
LARGE LANGUAGE MODELS ARE UNRELIABLE FOR CYBER THREAT INTELLIGENCE

Emanuele Mezzi
e.mezzi@vu.nl

Fabio Massacci
fabio.massacci@ieee.org

Katja Tuma
k.tuma@tue.nl

ABSTRACT

Several recent works have argued that Large Language Models (LLMs) can be used to tame the data deluge in the cybersecurity field, by improving the automation of Cyber Threat Intelligence (CTI) tasks. This work presents an evaluation methodology that other than allowing to test LLMs on CTI tasks when using zero-shot learning, few-shot learning and fine-tuning, also allows to quantify their consistency and their confidence level. We run experiments with three state-of-the-art LLMs and a dataset of 350 threat intelligence reports and present new evidence of potential security risks in relying on LLMs for CTI. We show how LLMs cannot guarantee sufficient performance on real-size reports while also being inconsistent and overconfident. Few-shot learning and fine-tuning only partially improve the results, thus posing doubts about the possibility of using LLMs for CTI scenarios, where labelled datasets are lacking and where confidence is a fundamental factor.

Keywords Large Language Models · Cyber Threat Intelligence · Unreliability · Consistency quantification · Calibration

1 Introduction

The number of vulnerabilities is becoming overwhelming: De Smale et al. [1] report that companies have reduced the raw intake of vulnerability information by 95%. In the quest to only do the work that really matters [2], Cyber threat intelligence (CTI) seems to be the new coping strategy. Unfortunately, despite standardization efforts such as STIX [3], TAXII [4], and MISP [5], CTI still requires humans to manage the massive amount of natural language information [6].

Large Language Models (LLMs) seem to be the solution to tame the CTI data deluge [7, 8] compared to pre-trained language models (PLM) such as BERT, used to automate the identification of Advanced Persistent Threat (APT) attack events [9, 10] and the extraction of knowledge graph (KG) [11]. Recent papers report high accuracy levels of LLMs processing CTI. Patsakis et al. [12] report 89% accuracy for extracting indicators of attack compromise and Hu et al. [13] used few-shot learning and fine-tuning in both named entity recognition and MITRE’s Tactics, Techniques and Procedures classification, with precision of 88% and 97%, respectively. Prompt engineering [14] and model benchmarking for CTI tasks [15] also produced promising results. CTI extraction is being augmented to use the LLM as a chatbot and requests to furnish information regarding a specific subject [16]. LLMs are known to hallucinate [17] but their CTI applications seem immune.

Unfortunately, *such promises are not based on real CTI reports*. Table 1 shows the lengths of inputs employed to perform CTI extraction, and compares them to the CISA’s emergency directive on SolarWinds. Previous claims analyzed sentences instead of reports or at best small paragraphs, shorter than the abstract of this paper. Hence our first question:

RQ1: How do LLMs perform extraction on real CTI reports?

We use the open access datasets of reports by Di Tizio et al. [2] which includes 350 reports on all Advanced Persistent Threats (APT)s reported by MITRE up to 2020 (and therefore should be known to the LLMs). Real, raw reports are two orders of magnitude larger than what is tested on the considered papers. The last column of Table 1 leaks the answer to our evaluation: LLMs are not great. Few-shot learning, prompt engineering, and fine-tuning do not significantly improve the results.

A different dimension of analysis is the consistency in the presence of repeated questions which is critical if the LLM is to be used as a chatbot. LLMs’ output generation is potentially not deterministic [23, 24] and the possibility of receiving

Table 1: LLMs extracting CTI info: reality vs papers

LLMs results look great because they use as input texts shorter than the abstract of this paper. Some work, such as Liu et al. [18] (average input 163 words), misuse the LLM to recognize the title of the report and not the actual entities of interest (e.g. APTs, CVEs). Real-size reports such as the Emergency Directive 21-01 (SolarWinds) [19], are two orders of magnitude larger. On the real 350 reports from MITRE’s APTs from [2], on average 3009 words, LLMs are not doing well . . .

Dataset	Input type	Words	Precision
Wang et al. [20]	Sentence	20	0.89
Wang et al. [21]	Sentence	18	0.83
Fieblinger et al. [22]	Sentence	54	miss.
Hu et al. [13]	Paragraph	106	0.88
<i>The abstract of this paper</i>	Paragraph	174	-
<i>Emergency Directive 21-01 (SolarWinds)</i> [19]	Report	1764	-
Our evaluation	Report	3009	0.76

different results when extracting information from the same CTI report, poses severe risks in CTI such as for patch management.

RQ2: How to evaluate the CTI consistency of LLMs?

Finally, even if the LLM may allucinate overconfidently [25] we would like to have an estimate of its uncertainty which is normally required when reporting risks [26]. This analysis is also not reported in related work [24, 13, 22, 27].

RQ3: Are LLMs over(under)confident when making predictions in CTI?

To address these challenges we: 1) design and deploy a novel evaluation pipeline to test the effectiveness, consistency and confidence calibration of LLMs for pre-attack CTI practices by extracting information from CTI reports and generating information regarding APTs, 2) run a validation experiment leveraging an existing open-source dataset consisting of 350 threat intelligence reports structured in STIX standard [2] as ground truth where we evaluate OpenAI, Google, and Mistral LLMs, 3) presents the result of the analysis on (i) the ineffectiveness of few-shot learning and fine-tuning for CTI, (ii) the inconsistency of LLMs, when used for information generation, and (iii) the low LLM confidence calibration in the extraction and generation of CTI information.

1.1 Threat Model and Non-Goals

A current focus of LLMs security research [28] is to implement adversarial attacks on LLMs such as prompt injection [29] and data poisoning [30] whose goal is to alter LLMs’ behaviour by deceiving them into making incorrect predictions. The corresponding mitigation measures to prevent adversarial attacks [31], include adding guardrails that can control LLMs output avoiding harmful outputs [32].

While interesting, we consider it a *non-goal* of our study, because there are already big security problems without calling attackers into play. We show that (uncompromised) LLMs may jeopardize the security of organizations using them to summarize CTI due to lack of consistency [25] and calibration [24]. The intuitive cause is that real reports contain additional information besides the entities which the LLM must retrieve. While a report concerns a particular attack scenario which involved a specific APT, the vulnerabilities they exploited, and the attack vectors they used in that particular scenario, it often contains information about other APTs that might have used the same vectors in other occasions or other attack vectors used by the same APT in other scenarios. The irrelevant information in the report is easily confused as relevant and thus raises the number of False Positives (FPs) and False Negatives (FN).

2 Examples of Unreliable LLMs in CTI

Longer reports, worse output. The analyst prompts the LLM with the following instructions: *Given the following CTI report, extract the name of the APT, the starting date of the campaign, the CVE of the vulnerabilities exploited and the attack vector employed.* However, the final result can be compromised by the ambiguity of natural language CTI reports, shown in Figure 1, that makes automated extraction of CTI challenging. The correct (human) interpretation of

this information is that spear-phishing links are the attack vector (as shown in the STIX format), but if taken out of context, the CTI report can be contradictory. This ambiguity, which can be reduced on small-size reports, forces to evaluate LLMs on real-size threat reports, to avoid the capacity of the models being overestimated due to the evaluation on small reports that do not represent the complexity of real CTI documents. Figure 2 shows the effects of report length

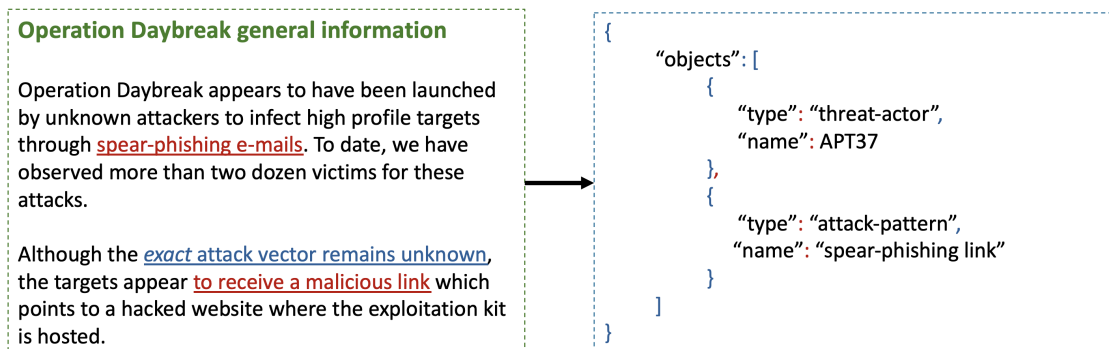


Figure 1: Excerpt of a CTI report (top) and STIX format (bottom). CTI reports can deceive tools by conveying contradictory information: **spear-phishing link** or **unknown attack vector**? The text fragment shown in the picture is extracted from a report about a campaign by APT37.

on the LLM performance. The user prompts the LLM first by giving in input a paragraph of a threat report. The result reported by the LLM is perfect, as all the requested entities are extracted. The user then asks to repeat the same task with a complete report, and the result is disappointing, as the APT and starting date of the attack campaign should be one, thus raising the FP.

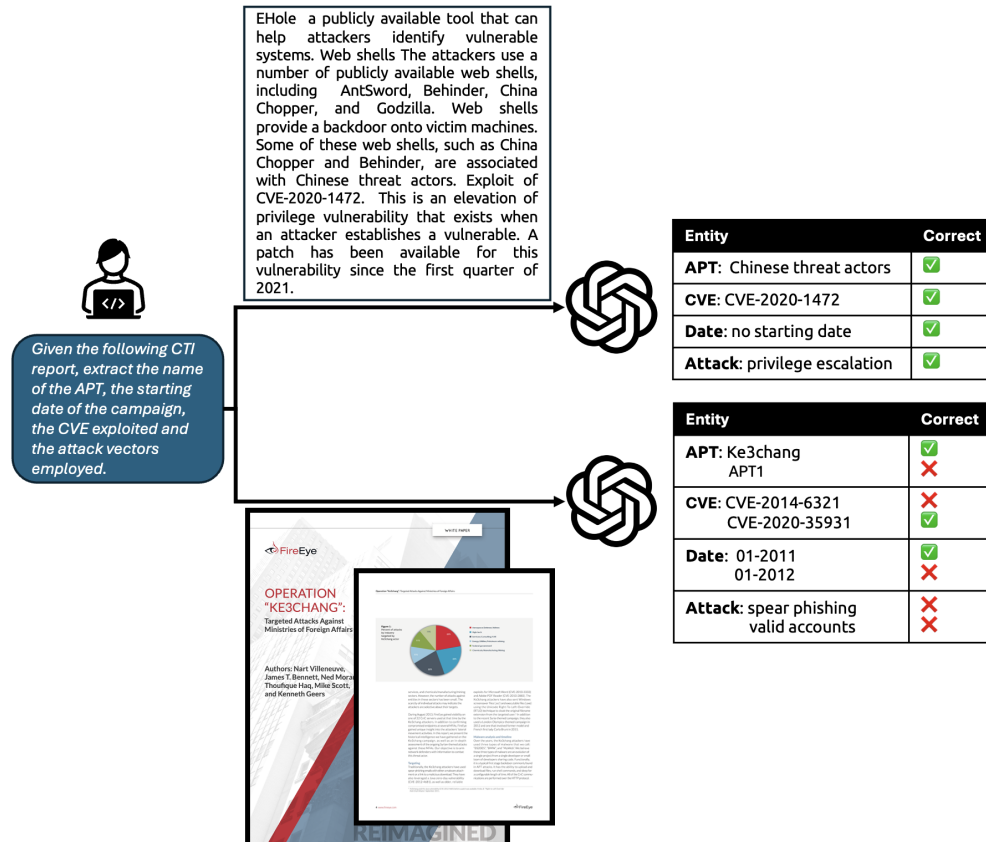
Ask twice and patch two different CVEs. It is possible that when the user prompts the LLM twice with the same instructions the information extracted differs between the two iterations [24, 33], creating uncertainty in the patching steps. If we focus on the information regarding the vulnerabilities, this lack of determinism brings uncertainty concerning the CVE that is necessary to patch, thus delaying the fixing process. However, the delay in patching the right vulnerabilities will result in a higher chance of being attacked [34]. A key information is assessing the type of threat actor (e.g. nation-state actor or criminal organization), as that implies different attacker resources, and requires proportional defences [35]. Lack of consistency in this type of information can have ominous consequences.

(Un)sure about the (right)wrong CTI. When an LLM is employed it associates a number corresponding to the confidence about the predicted token(s). For example, the LLM may associate confidence of 0.90 for the *APT K3chang* and 0.20 to the *CVE CVE-2014-6321*. Model confidence is the parameter used to decide whether to accept or not the prediction when no dataset is available. Model confidence can only be trusted if the model is calibrated: the estimated probabilities are then representative of the true correctness likelihood [36]. When automated CTI pipelines blindly rely on the model confidence level, we do not know whether a high/low confidence value is justified or the model is over or underconfident. The model could return medium to high confidence levels about a specific APT extracted from a CTI report. An automated CTI pipeline does not know if the model is overconfident, and thus if the high confidence expressed reflects the true correctness likelihood. The model might be attributing an attack campaign from a report to a threat actor that probably did not implement it, *introducing false positives*. The opposite case is also possible. If the model returns low confidence regarding an extracted vulnerability, an automated CTI pipeline will likely discard the prediction. If the model is underconfident, the confidence level does not reflect the true correctness likelihood. Thus, an automated pipeline may discard the prediction regarding a vulnerability that could be correct with a higher chance, *raising the number of false negatives* and thus leading to failing to repair a vulnerability that can be attacked in the future [37].

3 Related Work

We searched on Scopus for all papers with CTI keywords and either NLP or LLM over the last five years.

NLP and LLMs in CTI. The increasing scale and complexity of cyber attacks lead to the necessity to automate CTI practices and sharing [38, 39, 40]. Thus, researchers invented formats to structure datasets containing CTI, which can be used to test new emerging methodologies for the automation of CTI analysis. Some datasets can be employed for specific tasks such as named-entity recognition (NER) in CTI, such as [41] and [21]. Other general-purpose datasets

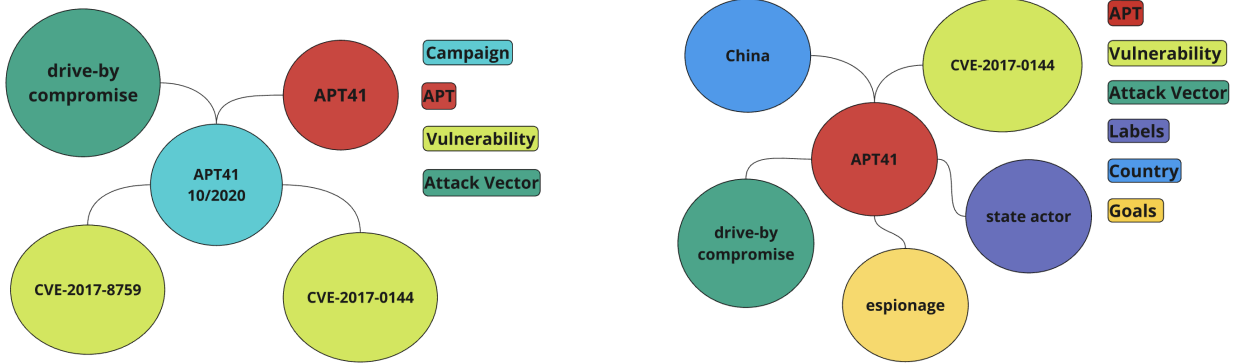


First, the analyst inputs the LLM with a paragraph of a threat report and the LLM extracts all the information contained in it. Then, the input consists of the entire report. Even though the input is better in terms of completeness, the result is worse as the number of FP and FN rises. The LLM is thus unable to distinguish the important information from the irrelevant ones.

Figure 2: More Information, worse output

can be employed for KG extraction [20, 42]. Husari et al. [43], propose an NLP-based analysis to reconstruct the chains of the attacks composed by the actions performed by the APT. In contrast, Zhang et al. [44] propose EX-Action a framework to extract threat actions from cyber threat reports. Abdi et al. [45] propose an NLP-based system to automatically highlight a CTI report with the responsible actor. Important tools offered by the world of NLP to researchers consist of PLMs and LLMs. One of the first activities to which PLMs were applied was called entity recognition in cyber threat reports. Quiao et al. [46] employs BERT-based models to annotate nine different categories of entities. Researchers also employed LLMs to help analysts respond to incoming attacks or threat campaigns, for instance automating the labelling of network intrusion detection systems rules, and the consultation of frameworks such as ATT&CK and D3FEND [47]. Researchers also tested the capacity of LLMs to automate the generation of threat reports [48] or to retrieve recovery steps from threat reports [27]. Finally, researchers employ PLMs [49, 50, 51] and LLMs [18] to extract entities and connections from threat reports, thus automating the operation of KG extraction. Fieblinger et al. [22], Hu et al. [13], and [18] tests LLMs on the task of KGs from cyber threat reports reporting promising performances. LLMs are also starting to be employed as CTI assistants to assign the correct CWE to a specific CVE [16, 52]. Current approaches that employ LLMs for CTI extraction from threat reports show promising performance. However, they test LLMs on single sentences and report paragraphs that are less complex than real-size reports. Moreover, even though LLMs are starting to be used as CTI assistants, there is still the necessity to perform proper evaluation. In this work, we address this gap by evaluating LLMs in the task of information extraction on real-size reports and in the task of CTI assistants that will build the APT profile given the given APTs' names.

Consistency quantification in LLMs. Considering the lack of determinism of LLMs [17, 24], a complete branch of research is dedicated to quantifying the consistency characterising their outputs [53, 54, 33]. It is possible to highlight two categories of consistency quantification in the realm of LLMs, where the first concerns the general characteristics of LLMs as free-form generation, while the second looks at the capacities of the LLMs regarding more limited and close-ended applications such as classification and information extraction. Examples of the first case are the



(a) KG extracted from CTI reports.

(b) KG from information generation.

Figure 3a and Figure 3b picture the graphs created when extracting information from cyber threat reports and generation information from the names of APTs. The graph resulting from information extraction is characterized by the entity *Campaign* that is composed by the name of the APT and by the starting date of the campaign, while the graph resulting from information generation is characterized by the *country* of origin of the APT, by its *goals* and by its *labels* which correspond to the type of APT.

Figure 3: KG extracted (a) and generated (b) when performing CTI tasks.

scientific research published by Kuhn et al. [55] who introduce the concept of semantic entropy, helping to measure the uncertainty regarding sentences and outputs which have the same meaning, and Lin et al. [23] who propose different metrics to evaluate the uncertainty in black-box LLMs. Examples of the second path of research can be seen in the scientific articles written by Wang et al. [56], concerned with uncertainty quantification in the realm of text-regression, and in the research carried on by Jiang et al. [57] and Kamath et al. [58], concerning uncertainty quantification in the case of PLMs employed for question answering. Current approaches in CTI measure performance with point value metrics and do not consider issues due to a lack of determinism. We address this gap and present an evaluation step that quantifies LLM consistency.

Calibration in LLMs. Lack of calibration is a characteristic of contemporary Deep Learning (DL) models [36]. However, even though new methods to improve calibration for LLMs have been suggested [25], in security this is still an unexplored problem. To the best of our knowledge, only one article explored LLMs' calibration and how to improve it for code generation [59]. Lack of calibration can have consequences in case LLMs are deployed in the absence of a labelled dataset, where the prediction confidence is used to choose whether to consider or discard their predictions. No previous work analysed calibration when LLMs are involved in CTI tasks. We address this gap by performing a calibration analysis to determine whether LLMs are (under)overconfident.

4 Overview of the Approach

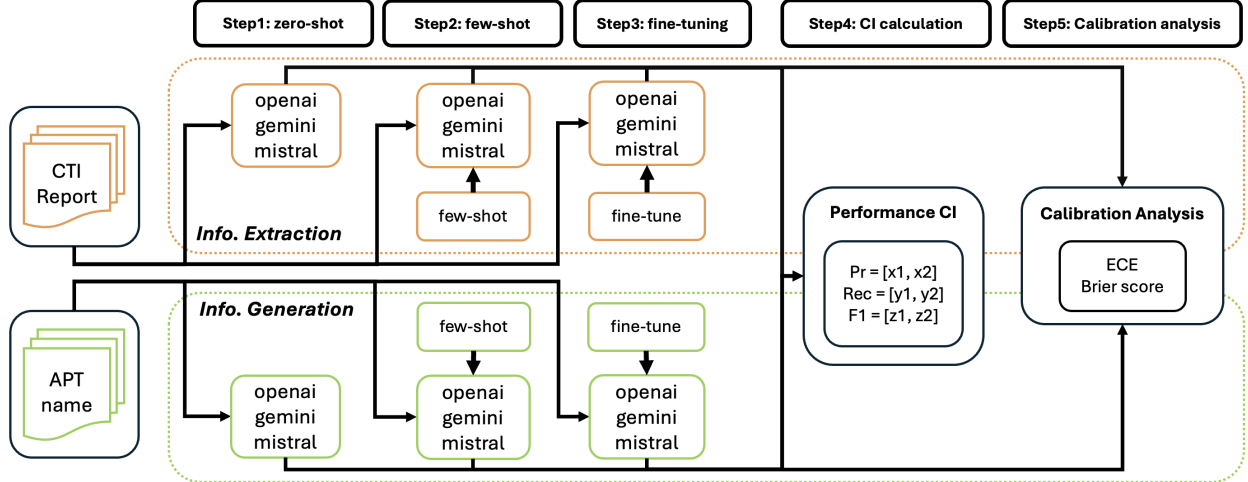
We design our five-step evaluation pipeline for two critical tasks.

Information extraction. When the LLM is given a *single* unstructured CTI report it is asked to extract entities from this report such as the APT the report is about or the attack vector it has exploited according to the report. This task is essential to transform natural language text into a structured format such as STIX. Figure 3a, shows an example of the expected outcome.

Information generation. When the LLM is asked to provide information related to an APT (e.g., the used attack vectors, the country of origin etc.) based on the knowledge embedded in its weights. Figure 3b shows an example output.

Figure 4 shows the evaluation steps. The first step consists of evaluating each LLM with zero-shot learning. The second step uses few-shot learning. To control for *prompt overfitting*¹, we evaluate the LLMs on the part of the dataset from which we do not gather any few-shot examples, mimicking the division between train and test datasets, and checking whether the model can generalize from given examples. The same approach is performed for fine-tuning: we fine-tune the LLM on the dedicated dataset section and then test it on the remaining part of the dataset. The fourth step calculates

¹With prompt overfitting, we refer to the case in which an LLM can perform well on the data from which the few-shot examples that enrich the prompt are gathered, but it is not able to use those few-shot examples to generalize its performance to data samples that do not contain those same few-shot examples [60].



Each evaluation step involves two selected CTI tasks: information extraction from threat reports and information generation from APT names, respectively represented by the orange and the green frame. The 1st step evaluates LLMs with zero-shot learning, the 2nd with few-shot learning, and the 3rd after fine-tuning the model. Performance is registered as point values. In the 4th step, we quantify LLM performance consistency by calculating performance confidence intervals (CI). During the 5th step, we compute the expected calibration error (ECE) and Brier Score (BS)

Figure 4: Illustration of the LLMs evaluation.

the confidence intervals (CIs) regarding the performance, and quantifies the oscillation in model performance, an aspect overlooked in previous work evaluating LLMs on CTI tasks [13, 22].

Precision and recall are sufficient only when models can be evaluated with a labelled dataset. When models must be used in real-world scenarios, knowledge regarding calibration is needed to know whether the predictions can be trusted [59]. The fifth step consists of deriving model confidence by analyzing the log probabilities generated by the LLMs. We can thus check whether the models are calibrated and whether few-shot learning and fine-tuning improve the calibration level, an aspect overlooked in previous work on LLMs for CTI [13, 22, 27, 61].

5 Evaluation Metrics

Traditional Metrics of Performance. For evaluating the impact of the size of the report we use the traditional metrics used in the literature (See Section 3 and Table 1): *Precision* $P = \frac{TP}{TP+FP}$ is the portion of extracted elements of a specific class which were correctly extracted; *Recall* $R = \frac{TP}{TP+FN}$ is the portion of elements of a specific class which have been extracted by the model; *F1* is the harmonic mean of precision and recall.

Confidence intervals. We measure the consistency of LLM output when prompted several times with identical inputs. To obtain measures beyond point estimates, we build confidence intervals from our observations by relying on a multi-sample method [55]. We draw sample values from the population with bootstrapping [62] with replacement, and calculate the sample mean. We repeat this process where n is the population size and k is the size of the sample drawn. Then, we create a list with the sample means, the 5th and 95th percentile (lower and upper bound of the interval).

Calibration metrics. We evaluate the LLMs calibration with expected calibration error (ECE) and Brier score (BS). The ECE [63] and BS [64], are two measures of calibration that quantify the deviation from perfect calibration. From Jiang et al. [57], given an input X and true output Y , a model output \hat{Y} , and a probability $PN(\hat{Y}|X)$ calculated over this output, a perfectly calibrated model satisfies the condition:

$$P(\hat{Y} = Y | PN(\hat{Y}|X) = p) = p, \forall p \in [0, 1] \quad (1)$$

which is that the confidence p of the model, corresponding to the calculated probability of its prediction that the output is \hat{Y} , equals the empirical fraction p of the cases where the actual output Y correctly matches the prediction \hat{Y} [59]. The best obtainable value for ECE and BS is zero, indicating perfect calibration. The higher the value is, the lower the calibration.

Table 2: Prompt engineering techniques employed taken from [14, 65]

Prompt Technique	Description	Example
Role specification	Instructing the LLM on the role.	You are a Cyber Threat Intelligence (CTI) analyst.
Step specification	Specify the steps required to accomplish a task.	Step 1 - Extract the starting date of the campaign, the Advanced Persistent Threat (APT), the CVE codes of the vulnerabilities exploited by the APT ...
Input subdivision	Split the steps into different and separated sections.	Step 2 - Return the information filling in this JSON format: "nodes": { "APT": [{"name": ""}], "attack_vector": [{"name": ""}], ...
World closing	Reducing the values an LLM can assign to an entity, indicating the possible alternatives.	The name of the attack vector can only be one of the following: drive-by compromise, supply chain compromise, spear-phishing via service, spear-phishing attachment ...
Few-shot learning	Providing a small number of labelled examples to the LLM from which it can generalize.	Examples to understand which attack vector the APT used. - ... employed legitimate user credentials to access its targets: valid accounts. - ... has been linked to a watering hole attack: drive-by compromise. ...

6 Experimental Pipeline

Prompt preparation. To ensure the quality of the prompts we rely on the practices from prompt engineering (Table 2). The techniques used are role specification, which instructs the LLM to embody the role of CTI analyst, input subdivision and step specification, which allows to specify a different task for each section of the prompt. We apply world closing by reducing, based on the dataset in [2], the possible values that entities can be assigned to by the LLM during information extraction and generation.

Zero-shot learning. We evaluate the capacity of the LLMs to perform the selected tasks without providing any examples through the prompt to help the LLM to extract and classify the information contained in the threat reports or the description of the APTs [66]. To perform information extraction, the LLM receives in input the prompt which instructs it regarding the entities to extract, and the cyber threat report from which they are to be extracted. For information generation, the LLM receives in input the name of the APT, the description of the APT, and the instructions indicating the information to recover.

Few-shot learning and fine-tuning. In the second and third steps, we repeat the tasks, by applying few-shot learning [67] and fine-tuning [68], measuring to what extent these techniques affect the performance of the LLMs for information extraction and generation. To implement few-shot learning for information extraction and generation, we respectively extract few-shot examples from threat reports and APTs' descriptions. We evaluate the LLMs on the dataset section from which few-shot examples were not extracted, to ensure that prompt overfitting is avoided. To fine-tune LLMs we randomly select a portion of the threat reports and APTs' descriptions and create the training dataset composed of the prompts and the correct answers to be generated. At the end of the fine-tuning, we evaluate the models on the test dataset, consisting of the dataset section not used for the training. We rely on supervised fine-tuning (SFT) [69] instead of more recent techniques based on reinforcement learning, such as Reasoning with REinforced Fine-Tuning (REFT) [70] as they suffer from reward hacking.

Generation of confidence intervals (CI). Since LLMs suffer from a lack of determinism [24], our evaluation measures their capacity to generate consistent results over multiple iterations. To this aim, we rely on a multi-sample approach [55, 33], by repeating the process of information extraction and information generation multiple times on the same input, registering the performance and then generating CIs. For each task, we re-prompt each LLM ten times on the same input, with $temperature=0$ and the same seed to guarantee maximum determinism, derive precision, recall, and f1-score, and then employ the bootstrapping method [62] to empirically build CIs. The bootstrapping method avoids assumptions over the distribution of data samples represented by metrics values gathered during the iterations. We re-prompt the LLM ten times to balance the necessity of checking for lack of consistency in model output with expense limits that would allow researchers to reproduce experiments.

Table 3: Dataset [2] key Indicators.

Entities	Num	Report Size	
		Mean	Max
# of reports	350		
# Campaign	350	# words	3 009 21 569
# APT	86	# tokens	4 002 27 794
# Vulnerability	123		
# Attack Vector	170		
# Country	17		
Max vulns/campaign	6		
Max attack vectors/campaign	4		

LLM calibration analysis. In the last step, we analyse whether LLMs are calibrated regarding the selected CTI tasks. We perform this analysis for zero-shot learning, few-shot learning and the fine-tuned model, to check whether few-shot learning and fine-tuning help in improving model calibration in CTI tasks [59, 71]. We implement calibration analysis by extracting the log probabilities that the LLM assigns to the tokens composing the generated response. Considering that the response is in the JSON format, and each field of the JSON file corresponds to a different entity, we isolate the tokens and the probabilities of each section of the JSON file corresponding to a specific entity, multiply the token probabilities of that JSON file section, and derive the overall confidence for each entity. Finally, we calculate the ECE [63] and BS [64], to measure the deviation from the ideal confidence level.

7 Dataset and Models

Dataset selection. For our evaluation, we use the APT dataset from Di Tizio et al. [2]. We report in Table 3 the characteristics of interest for our study.

We have selected the dataset as it is open source, it has been manually curated for correctness and presents characteristics in terms of length of reports and heterogeneity that allow us to test the LLMs on real-size CTI reports. The dataset guarantees heterogeneity by selecting 86 of the 163 APT groups on the MITRE ATT&CK. The selected APTs are the ones that launched at least one campaign during 2008 and 2020. The amount of attack campaigns and thus of CTI reports is 350. Each CTI report is associated with the following entities that can be extracted: *APT*, *Campaign*, *Vulnerability*, and *Attack vector*. The *APT* is the actor responsible for the attacks, the *Campaign* is composed of the name of the APT responsible for the attacks and its starting date, the *Vulnerability* corresponds to the *CVE* codes of the vulnerabilities exploited by the APT, and the *Attack Vectors* are the techniques employed by the threat actor.

The origin and structure of the textual sources widely vary: from cyber threat reports written by cybersecurity providers to blog posts shared by cybersecurity enthusiasts. Other than the threat reports, the dataset is equipped with information regarding the 86 APTs. Each APT is associated with the *country* of origin of the APT, the *label* of the APT which refers to whether the APT is a criminal, a nation-state actor or a spy, its *goals* such as espionage, the *vulnerabilities* they exploited over their campaigns, and the *attack vectors* employed.

Dataset division. We randomly select 70% of the dataset, thus mimicking the splitting which is typical of ML and DL model validations. In ML and DL, the splitting would have been in training, validation, and testing respectively in portions of 70%, 20%, and 10%. Since we do not train from scratch our model we assign 70% to the few-shot examples section and fine-tuning and 30% for the testing. We employ this division to avoid overfitting both for fine-tuning and for few-shot learning, if the few-shot examples are gathered directly from the dataset on which testing is conducted [60].

Models. We implement all the experiments by employing LLMs which are state-of-the-art at the time of writing, that can be fine-tuned, and for which it is available the JSON modality: *gpt4o* from OpenAI, *gemini-1.5-pro-latest* from Google, and *mistral-large-2* from Mistral, that can digest even the longest threat report given their large context windows (128k, 2M, and 128k tokens).

8 RQ1. Performance of LLMs in CTI

Table 4 and Table 7 (in the Appendix) show the results related to the first research question for entity extraction and entity generation respectively. We test LLMs with zero-shot learning, few-shot learning, and after fine-tuning.

Information extraction. Focusing on information extraction, Table 4 shows that zero-shot learning does not bring positive performance. Focusing on the recall in the best case is equal to 0.90, when *gemini* and *mistral* retrieve

Table 4: Results for information extraction.

	Model	zero-shot			few-shot			fine-tuning		
		P	R	F1	P	R	F1	P	R	F1
campaign	gpt4o	0.72	0.72	0.72	0.72	0.72	0.72	0.58	0.58	0.58
	gemini	0.77	0.77	0.77	0.73	0.73	0.73	0.61	0.61	0.61
	mistral	0.74	0.74	0.74	0.69	0.69	0.69	0.58	0.58	0.58
APT	gpt4o	0.87	0.87	0.87	0.84	0.84	0.84	0.68	0.68	0.68
	gemini	0.89	0.89	0.89	0.82	0.82	0.82	0.80	0.80	0.80
	mistral	0.89	0.89	0.89	0.82	0.82	0.82	0.68	0.68	0.68
CVE	gpt4o	0.67	0.87	0.76	0.74	0.92	0.82	0.71	0.69	0.70
	gemini	0.69	0.90	0.78	0.75	0.89	0.81	0.81	0.63	0.71
	mistral	0.72	0.90	0.80	0.79	0.91	0.85	0.71	0.69	0.70
attack vector	gpt4o	0.53	0.75	0.62	0.44	0.77	0.56	0.69	0.65	0.67
	gemini	0.68	0.74	0.71	0.71	0.78	0.74	0.89	0.84	0.87
	mistral	0.67	0.83	0.74	0.67	0.85	0.75	0.69	0.65	0.67

CVEs, meaning that 10% of the vulnerabilities will be overlooked. Recall can be as low as 0.72 when *gpt4o* is employed to retrieve *campaigns* entities, thus overlooking 28% of the campaigns. Even more worrying is the lack of benefits of few-shot learning and fine-tuning. When applying few-shot learning the performance can decrease below the performance obtained with zero-shot learning. The maximum decrease is obtained when considering the *APT* specifically for the models *gemini* and *mistral* for which the performance decreases by 7.87% from 0.89 to 0.82. The decrease in performance means that employing LLMs for CTI tasks will cause overlooking *APT* names, *CVEs*, and *campaign* and thus to an error in comprehension of the scenario and in the attack attribution. The same and even worsened negative trend is evident when analyzing the performance model after fine-tuning. The table shows that the recall decreases as low as 0.58 from 0.72 when *gpt4o* is applied to retrieve *campaign*, leading to overlooking the 42% of the *campaign* entities. However, the worst decrease is obtained when *gpt4o* and *mistral* are used to extract *APT* names with precision and recall decreasing from 0.87 to 0.68, with a decrease of 21.84%.

Information generation. Table 7 (in the Appendix) shows the results for information generation. Generally, the recall is considerably low, with the lowest level registered when the LLM is used to generate the type of *APT*, with 0.02 reached by *gemini* and *mistral*. The effect of few-shot learning and fine-tuning is limited and can also be detrimental, with the only exception represented by the *goals* entity. For few-shot learning, the maximum decrease in performance is obtained when *mistral* is used to generate the country of the threat actor, starting with precision and recall of 0.78 and ending with precision and recall of 0.64, and thus 36% of *APTs*' countries are wrongly predicted. Low performance is obtained when the LLM is fine-tuned. The worst case is represented by the application of *gpt4o* and *mistral* to generate *CVEs*, when the registered recall is 0.00, meaning that after fine-tuning the LLM is completely unable to generate and highlight the *CVEs* exploited by threat actors.

9 RQ2. Consistency of LLMs output

Here we analyze the results concerning the consistency quantification for LLMs generated output. The larger the CI computed the lower the consistency. Thus, we attribute an LLM perfect determinism or perfect consistency when the width of the CI is zero, indicating that the LLM returns the same output given the same input over multiple iterations.

Table 5 and Table 8 (in the Appendix) show the performance CI for the two tasks. For both tasks and most entities the LLMs are not consistent. Between the two tasks, information generation shows a greater lack of consistency. For information extraction in the worst case, the difference between the lower and upper bound of the CI is 0.02. The highest width can only be found between fine-tuned models, for instance when *gemini* is used to retrieve *campaign* with a precision and recall CI of [0.58, 0.60] with a difference between the lower and upper bound of 3.39% or when *mistral* is used to retrieve *CVE* codes, with precision and recall of [0.72, 0.74] and thus a difference between lower and upper bound of 2.74%.

For the task of information generation, the maximum difference between the lower and upper bound is 0.06, when few-shot learning is applied to *gemini* to generate *CVE* codes with a recall CI of [0.19, 0.25] and thus a 27.27% percentage difference between the lower and upper bound.

Table 5: Performance CI calculated with the bootstrapping techniques, for information extraction. LLMs lack complete determinism.

	Models	few-shot			fine-tuning		
		P	R	F1	P	R	F1
campaign	gpt4o	[0.69, 0.70]	[0.69, 0.70]	[0.69, 0.70]	[0.55, 0.56]	[0.55, 0.56]	[0.55, 0.56]
	gemini	[0.73, 0.74]	[0.73, 0.74]	[0.73, 0.74]	[0.58, 0.60]	[0.58, 0.60]	[0.58, 0.60]
	mistral	[0.67, 0.68]	[0.67, 0.68]	[0.67, 0.68]	[0.55, 0.56]	[0.55, 0.56]	[0.55, 0.56]
APT	gpt4o	[0.84, 0.84]	[0.84, 0.84]	[0.84, 0.84]	[0.66, 0.68]	[0.66, 0.68]	[0.66, 0.68]
	gemini	[0.84, 0.84]	[0.84, 0.84]	[0.84, 0.84]	[0.78, 0.79]	[0.78, 0.79]	[0.78, 0.79]
	mistral	[0.82, 0.82]	[0.82, 0.82]	[0.82, 0.82]	[0.66, 0.68]	[0.66, 0.68]	[0.66, 0.68]
CVE	gpt4o	[0.74, 0.74]	[0.89, 0.90]	[0.81, 0.82]	[0.72, 0.74]	[0.71, 0.74]	[0.72, 0.74]
	gemini	[0.74, 0.75]	[0.89, 0.89]	[0.80, 0.81]	[0.80, 0.81]	[0.62, 0.63]	[0.70, 0.71]
	mistral	[0.80, 0.81]	[0.91, 0.91]	[0.85, 0.85]	[0.72, 0.74]	[0.72, 0.74]	[0.72, 0.74]
attack vector	gpt4o	[0.42, 0.43]	[0.77, 0.77]	[0.54, 0.55]	[0.71, 0.73]	[0.67, 0.69]	[0.69, 0.71]
	gemini	[0.72, 0.73]	[0.77, 0.78]	[0.74, 0.75]	[0.90, 0.91]	[0.85, 0.86]	[0.87, 0.88]
	mistral	[0.66, 0.66]	[0.83, 0.84]	[0.73, 0.74]	[0.71, 0.73]	[0.71, 0.73]	[0.71, 0.73]

10 RQ3. Analysis of LLMs calibration level

The calibration analysis is performed by reporting the values of ECE and BS, which indicate the deviation of the model from perfect calibration. The best possible value for ECE and BS is zero, thus the higher their value the lower the calibration. Table 6 shows the results related to the calibration analysis.

Information extraction. Few-shot learning is detrimental to the model calibration, as seen by the variations in the ECE and BS. An increase in ECE and BS is seen when considering *campaign* entity in which the ECE and BS increase from 0.25 and 0.26 to 0.26 and 0.28 respectively. The same can be seen for the *attack vector* where ECE and BS increase from 0.13 and 0.46 to 0.19 and 0.49. The *CVE* entity presents the worst increase for ECE and BS, as they increase from 0.28 and 0.32 when zero-shot learning is used to 0.35 and 0.37 when employing few-shot learning.

Also fine-tuning worsens the LLM calibration. The only exception that shows improvement is the *CVE* entity, which diminishes both the ECE and BS when evaluating the fine-tuned LLM. All the other entities worsen their performance, with the maximum increase reached by the *campaign* entity, which registers an ECE and BS of 0.48 when calculated on the fine-tuned model performance. An increase of 92% in ECE and 84.62% for BS, compared to the ECE and BS registered when zero-shot learning is applied.

Information generation. The LLM employed for information generation offers a similar scenario. The only exception is the *goals* entity, in which few-shot learning and fine-tuning improve the ECE and BS. All other entities worsen their performance. An example of this pattern is the *labels* entity, which registers an ECE and BS of 0.45 and 0.44 when zero-shot learning is used, which rises to 0.57 and 0.53 when few-shot learning is used. This pattern is evident also for *attack vector* where the ECE raises from 0.47 to 0.48.

More concerning are the effects of fine-tuning on model calibration. The worst case is highlighted by the *CVE* and *attack vector* entity. For *CVE* the fine-tuned LLM registers an ECE and BS of 0.91 and 0.98 signalling a complete misalignment between performance and confidence level. For generation *attack vector*, the ECE and BS are 0.87 and 1.00 respectively rising from 0.47 and 0.43.

Table 6: ECE and BS with *gpt4o* for information generation and extraction.

	Information extraction						Information generation						
	zero-shot		few-shot		fine-tuning		zero-shot		few-shot		fine-tuning		
	ECE	BS	ECE	BS	ECE	BS	ECE	BS	ECE	BS	ECE	BS	
campaign	0.25	0.26	0.26	0.28	0.48	0.48	0.13	0.14	0.04	0.03	0.08	0.05	
APT	0.16	0.15	0.17	0.15	0.25	0.23	0.45	0.44	0.57	0.53	0.48	0.49	
CVE	0.28	0.32	0.35	0.37	0.18	0.21	0.19	0.22	0.38	0.27	0.35	0.29	
attack vector	0.13	0.46	0.19	0.49	0.27	0.58	CVE	0.15	0.29	0.13	0.22	0.91	0.98
							attack vector	0.47	0.43	0.48	0.42	0.87	1.00

11 Discussion

Regarding information extraction, we could see how for *campaign* entity the recall of each model is below 80%, meaning that the *campaign* entity is in more than 20% of the cases overlooked. Regarding the *CVE* codes we could see how the maximum recall is 0.90 (*gemini* and *mistral*) and the minimum is 0.87 (*gpt4o*), meaning that LLMs overlook at least the 10% of the vulnerabilities contained in the threat report and exploited during threat campaigns. The same phenomena concern *attack vector*, in which the maximum recall is 0.83 and the minimum is 0.74, meaning that 26% of the techniques used to exploit CVEs are overlooked. Few-shot learning and fine-tuning do not improve LLM performance and present detrimental effects after their use, as can be seen in Table 4 focusing on the maximum decrease brought by few-shot learning and fine-tuning of respectively 0.07 and 0.19 related to the *APT*.

The same pattern can be seen in information generation, when the LLM is used to recover the APT profile starting from its name and description, allowing to gather strategic and tactical information. The lowest recall and precision registered are 0.02 when the LLM is used to generate the label of the APT, meaning that the lack of knowledge regarding the type of APT (e.g., nation-state actor, criminal organization, etc.) is almost complete. Low and precision are also present when dealing with *CVE* entity as the minimum precision and recall are 0.10 and 0.06 (*gpt4o*) and the maximum precision and recall are 0.21 and 0.17 (*mistral*). With *country* entity the recall oscillates between 0.70 (*gpt4o*) and 0.78 (*mistral*), meaning that the best LLM in the 20% of the cases wrongly derives the country of an APT, leading to errors in the attribution of the cyber attacks. As for information extraction, the performance is rarely improved by the application of few-shot learning and fine-tuning, and their use can also be detrimental as shown by Table 7 (in the Appendix), where in the worst case few-shot learning decreases precision and recall of 0.14 (*labels* entity) and fine-tuning of 0.21 for precision and 0.17 for recall (*CVE* entity) leading to recall and precision of 0.00 when *mistral* generates CVEs exploited by APTs.

Our analysis highlights the (in)consistency of LLMs outputs. The repeated iteration and bootstrapping techniques are effective in building performance CIs that quantify LLM consistency. The computed CI show that in the task of information extraction from unstructured threat reports, LLMs lack complete determinism, with a maximum CI width of 0.02 (Table 5), posing an additional burden to their use in real-world environments, where repeated analysis could generate contradictory answers regarding the same threat report. Even larger lack of consistency, involves information generation, with a maximum CI width of 0.06, posing limitations to the use of LLMs as CTI chatbots and assistants.

Finally, our work shows that LLMs are not calibrated for CTI as shown by ECE and BSs. Few-shot learning and fine-tuning do not improve model calibration and can also be detrimental. The lack of calibration poses a third burden to using LLMs in real-world scenarios, where a labelled dataset is not available and thus the parameter used to accept or reject the model prediction consists of the confidence assigned to the information extracted or generated.

12 Limitations

Regarding the performance evaluation, a limitation is posed by evaluating LLMs on only one dataset, as evaluating LLMs on more intelligence reports would allow us to gain a deeper comprehension of the limitations posed by different types of CTI reports. We tried to handle this limitation by selecting a dataset characterized by great heterogeneity in terms of length and report source to reproduce a real-world scenario on which to test LLMs. Another limitation is posed by the limited number of LLMs, which we mitigated by choosing three state-of-the-art LLMs characterized by different architecture and from different providers.

Second, we consider limitations of the consistency quantification, where we highlight the cost-precision trade-off, as the higher the number of re-prompting the higher the precision and the higher the computational and economical costs. As we quantify the LLM consistency by re-prompting them ten times there is a risk of obtaining unreliable confidence intervals due to bootstrapping with a relatively small sample size. Raising the number of re-prompting would cause higher costs which would lead to the impossibility of reproducing the experiments. However, since the distribution of the samples is unknown, this is, to the best of our knowledge, the most appropriate method.

We perform the calibration analysis only with *gpt4o* as it is the only closed-source model between the ones used that allows extracting log probabilities and thus calculating the confidence level assigned by the LLM to tokens. This characteristic of closed-source models also limits the application of post-processing calibration methods such as Platt Scaling which requires access to logits [36, 72]. We chose closed-source models as they can be run directly on the cloud of proprietary companies, simplifying experiment reproducibility and reducing computational time and economic cost as open-source models need powerful and expensive hardware to be run [22].

13 Conclusion and future work

We show that LLMs are not ready for real-world CTI tasks. Regarding information extraction, performed on a dataset of real-size reports, given the low precision and recall, the performance cannot guarantee a faithful reconstruction of an attack scenario. Few-shot learning and fine-tuning does not seem to help. We observe the same pattern for information generation, used to build the profile of an APT, indicating that it would be risky to use LLMs as CTI assistants. A further worrying aspect consists in the performance oscillations measured and in the lack of calibration which further limits the trustworthiness posed in LLM predictions in a context where few evaluation datasets are available thus model confidence is a fundamental factor in choosing whether to rely on or not on model predictions.

In future work we plan to extend our experiments to other datasets, such as vulnerability databases, and to more LLMs, to generalize our results regarding the performance limitations that characterize them. Moreover, we plan to experiment with other prompting techniques such as Chain-of-Thought (CoT) [73], with other AI frameworks such as Retrieval Augmented Generation (RAG) [74] to improve information generation, and with approaches involving more than one LLM contemporarily, such as LLMs-based multi-agent systems [75]. Regarding measuring the consistency of LLMs, we plan to improve the empirical analysis by relying on the combination of different consistency quantification methods [33] gaining more insights regarding the model's consistency. At the same time, we also plan to integrate the empirical work with a formal analysis, thus complementing the empirical investigation.

References

- [1] Stephanie de Smale, Rik van Dijk, Xander Bouwman, Jeroen van der Ham, and Michel van Eeten. No one drinks from the firehose: How organizations filter and prioritize vulnerability information. In *2023 IEEE Symposium on Security and Privacy (SP)*, pages 1980–1996. IEEE, 2023.
- [2] Giorgio Di Tizio, Michele Armellini, and Fabio Massacci. Software updates strategies: A quantitative evaluation against advanced persistent threats. *IEEE Transactions on Software Engineering*, 49(3):1359–1373, 2022.
- [3] Sean Barnum. Standardizing cyber threat intelligence information with the structured threat information expression (stix). *Mitre Corporation*, 11:1–22, 2012.
- [4] Julie Connolly, Mark Davidson, and Charles Schmidt. The trusted automated exchange of indicator information (taxii). *The MITRE Corporation*, pages 1–20, 2014.
- [5] Cynthia Wagner, Alexandre Dulaunoy, Gérard Wagener, and Andras Iklody. Misp: The design and implementation of a collaborative threat intelligence sharing platform. In *Proceedings of the 2016 ACM on Workshop on Information Sharing and Collaborative Security*, pages 49–56, 2016.
- [6] Muhammad A Nainna, Julian M Bass, and Lee Speakman. Factors amplifying or inhibiting cyber threat intelligence sharing. In *European, Mediterranean, and Middle Eastern Conference on Information Systems*, pages 204–214. Springer, 2023.
- [7] Siva Sai, Utkarsh Yashvardhan, Vinay Chamola, and Biplab Sikdar. Generative ai for cyber security: Analyzing the potential of chatgpt, dall-e and other models for enhancing the security space. *IEEE Access*, 2024.
- [8] Venkata Ramana Saddi, Santhosh Kumar Gopal, Abdul Sajid Mohammed, S Dhanasekaran, and Mahaveer Singh Naruka. Examine the role of generative ai in enhancing threat intelligence and cyber security measures. In *2024 2nd International Conference on Disruptive Technologies (ICDT)*, pages 537–542. IEEE, 2024.
- [9] Ga Xiang, Chen Shi, and Yangsen Zhang. An apt event extraction method based on bert-bigru-crf for apt attack detection. *Electronics*, 12(15):3349, 2023.
- [10] JunJun Chen, Chengliang Gao, Fei Tang, Jiaxu Xing, Qianlong Xiao, Dongyang Zheng, and Jing Qiu. Automatically identifying sentences with attack behavior from cyber threat intelligence reports. In *2023 8th International Conference on Data Science in Cyberspace (DSC)*, pages 491–498. IEEE, 2023.
- [11] Gaosheng Wang, Peipei Liu, Jintao Huang, Haoyu Bin, Xi Wang, and Hongsong Zhu. Knowcti: Knowledge-based cyber threat intelligence entity and relation extraction. *Computers & Security*, 141:103824, 2024.
- [12] Constantinos Patsakis, Fran Casino, and Nikolaos Lykousas. Assessing llms in malicious code deobfuscation of real-world malware campaigns. *Expert Systems with Applications*, 256, 2024. Cited by: 0; All Open Access, Green Open Access.
- [13] Yuelin Hu, Futai Zou, Jijia Han, Xin Sun, and Yilei Wang. Llm-tikg: Threat intelligence knowledge graph construction utilizing large language model. *Computers & Security*, 145:103999, 2024.
- [14] Banghao Chen, Zhaofeng Zhang, Nicolas Langrené, and Shengxin Zhu. Unleashing the potential of prompt engineering in large language models: a comprehensive review. *arXiv preprint arXiv:2310.14735*, 2023.

- [15] Hangyuan Ji, Jian Yang, Linzheng Chai, Chaoren Wei, Liqun Yang, Yunlong Duan, Yunli Wang, Tianzhen Sun, Hongcheng Guo, Tongliang Li, et al. Sevenllm: Benchmarking, eliciting, and enhancing abilities of large language models in cyber threat intelligence. *arXiv preprint arXiv:2405.03446*, 2024.
- [16] Mohamed Amine Ferrag, Fatima Alwahedi, Ammar Battah, Bilel Cherif, Abdechakour Mechri, and Norbert Tihanyi. Generative ai and large language models for cyber security: All insights you need. *arXiv preprint arXiv:2405.12750*, 2024.
- [17] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232*, 2023.
- [18] Jiehui Liu and Jieyu Zhan. Constructing knowledge graph from cyber threat intelligence using large language model. In *2023 IEEE International Conference on Big Data (BigData)*, pages 516–521. IEEE, 2023.
- [19] CISA. Emergency directive 21-01: Mitigate solarwinds orion code compromise. Technical report, "Cybersecurity and Infrastructure Security Agency (CISA)", 2020. available at <https://www.cisa.gov/news-events/directives/ed-21-01-mitigate-solarwinds-orion-code-compromise>.
- [20] Xuren Wang, Xinpei Liu, Shengqin Ao, Ning Li, Zhengwei Jiang, Zongyi Xu, Zihan Xiong, Mengbo Xiong, and Xiaoqing Zhang. Dnrti: A large-scale dataset for named entity recognition in threat intelligence. In *2020 IEEE 19th International Conference on Trust, Security and Privacy in Computing and Communications (TrustCom)*, pages 1842–1848. IEEE, 2020.
- [21] Xuren Wang, Songheng He, Zihan Xiong, Xinxin Wei, Zhengwei Jiang, Sihao Chen, and Jun Jiang. Aptner: A specific dataset for ner missions in cyber threat intelligence field. In *2022 IEEE 25th International Conference on Computer Supported Cooperative Work in Design (CSCWD)*, pages 1233–1238. IEEE, 2022.
- [22] Romy Fieblinger, Md Tanvirul Alam, and Nidhi Rastogi. Actionable cyber threat intelligence using knowledge graphs and large language models. In *2024 IEEE European Symposium on Security and Privacy Workshops (EuroS&PW)*, pages 100–111. IEEE, 2024.
- [23] Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. Generating with confidence: Uncertainty quantification for black-box large language models. *arXiv preprint arXiv:2305.19187*, 2023.
- [24] Yifan Song, Guoyin Wang, Sujian Li, and Bill Yuchen Lin. The good, the bad, and the greedy: Evaluation of llms should not ignore non-determinism. *arXiv preprint arXiv:2407.10457*, 2024.
- [25] Haoyan Yang, Yixuan Wang, Xingyin Xu, Hanyuan Zhang, and Yirong Bian. Can we trust llms? mitigate overconfidence bias in llms through knowledge transfer. *arXiv preprint arXiv:2405.16856*, 2024.
- [26] Seth Guikema. Artificial intelligence for natural hazards risk analysis: Potential, challenges, and research needs. *Risk Analysis*, 40(6):1117–1123, 2020.
- [27] Zsolt Levente Kucsván, Marco Caselli, Andreas Peter, and Andrea Continella. Inferring recovery steps from cyber threat intelligence reports. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 14828 LNCS:330 – 349, 2024. Cited by: 0.
- [28] Erfan Shayegani, Md Abdullah Al Mamun, Yu Fu, Pedram Zaree, Yue Dong, and Nael Abu-Ghazaleh. Survey of vulnerabilities in large language models revealed by adversarial attacks. *arXiv preprint arXiv:2310.10844*, 2023.
- [29] Kai Greshake, Sahar Abdelnabi, Shailesh Mishra, Christoph Endres, Thorsten Holz, and Mario Fritz. Not what you’ve signed up for: Compromising real-world llm-integrated applications with indirect prompt injection. In *Proceedings of the 16th ACM Workshop on Artificial Intelligence and Security*, pages 79–90, 2023.
- [30] Xiaoyi Chen, Siyuan Tang, Rui Zhu, Shijun Yan, Lei Jin, Zihao Wang, Liya Su, Zhikun Zhang, Xiaofeng Wang, and Haixu Tang. The janus interface: How fine-tuning in large language models amplifies the privacy risks. *arXiv preprint arXiv:2310.15469*, 2023.
- [31] Arijit Ghosh Chowdhury, Md Mofijul Islam, Vaibhav Kumar, Faysal Hossain Shezan, Vinija Jain, and Aman Chadha. Breaking down the defenses: A comparative survey of attacks on large language models. *arXiv preprint arXiv:2403.04786*, 2024.
- [32] Traian Rebedea, Razvan Dinu, Makesh Sreedhar, Christopher Parisien, and Jonathan Cohen. Nemo guardrails: A toolkit for controllable and safe llm applications with programmable rails. *arXiv preprint arXiv:2310.10501*, 2023.
- [33] Miao Xiong, Andrea Santilli, Michael Kirchof, Adam Golinski, and Sinead Williamson. Efficient and effective uncertainty quantification for llms. In *Neurips Safe Generative AI Workshop 2024*, 2024.

- [34] Nesara Dissanayake, Mansooreh Zahedi, Asangi Jayatilaka, and Muhammad Ali Babar. Why, how and where of delays in software security patch management: An empirical investigation in the healthcare sector. *Proceedings of the ACM on Human-computer Interaction*, 6(CSCW2):1–29, 2022.
- [35] Vasileios Mavroeidis, Ryan Hohimer, Tim Casey, and Audun Jesang. Threat actor type inference and characterization within cyber threat intelligence. In *2021 13th International Conference on Cyber Conflict (CyCon)*, pages 327–352. IEEE, 2021.
- [36] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International conference on machine learning*, pages 1321–1330. PMLR, 2017.
- [37] Amit Seal Ami, Kevin Moran, Denys Poshyvanyk, and Adwait Nadkarni. " false negative-that one is going to kill you": Understanding industry perspectives of static analysis based security testing. In *2024 IEEE Symposium on Security and Privacy (SP)*, pages 3979–3997. IEEE, 2024.
- [38] Siri Bromander, Morton Swimmer, Lilly Pijnenburg Muller, Audun Jøsang, Martin Eian, Geir Skjøtskift, and Fredrik Borg. Investigating sharing of cyber threat intelligence and proposing a new data model for enabling automation in knowledge representation and exchange. *Digital Threats: Research and Practice (DTRAP)*, 3(1):1–22, 2021.
- [39] Cristoffer Leite, Jerry Den Hartog, Daniel R Dos Santos, and Elisa Costante. Automated cyber threat intelligence generation on multi-host network incidents. In *2023 IEEE International Conference on Big Data (BigData)*, pages 2999–3008. IEEE, 2023.
- [40] Amit Kumar Bairwa, Rohan Khanna, Sandeep Joshi, and Pljonkin Anton Pavlovich. Enhancing cyber threat intelligence and security automation: A comprehensive approach for effective protection. In *World Conference on Information Systems for Business Management*, pages 297–306. Springer, 2023.
- [41] Yuan Liu, Rong Shi, Yahe Chen, Xiaorui Gong, Qingli Guo, and Xiu Zhang. Apttoolner: A chinese dataset of cyber security tool for ner task. In *2023 3rd Asia-Pacific Conference on Communications Technology and Computer Science (ACCTCS)*, pages 368–373. IEEE, 2023.
- [42] Md Tanvirul Alam, Dipkamal Bhusal, Youngja Park, and Nidhi Rastogi. Looking beyond iocs: Automatically extracting attack patterns from external cti. In *Proceedings of the 26th International Symposium on Research in Attacks, Intrusions and Defenses*, pages 92–108, 2023.
- [43] Ghaith Husari, Ehab Al-Shaer, Bill Chu, and Ruhani Faiheem Rahman. Learning apt chains from cyber threat intelligence. In *Proceedings of the 6th Annual Symposium on Hot Topics in the Science of Security*, pages 1–2, 2019.
- [44] Huixia Zhang, Guowei Shen, Chun Guo, Yunhe Cui, and Chaohui Jiang. Ex-action: Automatically extracting threat actions from cyber threat intelligence report based on multimodal learning. *Security and Communication Networks*, 2021:1–12, 2021.
- [45] Hamza Abdi, Steven R Bagley, Steven Furnell, and Jamie Twycross. Automatically labeling cyber threat intelligence reports using natural language processing. In *Proceedings of the ACM Symposium on Document Engineering 2023*, pages 1–4, 2023.
- [46] Zhe Qiao, Chen Zhang, and Gang Du. Improving cybersecurity named entity recognition with large language models. In *2023 6th International Conference on Software Engineering and Computer Science (CSECS)*, pages 01–06. IEEE, 2023.
- [47] Nir Daniel, Florian Kaiser, Anton Dzega, Aviad Elyashar, and Rami Puzis. *Labeling NIDS Rules with MITRE ATT&CK Techniques Using ChatGPT*, pages 76–91. 03 2024.
- [48] Filippo Perrina, Francesco Marchiori, Mauro Conti, and Nino Vincenzo Verde. Agir: Automating cyber threat intelligence reporting with natural language generation. In *2023 IEEE International Conference on Big Data (BigData)*, pages 3053–3062. IEEE, 2023.
- [49] Hyeonseong Jo, Yongjae Lee, and Seungwon Shin. Vulcan: Automatic extraction and analysis of cyber threat intelligence from unstructured text. *Computers & Security*, 120:102763, 2022.
- [50] Mohit Sewak, Vamsi Emani, and Annam Naresh. Crush: Cybersecurity research using universal llms and semantic hypernetworks. 2023.
- [51] Francesco Marchiori, Mauro Conti, and Nino Vincenzo Verde. Stixnet: A novel and modular solution for extracting all stix objects in cti reports. In *Proceedings of the 18th International Conference on Availability, Reliability and Security*, pages 1–11, 2023.
- [52] Chris Madden. Using genai for efficient cve mitigation. presentation at EPSS Sig, 2024.

- [53] Ekaterina Fadeeva, Roman Vashurin, Akim Tsvigun, Artem Vazhentsev, Sergey Petrakov, Kirill Fedyanin, Daniil Vasilev, Elizaveta Goncharova, Alexander Panchenko, Maxim Panov, et al. Lm-polygraph: Uncertainty estimation for language models. *arXiv preprint arXiv:2311.07383*, 2023.
- [54] Fanghua Ye, Mingming Yang, Jianhui Pang, Longyue Wang, Derek F Wong, Emine Yilmaz, Shuming Shi, and Zhaopeng Tu. Benchmarking llms via uncertainty quantification. *arXiv preprint arXiv:2401.12794*, 2024.
- [55] Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. *arXiv preprint arXiv:2302.09664*, 2023.
- [56] Yuxia Wang, Daniel Beck, Timothy Baldwin, and Karin Verspoor. Uncertainty estimation and reduction of pre-trained models for text regression. *Transactions of the Association for Computational Linguistics*, 10:680–696, 2022.
- [57] Zhengbao Jiang, Jun Araki, Haibo Ding, and Graham Neubig. How can we know when language models know? on the calibration of language models for question answering. *Transactions of the Association for Computational Linguistics*, 9:962–977, 2021.
- [58] Amita Kamath, Robin Jia, and Percy Liang. Selective question answering under domain shift. *arXiv preprint arXiv:2006.09462*, 2020.
- [59] Claudio Spiess, David Gros, Kunal Suresh Pai, Michael Pradel, Md Rafiqul Islam Rabin, Amin Alipour, Susmit Jha, Prem Devanbu, and Toufique Ahmed. Calibration and correctness of language models for code. *arXiv preprint arXiv:2402.02047*, 2024.
- [60] Youngjae Cho, HeeSun Bae, Seungjae Shin, Yeo Dong Youn, Weonyoung Joo, and Il-Chul Moon. Make prompts adaptable: Bayesian modeling for vision-language prompt learning with data-dependent prior. *arXiv preprint arXiv:2401.06799*, 2024.
- [61] Norbert Tihanyi, Mohamed Amine Ferrag, Ridhi Jain, Tamas Bisztray, and Merouane Debbah. Cybermetric: A benchmark dataset based on retrieval-augmented generation for evaluating llms in cybersecurity knowledge. page 296 – 302, 2024. Cited by: 0; All Open Access, Green Open Access.
- [62] Bradley Efron. Bootstrap methods: another look at the jackknife. In *Breakthroughs in statistics: Methodology and distribution*, pages 569–593. Springer, 1992.
- [63] Mahdi Pakdaman Naeini, Gregory Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 29, 2015.
- [64] Glenn W Brier. Verification of forecasts expressed in terms of probability. *Monthly weather review*, 78(1):1–3, 1950.
- [65] OpenAI. Prompt engineering. 2024.
- [66] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213, 2022.
- [67] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Nee-lakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [68] Marius Mosbach, Tiago Pimentel, Shauli Ravfogel, Dietrich Klakow, and Yanai Elazar. Few-shot fine-tuning vs. in-context learning: A fair comparison and evaluation. *arXiv preprint arXiv:2305.16938*, 2023.
- [69] Aldo Pareja, Nikhil Shivakumar Nayak, Hao Wang, Krishnateja Killamsetty, Shivchander Sudalairaj, Wenlong Zhao, Seungwook Han, Abhishek Bhandwaldar, Guangxuan Xu, Kai Xu, et al. Unveiling the secret recipe: A guide for supervised fine-tuning small llms. *arXiv preprint arXiv:2412.13337*, 2024.
- [70] Trung Quoc Luong, Xinbo Zhang, Zhanming Jie, Peng Sun, Xiaoran Jin, and Hang Li. Reft: Reasoning with reinforced fine-tuning. *arXiv preprint arXiv:2401.08967*, 2024.
- [71] Sanyam Kapoor, Nate Gruver, Manley Roberts, Arka Pal, Samuel Dooley, Micah Goldblum, and Andrew Wilson. Calibration-tuning: Teaching large language models to know what they don’t know. In *Proceedings of the 1st Workshop on Uncertainty-Aware NLP (UncertainNLP 2024)*, pages 1–14, 2024.
- [72] John Platt et al. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. *Advances in large margin classifiers*, 10(3):61–74, 1999.
- [73] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

- [74] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474, 2020.
- [75] Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xiangliang Zhang. Large language model based multi-agents: A survey of progress and challenges. *arXiv preprint arXiv:2402.01680*, 2024.

A Additional tables

Table 7: Results for information generation. Few-shot learning and fine-tuning have limited performance improvement and their effect can be detrimental, as we can see for the entities *country* and *labels*. When a fine-tuned model the performance considerably decreases, even under the performance obtained with zero-shot learning, as it happens with *CVE*, and *labels*.

	Model	zero-shot			few-shot			fine-tuning		
		P	R	F1	P	R	F1	P	R	F1
goals	gpt4o	0.85	0.85	0.85	0.96	0.96	0.96	0.96	0.96	0.96
	gemini	0.72	0.72	0.72	0.83	0.83	0.83	0.84	0.84	0.84
	mistral	0.77	0.77	0.77	0.92	0.92	0.92	0.96	0.96	0.96
labels	gpt4o	0.50	0.50	0.50	0.44	0.44	0.44	0.44	0.44	0.44
	gemini	0.02	0.02	0.02	0.54	0.54	0.54	0.40	0.40	0.40
	mistral	0.02	0.02	0.02	0.36	0.36	0.36	0.44	0.44	0.44
country	gpt4o	0.70	0.70	0.70	0.56	0.56	0.56	0.60	0.60	0.60
	gemini	0.73	0.73	0.73	0.81	0.71	0.76	0.68	0.68	0.68
	mistral	0.78	0.78	0.78	0.64	0.64	0.64	0.60	0.60	0.60
CVE	gpt4o	0.10	0.06	0.08	0.08	0.07	0.07	0.00	0.00	0.00
	gemini	0.13	0.13	0.13	0.23	0.19	0.21	0.17	0.36	0.23
	mistral	0.21	0.17	0.19	0.24	0.24	0.24	0.00	0.00	0.00
attack vector	gpt4o	0.37	0.52	0.43	0.37	0.51	0.43	1.00	0.09	0.16
	gemini	0.24	0.54	0.33	0.27	0.56	0.36	0.52	0.84	0.64
	mistral	0.22	0.58	0.32	0.20	0.75	0.32	1.00	0.09	0.16

Table 8: Performance CI for the task of information generation.

	Models	few-shot			fine-tuning		
		P	R	F1	P	R	F1
goals	gpt4o	[0.96, 0.96]	[0.96, 0.96]	[0.96, 0.96]	[0.96, 0.96]	[0.96, 0.96]	[0.96, 0.96]
	gemini	[0.87, 0.90]	[0.87, 0.90]	[0.87, 0.90]	[0.84, 0.84]	[0.84, 0.84]	[0.84, 0.84]
	mistral	[0.92, 0.92]	[0.92, 0.92]	[0.92, 0.92]	[0.96, 0.96]	[0.96, 0.96]	[0.96, 0.96]
labels	gpt4o	[0.44, 0.44]	[0.44, 0.44]	[0.44, 0.44]	[0.44, 0.44]	[0.44, 0.44]	[0.44, 0.44]
	gemini	[0.54, 0.56]	[0.54, 0.56]	[0.54, 0.56]	[0.39, 0.40]	[0.39, 0.40]	[0.39, 0.40]
	mistral	[0.36, 0.36]	[0.36, 0.36]	[0.36, 0.36]	[0.44, 0.44]	[0.44, 0.44]	[0.44, 0.44]
country	gpt4o	[0.57, 0.59]	[0.57, 0.59]	[0.57, 0.59]	[0.60, 0.61]	[0.60, 0.61]	[0.60, 0.61]
	gemini	[0.82, 0.86]	[0.71, 0.74]	[0.76, 0.79]	[0.68, 0.68]	[0.68, 0.68]	[0.68, 0.68]
	mistral	[0.64, 0.64]	[0.64, 0.64]	[0.64, 0.64]	[0.60, 0.61]	[0.60, 0.61]	[0.60, 0.61]
CVE	gpt4o	[0.08, 0.09]	[0.07, 0.08]	[0.07, 0.08]	[0.00, 0.00]	[0.00, 0.00]	[0.00, 0.00]
	gemini	[0.21, 0.26]	[0.19, 0.25]	[0.20, 0.25]	[0.16, 0.19]	[0.32, 0.37]	[0.21, 0.24]
	mistral	[0.23, 0.24]	[0.24, 0.24]	[0.23, 0.24]	[0.00, 0.00]	[0.00, 0.00]	[0.00, 0.00]
attack vector	gpt4o	[0.37, 0.37]	[0.51, 0.52]	[0.43, 0.44]	[1.00, 1.00]	[0.09, 0.10]	[0.16, 0.18]
	gemini	[0.24, 0.27]	[0.55, 0.58]	[0.34, 0.36]	[0.49, 0.50]	[0.82, 0.84]	[0.61, 0.63]
	mistral	[0.20, 0.20]	[0.74, 0.75]	[0.31, 0.31]	[1.00, 1.00]	[0.00, 0.00]	[0.00, 0.00]