

# Towards a standardized methodology and dataset for evaluating LLM-based digital forensic timeline analysis

By: Hudan Studiawan, Frank Breitinger, Mark Scanlon

From the proceedings of
The Digital Forensic Research Conference **DFRWS APAC 2025**Nov 10-12, 2025

**DFRWS** is dedicated to the sharing of knowledge and ideas about digital forensics research. Ever since it organized the first open workshop devoted to digital forensics in 2001, DFRWS continues to bring academics and practitioners together in an informal environment.

As a non-profit, volunteer organization, DFRWS sponsors technical working groups, annual conferences and challenges to help drive the direction of research and development.

https://dfrws.org

FISEVIER

Contents lists available at ScienceDirect

# Forensic Science International: Digital Investigation

journal homepage: www.elsevier.com/locate/fsidi



DFRWS APAC 2025 - Selected Papers from the 5th Annual Digital Forensics Research Conference APAC



# Towards a standardized methodology and dataset for evaluating LLM-based digital forensic timeline analysis

Hudan Studiawan a,\*, Frank Breitinger b, Mark Scanlon c

- <sup>a</sup> Department of Informatics, Institut Teknologi Sepuluh Nopember, Indonesia
- <sup>b</sup> Chair for Cybersecurity, University of Augsburg, Augsburg, Germany
- <sup>c</sup> Forensics and Security Research Group, School of Computer Science, University College Dublin, Ireland

# ARTICLE INFO

# Keywords: LLM evaluation Forensic timeline analysis Large language models ChatGPT log2timeline/plaso

#### ABSTRACT

Large language models (LLMs) have widespread adoption in many domains, including digital forensics. While prior research has largely centered on case studies and examples demonstrating how LLMs can assist forensic investigations, deeper explorations remain limited, i.e., a standardized approach for precise performance evaluations is lacking. Inspired by the NIST Computer Forensic Tool Testing Program, this paper proposes a standardized methodology to quantitatively evaluate the application of LLMs for digital forensic tasks, specifically in timeline analysis. The paper describes the components of the methodology, including the dataset, timeline generation, and ground truth development. In addition, the paper recommends the use of BLEU and ROUGE metrics for the quantitative evaluation of LLMs through case studies or tasks involving timeline analysis. Experimental results using ChatGPT demonstrate that the proposed methodology can effectively evaluate LLM-based forensic timeline analysis. Finally, we discuss the limitations of applying LLMs to forensic timeline analysis.

## 1. Introduction

Forensic investigations often require the reconstruction of a timeline of events and activities related to a digital device or users (Hargreaves and Patterson, 2012). Such timelines can provide valuable insights into various criminal activities, including malware, brute-force attacks, or attacker post-exploitation activities. The timeline analysis process is complex and time-consuming, particularly when dealing with large amounts of digital data from multiple sources (Breitinger et al., 2025). Traditional methods for timeline analysis are based on manual analysis, which can be subjective and prone to errors, and can lead to missing critical information (Studiawan et al., 2020).

The development of large language models (LLMs), such as OpenAI's GPT-3 (Brown et al., 2020) has opened up many possibilities, including in digital forensic research. The model has been implemented in the ChatGPT application and instantly gained many users (Buchholz, 2023). LLMs performed remarkably in various natural language processing tasks, including language generation and question-answering. Leveraging these capabilities, natural language processing techniques can be applied to digital data sources to analyze temporal information and investigate timelines of events. Other studies also suggest that

artificial intelligence should provide more assistance in forensic investigation (Hall et al., 2022; Studiawan et al., 2019).

An editorial article by Scanlon et al. (2023b) discusses the increasing demand for expert digital forensic analysts and the potential use of LLMs such as ChatGPT in this domain. They emphasize the importance of maintaining the "AI-assisted investigation" and "human-in-the-loop" mantras when using LLMs in digital forensics. The article suggests that LLMs could lead to a new career specialization of digital forensic prompt engineers. Wickramasekara et al. (2025) provides a comprehensive overview of where LLMs may assist digital forensics.

Various studies explored the application of LLMs in digital forensics. For instance, Scanlon et al. (2023a) assessed ChatGPT's impact on tasks such as understanding artifacts, evidence searching, and anomaly detection. Although ChatGPT shows promise in several low-risk forensic applications, concerns arise about evidence security and the model's occasional inaccuracies. Experts must exercise caution and have a deep understanding of the subject to effectively use ChatGPT in forensic scenarios. Furthermore, ChatGPT has been explored for digital evidence investigations (Henseler and van Beek, 2023), virtual forensic assistants (Dinis-Oliveira and Azevedo, 2023), and report writing (Michelet and Breitinger, 2024). Based on our literature review, existing work in this

E-mail addresses: hudan@its.ac.id (H. Studiawan), frank.breitinger@uni-a.de (F. Breitinger), mark.scanlon@ucd.ie (M. Scanlon).

https://doi.org/10.1016/j.fsidi.2025.301982

<sup>\*</sup> Corresponding author.

area has not discussed standardized evaluation for LLM-based digital investigation.

# 1.1. Contribution

The contributions of this paper are as follows:

- This paper proposes a standardized methodology to quantitatively evaluate the performance of LLMs in forensic timeline analysis tasks, such as event summarization.
- 2. This study presents a case study of forensic timeline analysis using LLM, e.g., ChatGPT.
- We created forensic timeline datasets and ground truth from Windows 11 using Plaso and these are publicly available on Zenono<sup>1</sup> for research and education purposes.

The proposed methodology is not fully autonomous. Instead, it functions as an LLM-assisted framework that depends on structured prompts, external libraries, and curated inputs. The current implementation emphasizes feasibility and standardization, rather than full automation.

The remainder of the paper is organized as follows: Sec. 2 provides related research. Sec. 3 describes the proposed approach for standard methodology and quantitative evaluation for LLM-based timeline analysis. Sec. 4 presents the case study that demonstrates the application of the proposed method and a discussion of the results. Finally, Sec. 5 concludes this study.

#### 2. Related work

# 2.1. Forensic tool testing and validation

To effectively validate digital forensic tools and methods, a proper validation test plan should include laboratory use in the real world, controlled internal tests based on scientific principles, and peer review. Brunty (2023) provides an overview of the foundational scientific aspects of forensic validations and describes the recommended steps to conduct a forensically sound validation method.

The Computer Forensics Tool Testing (CFTT) Program at NIST aims to establish a methodology for testing computer forensic tools, including developing specifications, test procedures, and criteria (NIST, 2019). The program helps to ensure the reliability of forensic software tools, helping tool makers, users, and interested parties. CFTT methodology involves breaking down forensic tasks into discrete functions and creating test methodologies for each.

Hughes and Karabiyik (2020) discuss the need for rigorous validation practices in digital forensics to establish accuracy and reliability. They highlight challenges in developing statistical confidence for forensic tools, such as the lack of reference data, validation methods, and precise definitions of measurement. The authors propose a method for generating data procedures, virtual machine-based validation, and empirical models to guide the analysis.

Another study discusses the challenges of scientifically validating digital forensic evidence (Arshad et al., 2018). The authors emphasize the lack of standard datasets, formal testing procedures, and established error rates. Horsman (2019) examines the challenges of ensuring reliability in digital forensic tools. The paper discusses the lack of standardized validation methods and the issues of transparency from software vendors. A survey of practitioners reveals widespread concerns about tool reliability and a need for improved testing standards and error rate reporting.

The related study on tool testing and validation shows a research gap where we need a method to evaluate and validate LLMs as tools in digital

forensics. This paper aims to fulfill this need specifically for the task of forensic timeline analysis.

# 2.2. Forensic timeline analysis

Forensic timeline analysis involves reconstructing the sequence of events and activities related to a user or a system. Therefore, a variety of artifacts, such as browsing history, log files, or file metadata, are being parsed, and relevant information is extracted (Palmbach and Breitinger, 2020). The analysis of the timeline is then conducted using tools and data visualization techniques (Inglot and Liu, 2014). If tools do not yield expected results, a manual examination of data sources may be required. However, this approach can be time-consuming, labor-intensive, and prone to errors (Breitinger et al., 2025).

Timeline generation tools, such as log2timeline/Plaso, Autopsy, and Magnet AXIOM, can automate the timeline analysis process to some extent by extracting relevant temporal information from digital data sources. However, these tools are limited by the quality of the extracted data and may not be able to capture all relevant events and activities from acquired artifacts (Studiawan et al., 2022a).

The approach by Hargreaves and Patterson (2012) can automatically reconstruct or summarize high-level events from low-level events. Previous techniques focus on extracting times from a disk image into a timeline, which can produce several million "low-level" events (e.g., file modification or Registry key update) for a single disk. In contrast, this approach can automatically reconstruct high-level events (e.g., connection of a USB stick) from this set of low-level events. The knowledge representation model presented in Chabot et al. (2014) allows a semantically rich representation of events related to the incident. It includes the identification of correlated events that can highlight valuable information for investigators.

The construction of a timeline array using time information from web browser log files is one way to perform forensic timeline analysis (Nalawade et al., 2016). Different data types of timelines can be constructed from web browser artifacts such as web history, cache, cookies, download history, and search term timelines. Furthermore, Bhandari and Jusas (2020) propose an abstraction-based approach to reconstruct a timeline of events and artifacts. The method enhances the relevance of the timeline by reconstructing it into four levels of depth, from general to specific, to reduce complexity and extract information.

The use of deep learning techniques, e.g., autoencoders, improves anomaly detection in a forensic timeline by establishing a baseline for normal activities (Studiawan and Sohel, 2021). Another tool, namely Drone Timeline, constructs a timeline from a drone device and considers time extracted not only from file metadata, but also from various source artifacts of a drone or its control devices (Studiawan et al., 2022b).

# 2.3. LLMs for digital forensics

In the case of LLM application for digital forensics, Henseler and van Beek (2023) discuss how ChatGPT can assist investigators by writing structured queries, summarizing and evaluating large volumes of communication data, and analyzing search results. The authors highlight that ChatGPT can transform natural language queries into structured formats, summarize, and visualize chat logs to reveal key relationships. The study notes limitations, such as hallucinations and the need for expert guidance. Another work explores the potential of using LLMs, e.g., ChatGPT and Llama, to assist in the generation of forensic reports in digital investigations (Michelet and Breitinger, 2024). The authors assess the ability of LLMs to automate parts of the report writing process, focusing on sections such as the introduction, items received, methodology, and results. They found that while ChatGPT performs well and generates relatively accurate drafts, Llama struggles with accuracy and completeness. The results show that LLM outputs still require proofreading and corrections.

Dinis-Oliveira and Azevedo (2023) also explore the potential and

<sup>1</sup> https://zenodo.org/records/15493424.

challenges of using ChatGPT in forensic sciences. The authors highlight the advantages of ChatGPT, such as assisting forensic professionals in drafting reports, analyzing forensic data, performing literature searches, and serving as a virtual forensic assistant. However, the paper also raises concerns about the ethical and legal challenges associated with using AI in this field, such as credibility issues, inaccuracies, plagiarism, and the risk of overreliance on AI in judicial decisions.

Finally, Scanlon et al. (2023a) describes the potential applications of ChatGPT and LLMs in digital forensics. The authors assess how ChatGPT can assist in various forensic tasks, such as identifying digital artifacts, generating code for forensic activities, and detecting anomalies in logs. LLMs present challenges, including issues with hallucinations, inaccuracies, and limitations when dealing with sensitive data. The study shows that ChatGPT can be a useful tool for investigators when used with caution, but human expertise remains important to ensure reliability in forensic investigations.

The application of LLMs in digital forensics has the potential to enhance investigators' capabilities to handle digital evidence and help solve cases with greater accuracy. However, it is important to remember that LLMs are not a replacement for human expertise, but rather a valuable tool that complements and assists forensic professionals. Therefore, we need a methodology and a dataset to evaluate LLMs as a forensic tool, particularly for timeline analysis, as discussed in this paper.

# 3. Proposed standardized methodology

To assess the performance of an LLM for timeline analysis, several aspects are important as depicted in Fig. 1. Note that we are not intending to replace forensic tools, but rather to assess the LLM's ability to follow structured instructions in a forensic context. We must define one or more tasks (Sec. 3.2) that we expect the LLM to perform. This involves designing a prompt to interact with the system, such as summarizing events into high-level insights or identifying indicators of compromise. In addition, a ground-truth dataset is needed that can be used to assess the outcome of an LLM (Sec. 3.3). Lastly, evaluation metrics are required that allow us to compare the ground truth with LLM output (Sec. 3.1). While starting with the tasks may seem natural, we recommend beginning with the evaluation metric, as it defines the required output, which in turn influences the task and prompt. This proposed methodology is designed to be flexible so that additional tasks and evaluation metrics can be incorporated in the future.

# 3.1. Evaluation metrics

We decided to use BLEU (Bilingual Evaluation Understudy) and ROUGE (Recall-Oriented Understudy for Gisting Evaluation). They were selected due to their widespread acceptance and established methodologies in machine translation and summarization. These metrics provide a way to quantify the quality of generated text and allow for comparisons across different models and tasks.

# 3.1.1. BLEU – Bilingual Evaluation Understudy

BLEU assesses the quality of machine-generated outputs by comparing them to human-curated reference texts (ground truth) (Papineni et al., 2002). The score focuses on how accurately and completely the machine or LLM has replicated the human ground truth. It is calculated as follows:



Fig. 1. The proposed methodology for quantitative evaluation of LLM-based timeline analysis.

$$BLEU = BP \times exp\left(\sum_{n=1}^{N} w_n log p_n\right)$$
 (1)

where  $p_n$  is the precision for each n-gram,  $w_n$  is the weight for each n-gram, and BP is the brevity penalty (BP). BP is designed to penalize generated text that is too short. The idea is that shorter text might artificially increase precision, but may not capture the full meaning of the original text. The brevity penalty is calculated as:

$$BP = \begin{cases} 1 & \text{if } c > r \\ e^{(1-r/c)} & \text{if } c \le r \end{cases}$$
 (2)

where c is the length of the candidate (machine) translation and r is the reference length.

# 3.1.2. ROUGE - Recall-Oriented Understudy for Gisting Evaluation

ROUGE is a collection of metrics designed to evaluate automatic summarization and machine translation systems (Lin, 2004). It focuses on the quality of the output generated by these systems. In our case, the essence of ROUGE is to provide a quantitative measurement of the quality of an automatically generated text from an LLM by comparing it with reference data or ground truth created by humans.

ROUGE includes several metrics, each serving a unique purpose in evaluating text. Two of the key metrics are ROUGE-N and ROUGE-L. ROUGE-N assesses the overlap of n-grams between the machine-generated text and the reference, where n is 1 and 2 in our experiments. ROUGE-L focuses on the longest common subsequence (LCS) between the LLM-generated n and the reference.

ROUGE-N is based on the *n*-gram overlap between the machinegenerated text and the reference as follows:

$$ROUGE - N = \frac{\sum_{s \in \{Reference\}} \sum_{n-gram \in s} Count_{match}(n - gram)}{\sum_{s \in \{Reference\}} \sum_{n-gram \in s} Count(n - gram)}$$
(3)

where  $Count_{match}(n\text{-gram})$  is the count of n-grams in the machinegenerated text that matches the ground truth. Count(n-gram) is the count of n-grams in the ground truth. On the other hand, ROUGE-L evaluates the LCS between the machine-generated text and the reference as follows:

$$ROUGE - L = \frac{\sum_{s \in \{Reference\}} LCS(s, Machine)}{\sum_{s \in \{Reference\}} Length(s)}$$
(4)

where LCS(s, Machine) refers to the length of the LCS between the system-generated text and the reference s. Finally, Length(s) is the length of the reference text. For both, BLEU and ROUGE, the higher the score, the better. We implemented both metrics using Hugging Face evaluate library (Hugging Face, 2024a,b). Note that we evaluate the text from the LLM's answer that is generated in a downloadable file, not from the text-based responses.

# 3.1.3. Considerations

Achieving high BLEU and ROUGE scores requires a significant overlap between the LLM's output and our ground truth data, where 'overlap' means identical wording. These metrics do not assess meaning but only textual similarity. For example, the sentences 'He is 25' and 'He was born in 2000' would yield low scores despite conveying the same information. To mitigate hallucinations, we must ensure that the LLM returns data in a specific format, which we also use in our ground truth. To achieve this, we designed largely deterministic tasks (solvable by traditional software) and provided examples within the prompt to guide the LLM. In the future, we plan to explore fine-tuning an LLM, which could enhance the user experience. However, this study focuses on feasibility, and fine-tuning is beyond its scope.

#### 3.2. Common tasks for forensic timeline analysis

Given the considerations and to quantitatively evaluate the capabilities of an LLM, we selected the following four tasks:

- Running grep for specific terms, i.e., assess how well the LLM handles a straightforward task such as running grep.
- Rule-based anomaly detection, i.e., looking for patterns that could also be identified using rules, such as multiple failed login attempts, could mean a brute-force attack.
- Event summarization, i.e., combining several low-level events into a more meaningful event, such as if events A and B are found (lowlevel), this means a new user was created (meaningful event).
- 4. Exploratory data analysis.

Tasks have been carefully chosen to be realistic and deterministic (besides task four), but also allow for validation, e.g., for running grep, we can develop our own grep expression. Note that only the first three tasks require a ground truth. With respect to the prompts, we follow the prompt style of Scanlon et al. (2023b) and the OpenAI prompt engineering guides (OpenAI, 2024b). More details are provided in the subsequent sections.

# 3.2.1. Prompts for running grep of specific terms

This task simulates a grep command to ensure that it can handle basic tasks without making critical errors. The example prompts are shown below:

- 1. "I am a forensic investigator. I need to find these terms: \b[A-Za-z0-9\_\\:.]+\.exe\b in the given CSV file to get all entries related to executable files (.exe). The CSV file is a forensic timeline generated from the log2timeline/Plaso tool."
- 2. "For your references, the grep command is: grep -E "\b[A-Za-z0-9\_\\:.]+\.exe\b" timeline.csv."
- "Do not include the first line of the file containing column names. Include all columns in the results, not only the message column. Export the results into plain text."

The prompt asks an LLM to replicate the functionality of a <code>grep</code> command, which is commonly used to search for patterns in the text. The goal is to search the CSV file for all entries that contain executable files with the <code>.exe</code> extension. In total, five terms need to be found, the system is expected to identify these entries and save the results in plain text format. In addition, we ask the system to exclude the header row and include all column values in the results. This task checks whether the LLM can effectively search through the forensic timeline using a regular expression to filter out relevant entries.

# 3.2.2. Prompts for rule-based anomaly detection

The goal of this task is to enable more natural queries against the timeline. This simulates providing a timeline and then asking about specific aspects, such as 'Have there been failed login attempts?' or 'Was registry.exe executed?' Rather than posing these queries one by one, we opted to include multiple elements of interest in a file (keyword list), which the user uploads. This approach effectively cross-references a keyword list with the CSV-based timeline.

Specifically, we provide the following prompt: 'I am a forensic investigator. Read this list of keywords to find suspicious events.' The user uploads a keyword list, allowing the system to focus on specific patterns or terms that may indicate anomalous or potentially adversarial behavior within the timeline.

As we require the output in a specific format, the uploaded file is in reality a JSON file which includes elements of the prompt (event) as well as what to look for (keyword). This helps the LLM to detect suspicious events in a timeline CSV file. Note, the keyword is extracted from the message column from the timeline data, i.e., it exists in the timeline

CSV. The event is our creation.

```
{
   "event": "Registry launch with prefetch file",
   "keyword": "Prefetch [REGEDIT.EXE] was executed'
}
```

The LLM is expected to return a JSON-formatted response that includes the timestamp of the detected event (date time), the name of the matched event (event), the keyword that triggered the match (keyword), and the full log message (message) from the timeline. This structured format facilitates automated comparison with ground truth data and supports downstream forensic analysis. This format also maintains consistency and interpretability to allow for an accurate evaluation using BLEU and ROUGE. An example output structure is shown below:

```
[{
    "datetime": "2023-12-27T00:37:14.609465+00:00",
    "event": "Registry launch with prefetch file",
    "keyword": "Prefetch [REGEDIT.EXE] was executed",
    "message": "Prefetch [REGEDIT.EXE] was executed -
    run count 3 path hints: \\WINDOWS\\REGEDIT.EXE
    hash: 0x246AC210 volume: 1 [serial number:
    0x5CE1DF5A device path: \\VOLUME{01da182ce1985a
    64-5ce1df5a}]"
}
```

### 3.2.3. Prompts for event summarization

A user action (high-level event) causes many entries in a timeline (low-level events). This task looks at the possibility of summarizing low-level tasks to high-level tasks. To solve this task without fine-tuning, we provide a code (a Python library) that can be used (executed) by the LLM.

The interaction between the user and ChatGPT is outlined in Fig. A.1. We provide a persona, such as stating a role (e.g., forensic investigator), including detailed information about the task (e.g., event type or data format), and offer additional tools to improve accuracy. These steps help the system manage responses more accurately. The prompt uses a space delimiter to provide suitable spacing to separate key pieces of information. In the third-to-last box ("Specify steps to run an event summarization"), the user outlines the exact procedure for summarizing events. This involves uploading the CSV file, selecting the type of event (such as "last-shutdown"), and executing the summarization using the given libraries. The expected return value for this task is as follows:

```
{"0": {
   "id": 1002,
   "date_time_min": "2023-12-26 00:34:47.890403+00:00",
   "date_time_max": "2023-12-26 00:34:47.890403+00:00",
   "date_time_max": "2023-12-26 00:34:47.890403+00:00",
   "evidence_source": "[9707 / Ox25eb] Provider identifier:
   {...} Source Name: Microsoft-Windows-Shell-Core Strings:
   ['msedge.exe\" --no-startup-window --win-session-start']
   Computer Name: WinDev2311Eval Record Number: 2249
   Event Level: 4",
   "type": "Process Creation",
   "description": "Process creation of 'msedge.exe'",
   "category": "Windows",
   "plugin": "EVT-WinEVTX-winevtx",
   "files": "NTFS:\Windows\System32\\winevt\\Logs\\
   Microsoft-Windows-Shell-Core%40perational.evtx",
   "keys":
   "Windows Event ID": "9707",
   "Windows Event ID (hex)": "0x25eb",
   "Executable name": "msedge.exe"
   },
   "supporting": { . . . },
   "trigger": { . . . }
} . . . }
} . . . }
} . . . }
```

# 3.2.4. Prompts for exploratory data analysis

This task analyzes exploratory data analysis (EDA), which allows gaining valuable insights into the dataset as a whole. This task is a qualitative module and is not part of the core quantitative benchmark. For instance, EDA may help investigators grasp the structure, distribution, and key features. It may also enable the identification of patterns and relationships between events, such as how user behaviors might be interconnected. In addition, it facilitates the visualization of temporal data, which is an important aspect of timeline analysis. Using diagrams

such as histograms and heatmaps, investigators can acquire a clearer understanding of trends and cycles in the data. These visualizations pinpoint periods of interest and aid in the identification of suspicious activities for further investigation.

The example EDA prompt is: "Explore patterns of event occurrences based on the datetime field per second (e.g., busiest times, significant gaps), use a bar chart. Write the hour:minute:second in the x axis". An LLM will generate Python code to create the bar chart, and we can download the chart as a PNG file.

In this proposed standardized LLM evaluation, event summarization comprises two scenarios: summarizing a single event or multiple events. Summarizing a single event means the method extracts one specific event from the provided timeline, such as a Google search (full list see Sec. 3.3.5). Consequently, multiple events mean the LLM is tasked with summarizing all defined events.

# 3.3. Ground truth

To assess the quality of output (LLM response), we require a ground truth dataset, i.e., documentation of the underlying dataset (Göbel et al., 2023; Breitinger and Jotterand, 2023). A peculiarity in our scenario is that we need the ground truth in a specific format so that it is comparable with the output of an LLM (automated). Specifically, there is no easy way to compare a disk image or its corresponding timeline against the LLM output. Consequently, the underlying dataset must be converted into a text-based format (ground truth), allowing an automated comparison with the LLM output.

To accomplish this, we first must create a dataset (Sec. 3.3.1) where the creation process is documented or recorded. Next, we generate a timeline of the disk image (Sec. 3.3.2) which serves as an input for the LLM. Lastly, using the documentation and timeline, we manually create the *expected outcome* which represents our ground truth (Sec. 3.3.5 to 3.3.3).

# 3.3.1. Scenario and dataset generation

The first step was to create a dataset, as no appropriate dataset was available. The procedure is illustrated in Fig. 2, and the dataset is shared through the Zenodo. Our test bed was a Windows 11 machine within a virtual environment simulating regular computer usage. All activities were recorded using screen capture (video) and are documented (written notes). We emphasize that this dataset was deliberately designed to be synthetic and single-user to serve as a controlled feasibility testbed. The goal was to validate the operation and reliability of the proposed standardized evaluation methodology using a timeline with clear ground truth and minimal ambiguity.

The scenario follows a sequence of opening applications, downloading software, and accessing websites. The user begins by opening the Edge browser and then navigates to Bing. They perform a search query for "Mozilla Firefox download" on Bing and visit Mozilla's official website to download the Firefox browser. After that, the user opens the File Explorer to navigate the downloaded installer. The user runs the Firefox installer and opens the newly installed Firefox browser. Afterward, they navigate to Google, perform a search related to SQL injection, and open a tutorial on the W3Schools website. The session ends with a system shutdown, indicating that the user has completed all activities.

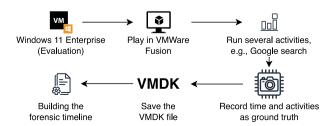


Fig. 2. Building ground truth for LLM evaluation.

There is no strict rule determining the number of ground truth entries used. The number provides sufficient variation to evaluate the model's ability to run different types of tasks while keeping the dataset size manageable for manual validation.

# 3.3.2. Timeline generation

To generate the timeline, we ran log2timeline/Plaso (Metz et al., 2024) on the vmdk file. The tool (Plaso) analyzes all known artifacts<sup>2</sup> and compiles them into a single unified timeline. By default, the tool processes all partitions from a vmdk file and generates a Plaso storage file (\*.plaso, a database file) containing the forensic timeline. To convert plaso file to a CSV timeline file, we ran psort.

# 3.3.3. Ground truth for task 1: running grep for specific terms

Building the ground truth requires manually running grep on the dataset and taking note of the output. This was done for the following five keywords:

- RegisteredApplications: obtaining events related to registered applications in the Windows registry.
- (OneDrive OneDrive\.exe): finding events related to Microsoft OneDrive application.
- 3.  $\b [A-Za-z0-9_{\cdot}: ]+\ exe\b : looking for all entries related to executable files (.exe).$
- 4616 /: finding Windows event ID 4616, which is related to system time change without regex.
- 5. \[4616 / 0x1208\].\*Microsoft-Windows-Security -Auditing.\*svchost.exe: finding Windows event ID 4616 with regex.

The command to generate this ground truth is grep -E keyword timeline.csv, where -E indicates that extended regular expressions are being used with the grep command.

# 3.3.4. Ground truth for task 2: rule-based anomaly detection

The second ground truth requires matching keywords (or phrases) with events. We create the keywords as a rule-based approach by first checking the date and time of the event we performed earlier in the Windows test-bed. Next, we manually look for related entries in the timeline CSV file. Once we find the relevant entry, such as registry launch, we extract the keywords linked to the event. Finally, we format these keywords into a JSON format as shown below:

```
"event": "Registry launch with prefetch file",
    "keyword": "Prefetch [REGEDIT.EXE] was executed"
}
```

In the evaluation, we can ask questions in natural language because the event and the keyword the LLM searches for are already defined. Unlike an event summarization task, no script or library is provided, and the LLM handles the matching on its own. These keywords collected are a useful technique to identify suspicious events in the forensic timeline. There are seven keywords in total, and the full list of keywords in JSON format is available on the Zenodo.

# 3.3.5. Ground truth for task 3: event summarization

Event summarization aims at combining low-level events to obtain high-level events as proposed by Hargreaves and Patterson (2012). Forming the ground truth was accomplished by implementing the dftpl tool<sup>3</sup> as described by the authors. Given a CSV timeline, our prototype can extract certain high-level events and return a JSON file.

 $<sup>^{2}\,</sup>$  Plaso consists of various parsers for different artifacts. Artifacts unknown to Plaso are ignored.

<sup>&</sup>lt;sup>3</sup> https://github.com/studiawan/dftpl.

There are eight predefined events, grouped into three categories:

- 1. Web: Google search, Bing search, and web visit
- 2. Windows: last shutdown, process creation, and program opened
- 3. User activity: file download, and recent file access

We chose JSON because it is human-readable, making it easier for investigators to interpret and manually validate results. JSON also facilitates straightforward comparison with evaluation metrics due to its structured nature for efficient parsing. Moreover, its compatibility with various programming languages and tools further supports automation and quantitative evaluation in forensic analysis workflows.

To create the high-level events, we ran the dftpl command as follows: dftpl -i timeline-input.csv -o summarization-output.json -t last-shutdown, where -i is a Plaso CSV file, -o is the output (in JSON), and -t specifies the event of interest. The -t option can be omitted to summarize multiple events. The list of high-level events was then manually validated to ensure it was correct.

A sample ground truth in JSON format is provided below:

The sample includes the following high-level activities:

- 1. *id*: A unique identifier for the event, which is a number that differentiates this event from others.
- date\_time\_min: The earliest possible timestamp for when the event could have occurred.
- 3. *date\_time\_max*: The latest possible timestamp for when the event could have occurred.
- evidence\_source: Refers to the Plaso message that provides information about the event.
- 5. *type*: The nature of the event, such as Google Search, File Download, or any other high-level event type.
- description: A human-readable explanation or summary of the event.
- category: A higher-level classification or tag for filtering or organizing events.
- 8. *plugin*: Identifies the Plaso plugin used to parse the source file from which the event was extracted.
- 9. *files*: Refers to the file(s) related to the event, such as the log file, binary file, or any other data source.
- keys: Stores additional key-value pairs related to the event, such as specific attributes or metadata.
- 11. *supporting*: Stores a list of five low-level events before and after the main event for context.
- 12. *trigger*: Refers to the reasoning artifact or piece of evidence that caused the event to be recognized.

#### 4. Experimental results and analyses

This section details the experimental settings, along with the analysis, results, and discussion of our case study.

# 4.1. Experimental settings

We used the version of log2timeline/Plaso, which was the Docker image version 20230717. The target operating system was Microsoft Windows 11 Enterprise. The OS was sourced from the Microsoft Developer Network, specifically the evaluation virtual machine (VM) version 2311 (Microsoft Developer, 2024). For virtualization, we opted for VMWare Fusion 13.5.0. For the LLM, we selected ChatGPT-40, one of the most advanced models available at the time of writing this paper. This study focuses on ChatGPT-40 to demonstrate the feasibility and structure of a standardized evaluation methodology, rather than to benchmark performance across multiple LLMs. To facilitate containerized environments, we use Docker Desktop version 4.22.1 (118664).

The extracted full timeline is too large to be handled by ChatGPT due to token limitations. Consequently, we only provided ChatGPT with about 2000 lines of Plaso entries as a timeline of interest. We have experimented with different sizes (e.g., 1000, 2000, 3000 lines) and found 2000 lines to be a manageable amount that balances input size and processing efficiency. We understand that real-world forensic timelines often span millions of lines. In practice, deployment would require preprocessing techniques, such as time-window-based segmentation, event-type filtering, or sliding windows, to divide large timelines into manageable subunits for LLM processing.

# 4.2. Timeline analysis with ChatGPT

The Advanced Data Analysis feature of ChatGPT, previously called Code Interpreter, is now integrated into ChatGPT versions 4 and 40 (OpenAI, 2024a). This feature allows users to analyze data and interpret code directly within the platform. This enhances the user experience by supporting data uploads, where users can write, test, and execute code seamlessly. The supported file formats include text, image files, PDFs, and Word documents, code or other data files, as well as audio and video. In this study, we used the CSV file generated by Plaso. Once the data is uploaded, we can use the prompts to instruct ChatGPT to read or analyze the timeline.

We employ ChatGPT in two scenarios: with and without additional knowledge. In the first scenario, we provided ChatGPT with specific information related to the task, such as a library for event summarization (Sec. 3.3.5) or a list of keywords to detect suspicious activities (Sec. 3.3.4). In the latter scenario, we did not provide any additional information and relied solely on ChatGPT's existing language model to analyze the timeline.

# 4.3. Results and analysis

To quantitatively evaluate ChatGPT for forensic timeline analysis, we developed four tasks, including ground-truth data. For example, the event summarization task has 14 event types, the rule-based anomaly detection task has seven rules, and the search task for specific terms has five keywords. Note that the exploratory data analysis task does not have evaluation metrics because there is no ground truth data for this task.

A sample result of the given prompts and the ChatGPT answers is depicted in Fig. A.1. The evaluation results for the used datasets are shown in Table 1, where the metric values represent the mean values for each task.

# 4.3.1. Results of running grep for specific terms

It is important to note that when asked to search for specific terms, ChatGPT does not run the grep command. Instead, it generates Python code to perform the search. The results of this task are shown in Table 1.

**Table 1**Evaluation results of various tasks given to ChatGPT for forensic timeline analysis.

Task	BLEU	ROUGE- 1	ROUGE- 2	ROUGE- L	Mean
Without additional knowled	dge				
Event summarization (single)	0.077	0.192	0.129	0.136	0.134
Event summarization (multiple)	0.001	0.171	0.120	0.132	0.106
Rule-based anomaly detection	0.147	0.144	0.075	0.141	0.127
Run grep for specific terms	0.847	1.000	1.000	1.000	0.962
With additional knowledge					
Event summarization (single)	0.999	1.000	1.000	1.000	1.000
Event summarization (multiple)	0.743	0.786	0.786	0.786	0.775
Rule-based anomaly detection	0.945	0.997	0.996	0.997	0.984
Run grep for specific terms	0.847	1.000	1.000	1.000	0.962

The results indicate that ChatGPT performs this task effectively, especially when provided with additional knowledge, i.e., the corresponding grep command. Without additional knowledge, the BLEU score is 0.847, and both ROUGE-1 and ROUGE-L are 1.000. The results suggest that the system accurately identifies specific terms most of the time, but with minor variations that affect the BLEU score. With additional knowledge, the BLEU, ROUGE-1, and ROUGE-L scores all reach 1.000, and they demonstrate that the model can perfectly match the specific terms when it has more context or knowledge about the data. These findings imply that the performance of ChatGPT in conducting targeted searches is enhanced when it is given relevant prior information. Therefore, it produces consistent and fully accurate results.

ChatGPT can detect all entries correctly when provided with additional knowledge or information. However, the grep output from ChatGPT does not contain commas, whereas the ground truth does, as the timeline is a comma-separated file. In addition, the model's output has extra spaces that are not present in the original data. Furthermore, it gives inconsistent output when no additional knowledge is provided. In several cases, it only produces incomplete results, displaying only the "message" column without including all other columns.

# 4.3.2. Results of rule-based anomaly detection

As mentioned in Sec. 4.2, there are two scenarios: one with additional knowledge and one without. In the case without additional knowledge, the prompt is slightly different because it does not include instructions to read the uploaded keywords file. In this task, we can instruct ChatGPT to format the answers in a specific format, such as JSON. The prompt would be "Format your answer using this JSON format:" and we can give an example format as follows:

```
{
  "datetime": "2023-12-26T23:26:57.492844+00:00",
  "event": "Disabling Windows Firewall recorded in Windows logs",
  "keyword": "[2082 / 0x0822] Provider identifier",
  "message": "[2082 / 0x0822] Provider identifier: {d1bc9aff-...}
  Source Name: Microsoft-Windows-Windows Firewall
  With Advanced Security Strings"
```

Moreover, we instruct the system to export all results to a down-loadable file, with "I need all entries of suspicious entries. Export to a JSON file for all of the results".

In the task of rule-based anomaly detection without additional knowledge, the performance was poor: The BLEU score is 0.147, and the ROUGE scores range from 0.141 to 0.192, indicating that the model's output is significantly different from the expected output. The keywords

generated by ChatGPT are as follows: 'delete', 'clear', 'wipe', 'remove', 'malware', and 'unauthorized'. These low scores reflect minimal overlap between the system's output and the expected results, both in terms of individual words and word sequences. However, it is important to note that these evaluation metrics are based on word matching and do not account for semantic similarity. Although the wording used by ChatGPT may differ from the predefined ground truth, the underlying interpretation or intent of the result may still be forensically relevant or correct.

In contrast, the results improve when additional knowledge is provided. Specifically, the BLEU score rises to 0.945, and the ROUGE scores increase to nearly perfect values (ranging from 0.996 to 0.997). This means that the generated outputs closely match the expected results. This highlights the importance of providing context or specialized knowledge to improve performance in more complex forensic analysis tasks.

Even with additional information or knowledge, ChatGPT can still make mistakes. The errors are mainly due to differences in how characters are escaped. For example, the ground truth uses two backslashes to escape regular expressions (regex), while ChatGPT's output uses four backslashes to escape the "\" character.

#### 4.3.3. Results of event summarization

Summarizing a single event means the method extracts one specific event from the provided timeline, such as a Google search (see full list in Sec. 3.3.5). Consequently, multiple events mean the LLM is tasked with summarizing all defined events.

Our research indicates that ChatGPT uses a virtual environment to run Python code when responding to user prompts. This means we can install the <code>dftpl</code> Python wheel installer within that virtual environment. Note that in the "with additional knowledge" setting, the LLM is instructed to use the <code>dftpl</code> library to perform the summarization. This means the LLM is not independently interpreting the event semantics, but rather acting as an automation agent that follows structured instructions to run a pre-existing summarization tool. This setup allows us to evaluate the model's ability to accurately execute forensic tools and follow code-level prompts, rather than to perform semantic reasoning or event abstraction on its own.

To respond to the user prompts, ChatGPT generates Python code as shown in Fig. A.1. For example, if the parser example is designed to work for all supported events, ChatGPT can summarize a specific event, such as the last shutdown event on Windows. One can click the '[>\_]'-button to view the generated Python source code. Thus, experienced investigators may validate the code and, with it, the answer. Finally, the results can be downloaded in a JSON format, and this file will be quantitatively evaluated based on the ground truth from Sec. 3.3.5.

The result of event summarization on single and multiple events without additional knowledge shows a low performance, with a BLEU score of 0.077, indicating limited precision in generating a summarization that closely matches the expected events. The ROUGE-1 score of 0.192 suggests that around 19.2 % of single words in the generated output matched the reference, while the ROUGE-2 score of 0.129 shows even lower overlap in bigrams (two-word sequences). The ROUGE-L score of 0.136 reflects a moderate match in terms of the longest sequence of matching words. However, we conclude that without additional knowledge, the system cannot accurately summarize events.

In contrast, the result for a single event with additional knowledge, i. e., using the dftpl library, shows near-perfect performance, with a BLEU score of 0.999 and ROUGE-1, ROUGE-2, and ROUGE-L scores, all at 1.000. This indicates that the ChatGPT output almost exactly matched the reference in terms of precision, word overlap, and sequence structure. The high scores suggest that, with additional knowledge, the system was able to mimic the expected results. The reason is that we gave a Python library that can summarize events based on the method described in Hargreaves and Patterson (2012) to ChatGPT (Fig. A.1). Although we did not explicitly instruct ChatGPT to follow a particular order, the ground truth output produced by the dftpl library is

chronologically ordered by timestamp. For the multiple event summarization task, the evaluation scores were lower because ChatGPT generated the correct events but in a different order than the ground truth. The beginning of the file displays timestamps that increase or remain the same, indicating a mostly sorted order. Similarly, the end of the file follows a chronological pattern. However, the middle sections break this order, with some events appearing earlier than preceding ones. This discrepancy in ordering affected the BLEU and ROUGE scores, which are sensitive to the sequence of words or structures. Importantly, while the order differed, the extracted content was sometimes semantically correct and forensically valid. Future work may include implementing order-invariant evaluation metrics or normalizing the output order before comparison to address this issue.

# 4.3.4. Results of exploratory data analysis

This section aims to explore how ChatGPT can assist forensic investigators in identifying patterns or anomalies within large timelines through exploratory data analysis (EDA). Specifically, we evaluate the model's ability to generate useful visualizations that support investigative tasks. The example of a generated bar chart is shown in Fig. 3. The chart shows the number of event occurrences per second within a specific time range, where each bar corresponds to a second in the format: hour:minute:second. The data reveal variability in event activity, with most seconds seeing between 50 and 150 events. However, there is a noticeable spike at 00:45:55, where the event count exceeds 250, which indicates a sudden surge in activity during that particular second. The concentration of events at specific seconds may point to important actions or incidents that require further investigation, especially during periods of relatively low activity that are punctuated by intense bursts (Studiawan and Sohel, 2021).

Another chart generated by ChatGPT using Python is a heatmap shown in Fig. 4. The heatmap illustrates the flow of the event sequence, showing the transitions between various types of events based on their timestamps. The rows represent the current event types, while the columns represent the next event types, with each cell indicating how often a specific event type is followed by another. The color intensity, as shown by the legend, reflects the frequency of these transitions, with darker shades showing more frequent sequences. The heatmap highlights common flows in the event timeline and provides valuable insight into which events tend to trigger others. Therefore, it can help to understand the sequences of events within the forensic timeline analysis.

Key patterns can be observed in this visualization. For example, 'Metadata Modification Time' transitions into itself 1079 times, suggesting that it frequently repeats or is followed by itself in the sequence of events. There are also noticeable transitions from 'Creation Time' to 'Metadata Modification Time' (118 times) and from 'Last Access Time' to 'Metadata Modification Time' (98 times).

The heatmap reveals typical patterns in the event timeline by showing how certain events frequently follow others. This visualization helps investigators better understand the sequence and relationship between events during forensic timeline analysis. In short, EDA can be



Fig. 3. A bar chart generated by ChatGPT in an exploratory data analysis task.

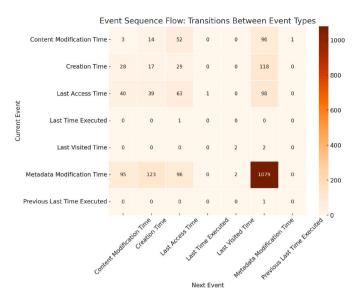


Fig. 4. A heatmap generated by ChatGPT for event sequence flow.

done by a human investigator, but using ChatGPT can help speed up this manual work.

#### 4.4. Discussion

#### 4.4.1. Overall quantitative evaluation

Without additional knowledge, tasks such as 'Event summarization (single)' and 'rule-based anomaly detection' have mean scores of 0.134 and 0.127, respectively, indicating limited accuracy. However, 'run grep for specific terms' achieves a much higher mean score of 0.962, suggesting that ChatGPT can handle searches for specific terms relatively well, even without prior information. With additional knowledge, mean scores improve across all tasks. Single event summarization tasks achieved a perfect mean score of 1.000, while the multiple event one obtained 0.775. The results demonstrate inconsistent accuracy scores, even when provided with relevant context. The mean score for the rule-based anomaly detection task also increases to 0.984. The consistent mean score of 0.962 for "run grep for specific terms" shows that the task is already handled effectively regardless of additional knowledge. In the grep task, providing prior information does not lead to further improvement. Future work will incorporate secondary human expert review of LLM outputs that differ lexically but appear plausible.

# 4.4.2. API-based and BERTScore implementations

Our experiments were conducted using the ChatGPT web interface to evaluate feasibility and task alignment under realistic investigator interactions. However, forensic analysts aiming to operationalize or scale these evaluations would rely more on API-based implementations. We acknowledge that BLEU and ROUGE are limited in capturing semantic equivalence, particularly in forensic tasks involving reasoning and inference. Their use in this study was motivated by the need for an initial, automated, and deterministic evaluation protocol. To address their limitations, we are incorporating semantic-aware metrics such as BERTScore (Zhang et al., 2020), which better account for meaning despite syntactic variation. API-based and BERTScore implementations are provided as work in progress on a GitHub repository.

# 4.4.3. CSV file size of a forensic timeline

While ChatGPT is advertised as being capable of handling CSV files

<sup>&</sup>lt;sup>4</sup> https://github.com/studiawan/llm-forensic-timeline.

up to 50 MB in size, <sup>5</sup> we found that in practice, it struggles to process files of that size. Throughout our experiments, we observed that ChatGPT could successfully analyze smaller CSV files, but when attempting to work with larger files (10 MB or more), the model often encountered errors or failed to provide results. This discrepancy suggests that, despite the claims in the documentation, there are practical limitations when analyzing larger datasets. LLMs are limited by input size constraints and context window lengths, which restrict their ability to reason over extensive forensic timelines without segmenting the data or relying on iterative processing techniques.

#### 5. Conclusion and future work

The proposed methodology and dataset have demonstrated their potential for quantitative evaluation of timeline analysis using LLMs. Using the proposed standardized methodology and dataset, researchers can apply and expand the test and evaluation of LLM-based forensic timeline analysis. By employing the advantages of natural language processing on LLMs, e.g., ChatGPT, the presented case studies show that it can assist in analyzing events and temporal information from a forensic timeline. It also provides valuable information for forensic

investigators, particularly in the task of exploratory data analysis. However, based on the quantitative evaluation, ChatGPT performs worse than a rule-based approach or a regular expression-based approach accompanied by a human investigator.

For future work, we plan to add more datasets and tasks, such as real-world malware scenarios. We also plan to explore the use of other commercial LLM services, such as Google Gemini and Claude, to evaluate the robustness of our approach. In addition, to address concerns about the confidentiality of digital evidence, we can deploy open-source LLMs, such as LLaMA (Touvron et al., 2023) and Mixtral (Jiang et al., 2024) on a local device as shown by Sharma et al. (2025). By keeping the forensic timeline on the local computer, we aim to avoid the need to upload sensitive data to cloud-based LLM services, thus ensuring the privacy of the investigation.

# Acknowledgments

We would like to thank Christopher Hargreaves for his valuable comments and feedback. This research is supported by the Institut Teknologi Sepuluh Nopember, Indonesia through the ITS Research Funding scheme.

Appendix A. A conversation sample between an investigator and ChatGPT

<sup>&</sup>lt;sup>5</sup> https://help.openai.com/en/articles/8555545-file-uploads-faq.

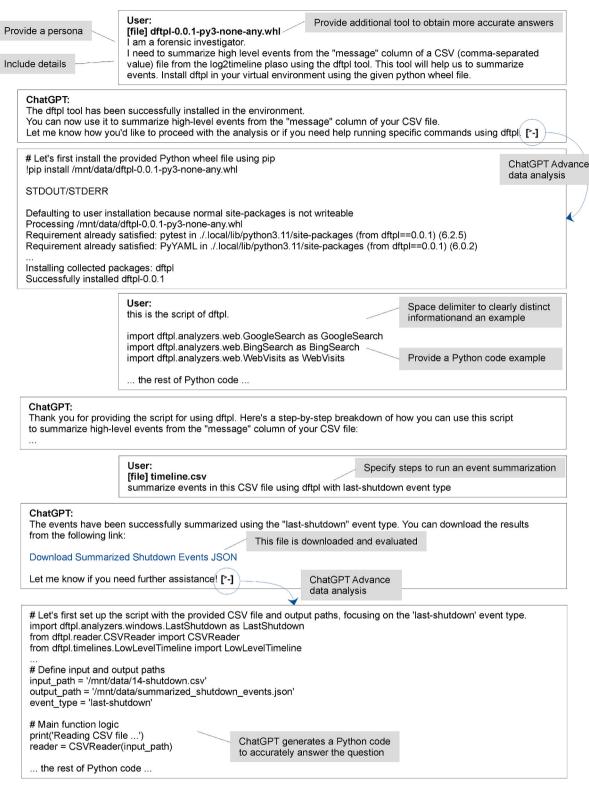


Fig. A.1. Example of ChatGPT interaction for summarizing high-level events from a forensic timeline.

# References

Arshad, H., Jantan, A.B., Abiodun, O.I., 2018. Digital forensics: review of issues in scientific validation of digital evidence. J. Informat. Process. Syst. 14, 346–376. Bhandari, S., Jusas, V., 2020. An abstraction based approach for reconstruction of timeline in digital forensics. Symmetry 12, 104.

Breitinger, F., Jotterand, A., 2023. Sharing datasets for digital forensic: a novel taxonomy and legal concerns. Forensic Sci. Int.: Digit. Invest. 45, 301562.

Breitinger, F., Studiawan, H., Hargreaves, C., 2025. Sok: timeline based event reconstruction for digital forensics: terminology, methodology, and current challenges. URL. https://arxiv.org/abs/2504.18131,arXiv:2504.18131.

- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al., 2020. Language models are few-shot learners. Adv. Neural Inf. Process. Syst. 33, 1877–1901.
- Brunty, J., 2023. Validation of Forensic Tools and Methods: a Primer for the Digital Forensics Examiner, 5. Wiley Interdisciplinary Reviews: Forensic Science, e1474
- Buchholz, K., 2023. One million users: threads shoots past one million user mark at lightning speed. https://www.statista.com/chart/29174/time-to-one-million-users/.
- Chabot, Y., Bertaux, A., Nicolle, C., Kechadi, M.T., 2014. A complete formalized knowledge representation model for advanced digital forensics timeline analysis. Digit. Invest. 11, S95–S105.
- Dinis-Oliveira, R.J., Azevedo, R.M., 2023. ChatGPT in forensic sciences: a new Pandora's box with advantages and challenges to pay attention. Forens. Sci. Res. 8, 275–279.
- Göbel, T., Baier, H., Breitinger, F., 2023. Data for digital forensics: why a discussion on "how realistic is synthetic data" is dispensable. Digit. Threat.: Res. Pract. 4, 1–18.
- Hall, S.W., Sakzad, A., Choo, K.K.R., 2022. Explainable Artificial Intelligence for Digital Forensics, 4. Wiley Interdisciplinary Reviews: Forensic Science, e1434.
- Hargreaves, C., Patterson, J., 2012. An automated timeline reconstruction approach for digital forensic investigations. Digit. Invest. 9 (Suppl. m), 869–879.
- Henseler, H., van Beek, H., 2023. ChatGPT as a copilot for investigating digital evidence. In: Proceedings of the Third International Workshop on Artificial Intelligence and Intelligent Assistance for Legal Professionals in the Digital Workplace, pp. 58–69.
- Horsman, G., 2019. Tool testing and reliability issues in the field of digital forensics. Digit. Invest. 28, 163–175.
- Hugging Face, 2024a. Metric: bleu. https://huggingface.co/spaces/evaluate-metric/ble
- $\label{prop:hugging} \mbox{ Hugging Face, 2024b. Metric: rouge. $https://huggingface.co/spaces/evaluate-metric/rouge. $https://huggingface.go/spaces/evaluate-metric/rouge. $https://huggingface.go/spaces/evalu$
- Hughes, N., Karabiyik, U., 2020. Towards reliable digital forensics investigations through measurement science. Wiley Interdisciplin. Rev.: Forensic Sci. 2, e1367.
- Inglot, B., Liu, L., 2014. Enhanced timeline analysis for digital forensic investigations. Inf. Secur. J. A Glob. Perspect. 23, 32–44.
- Jiang, A.Q., Sablayrolles, A., Roux, A., Mensch, A., Savary, B., Bamford, C., Chaplot, D.S., Casas, D.d.l., Hanna, E.B., Bressand, F., et al., 2024. Mixtral of Experts arXiv: 2401.04088.
- Lin, C.Y., 2004. ROUGE: a package for automatic evaluation of summaries. In: Text Summarization Branches out, pp. 74–81.
- Metz, J., Gudjonsson, K., White, D., et al., 2024. log2timeline Plaso: super timeline all the
- things. https://github.com/log2timeline/plaso.
  Michelet, G., Breitinger, F., 2024. ChatGPT, Llama, can you write my report? An experiment on assisted digital forensics reports written using (local) large language models. Forensic Sci. Int.: Digit. Invest. 48, 301683.
- Microsoft Developer, 2024. Get a windows 11 development environment. https://developer.microsoft.com/en-us/windows/downloads/virtual-machines/.

- Nalawade, A., Bharne, S., Mane, V., 2016. Forensic analysis and evidence collection for web browser activity. In: 2016 International Conference on Automatic Control and Dynamic Optimization Techniques (ICACDOT), pp. 518–522.
- NIST, 2019. Computer forensics tool testing program (CFTT). https://www.nist.gov/it l/ssd/software-quality-group/computer-forensics-tool-testing-program-cftt.
- OpenAI, 2024a. Data analysis with ChatGPT. https://help.openai.com/en/articles/8437071-data-analysis-with-chatgpt.
- OpenAI, 2024b. Prompt engineering. https://platform.openai.com/docs/guides/promp t-engineering.
- Palmbach, D., Breitinger, F., 2020. Artifacts for detecting timestamp manipulation in ntfs on windows and their reliability. Forensic Sci. Int.: Digit. Invest. 32, 300920.
- Papineni, K., Roukos, S., Ward, T., Zhu, W.J., 2002. BLEU: a method for automatic evaluation of machine translation. In: Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, pp. 311–318.
- Scanlon, M., Breitinger, F., Hargreaves, C., Hilgert, J.N., Sheppard, J., 2023a. ChatGPT for digital forensic investigation: the good, the bad, and the unknown. Forensic Sci. Int.: Digit. Invest. 46, 301609.
- Scanlon, M., Nikkel, B., Geradts, Z., 2023b. Digital forensic investigation in the age of ChatGPT. Forensic Sci. Int.: Digit. Invest. 44, 301543.
- Sharma, B., Ghawaly, J., McCleary, K., Webb, A.M., Baggili, I., 2025. Forensicllm: a local large language model for digital forensics. Forensic Sci. Int.: Digit. Invest. 52, 201272
- Studiawan, H., Ahmad, T., Santoso, B.J., Pratomo, B.A., 2022a. Forensic timeline analysis of iOS devices. In: 2022 International Conference on Engineering and Emerging Technologies (ICEET), pp. 1–5.
- Studiawan, H., Ahmad, T., Santoso, B.J., Shiddiqi, A.M., Pratomo, B.A., 2022b.

  DroneTimeline: forensic timeline analysis for drones. SoftwareX 20, 101255.
- Studiawan, H., Sohel, F., 2021. Anomaly detection in a forensic timeline with deep autoencoders. J. Inf. Secur. Appl. 63, 103002.
- Studiawan, H., Sohel, F., Payne, C., 2019. A survey on forensic investigation of operating system logs. Digit. Invest. 29, 1–20.
- Studiawan, H., Sohel, F., Payne, C., 2020. Sentiment analysis in a forensic timeline with deep learning. IEEE Access 8, 60664–60675.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al., 2023. Llama: Open and Efficient Foundation Language Models, 13971 arXiv:2302.
- Wickramasekara, A., Breitinger, F., Scanlon, M., 2025. Exploring the potential of large language models for improving digital forensic investigation efficiency. Forensic Sci. Int.: Digit. Invest. 52, 301859.
- Zhang, T., Kishore, V., Wu, F., Weinberger, K.Q., Artzi, Y., 2020. BERTScore: evaluating text generation with BERT. In: International Conference on Learning Representations (ICLR).