

CTINEXUS: Leveraging Optimized LLM In-Context Learning for Constructing Cybersecurity Knowledge Graphs Under Data Scarcity

Yutong Cheng Osama Bajaber Saimon Amanuel Tsegai Dawn Song Peng Gao
Virginia Tech Virginia Tech Virginia Tech UC Berkeley Virginia Tech
yutongcheng@vt.edu obajaber@vt.edu saimon.tsegai@vt.edu dawnsong@berkeley.edu penggao@vt.edu

Abstract—Textual descriptions in cyber threat intelligence (CTI) reports, such as security articles and news, are rich sources of knowledge about cyber threats, crucial for organizations to stay informed about the rapidly evolving threat landscape. However, current CTI extraction methods lack flexibility and generalizability, often resulting in inaccurate and incomplete knowledge extraction. Syntax parsing relies on fixed rules and dictionaries, while model fine-tuning requires large annotated datasets, making both paradigms challenging to adapt to new threats and ontologies. To bridge the gap, we propose CTINexus, a novel framework leveraging optimized in-context learning (ICL) of large language models (LLMs) for data-efficient CTI knowledge extraction and high-quality cybersecurity knowledge graph (CSKG) construction. Unlike existing methods, CTINexus requires neither extensive data nor parameter tuning and can adapt to various ontologies with minimal annotated examples. This is achieved through: (1) a carefully designed automatic prompt construction strategy with optimal demonstration retrieval for extracting a wide range of cybersecurity entities and relations; (2) a hierarchical entity alignment technique that canonicalizes the extracted knowledge and removes redundancy; (3) an ICL-enhanced long-distance relation prediction technique to further complete the CSKG with missing links. Our extensive evaluations using 150 real-world CTI reports collected from 10 platforms demonstrate that CTINexus significantly outperforms existing methods in constructing accurate and complete CSKGs, highlighting its potential to transform CTI analysis with an efficient and adaptable solution for the dynamic threat landscape.

Index Terms—Threat Intelligence, Large Language Model, In-Context Learning, Cybersecurity Knowledge Graph

1. Introduction

Modern cyberattacks are becoming increasingly complex and rapidly evolving. Many public and commercial organizations extensively record and share cyber threat intelligence (CTI) on their platforms to combat evolving threats. According to Gartner, CTI is defined as “evidence-based knowledge, including context, mechanisms, indicators, implications and actionable advice, about an existing or emerging menace or hazard to assets that can be used to inform decisions regarding the subject’s response to that menace or hazard” [60]. Such knowledge is crucial for organizations to monitor the rapidly evolving threat landscape, promptly detect early signs of an attack, and

effectively contain the attack with proper measures. Given its importance, CTI has been increasingly collected and exchanged across organizations, often in the form of Indicators of Compromise (IOC) [56]. IOCs are forensic artifacts of an intrusion such as virus signatures, IPs/domains of botnets, MD5 hashes of attack files, etc. However, recent studies [56], [77] showed that knowledge offered by IOCs is rather limited, which covers only a limited set of knowledge and has a short lifespan.

Recognizing the limitations of IOCs, recent research has shifted towards automatically extracting richer knowledge from textual threat descriptions in CTI reports (i.e., CTI text). These reports, such as security blog articles [22], [6] and news [10], [20], are produced by security researchers and practitioners and published on websites, summarizing threat behaviors in natural language. Besides IOCs, these reports contain various other cybersecurity entities, such as malware, vulnerabilities, and attack techniques, as well as their relationships, illustrating their interactions and dependencies. This knowledge is crucial for building a comprehensive cyber threat profile.

Several approaches have been proposed for automatically extracting security knowledge from CTI text and constructing a *cybersecurity knowledge graph (CSKG)*. Syntax parsing-based approaches [43], [39], [56] leverage fixed dependency rules and hand-crafted dictionaries to parse the grammatical structure of sentences and extract key subject-verb-object triplets. Fine-tuning-based approaches [71], [55], [30] leverage pre-trained transformer models and fine-tune them on labeled CTI text datasets to identify semantic roles and extract entities and relations. However, all these existing methods suffer from several *key drawbacks* (see Section 2.2 for details), particularly when facing the evolving threat landscape: (1) **Lack of flexibility and generalizability**: Many of these methods are tailored to specific cybersecurity ontologies, focusing on a fixed set of entities and relation types. They are difficult to generalize to new ontologies and emerging threats and terminologies. Fixed rules have limited flexibility to adapt to new patterns and require manual creation and maintenance. Model fine-tuning, however, requires a large amount of labeled CTI data. Such data is scarce in security, especially for new threats that lack annotations. (2) **Information inaccuracy and incompleteness**: Due to the peculiarities of the security context and the lack of deep analysis, these methods often generate low-quality CSKGs that are incomplete, inaccurate, and disconnected. Fig. 1 shows example CSKGs generated by two representative methods for a real-world CTI report. We can observe

several issues, including incomplete entities, misidentified entity boundaries, misaligned entities, missing links, etc. These low-quality CSKGs limit the ability to obtain a comprehensive threat profile, hindering the effective use of CTI to enhance other defensive measures.

These limitations highlight the need for a paradigm shift in CTI knowledge extraction that enables accurate knowledge capture in data-limited environments while adapting to evolving threats. Recent advancements in LLMs have demonstrated strong capabilities in various natural language tasks [33], shifting the focus from fine-tuning to *in-context learning (ICL)*, which requires minimal annotated data and no parameter updates. However, ICL strategies vary in performance, from state-of-the-art to suboptimal [57]. To address this, we conducted thorough experiments to identify optimal ICL settings for CSKG construction. With the optimized ICL strategy, LLMs can effectively learn from a few examples and adapt to new tasks with stability and high performance without requiring model weight updates.

Contributions. We present CTINEXUS, an LLM-powered framework for automated CTI knowledge extraction and CSKG construction from CTI reports. Unlike existing methods limited by generalizability and data demands, CTINEXUS introduces an *optimized-ICL-based pipeline for data-efficient inference*, enabling precise extraction of diverse cybersecurity entities and relations while adapting to various ontologies. By leveraging an optimal ICL strategy, CTINEXUS generalizes from selected demonstrations to perform versatile CSKG construction tasks. In addition to cybersecurity triplet extraction, CTINEXUS also refines extracted knowledge to enhance the canonicalization and completeness of the resulting knowledge graph. As shown in Fig. 1, the CSKG constructed by CTINEXUS has significantly higher quality compared to existing approaches.

CTINEXUS leverages the ICL paradigm of LLMs to directly extract entity-relation triplets (i.e., (head entity, relation, tail entity)) by analogizing similar demonstration examples in the prompt construction, eliminating the need for large amounts of training data or extensive model tuning. Unlike multi-round dialogue approaches, CTINEXUS performs end-to-end extraction of triplets in a single step (see Fig. 3), significantly reducing inference token costs. To ensure the high quality of the extracted knowledge, CTINEXUS employs a carefully designed *prompt template* and an optimal *demonstration retrieval strategy* for automatic prompt construction. This prompt construction also incorporates the defined ontology for the task domain. Different ontologies can be easily swapped in, and with just a few demonstration examples, CTINEXUS can automatically bootstrap and adapt to new threats and tasks.

To canonicalize the knowledge and remove redundancy in entities, we designed a *hierarchical entity alignment* technique, which consists of two phases. In coarse-grained entity grouping, CTINEXUS assigns entity types to each entity in the extracted triplets using LLM’s ICL and groups entities within the same type. This ensures preliminary categorization and prevents the merging of textually similar entities that belong to different types. In fine-grained entity merging, CTINEXUS calculates the semantic similarity among the grouped entities and merges those with high similarity. With this hierarchical approach,

CTINEXUS avoids the high costs of querying LLMs for each entity pair’s similarity.

To further complete the CSKG with implicit relations for distant entities, we designed a *long-distance relation prediction* technique. First, entities with the highest degree centrality in a subgraph are selected as the central nodes of that subgraph. Then, CTINEXUS performs ICL-enhanced implicit relation prediction on these central nodes, conditioned on the input context, to infer connections among the disjoint subgraphs.

Evaluation. We conducted comprehensive evaluations using 150 CTI reports from 10 well-recognized CTI sharing platforms [3], [4], [6], [9], [20], [21], [10], [22], [23], [25]. CTINEXUS achieved F1 scores of 87.65% in cybersecurity triplet extraction, 89.94% in coarse-grained entity grouping, 99.80% in fine-grained entity merging, and 90.99% in long-distance relation prediction. Qualitative analysis showed that CTINEXUS constructs more comprehensive and interconnected CSKGs compared to TTPDrill [43], EXTRACTOR [71], and LADDER [30]. Quantitatively, CTINEXUS outperforms EXTRACTOR by 25.36% in F1 score for cybersecurity triplet extraction and LADDER by 19% in cybersecurity entity extraction. We also explored various prompting strategies and four backbone models (closed-source models: GPT-3.5, GPT-4; open-source models: Llama3 and QWen2.5) to identify the optimal ICL paradigm for CTI knowledge extraction, providing valuable insights for future research.

2. Background and Motivating Example

2.1. Cyber Threat Intelligence

Although crowd-sourced CTI reports provide valuable information, their unstructured format significantly hinders their effectiveness. As the number and complexity of cyberattacks increase, the textual CTI descriptions have also expanded, creating an urgent need for automated information extraction from CTI [68]. The extracted knowledge can be used to construct cybersecurity knowledge graphs (CSKGs), where nodes represent entities and edges represent relations. Compared to unstructured CTI text, CSKGs provide a holistic profile for cyber threats, offer better visualization, and are more amenable to integration into downstream applications. The construction of a CSKG typically follows an ontology, which specifies the entity types and their allowed relations. However, despite efforts to create multiple security ontologies [44], [69], [75] covering different aspects of threats, keeping up with the rapidly evolving threat landscape remains challenging as new threats, techniques, and tools constantly emerge. It is nearly impossible to develop a universal ontology that encompasses all current and future threats. This underscores the need for CTI knowledge extraction approaches that can flexibly adapt to different ontologies and emerging threats with minimal transition effort.

2.2. Limitations of Existing Approaches

Existing CTI knowledge extraction approaches face several fundamental challenges in adapting to the rapidly evolving threat landscape. Existing approaches follow two

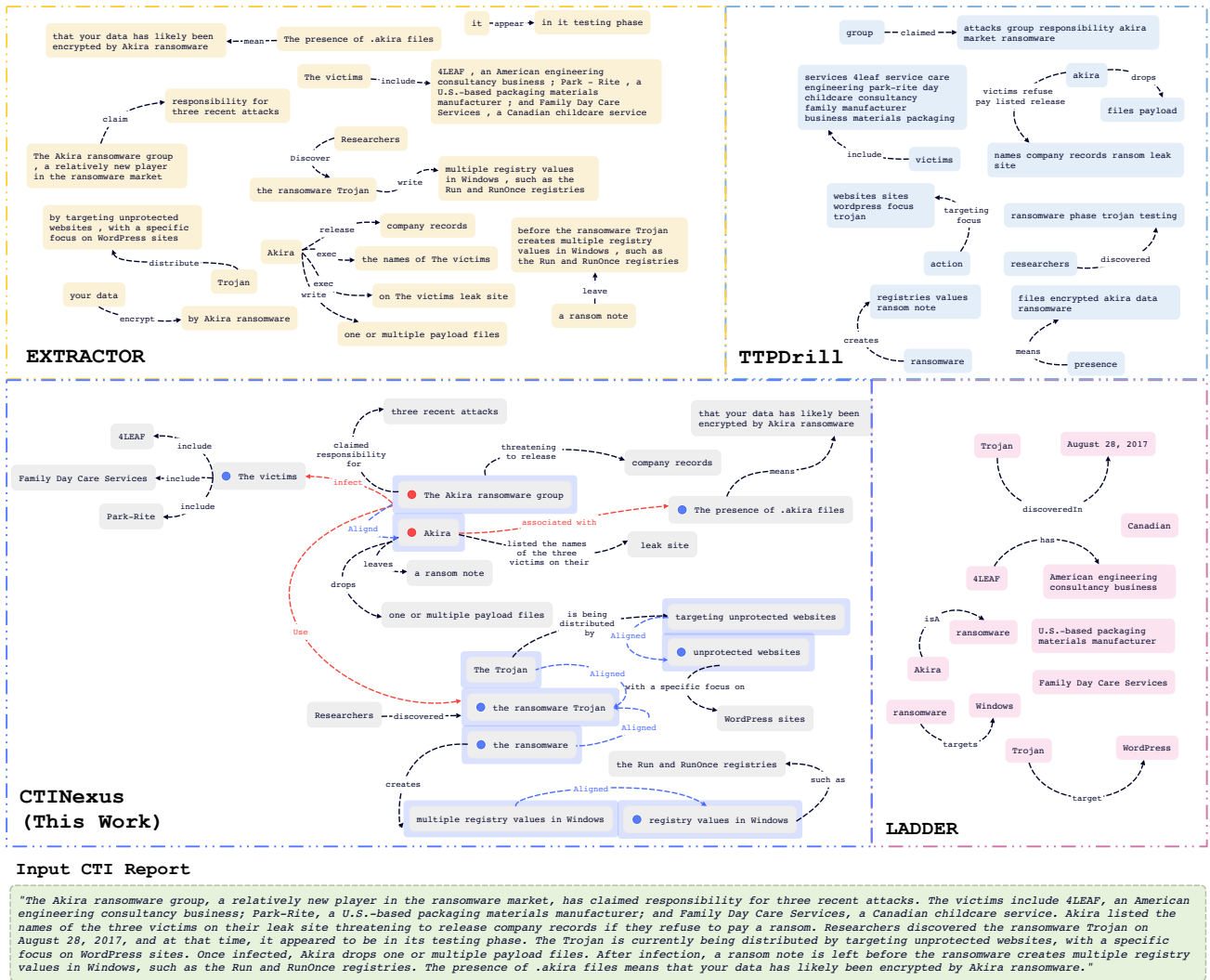


Fig. 1: CSKGs extracted by EXTRACTOR, TTPDrill, LADDER, and CTINEXUS for a real-world CTI report. EXTRACTOR, TTPDrill, and LADDER tend to produce incomplete and fragmented subgraphs, lacking comprehensive contextual connections. In contrast, CTINEXUS constructs a more integrated and comprehensive CSKG, with key information extracted and entities linked, providing a clearer and more complete representation of the threat profile.

paradigms: syntax parsing-based and fine-tuning-based. *Syntax parsing-based* methods leverage typed dependency rules to analyze the grammatical structure of a sentence and extract subject-verb-object (SVO) triplets. For example, TTPDrill [43] extracts subject entities and verb relations in CTI-related sentences as threat actions. iACE [56] extracts verb relations between IOCs and context terms. ThreatRaptor [39] extracts verb relations between subject IOC and object IOC. However, syntax parsing-based methods have *two main drawbacks*:

- **Domain complexity:** The grammatical rules can apply to any domain. However, CTI text in the security domain has several peculiarities that can confuse syntax parsing, leading to inaccurate extraction. Cybersecurity entities can contain special characters, such as dots in IPv4 addresses, underscores in file names, and slashes in file paths. These special characters can confuse basic NLP modules, like sentence segmentation and tokenization, which syntax parsing relies upon.

- **Static nature:** These methods rely on fixed syntax rules and predefined dictionaries to filter out irrelevant information and canonicalize extracted information. For example, TTPDrill maps extracted SVOs to a curated list of threat action terms, while ThreatRaptor uses a dictionary to canonicalize the extracted relation verbs. Keeping up with the evolving threat landscape requires continuous updates and maintenance of these rules and dictionaries, which is hard to scale.

On the other hand, *fine-tuning-based* approaches fine-tune pre-trained neural networks on annotated CTI domain datasets to perform named entity recognition (NER) and relation extraction (RE). For example, EXTRACTOR [71] fine-tunes a pre-trained BERT [72] model with thousands of annotated CTI sentences, to perform semantic role labeling to extract subjects, objects, and verb actions. AttackKG [55] fine-tunes a pre-trained model in the SpaCy library [41] to recognize entities and extract dependencies. LADDER [30] fine-tunes different pre-trained transformer models, including BERT, RoBERTa, and XML-RoBERTa, on their custom datasets annotated according to their own

ontology for performing NER and RE. However, fine-tuning-based methods also have *several drawbacks*:

- *Resource requirement*: Fine-tuning requires large amounts of labeled data (i.e., annotated CTI text corpora), and the labeling needs to be aligned with the targeted ontology. Such annotations are expensive to obtain, especially for emerging threats. Additionally, fine-tuning can be computationally expensive if the backbone model contains lots of parameters.
- *Ontology lock-in*: Since the models are fine-tuned on datasets annotated using a specific ontology, they are difficult to generalize to new ontologies that cover different entities and relations. Transferring to other ontologies would require reannotating vast data and retraining the models.

2.2.1. Motivating Example. We further investigate the quality of the constructed CSKG by existing approaches using a real-world CTI report. Fig. 1 illustrates a snippet of the report titled “RANSOMWARE - AKIRA AND RAPTURE” published on May 9, 2023, by Avertium [3]. The report provides rich information about the new Akira ransomware group. We run this CTI text snippet with three representative approaches, TTPDrill, EXTRACTOR, and LADDER using their released implementations [24], [7], [13]. Fig. 1 shows their constructed CSKGs. **We observe that the quality of CSKGs is very low.**

- *Some triplets have wrong directions.* For example, in EXTRACTOR, “ransom note” is extracted as the subject of “leave”, whereas it should be the object.
- *Many extracted entities have poor quality.* Some are not meaningful, such as “presence” extracted by TTPDrill. Others include unnecessary words or combine multiple distinct entities; for example, TTPDrill incorrectly extracts “registry values” and “ransom note” together when they should be separate. Similarly, in EXTRACTOR, the victim entities are not properly distinguished and should be individually separated. Although LADDER’s extracted content is of higher quality compared to TTPDrill and EXTRACTOR, it often lacks completeness. For instance, in the context where a “Trojan” targets “WordPress sites”, LADDER only extracts “WordPress” thereby omitting contextual information from the original phrase.
- *Entities are not aligned.* For example, in EXTRACTOR, “Trojan” and “the ransomware Trojan” refer to the same object and should be merged or associated. The same issue is observed in TTPDrill and LADDER.
- *Some critical relations are missing.* In the text, “the Akira ransomware group” uses the “ransomware Trojan” to launch the attack. However, since these two entities are mentioned in different sentences without explicit relational indicators, all approaches fail to infer the relationship between them.

As shown in Fig. 1, the CSKG constructed by CTINEXUS is comprehensive, well-connected, and of much higher quality, addressing all previous drawbacks. By leveraging the in-context learning of LLMs, the construction of such a CSKG does not rely on large amounts

of training data and can flexibly adapt to different ontologies. We describe our approach in Section 4.

2.3. Large Language Models

Recently, LLMs have shown emergent abilities to learn from just *a few demonstration examples* in the prompt, a paradigm known as in-context learning (ICL) [37]. In the ICL paradigm, the prompt input to the LLM typically includes three components: (1) an instruction specifying the task, (2) several demonstration examples containing ground truth to provide task-specific knowledge, and (3) a query to the LLM with the expectation of an appropriate answer. This allows LLMs to adapt to new tasks with minimal cost using task-specific prompts and demonstration examples. Multiple studies have shown that LLMs perform well in various tasks under ICL, such as fact retrieval [79] and mathematical reasoning [45], [28]. Additionally, LLMs have shown promise in different cybersecurity tasks, such as vulnerability detection [38], [59], patch generation [50], and software fuzzing [81], [61]. However, the use of LLMs for CTI knowledge extraction and CSKG construction remains largely underexplored.

3. Overview

Fig. 2 illustrates CTINEXUS. CTINEXUS introduces a novel ICL-based approach for data-efficient CTI knowledge extraction and CSKG construction. Unlike previous methods, CTINEXUS eliminates the need for extensive data annotations and parameter tuning, facilitating easy generalization to different ontologies. CTINEXUS focuses on constructing a connected and comprehensive CSKG, enabling entity alignment and long-distance relation inference. CTINEXUS consists of three phases.

Phase 1: Given a CTI report, CTINEXUS first extracts entity-relation triplets that align with the task ontology. The kNN-based demonstration retriever embeds the report and the candidate reports in the demonstration set into a high-dimensional latent space. The retriever then selects the top- k candidates with the highest similarity scores. The selected demonstrations are fed into an automatic prompt construction module to create a customized prompt for the current report. As illustrated in Fig. 2, our prompt template consists of three sections: an instruction describing the task, a query containing the input CTI report, and demonstration examples arranged in a specific order. Fig. 3 illustrates our carefully designed instruction. *Note that the ontology is incorporated into the instruction.* This design allows different ontologies to be easily switched, and our automatic prompt construction module will create a prompt specifically for this ontology and report, enhancing knowledge extraction performance.

Phase 2: With the extracted triplets, CTINEXUS removes redundancy by merging entities that refer to the same cybersecurity object using a hierarchical approach. The coarse-grained entity grouping module assigns types to entities using an automatically populated ICL prompt template, as illustrated in Fig. 4. The instruction incorporates the ontology that defines possible entity types. The demonstration examples show how to label each entity in the triplet. The query includes all the triplets to be typed. Entities assigned the same type are grouped together.

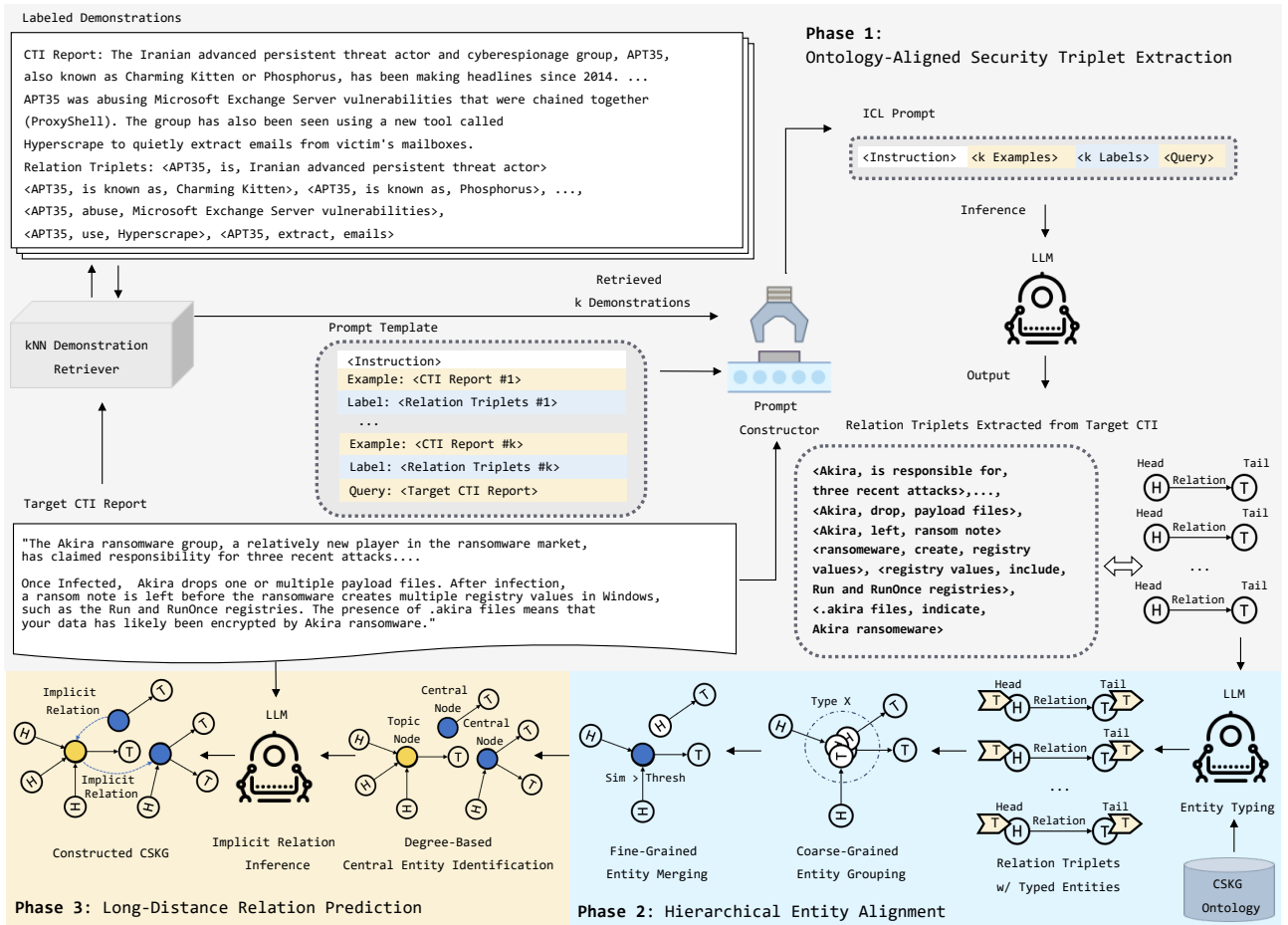


Fig. 2: Overview of CTINEXUS. CTINEXUS comprises three phases. Phase 1, *Security Triplet Extraction*, enables end-to-end extraction of cybersecurity triplets using in-context learning of LLM. Phase 2, *Hierarchical Entity Alignment*, reduces the redundancy of CSKG through coarse-grained grouping and fine-grained clustering. Phase 3, *Long-Distance Relation Prediction*, connects disjoint subgraphs by identifying central nodes and performing relation inference.

Next, the fine-grained entity merging module embeds all entities within each group and merges those that exceed a predefined similarity threshold into a single entity.

Phase 3: To infer missing links between distant entities, CTINEXUS performs long-distance relation prediction. The central entity identification module selects a central node in each connected subgraph based on the node’s degree centrality. Among central nodes, the module then selects a topic node with the highest importance, which serves as the main subject of the report. The central nodes and the topic node are passed to the ICL-enhanced relation prediction module to infer their implicit relationships. CTINEXUS automatically constructs an ICL prompt (illustrated in Fig. 5) to perform this inference.

4. Design of CTINEXUS

4.1. CSKG Ontology

We choose MALOnt for the current implementation, as MALOnt [69] is the most comprehensive among open-source ontologies, featuring 33 entity types (17 types and 16 sub-types) and 27 relation types. MALOnt covers a broad range of entities, such as Account, Action, Threat

Actor, Campaign, Event, Exploit Target, Host, Information, Infrastructure, Location, Malware, Person, Software, System, and Vulnerability, with detailed sub-types under Indicator and Malware Characteristics. However, note that CTINEXUS’s ICL-based pipeline eliminates the need for parameter tuning on large, ontology-specific training sets, largely simplifying generalization to other ontologies. If downstream applications require ontologies not covered by MALOnt, CTINEXUS can easily switch to a different ontology. This only requires a few demonstration examples aligned with the new ontology for each ICL task, and the ontology defined in a JSON format incorporated in the prompts (illustrated in Figs. 3 and 4). If the new ontology is a subset of MALOnt (which is already quite comprehensive), CTINEXUS can directly adapt by simply removing unrequired entity types without further actions.

4.2. Cybersecurity Triplet Extraction

Given that CTI text may contain diverse relations and we want the approach to be adaptable to emerging threats, we formulate the cybersecurity triplet extraction module in our pipeline as a *semi-open* extraction problem: Entity types are prescribed using MALOnt, as its coverage is already comprehensive, while relation extrac-

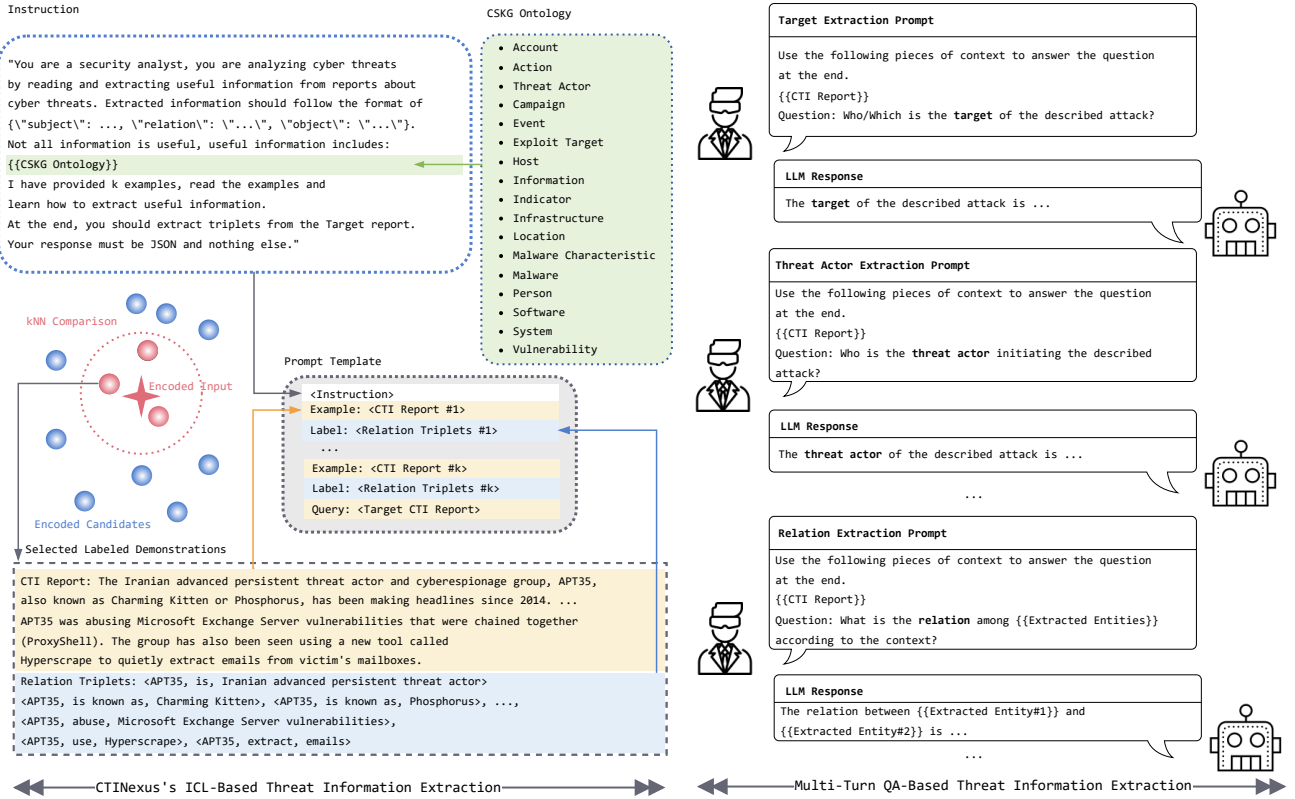


Fig. 3: Comparison of CTINEXUS’s ICL-based CTI knowledge extraction (left) and a multi-turn QA-based extraction (right). CTINEXUS consolidates task descriptions (including applied ontology), k selected demonstrations, and query into a single instruction for efficient cybersecurity triplet extraction. In contrast, the multi-turn QA paradigm requires multiple rounds of conversations with multiple prompts to extract different entities and relations, resulting in inefficiency.

tion is modeled as open RE to maximize the coverage. These approaches transform information extraction tasks into multi-turn question-answering, leveraging the conversational capabilities of LLMs. Fig. 3 illustrates this paradigm. This method involves creating multiple questioning prompts for each information type and refining the responses. However, applying this multi-turn QA formulation to cybersecurity entity and relation extraction requires numerous lengthy prompts due to the extensive cybersecurity ontology that could contain many entity classes. For N entities in the input CTI, $\frac{N(N-1)}{2}$ prompts are needed to extract relations between identified entities, leading to repetitive content and significant token waste, hindering scalability. Additionally, the multi-turn paradigm suffers from confirmation bias [34], as LLMs may confirm with a non-existing relation after several rounds of dialogue. These erroneous links can be particularly harmful in the CTI domain, negatively affecting downstream defense solutions by producing false alarms.

ICL prompt template. To improve efficiency and reduce confirmation bias, we develop a kNN-enhanced ICL paradigm that completes the cybersecurity triplet extraction process with only one LLM query. As illustrated in Fig. 3, CTINEXUS extracts all cybersecurity triplets by automatically populating a comprehensive ICL prompt template, which comprises the following components:

(1) *Instruction:* The instruction specifies the task, the applied ontology, and the required format for the ex-

tracted triplets. Instruction design is critical in LLMs, as an unclear definition of the task can severely degrade the performance. We carefully designed several versions of the instruction and identified the one presented in Fig. 3 as the most effective.

(2) *Demonstrations:* Top- k most relevant examples are retrieved using the demonstration retriever. Each example consists of a CTI report annotated with the security triplets. These examples are ordered in ascending similarity to the input query based on findings described in Section 5.3.

(3) *Query:* The input CTI text that needs to be analyzed.

kNN-based demonstration retriever. Multiple studies [67], [57] have shown that prompt examples selection can significantly affect LLM’s ICL capacity. One approach for selecting demonstration examples involves training a proxy LM to score candidates in the demonstration set [82]. However, this method requires large amounts of labeled data, which conflicts with our goal of designing a data-efficient solution. Recently, a k-nearest neighbors (kNN) method for selecting the most relevant demonstration examples based on semantic similarity has proven effective [57]. This method requires no dataset annotation or model tuning, making it ideal for our purposes. Specifically, we compute high-dimensional embeddings for the query and all candidate demonstrations using a pre-trained embedding model. Among the models explored,

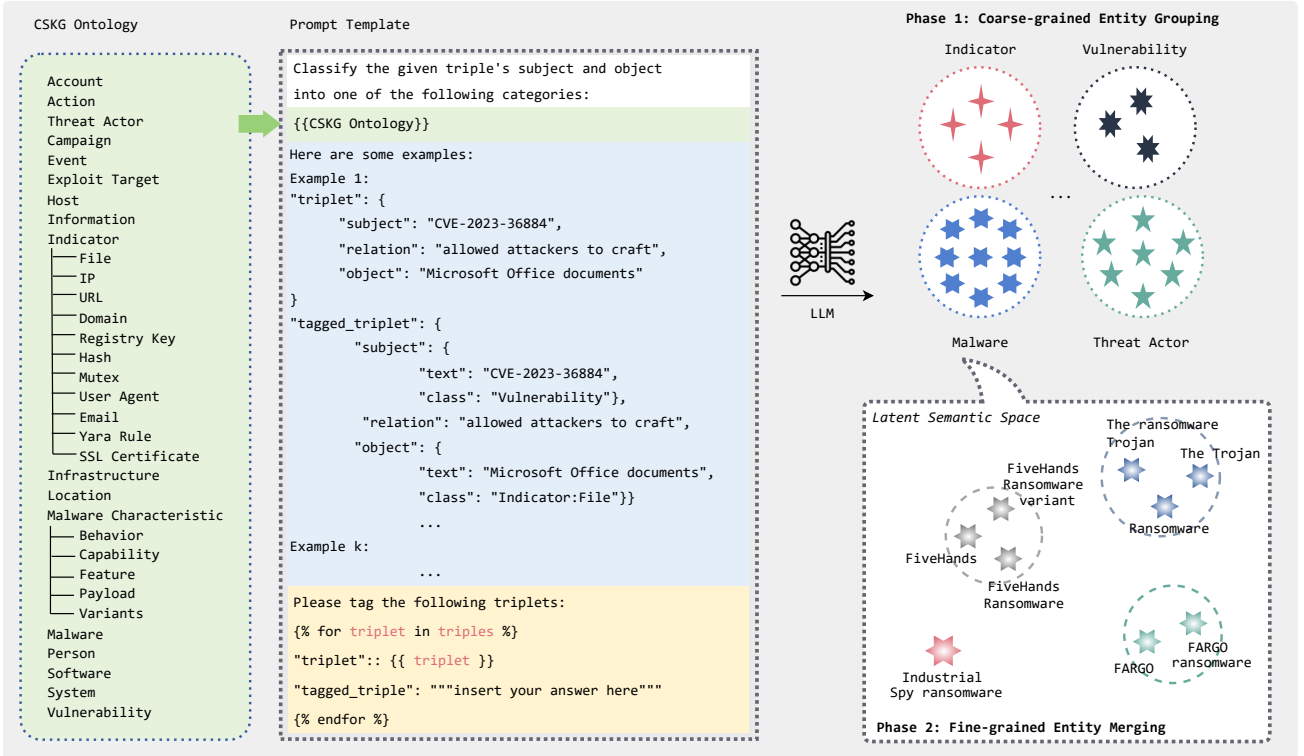


Fig. 4: The design of CTINEXUS’s hierarchical entity alignment. The coarse-grained entity grouping module populates an ICL prompt to assign entity types to the extracted triplets according to the applied ontology. Entities with the same type are grouped together. The fine-grained entity merging module then uses an embedding-based technique to merge semantically similar entities within each group based on a predefined similarity threshold.

text-embedding-3-large yielded the best performance. We then calculate the cosine similarity between the query embedding and each candidate demonstration’s embedding, selecting the top- k most similar candidates.

Several studies [57], [37], [62] have pointed out that the order of demonstration examples can also affect the performance of ICL. In particular, the model’s prediction often exhibits a recency bias [58], meaning that LLMs tend to pay more attention to the demonstration placed near the query. Also, kNN similarity indicates that if the demonstration is more similar to the query, LLMs can better analogize it. To investigate the impact of demonstration order in the CTI domain, we evaluated various permutations, including random, ascending, and descending orders (Section 5.3). Our findings indicate that arranging the demonstration examples in ascending order of their similarity to the query yields the best performance. This confirms the recency bias phenomenon in our scenario, as the demonstration example most similar to the query is placed at the bottom of the list, closest to the query.

4.3. Hierarchical Entity Alignment

Entity alignment identifies entities with different mentions that refer to the same real-world object, a key area in knowledge graph research [26]. Aligning these mentions integrates sub-graphs containing complementary knowledge, enhancing the comprehensiveness of the knowledge graph. Traditional entity alignment techniques rely on heuristics like string matching and structural similarities,

which fail to capture the underlying semantics or context of entities and have limited accuracy.

Recent studies [74], [78], [65] have adopted deep learning-based methods to learn vector representations (i.e., embeddings) of entities, achieving better accuracy. However, embedding-based techniques face unique challenges in our problem domain. In CTI text, entities with similar embeddings may refer to different concepts, e.g., “.akira files” (an IOC) and “Akira” (a threat actor). Besides, comparing the semantic distance between every pair of entities has a computational complexity of n^2 , where n is the total number of entities. This is inefficient when n becomes large.

To address these challenges, we perform entity alignment in a hierarchical way. The coarse-grained entity grouping module leverages LLM’s ICL ability to assign types to entities. Entities assigned the same type are then grouped together as potential candidates for alignment, narrowing the scope for later fine-grained merging. Fig. 4 illustrates our prompt template. CTINEXUS automatically creates a customized prompt by assembling k carefully annotated demonstration examples. Each demonstration example contains an untagged triplet and a tagged triplet with subject and object entities assigned type labels. The query part automatically traverses all triplets generated by the triplet extraction phase. For each triplet, we add a placeholder, “tagged_triplet”: “insert your answer here” to follow the format provided in the demonstration examples, better guiding the LLM to correctly fill in the answers.

For entities within each group, the fine-grained entity merging module uses an embedding-based technique to

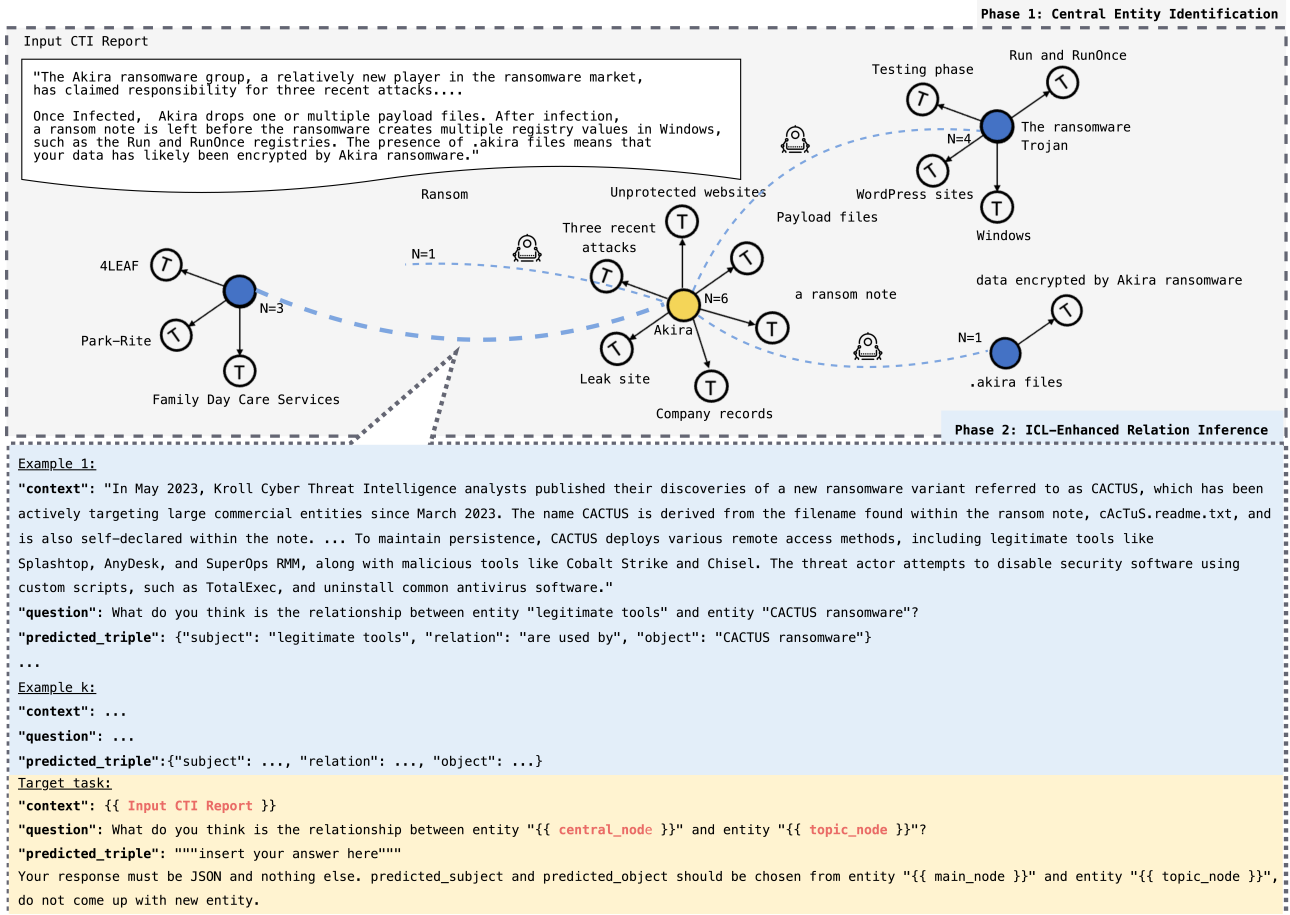


Fig. 5: The design of CTINEXUS's long-distance relation prediction. Phase 1 selects central entities (blue) and the topic entity (yellow) from separate subgraphs based on their degree centrality. Phase 2 populates an ICL prompt to infer implicit relations between each central entity and the topic entity.

merge entities with similar semantic representations. The embedding model is central to this procedure, as its generated embeddings are used to determine the semantic closeness of entities. We evaluated state-of-the-art, general-purpose text embedding models of various sizes (i.e., text-embedding-3-small, text-embedding-3-large) for this task. Since these models are not specifically pre-trained on a cybersecurity corpus, we also experimented with a security-specific embedding model, SecureBERT [27], which has been pre-trained on millions of cybersecurity websites, articles, and books. Another aspect to consider is the similarity threshold (degree of closeness) for determining alignment. To find the optimal threshold value, we experimented with common threshold values in semantic similarity comparison [66]: 0.4, 0.5, 0.6, and 0.7. Our results indicated that 0.6 is the most effective value. Detailed results and discussions are in Section 5.4.1.

4.4. Long-Distance Relation Prediction

After entity alignment, the triplets form a set of disconnected subgraphs, leaving implicit relations between distant entities unidentified. Previous methods primarily rely on graph structure learning and graph neural networks [86], [48] to perform link prediction. However, these methods require large amounts of annotated graph

data for model training. Additionally, in the CTI analysis domain, establishing relationships between distant cybersecurity entities requires a deep natural language understanding of their corresponding context. To make the procedure more data-efficient, we develop a long-distance relation prediction technique leveraging the ICL ability of LLMs. Fig. 5 illustrates our design.

Creating links for every pair of distant entities would introduce excessive connections, complicating the CSKG and consuming significant computational resources. Thus, CTINEXUS first runs a depth-first search to find all connected subgraphs. Then, CTINEXUS leverages graph structure information to identify a central entity for each subgraph. A central entity represents the most important entity in the subgraph and will be the head for inter-subgraph connections. In our design, we identify central entities based on their degree centrality [84], which is the most widely used measure of a node's importance in a graph. It is easy to calculate, by counting the total number of edges that a node has to other nodes. The intuition is that an entity with the most explicit relations with other entities is more likely to be the core subject of that part of CTI text. Among all identified central entities, we further identify a topic entity, which is the one with the highest degree, representing the core subject of the entire CTI report. Specifically, we consider both incoming

and outgoing edges when calculating degree centrality to identify the central identity. If multiple entities have the same highest score, we further prioritize out-degree over in-degree, as subjects in triplets (e.g., “AndroXgh0st” in <“AndroXgh0st”, “targets”, “.env files”>) are generally more important than objects. If there is still a tie, they are all determined as central entities. We follow the same procedure for identifying the topic entity. In the example shown in Fig. 5, there are five subgraphs. We identify the following central entities: “Victim”, “Akira”, “the ransomware Trojan”, “Akira ransomware group”, and “.akira files”. We select “Akira” as the topic entity, which has the highest degree centrality score of 6. These central entities and the topic entity are then fed into the next module for relation inference.

The ICL-enhanced relation prediction module leverages ICL of LLM to infer implicit relations between each central entity and the topic entity, creating inter-subgraph connections. Fig. 5 illustrates our prompt template. For each central entity, CTINEXUS automatically creates a customized prompt by assembling k fixed, carefully annotated demonstration examples, similar to the entity alignment process. The prompt template consists of two sections: a demonstration section (blue) and a query section for the target task (yellow). Both sections include “context”, “question”, and “predicted_triple” components. The “context” component presents the CTI report, while the “question” component asks the LLM about the relation between the queried central entity and the topic entity. The “predicted_triple” component contains the annotated relations for the demonstration examples and a placeholder, “insert your answer here”, for the queried task. This consistent design across the three components in both the query and demonstration sections helps the LLM effectively analogize the demonstration examples, facilitating better relation inference.

5. Evaluation

To comprehensively study the performance of CTINEXUS

in various phases of CSKG construction, we set the following four research questions:

- RQ1:** How does CTINEXUS’s performance in extracting cybersecurity entities and relations compare to existing baseline methods?
- RQ2:** How well does CTINEXUS perform in cybersecurity triplet extraction?
- RQ3:** How well does CTINEXUS perform in knowledge graph construction?
- RQ4:** What is the efficiency of CTINEXUS?

5.1. Dataset Annotation

Existing datasets benchmark triplet extraction but do not encompass other procedures in our pipeline, and the CTI reports they include are mostly outdated. For example, the dataset constructed by LADDER is limited to reports from 2010 to 2021 [30]. To address this, we created a dataset to evaluate CTINEXUS across cybersecurity triplet extraction, hierarchical entity alignment,

and long-distance relation prediction phases. Our dataset includes 150 reports from May 2023 onwards, Published by organizations like Trend Micro, Symantec, and The Hacker News, it includes 10 sources, averaging 15 reports each. Annotators with over three years of threat analysis expertise followed a four-phase process: *Phase I* involved annotating all cybersecurity entities and selecting entity types; *Phase II* identified explicit relations among entities and organized them into JSON-formatted triplets; *Phase III* grouped entities by type and merged those referring to the same threat-related concept; *Phase IV* selected central and topic entities based on degree and summarized implicit relations. This process resulted in 59,776 mentions, 35,258 entities, and 34,876 relations, enabling an effective evaluation of CTINEXUS’s performance in constructing cohesive and complete CSKGs.

5.2. RQ1: How does CTINEXUS compare against existing CTI extraction methods?

We evaluate CTINEXUS against two state-of-the-art baselines: EXTRACTOR [71] and LADDER [30], representing syntactic analysis-based and fine-tuning-based approaches, respectively. Several methodological challenges were addressed to enable fair comparison. For EXTRACTOR, we adapted its output to our broader ontology using CTINEXUS’s coarse-grained entity grouping module. For LADDER, we addressed two key differences: (1) LADDER uses a word-level annotation format, where each token is labeled with its target class. In contrast, our dataset follows an end-to-end report-to-triplet format, where the entire report is input, and the label is a set of extracted triplets. (2) LADDER uses a simplified ontology derived from MALONT, which includes only 10 entity types, a subset of the entity types used in our ontology. To facilitate comparison, we developed scripts to tokenize our data and convert our manually annotated datasets into LADDER’s word-level format. To ensure a fair comparison with LADDER, we merged our training set with LADDER’s in a 5:1 ratio, maintaining their original training/validation split. We also replaced LADDER’s test set with ours to ensure consistent evaluation on the same data. We compare with LADDER solely on named entity extraction performance. The reason is that our method focuses on open relation extraction, while LADDER targets relation classification within fixed categories.

Table I demonstrates that CTINEXUS outperforms EXTRACTOR in all metrics in cybersecurity triplet extraction. The evaluation results in Section 5.3 showed that GPT-4 outperforms all other backbone models. Thus, we use GPT-4 as the default backbone model for CTINEXUS’s implementation. This superior performance can be attributed to several factors. First, CTINEXUS leverages the robust context understanding and instruction-following capabilities of LLMs and enhances specificity with kNN-selected demonstration examples for extracting triplets. In contrast, EXTRACTOR employs general fine-tuning to extract semantic roles not specific to any ontology, reducing its accuracy in triplet extraction. Also, the CTI context introduces peculiarities that lead to errors in EXTRACTOR’s semantic role labeling module, which relies on a simple BERT model. For instance, EXTRACTOR might extract a triplet like

TABLE I: Performance comparison of CTINEXUS and EXTRACTOR on cybersecurity triplet extraction.

Method	F1 Score	Precision	Recall
EXTRACTOR	62.29	51.62	78.53
CTINEXUS	87.65	93.69	82.34

TABLE II: Performance comparison of CTINEXUS and LADDER on cybersecurity entity extraction.

Method	F1 Score	Precision	Recall
LADDER	71.13	78.31	73.94
CTINEXUS	90.13	92.00	88.35

⟨“Androxxgh0st malware”, “support”, “numerous functions capable of abusing the Simple Mail Transfer Protocol (SMTP), such as scanning and exploiting exposed credentials and application programming interfaces (APIs), and web shell deployment”⟩, where the object is a long sentence not suitable as a single entity. The object contains multiple entities due to misidentified boundaries. Conversely, CTINEXUS can comprehend implicit meanings and transform phrases to be more suitable as entities, resulting in a triplet like ⟨“Androxxgh0st malware”, “supports”, “functions abusing SMTP”⟩.

Table II demonstrates that CTINEXUS outperforms LADDER in F1 Score, Precision, and Recall by 26.7%, 17.5%, and 19.5%, respectively. Specifically, LADDER achieved an F1 Score of 71.13%, Precision of 78.31%, and Recall of 73.94%, which are slightly lower than the numbers reported in LADDER’s original evaluation (75.32%, 79.06%, and 76.98%, respectively). LADDER’s lower performance on our test data compared to its reported values is likely due to a distribution shift. LADDER’s dataset spans 2015 to 2021, while our data is from May 2023 onward. This temporal gap may introduce new patterns, terminologies, or threat vectors that LADDER’s model struggles to generalize to, even when retrained on a mix of old and new data. The performance disparity between LADDER and CTINEXUS can be attributed to several factors. First, fine-tuning the model in LADDER may lead to overfitting on the training set. Consequently, when confronted with unseen entities in the test set, the model may struggle to recognize them accurately, potentially misclassifying them or recognizing only parts of the entities. For example, in the sentence “... with a specific focus on WordPress sites”, LADDER extracts only “WordPress” as an application, resulting in ambiguous content. In contrast, CTINEXUS correctly extracts “WordPress sites”, which more accurately reflects the original context. Second, similar to EXTRACTOR, the LADDER model lacks sufficient contextual understanding. For instance, in the sentence “The victims include Family Day Care Services, a Canadian childcare service”, LADDER incorrectly identifies “Canadian” as a “B-Location”, whereas it should be recognized as a descriptive term for the childcare service. Furthermore, LADDER models relation extraction as relation classification, limited to ten relation classes (i.e., closed-world setting). This constraint restricts the contextual information in the extracted content and hinders the model’s ability to generalize to new CTI data containing different or additional relation classes.

TABLE III: Impact of example numbers on CTINEXUS’s cybersecurity triplet extraction.

Demo. Num.	F1 Score	Precision	Recall	Input _{Len}
1	85.05	94.39	77.40	949.95
2	87.65	93.69	82.34	1539.68
3	87.04	93.62	81.31	2138.41
4	86.73	89.55	84.07	2761.38

TABLE IV: Impact of example permutation on CTINEXUS’s cybersecurity triplet extraction.

Permutation	F1 Score	Precision	Recall
kNN-ascend	87.65	93.69	82.34
kNN-descend	85.82	90.58	81.53
random	84.96	90.29	80.22

5.3. RQ2: How well does CTINEXUS perform in cybersecurity triplet extraction

To demonstrate the effectiveness of CTINEXUS in cybersecurity triplet extraction, stemming from the superiority of the ICL paradigm and our specific prompt design, we conducted experiments on different ICL configurations, focusing on three aspects: (1) the number of demonstration examples, (2) the permutation of these examples, and (3) the backbone model types. By default, CTINEXUS uses GPT-4 as the model backbone, selects the k most similar prompt examples sorting in ascending order of query similarity.

Impact of prompt example numbers. To investigate the impact of prompt example numbers, we evaluated 4 configurations: 1, 2, 3, and 4 examples. Our observations show effectiveness plateau when using 2 or 3 examples, while input ICL prompt size increases significantly with more examples. As shown in Table III, increasing the prompt example number from 1 to 2 improves all metrics, particularly recall. However, with 3 examples, precision and F1-score plateau, and recall drops by 1%. With 4 examples, recall improves from 82 to 84%, but precision drops from 93 to 89%. This contradicts the heuristic that more examples always improve ICL performance but aligns with Chandra et al. [32], noting that each scenario has an optimal number of examples. Additionally, each additional example increases the input length by an average of 603 tokens, slowing inference speed and increasing computational costs. Thus, our implementation uses two examples in the cybersecurity triplet extraction phase, balancing effectiveness and efficiency.

Impact of prompt example permutations. To analyze the effect of the permutation method for selected examples, we examined three strategies: (1) random selection and sorting (random), (2) selection based on kNN similarity and sorting in ascending order (kNN-ascend), and (3) selection based on kNN similarity and sorting in descending order (kNN-descend). These methods were chosen to explore the impact of recency bias in LLMs, which suggests that models give more weight to examples placed nearer to the query [58]. The random method serves as a baseline, while kNN-ascend and kNN-descend test the influence of example order based on similarity. As shown in Table IV, kNN-ascend outperforms other methods across

TABLE V: Impact of backbone models on CTINEXUS’s cybersecurity triplet extraction.

Backbone	F1 Score	Precision	Recall
GPT-4	87.65	93.69	82.34
GPT-3.5	76.97	82.37	72.24
Qwen2.5-72B	78.18	80.83	75.71
Llama3-70B	77.85	81.74	74.32

TABLE VI: Impact of example numbers on CTINEXUS’s coarse-grained entity grouping.

Model Config.	Acc	Micro-F1	Macro-F1
GPT-3.5 (1-shot)	61.50	74.71	78.50
GPT-3.5 (4-shot)	66.18	78.45	79.86
GPT-3.5 (8-shot)	69.52	80.99	82.16
GPT-3.5 (12-shot)	69.68	81.11	81.95
GPT-4 (1-shot)	76.98	86.27	86.10
GPT-4 (4-shot)	81.02	88.94	87.87
GPT-4 (8-shot)	82.58	89.94	89.24
GPT-4 (12-shot)	81.18	89.05	88.28

all metrics, indicating the presence of recency bias and its potential for improving results. Consequently, we adopted kNN-ascend for CTINEXUS and recommend arranging prompt examples in ascending order of similarity as a universal strategy for other ICL applications.

Impact of backbone models. The emergence of ICL is closely associated with the substantial parameter counts of LLMs. To assess CTINEXUS’s generalizability across different backbone models, we evaluate its performance on representative closed-source LLMs, GPT-3.5 and GPT-4, and leading open-source LLMs, Llama3 and Qwen2.5. As shown in Table V, CTINEXUS achieves over a 10% improvement in both recall and precision when using GPT-4 compared to GPT-3.5-turbo. This underscores the importance of leveraging larger models to fully exploit ICL’s potential within CTINEXUS’s framework. For Qwen2.5 and Llama3, due to computational resource limitations, we deployed their 72B and 70B parameter versions, respectively. As shown in Table V, both Qwen2.5 and Llama3 demonstrate performance generally comparable to GPT-3.5-turbo. Specifically, Qwen2.5 exhibits a 1.4% higher recall but a 0.9% lower precision compared to Llama3. GPT-4 excelled in both precision and recall among all evaluated backbones. Therefore, all subsequent experiments will employ GPT-4 as the default base model.

5.4. RQ3: How well does CTINEXUS perform in knowledge graph construction?

5.4.1. Hierarchical Entity Alignment. As described in Section 4, for entity alignment, we first apply ICL to perform coarse-grained grouping of entities based on their types. We then vectorize these entities into high-dimensional embeddings and conduct fine-grained merging based on their semantic similarity. In the following, we present a series of experiments to investigate the impact of different configurations in entity grouping and entity merging, aiming to identify the optimal combination.

Impact of demonstration numbers.

TABLE VII: Impact of merging threshold values on CTINEXUS’s fine-grained entity merging.

Threshold	F1 Score	Precision	Recall	EntityNum
0.4	90.10	81.99	100	13.32
0.5	95.18	90.80	100	15.13
0.6	99.80	99.61	100	16.62
0.7	96.29	99.58	93.21	17.50

TABLE VIII: Impact of embedding models on CTINEXUS’s fine-grained entity merging.

Model	F1 Score	Precision	Recall	EntityNum
SecureBERT	79.15	65.50	100	8.11
text-embedding-3-small	98.10	97.54	98.66	16.50
text-embedding-3-large	99.80	99.61	100	16.62

We assessed the impact of demonstration example numbers on ICL through comparative experiments with four example quantities: 1, 4, 8, and 12. Additionally, we evaluated the performance of two LLMs, GPT-4 and GPT-3.5, across different model sizes. Notably, the performance showed no significant improvement once the number of examples exceeded 12, so these results are excluded from the table. Our evaluation methodology used accuracy, macro-F1, and micro-F1 metrics, consistent with previous text classification studies [63]. The experimental results, shown in Table VI, indicate that GPT-4 consistently outperforms GPT-3.5 across all demonstration number hierarchies. Remarkably, GPT-4 with 1 demonstration yields better results than GPT-3.5 with 12 demonstrations. Both models show substantial improvements when increasing from one to eight demonstrations, but a saturation trend appears when the number of examples exceeds eight. This trend is especially evident in GPT-4, where all three metrics slightly decrease as demonstration numbers increase from eight to twelve.

Impact of embedding models and merging threshold. The entity merging module applies a text embedding model to vectorize candidate entities grouped by the entity grouping module and uses a merging threshold to identify equivalent entities. Our evaluation focuses on the selection of the embedding model and the determination of the merging threshold. We use OpenAI’s third-generation embedding models, text-embedding-3-small and text-embedding-3-large, which differ in vector size and represent the latest state-of-the-art general-purpose models. In addition, we also compare with SecureBERT[27], a cybersecurity-specific embedding model based on the RoBERTa architecture pre-trained on a large corpus of cybersecurity data. We consider merging thresholds of 0.4, 0.5, 0.6, and 0.7. Besides common metrics for entity alignment, we introduce *Num_ent*, which records the number of entities after alignment.

The experimental results are shown in Table VII and Table VIII. Threshold values of 0.4, 0.5, and 0.6 all achieve a 100% recall rate, indicating the algorithm’s ability to detect all entities that should be merged. However, lower thresholds can erroneously merge non-equivalent entities based on the *Num_ent* and precision metrics. The highest precision is observed when the merging threshold is 0.6. Increasing the threshold to 0.7 maintains preci-

TABLE IX: Impact of example numbers on CTINEXUS’s relation prediction.

Model Config.	F1 Score	Precision	Recall
GPT-3.5 (0-shot)	65.95	51.26	92.42
GPT-3.5 (1-shot)	70.21	55.46	95.65
GPT-3.5 (2-shot)	76.87	63.31	97.84
GPT-3.5 (3-shot)	74.83	61.06	96.46
GPT-4 (0-shot)	85.76	75.07	100
GPT-4 (1-shot)	89.13	80.39	100
GPT-4 (2-shot)	90.99	83.47	100
GPT-4 (3-shot)	89.00	80.11	100

sion but significantly reduces recall, suggesting overly fine granularity that misclassifies equivalent entities as distinct. Regarding embedding models, text-embedding-3-large demonstrates the best performance, with text-embedding-3-small showing similar results. SecureBERT, despite its high recall, struggles to correctly cluster entities, as reflected in its low precision and *Num_ent* scores. This may be due to the smaller size of RoBERTa compared to the text-embedding-3 models, leading to less accurate entity distinction.

5.4.2. ICL-Enhanced Relation Prediction. As mentioned in Section 4, we compose ICL prompts to guide LLMs in inferring relations between disconnected subgraphs using the provided examples and context. We evaluated different ICL settings by varying the number of demonstration examples (1, 2, and 3) and the sizes of backbone models. Additionally, we examined the effectiveness of zero-shot learning, where the LLM infers relationships of given entities *without* demonstration examples. Zero-shot learning results are excluded from previous ICL experiments due to poor performance. The better performance in implicit relation inference compared to other tasks in CTINEXUS could be that relation prediction aligns more closely with general NLP tasks. Unlike triplet extraction or entity alignment, which require domain-specific knowledge in the cybersecurity context, relation prediction relies more on LLMs’ general ability to infer connections between entities based on linguistic cues in the text. This makes relation prediction less dependent on specialized domain knowledge and more aligned with the LLM’s general language understanding capabilities.

Experimental results, shown in Table IX, indicate that GPT-4 outperforms GPT-3.5 in every setting by a large margin, achieving a 100% recall rate compared to 92%-96% for GPT-3.5. The reason for this discrepancy is that GPT-3.5 has a higher tendency to produce hallucinated answers, either by not following the required instructions for the task (e.g., generating relations between entities not present in the queries) or by not adhering to the required format (e.g., generating a string instead of the requested JSON format). Both models show suboptimal performance with zero-shot learning. Increasing the number of demonstration examples from 1 to 2 significantly improves results, but a slight decline is observed with 3-shot examples. This suggests that while some examples can enhance performance, too many examples may introduce additional complexity or noise.

5.5. RQ4: What is the efficiency of CTINEXUS?

In this RQ, we assess the average token and time costs of three modules within CTINEXUS, using GPT-3.5 and GPT-4 as backbone models. The results, shown in Table X, indicate that using GPT-4 as the backbone results in token costs 20-30 times higher than those of GPT-3.5. Additionally, the time cost of using GPT-4 is approximately twice as high compared to GPT-3.5 for each module and the overall pipeline. The ICL-enhanced relation prediction module is the most computationally expensive, requiring multiple inferences for each input CTI. In contrast, the cybersecurity triplet extraction and hierarchical entity alignment modules have similar token costs, approximately half that of the long-distance relation prediction module, as they adhere to the “one input, one inference” principle, making them more economical. Specifically, for the hierarchical entity alignment module, the token and time costs are mainly attributed to the coarse-grained entity grouping module. The fine-grained entity merging module, which uses the text-embedding-3-large model, incurs minimal costs (\$0.13 per 1M tokens), resulting in the entire experiment costing less than \$0.30.

6. Discussion

Limitations. In CTINEXUS, the demonstrations must be carefully chosen and of high quality, with correct answers and the required prompt format. This ensures that CTINEXUS can fully utilize the ICL capability to infer the correct answers from the provided examples. According to Zhao et al. [85], CTINEXUS’s performance degrades significantly if the demonstration set contains incorrect or misformatted samples. Additionally, although CTINEXUS can operate in a data-constrained manner, it still requires a certain amount of labeled data, with a recommended minimum of 100 samples. Data imbalance within the demonstration set also affects CTINEXUS’s performance, as an imbalanced label distribution leads to less diverse retrieved examples, increasing the likelihood of biased content generation and reducing overall effectiveness.

Hallucinations in LLMs. Large Language Models (LLMs) can generate hallucinations, which are plausible yet factually inaccurate outputs [42], [76]. For instance, CTINEXUS with GPT-3.5 extracted the incorrect triplet ⟨“July 2022”, “threat actors behind FARGO attacks were hijacking”, “vulnerable Microsoft SQL servers”⟩ instead of ⟨“vulnerable Microsoft SQL servers”, “are hijacked by”, “July 2022”⟩, leading to a complete misplacement of the subject and object and an incoherent relation. This issue is more prevalent in smaller models like GPT-3.5, LLaMA3-70B, and QWen2.5-72B. While potential solutions include fine-tuning hallucination detection classifiers or using stronger LLMs for verification, we leave these challenges for future work. Our current focus is on CSKG construction under data scarcity, where GPT-4 has demonstrated reliable performance.

Empowering downstream defenses. Various applications can be potentially empowered by CTINEXUS. For example, the extracted CTI knowledge can be converted into open formats like STIX [49] (e.g., also via LLMs), and exchanged in platforms like AlienVault OTX [2], and

TABLE X: Token and time costs of CTINEXUS across different modules.

	Ontology-Aligned Security Triplet Extraction				Hierarchical Entity Alignment				ICL-Enhanced Long-Distance Relation Prediction				Pipeline		
	Input	Output	Overall	Time	Input	Output	Overall	Time	Time	Input	Output	Overall	Time	Token Cost	Time
CTINEXUS w/ GPT-4	0.0246	0.0117	0.0364	11.0905	0.0158	0.0236	0.0393	26.1590	5.9887	0.0644	0.0083	0.0728	24.2483	0.1485	67.4865
CTINEXUS w/ GPT-3.5	0.0007	0.0006	0.0013	5.9824	0.0008	0.0010	0.0018	10.6606	5.9887	0.0033	0.0005	0.0038	9.5013	0.0069	32.1330

integrated into intrusion detection systems [52], [31]. A question-answering system can be developed upon the constructed CSKG using LLM’s retrieval-augmented generation [73], to provide grounded answers to threat-related questions. Cyber threat hunting [83] can also potentially be enhanced. For example, the effort required for manually constructing threat hunting queries can be reduced by using LLMs to synthesize or suggest next steps based on the constructed CSKG and partial user input. We leave the exploration of these applications for future work.

7. Related Work

In Section 2, we discussed CTI knowledge extraction works in detail. Here, we discuss additional related work.

CTI services and platforms. There exist several services that regularly publish updated CTI feeds. For example, PhishTank [19] and OpenPhish [18] focus on phishing URLs. Abuse.ch [1] focuses on malware samples and botnet C&C servers. A key limitation is that they only provide isolated IOC feeds. There are also several comprehensive platforms that allow users to (1) share CTI data with other members of the community to benefit from the crowd-sourced knowledge, or (2) systematically manage their CTI data. These systems often provide web interfaces for user exploration and APIs for system integration. For example, AlienVault OTX [2] and IBM X-Force Exchange [11] are company-owned crowd-sourced platforms for sharing and searching threat data like IOCs, malware, and vulnerabilities. MISP [14] is an open-source platform for sharing, storing, and correlating IOCs of targeted attacks. OpenCTI [17] is an open-source platform that allows users to structure, store, organize, and visualize their CTI knowledge and observables. Unlike CTINEXUS’s automated approach, these platforms require users to actively participate in the sharing process and *manually* contribute CTI data.

Cybersecurity knowledge bases. Several comprehensive cybersecurity knowledge bases have been created by the industry. CVE [5] and NVD [16] are two most widely used vulnerability databases. Several threat encyclopedias exist (Trend Micro [23], Kaspersky [12], F-Secure [8]) for malware and vulnerabilities. MITRE ATT&CK [15] is a knowledge base for cyber adversary tactics and techniques based on real-world observations. These knowledge bases are manually created by security experts, and hence their update frequency is typically low. The scope of CTINEXUS differs from these systems. Nevertheless, since these knowledge bases also contain textual CTI descriptions about malware and vulnerabilities, CTINEXUS can be applied to further structuralize such knowledge.

LLMs for cybersecurity. Recent works have explored applying LLMs to cybersecurity challenges. PentestGPT [35] investigates LLM capabilities in penetration testing, revealing that while LLMs can handle fundamental tasks and use testing tools competently, they struggle

with context loss and attention issues. TitanFuzz [36] introduces an innovative approach for fuzzing deep-learning libraries using LLMs. It employs a generative LLM for high-quality seed programs and an infilling LLM for mutations, significantly improving API and code coverage, and detects numerous previously unknown bugs. Recent studies have also explored the use of LLMs in tasks such as vulnerability detection [38], [59], patch generation [50], malware detection [29], [70], botnet traffic analysis [53], [40], and phishing and scam detection [46], [51]. Unlike these works, CTINEXUS leverages the ICL paradigm of LLMs for comprehensive CTI knowledge extraction and CSKG construction.

Other CTI research. Several studies have empirically examined various aspects of CTI, including understanding vulnerability reproducibility [64], evaluating the quality of CTI feeds in terms of volume, timeliness, and coverage [47], [54], and analyzing information inconsistencies [80]. These works offer valuable insights into the current state of CTI data. In contrast to these empirical efforts, CTINEXUS focuses on designing an LLM-empowered approach for automated extraction of CTI knowledge from CTI reports. The scope is different.

8. Conclusion and Future Work

We proposed CTINEXUS, a new framework leveraging ICL of LLMs for efficient and adaptive CTI extraction and CSKG construction. Unlike existing methods, CTINEXUS requires minimal data and parameter tuning and can adapt to various ontologies with minimal data annotation. Extensive evaluations demonstrated CTINEXUS’s effectiveness in extracting comprehensive knowledge, highlighting its potential to transform CTI analysis into a data-efficient and adaptable paradigm.

Future directions include integrating CTINEXUS with downstream applications, such as intrusion detection systems, penetration testing tools, or cybersecurity question-answering systems, to enable timely knowledge updates and ensure factual accuracy in alignment with the evolving threat landscape.

References

- [1] “Abuse.ch,” <https://abuse.ch/>.
- [2] “Alienvault open threat exchange,” <https://otx.alienvault.com/>.
- [3] “Avertium,” <https://www.avertium.com/>.
- [4] “Bleeping computer,” <https://www.bleepingcomputer.com/>.
- [5] “Cve - common vulnerabilities and exposures,” <https://cve.mitre.org/>.
- [6] Dark reading. <https://www.darkreading.com/>.
- [7] “Extractor,” <https://github.com/ksatvat/EXTRACTOR>.
- [8] “F-secure threat descriptions,” <https://www.f-secure.com/en/business/security-threats/threat-descriptions>.

- [9] "Google threat analysis group," <https://blog.google/threat-analysis-group/>.
- [10] The hacker news. <https://thehackernews.com/>.
- [11] "Ibm x-force exchange," <https://exchange.xforce.ibmcloud.com/>.
- [12] "Kaspersky threats," <https://threats.kaspersky.com/>.
- [13] "Ladder," <https://github.com/aiforsec/LADDER>.
- [14] "Misp-open source threat intelligence platform & open standards for threat information sharing," <https://www.misp-project.org/>.
- [15] "Mitre att&ck@," <https://attack.mitre.org/>.
- [16] "National vulnerability database," <https://nvd.nist.gov/>.
- [17] "Opencti - open cyber threat intelligence platform," <https://www.opencti.io/>.
- [18] "Openphish," <https://openphish.com/>.
- [19] "Phishtank," <https://www.phishtank.com/>.
- [20] "Securityweek," <https://www.securityweek.com/>.
- [21] "Symantec security center," <https://symantec-enterprise-blogs.security.com/>.
- [22] Threatpost. <https://threatpost.com/>.
- [23] "Trend micro," <https://www.trendmicro.com/vinfo/us/security/news/>.
- [24] "Ttpdrill 0.5," <https://github.com/KaiLiu-Leo/TTPDrill-0.5>.
- [25] "Unit42," <https://unit42.paloaltonetworks.com/>.
- [26] B. Abu-Salih, "Domain-specific knowledge graphs: A survey," *Journal of Network and Computer Applications (JNCA)*, vol. 185, p. 103076, 2021.
- [27] E. Aghaei, X. Niu, W. Shadid, and E. Al-Shaer, "Securebert: A domain-specific language model for cybersecurity," in *International Conference on Security and Privacy in Communication Systems (SecureComm)*, 2022, pp. 39–56.
- [28] J. Ahn, R. Verma, R. Lou, D. Liu, R. Zhang, and W. Yin, "Large language models for mathematical reasoning: Progresses and challenges," *arXiv preprint arXiv:2402.00157*, 2024.
- [29] J. Al-Karaki, M. A.-Z. Khan, and M. Omar, "Exploring llms for malware detection: Review, framework design, and countermeasure approaches," *arXiv preprint arXiv:2409.07587*, 2024.
- [30] M. T. Alam, D. Bhusal, Y. Park, and N. 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 (RAID)*, 2023, pp. 92–108.
- [31] P. V. Bro, "A system for detecting network intruders in real-time," in *Proc. 7th USENIX security symposium (USENIX Security)*, 1998.
- [32] M. Chandra, D. Ganguly, Y. Li, and I. Ounis, "'one size doesn't fit all': Learning how many examples to use for in-context learning for improved text classification," *arXiv preprint arXiv:2403.06402*, 2024.
- [33] Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, H. Chen, X. Yi, C. Wang, Y. Wang *et al.*, "A survey on evaluation of large language models," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 15, no. 3, pp. 1–45, 2024.
- [34] Y.-S. Chuang, A. Goyal, N. Harlalka, S. Suresh, R. Hawkins, S. Yang, D. Shah, J. Hu, and T. T. Rogers, "Simulating opinion dynamics with networks of llm-based agents," *arXiv preprint arXiv:2311.09618*, 2023.
- [35] G. Deng, Y. Liu, V. Mayoral-Vilches, P. Liu, Y. Li, Y. Xu, T. Zhang, Y. Liu, M. Pinzger, and S. Rass, "Pentestgpt: An llm-empowered automatic penetration testing tool," *arXiv preprint arXiv:2308.06782*, 2023.
- [36] Y. Deng, C. S. Xia, H. Peng, C. Yang, and L. Zhang, "Large language models are zero-shot fuzzers: Fuzzing deep-learning libraries via large language models," in *Proceedings of the 32nd ACM SIGSOFT international symposium on software testing and analysis (ISSTA)*, 2023, pp. 423–435.
- [37] Q. Dong, L. Li, D. Dai, C. Zheng, Z. Wu, B. Chang, X. Sun, J. Xu, and Z. Sui, "A survey for in-context learning," *arXiv preprint arXiv:2301.00234*, 2022.
- [38] R. Fang, R. Bindu, A. Gupta, and D. Kang, "Llm agents can autonomously exploit one-day vulnerabilities," *arXiv preprint arXiv:2404.08144*, 2024.
- [39] P. Gao, F. Shao, X. Liu, X. Xiao, Z. Qin, F. Xu, P. Mittal, S. R. Kulkarni, and D. Song, "Enabling efficient cyber threat hunting with cyber threat intelligence," in *2021 IEEE 37th International Conference on Data Engineering (ICDE)*, 2021, pp. 193–204.
- [40] M. Guastalla, Y. Li, A. Hekmati, and B. Krishnamachari, "Application of large language models to ddos attack detection," in *International Conference on Security and Privacy in Cyber-Physical Systems and Smart Vehicles*. Springer, 2023, pp. 83–99.
- [41] M. Honnibal, I. Montani, S. Van Landeghem, and A. Boyd, "spaCy: Industrial-strength Natural Language Processing in Python," 2020.
- [42] L. Huang, W. Yu, W. Ma, W. Zhong, Z. Feng, H. Wang, Q. Chen, W. Peng, X. Feng, B. Qin *et al.*, "A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions," *arXiv preprint arXiv:2311.05232*, 2023.
- [43] G. Husari, E. Al-Shaer, M. Ahmed, B. Chu, and X. Niu, "Ttpdrill: Automatic and accurate extraction of threat actions from unstructured text of cti sources," in *Proceedings of the 33rd annual computer security applications conference (ACSAC)*, 2017, pp. 103–115.
- [44] M. Iannacone, S. Bohn, G. Nakamura, J. Gerth, K. Huffer, R. Bridges, E. Ferragut, and J. Goodall, "Developing an ontology for cyber security knowledge graphs," in *Proceedings of the 10th Annual Cyber and Information Security Research Conference (CISRC)*, 2015, pp. 1–4.
- [45] S. Imani, L. Du, and H. Shrivastava, "Mathprompter: Mathematical reasoning using large language models," *arXiv preprint arXiv:2303.05398*, 2023.
- [46] L. Jiang, "Detecting scams using large language models," *arXiv preprint arXiv:2402.03147*, 2024.
- [47] B. Jin, E. Kim, H. Lee, E. Bertino, D. Kim, and H. Kim, "Sharing cyber threat intelligence: Does it really help?" in *Proceedings of the 31st Annual Network and Distributed System Security Symposium (NDSS)*, 2024.
- [48] W. Jin, Y. Ma, X. Liu, X. Tang, S. Wang, and J. Tang, "Graph structure learning for robust graph neural networks," in *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining (KDD)*, 2020, pp. 66–74.
- [49] B. Jordan, R. Piazza, and T. Darley, "STIX version 2.1," <https://docs.oasis-open.org/cti/stix/v2.1/os/stix-v2.1-os.html>, June 2021.
- [50] U. Kulsum, H. Zhu, B. Xu, and M. d'Amorim, "A case study of llm for automated vulnerability repair: Assessing impact of reasoning and patch validation feedback," *arXiv preprint arXiv:2405.15690*, 2024.
- [51] J. Lee, P. Lim, B. Hooi, and D. M. Divakaran, "Multimodal large language models for phishing webpage detection and identification," *arXiv preprint arXiv:2408.05941*, 2024.
- [52] W. Lee and S. Stolfo, "Data mining approaches for intrusion detection," 1998.
- [53] Q. Li, Y. Zhang, Z. Jia, Y. Hu, L. Zhang, J. Zhang, Y. Xu, Y. Cui, Z. Guo, and X. Zhang, "Dollm: How large language models understanding network flow data to detect carpet bombing ddos," *arXiv preprint arXiv:2405.07638*, 2024.
- [54] V. G. Li, M. Dunn, P. Pearce, D. McCoy, G. M. Voelker, and S. Savage, "Reading the tea leaves: A comparative analysis of threat intelligence," in *28th USENIX Security Symposium (USENIX Security)*, 2019, pp. 851–867.
- [55] Z. Li, J. Zeng, Y. Chen, and Z. Liang, "Attackg: Constructing technique knowledge graph from cyber threat intelligence reports," in *European Symposium on Research in Computer Security (ESORICS)*, 2022, pp. 589–609.
- [56] X. Liao, K. Yuan, X. Wang, Z. Li, L. Xing, and R. Beyah, "Acing the ioc game: Toward automatic discovery and analysis of open-source cyber threat intelligence," in *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security (CCS)*, 2016, pp. 755–766.

- [57] J. Liu, D. Shen, Y. Zhang, B. Dolan, L. Carin, and W. Chen, "What makes good in-context examples for gpt-3?" *arXiv preprint arXiv:2101.06804*, 2021.
- [58] N. F. Liu, K. Lin, J. Hewitt, A. Paranjape, M. Bevilacqua, F. Petroni, and P. Liang, "Lost in the middle: How language models use long contexts," *Transactions of the Association for Computational Linguistics (ACL)*, vol. 12, pp. 157–173, 2024.
- [59] G. Lu, X. Ju, X. Chen, W. Pei, and Z. Cai, "Grace: Empowering llm-based software vulnerability detection with graph structure and in-context learning," *Journal of Systems and Software (JSS)*, vol. 212, p. 112031, 2024.
- [60] R. McMillan, "Definition: Threat intelligence," <https://www.gartner.com/en/documents/2487216>, May 2013.
- [61] R. Meng, M. Mirchev, M. Böhme, and A. Roychoudhury, "Large language model guided protocol fuzzing," in *Proceedings of the 31st Annual Network and Distributed System Security Symposium (NDSS)*, 2024.
- [62] S. Min, X. Lyu, A. Holtzman, M. Artetxe, M. Lewis, H. Hajishirzi, and L. Zettlemoyer, "Rethinking the role of demonstrations: What makes in-context learning work?" *arXiv preprint arXiv:2202.12837*, 2022.
- [63] S. Minaee, N. Kalchbrenner, E. Cambria, N. Nikzad, M. Chenaghlu, and J. Gao, "Deep learning–based text classification: a comprehensive review," *ACM computing surveys (CSUR)*, vol. 54, no. 3, pp. 1–40, 2021.
- [64] D. Mu, A. Cuevas, L. Yang, H. Hu, X. Xing, B. Mao, and G. Wang, "Understanding the reproducibility of crowd-reported security vulnerabilities," in *27th USENIX Security Symposium (USENIX Security)*, 2018.
- [65] S. Pei, L. Yu, R. Hoehndorf, and X. Zhang, "Semi-supervised entity alignment via knowledge graph embedding with awareness of degree difference," in *The world wide web conference (WWW)*, 2019, pp. 3130–3136.
- [66] M. T. Pilehvar, D. Jurgens, and R. Navigli, "Align, disambiguate and walk: A unified approach for measuring semantic similarity," in *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (ACL) (Volume 1: Long Papers)*, 2013, pp. 1341–1351.
- [67] C. Qin, A. Zhang, A. Dagar, and W. Ye, "In-context learning with iterative demonstration selection," *arXiv preprint arXiv:2310.09881*, 2023.
- [68] M. R. Rahman, R. M. Hezaveh, and L. Williams, "What are the attackers doing now? automating cyberthreat intelligence extraction from text on pace with the changing threat landscape: A survey," *ACM Computing Surveys (CSUR)*, vol. 55, no. 12, pp. 1–36, 2023.
- [69] N. Rastogi, S. Dutta, M. J. Zaki, A. Gittens, and C. Aggarwal, "Malont: An ontology for malware threat intelligence," in *International workshop on deployable machine learning for security defense (MLHat)*, 2020, pp. 28–44.
- [70] P. M. S. Sánchez, A. H. Celdrán, G. Bovet, and G. M. Pérez, "Transfer learning in pre-trained large language models for malware detection based on system calls," *arXiv preprint arXiv:2405.09318*, 2024.
- [71] K. Satvat, R. Gjomemo, and V. Venkatakrishnan, "Extractor: Extracting attack behavior from threat reports," in *2021 IEEE European Symposium on Security and Privacy (EuroS&P)*, 2021, pp. 598–615.
- [72] P. Shi and J. Lin, "Simple bert models for relation extraction and semantic role labeling," *arXiv preprint arXiv:1904.05255*, 2019.
- [73] J. Sun, C. Xu, L. Tang, S. Wang, C. Lin, Y. Gong, H.-Y. Shum, and J. Guo, "Think-on-graph: Deep and responsible reasoning of large language model with knowledge graph," *arXiv preprint arXiv:2307.07697*, 2023.
- [74] Z. Sun, W. Hu, Q. Zhang, and Y. Qu, "Bootstrapping entity alignment with knowledge graph embedding," in *International Joint Conference on Artificial Intelligence (IJCAI)*, vol. 18, no. 2018, 2018.
- [75] Z. Syed, A. Padia, T. Finin, L. Mathews, and A. Joshi, "Uco: A unified cybersecurity ontology," *UMBC Student Collection*, 2016.
- [76] S. Tonmoy, S. Zaman, V. Jain, A. Rani, V. Rawte, A. Chadha, and A. Das, "A comprehensive survey of hallucination mitigation techniques in large language models," *arXiv preprint arXiv:2401.01313*, 2024.
- [77] W. Tounsi, "What is cyber threat intelligence and how is it evolving?" *Cyber-Vigilance and Digital Trust: Cyber Security in the Era of Cloud Computing and IoT*, pp. 1–49, 2019.
- [78] B. D. Trisedya, J. Qi, and R. Zhang, "Entity alignment between knowledge graphs using attribute embeddings," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, no. 01, 2019, pp. 297–304.
- [79] C. Wang, X. Liu, Y. Yue, X. Tang, T. Zhang, C. Jiayang, Y. Yao, W. Gao, X. Hu, Z. Qi *et al.*, "Survey on factuality in large language models: Knowledge, retrieval and domain-specificity," *arXiv preprint arXiv:2310.07521*, 2023.
- [80] J. Wunder, A. Kurtz, C. Eichenmüller, F. Gassmann, and Z. Benenson, "Shedding light on cvss scoring inconsistencies: A user-centric study on evaluating widespread security vulnerabilities," *arXiv preprint arXiv:2308.15259*, 2023.
- [81] C. S. Xia, M. Paltenghi, J. Le Tian, M. Pradel, and L. Zhang, "Fuzz4all: Universal fuzzing with large language models," in *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering (ICSE)*, 2024, pp. 1–13.
- [82] X. Xu, Y. Liu, P. Pasupat, M. Kazemi *et al.*, "In-context learning with retrieved demonstrations for language models: A survey," *arXiv preprint arXiv:2401.11624*, 2024.
- [83] F. Yang, Y. Han, Y. Ding, Q. Tan, and Z. Xu, "A flexible approach for cyber threat hunting based on kernel audit records," *Cybersecurity*, vol. 5, no. 1, p. 11, 2022.
- [84] J. Zhang and Y. Luo, "Degree centrality, betweenness centrality, and closeness centrality in social network," in *2017 2nd international conference on modelling, simulation and applied mathematics (MSAM)*. Atlantis press, 2017, pp. 300–303.
- [85] H. Zhao, M. Andriushchenko, F. Croce, and N. Flammarion, "Is in-context learning sufficient for instruction following in llms?" *arXiv preprint arXiv:2405.19874*, 2024.
- [86] Y. Zhu, W. Xu, J. Zhang, Y. Du, J. Zhang, Q. Liu, C. Yang, and S. Wu, "A survey on graph structure learning: Progress and opportunities," *arXiv preprint arXiv:2103.03036*, 2021.

Appendix A. Data Availability

Upon acceptance of this paper, all code and datasets necessary to reproduce the results will be made publicly available under an open-source license. The code and data will be shared via a publicly accessible repository to ensure transparency and facilitate further research in the field.