision-making is shown in 413 Figure 2, with the ground truth label being 'reconnaissance activities'. Despite 3 out of 6 machine 414 learning models predicting the traffic as benign, IDS-Agent accurately labels it as a reconnaissance 415 attack. This decision is rooted in the IDS-Agent's understanding that both Host Discovery and OS 416 Scanning belong to reconnaissance attacks. By considering the relationship between these types of 417 attacks, the IDS-Agent demonstrates its ability to detect subtle patterns that individual classifiers 418 may overlook, thus improving the overall detection accuracy. Additional examples are provided in 419 the appendix to demonstrate its decision-making process further. 420

4.5 ZERO-DAY ATTACK 421

422 We find that IDS-Agent not only excels in detecting known attacks but also shows promise in 423 identifying zero-day threats. Table 2 presents the detection results on the CIC-IoT'23 dataset, 424 comparing IDS-Agent with two recent zero-day intrusion detection methods: ACGAN (Zhao 425 et al., 2022) and RealNVP (Matejek et al., 2024). We reproduced ACGAN using its official code 426 and RealNVP based on the FrEIA framework, training separate normalizing flows for benign and 427 malicious samples. We set the detection threshold of RealNVP at the 90th percentile of the negative 428 log-likelihood from the INN model; samples exceeding this threshold are considered unknown attacks 429 or out-of-distribution (OOD). We selected nine attack types from the CIC-IoT'23 dataset as 'unseen' for evaluation, excluding them from training data. Each attack type was evaluated with 50 samples. 430 IDS-Agent's classifiers were trained on data from other known attack categories. Additionally, we 431 enhanced the system prompt with specific instructions to classify ambiguous samples as "Unknown"

32	
33	Thought: I have obtained the predictions from multiple classifiers. Here are the results:
84	1. **Random Forest**: - Recon-HostDiscovery: 0.2800, MITM-ArpSpoofing: 0.2500, BenignTraffic: 0.2400
_	2. **K-Nearest Neighbors**: BenignTraffic: 0.4000, Recon-HostDiscovery: 0.4000, MITM-ArpSpoofing: 0.2000
S	3. **Logistic Regression**: Recon-OSScan: 0.3909, Recon-PortScan: 0.2442, Recon-HostDiscovery: 0.1493
	4. ***/VIP***: beingn iramic: 0.5555, WITW-Arpspooling: 0.3421, Recon-USScatt: 0.0386 5. ***Support Vertar Classifier**: beingerTraffic: 0.6321, Becon-USScatt: 0.2366 Becon-USScatt: 0.0855
	support vector lessment : Dengrimaria: 0.0220, Recon-PortScan: 0.0000, DDoS-UDP Flood: 0.0000 secon-HostDiscovery: 1.0000, Recon-PortScan: 0.0000, DDoS-UDP Flood: 0.0000
	Now, I will analyze the results to determine the final classification. The predictions show a mix of benign and attack classifications, I will search for
	more information regarding the features of the attacks and benign traffic to make a more informed decision.
)	Action: knowledge_retriever
	Action Input: {"query": "what is Recon-HostDiscovery and how to detect it"
	Observation: The provided context does not specifically mention ("Recon-HostDiscovery(" or provide details on what it is or how to detect it. Based on provide details on what it is or how to detect it. Based on the specification of the specific details on what it is or how to detect it. Based on the specification of the s
	general knowledge, necton-noscibicovery typically relevant to identify active devices open ports and services running on those devices
)	Final Answer:
_	{
5	"line_number": 229,
L	"analysis": "The traffic features were classified by multiple models with discrepancies. The Random Forest and Decision Tree models indicated Recon-
	HostDiscovery, while Logistic Regression suggested Recon-OSScan. The K-Nearest Neighbors, MLP, and Support Vector Classifier models leaned towards begins the centre to find ensure and the active of the predictions the find elevision of the predictions the find elevision of the predictions of the predictions the find elevision of the predictions of the p
	"inredicted label for 1". "Recon-HostDiscovery"
	"predicted label top 2": "Recon-OSScan",
7	"predicted_label_top_3": "BenignTraffic",
2	}
1	
)	Figure 2: A case study of the IDS-Agent detecting reconnaissance activities. Despite most machine

learning models classifying the traffic as benign, the IDS-Agent ultimately predicts the final label 451 as a reconnaissance attack. This decision is based on its understanding that both Host Discovery and 452 OS Scan belong to reconnaissance activities. 453

456

457

458

459

461

463 464

455 if most classifiers' confidence was below a certain threshold. We set the threshold as 0.7 in our experiments. The recall for IDS-Agent in predicting an 'unknown attack' from these unseen classes was measured, achieving a top-1 recall of 0.61. Notably, IDS-Agent showed high recall for Vulnerability Scan and SQL Injection attacks, likely due to their distinct deviation from the training data distribution, resulting in low classification confidence. This indicates IDS-Agent's capability 460 to detect unknown attacks, particularly when diverse machine learning models yield divergent results with low confidence or when traffic features are anomalous. The recall for benign examples decreased from 0.91 to 0.86 after introducing the zero-day detection prompt. This drop is attributed to the 462 tendency of the model to classify OOD benign examples as unknown attacks.

Table 2: Comparisons of the recall of different zero-day attack detection methods.

Methods	Backdoor	DNS Spoofing	Uploading Attack	XSS	Dictionary BruteForce	Command Injection	VulScar	n Browser Hijacking	SQL	Avg Recall
ACGAN	0.38	0.35	0.59	0.32	0.36	0.42	0.88	0.05	0.38	0.41
RealNVP	0.45	0.37	0.68	0.49	0.42	0.45	0.89	0.03	0.45	0.47
IDS-Agent	0.64	0.46	0.74	0.65	0.45	0.55	0.95	0.15	0.86	0.61

4.6 ABLATION STUDY

474 We conduct an ablation study focusing on two key modules: the Knowledge Retrieval Module and 475 the Long-Term Memory Module. We evaluate how each module affects the detection performance on 476 two settings: 1) 'in-distribution' setting with all attack labels known and 2) zero-day attack setting 477 in Section 4.5 with a subset of unknown attacks. We also evaluate the performance under different 478 detection sensitivities.

479

473

480 Effect of Knowledge Retrieval Module The Knowledge Retrieval Module is designed to augment 481 the classifiers' outputs with relevant information from an external knowledge base, enhancing the 482 agent's understanding of complex attack patterns. To assess its impact, we disable this module and compare the performance of the modified agent with the original IDS-Agent. Table 3 summarizes 483 the detection performance with and without the Knowledge Retrieval Module. When the module is 484 removed, we observe significant decrements in the detection accuracy for both settings. Specifically, 485 the recall for detecting zero-day attacks drops from 0.61 to 0.42. This indicates that without access to

486 external knowledge, the agent has a diminished ability to recognize patterns that are not represented 487 in the training data. 488

 Table 3: Effect of Knowledge Retrieval Module

Table 4: Effect of Long-Term Memory Module

In-Distribution	0.733	0.702
Zero-Day	0.610	0.560
	In-Distribution Zero-Day	In-Distribution0.733Zero-Day0.610

496 Effect of Long-Term Memory Module The Long-Term Memory Module allows IDS-Agent to 497 maintain a history of previous observations and decisions, which is essential for detecting attacks 498 exhibiting temporal dependencies and utilizing previous success experiences. To evaluate its effect, we disable the Long-Term Memory Module and assess the agent's performance on the same evaluation 499 set. We set the λ_1 and λ_2 as 0.5 in Eq. 1 to balance the recency similarity to retrieve the most relevant 500 past examples from the agent's LTM. 501

502 As presented in Table 4, the removal of the Long-Term Memory Module leads to a degradation in detection performance, particularly for attacks that unfold over time, such as Brute Force and 504 Distributed Denial of Service (DDoS) attacks. The overall detection accuracy decreases from 0.733 505 to 0.702, and the recall for zero-day attacks drops to 0.56. The decrease in performance underscores the importance of the Long-Term Memory Module in capturing temporal features and improving 506 the agent's ability to detect attacks that evolve over time. By retaining historical information, 507 IDS-Agent can identify suspicious patterns that may not be apparent when considering individual 508 events in isolation. 509

510 **Detection Sensitivity.** Sensitive configuration is a critical function of intrusion detection systems (IDS). In signature-based IDS, experts need to manually adjust detection sensitivity, which can be 511 both costly and time-consuming (Díaz-Verdejo et al., 2022). In contrast, the detection sensitivity 512 of our IDS-Agent can be easily adjusted through input prompts. Here, we can optionally instruct 513 the IDS-Agent to operate under three different sensitivity levels: aggressive, balanced, and 514 conservative. This sensitivity level controls the trade-off between the false alarm rate and the missed 515 detection rate. The results are shown in Table 5. For the aggressive detection, the IDS-Agent 516 achieves high recall for attacks (0.97) but lower recall for benign examples (0.90). Conversely, 517 conservative detection shows relatively lower recall for attacks (0.85) but higher recall for benign 518 examples (0.98). The complete prompts and detection results are shown in Appendix A.4. 519

520 Table 5: The classification results of different detection sensitivities. The core LLM is GPT-40. We compute the precision, recall and F1-score under different sensitivities on the ACI-IoT'23 dataset.

Sensitivity	Aggressive				Balance		Conservative		
Metrics	Precision	Recall	F1-score	Precision	Recall	F1-score	Precision	Recall	F1-score
Benign	0.96	0.90	0.92	0.87	0.96	0.91	0.60	0.98	0.75
Attack (Macro Avg)	0.97	0.97	0.97	0.98	0.95	0.96	0.95	0.85	0.87

527 528

521

489

490

529

530

5 CONCLUSION

In this paper, we propose IDS-Agent, the first LLM-powered agent for intrusion detection. We 531 design an iterative reasoning-followed-by-action pipeline for IDS-Agent to extract data from the 532 network traffic, preprocess the data, consult different machine learning models for classification results and details, retrieve both internal and external knowledge, and summarize the final detection inference. 534 These agent actions are facilitated by a memory module and a wide array of tools for intrusion detection and general purposes. Empirically, IDS-Agent outperforms diverse types of SOTA IDSs 536 on ACI-IoT'23 and CIC-IoT'23. We find that when classifiers produce discrepancy predictions, IDS-Agent can utilize inherent and external knowledge to assist decision-making. Moreover, IDS-Agent can be easily adapted to detect zero-day attacks, exhibiting better performance than 538 existing methods. Finally, we find that IDS-Agent effectively follows the sensitivity instructions in detection without requiring expert intervention or additional tuning.

540 REFERENCES

553

554

555

556

559

560

561

562

566

567

568

569

575

580

- 542 Mahyar Abbasian, Iman Azimi, Amir M. Rahmani, and Ramesh Jain. Conversational health agents:
 543 A personalized llm-powered agent framework, 2024.
- Tanveer Ahmad and Dongdong Zhang. Using the internet of things in smart energy systems and networks. *Sustainable Cities and Society*, 68:102783, 2021.
- 547 Yara Alghofaili and Murad A. Rassam. A trust management model for iot devices and services
 548 based on the multi-criteria decision-making approach and deep long short-term memory technique.
 549 Sensors, 22(2), 2022.
- Mutep Y AlYousef and Nabih T Abdelmajeed. Dynamically detecting security threats and updating
 a signature-based intrusion detection system's database. *Procedia Computer Science*, 159:1507–1516, 2019.
 - Karan Bajaj and Amit Arora. Dimension reduction in intrusion detection features using discriminative machine learning approach. *International Journal of Computer Science Issues (IJCSI)*, 10(4):324, 2013.
- Nathaniel Bastian, David Bierbrauer, Morgan McKenzie, and Emily Nack. ACI IoT Network Traffic
 Dataset 2023. IEEE Dataport, December 29 2023.
 - Luca Catarinucci, Danilo de Donno, Luca Mainetti, Luca Palano, Luigi Patrono, Maria Laura Stefanizzi, and Luciano Tarricone. An iot-aware architecture for smart healthcare systems. *IEEE Internet of Things Journal*, 2(6):515–526, 2015.
- Robin Chataut, Alex Phoummalayvane, and Robert Akl. Unleashing the power of iot: A comprehensive review of iot applications and future prospects in healthcare, agriculture, smart homes, smart cities, and industry 4.0. *Sensors*, 23(16), 2023.
 - Patrick J Davis, Sean Coffey, Lubjana Beshaj, and Nathaniel D Bastian. Quantum machine learning for feature selection in internet of things network intrusion detection. In *Quantum Information Science, Sensing, and Computation XVI*, volume 13028, pp. 78–92. SPIE, 2024.
- Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and
 Yu Su. Mind2web: Towards a generalist agent for the web, 2023.
- Jesús Díaz-Verdejo, Javier Muñoz-Calle, Antonio Estepa Alonso, Rafael Estepa Alonso, and Germán
 Madinabeitia. On the detection capabilities of signature-based intrusion detection systems in the
 context of web attacks. *Applied Sciences*, 12(2), 2022. ISSN 2076-3417.
- Qiuhan Gu. Llm-based code generation method for golang compiler testing. In *Proceedings of the* 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2023, pp. 2201–2203, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400703270. doi: 10.1145/3611643.3617850. URL https://doi.org/10.1145/3611643.3617850.
- Azidine Guezzaz, Said Benkirane, Mourade Azrour, and Shahzada Khurram. A reliable network
 intrusion detection approach using decision tree with enhanced data quality. *Security and Communication Networks*, 2021(1):1230593, 2021.
- Yuejun Guo, Constantinos Patsakis, Qiang Hu, Qiang Tang, and Fran Casino. Outside the comfort
 zone: Analysing llm capabilities in software vulnerability detection. In *European Symposium on Research in Computer Security*, pp. 271–289. Springer, 2024.
- Izzeddin Gur, Hiroki Furuta, Austin V Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and Aleksandra Faust. A real-world webagent with planning, long context understanding, and program synthesis. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=9JQtrumvg8.
- Philokypros Ioulianou, Vasileios Vasilakis, Ioannis Moscholios, and Michael Logothetis. A signature based intrusion detection system for the internet of things. In *Information and Communication Technology Form*, June 2018.

594 595 596	Huaping Jia, Jun Liu, Min Zhang, Xiaohu He, and Weixi Sun. Network intrusion detection based on ie-dbn model. <i>Computer Communications</i> , 178:131–140, 2021.
597 598 599 600	Ansam Khraisat, Iqbal Gondal, and Peter Vamplew. An anomaly intrusion detection system using c5 decision tree classifier. In <i>Trends and Applications in Knowledge Discovery and Data Mining: PAKDD 2018 Workshops, BDASC, BDM, ML4Cyber, PAISI, DaMEMO, Melbourne, VIC, Australia, June 3, 2018, Revised Selected Papers 22</i> , pp. 149–155. Springer, 2018.
601 602	Ansam Khraisat, Iqbal Gondal, Peter Vamplew, and Joarder Kamruzzaman. Survey of intrusion detection systems: techniques, datasets and challenges. <i>Cybersecurity</i> , 2(1):1–22, 2019.
603 604 605 606	Gisung Kim, Seungmin Lee, and Sehun Kim. A novel hybrid intrusion detection method integrating anomaly detection with misuse detection. <i>Expert Systems with Applications</i> , 41(4):1690–1700, 2014.
607 608	Junkai Li, Siyu Wang, Meng Zhang, Weitao Li, Yunghwei Lai, Xinhui Kang, Weizhi Ma, and Yang Liu. Agent hospital: A simulacrum of hospital with evolvable medical agents, 2024.
609 610 611 612	Wenchao Li, Ping Yi, Yue Wu, Li Pan, and Jianhua Li. A new intrusion detection system based on knn classification algorithm in wireless sensor network. <i>Journal of Electrical and Computer Engineering</i> , 2014(1):240217, 2014.
613 614 615	Ziadoon Kamil Maseer, Robiah Yusof, Nazrulazhar Bahaman, Salama A Mostafa, and Cik Fer- esa Mohd Foozy. Benchmarking of machine learning for anomaly based intrusion detection systems in the cicids2017 dataset. <i>IEEE access</i> , 9:22351–22370, 2021.
616 617 618	Brian Matejek, Ashish Gehani, Nathaniel D Bastian, Daniel J Clouse, Bradford J Kline, and Susmit Jha. Safeguarding Network Intrusion Detection Models from Zero-day Attacks and Concept Drift. <i>AAAI Workshop on Artificial Intelligence for Cyber Security (AICS)</i> , 2024.
619 620 621 622 623	Mokhtar Mohammadi, Tarik A Rashid, Sarkhel H Taher Karim, Adil Hussain Mohammed Aldalwie, Quan Thanh Tho, Moazam Bidaki, Amir Masoud Rahmani, and Mehdi Hosseinzadeh. A compre- hensive survey and taxonomy of the svm-based intrusion detection systems. <i>Journal of Network</i> <i>and Computer Applications</i> , 178:102983, 2021.
624 625 626	S. Nagaraju, B. Shanmugham, and K. Baskaran. High throughput token driven fsm based regex pattern matching for network intrusion detection system. <i>Materials Today: Proceedings</i> , 47: 139–143, 2021.
627 628 629	E. C. P. Neto, S. Dadkhah, R. Ferreira, A. Zohourian, R. Lu, and A. A. Ghorbani. CICIoT2023: A real-time dataset and benchmark for large-scale attacks in IoT environment. <i>Sensor</i> , 2023. (submitted to Journal of Sensors).
630 631 632	Souradip Roy, Juan Li, Bong-Jin Choi, and Yan Bai. A lightweight supervised intrusion detection mechanism for iot networks. <i>Future Generation Computer Systems</i> , 127:276–285, 2022.
633 634 635	Wenqi Shi, Ran Xu, Yuchen Zhuang, Yue Yu, Jieyu Zhang, Hang Wu, Yuanda Zhu, Joyce Ho, Carl Yang, and May D. Wang. Ehragent: Code empowers large language models for few-shot complex tabular reasoning on electronic health records, 2024.
636 637 638 639	Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning, 2023. URL https://arxiv.org/abs/2303.11366.
640 641 642	Ivan Studnia, Eric Alata, Vincent Nicomette, Mohamed Kaâniche, and Youssef Laarouchi. A language-based intrusion detection approach for automotive embedded networks. <i>International Journal of Embedded Systems</i> , 10(1):1–12, 2018.
643 644 645	Farhan Ullah, Hamad Naeem, Sohail Jabbar, Shehzad Khalid, Muhammad Ahsan Latif, Fadi Al- Turjman, and Leonardo Mostarda. Cyber security threats detection in internet of things using deep learning approach. <i>IEEE access</i> , 7:124379–124389, 2019.
647	Abhishek Verma and Virender Ranga. Machine learning based intrusion detection systems for iot applications. <i>Wireless Personal Communications</i> , 111(4):2287–2310, 2020.

- 648 Weizhi Wang, Li Dong, Hao Cheng, Xiaodong Liu, Xifeng Yan, Jianfeng Gao, and Furu Wei. 649 Augmenting language models with long-term memory. Advances in Neural Information Processing 650 Systems, 36, 2024. 651 652 Zhen Yang, Xiaodong Liu, Tong Li, Di Wu, Jinjiang Wang, Yunwei Zhao, and Han Han. A 653 systematic literature review of methods and datasets for anomaly-based network intrusion detection. Computers Security, 116:102675, 2022. ISSN 0167-4048. doi: https://doi.org/10.1016/j.cose. 654 2022.102675. URL https://www.sciencedirect.com/science/article/pii/ 655 S0167404822000736. 656 657 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 658 React: Synergizing reasoning and acting in language models. In International Conference on 659 Learning Representations (ICLR), 2023. 660 661 Yangyang Yu, Haohang Li, Zhi Chen, Yuechen Jiang, Yang Li, Denghui Zhang, Rong Liu, Jordan W. 662 Suchow, and Khaldoun Khashanah. Finmem: A performance-enhanced llm trading agent with 663 layered memory and character design, 2023. 664 665 Han Zhang, Akram Bin Sediq, Ali Afana, and Melike Erol-Kantarci. Large language models in 666 wireless application design: In-context learning-enhanced automatic network intrusion detection, 667 2024. URL https://arxiv.org/abs/2405.11002. 668 669 Oingling Zhao, Minggiang Chen, Zonghua Gu, Siyu Luan, Haibo Zeng, and Samarjit Chakrabory. 670 Can bus intrusion detection based on auxiliary classifier gan and out-of-distribution detection. 671 ACM Transactions on Embedded Computing Systems (TECS), 21(4):1–30, 2022. 672 Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. Gpt-4v(ision) is a generalist web 673 agent, if grounded. arXiv preprint arXiv:2401.01614, 2024. 674 675 Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. Memorybank: Enhancing large 676 language models with long-term memory. In Proceedings of the AAAI Conference on Artificial 677 Intelligence, volume 38, pp. 19724–19731, 2024. 678 679 680 А APPENDIX A 681
 - A.1 GENERAL PROMPT AND REASONING TRACE OF IDS-AGENT

The general prompt of the IDS-Agent is illustrated in Figure 3. The process begins with instructing the IDS-Agent to load network traffic data and perform feature preprocessing. Afterward, we utilize a range of classifiers to analyze the data. To enhance decision-making, the IDS-Agent retrieves prior successful examples from its knowledge base for comparison. In cases where discrepancies arise between the predictions of different models, we prompt the IDS-Agent to consult internal or external knowledge bases for additional insights to resolve the conflict. Finally, the IDS-Agent consolidates the findings and presents the result in a structured JSON format. Figure 4 provides an example of the reasoning trace produced by the IDS-Agent during this process.

682

683 684

685

686

687

688

689

690

A.2 ADDITIONAL CASE STUDIES

The cases in Figure 5 and Figure 6 highlight the enhanced reasoning ability of IDS-Agent with the knowledge retriever. It is shown that IDS-Agent not only considers the top-1 predictions but also the second and third predictions and their confidences. Moreover, in these examples, when the models have discrepancies in the predictions, the IDS-Agent automatically accesses external databases to extract additional knowledge, aiding in feature analysis and supporting its final decision. By leveraging these external knowledge sources, the IDS-Agent gains a deeper understanding of complex attack patterns and anomalies, enhancing both accuracy and decision-making. This dynamic capability allows the IDS-Agent to better adapt to new or evolving threats in the IoT environment.

/02	General Prompt:
703	You are a helpful assistant that can implement multi-step tasks, such as intrusion detection. I will give you the traffic features, you are asked to classifier results. You should plan
04	are asked to classify it using cools. The man output must be in a soon format according to the classifier results. For should plan first such as:
05	1. Load the traffic features from the CSV file. You can use the load_data_line tool to obtain the complete traffic.
06	2. Preprocessing the feature. This can be done using the data_preprocessing tool. Input the traffic in the original format.
07	load classifiers for classification. This can be done using the classifier tool. You can use multiple classifiers. The tool params include a classifier name, which must be one from {model_names} and the preprocessed features.
08	4. Retrieve previous successful reasonings to help you predict. This can be done using the memory_retrieve tool with the
)9	classifier's names and their classification results as input.
0	 When there are discrepancies/disagreements for different models, you can search from vector database/google/wiki to get more information about the difference of attacks to help you make decisions.
1	6. At the end, you should summarize the results from these classifiers and provide a final result. Summarize the classification
2	with Balance sensitivity, which means balancing the false alarm rate and the missing alarm rate. The predicted label should be the original format of classifier prediction. The final output format **must** here
3	Final Answer:
4	```json
5	{ 'line_number': \line_number,
6	'analysis': str, \nere is the Analysis, 'nredicted label ton 1': str
7	'predicted_label_top_2': str,
3	'predicted_label_top_3': str,
)	}
-)	Now, classify the traffic from file name {file_name} with index {line_number}
1	
>	Figure 3: General prompt and user input
2	
	· · · · · · · · · · · · · · · · · · ·
-	Thought: I have obtained the classification results from all five classifiers. Here
5	Thought: I need to follow the plan to classify the traffic features
6	from line number 230. I will start by loading the traffic features have don the highest confidence scores from the classifiers
7	Action: load_data_line ### Summary of Predictions:
3	Action Input: {"line_number": 230} Observation: flow, duration: 4.198554983, Header, Length: Neighbors (1.0000), and MLP (0.8543).
9	-**DDoS-HTTP_Flood**: Predicted by Logistic Regression (0.2586), MLP
)	Thought: I have successfully loaded the traffic features from line number 230. Now, I will proceed to preprocess the loaded traffic - **Recon-PortScan**: Predicted by Logistic Regression (0.1300) and K-Nearest
	features to prepare them for classification. **BenjanTraffic** Predicted by Support Vector Classifier (0.2320)
	Action: data_preprocessing Action Input: {"traffic_features": "flow_duration: 4.198554983, ### Final Decision:
	Header_Length: 10051.49,"} The most consistent prediction across the classifiers is **DDoS-SlowLoris**, which received the highest confidence scores from multiple models.
ļ.	Thought: I have successfully preprocessed the traffic features. Now, I will prepare the final output in the required JSON format.
5	Now, I will proceed to classify the preprocessed features using multiple classifiers to determine if the traffic record is an attack. Final Answer:
6	Action: classifier
,	[-0.355, -0.3088,]}
ξ	Observation: Top predictions: DDoS-SlowLoris: 0.9100, DDoS- HTTP: Flood: 0.0500, DoS-HTTP: Flood: 0.0200 SlowLoris attack based on the predictions from multiple classifiers, particularly
о О	Thought: I have obtained the classification results from the
2	Random Forest model. Now, I will classify the same preprocessed features using additional classifiers to gather more predictions. "predicted_label_top_1": "DDOS-SlowLong",
)	"predicted_label_top_3": "Recon-PortScan"
-	
3	Figure 4: An example of the reasoning trace and final answer.
4	
5	
16	A.3 THE PERFORMANCE OF DIFFERENT ML CLASSIFIERS

Table 6 shows the F-score of different ML classifiers on the ACI-IoT'23 dataset as well as our method.
Among the six classifiers, MLP achieves the highest F-score of 0.96. The IDS-Agent outperforms all six classifiers and the majority vote method. Moreover, we achieve a high F-score on the UDP Flood attack, while the majority method only has an F-score of 0.55.

747

Table 7 shows the F-score of different ML classifiers on the CIC-IoT'23 dataset as well as our method. Among the six classifiers, Random Forest achieves the highest F-score of 0.75. For the IDS-Agent, we use the GPT-40 as the core LLM. It is shown that our attack achieved a higher F-score compared with the majority vote classifier. Moreover, we achieved a higher F-score on the benign traffic compared with six classifiers and the majority vote method, which means our method

56	Thought: I have obtained predictions from all classifiers. Now, I will summarize the results to determine the final classification.
57	The predictions from the classifiers are as follows:
58	- Random Forest: DDoS-SYN_Flood (0.6458)
59	 k-nearest neighbors: Do->1n_riooa (1.0000) Logistic Pagressian: Does SNL Eload (0.6473)
60	- Eugratic Regression: DD0-5-010_1000 (0.0475)
61	- Support Vector Classifier: DDoS-SYN_Flood (0.6329)
62	The majority of classifiers predict "DDoS-SYN_Flood" as the top label. However, K-Nearest Neighbors predicts "DOS-SYN_Flood"
63	with a certainty of 1.0000. To resolve this discrepancy, I will search for more information regarding the features of DDoS and DoS attacks to understand the differences better.
64	Action: knowledge retriever
65	Action Input: {"query": "What are the features of DoS and DDoS attacks and their differences?"
66	Observation: DoS (Denial-of-Service) and DDoS (Distributed Denial-of-Service) attacks have several features and key differences:
~~	**Features of DoS Attacks:**
11	**Features of DDoS Attacks:**
8	**Key Differences:**
9	 **Source of Attack:** DoS attacks come from one source, while DDoS attacks come from multiple sources. **Volume of Traffic:** DDoS attacks typically generate a larger volume of traffic, making them more damaging.
0	 **Detection and Mitigation:** DDoS attacks are harder to detect and eliminate quickly due to their distributed nature, whereas
1	DoS attacks can be easier to manage since they originate from a single source.
2	Overall, Dubo attacks are generally considered more severe and challenging to defend against compared to Dos attacks.
5	{
	line number": 83,
	"analysis": "The majority of classifiers predict 'DDoS-SYN_Flood', indicating a distributed denial-of-service attack. The K-Nearest
	Neighbors classifier predicts 'DoS-SYN_Flood' with high certainty, but given the nature of the features and the consensus among other classifiers, 'DDOS-SYN_Flood' is deemed the more appropriate classification.",
	"predicted_label_top_1": "DDoS-SYN_Flood",
	"predicted_label_top_2": "DoS-SYN_Flood",
3	"predicted_label_top_3": "DDoS-SynonymousIP_Flood"
	}
	Figure 5: Additional case study on the DDoS attacks of the IDS-Agent. From the final analysis, it is
	shown that the IDS-Agent not only considers the predicted labels but also considers the confidence
2	of different classifiers
}	
ļ.	
	Thought: I have obtained the predictions from multiple classifiers. Here are the results:
C	 **Kanaom Forest**: wii IM-ArpSpooting: U.4200 - Benign Irafiit: U.2200 - Recon-PortScan: U.1300 **K-Naarset Neighbors**: Becon-PortScan: 0.4000 - MUTM-ArpSpoofing: 0.4000 - RenignTaffic: 0.2000
7	3. **Logistic Regression**: - BenignTraffic: 0.3689 - MITM-ArpSpoofing: 0.3105 - DoS-UDP Flood: 0.1391
3	4. **MLP**: - MITM-ArpSpoofing: 0.4019 - BenignTraffic: 0.3296 - Recon-HostDiscovery: 0.1095
	5. **Support Vector Classifier**: - BenignTraffic: 0.6327 - MITM-ArpSpoofing: 0.2328 - Recon-HostDiscovery: 0.0855
	6. **Decision Tree**: - Recon-OSScan: 1.0000 - Recon-PortScan: 0.0000 - DDoS-UDP_Flood: 0.0000 New Lyuill apply to the results to determine the final classification. The predictions show a mix of benign and attack electrifications with some
	classifiers leaning towards "MITM-ArpSpoofing" and others towards "BenignTraffic." The Decision Tree classifications, with some
	classifying it as "Recon-OSScan."To resolve the discrepancies, I will search for more information regarding the features of the attacks and benign
	traffic to make a more informed decision.
2	Action: knowledge_retriever
	A THE A THE A THE DATION OF A THE DATION OF A THE ATTRACT AND A THE ATTRACT ATTRACT AND A THE ATTRACT

- 794
- Observation: The features of MITM-ARP-Spoofing include: 1. **Interception of Communication**: The attacker positions themselves between two network devices, allowing them to intercept and potentially modify the data being transmitted. 795
- 2. **ARP Cache Poisoning**: The attacker sends falsified ARP messages to both devices, associating their MAC address with the IP address of the other device, thus redirecting traffic through their machine.... 796 **Final Answei**

797

- "line_number": 83, 798
- "analysis": "The traffic features indicate a potential MITM-ArpSpoofing attack based on the predictions from multiple classifiers, particularly 799 Random Forest and MLP. The Decision Tree classifier also suggests a reconnaissance attack, but the overall consensus leans towards MITM.", "predicted_label_top_1": "MITM-ArpSpoofing", 800 "predicted_label_top_2": "BenignTraffic", "predicted_label_top_3": "Recon-PortScan", 801
- 802 803

Figure 6: Additional case study on the MITM attacks of the IDS-Agent. From the final analysis, it is shown that the IDS-Agent not only considers the top 1 prediction but also considers the second and third predictions.

805 806 807

804

- 808
- has a lower false alarm rate, which is an important metric for intrusion detection. Figure 7 shows the 809 confusion matrix of the majority voting classifier and IDS-Agent.

Table 6: The F-score of different ML classifiers on the ACI-IoT'23 dataset. For the IDS-Agent, we use the GPT-40 as the core LLM.

Model	RF	LR	KNN	MLP	DT	SVC	Majority Vote	IDS-Agent
Benign	0.90	0.59	0.91	0.91	0.91	0.80	0.91	0.91
DNS Flood	0.95	0.10	0.80	0.95	0.91	0.91	1.00	0.95
Dictionary Attack	1.00	0.71	0.98	0.95	1.00	0.92	1.00	1.00
ICMP Flood	1.00	0.98	0.98	1.00	0.95	0.98	0.98	0.98
OS Scan	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
Ping Sweep	0.98	0.98	0.97	0.98	0.97	0.98	1.00	1.00
Port Scan	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
SYN Flood	1.00	1.00	1.00	1.00	0.98	1.00	1.00	1.00
Slowloris	1.00	0.43	1.00	1.00	1.00	0.97	1.00	1.00
UDP Flood	0.60	0.00	0.45	0.74	0.50	0.00	0.55	0.80
Vulnerability Scan	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
Macro Avg	0.95	0.71	0.92	0.96	0.93	0.87	0.96	0.97

Table 7: The F-score of different ML classifiers on the CIC-IoT'23 dataset. For the IDS-Agent, we use the GPT-40 as the core LLM.

Model	DT	KNN	LR	MLP	RF	SVC	Majority Vote	IDS-Agent
BenignTraffic	0.79	0.77	0.79	0.75	0.75	0.73	0.74	0.84
DDoS-ACK_Fragmentation	0.98	0.95	0.95	0.93	0.95	0.98	0.95	1.00
DDoS-HTTP_Flood	0.58	0.53	0.24	0.79	0.68	0.38	0.69	0.70
DDoS-ICMP_Flood	0.98	0.95	0.98	0.95	1.00	1.00	1.00	1.00
DDoS-ICMP_Fragmentation	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
DDoS-PSHACK_Flood	1.00	1.00	0.98	1.00	1.00	1.00	1.00	0.95
DDoS-RSTFINFlood	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
DDoS-SYN_Flood	0.72	0.08	0.63	0.73	0.76	0.64	0.75	0.75
DDoS-SlowLoris	0.76	0.74	0.00	0.89	0.79	0.44	0.79	0.82
DDoS-SynonymousIP_Flood	0.70	0.74	0.70	0.72	0.74	0.65	0.74	0.78
DDoS-TCP_Flood	0.69	0.08	0.66	0.67	0.69	0.67	0.67	0.71
DDoS-UDP_Flood	0.71	0.23	0.77	0.75	0.71	0.74	0.71	0.74
DDoS-UDP_Fragmentation	0.97	0.95	0.95	0.97	0.97	0.97	0.97	0.95
DoS-HTTP_Flood	0.65	0.79	0.73	0.85	0.76	0.89	0.83	0.84
DoS-SYN_Flood	0.62	0.57	0.31	0.72	0.76	0.37	0.72	0.80
DoS-TCP_Flood	0.32	0.60	0.00	0.00	0.32	0.00	0.17	0.33
DoS-UDP_Flood	0.53	0.65	0.55	0.52	0.48	0.46	0.48	0.46
MITM-ArpSpoofing	0.54	0.58	0.09	0.60	0.58	0.58	0.62	0.67
Mirai-greeth_flood	0.95	0.98	0.90	0.97	0.97	1.00	0.97	0.95
Mirai-greip_flood	0.98	0.97	0.89	0.98	0.98	1.00	0.98	0.95
Mirai-udpplain	0.95	0.98	1.00	0.98	0.98	1.00	0.98	1.00
Recon-HostDiscovery	0.54	0.45	0.40	0.48	0.55	0.39	0.47	0.53
Recon-OSScan	0.30	0.30	0.08	0.00	0.15	0.10	0.17	0.15
Recon-PortScan	0.41	0.36	0.24	0.31	0.44	0.31	0.37	0.31
Macro Avg	0.74	0.68	0.62	0.73	0.75	0.68	0.70	0.75

864 A.4 DETAILS FOR SENSITIVITY CUSTOMIZATION 865

866 We adjust the detection sensitivity by prompting the core LLM with "Summarize the classification with {sensitivity} sensitivity, {sensitivity details}". Here, 'sensitivity details' will be 'discover the 867 attack as the priority', 'balance the false alarm rate and the missing alarm rate', and 'do not alert 868 unless you are very sure', for 'sensitivity' being 'aggressive', 'balanced' and 'conservative'. The detection performances of IDS-Agent for different detection sensitivities are shown in Table 8. It is 870 shown that the 'Aggressive' command achieves a higher recall on the attacks while the 'Conservative' 871 command achieves a higher recall on the benign examples. The classification results, detailed in 872 Table 8 of the appendix, show that the IDS-Agent effectively follows these sensitivity instructions 873 without requiring expert intervention or additional tuning. 874

875 876

877 878 879

882 883

885

893

897

900

901

911

912

Sensitivity	Aggressive				Balance		Conservative			
Metrics	Precision	Recall	F1-score	Precision	Recall	F1-score	Precision	Recall	F1-score	
Benign	0.96	0.90	0.92	0.87	0.96	0.91	0.60	0.98	0.75	
DNS Flood	0.91	1.00	0.95	0.91	1.00	0.95	0.94	0.80	0.86	
Dictionary Attack	0.91	1.00	0.95	1.00	1.00	1.00	1.00	0.65	0.79	
ICMP Flood	0.95	1.00	0.89	0.95	1.00	0.98	0.95	1.00	0.98	
OS Scan	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
Ping Sweep	0.95	1.00	0.98	1.00	1.00	1.00	1.00	1.00	1.00	
Port Scan	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
SYN Flood	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
Slowloris	0.95	1.00	0.98	1.00	1.00	1.00	1.00	0.40	0.57	
UDP Flood	1.00	0.80	0.89	1.00	0.53	0.69	1.00	0.47	0.64	
Vulnerability Scan	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
Macro Avg	0.97	0.97	0.97	0.98	0.95	0.96	0.95	0.85	0.87	

A.5 THE ZERO-DAY ATTACK DETECTION DETAILS

894 We prompt GPT-40 to classify an example as an unknown attack if multiple classifiers output low 895 confidence for their top predictions or if there are conflicting predictions among different classifiers. This is based on our observation that, for unknown attacks, machine learning models typically exhibit 896 relatively low confidence levels, as shown in Figure 8. Specifically, we instruct the LLM to consider an example as a potential unknown attack if more than two models have low confidence (e.g., below 898 a threshold of 0.7). Moreover, if more than two models have low confidence or if different models 899 produce significantly divergent predictions, we direct IDS-Agent to search the knowledge base for characteristics of the most probable predicted attacks. If the traffic features do not match these attack characteristics, we confirm the example as an unknown attack and provide this as the final output. 902

903 A.6 THE INFLUENCE OF HYPERPARAMETERS 904

905 To assess the influence of different values of λ_1 and λ_2 in Eq. 1, we conducted experiments by 906 varying these parameters and measuring the impact on retrieval effectiveness and overall classification 907 performance. Table 9 summarizes the results of our experiments. The experimental results indicate that both recency and content similarity are crucial for effective LTM retrieval. A balanced approach, 908 where λ_1 and λ_2 are equal, provides the best performance, suggesting that the agent benefits from 909 considering both embedding similarity and recency. 910

Table 9: Performance metrics for different values of λ_1 and λ_2 .

3					
4	λ_1	λ_2	Accuracy (%)	Precision (%)	Recall (%)
5	0.1	0.9	97.2	97.2	96.5
	0.5	0.5	98.0	98.2	97.2
	0.9	0.1	97.3	97.1	96.1

918 A.7 EXCUTION TIME OF IDS-AGENT

927

In this section, we evaluate the execution time of the proposed IDS-Agent and compare it with
 the in-context-learning-based GPT-4 approach. We conducted the execution time experiments with
 the Intel Core i7 CPU of 3.8GHz. The operating system is MacOS 14.6. As shown in Table 10, the
 IDS-Agent balances performance and efficiency, averaging 8.65 seconds per instance. We use
 GPT-40 API as the core LLM of IDS-Agent. The additional time compared to the GPT-4 method
 is due to the knowledge retrieval and aggregation process, but it remains well within acceptable limits
 for real-time applications.

Table 10: Execution time comparison between different methods.

Table 10: Zheetalon time company		
Method	GPT-4	TDS-laent
		100 ngene
Average Time per Instance (s)	3.36	8.65





Figure 8: The confidence distributions of difference classifiers on the in-distribution dataset and out-of-distribution dataset.