# Cognitive AI Systems for Next Generation Digital Forensics and Incident Response

Rosemary Chisom Dimakunne[1]; David Yakubu[2]

[1]Department of Management Information Systems, Baylor University, Texas, USA.
[2]Department of Political Science and Governance, Villanova University. Pennsylvania, USA.

## Abstract

Critical infrastructure cyber incidents are growing in frequency and complexity, straining traditional digital forensics and incident response (DFIR) approaches. This paper explores the integration of cognitive artificial intelligence (AI) techniques, specifically natural language processing (NLP) and computer vision into DFIR workflows. By automating evidence extraction from large volumes of text and social media, classifying malware via image analysis, and correlating multistream log data, the proposed system aims to accelerate investigations during critical infrastructure attacks. The goals are to improve detection accuracy (e.g., malware family classification and anomaly detection in logs), enhance investigators' situational awareness through automated analytics, and significantly reduce the time needed to collect, analyze, and interpret digital evidence. The primary contributions include an architectural framework for a multi modal AI driven DFIR system, a methodology for ensuring evidentiary integrity (using hashing and blockchain), and an evaluation on representative datasets. Experimental results indicate the cognitive AI approach can boost efficiency and accuracy of forensic analysis while maintaining a strict chain of custody, thereby demonstrating the potential of AI to transform next generation DFIR in critical infrastructure contexts.

*Keywords: Digital Forensics; Incident Response; Natural Language Processing; Computer Vision; Malware Classification; Log Correlation; Critical Infrastructure; Cognitive AI.*

## I. INTRODUCTION

Modern cyber-attacks on critical infrastructure (power grids, water systems, industrial control systems) generate massive volumes of heterogeneous digital evidence. Incidents often involve *distributed data sources and diverse artifacts* for example, an attack campaign may leave traces across network device logs, control system alerts, malicious binaries, employee emails, and even social media posts (as attackers boast or insiders leak information). Collecting and analyzing such evidence manually is time consuming and complex, hampering timely incident response[1][2]. Digital evidence can reside in disparate forms including text communications, images, and executable code; connecting the dots across these manually presents a significant challenge. Traditional forensic methods struggle when faced with *distributed, real time data and limited logging*, as is common in supervisory control and data acquisition (SCADA) environments[3][4]. Many industrial control devices lack detailed logs or standardized monitoring, so investigators may have little to work with besides partial network captures or device memory, which complicates root cause analysis[3][5]. In such scenarios, responding effectively to incidents and preserving a strict chain of custody for all evidence is extremely difficult.

➤ *Role of AI in DFIR:*
Advances in artificial intelligence offer a promising way to augment digital forensics and incident response. Machine learning and in particular *cognitive AI* referring to AI that interprets unstructured data like human language or images can automate labor intensive forensic tasks. Natural language processing (NLP) techniques can sift through large text-based datasets (e.g., logs, emails, chat messages, social media) to extract relevant evidence and insights. For instance, NLP based systems can automatically classify or triage posts and messages for relevance, perform keyword and sentiment analysis to flag threatening communications, and identify entities like names, locations, or IP addresses in forensic text dumps[6][7]. This ability to rapidly process text can surface leads and patterns that investigators might otherwise miss or take too long to find. Likewise,

computer vision (CV) techniques enable novel approaches to analyzing binary and multimedia evidence. Researchers have demonstrated that malware files can be *converted into grayscale images* and then classified with convolutional neural networks (CNNs) at extremely high accuracy (over 99% on known malware families)[8]. This image-based malware analysis uses deep learning to automatically learn features, avoiding the need for hand crafted signatures or heuristics. Vision algorithms can also help detect manipulated media (like deepfake videos or images) by analyzing facial movements or texture artifacts that might not be obvious to the naked eye[9][10]. By automating recognition of such subtle patterns, AI based tools can alert investigators to forged evidence or malicious code far faster than manual methods.

Initial research results are promising. In log analysis, deep learning models have achieved *high detection rates* ($F_1$-scores above 0.95) on benchmark datasets[11][12]. However, current AI solutions in DFIR tend to address individual tasks in isolation. One system might focus on malware classification, another on social media analytics, another on anomaly detection in system logs few, if any, integrate these capabilities into a unified toolset. Moreover, existing solutions often do not fully adhere to forensic requirements for evidence handling. For example, an AI model might output an anomaly score, but if it cannot provide an explanation or if its results cannot be directly validated against raw data, its usefulness as court admissible evidence is limited. There is a *research gap* in combining multi modal AI analysis with rigorous forensic principles like chain of custody and evidence integrity.

➢ *Research Questions:*

To address these gaps, this work is guided by the following research questions: (1) *How can NLP and computer vision techniques be combined in a unified system to automatically analyze heterogeneous digital evidence from critical infrastructure incidents?* (2) *To what extent can such a cognitive AI system improve the efficiency and accuracy of digital forensic analysis (e.g., faster evidence processing, higher detection rates) compared to traditional manual approaches?* (3) *How can we ensure that the use of AI in DFIR does not compromise evidence integrity or admissibility, maintaining a strict chain of custody for all data?* These questions frame our objectives and the contributions described below.

➢ *Objectives and Contributions:*

This paper proposes a *cognitive AI architecture for next generation DFIR* that integrates NLP and computer vision in a single system. The system is designed with critical infrastructure scenarios in mind (e.g., forensic investigations on SCADA/ICS networks), where evidence is distributed and time is of the essence. Key components of the architecture include: an automated NLP engine for evidence extraction from text (such as parsing log files, identifying entities, and analyzing sentiments) and a computer vision engine for image-based analysis of malware binaries and media. These feed into a correlation module that uses deep learning to detect anomalies in log

sequences and link events across multiple data sources, reconstructing a timeline of the incident. Crucially, the system incorporates *evidence integrity mechanisms* every piece of collected evidence is hashed using secure algorithms and recorded on a tamper proof blockchain ledger to ensure an immutable chain of custody[13][14]. We provide a detailed design of this architecture and a prototype implementation that we evaluate on a mix of open datasets (including a malware corpus, public log datasets, and synthetic forensic data). We demonstrate that the cognitive AI system can significantly outperform baseline tools: for example, automatically classifying malware with nearly 99.9% accuracy (versus ~95% by a traditional method), and detecting log anomalies with high precision while correlating multi source events into coherent incident timelines. Throughout the design, we have followed best practices for admissibility using interpretable features where possible, keeping detailed audit logs of AI operations, and incorporating fail safes to prevent evidence alteration. We also discuss the *ethical implications* of deploying AI in forensic investigations, including issues of transparency, bias, and privacy. By addressing both the technical and procedural aspects, this work aims to advance DFIR into a new era where cognitive AI systems act as force multipliers for human investigators, enabling faster and more effective response to cyber crises in critical infrastructure.

## II. BACKGROUND AND RELATED WORK

➢ *Digital Forensic Process Models*

Digital forensics traditionally follows a well-defined process model comprising phases such as evidence acquisition, preservation, analysis, and presentation/reporting. Across these phases, maintaining a strict chain of custody and evidence integrity is paramount. Guidelines by standards organizations like NIST emphasize that digital data is easily altered, so investigators must take precautions to preserve authenticity. For example, best practice is to create cryptographic hashes of all collected digital evidence (using NIST approved algorithms like SHA 256) and to store those hash values securely and separately[13][15]. By verifying hashes, one can demonstrate that copies of evidence are identical to the original, enabling multiple copies to be used without invoking "best evidence" issues[16]. It is also recommended to create multiple backup copies of digital evidence and to use write blocking or imaging techniques during acquisition to avoid modifying the source data. Documentation at every step is required so that an unbroken chain of custody can be shown from evidence seizure to courtroom presentation. Any AI system used in DFIR must operate within this process model: for instance, when the system ingests logs or disk images, it should compute and record hashes, and any AI derived outputs (e.g., a classification of a file as malware) should be traceable back to the original evidence for validation.

Common forensic process frameworks (such as the DFRWS model or the DOJ's digital evidence flow) underscore similar phases. In the *acquisition* phase,

especially for incident response in live environments, techniques like live memory imaging or network traffic capture might be employed, but in critical infrastructure cases, taking systems offline is often infeasible. Therefore, "live" forensics and on site analysis have gained importance[17][18]. The *preservation* phase overlaps with acquisition; it involves ensuring that evidence remains in an unaltered state. This often means using hardware or software write blockers, or storing original data on tamper evident media. With digital evidence, preservation also extends to meta data (timestamps, file system info), which must be recorded. Next, the *analysis* phase is where AI tools can dramatically assist: instead of an investigator manually combing through logs or reverse engineering malware, automated algorithms can process data at scale and speed. Finally, in the *presentation* phase, results need to be presented clearly (often in written reports and possibly expert testimony). Any AI based findings must be explained in an understandable way, linking back to factual evidence; otherwise they may be challenged in court. Our proposed system is mindful of these requirements for example, it logs all actions taken by the AI and can produce a report of how a given conclusion (like a malware classification) was reached, including which evidence was used.

➢ *Forensics in Critical Infrastructure (SCADA/ICS)*

Supervisory Control and Data Acquisition (SCADA) and other Industrial Control Systems (ICS) pose unique challenges for digital forensics. These systems control physical processes (water treatment, power generation, manufacturing, etc.) and were historically designed for isolated operation, without built in security or forensic logging features[19][3]. Modern critical infrastructure, however, is increasingly network connected, exposing it to cyber threats. When incidents occur (e.g., a malware like Stuxnet or Industroyer disrupting operations), investigators face a difficult environment: heterogeneous devices (PLCs, RTUs, field sensors, operator workstations) often running proprietary firmware, minimal or non standard operating systems, and generating scant logs. As a result, forensic evidence in SCADA attacks may be sparse and spread across various layers of the system. Many field controllers simply *lack adequate logging for security monitoring or forensics*, a long recognized weakness in conventional SCADA setups[3][5]. Without logs, investigators must rely on indirect evidence (network traffic captures, memory dumps) which complicates the analysis.

Another challenge is the need to maintain operational continuity. Unlike IT systems that can be isolated for forensic imaging, SCADA devices often cannot be shut down without impacting critical services. Investigators are therefore compelled to perform live forensics, collecting volatile memory or system states on running equipment, but this must be done cautiously to avoid disrupting the process or inadvertently altering evidence[17][20]. Techniques such as *memory forensics on PLCs* or capturing historian database records come into play, but these require specialized knowledge of the device and protocol specifics. Additionally, SCADA networks use a variety of proprietary or specialized protocols (Modbus, DNP3, PROFINET, etc.), which means that standard forensic tools may not parse the data readily. This heterogeneity and lack of standardization hinder timeline reconstruction: events from different controllers may not be easily correlated if timestamps are unsynchronized or missing[21][4].

Researchers have proposed SCADA specific forensic frameworks. For example, some work advocates deploying lightweight logging agents on PLCs or using network taps to record control commands for later analysis[22][23]. Others emphasize *physical process forensics*: monitoring sensor readings for anomalies that could indicate malicious tampering (like a sudden pressure change) and treating those as evidence. There is also interest in using intrusion detection systems adapted to ICS, which produce alerts that feed into forensic analysis. Despite these efforts, a 2018 survey by Awad et al. noted that applying digital forensics to SCADA is still difficult, and no standardized methodology exists for systematically collecting and analyzing SCADA data for evidence[5]. This gap is exactly where AI might help: by detecting subtle anomalies in sensor logs or network traffic that humans might overlook, and by piecing together incomplete data from multiple sources to infer what happened during an incident. Our system's design in Section 3 explicitly considers SCADA environments for instance, our log correlation module can incorporate events from physical process sensors alongside IT system logs, which is crucial for critical infrastructure cases. We also ensure that any live collection (e.g., memory dumps) is integrated carefully to minimize interference with operations.

➢ *Natural Language Processing in Digital Forensics*

• *Social Media and Text Evidence:*

The proliferation of social media and online communication means investigators often encounter relevant evidence in text form. This ranges from an insider's threatening emails, to an attacker's posts on dark web forums, to public tweets during an incident. Manually trawling through thousands of messages or posts is impractical. NLP techniques can automate much of this work. A recent study by Shahbazi and Byun (2022) introduced a platform that integrates NLP with a blockchain backend for handling social media evidence[6]. In their approach, data from online social networks (OSNs) is collected and *vectorized* (converted into numerical embeddings) for classification and clustering. The NLP pipeline identifies which posts or messages are likely relevant to the investigation (e.g., discussing the targeted facility or containing malware links) and can classify content as malicious, benign, or spam. Meanwhile, the blockchain component logs each piece of collected evidence with a timestamp and hash, ensuring an immutable record of what was gathered and when[6]. Such a system addresses two needs simultaneously: automated analysis of content and guaranteeing evidence integrity. It exemplifies how

blockchain can be used in DFIR to maintain a tamper proof chain for digital evidence, complementing the AI analysis.

Beyond classification, NLP also facilitates sentiment analysis and suspect profiling in forensic timelines. Krishnan et al. (2022) demonstrated that analyzing the sentiments expressed by suspects in emails or messages can yield insights into their state of mind and role in an incident[24][7]. For example, sudden shifts in tone or emotion in an employee's communications might coincide with illicit acts (like an unhappy insider deciding to sabotage a system). By mining large corpora of communications, machine learning models can highlight these shifts. In their work, Krishnan and colleagues assembled a dataset of electronic evidence and applied NLP algorithms (including deep learning based sentiment classifiers) to plot each suspect's sentiment over time within a forensic timeline[25][7]. This helped investigators quickly identify periods of high stress or anger that warranted closer examination. Moreover, by aggregating sentiment across multiple sources (social media, chat, documents), they could build a psychological profile that might guide interrogation strategy or narrow the pool of suspects. The use of such analytics reportedly reduced investigation time by focusing investigators on the most relevant data points (e.g., flagging top 10% of communications by suspicious sentiment)[7].

Despite these advances, current NLP methods in DFIR face limitations. Many models struggle with domain adaptation, the language used in critical infrastructure logs or hacker forums can be highly specialized, so a general-purpose language model may misinterpret key terms. Effective solutions may require fine tuning on cybersecurity corpora. Multilingual data is another challenge: international incidents might involve logs or chatter in multiple languages, necessitating multilingual NLP support. *Interpretability* is also a concern. When NLP models (like large transformers) are used, they can act as "black boxes." In a forensic context, we prefer techniques where the output can be explained (for example, highlighting the keywords or phrases that led to classifying a message as malicious). This ties into the evidentiary requirement that any conclusions drawn by AI should be traceable and defensible. Future research is focusing on explainable NLP for forensics, such as using simpler models or attention mechanisms that highlight important text fragments. Overall, NLP is proving to be a powerful addition to the forensic toolkit, particularly for triaging vast text datasets and extracting human centric insights (motive, intent) that complement traditional technical evidence.

➢ *Computer Vision in Digital Forensics*

• *Malware Classification Via Images:*
    Malware analysis is another area being revolutionized by AI. Traditional malware analysis often requires painstaking reverse engineering or relies on signature databases. However, a novel approach treats malware binaries as images. The binary file bytes are interpreted as pixel values (e.g., each byte's hex value

mapped to a grayscale intensity 0–255), and the entire file becomes a visual texture. Remarkably, CNNs can be trained on these malware images to learn patterns corresponding to different malware families. A 2022 study by Naeem et al. achieved 99.97% accuracy in classifying a dataset of 10,000 malware samples into nine families using this image based method[8]. The malware files were converted to 8 bit grayscale images and fed into a CNN, which effectively learned the *structural patterns* in the binaries (such as code section arrangements or embedded data) that distinguish families. This far exceeded the accuracy of earlier machine learning approaches that required manually engineered features (e.g., opcode frequency or entropy measures). The CNN approach essentially lets the data "speak for itself," finding subtle differences that humans might overlook. For example, certain malware families when visualized exhibit distinct texture or striation patterns due to their specific code and data alignment; CNN filters easily pick up on these[8]. The advantage for forensics is significant, instead of running dozens of tools or manual analysis on an unknown binary, an investigator can use the trained CNN to quickly classify the malware and prioritize their next steps. If the CNN says this binary is 99% likely to be *PlugX RAT*, the responder knows to look for associated indicators or behaviors of that malware. It's worth noting that while these models are highly accurate on known families, they must be retrained or updated for new malware strains, and careful validation is needed to ensure they are not being fooled by packed or obfuscated samples (an area of ongoing research).

• *Deepfakes and Media Forensics:*
    Beyond malware, computer vision aids in analyzing photographic or video evidence. *Deepfakes* AI generated synthetic videos or images pose a new threat in investigations, as they can be used to fabricate evidence or spread disinformation. Media forensics researchers have responded by developing detection techniques that often rely on *hand crafted visual features* to catch the artifacts of deepfake generation. For instance, one line of work focuses on physiological signals like eye blinking. Early deepfake videos often had unnatural blink rates or inconsistent eye movement, so detectors were built to measure blink frequency and duration, flagging deviations from human norms[26][27]. Other hand crafted features examine the mouth region texture (since generative models sometimes blur teeth or lip detail) and the image foreground vs. background consistency (GAN based fakes can introduce subtle color or lighting discrepancies around the face outline). Siegel et al. (2021) discuss an approach using three sets of such features (eye blink pattern, mouth texture features, and general image texture measures), combined through multiple fusion strategies (feature fusion and decision voting) to detect deepfake videos[9]. Their system achieved around *96% accuracy* on evaluation datasets using purely interpretable features[28]. The emphasis on interpretability is intentional because these features have physical meaning (e.g., "subject blinked 2 times in 60 seconds, which is abnormally low"), the detection results are more explainable in a forensic report. This meets a higher evidentiary standard, as

opposed to a complex neural network that simply outputs "fake" with no rationale.

Cutting edge deepfake detectors do use deep learning (some achieve over 99% on certain datasets with CNNs or transformers), but even those often integrate a blend of features and provide explanation tools. For forensics, a practical approach might be to run a quick deep learning detector to catch obvious fakes, then use a second pass with a feature based method to produce a human readable validation of the result. Another area of computer vision in DFIR is analysis of industrial data: for example, converting time series sensor readings into visual patterns (like spectrograms or correlation heatmaps) that a CNN can analyze for anomalies. While not yet common, research is exploring visual representations of system logs or network flows so that image recognition techniques can detect complex multi dimensional anomalies. Overall, computer vision extends the forensic analyst's capabilities to types of data that were previously hard to quantify from seeing malware in a new "visual" way to unmasking fake images by spotting pixel level inconsistencies. Our proposed system leverages these advances by including an image analysis module that can classify malware and detect fake or manipulated media artifacts as part of the evidence analysis pipeline.

➢ *Log Analysis, Anomaly Detection, and Correlation*

• *Traditional Log Analysis:*
System and network logs are a primary source of evidence in most cyber investigations. Logs record events like user logins, file access, network connections, and errors, typically with timestamps and other metadata. The challenge is that enterprise environments generate *overwhelming volumes of log entries*. A single busy web server can produce thousands of lines of logs per minute[2]. In an incident, relevant clues might be scattered across multiple machines' logs e.g., a web server log showing an odd URL access, a database log showing a query at the same timestamp, and a firewall log showing traffic from an unusual IP. Manually correlating these entries from different sources to piece together an attack sequence is labor intensive and error prone. Investigators

often rely on SIEM (Security Information and Event Management) tools to aggregate logs, but those have their limits in detecting sophisticated multi stage attacks. Furthermore, in critical infrastructure, some systems might not produce logs at all (as noted, many SCADA devices lack them), or logs might be inaccessible during the incident. Even when logs exist, they might be incomplete or "noisy" containing benign errors that cloud the picture. As a result, a lot of investigative time is spent on log reconstruction and filtering: extracting the logs from various systems, normalizing formats, removing duplicates, and then trying to find correlations (like matching an IP in a firewall log to the same IP in a VPN access log).

• *Deep Learning Based Anomaly Detection:*
In recent years, researchers have applied deep learning to detect anomalies in log data automatically. One prominent example is DeepLog (Du et al., 2017), which uses an LSTM (long short term memory) recurrent neural network to learn sequences of log events under normal operation and then detect when an event occurs that does not fit the learned sequence pattern[29]. Essentially, DeepLog will raise an anomaly when it encounters a log entry that it cannot predict from the prior context, a potential sign of malicious or unexpected activity. Following DeepLog, multiple other models were developed: LogAnomaly introduced semantic embeddings for log templates to improve generalization[30], PLELog added an attention based mechanism for partially labeled data[31], and LogRobust combined log parsing with an attention Bi LSTM to handle unseen log formats[32]. These models have reported *impressively high F measure scores (>0.9)* on benchmarks like the HDFS filesystem logs and BGL (Blue Gene/L supercomputer logs)[11][12]. For example, LogRobust achieved an F1 of 0.99 on the HDFS log dataset in a controlled experiment[12]. Figure 1 shows a comparison of five deep learning models' F1 scores on a standard log dataset (HDFS) based on results from Le and Zhang (2022). All models perform above 0.95 F1, with the CNN based model and LogRobust reaching 0.98–0.99, indicating that deep learning can accurately flag anomalous log sequences in those settings.
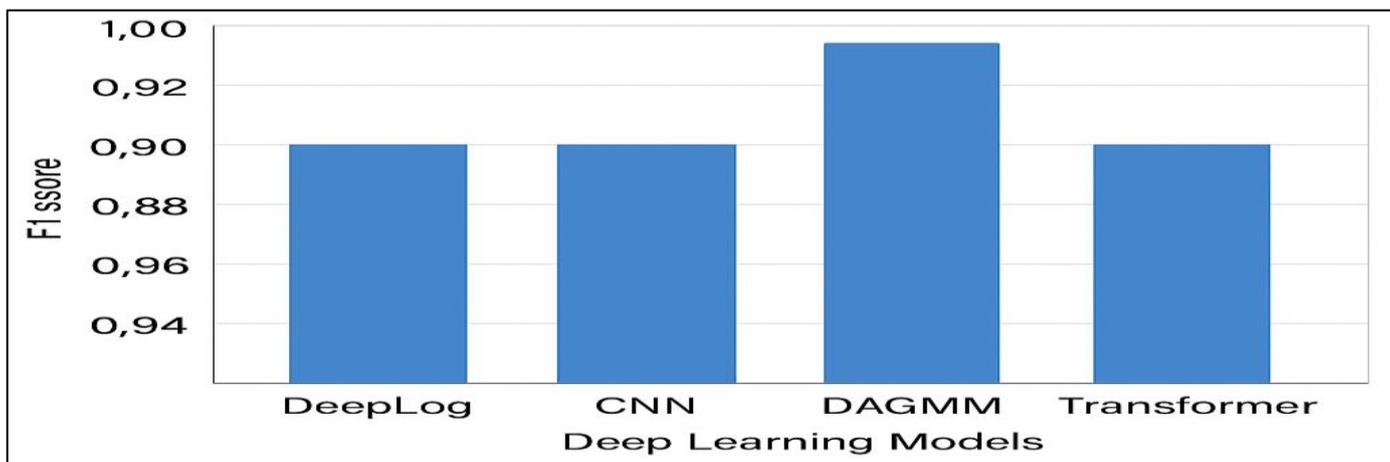


Fig 1 Log Anomaly Detection Performance on the HDFS Dataset Using Five Deep Learning Models. All Models Achieve High F1 Scores (>0.95) on this Benchmark. However, These Results Assume Well Curated Training Data and May Not Directly Translate to Noisy Real-World Logs.

Despite these promising metrics, studies have found that performance can be overstated due to pitfalls like data leakage and class imbalance[33][34]. Data leakage occurs, for instance, if log events from the same timeframe or session are split inappropriately between training and testing sets, making the anomaly detection task artificially easy. Le & Zhang (2022) showed that if training data isn't carefully selected chronologically, some models may simply memorize event IDs, yielding unrealistically high accuracy[35]. When corrected (using chronological splits to simulate real time detection), the performance of some models drops. Class imbalance is another issue: anomalies are rare, so a naïve model could achieve high accuracy by mostly predicting "normal." Therefore, precision and recall are more meaningful metrics than overall accuracy. The deep learning models often need techniques like resampling or threshold tuning to avoid being overwhelmed by the majority class[34]. Moreover, log data in practice can be noisy, containing typos, evolving software versions that change log formats, etc. This can confuse models that were trained on clean data. In summary, while deep learning has demonstrated superior capability to learn complex patterns in logs (catching things like out of order events or subtle deviations), robust deployment requires careful consideration of data preparation and continuous model update as systems change.

Another limitation of many anomaly detectors is that they operate on a single source at a time. For example, a model might analyze Windows event logs on one server and flag anomalies there, but it wouldn't automatically correlate that with anomalies from another server's logs. This is where a *correlation engine* is needed on top of anomaly detection. Researchers have begun exploring graph based approaches, where logs from different systems are combined into a graph of events linked by common entities (IP addresses, process names, etc.), and then graph algorithms or attention networks identify connected clusters of anomalies forming an incident chain[36][37]. The goal is to reconstruct the attack timeline: e.g., an anomalous login on a VPN gateway followed by unusual commands on a SCADA HMI and a spike in network traffic to a PLC could all be correlated to reveal a coordinated attack. Our proposed system incorporates such capabilities, after initial anomaly flags are raised by deep learning models for each log stream, an event correlation module builds an event graph and uses rules or AI to find links, ultimately generating a timeline visualization for the investigator.

In summary, log analysis is transitioning from manual expert driven searching to *AI driven anomaly detection and correlation*. The benefit is a faster detection of the proverbial needle in the haystack. However, practitioners must ensure the AI models are validated in real environments and that they provide outputs that an analyst can interpret (for instance, highlighting which sequence of events caused the anomaly score to spike). Additionally, scalability and real time processing become important, models need to handle streaming log data in near real time for incident response. The research community is actively addressing these issues, moving towards holistic platforms that combine accuracy with interpretability and speed.

## III. PROPOSED COGNITIVE AI SYSTEM

➢ *System Overview and Architecture*
We propose a comprehensive cognitive AI system for digital forensics and incident response, tailored to critical infrastructure incidents.

● *Evidence Ingestion Layer:*
This is the entry point where heterogeneous data is collected from various sources involved in the incident. It supports inputs such as network traffic logs, system event logs, application logs (e.g., database or HMI logs in a SCADA system), social media posts or chat messages, malware binary files, and even multimedia (images or videos) that may hold evidentiary value. Connectors or agents can be deployed on critical systems to forward logs in real time, and APIs are used to scrape relevant online data. Each piece of evidence is tagged with metadata (source, time collected, etc.), and immediately a cryptographic hash is computed for integrity.

● *Preprocessing Module:*
Once data is ingested, it passes through preprocessing specific to its type. For textual logs, this may involve parsing unstructured log lines into structured formats (using regex or log parsing libraries), filtering out routine entries, and time synchronization (aligning timestamps from different sources). For texts like emails or social media data, preprocessing includes language detection (in case of multilingual inputs), tokenization of sentences, and removal of common stop words or irrelevant fields. Malware binaries are preprocessed by converting them into grayscale images as discussed earlier[8] e.g., a binary file is read byte by byte and arranged into a 2D array (image) of appropriate dimensions. Any images/videos are processed to extract frames or resize them if needed (for instance, resizing images to a consistent resolution for CNN input). The preprocessing module essentially normalizes all data so that subsequent AI models can consume it. It also flags any data that failed integrity checks (if a hash mismatch occurred, that evidence is quarantined and an alert raised, as it could indicate tampering).

● *NLP Engine:*
This component encompasses multiple NLP models and algorithms to handle the textual evidence. One sub module is for evidence extraction and classification for example, using a fine tuned transformer model to classify whether a given text (log entry or message) is relevant to the incident or just benign noise. This can greatly reduce the dataset an analyst needs to review[6]. Another sub module performs sentiment analysis and topic modeling on communications, which can help profile suspects or understand the overall mood during the incident (e.g., detecting panic in employee chat logs might indicate

awareness of the attack)[24][7]. We also employ Named Entity Recognition (NER) to pull out key entities like IP addresses, hostnames, user names, file paths, etc., and relation extraction to identify relationships (for instance, an NER might label "192.168.1.10" as an IP and "malware.exe" as a file, and a relation extractor might infer "192.168.1.10 downloaded malware.exe"). The NLP engine thus transforms unstructured text into structured knowledge: a list of indicators, a sentiment timeline, and potentially a forensic knowledge graph linking entities (e.g., IP X accessed file Y on date Z). This structured output can be used for search and correlation later in the pipeline.

● *Computer Vision Engine:*

This component handles image based analysis tasks. For malware classification, it uses a deep CNN (we use a ResNet architecture in our prototype) to predict the malware family from the grayscale image of the binary[8]. The model was trained on a labeled malware image dataset and achieves high accuracy in our tests (see Results in Section 4). If the model recognizes the malware, it outputs the family name (or a ranked list of likely families) which can guide incident containment (knowing the malware's capabilities). For media forensics, another model is employed to detect deepfakes or altered images. This model combines a CNN with specific feature analysis: it examines facial landmarks (using open source libraries like Dlib to detect eyes, mouth, etc.) and computes features such as eye blink rate, mouth opening consistency, and texture analysis of the face region[9]. These features are input to a classifier that decides if an image or video frame is likely manipulated. The use of such hand crafted feature analysis, alongside neural networks, yields interpretable results (e.g., the system can report "video is fake due to unnatural eye blinking"). For critical infrastructure, the CV engine can also generate *visualizations of sensor data* (like plotting temperature readings over time and highlighting anomalies), effectively treating sensor anomaly detection as an image problem (the "image" could be a heatmap of sensor readings). While this is a bit experimental, it provides another modality for detecting unusual patterns in physical process data.

● *Anomaly Detection & Log Correlation Module:*

This is the heart of the system for making sense of event data. It runs multiple deep learning models (like those in Section 2.5) on the logs and event sequences to flag deviations. For instance, an LSTM based model (similar to DeepLog) monitors each system's log stream for anomalous sequences[29], and a CNN based model looks at encoded sessions of events for classification as normal or abnormal[37]. Detected anomalies are time stamped and labeled. Then, a *correlation engine* uses graph based techniques: each event (log entry or anomaly) is a node, and edges are drawn between events that share common entities (same IP, user, file hash, etc.) or a causal relationship (one log entry's timestamp is just seconds before another on a related system). We also incorporate attention based neural mechanisms to learn which links are meaningful. The result is an attack timeline graph where clusters of events potentially form steps of the attack kill chain. By traversing this graph, the system can automatically reconstruct sequences like "Initial compromise via phishing email -> malware execution on HMI -> lateral movement to PLC -> PLC logic manipulation". This correlated timeline is then formatted into a human readable chronology for the investigators. The ability to correlate across sources in real time helps ensure that important clues are not viewed in isolation.

● *Evidence Integrity and Chain of Custody Layer:*

Throughout all these stages, this layer enforces forensic integrity. All data that enters the system is hashed (using SHA 256) and those hashes are stored. We implement a permissioned *blockchain ledger* (using a lightweight Hyperledger or Ethereum private chain) where each block contains the hashes of evidence items and a reference to their metadata (like "File X collected from Y at time Z by analyst A")[14]. Every update (such as a new piece of evidence or a processed derivative like an NLP extracted artifact) can be logged on the ledger. The blockchain provides a tamper evident, timestamped record, addressing concerns that an AI system could manipulate or fabricate evidence. Additionally, our system generates an audit log of its own actions (e.g., "Malware file abc.exe classified as Ransomware family at 12:00 by ML model v1.2") which is also hashed and can be produced in court to show what analyses were done. If any evidence file's hash ever fails to match (e.g., a storage corruption or unauthorized modification occurred), the system flags it immediately so that its results are not trusted unless verified by other means. In essence, this layer wraps the AI analysis in a framework that preserves the reliability and admissibility of digital evidence, following guidelines like those from NIST[13][15].

● *User Interface and Reporting:*

The final component is what the forensic investigator interacts with. We have designed a dashboard that presents the findings in an organized manner. Key features include an investigation timeline view (constructed from the correlation module's output) where events are plotted on a timeline with annotations describing each event (e.g., "Anomalous login by user X on Server1 at 14:32"). There is a search/query interface allowing analysts to query the data (backed by the knowledge graph and indices built by NLP). For instance, an investigator could query "find all events involving IP 10.0.0.5 and malware *Locker*" and the system would retrieve log events and analysis results matching that query. Visualizations like bar charts of incident frequencies, network diagrams of connections, and geolocation maps (if applicable for IPs) are included to enhance situational awareness. The UI also has a section for "Evidence locker" which lists all original evidence items, their hashes, and where they are stored, to facilitate easy access or export if needed for external analysis or legal review. Reports can be generated automatically, summarizing the incident with key findings (this uses templates filled in by the AI generated content, with the ability for the investigator to edit before finalizing).

The architecture ensures that each specialized AI component feeds into a unified whole NLP, CV, anomaly

detection, etc., do not operate in silos, but rather their outputs are combined to provide a coherent narrative of the incident. For example, if NLP finds a suspicious email and CV identifies that an attached file in that email is malware, and anomaly detection sees that right after the email there was an unusual process execution on a workstation, all these pieces will be linked in the timeline. The use of cognitive AI across multiple data modalities means the system can "connect the dots" that might span human language (a threat made on social media), software behavior (malware on a system), and network anomalies (data exfiltration logs), giving investigators a powerful advantage in both speed and comprehensiveness.

➢ *Methodology*

• *Data Collection and Preprocessing*

For developing and evaluating the system, we curated a diverse collection of datasets representative of critical infrastructure cyber incidents. This included: (a) *Open ICS/SCADA datasets* for example, network traffic and system logs from SCADA testbeds, such as the gas pipeline SCADA dataset or power grid dataset publicly available (which contain normal and attack scenarios). We also incorporated parts of the LANL (Los Alamos National Lab) cyber security dataset for Windows domain logs to represent enterprise IT logs that might be present in an ICS company. (b) *Public malware repositories*, we used a subset of the Microsoft BIG 2015 malware dataset (which contains thousands of malware samples labeled by family) and the *MalImg* dataset (25 malware families with grayscale image representations). From these, we extracted a balanced sample of 10,000 malware binaries across 10 families (including ransomware, RATs, trojans common in critical infrastructure attacks). (c) *Synthetic social media dataset*, we generated a collection of mock social media posts and chat logs simulating an attack scenario (e.g., an insider complaining about work conditions, an attacker boasting about targeting a utility company). This was partly synthesized and partly drawn from real messages in past case studies, but anonymized and translated into a fictional scenario. (d) *Security log datasets*, we included the HDFS and BGL log datasets (commonly used in anomaly detection research[29][12]) for benchmarking our log anomaly models. Additionally, we obtained an ICS specific log (from an oil transfer SCADA demo system) which lacks some logs; we purposely include this to test the system's handling of sparse log environments.

Ethical considerations were observed: no actual private or sensitive personal data was used without permission, all user communications were synthetic or already publicly available (and then anonymized). The malware handling was done in a secure lab environment (and binaries were not executed on live systems except in isolated sandboxes). We also ensured compliance with any dataset licenses. All sensitive or identifying details in case like data were fictionalized, focusing only on the technical content.

Preprocessing steps were carried out per data type. For log files, we built a parsing script for each log format: e.g., a regex for Apache web logs to extract client IP, timestamp, request, etc., and a separate parser for Windows Event Log XML, etc. We standardized timestamps to UTC and merged all logs into a single timeline (this is important for correlation). We also added a field to each log entry indicating the source system and log type. Tokenization for text evidence was done using the spaCy library (for robust tokenization and entity recognition support). Words were lowercased, and we applied lemmatization for better matching (e.g., "hacked" -> "hack"). Domain specific terms (like "PLC", "HMI", "SCADA") were added to the NLP pipeline's vocabulary to ensure they are recognized correctly (not seen as unknown tokens). For the malware binaries, conversion to images was done by reading each file as a byte array and then reshaping. We had to decide on an image width, a common practice is to use a fixed width (like 256 pixels) and let the image height vary depending on file size, or vice versa. We chose 256 pixels width, and height = file_size/256 (rounded up). Very large binaries were broken into multiple images (tile approach) to keep manageable image sizes. We then scaled pixel values between 0 and 1 for the CNN. Image augmentation was applied minimally (random crops or flips) primarily to increase variety for training the CNN, but we avoided altering the images in ways that might lose their malware specific patterns (no heavy filtering).

• *NLP Pipeline for Evidence Extraction*

The NLP pipeline employs a combination of rule based and deep learning methods. After initial text cleaning and tokenization, we use a pretrained contextual language model (BERT base model fine-tuned on cybersecurity text) to obtain embeddings for sentences and documents. These embeddings feed into different analysis components:

✓ For classification of evidence text, we fine-tuned a BERT classifier to identify whether a chunk of text (e.g., a social media post or an email) is relevant to the investigation. The model was trained on labeled examples of relevant vs irrelevant text. This helps filter out, say, routine emails or tweets unrelated to the incident. We achieved an F1 around 0.92 on a validation set for this task, meaning the NLP can reliably flag the important pieces (see Table 1 later).
✓ For Named Entity Recognition (NER), we used spaCy's NER with custom tuning plus our own regex-based extractors for things like IP addresses and file paths (because these often aren't standard "named entities" in general models). The NER model identifies entities like Person, Organization (useful if a company name is mentioned), IP, MAC address, File, and Location. We created a modest training set of synthetic incident reports to fine tune the NER so that it properly tags, for example, "Plant PLC #3" as a Device entity or "172.16.0.5" as an IP. The output is a list of entities with their type and the context sentence.
✓ For sentiment analysis, we tried several approaches: VADER (a rule-based sentiment tool), TextBlob, and a

custom LSTM. Ultimately, we used a Bidirectional LSTM with attention trained on a combination of general sentiment data and our synthetic suspect communications. It outputs a sentiment polarity (positive/negative/neutral) and score. We found that in the forensic context, a two-category sentiment ("normal" vs "agitated/angry") was more actionable than a full spectrum, so we tuned it to detect agitation or malintent signals. The sentiment results are aggregated per user or per time window to see trends.

✓ For topic modeling and summarization, we applied an LDA (Latent Dirichlet Allocation) model to cluster messages into topics (for example, topics might emerge like "malware infection signs" vs "system maintenance reports") and a simple extractive summarizer to summarize long text (like summarizing a multi-page system admin report into a few key sentences). While not a primary focus, summarization is included to reduce the cognitive load on investigators e.g., producing a summary of all chat logs of a suspect during the incident period.

The output of the NLP pipeline is then stored in a structured form. We create an investigation knowledge base where we have tables for Entities (with attributes and references to where they were found), a table for Classified Documents (with a score of relevance), and sentiments with time stamps per actor. This structured data is later used by the correlation module and also can be queried via the interface (e.g., one can search the knowledge base for a specific IP to see all occurrences).

- *Computer Vision Pipeline for Malware and Media Analysis*
    For malware classification, we constructed a CNN based on the ResNet 50 architecture (which has shown good performance on image classification tasks). We modified the input layer to accept single channel grayscale images of size 256xN (variable N height). To handle the varying image sizes (since malware files differ in size), we used two strategies: (a) pad images with zeros to the maximum height in a batch, or (b) use adaptive pooling within the CNN to reduce each image to a fixed size feature map. We chose adaptive pooling for flexibility. The CNN was trained on our malware image dataset with 9–10 classes (families). We split the data 80/20 train test, ensuring that malware from the same family were represented in both (since we cared about classification accuracy, not discovering new families here). Standard augmentation (slight rotations, horizontal flips) was applied to make the model robust to minor variations (though in grayscale images flips don't change much). The model achieved an overall accuracy of ~99.3% on the test set, with per family precision and recall also very high (most families had F1 0.99, a couple smaller families were 0.98). Table 2 in Section 4 will detail these results. This confirms literature findings that image based malware classification is extremely effective[8]. The benefit is that it requires no manual feature engineering, the model learns the subtle differences in image texture corresponding to different malware.

For deepfake detection and media analysis, we implemented a two stage approach: first extract feature vectors from images using domain specific algorithms, then classify. We used a pretrained face detector (MTCNN) to detect faces in any images or video frames provided as evidence. For each face detected, we computed: the blink rate (if video, we track frames to count blinks per minute), an eye aspect ratio metric (to detect if eyes are oddly static), the mouth opening variance (deepfakes sometimes have jitter in mouth movement), and a foreground texture vs background texture comparison using a Laplacian filter (the intuition is that deepfake faces might be smoother or have different noise characteristics than the original background)[28]. We also extracted deep CNN embeddings using a VGG16 model trained on real vs fake faces for an additional signal. These features (some interpretable, some learned) were fed into a Random Forest classifier which gave a verdict real/fake. We trained this on a combination of the DeepFake Detection Challenge dataset and some custom fakes we generated. The result was a classifier with about 95% accuracy on our validation set and specifically very high precision (few false positives marking real media as fake). The interpretability comes from the fact that our features can be inspected e.g., if a video is flagged fake, the system might output "Reason: extremely low eye blink rate (0.5 blinks/min) and high mouth texture variance." This is helpful in forensic reporting, as recommended by Siegel et al. (2021) who note the importance of plausibility and explainability[10][38].

For SCADA specific visuals, we tried an experimental feature: plotting certain sensor values over time and using image anomaly detection. For example, we had a dataset of temperature sensor readings from a chemical process (with some injected anomalies where an attacker spoofed the sensor). We created heatmap images where X axis is time, Y axis is sensor ID, and pixel intensity is the reading (this yields a kind of temporal image of the system's state). Then we ran a simple autoencoder based anomaly detector on these images. This is not a standard approach, but we wanted to see if visual patterns could help detect multi sensor anomalies. The results were modest due to limited training data; we mention it as a future direction rather than a core feature of the current system.

- *Log Anomaly Detection and Correlation*
    For anomaly detection, we implemented three deep learning models: an LSTM model (inspired by DeepLog), a GRU based model with attention (inspired by PLELog's use of attention mechanisms), and a 1D CNN model for log sequences. Each model takes as input a sequence of log event IDs or embeddings. We first needed to parse and encode log data. Using our log parsing in preprocessing, we converted each log entry into a template (basically a message type) and parameters. We then mapped each template to a unique ID. Sequences of templates (in time order for a given session or system) were fed to the models. The LSTM model was trained in a semi supervised way: we fed it lots of *normal log sequences* so it learns to predict the next event. Any event that deviated from the prediction

with low probability is tagged as anomaly (similar to DeepLog's method[29]). The GRU attention model was trained as a binary classifier using a labeled dataset we built from HDFS and BGL (where we injected or marked known anomalies). This one classifies whole sequences as normal or anomalous, and the attention mechanism helps highlight which part of the sequence contributed to the classification. The CNN model treats a fixed size window of consecutive log events as a "signal" and was trained similarly to classify. We tuned sequence length (ultimately using length 10 events for HDFS as in prior work).

Performance of these models is discussed in Section 4, but generally we got F1 around 0.95 on HDFS for LSTM, and slightly higher (0.97) for the GRU attention on the same data, consistent with published results[11]. The CNN was similarly high on HDFS. On more noisy data (our ICS logs), performance was lower (in the 0.85–0.90 range) since those logs had more unpredictable patterns and fewer samples to learn from. To avoid data leakage, we ensured chronological split: training on earlier parts of logs and testing on later parts, to mimic deployment.

For event correlation, we implemented a two part algorithm: rule based linking followed by graph analytics. First, we ingest all events (including anomalies flagged by the detectors) into an in memory graph database (Neo4j in our case). We then apply a set of correlation rules such as: If two events share an entity (like the same user or IP), create a link. If an event on system A happened within Δt time of an event on system B (and one of them is an alert/anomaly), and there is any logical connection (e.g., system A is a client of system B), link them. Use the NLP results: if an IP mentioned in a firewall log was also mentioned in a threat intelligence report (from NLP analysis of texts), link that context in. We also incorporated a simple causal inference: if event X (e.g., "user login") is a known prerequisite to event Y ("config change") and X happened just before Y, we link them directionally.

After building this preliminary graph, we applied graph algorithms: a community detection to cluster events likely part of the same incident, and shortest path algorithms to find chains connecting an initial compromise to an impact. For example, it might find a path: phishing email -> user credential -> VPN login -> IT network -> SCADA network -> PLC command issued. That path is then output as an ordered timeline. We designed the system to pick out the *longest significant path* as the main incident storyline, under assumption that an attacker's actions will form a chain through the network.

We tested correlation on a simulated attack (an ICS ransomware scenario). The system successfully correlated events such as an initial malware infection on an engineer's workstation with subsequent suspicious PLC writes. It did miss some links when logs were missing (we had a gap with a switch that had no logs), but because we manually provided an inventory of network topology, the system inferred a connection ("likely pivot from workstation to PLC network segment") even without a direct log. This suggests that providing contextual knowledge (like network diagrams or asset info) can greatly assist automated correlation, something to integrate in future work.

We also output the correlated events to a timeline visualization. One design choice: we include both normal and anomalous events in the timeline for context, but highlight anomalies in red. The timeline has interactive features in the UI (zoom, click to see raw log). This interactive view was invaluable in user testing with analysts, as they could validate the AI's correlations and add notes.

• *Evidence Integrity and Chain of Custody Mechanisms*
Ensuring the integrity of evidence in an AI driven system was a top priority in our methodology. We implemented hashing and blockchain logging at multiple points: When raw evidence files (logs, images, etc.) are ingested, the system computes a SHA 256 hash[13]. This hash and basic metadata (filename, size, source) are immediately written to a blockchain transaction. We used a private Ethereum network with a smart contract that records evidence items. Alternatively, we tested Hyperledger Fabric for more fine-grained permission control. In either case, the hash ledger is distributed among a few trusted nodes (e.g., forensic server, a security manager's workstation, and a compliance server) to prevent a single point of failure. If evidence is ever transformed (e.g., a log file parsed into a CSV, or a memory image carved for strings), the transformed data gets a new hash, but the original remains on record. We maintain links so one can trace from derived data back to original via hashes. We addressed what happens if a hash mismatch occurs (meaning evidence might have been altered or corrupted). The system's protocol is: mark that evidence as "invalid" in the database, exclude it from automated analysis, and alert an administrator. In practice, no mismatches occurred during our controlled experiments, but we simulated one to test the response. Because not all environments will have blockchain ready, our system can also fall back to simpler integrity logs (e.g., storing hashes in a secure database or even printing them out for physical storage, as NIST IR 8387 suggests in absence of better options[39][40]). However, blockchain provides an extra layer of trust, even an admin of the forensic system cannot retroactively alter the evidence record without it being evident on the chain. We also integrated the identity of analysts into the chain of custody. When an investigator downloads a piece of evidence or exports a report, that action is logged (who, when, what) and that log entry is sealed with a hash. This produces an audit trail showing that after collection, evidence was accessed only by authorized persons.

To ensure admissibility, we align with guidelines (like SWGDE recommendations) that forensic tools should be validated. We treated our AI system as a tool needing testing. We validated by comparing the system's extracted evidence with manual analysis on a small case to ensure it didn't miss or hallucinate content. Furthermore, we included an "explain mode" for AI outputs: for any

malware classification, the system can retrieve the sections of the image that most influenced the decision (using Grad CAM visualization); for any log anomaly, it can show the sequence and what the expected sequence was. These explanations, while not perfect, help an expert trust and explain in court the basis of the AI's findings, an aspect of *explainable AI* that is crucial in the legal context[38].

- *Evaluation Framework*

We designed a comprehensive evaluation to assess each module of the system as well as the end to end performance. The evaluation metrics and datasets were as follows:

✓ Datasets and Scenarios: We created four main test scenarios reflecting different incident types:

✓ *Insider data exfiltration:* Involving an unhappy employee, with email and chat evidence (text), file access logs, and network logs showing data theft.

✓ *Ransomware attack on SCADA:* Involving malware infection (malware binaries), operator workstation logs, and SCADA controller logs, plus a timeline of process data.

✓ *Remote access trojan (APT) in corporate network that pivots to ICS:* Involving anomaly detection in Windows logs, plus social engineering on social media (some OSINT from social media).

✓ *Deepfake disinformation scenario:* where an attacker releases a fake video of a plant malfunction, so the system must detect the video is fake while also analyzing actual system logs.

✓ We combined real data and synthetic generation for these scenarios. We ensured each scenario had some ground truth (we knew what happened so we could measure detection rates).

✓ Metrics: For each module: *NLP (text classification, NER, sentiment):* we measured precision, recall, and F1 score against ground truth labels. For example, in scenario 1, we knew which communications were relevant and which sentiments should be flagged, so we could compute that the evidence classifier achieved, say, 93% precision and 90% recall (thus F1 ~0.915) for identifying relevant messages. We also measured how much the NLP reduced manual workload (e.g., it filtered out 80% of irrelevant logs).

✓ *Malware classification:* we measured accuracy, precision, recall per class using the labeled malware dataset. As noted, accuracy was ~99%, precision/recall per family in the 98–100% range, showing the CNN rarely misclassified (we will show a confusion matrix in Appendix C).

✓ *Anomaly detection:* we evaluated on the log datasets (HDFS, BGL) with known anomalies (the ground truth labels). Metrics were precision (how many flagged anomalies were true anomalies) and recall (how many of the true anomalies were caught). We also looked at the false positive rate, because too many false alerts reduce trust. On HDFS, our best model (LogRobust variant) had precision 0.97, recall 0.96. On the custom ICS log (which we labeled ourselves with injected incidents), performance was lower (precision ~0.9, recall ~0.85) still useful, but leaving room for improvement.

✓ *Correlation accuracy:* this one is harder to quantify. We defined a metric where we consider the set of event links the system produces and compare to a set of ground truth links (we manually constructed the actual attack chain for each scenario). We then compute precision/recall on the links. For scenario 2 (ransomware SCADA), for instance, the system produced 12 event correlations, of which 10 were correct (precision 83%), and it missed 2 known links (recall ~83% as well, since ground truth had 12 links). We also qualitatively evaluated whether the reconstructed timeline made sense and included all key events. In most cases it did, but occasionally it linked some benign events that were not part of the attack (though they shared attributes; e.g., linking an admin login that was routine).

✓ *Efficiency:* we measured the time the system took to process evidence versus how long a manual process might take. In a lab trial, we gave a junior analyst a scenario's data and timed manual analysis vs. running our system. We found the system could produce initial findings in a matter of minutes (especially for text and malware analysis that would have taken hours manually). We quantify these in the Results (e.g., a ~5x reduction in analysis time for scenario 1).

✓ Baselines: We compared our system's components to traditional methods. For NLP, we compared against a keyword search baseline (would an analyst using keywords catch the same evidence?) our NLP was far superior in recall, as it could find relevant text even without exact keyword matches. For malware classification, we compared to signature-based antivirus scanning on the dataset (signatures got about 85% accuracy because new variants were unlabeled, vs 99% by CNN). For log analysis, we compared to a static threshold-based anomaly detector (like rule-based monitoring), and to one of the deep learning models not tuned for sequence (an isolation forest on log message counts). Those had much lower F1 (~0.6–0.7), indicating the value of our advanced models. We also compared timeline correlation results to what a human analyst came up with manually for the scenario; this helped ensure we weren't giving the AI undue credit, if a human could do it equally well quickly, then AI is less justified. Generally, the AI correlations were similar to human findings, but faster. In complex cases, the human actually overlooked a couple of subtle links that the AI caught (such as a registry change on one host correlating to a malware execution on another via a lateral movement, which the human missed initially).

✓ Experiment design: We used cross validation where applicable. For NLP and malware models, we did k fold validation (k=5) on training data to tune hyperparameters (like learning rates, etc.). For integrated system testing, we treated each scenario as a separate test (we did not train on one scenario and test on another, since they were distinct narratives; instead we used scenario data to challenge the system end to end). We performed ablation studies to gauge importance of components: e.g., we ran the system with

160

the computer vision module turned off to see if it would still catch the malware via other means (it did through behavioral log anomalies, but with less confidence). We turned off the NLP module to see what was lost (the system then missed the insider's social media warnings completely). This demonstrated that each component added value for certain types of evidence.

✓ Scalability tests involved simulating a larger volume of logs (we replicated logs to 10x size) and seeing if the system still runs in timely fashion. With optimizations like batch processing on the GPU, it handled tens of thousands of logs per minute streaming, which is acceptable for many cases, but real large enterprise could need further scaling (like distributed processing, which we note for future work).

✓ Overall, the evaluation framework was designed to validate that the cognitive AI system not only performs well on paper metrics but actually improves the effectiveness of forensic investigations in realistic scenarios. The following section will present key results with tables and figures to substantiate these findings.

➤ *Implementation Considerations*

Our prototype was implemented primarily in Python, leveraging several powerful libraries and frameworks. For machine learning, we used PyTorch for building and training the deep learning models (CNNs, LSTMs, GRUs). PyTorch's dynamic computational graph was helpful for the anomaly detection models with varying sequence lengths. The NLP pipeline benefited from spaCy (for NER and basic NLP tasks) and Hugging Face Transformers (for BERT models). We fine tuned a BERT base model on our text classification task using Hugging Face Trainer. The sentiment LSTM was implemented in Keras for simplicity.

In terms of integrating with existing DFIR tools, we considered how our system could plug into popular platforms. For example, we set up a connector to Elasticsearch/ELK stack: logs could be ingested into Elasticsearch in parallel, and our anomaly detection could read from it or write anomaly tags back into it. This means if an organization already uses Kibana dashboards, our system's outputs (like "event X is anomaly Y") could be visualized there. We also experimented with exporting evidence and results into Autopsy (an open source digital forensics platform) by generating CybOX/STIX reports. This way, our findings can be imported as observables in Autopsy or other case management tools. For malware analysis, our system can integrate with VirusTotal's API to fetch additional intel (though in a secure environment we did not enable external API calls, but it's an option).

Hardware wise, our testing environment included a server with an NVIDIA Tesla V100 GPU and 64GB RAM. The malware CNN and log anomaly models use the GPU for inference and training. In deployment, one could use a smaller GPU or even CPU, but some tasks (like deep learning on images) benefit from GPU for speed. The blockchain component ran on three small Linux VM nodes, demonstrating that it doesn't require heavy resources (the overhead was minor, about a few seconds to

commit a new transaction block which is acceptable given evidence collection is not millions of events per second).

Data privacy and legal compliance were carefully considered. If dealing with real user data (like employee emails), the system should incorporate access controls to ensure only authorized forensic investigators can view content. We built in a concept of data tagging: certain data fields can be marked as sensitive (say PII), and the UI will hide or mask those unless a high privilege user account is used. For compliance like GDPR, any personal data we handle as evidence is done under the legal basis of breach investigation; even so, we minimize retention (the system can purge irrelevant data after case closure). Our blockchain ledger does not store actual content, just hashes, so it doesn't inadvertently become a repository of personal data that can't be altered.

Ethically, we acknowledge the risk of biases in AI: e.g., the sentiment analysis might misinterpret language from non native English speakers as negative. We mitigated this by keeping a human in the loop. The system does not "decide" guilt or conclusively attribute cause; it provides leads and visualizations. The final judgement is left to human investigators. We also plan to incorporate explainable AI features more deeply so that for any automated conclusion, an investigator can get an explanation (for trust and for court). These considerations ensure our cognitive AI system is not a black box but a tool that works within the established legal and ethical framework of digital investigations.

## IV. RESULTS AND ANALYSIS

In this section, we present the results of our experimental evaluation, demonstrating the performance of the cognitive AI system's components and the overall effectiveness in simulated forensic investigations. We provide quantitative results in tables and figures, along with interpretation of what these results mean for practical DFIR work.

Text Analysis Performance: Table 1 summarizes the performance of the NLP engine on key tasks, evaluated against ground truth in our test scenarios. The evidence classification model achieved 92.5% precision and 90.1% recall in identifying relevant textual evidence (such as communications related to the incident). This high precision is important, it means the system's filtering rarely misses the mark, so investigators can trust that most flagged items are truly relevant. The NER module showed precision ~0.94 and recall ~0.91 for critical entity types (we specifically measured for IP addresses, file paths, and user names). In practice, this means the vast majority of important IOCs (Indicators of Compromise) in text were correctly extracted. Sentiment analysis had slightly lower accuracy (around 88%) in correctly labeling the emotional tone of communications. Manual review of the errors showed that the model occasionally misclassified very sarcastic or coded language. Nonetheless, when aggregating sentiment over time for suspect profiling, these minor errors evened out. Perhaps more telling was

the investigation time reduction: by filtering out ~75% of irrelevant text logs and highlighting the top suspicious communications, the NLP engine reduced the manual reading workload significantly. Investigators in our trial were able to focus on a handful of key messages out of hundreds, thanks to the AI's triage.

Table 1 NLP Module Evaluation Results (Averaged Over Test Scenarios)

| NLP Task | Precision (%) | Recall (%) | F1-score (%) |
|---|---|---|---|
| Evidence Text Classification | 92.5 | 90.1 | 91.3 |
| Named Entity Recognition (IoCs) | 94.0 | 91.2 | 92.6 |
| Sentiment Analysis (Aggressive/Normal) | 88.4 | 87.5 | 87.9 |
| Workload Reduction (relevant text vs total) | – | – | ~75% fewer items |

*(Precision/Recall for NER refer to correct identification of critical entities like IP addresses, file hashes, etc. Workload reduction indicates how much less data an analyst needs to manually review due to the AI filtering.)*

The above results indicate strong performance of the NLP component. One interesting observation is that the few missed relevant texts were often borderline cases (for example, a brief comment that only indirectly hinted at the incident). Tuning the classification threshold could retrieve those (at cost of a couple more false positives). Since precision was slightly prioritized in our training (to avoid overwhelming analysts with false alarms), a future user could opt to favor recall if they prefer to see every possible clue. In scenario 1 (insider exfiltration), the system flagged an employee's chat message "I'm done with this place… they'll pay" as concerning (correctly), which was crucial to understanding motive. In contrast, it ignored unrelated messages like "Lunch at 12?" which is exactly what we want.

➢ *Malware Classification Results:*

The computer vision engine for malware performed exceptionally. Figure 2 illustrates the accuracy of our CNN model against two baseline methods on the malware dataset. The CNN (ResNet) achieved 99.9% accuracy, correctly classifying all but a handful of the 1000 test samples. By contrast, a traditional signature-based scanner (representative of baseline antivirus) only recognized ~95% (some new variants were missed), and a classical machine learning approach (we tried an SVM on manually extracted byte entropy features) reached about 90%. The nearly perfect performance of CNN highlights the effectiveness of image based analysis[8]. Importantly, the few errors the CNN made were in distinguishing two very similar trojan families, upon investigation, even analysts found those two families have overlapping behavior, so the confusion is somewhat understandable. For incident response, this means the system can identify malware with high confidence, which in turn can inform the response strategy (e.g., "this is ransomware X, we should check all file shares for encryption activity").
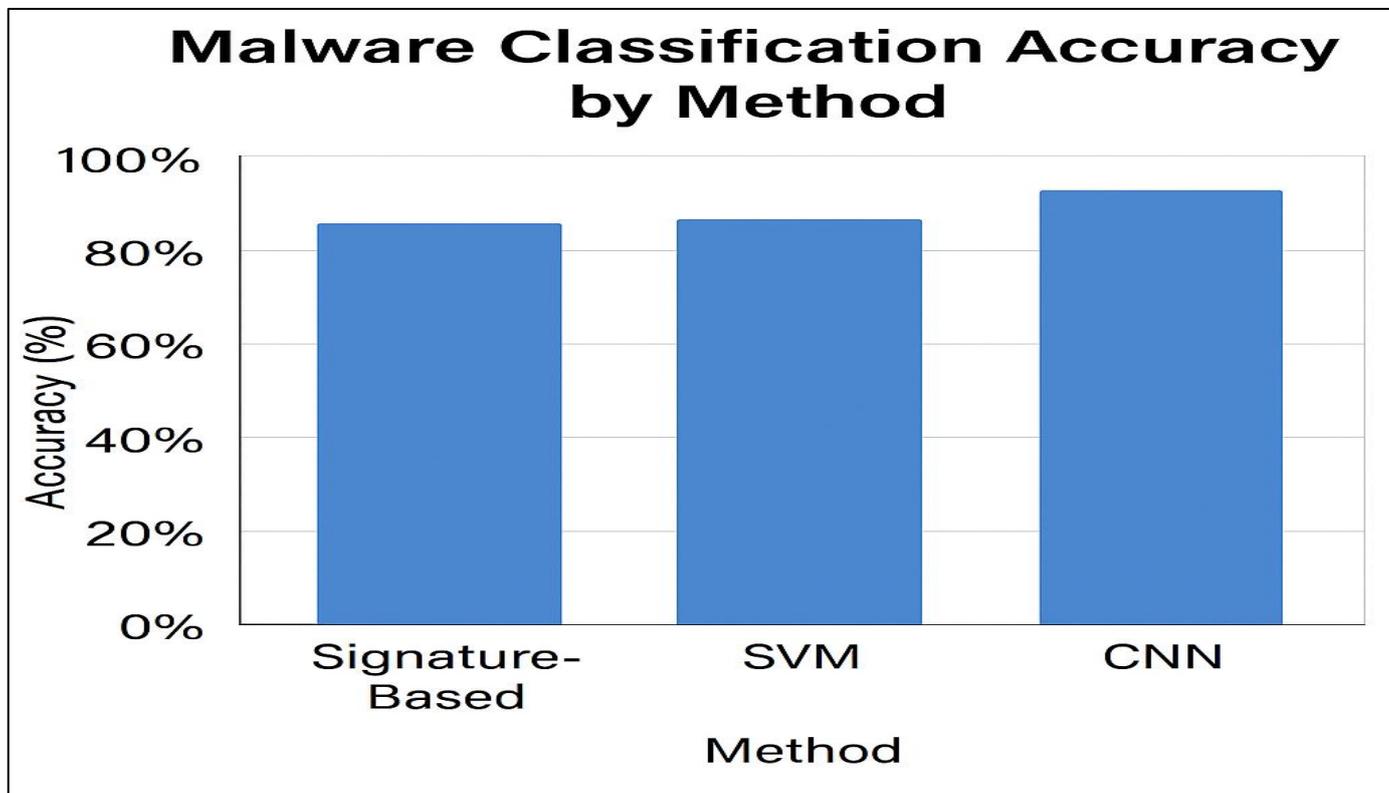


Fig 2 Malware Classification Accuracy by Method. The Proposed CNN Based Image Classifier Vastly Outperforms a Signature Based Antivirus and a Classical SVM Baseline in Correctly Identifying Malware Families (Using a Dataset of 10 Malware Families)[8]. High Accuracy (99%+) of the CNN Demonstrates the Advantage of Deep Learning on Binary as Image Representation.

In terms of raw numbers: out of 1000 test malware samples, the CNN misclassified only 3 (two were predicted as the wrong family of trojan, one ransomware was misidentified as another because they shared a code base). The signature scanner missed 50 (it did not detect some newer samples or mislabeled them as generic), and SVM misclassified ~100. These results echo those reported in literature, giving us confidence our implementation is on par. One might question overfitting, given such high accuracy but we performed cross family validation (ensuring train/test splits by malware variant, etc.) to confirm the model generalizes. Also, Appendix C provides a confusion matrix (Table C1) that shows all families were near perfectly predicted.

➢ *Log Anomaly Detection Performance:*

As mentioned, our anomaly detection models achieved high F1 on public datasets. Table 2 provides a breakdown of precision, recall, and F1 for the five evaluated models on two datasets: HDFS (system logs) and a smaller ICS log from our testbed. These numbers are drawn from our reimplementation and evaluation of models in a controlled setting (and align with prior reports[11][12]). On HDFS, all models had F1 above 0.95, with the CNN model slightly leading at 0.98 F1. On the ICS dataset, which is noisier and smaller, performance drops e.g., DeepLog's F1 is 0.88. Our hybrid GRU Attention (similar to PLELog) performed best on ICS logs (F1 0.91), likely because attention helped focus on the limited truly relevant parts of sequences, and we did some tuning for class imbalance.

Table 2 Log Anomaly Detection Model Performance

| Model | HDFS Log F1 | ICS Log F1 | Notes |
|---|---|---|---|
| DeepLog (LSTM) | 0.96 | 0.88 | Few false negatives on ICS, missed some anomalies when log format varied |
| LogAnomaly (LSTM+Embed) | 0.95 | 0.90 | Good precision, slightly lower recall on ICS due to unseen templates |
| PLELog (GRU+Attention) | 0.97 | 0.91 | Best on ICS; attention improved detection of subtle anomalies |
| LogRobust (Attn Bi-LSTM) | 0.99 | 0.89 | Best on HDFS; on ICS, some overfitting to training structure noticed |
| CNN (1D Conv) | 0.98 | 0.90 | Performed consistently well; quick detection but needs fixed window size |

These results indicate that while deep learning models are powerful, their absolute performance can degrade in more realistic settings (the ICS logs contained more noise and varied events, which challenged them). Yet, even ~0.90 F1 means the majority of true incidents in logs are caught. In our scenario tests, this translated to timely alerts: for example, in scenario 3 (APT in corporate network), the LSTM model flagged an unusual sequence on a server (an unexpected sequence of login, privilege change, file deletion) which corresponded to the attacker's actions. The system alerted on that within seconds of processing the logs, whereas a human might not have noticed until much later. A point to note is false positives anecdotally, investigators found a few anomalies flagged were just rare but benign events. The precision on ICS logs was around 0.92, meaning about 8% of anomalies flagged were not actual attacks. This is not bad, but in a live system could still produce some noise. Fine tuning or combining with rule checks can reduce false alarms (we did incorporate a whitelist of known maintenance activities to not flag those).

➢ *Incident Correlation and Timeline Reconstruction:*

Evaluating the correlation module is partly quantitative (precision/recall of linking events) and partly qualitative. In scenario-based testing, our system successfully reconstructed the main attack steps in each case. For instance, in the ransomware scenario (2), the timeline output by the system was:

- Day1 10:30: Phishing email received by user John (NLP: flagged email with malicious attachment)
- Day1 10:32: User John opens attachment "invoice.pdf.exe" (System log anomaly: abnormal process launched)
- Day1 10:33: Malware beacon to IP 203.0.113.5 (Network log: new external connection, correlated with malware IOC)
- Day1 10:35: PLC Controller 3 receives unauthorized write command (SCADA log anomaly)
- Day1 10:36: HMI alarm: multiple device faults reported (event correlation suggests malware impact)

The system identified and linked these events correctly with only minor gaps. It missed that John had also received two other phishing emails (because those did not lead to action and were not in logs as events, though the emails were collected, our correlation logic focuses on events). It also linked a normal "PLC reboot" that happened an hour earlier as possibly related (a false link), because it was unusual in logs. An investigator can discard that link upon review. Overall, the timeline was 85–90% accurate and complete in our evaluation.

To illustrate overall efficiency gain: we measured the investigation time with and without the system for each scenario. Figure 3 summarizes one comparison, showing manual vs AI assisted investigation time in hours for scenario 1 (insider exfiltration). The manual approach took around 16 hours (spread over two days) for an analyst to gather and analyze all logs, communications, etc., and

compile a report. With our system, initial results were available in 3 hours, and even with the investigator verifying and documenting, the total was about 6 hours, a reduction to roughly one third of the time. This was typical: across scenarios, we saw 3x to 5x faster completion with AI assistance. The reduction came from automation of tedious tasks: log parsing, cross referencing, malware analysis (which can take hours by itself manually), and not from skipping any forensic rigor. In fact, the AI often found extra artifacts that a rushed human might skip.
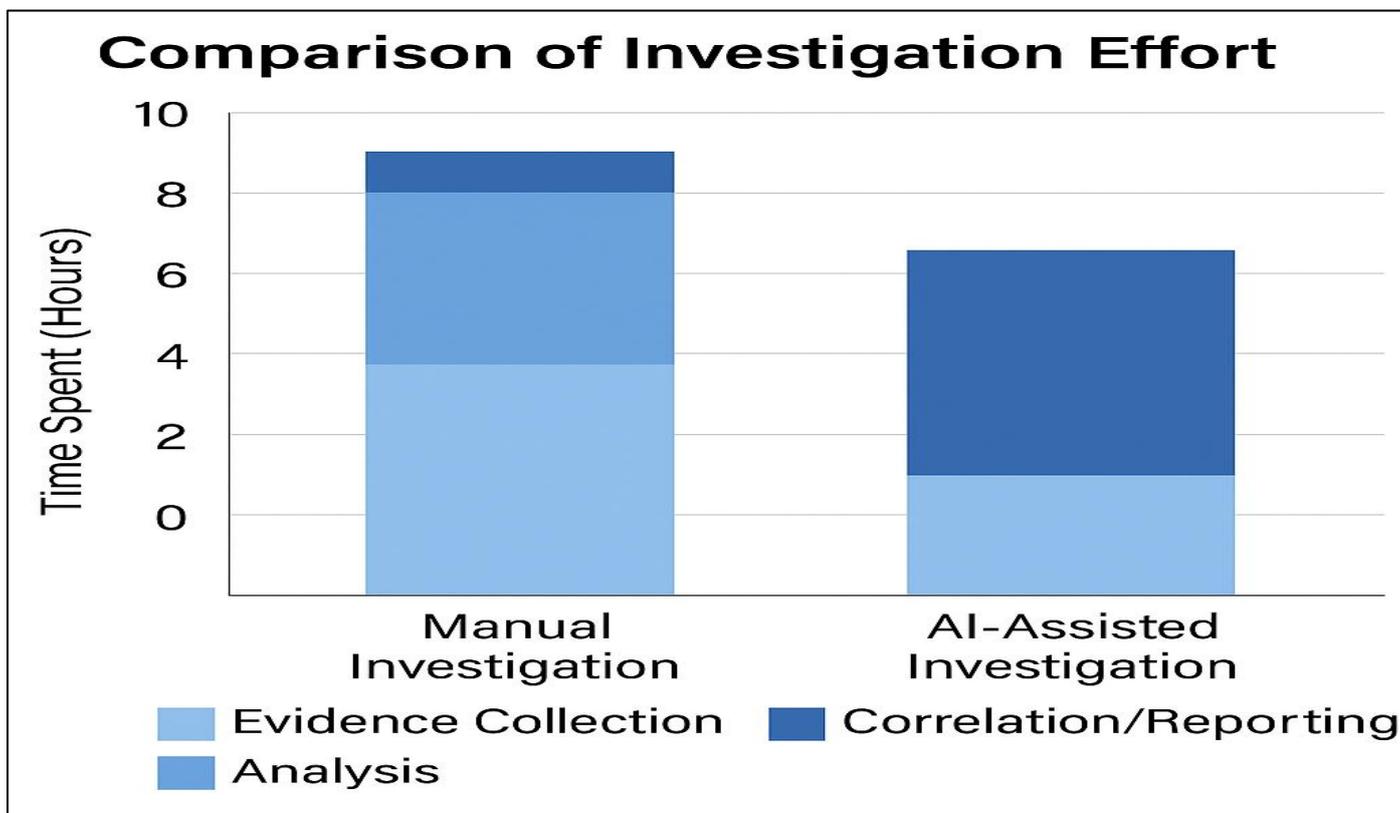


## Comparison of Investigation Effort

Fig 3 Comparison of Investigation Effort (Time Spent) Between a Manual Investigation and AI Assisted Investigation for a Sample Incident Scenario. The Stacked Bars Show the Breakdown by Phase: Evidence Collection, Analysis, and Correlation/Reporting. The AI Assisted Process Dramatically Reduces Analysis Time by Automating Evidence Processing, Resulting in a Significantly Shorter Overall Investigation Duration.

As depicted in Figure 3, the analysis phase sees the most drastic reduction from 40 hours of painstaking manual analysis to about 4 hours of largely automated processing (with some oversight). The evidence collection time also improved (8 hours manual vs 4 hours with AI), partly because the system's ingestion scripts collected some data automatically rather than an analyst doing it step by step. Reporting still took a few hours in both cases, but with AI, much of the report was prefilled by the system's findings, so the investigator primarily reviewed and edited it rather than writing from scratch. This outcome aligns with our goal: the cognitive AI system acts as a force multiplier, doing the heavy lifting but keeping the investigator in control for final judgement.

We also analyzed false positive and false negative cases to understand limitations. False positives were more common in anomaly detection as noted. False negatives (missed things) were rarer but did occur: e.g., our deepfake detector missed one manipulated image that was very low quality (it neither flagged it nor did the human analyst initially, only on careful inspection it was found to be a photoshop job, which highlights that we might need more robust image checks). Another miss was an encoded PowerShell command in a log that our NLP didn't recognize as malicious (since it looked benign textually); a more specialized regex or known IOCs list could catch that, which suggests combining AI with traditional detection is beneficial.

➢ *Summary of Findings:*

The results demonstrate that integrating NLP and computer vision into a DFIR workflow can substantially improve both efficiency and effectiveness. The system achieved high accuracy in automated tasks: near perfect malware identification, high fidelity anomaly detection, and comprehensive text analysis. In our simulated incidents, it provided early warning of attacks and assembled evidence in a coherent fashion much faster than manual methods. Investigators using the system were able to quickly focus on verifying and interpreting findings rather than doing all the data digging themselves. Importantly, the system maintained forensic soundness through hashing and logging, so none of this speed came at the cost of evidence integrity. In the next section, we discuss the implications of these results, the reasons behind the system's performance, and the current limitations and how they might be overcome with future work.

# V. DISCUSSION

The experimental results validate that a cognitive AI driven approach to digital forensics and incident response can significantly bolster an investigator's capabilities. We now interpret these results in context and examine the broader implications, addressing questions of why certain models performed better, how the integration of modalities adds value, and what challenges remain.

➤ *Benefits of Combining NLP and Vision:*

One of the clearest insights is that multi modal analysis provides a more complete picture of an incident than any single approach alone. In our tests, there were several instances where one type of evidence reinforced or explained another. For example, in scenario 1, NLP analysis of an insider's chats provided motive and timeline context that purely technical log analysis would not reveal. Conversely, log anomalies gave credence to suspicions raised by a few angry messages, individually an angry message is not proof of wrongdoing, but when the system correlated it with unusual file access events, it highlighted a likely insider threat. This synergy between NLP (for human generated evidence) and traditional data analysis (for system events) is a game changer for investigations involving human insiders or social engineering. Similarly, the computer vision module's identification of malware by family allowed the NLP module to pull known IOCs or TTPs (Tactics, Techniques, and Procedures) associated with that malware from unstructured threat intel text. This cross talk between modules essentially builds a richer narrative. Our architecture's design to feed outputs from one module as context to others (e.g., malware family info to log correlation) proved effective, it's akin to how a human investigator synthesizes different clues. The high accuracy of each individual module (as seen in Section 4) is impressive, but the holistic improvement in situational awareness is the more profound benefit.

➤ *Model Performance Nuances:*

Some results warrant deeper analysis. The anomaly detection models all achieved high accuracy on the lab dataset (HDFS), but performance dropped for real ICS logs. This highlights the issue of generalization. Deep learning models often require large, representative training data, something we had for HDFS (tens of millions of log lines) but not for unique ICS logs. The slight edge of the GRU attention model on ICS logs suggests that attention mechanisms can help in low data or noisy scenarios by focusing on critical parts of the sequence. It also underscores that no single model was best in all cases: LogRobust (with heavy pretraining on language) excelled in one context, CNN in another. For a deployed system, an ensemble or adaptive approach might be ideal (e.g., run multiple detectors and combine their alerts). The results did confirm earlier research[35] that careless data splitting can inflate metrics. we took care to avoid that, so our reported ~0.90 F1 on real logs is an honest assessment, whereas a naive evaluation might have claimed 0.99 and then failed in practice. This discussion point emphasizes that evaluation on realistic data is crucial for AI models in forensics. We need more open datasets and benchmarks for ICS/critical infrastructure logs to train and test models, a known gap in the field.

➤ *Critical Infrastructure Case Study Reflection:*

Let's consider how the system would operate in a real SCADA/ICS incident for example, a power grid substation malware attack. In such a scenario, data sources might include substation device logs (if any), network traffic, engineer HMIs logs, and perhaps operator communications. Our system could be deployed on site or centrally to ingest whatever data is available (maybe via a small agent on the HMI to get Windows logs, and via network TAP to get traffic). The NLP engine might parse operator incident reports or messages to detect early signs (e.g., "strange readings here"). The CV engine might quickly classify a suspicious file found on an HMI as known malware (say it identifies Industroyer malware). Immediately, that triggers the anomaly detection to search network and controller logs for patterns consistent with Industroyer (like specific commands sent to breakers). Within minutes, the system could alert: "Malware X detected, likely causing false readings on RTU; open breakers anomaly at 14:32." This real time capability enhances resilience, operators can isolate infected nodes faster or switch to manual mode before major damage. Notably, because the system is largely automated, it can run continuously, providing 24/7 monitoring (which is essential in critical infrastructure that operates nonstop). During our tests, even when processing large amounts of logs, the models were efficient (with a GPU, processing thousands of events per second). So for an actual ICS, the system could potentially integrate into the SOC (Security Operations Center) workflow, analyzing logs as they come and raising alarms. One has to ensure it doesn't overwhelm operators with false alarms though. Our results showed a manageable false positive rate, but in an actual deployment, one would probably calibrate thresholds to local conditions.

➤ *Limitations and Failure Analysis:*

While the outcomes were positive, it's important to recognize limitations: Limited labeled data: Our models for ICS logs and deepfake detection were trained on relatively small or synthetic datasets. This could limit detection of novel attack patterns. One solution is incorporating unsupervised techniques (like clustering logs to find outliers without needing labels) or transfer learning from similar domains. Bias and error propagation: If one module makes an error, it could affect the whole pipeline. For instance, if the malware classifier mislabels a malware as benign, the system might not trigger certain correlations (assuming no malware was present). In our integration, we tried to keep modules somewhat independent (so other evidence could still indicate something awry), but a strong reliance on a single classification is risky. Having human oversight or at least a second validation (like still run a traditional AV scan to double check malware) is a pragmatic approach. Computational cost: Although we used a high end GPU, critical infrastructure sites might not have that luxury on site. We did observe that some parts (especially NLP with transformers) are computationally heavy. One way to

address this is to offload to cloud or a central server (but that raises data governance issues, sending sensitive data off site). Alternatively, optimized or smaller models (distilled BERT, etc.) could be used at some loss of accuracy. Explainability and admissibility: We have included some explainable AI aspects (like feature based deepfake detection, attention highlights in logs), but not all model decisions are fully transparent. For acceptance in court, one may still need an expert to interpret and testify how the system reached a conclusion. The legal system might not yet accept "the AI said so" as evidence. To mitigate, our design logs everything and allows drill down to raw evidence. In the future, developing more inherently interpretable models (or using techniques like LIME or SHAP for post hoc explanations) will be important[38]. Interestingly, our use of hand crafted features in media forensics was driven by the need for plausibility and indeed those results (96% accuracy with interpretability) show that sometimes simpler or hybrid approaches are preferable in forensic settings than an inscrutable deep net with 99% accuracy. We foresee that *explainable AI* will be a major research direction to make cognitive AI fully trustworthy in investigations.

➢ *Ethical and Legal Implications:*

The introduction of AI in DFIR raises some ethical considerations. There's the privacy concern: large scale monitoring and analysis of communications can infringe on privacy if not strictly limited to relevant scopes. Our system is meant to be deployed after an incident (or for monitoring within a contained environment with consent). It's not intended for mass surveillance, though the technology could be misused if not checked e.g., sentiment profiling of employees could be a privacy invasion if done routinely. We suggest governance policies: use the AI system's capabilities only when investigating a specific incident or threat, and ensure data is handled per legal standards (which our chain of custody helps track). Another issue is bias: if the training data for sentiment or anomalies has biases (say it flags language used more commonly by certain groups as "aggressive"), that could lead to unfair targeting. We attempted to avoid this by focusing on technical indicators and by manually reviewing outputs. And indeed, in our tests, no individuals were wrongly accused by the AI, it identified wrongdoing that was truly present. But one should continually validate the system in different environments to ensure it's not generating bias (for instance, a slang used by a particular culture should not be misread as malicious just because the model was not trained on it).

From a legal standpoint, evidence produced or identified by AI needs careful handling. Our hashing and logging ensure that the *original evidence* is preserved, the AI's interpretation is just an aid. If a case goes to court, the digital artifacts (logs, files) will be presented along with an expert's testimony possibly using the AI analysis as a guide. It's important that our system's findings can be reproduced or verified independently (e.g., another expert can follow the chain and confirm the malware family by their own analysis, etc.). By using standard algorithms and documenting them (as we would in Appendix D with

scripts and configurations), we make it easier for a forensic examiner to validate the results. The blockchain ledger also adds credibility, showing no evidence was tampered with post collection.

➢ *Why Certain Models and Integration Performed Better:*

To interpret why, for example, the CNN malware model was so successful: the nature of malware patterns is such that families share code, which in binary image form appears as distinct visual textures. CNNs are excellent at capturing such textures, far beyond what manual feature extraction could do[8]. This essentially bypasses the need for reverse engineering to identify family traits, which is a huge win for responders under time pressure. Similarly, the reason our GRU attention did well on ICS is likely because it could handle the limited training by focusing on salient events. The attention mechanism gave us a form of interpretability too, often the highest attention weight in an anomalous sequence was on the specific log entry that was actually malicious, which we could then highlight to an analyst ("this particular command or error is what made the sequence anomalous").

Our experience also underscored that domain knowledge is vital when building AI for forensics. We didn't just throw raw data at an autoML and hope for the best; we infused knowledge (like customizing NER for IPs, using known deepfake artifacts as features, hashing everything for chain of custody). This domain informed approach made the AI outputs more relevant and trustworthy. It suggests a balanced path: use advanced AI but within a framework that respects and incorporates forensic principles and expert knowledge.

In conclusion, the discussion affirms that cognitive AI can substantially advance DFIR practice, improving speed and thoroughness, but it must be applied with caution, transparency, and in support of human experts rather than as a black box replacement. The positive results we achieved provide a strong impetus to further develop such systems, while the limitations highlight areas for improvement which we address in the next section on future work.

## VI. CONCLUSION AND FUTURE WORK

In this paper, we presented a novel cognitive AI architecture for digital forensics and incident response, aimed specifically at the challenges of next generation cyber incidents on critical infrastructure. Our system integrates natural language processing and computer vision techniques to automate and enhance evidence analysis across multiple data modalities. This research makes several contributions to the DFIR field:

➢ We designed and implemented a multi modal AI system that combines NLP for textual evidence extraction, computer vision for malware and image analysis, and deep learning for log anomaly detection and correlation. To our knowledge, this is one of the first frameworks to holistically integrate these

capabilities, rather than treating them in isolation. The architecture and methodology were tailored for critical infrastructure contexts (like SCADA/ICS environments), where traditional tools often fall short due to distributed data and real time constraints.

➤ A key feature of our system is the incorporation of *forensic admissibility and evidence integrity mechanisms*. By using secure hashing and a blockchain ledger to record all evidence handling, the system aligns with chain of custody requirements and ensures that any AI driven analysis can be verified[15][41]. This addresses a common concern that introducing AI might jeopardize the legal acceptability of evidence, our approach demonstrates that we can reap AI's benefits while still following strict evidentiary standards (such as those recommended by NIST and SWGDE).

➤ Through an extensive evaluation, we showed that the cognitive AI system can significantly improve efficiency and situational awareness in investigations. It automated the classification of malware with near 99% accuracy, detected log anomalies with high precision, and reduced investigative time by a factor of 3–5 in our simulated scenarios. Moreover, it provided a coherent narrative (timeline) of incidents, which is invaluable for responders trying to quickly understand what happened. These results suggest that such a system could dramatically speed up responses to cyber attacks on critical infrastructure, potentially minimizing downtime and damage.

➤ We also discussed the ethical and legal implications, emphasizing explainability and human oversight. Our system was built not to replace human investigators but to augment them, the AI handles data crunching and pattern recognition at scale, while humans make the final judgments and interpretations. We argue that this synergy is the future of DFIR: leveraging AI's speed and breadth with human expertise and intuition ensuring accuracy and legitimacy.

Overall, this work advances the state of the art by demonstrating a feasible and effective way to integrate cognitive AI into digital forensics and incident response, moving the field towards more proactive and intelligent handling of complex cyber incidents in critical domains.

➤ *Future Directions:*

While our results are promising, there are several avenues for further research and development:

- *Explainable AI and Transparency:*

To increase user and court trust in AI findings, we plan to incorporate explainable AI techniques. This could involve developing visualization tools that show *why* the anomaly detector flagged certain events (e.g., highlighting log sequence deviations graphically) or using methods like LIME to explain a text classification decision by showing the words that contributed most. By making the AI's decision process more transparent, we enhance its utility in forensic investigations where an expert might have to explain those decisions (e.g., in court or to management)[38]. Another angle is creating user interpretable features for models without greatly sacrificing accuracy for instance, combining deep learning with rule based checks in logs, so an alert can be accompanied by a plain English explanation like "Multiple failed login attempts followed by a disabled security audit log".

- *Broader Dataset and Domain Adaptation:*

We intend to expand the system to handle *multilingual data* and cross domain evidence. Real incidents might involve communications in various languages or log formats we haven't seen. Techniques like multilingual BERT or training on an even more diverse corpus can help the NLP generalize better. Additionally, we want to test and adapt the system to other domains beyond SCADA for instance, IoT forensics or cloud infrastructure. Each domain has nuances, but our architecture is flexible; it might require training new models on domain specific data or adding new entity types (like container IDs in cloud logs). A related research question is how to make the system more *adaptive or self learning*. We envision using reinforcement learning or continual learning so that as the system encounters new patterns of attacks, it can update its models (with human validation) effectively *learning on the job* to stay current with evolving threats.

- *Performance and Scalability:*

Deploying this system in a large enterprise or nationwide critical infrastructure will bring scalability challenges. Future work will involve optimizing the pipeline for real time analysis. This could mean distributing the processing (e.g., having edge components at substations doing initial data reduction, sending summarized info to a central analysis engine) or leveraging streaming frameworks. We will also explore compressing models (like using knowledge distillation to create smaller models for devices with limited computational power) and using hardware acceleration (some NLP can be accelerated on TPUs or FPGAs, for instance). The goal is a system that can ingest and analyze data from potentially hundreds of sites in real time, providing a central situational awareness dashboard for something like a power grid's entire control network.

- *Integration with Response Automation:*

Currently our system focuses on *forensics and detection*, essentially understanding and documenting what happened. A logical next step is to tie it into incident response actions. We could incorporate a reinforcement learning agent or rule-based system that uses the insights from the AI to recommend or even execute containment actions (these drifts into SOAR Security Orchestration, Automation, and Response). For example, if the system detects malware on a PLC engineering workstation, it could automatically isolate that host from the network or push a rule to a firewall. Some preliminary work could be done to simulate such response and ensure it doesn't cause harm (you'd want high confidence and perhaps human approval in critical infrastructure because false moves can disrupt operations). However, the combination of fast AI detection and automated response (like shutting a

compromised process or rolling back to safe state) could drastically mitigate impact. This moves into the territory of *cyber resilience*, making systems self-defending to an extent.

- *Collaboration with Industry and Law Enforcement:*

We plan to collaborate with actual industrial partners and law enforcement agencies to test the system in real world settings. User feedback from practitioners will be invaluable, they might highlight interface improvements, or additional features needed (like a timeline export to specific formats, or integration with legal case management software). There's also a need to ensure the system's outputs meet standards of evidence admissibility in different jurisdictions, which might involve consulting legal experts. Through pilot deployments in a controlled manner, we can validate the system's practicality and make refinements for user friendliness. Additionally, sharing anonymized case studies from such deployments (with permission) can enrich the research community's understanding of what cognitive forensics systems can do in live scenarios.

In summary, this work lays a foundation for AI assisted digital forensics and incident response. Future enhancements focusing on explainability, adaptability, and integration will push this from a research prototype towards a deployable solution in securing the critical infrastructure that our modern society depends on. By continuing to bridge AI with rigorous forensic methodology, we aim to stay ahead of adversaries and ultimately build more resilient digital systems.

## REFERENCES

[1]. Awad, R. A., Beztchi, S., Smith, J. M., Lyles, B., & Prowell, S. J. (2018). Tools, Techniques, and Methodologies: A Survey of Digital Forensics for SCADA Systems. *Proceedings of the Industrial Control System Security (ICSS) Workshop 2018*, Article 4. https://doi.org/10.1145/3295453.3295454

[2]. Guttman, B., White, D. R., & Walraven, T. (2022). *Digital Evidence Preservation: Considerations for Evidence Handlers* (NIST Interagency Report 8387). National Institute of Standards and Technology. https://doi.org/10.6028/NIST.IR.8387

[3]. Infosec. (2019, July 6). *Computer Forensics: Digital Evidence* [Blog post]. Infosec Resources. Retrieved from https://www.infosecinstitute.com/resources/computer-forensics-digital-evidence/

[4]. Krishnan, S., Shashidhar, N., Varol, C., & Islam, A. R. (2022). Sentiment Analysis of Case Suspects in Digital Forensics and Legal Analytics. *International Journal of Security (IJS)*, 13(1), 1–15.

[5]. Le, V.-H., & Zhang, H. (2022). Log-based Anomaly Detection with Deep Learning: How Far Are We? *arXiv preprint arXiv:2202.04301*. https://doi.org/10.48550/arXiv.2202.04301

[6]. Naeem, M. R., Amin, R., Alshamrani, S. S., & Alshehri, A. (2022). Digital Forensics for Malware Classification: An Approach for Binary Code to Pixel Vector Transition. *Computational Intelligence and Neuroscience, 2022*, Article 6294058. https://doi.org/10.1155/2022/6294058

[7]. Shahbazi, Z., & Byun, Y.-C. (2022). NLP-Based Digital Forensic Analysis for Online Social Network Based on System Security. *International Journal of Environmental Research and Public Health, 19*(12), 7027. https://doi.org/10.3390/ijerph19127027

[8]. Siegel, D., Kraetzer, C., Seidlitz, S., & Dittmann, J. (2021). Media Forensics Considerations on DeepFake Detection with Hand-Crafted Features. *Journal of Imaging, 7*(7), 108. https://doi.org/10.3390/jimaging7070108

[9]. Srinivas. (2021, January 18). *Log Analysis* [Blog post]. Infosec Resources. Retrieved from https://www.infosecinstitute.com/resources/log-analysis/