



Exploring the potential of large language models for improving digital forensic investigation efficiency

Akila Wickramasekara^{a, id}, Frank Breitinger^{b, c, id}, Mark Scanlon^{a, id, *}

^a Forensics and Security Research Group, School of Computer Science, University College Dublin, Belfield, Dublin 4, Ireland

^b School of Criminal Justice, University of Lausanne, Lausanne, Switzerland

^c Institute of Computer Science, University of Augsburg, Augsburg, Germany¹

ARTICLE INFO

Keywords:

Digital forensics
Large language models
LLM
Investigative process
Challenges

ABSTRACT

The ever-increasing workload of digital forensic labs raises concerns about law enforcement's ability to conduct both cyber-related and non-cyber-related investigations promptly. Consequently, this article explores the potential and usefulness of integrating Large Language Models (LLMs) into digital forensic investigations to address challenges such as bias, explainability, censorship, resource-intensive infrastructure, and ethical and legal considerations. A comprehensive literature review is carried out, encompassing existing digital forensic models, tools, LLMs, deep learning techniques, and the use of LLMs in investigations. The review identifies current challenges within existing digital forensic processes and explores both the obstacles and the possibilities of incorporating LLMs. In conclusion, the study states that the adoption of LLMs in digital forensics, with appropriate constraints, has the potential to improve investigation efficiency, improve traceability, and alleviate the technical and judicial barriers faced by law enforcement entities.

1. Introduction

With the widespread growth of information and communication technology (ICT) and information systems, cybercrimes have seen a significant increase in recent years (Ali, 2019).² As a further compounding factor, the number of "traditional" police investigations that include digital evidence is also increasing (Du and Scanlon, 2019). Addressing and investigating this volume of cases presents substantial challenges.

Generative AI (GenAI) and Large Language Models (LLMs) have become prominent topics of global discussion, prompting researchers to intensify their investigations by leveraging the capabilities of LLMs. The usage of LLMs within the scientific community experienced a rapid surge after 2022, notably with the advent of OpenAI's ChatGPT platform. In a relatively short period of time, this topic has attracted the attention of academia, industry, and the research community at large (Ray, 2023). Simultaneously, researchers are exploring the potential of LLMs in various domains and assessing their impact on the future of science and society. This inquiry also includes an examination of the potential harmfulness associated with the deployment of LLMs (Brown et al., 2020; Nozza et al., 2022). In other words, the use of LLMs in various tasks can

be a double-edged sword, necessitating careful consideration depending on the specific situations and contexts.

Given the rapidly evolving landscape of LLMs, it is prudent to look into various types and their unique capabilities. A nuanced understanding of the strengths and characteristics of different LLMs can contribute to more informed and effective applications within the dynamic field of digital forensics (DF).

This paper reviews recent advances in the application of LLMs within digital forensics, focussing on established models, methods, and key challenges. By examining contemporary studies from 2019 onwards, the survey highlights core areas, such as automation, investigative methods, and efficiency improvements facilitated by LLMs. In addition, it explores the literature that addresses challenges in both DF and LLMs, covering limitations, ethical considerations, and forensic-specific risks. This comprehensive review synthesises current insights and emerging trends, offering a foundation for understanding the potential and limitations of LLMs in DF contexts.

In light of fast-paced advancements and the recent explosion in LLM-focused research, a substantial influx of LLM-focused research papers has occurred since the launch of ChatGPT in late 2022. Due to this

* Corresponding author.

E-mail addresses: akila.wickramasekara@ucdconnect.ie (A. Wickramasekara), frank.breitinger@uni-a.de (F. Breitinger), mark.scanlon@ucd.ie (M. Scanlon).

¹ Present affiliation.

² <https://go.crowdstrike.com/rs/281-OBQ-266/images/GlobalThreatReport2024.pdf>.

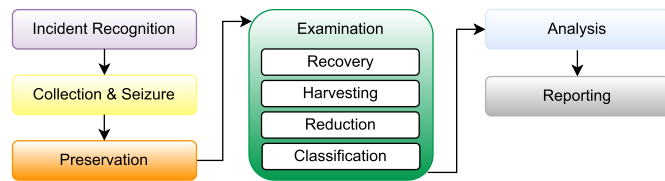


Fig. 1. Traditional digital forensic process model (Casey, 2011).

fast pace, many seminal research articles exist solely as preprints on preprint services, e.g., arXiv. For example, the official papers for GPT-4 (Achiam et al., 2023), LLaMA (Touvron et al., 2023) and DeepSeek-V3 (DeepSeek-AI et al., 2024) are only published on arXiv, but have garnered thousands of citations. Despite their preprint status, these articles offer essential insights critical for contemporary research and dialogue within the domain, making their incorporation into this article necessary to provide the most up-to-date knowledge and perspectives.

The paper is structured as follows: Section 1.1 provides a comprehensive background for the review, delving into existing DF process models, the challenges inherent in DF, and a detailed overview of the current work conducted with the use of LLMs within DF. In Section 2 the paper delves into the realm of Natural Language Processing (NLP), elucidating the working principles of LLMs, their architectural foundations, and the specifics surrounding specially trained LLMs. Section 3 provides an in-depth review that focusses on the capabilities and benchmark information of LLMs trained for coding tasks, as well as those tailored for vision assistance. Section 3 explores the synergy between DF and LLMs, detailing how LLMs can be effectively employed in each phase of the DF process model. Finally, in Section 6, the paper summarises the future challenges associated with integrating LLMs with automated agents within the DF domain. The conclusion outlines potential avenues for future research and development, shedding light on the path of future DF investigations employing LLMs. The discussion covers not only the potential negative impacts but also the practical difficulties and risks in real world environments.

1.1. Digital forensic context

DF is a process for identifying, preserving, analysing, and documenting digitally recorded data, which originate in electronic devices such as computers, servers, smartphones, and IoT devices (Baryamureeba and Tushabe, 2004). This exercise is required in most criminal cases. Data collected in this process are kept unchanged and safe for presentation in a court case or to support future investigations conducted by law enforcement agencies (Mukherjee and Haque, 2018).

1.1.1. Digital forensic process models

DF process models consist of a series of activities that help standardise the investigative process (Du et al., 2017) and outline the phases: collection, preservation of evidence, examination or analysis, and reporting.

DF encompasses various subdisciplines such as computer forensics, mobile device forensics, memory forensics, network forensics, and cloud forensics, each employing distinct processes reflected in a plethora of models within the literature (Scanlon et al., 2023a; Wu et al., 2020). These models often share phases but differ in their focus and execution. For example, Al-Dhaqm et al. (2020) proposed a mobile forensic model that adds a preparatory phase and bifurcates the analysis stage into examination and analysis phases. To accommodate the complexities of computer, network, cloud, and smart device forensics, Lutui (2016) introduced a multidisciplinary model that requires diverse skills for effective investigation, ranging from incident detection to evidence storage.

Casey's model, as shown in Fig. 1, includes phases such as incident recognition, evidence collection, preservation, and presentation, with the examination phase detailed in recovery, harvesting, reduction, and

classification (Casey, 2011). During incident recognition, the focus is on identifying the incident itself, possible evidentiary sources, and expected digital evidence types, as well as delineating the scope of the ensuing investigation. Conversely, investigators systematically acquire pertinent evidence from various sources encompassing computers, smartphones, storage media, and networks during collection and seizure. Preservation is of paramount importance in upholding the integrity of evidence, necessitating specific and accurate measures to ensure the unmodified condition of the collected data throughout the investigative process. The overarching objective remains the meticulous safeguarding of evidence integrity.

The subsequent phase entails examination, in which analysts rigorously scrutinise the gathered data to extract pertinent information. This endeavour may involve the use of various forensic hardware and software tools and techniques. The examination process includes the interpretation of the information extracted to draw conclusive inferences about the events under scrutiny. This phase often demands a profound understanding of both the technology employed and the context surrounding the evidence.

Next is the reporting phase, where the findings derived from the analysis are presented systematically in a format suitable for legal adjudication. This may involve preparing comprehensive reports and providing expert testimony in a court of law. This model emphasises the critical nature of maintaining the integrity and provenance of the evidence and the requirement for expert analysis to an automated and interpret pertinent information. This culminates in a report suitable for legal scrutiny. In particular, the analysis or examination stage is crucial in all models, demanding specialised knowledge in the relevant DF area (Mir et al., 2016; Prayudi et al., 2020).

The advent of cloud computing has led to the Digital Forensic as a Service (DFaaS) model by van Baar et al. (2014), which integrates evidence preservation and analysis into an automated and secure software service, marking a significant evolution in forensic methodologies (van Beek et al., 2015; Du et al., 2017; van Beek et al., 2020).

1.1.2. Existing challenges in digital forensics

DF is an evolving field, yet the literature highlights that it still undergoes changes to address ongoing challenges and advancements. Dubey et al. (2023) assert that DF faces key challenges, including the complexity of data and its volume, a lack of standardisation, inadequacies in the power of existing tools to support investigations, and issues related to timelines.

In addition to the previously mentioned challenges, other issues persist including scope creep in cases due to complexity and the vast data involved, selecting and prioritising the right set of evidence, and efficiently allocating time and investigators for the chosen evidence (Kalaimannan et al., 2013). Koper et al. (2014) focus on a number of these issues from the investigator's perspective, including challenges in adapting to a system, unexpected time sinks, and frustrations among officers arising from expected operational timelines and the adoption of complex systems. The contemporary landscape of forensic science is characterised by a notable deficiency in proficient agents, exacerbated by the swiftly evolving standards, practices, tools, and techniques within the field. Moreover, the predominant emphasis of law enforcement roles on fieldwork, as opposed to dedicated DF duties, has further contributed to the scarcity of adept human expertise in this domain (Vincze, 2016).

Automating the DF process using existing technology appears to be a promising solution to address issues related to time management and effectiveness (Michelet et al., 2023). However, an ongoing challenge revolves around measuring the accuracy of investigations and ensuring the verification of the automated process. This aspect remains an open area that requires more attention and resolution (Jarrett and Choo, 2021).

1.1.3. Existing work with LLMs in digital forensics

Scanlon et al. (2023a) analysed using ChatGPT for DF. In their assessment, the authors evaluated the programming, incident narration,

keyword list creation, and DF teaching abilities of ChatGPT. Their conclusion highlighted that while ChatGPT exhibited some hallucinations in the output results, it still serves as an effective assistant for code generation. Wickramasekara et al. (2025) introduced the AutoDFBench benchmarking framework, and corresponding score, to evaluate code generation for DF specific tasks against the tests and datasets used as part of NIST's Computer Forensics Tool Testing Program (CFTT).³

Timeline reconstruction helps investigators deduce the chronological "story" of an event. In line with timeline regeneration, Silalahi et al. (2023) proposed a method to detect anomalies in a drone flight by employing sentiment analysis with the assistance of a pre-trained LLM. Their approach successfully discerned the differences between normal and abnormal events with an accuracy of 92.5%.

Hansken, a DFaaS platform created by the Netherlands Forensic Institute, is designed to help investigators handle evidence and conduct investigations more efficiently (van Beek et al., 2015). ChatGPT has been used as an assistance for the Hansken DFaaS system using its bespoke query language, contributing to streamlined processes and improved support for investigators. In these experiments by Henseler and van Beek (2023), it was tasked with analysing evidence using Hansken's trace model. This work demonstrated the potential for ChatGPT in helping with analytical aspects of investigations, highlighting its ability to process and interpret evidence data.

While DF is the main focus of this paper, it is important to recognise the broader application of LLMs in adjacent areas within cybersecurity, many of which overlap with DF. LLMs are proving to be valuable tools in fields such as malware analysis, security log analysis, code security reviews, and intrusion detection areas that bridge the gap between cybersecurity and DF (Yu et al., 2024a; Karlsen et al., 2024; Lira et al., 2024). In malware analysis, LLMs can identify patterns in malicious code, while in log analysis, they assist in detecting anomalies across large datasets, thereby improving response times. In code-related security reviews, LLMs like GRACE have demonstrated the ability to identify vulnerabilities in software, achieving a detection rate of 28.65% of vulnerabilities (Lu et al., 2024). These applications contribute to DF investigations by improving and improving evidence collection and analysis.

2. Background of large language models

This section explores LLMs, concentrating on three principal aspects. Initially, it explores the architecture of LLMs, detailing their design and function. Then it assesses the usability of LLMs, underscoring the features and capabilities that make them suitable for a wide range of tasks. Finally, it showcases the versatility of LLMs by discussing their applications across various fields, demonstrating their wide-reaching impact and the extensive scope of their applications.

2.1. Natural language processing

Popular LLMs such as Generative Pre-trained Transformer (GPT) (Dale, 2021), Language Model for Dialogue Applications (LaMDA) (O'Leary, 2023), Pathways Language Model (PaLM) (Chowdhery et al., 2023), Bidirectional Encoder Representations from Transformers (BERT) (Devlin et al., 2019), and Large Language Model Meta AI (LLaMA) (Touvron et al., 2023) stem from advances in Natural Language Processing (NLP). NLP, which focusses on language-based tasks, uses traditional and deep learning models to enable applications such as language translation, text processing, and speech recognition (Reshamwala et al., 2013).

Deep learning, a branch of machine learning, uses complex computational layers and adaptive weights to improve prediction accuracy, offering a more refined analysis than conventional machine learning (Lecun et al., 2015). It has excelled in image and speech recognition and

natural language understanding, mimicking the decision making process of the human brain through artificial neurons. These neurons form networks capable of intricate pattern recognition and data analysis. Central to deep learning are neural networks with multiple hidden layers that autonomously learn and extract features from data, bypassing the need for manual variable selection. This automatic feature extraction makes them exceptionally adept at handling complex tasks (Dong et al., 2021).

2.2. LLMs

An LLM is a language model that employs neural networks with billions of parameters, trained on extensive text data. These models are engineered to comprehend and generate human language. Fundamentally, they rely on multiple neural network architectures, enabling them to recognise the relationships between words and phrases within sentences (Yang et al., 2024a; Shen et al., 2023). These architectures have been a transformative force in natural language processing. Its capability to excel across a diverse array of language-related tasks distinguishes it as a game-changer, in contrast to being tailored for a singular, specific task.

2.3. Architecture of LLMs

LLMs utilise deep learning, particularly neural networks, to process and produce human language. Fundamentally, a language model operates with letters or words, but since machine learning algorithms and neural networks require vector inputs, words are vectorised. Each word in the vocabulary is assigned a unique numerical value for input into neural networks. Through initial random weight assignments and subsequent backpropagation, words acquire numerical positions reflecting their semantic similarity, culminating in a word embedding model (Lai et al., 2016).

2.3.1. Word to vectors

Word embeddings, as introduced by Mikolov et al., entail precise and high-dimensional vector representations for words, particularly suited for extensive datasets comprising billions of text entries. The authors explored model architectures for word vectorisation, achieving substantial improvements in accuracy while requiring fewer computational resources and reduced training time (Mikolov et al., 2013). In the realm of LLMs, the primary objective is to generate new text based on the extensive dataset on which it was trained. For this purpose, Vaswani et al. (2017) introduced the transformer model, assisted by the word-to-vector model. This architecture incorporates a self-attention mechanism, as well as encoder and decoder processes, enabling the model to quickly and simultaneously focus on pertinent information.

2.3.2. Transformer models

The transformer model initially aimed at machine translation, translating input words into another language, begins with word embedding, where input, termed tokens, are vectorised. Recognising word order is achieved through positional encoding, with two main techniques: absolute and relative. Absolute positional encoding assigns unique vectors to each position, enhancing the model's ability to recognise word placement and facilitate position-specific attention (Ke et al., 2021). Relative positional encoding, on the other hand, calculates the relative positions of words by introducing a bias term that quantifies the distances between positions, improving the model's ability to understand the relationships between words within a sequence (Ke et al., 2021).

Self-attention, a core mechanism within the transformer, calculates the relationship among words in a sentence, allowing the model to assess the similarity of each word with others and generate unique representations for each (Al-Rfou et al., 2019). The decoder mirrors the encoder's steps but uses different weights, starting with positional encoding and computing self-attention values to identify the sentence's initial translation word.

³ <https://www.nist.gov/itl/ssd/software-quality-group/computer-forensics-tool-testing-program-cftt>.

This transformer process, which leverages stacking self-attention and unique positional encoding, has significantly advanced NLP tasks, including machine translation, text generation, and summarisation, by executing these processes in parallel and optimising the weights of both the encoder and decoder (Vaswani et al., 2017; Acheampong et al., 2021).

2.4. Specifically trained LLMs

The transformer model and the self-attention mechanism have paved the way for researchers to train language models on trillions of tokens with billions of parameters. Several LLMs have been trained and harnessed, each designed with specific capabilities for diverse fields such as security (Lira et al., 2024), chemistry (Bran et al., 2024; Tsai et al., 2023), engineering (Hou et al., 2024), medicine (Chang et al., 2024), business (Vidgof et al., 2023), tourism (Carvalho and Ivanov, 2024), and language-related applications (Chang et al., 2024). These models are used in tasks ranging from detecting security threats (Lira et al., 2024), analysing data and generating synthetic actions to teaching (Moore et al., 2023), code generation (Hou et al., 2024; Liu et al., 2023a), structured query generation (Vidgof et al., 2023; Li et al., 2024a), planning (Vidgof et al., 2023), assisting in medical education (Chang et al., 2024), clinical decision making (Eggmann et al., 2023), leveraging clinical settings (Thirunavukarasu et al., 2023), clinical validation (Karabacak and Margetis, 2023), understanding general patterns and decision making (Vidgof et al., 2023), bias detection (Moore et al., 2023), addressing ethical issues (Bonner et al., 2023), language translations (Caines et al., 2023), question answering (Moore et al., 2023), information extraction (Liu et al., 2024a), and business process automation (Vidgof et al., 2023), among others (Chang et al., 2024). The fine-tuning and retraining capabilities of LLMs enable them to be adapted to specific tasks or behaviours in a predefined manner. Fine-tuning involves taking an already trained language model and retraining its existing weights and bias values using a new dataset specific to a particular domain. This process allows the LLM to be customised and refined for tasks beyond its original training, enhancing its applicability in specific contexts (Bakker et al., 2022). This process results in a new model that is more tailored and focused on the specified domain. In the existing literature, it is frequently observed that LLMs are fine-tuned with a particular emphasis on engineering and research-related fields. This targeted fine-tuning ensures that the model is adept at handling tasks and generating content specifically relevant to the intricacies of these domains (Chang et al., 2024).

2.5. Multi-modal large language models (MLLMs)

Unlike traditional LLMs, which are trained on text data, Multi-modal Large Language Models (MLLMs) are designed to process and interpret image-based data alongside text. These models can extract and analyse information within images or videos, integrating visual and textual data to enhance comprehension and analysis (Tan et al., 2024). The application of MLLMs is rapidly expanding across fields such as digital forensics, where they can be used to analyse images of documents or identify visual anomalies. Section 3.2 provides a discussion on the potential and diverse applications of MLLMs in various domains.

2.6. Large action models

While LLMs excel in text generation and processing, they struggle with complex task manipulation and operational control, especially when moving from language understanding to action execution. This limitation arises because their core design emphasises prediction and generation over direct task execution. To overcome this shortcoming, recent research has introduced innovative approaches, such as the Large Action Model (LAM) developed by the Rabbit research team. LAM

extends the capabilities of LLMs by incorporating action-based operations.⁴ These LAMs can mimic human routines such as scheduling meetings with given instructions, sending emails, ordering taxis, and handling complex tasks such as making reservations for a whole trip. In this approach, the base model is trained to comprehend sequences of human-provided actions and commands, allowing it to execute these actions and tasks accordingly. Similarly, Microsoft introduced the concept of Visualisation-of-Thought (VoT) aimed at integrating human cognitive abilities, specifically the creation of mental images, into the model (Wu et al., 2024). Through this approach, it has been demonstrated that MLLMs excel in visual tasks, thereby enabling the extension of action capabilities within an LAM to any LLM. These advancements signify promising directions toward enhancing the practical applicability and versatility of language models across various domains.

3. Capabilities of large language models

This section focusses on the abilities and capabilities of Language Model Models (LLMs) as outlined in Section 2.4. This section also discusses the currently available fine-tuned LLMs that exhibit potential for application in DF. Although considered too broad for this article, Zhao et al. (2024) provide a detailed generic overview of LLMs, their operation, and how they are trained and fine-tuned.

3.1. Programming/scripting/code generation

The ability to generate source code within a specific context is a crucial skill inherent in a language model (Alon et al., 2020). Xu et al. (2022) conducted a systematic evaluation of six LLMs for code generation in 12 different programming languages. The benchmarking process employed the HumanEval benchmark, along with a tailored evaluation dataset designed to assess the functional correctness of the programs generated by an LLM (Chen et al., 2021). The benchmark comprises a set of coding problems in which the model is tasked with generating Python functions. Each problem is accompanied by a prompt and a set of unit tests that verify whether the generated code produces the expected output. This facilitates a systematic evaluation to generate both syntactically correct and functionally accurate programs. Using this dataset, it is possible to measure performance on real-world coding tasks, as well as its ability to generate solutions that satisfy functional requirements.

The Mostly Basic Programming Problems (MBPP) is another benchmark comprising 974 programming tasks. It serves as a frequently used evaluation dataset for LLMs specialising in code-related tasks (Wei et al., 2022). Several LLMs explicitly trained for code generation include Code LLaMA (Rozière et al., 2023), CodeGen (Nijkamp et al., 2022), StarCoder (Li et al., 2023a), PanGu-Coder (Yu et al., 2024b), PanGu-Coder2 (Yu et al., 2024b), WizardCoder (Luo et al., 2024), InCoder 6B (Fried et al., 2023), CodeGen-Mono 16B (Nijkamp et al., 2022), Code-Davinci-001 (Zhou et al., 2022), Code-Davinci-002 (Zhou et al., 2022), PaLM-Coder-540B (Chowdhery et al., 2023), CodeT5+ (Wang et al., 2023a, 2021), InstructCodeT5+ (Chen et al., 2023; Wang et al., 2021), GPT-4 with Reflexion (Shinn et al., 2023), CodeGeeX (Zheng et al., 2023), AlphaCode (Li et al., 2022), Santa-Coder (Allal et al., 2023), CodeFuse-13B (Di et al., 2024), Codex (Kalyan, 2024), WaveCoder (Yu et al., 2024c).

A higher score for both HumanEval and MBPP indicates high precision in code generation for a given task. Table 1 presents the scores for HumanEval and MBPP for each of the code generation LLMs mentioned above, along with the trained parameter size for each LLM. Four generic LLMs, or Mixture-of-Experts (MoE) models, language models are also included in Table 1: o1-mini (Yu et al., 2024d), GPT-4 with Reflexion, DeepSeek-V3 and DeepSeek-V3-Base (DeepSeek-AI et al., 2024). These are included as these are the top 4 best performing models for HumanEval despite them being MoE models.

⁴ <https://www.rabbit.tech/research>.

Table 1

Trained parameter count, HumanEval and MBPP scores for LLM based code generation (ordered by HumanEval score).

Model	Parameters	HumanEval	MBPP
o1-mini ¹	100B	97.6	93.9
GPT-4 with Reflexion	1.76T ²	91.0	77.1
DeepSeek-V3	671B	85.6	-
DeepSeek-V3-Base	671B	65.2	75.4
Code LLaMA	34B	62.2	61.2
PanGu-Coder2	15B	61.64	-
WizardCoder	15B	57.3	51.8
Code-Davinci-002 (GPT3.5)	175B	47.0	58.1
StarCoder	15.5B	40.8	49.5
Code-Davinci-001 (GPT3)	175B	39.0	51.8
PaLM-Coder	540B	36.0	47.0
InstructCodeT5+	16B	35.0	-
code-cushman-001	12B	33.5	45.9
CodeT5+	16B	30.9	-
CodeGen-MONO	16.1B	29.7	42.4
CodeGen	16.1B	29.28	35.28
Codex-12B	12B	28.81	-
PanGu-Coder	2.6B	27.78	23
Sanata-Coder	1.1B	18	35
AlphaCode	1.1B	17.1	-
InCoder 6B	6.7B	16.4	21.3

¹ Yu et al. (2024d).

² Estimated parameter count as value is not officially released (Rizzo et al., 2024).

3.2. Vision assistance

Traditional vision assistant systems face limitations in image processing or recognition, as they are typically trained on fixed types of datasets. However, with the emergence of LLMs, this paradigm has changed to the use of raw text as a source of supervision (Radford et al., 2021b; Tewel et al., 2022). Research on visual recognition language models is experiencing exponential growth, with the number of models exceeding 1,500 in 2023 (Zhang et al., 2024). Radford et al. (2021b) introduced a novel method called Contrastive Language Image Pre-training (CLIP). This method is efficient and capable of performing a wide range of tasks during pre-training. It enables a model to learn a shared representation space for both images and text, facilitating a deeper understanding of the relationships between the two modalities. Ramesh et al. (2021) proposes a model for text-to-image generation, capable of generating images as combinations derived from textual input or sentences. Moreover, with the model named Generating Images with Large Language Models (GILL), it becomes feasible to generate text, retrieve images, generate novel images, and interleave the results into coherent multimodal dialogues (Koh et al., 2023). VisionLLM is a framework leveraging LLMs for diverse vision tasks with unified language instruction, demonstrating generality and flexibility (Wang et al., 2023b). It incorporates a language-guided image tokeniser and an LLM-based task decoder, capable of handling open-ended tasks based on provided language instructions (Wang et al., 2023b).

Visual instruction tuning leverages language-only models, such as GPT-4, to generate multimodal language-image instruction following data. This data is then utilised to instruction-tune large multimodal models, such as Large Language and Vision Assistant (LLaVA) (Liu et al., 2023b, 2024b; Achiam et al., 2023). The open source LLaVA project introduces an end-to-end trained model, integrating a vision encoder with an LLM. Notably, LLaVA showcases multimodal chat capabilities. LLaVA has the capability to interact with images, provide detailed descriptions and respond to queries with a reported accuracy of 92.53% (Liu et al., 2023b). This shows its effectiveness in understanding and generating contextually relevant information on visual content. MiniGPT-4 is an open-source, powerful visual instruction-tuned LLM, and it demonstrates versatility by generating stories and poems inspired by provided images and teaching users how to cook based on visual cues from food

photos. This showcases its ability to understand and respond creatively to various visual stimuli (Zhu et al., 2024).

Position-Enhanced Visual Instruction Tuning (PVIT) represents an extended version of Multimodal Large Language Models (MLLMs). It facilitates region-level encoding in an image, enabling the model to discern and identify information within specific regions (Guo et al., 2024a). This model enables users to interact with both the language and drawing the bounding boxes to indicate the area of interest within an image (Guo et al., 2024a). Other MLLMs, such as Visual ChatGPT (Radford et al., 2021b), InternGPT (Thapa and Patil, 2024), Flamingo (Alayrac et al., 2022), BLIP-2 (Li et al., 2023b), and Kosmos (Wu et al., 2023a), are noted in the literature for their ability to assist users in visual-related information.

Video information is gaining prominence in vision assistance, and Zhao et al. (2023) has introduced a novel approach to automatically narrate lengthy videos using LLMs. UniViLM is another language pre-trained model designed for both multimodal understanding and generation. It is capable of retrieving a video segment based on text descriptions, generating captions for given video clips, segmenting a video according to a provided text input, and performing multimodal sentiment analysis of a video segment (Park and Kim, 2023). VidIL and LLaViLo are additional MLLMs with similar capabilities, demonstrating proficiency in video classification and video-language operations such as video captioning, video question answering, video caption retrieval, and prediction of future video events (Wang et al., 2022; Ma et al., 2023).

These MLLMs adhere to a shared task set, that includes visual question answering, visual captioning, visual common-sense reasoning, visual generation, multimodal affective computing, visual retrieval, vision language navigation, multimodal machine translation, visual question generation and visual dialoguing, as summarised in Table 2 (Kiros et al., 2014; Uppal et al., 2022).

3.3. Conversation

Specific LLMs are trained explicitly for meaningful and coherent dialogues with humans. An example is Dialogue Generative Pre-trained Transformer (DialogPT), a fine-tuned model trained on 174 million Reddit conversations (Zhang et al., 2020). DialogPT exhibits the ability to provide human-like answers in tested conversations (Zhang et al., 2020). Dettmers et al. (2023) introduced a fine-tuning mechanism for LLMs named Quantised Pre-trained Language Model into Low-Rank Adapters (QLoRA). This allows for the fine-tuning of large-parameter LLMs with low training costs. They introduced Guanaco, a fine-tuned LLM with 65 billion parameters, which achieved a performance level of 99.3%. Falcon-180B and Falcon-40B represent another set of open-source LLMs with 180 billion and 40 billion parameters. These models are trained to communicate in multiple languages, allowing users to engage in conversations in languages other than English (Penedo et al., 2023). To evaluate the accuracy of human-like dialogue systems, Ou et al. (2024) proposed a dialogue evaluation benchmark named DialogBench, which consists of 12 dialogue tasks to assess the capabilities of LLMs. In their evaluation, they assessed 28 pre-trained and instruction-tuned LLMs, demonstrating that GPT-4, ChatGPT, and KwaiYi-13B-Chat emerged as the top three models for conversations in domains related to daily life and professional knowledge.

In DF chat conversations, the significance lies in facilitating non-technical investigators to elucidate terminologies and areas lacking understanding. This serves a dual purpose, acting as an interactive teacher to enhance comprehension in discussions (Scanlon et al., 2023a).

3.4. Prompt engineering

Achieving quality outputs from LLMs often relies on providing well-crafted, meaningful, and precise input queries, known as input prompts. However, even human-defined natural language instructions may not

Table 2
Capabilities of MLLMs trained for vision assistance.

Capabilities		GILL	VisionLLM	GPT-4	LLaVa	MiniGPT-4	Visual ChatGPT	InternGPT with Husky	Flamingo	Kosmos	Uni VL	VidLL
Visual question answering (Task of providing an answer to a visual input.)	Image	✓	✓	✓	✓	✓	✓	✓	✓	✓		
Visual captioning (Task of generate fitting visual descriptions.)	Image	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Visual common-sense reasoning (Task of infer understanding from images or video clip.)	Image	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Visual generation (Task of generating image or video from a given textual input.)	Image	✓					✓	✓				
Multimodal affective computing (Task of automatically recognition of emotions and causes.)	Image	✓	✓	✓	✓	✓	✓	✓	✓	✓		
Visual retrieval (Task of language and vision understanding and retrieval.)	Image	✓	✓				✓	✓				
Vision-language navigation (Task of navigation based on linguistic instructions.)	Image	✓					✓	✓			✓	
Multimodal machine translation (Task of translation from a video or an image)	Image		✓	✓	✓	✓	✓		✓	✓		
Visual question generation (Task of generating questions for given image or video)	Image		✓	✓	✓	✓	✓	✓	✓	✓		
Visual dialoguing (Task of automating a conversation about a video or image)	Image	✓		✓	✓	✓	✓	✓	✓	✓		✓
	Video											✓

consistently yield the best results. Prompt engineering is a methodology that involves carefully defining and instructing LLMs to generate more accurate and desirable outputs. Through thoughtful refinement of input prompts, prompt engineering aims to enhance the performance and effectiveness of LLMs in generating output that align more closely with user expectations and requirements (Marvin et al., 2024; Zhou et al., 2023). This plays a crucial role in biasing LLMs toward specific domains or topics, enabling a more targeted and nuanced response. By carefully crafting prompts, users can guide LLMs to dive deeper into the nuances of their queries, leading to more accurate and relevant outputs. This approach enhances the model's responsiveness to specific areas of interest, allowing users to fine-tune and tailor their interactions with the LLM for more precise and meaningful outcomes. Prompt engineering with LLMs is employed in various sectors, including, but not limited to, medical, engineering, construction, and healthcare (Meskó, 2023; Polak and Morgan, 2024).

ChainForge is an open source Graphical User Interface (GUI) tool developed specifically for prompt engineering and hypothesis testing derived from LLMs that can be used in the aforementioned fields to generate accurate and quick output (Arawjo et al., 2024).

Although prompt engineering is a necessity in generating the desired output from an LLM, the results can still be biased based on the wording and phrasing provided by the user. Since the effectiveness of prompts depends on the user's proficiency in English, the output may vary significantly depending on the exact requirements of the user and how the LLM interprets these prompts. In addition, the prompts can unintentionally reinforce existing biases within the model's training data, potentially skewing the results. Therefore, prompt engineering must be approached carefully and methodically to minimise misinterpretation and maximise output relevance.

3.5. Autonomous agents

The evolution of LLMs, with their ability to generate information and communicate in a manner that resembles human interaction, has led to the development of autonomous agents. The expectation is that these agents will effectively perform a wide range of tasks, taking advantage of the human-like capabilities inherent in LLMs (Wang et al., 2024a). These autonomous agents follow a four-stage architecture that includes profiling, memory, planning, and action. Profiling defines the agent's role, privileges, domain, and expertise (Wang et al., 2024a). Memory stores

information on tasks and profile data relevant to the environment. Planning involves breaking down given tasks into subtasks and solving them individually. The action stage is the final phase where all decisions and subtasks are translated into actions executed by the agent. Zhang et al. (2023) developed a framework designed to facilitate collaboration between GenAI agents and humans. This framework enables planning and communication for specific tasks, leveraging the capabilities of LLMs. Similarly, AgentSims (Li et al., 2024b), ToolBench (Qin et al., 2024b), GameGPT (Sweetser, 2024), ChatDev (Qian et al., 2024), Voyager (Qin et al., 2024a), and RecMind (Wang et al., 2024b) represent a diverse array of autonomous GenAI agents developed with distinct goals and objectives. Certainly, AutoGen stands out as a multiagent framework with the capability to autonomously perform tasks or collaborate with human feedback. This flexibility makes it a versatile tool for various applications (Wu et al., 2023b).

In addition to the AutoGen framework, AgentLite (Bajwa et al., 2007), Camel (Li et al., 2023c) and CrewAI (Barbarroxa et al., 2025) are each similar LLM-based agent framework architectures. These platforms are distinguished by their support for task decomposition, multi-agent orchestration, and adaptable reasoning. In particular, AgentLite and CrewAI facilitate work delegation functionalities, increasing their utility in various operational contexts.

3.6. Retrieval-augmented generation

The content generated by LLMs is highly dependent on the extensive text-based datasets on which they are trained. These datasets may contain vast amounts of information utilised by the LLMs. However, maintaining up-to-date knowledge within LLMs is challenging, as fine-tuning or retraining a model is often extremely costly and resource-intensive. Retrieval-augmented Generation (RAG) is a technique designed to address this LLM knowledge gap by retrieving information from external sources and integrating it with the model's internal representations (Fan et al., 2024).

A major advantage of RAG is its ability to reduce the hallucination problem in LLMs, allowing them to generate more accurate and current information (Jiang et al., 2023; Lewis et al., 2020b). The architecture of RAG consists of a knowledge base and a retriever model. The retriever model converts input prompts and content from the knowledge base to vectors. The user prompts are then appended with the most relevant content from the knowledge base, and this augmented prompt is sent

to the base model to generate a more accurate response (Lewis et al., 2020b).

3.7. Limitations and risks

As explored in the preceding sections, LLMs appear to possess a vast range of capabilities. However, it is crucial to acknowledge that they are not without limitations and risks. In multimodal LLMs, it is a common problem that they are over-reliant (Subramonyam et al., 2024). There are also potentially significant drawbacks associated with LLMs, including issues such as bias, explainability challenges, reasoning errors, logical errors, hallucinations, vulnerability to prompt injections, and spelling and grammar errors. These limitations underscore the importance of a cautious and critical approach when using LLMs in various applications. Furthermore, the literature shows limitations in LLMs, including statistical inconsistency, the absence of emotional attributes in linguistic responses, and challenges related to fact verification (Tang et al., 2024; Fröhling and Zubiaga, 2021). These factors contribute to a comprehensive understanding of the constraints and potential shortcomings when working with LLMs. Thapa et al. (2023) contend that while LLMs can indeed reduce the time and costs associated with annotation tasks, they are not completely supplanting human annotation. This is because they struggle with intricate linguistic constructions, such as idioms, irony, sarcasm, and metaphor, which can potentially impact the precision of annotations.

Similar limitations are associated with MLLMs, such as over-reliance on training data, sensitivity to word order in input prompts, and vulnerability to prompts containing additional knowledge (Qi et al., 2023). There are several more concerns associated with LLMs, including restricted text input and output lengths, limited comprehension of syntax, ethical considerations with the generated information, constraints with multilingual capabilities, elevated costs associated with training and maintenance, inadequate understanding of human behaviours and limited ability to learn incrementally (Chiang and Lee, 2023; Alawida et al., 2023).

Despite their considerable capabilities, LLMs are not without risks. Wiggins and Tejani (2022); Lund and Wang (2023); Rahman and Santacana (2023) provide comprehensive overviews of risks linked to LLMs. These include the homogenisation of results, whereby defects or biases from the foundation model are inherited by all downstream models. There is also the risk of monopolistic control by foundation model owners, potentially concentrating decision making power, resource access, and influence over model usage in a single entity. Ethical and legal concerns are intertwined with concerns about privacy and intellectual property. In addition, there are economic and environmental impacts that raise concerns about the potential displacement of human workers. Furthermore, inequity and misuse of LLMs, such as the creation of deepfakes and their application in criminal and unethical activities, pose additional challenges. Given that LLMs do not inherently prioritise the precision of information, Bender et al. (2021) have highlighted the risk of generating social turbulence, especially when used on social media platforms. Furthermore, the use of LLMs is associated with significant costs, leading to a direct environmental impact due to their substantial energy consumption (Rillig et al., 2023).

The risks associated with LLMs are predominantly emphasised within Information Communication and Technology (ICT) and cyberspace. Primary concerns include the disclosure of personal information, the generation of malicious text, and the creation of malicious code (Rao, 2023).

The Beyond the Imitation Game benchmark (BIG-bench), serves as an evaluation framework for LLMs. It encompasses 204 distinct language-related tasks. These span contextual and context-free question-answering, reading comprehension, logical reasoning, etc. (Srivastava et al., 2023). It is acknowledged that the challenge of social biases and dependency on the English language persists in almost all LLMs.

When employing LLM-based agents, it is imperative to address challenges associated with LLM-based multi-agent frameworks as well. The handling of many defined agents may necessitate substantial computational resources and memory, thus mandating high-end computing infrastructure for seamless operations. Furthermore, the absence of a standardised comprehensive benchmarking system to evaluate the behaviour of such agents underscores the limitations inherent in the development of LLM-based multi-agent systems. These challenges underscore the need for further research and refinement in this domain to enhance the efficiency and effectiveness of LLM-based multi-agent frameworks (Guo et al., 2024b).

Similarly to these risks, the Open Web Application Security Project (OWASP) has identified ten major risk factors related to LLMs. These risks include training data poisoning, prompt injection, denial of service, insecure output handling, supply chain vulnerabilities, sensitive information leakage, excessive agency, insecure plugins, overreliance, and model theft of data. Despite these threat factors, OWASP also stressed the need for regulatory bodies to supervise LLMs in various domains and recommended the implementation of risk management programmes that incorporate the checklist provided by OWASP.⁵

The EU's Artificial Intelligence Act (AIA) proposes a framework for categorising AI applications based on their associated risk levels, with the primary aim of safeguarding human rights and maintaining ethical standards in AI deployment (Neuwirth, 2023). Within the AIA, AI applications are divided into categories such as "unacceptable risk", which includes practices such as exploiting vulnerabilities and social-ranking techniques due to their potential for individual manipulation and impact on fundamental rights. These categories have relevance to DF, where the ethical application of LLMs must balance the advantages of automation with the imperative to address privacy concerns. Since DF involves sensitive data and influences legal outcomes, the use of LLMs must align with AIA's risk-based principles, ensuring transparency, accountability, and fair application. In future AI applications within DF, it is essential to implement appropriate measures to prevent biases and hallucinations to mitigate the risk of misuse of AI.

4. Large language models for digital forensics

Section 4 summarises existing work with LLMs in DF, the feasibility of employing them, and potential future directions. As discussed in Sections 2 and 3, despite the widespread use of LLMs in various fields to improve the efficiency and accuracy of tasks within specific domains, their application in the field of DF is still relatively new.

Conducting a thorough analysis of LLMs use in conjunction with the stages of the DF process model, as highlighted in Section 1.1.1, proves to be a valuable undertaking.

4.1. Incident recognition phase

In the initial phase of Casey's DF process model, which delineates the recognition of an incident, LLMs can serve as a valuable detection mechanism (Goel et al., 2024). In cybercrime cases, the primary artefacts often involve data logs, data dumps and network dumps. Fine-tuning an LLM to monitor text-based logs and related files enables it to discern and identify potential or ongoing incidents within the environment. In network-related activities, anomaly detection plays a pivotal role in initiating an incident response. Various existing anomaly detection techniques are employed in systems for this purpose. Using their ability to identify patterns in a series of text data sets, LLMs exhibit potential as an Intrusion Detection System (IDS) within such systems (Lira et al., 2024; Yang et al., 2024b). For instance, Kan et al. (2024) introduces Mobile-LLaMA, an open source mobile network-specialised LLM,

⁵ <https://owasp.org/www-project-top-10-for-large-language-model-applications>.

fine-tuned through instructional data to enhance their capabilities for network analysis tasks within 5G environments. Mobile-LLaMA supports three primary functions: IP routing analysis, packet analysis, and performance analysis.

4.2. Collection phase

Although evidence collection or seizure traditionally involves physical tasks that require human interaction, LLMs can play a role in identifying and listing potential pieces of evidence at a crime scene. For example, in the examination of photographs or video records from a crime scene, an investigator can enlist the help of a MLLM such as LLaVa, GPT-4, or VisionLLM. These models are capable of processing information within the images and generating a text-based output, facilitating the interpretation and categorisation of visual data. Although this task may seem simple and within the capabilities of a human agent, the efficiency becomes particularly evident when dealing with a massive-scale investigation involving thousands of collected artefacts and photographs. Using an MLLM for initial processing can significantly save time, with human agents then focussing on the crucial task of verification and validation.

4.3. Preservation/acquisition phase

The preservation of evidence is centred on maintaining integrity. To achieve this, various tools such as EnCase and FTK Imager have been used, helping investigators streamline their work processes (Shah et al., 2017). In the context of non-technical DF stakeholders being able to interrogate the evidence, it becomes feasible for a user to articulate their requirements/query in natural language. Subsequently, the LLM generates source code tailored to the specific need, executes the code on the data, and returns the result in consumable natural language.

LLMs specialised in code generation, such as StarCoder and Code LLaMA, can be fine-tuned and retrained for domain-specific tasks, including the preservation of disk evidence through customised code and script generation. These LLMs are capable of generating scripts or code snippets that create secure copies of disk images, metadata, and partition information, as well as automating cryptographic hashing and verification routines to maintain the evidence's integrity through checksums. Additionally, LLMs can assist in documenting preservation steps by generating logs and summaries for each stage of the disk preservation process, thereby supporting the chain of custody during acquisitions. However, despite these capabilities, human expertise remains essential for identifying and collecting potential sources of evidence during the preservation phase, as the application of LLMs in this stage is currently limited to lower-potential tasks.

In certain instances, the gathering of live data for forensic investigations becomes crucial, particularly data collected at the crime scene. For this purpose, investigators can use DFaaS platforms such as Hansken. Hansken possesses the ability to amalgamate custom extraction APIs for data extractions, and these APIs can be developed using code-generative LLMs (van Beek et al., 2015). This approach improves the adaptability and efficiency of the investigative process.

As stated in Section 3.5, the automation of code generation and unit testing can be facilitated by autonomous agents that use LLMs as their core. AutoGen, being an open source framework, provides the means to develop AI agents tailored for specific tasks. These AutoGen agents are not only customisable and conversational, but can also operate in various modes, employing combinations of LLMs, human inputs, and various tools (Wu et al., 2023b). Automated agents, particularly those developed within frameworks such as AutoGen, can be used in the preservation phase of investigations. These agents can be assigned specific tasks, such as acquiring disk images, generating disk hashes, retrieving disk metadata, and compiling acquisition reports. By defining precise roles and tasks for AI agents, it is possible to streamline and standardise

these preservation actions, improving the management of digital evidence (Wu et al., 2023b; Wickramasekara and Scanlon, 2024).

4.4. Examination phase

This phase constitutes a crucial component of the investigation, playing a crucial role in elucidating the case through activities such as data recovery, collection, reduction, and classification. For each of these components, LLMs fine-tuned for scripting can significantly assist, especially at a larger scale. Within these components, tasks such as keyword search, file recovery, pattern matching, and fragment reassembly can be achieved with minimal technical knowledge using LLMs. LLMs can provide valuable assistance in these tasks by generating new codes, crafting regular expressions, generating passwords and/or password hash lists for decryption, and creating sample logs or files. LLMs can generate a set of instructions, queries, and Application Programming Interface (API) validations from natural language provided by a human. This opens up the possibility of integrating third-party tools like Scapy, tshark, John the Ripper, and others seamlessly into the investigative process, enhancing the toolkit available for DF investigations, and the ability to automate these processes enhances efficiency and effectiveness in the examination phase of the investigation.

The use of LAMs and VoT techniques in the examination phase can significantly enhance the efficiency of an investigation. Since LAMs and VoT specifically focus on task manipulation, investigators can offload some examination work to an LAM, which will then generate the final results from a series of subtasks. This approach can allow investigators to focus on higher-level analysis and decision making, thus streamlining the overall investigative process.

4.5. Analysis phase

The analysis phase involves understanding the incident and obtaining a conclusive understanding based on the information collected during the examination phase. As also highlighted in Section 1.1.3, it has been demonstrated that LLMs are effective in case analysis (Henseler and van Beek, 2023). The use of MLLMs, which possess the capability to interpret images, broadens the scope for analysing a crime case more comprehensively. Using Gemini 1.5, Xu et al. (2024a) presented a tutorial on profiling a suspect's web history through an LLM. This case study demonstrates how an LLM can help identify potential motivations, personal interests, and psychological characteristics of the suspect. In conclusion, the authors suggest that such mechanisms could power AI-assisted tools, enabling law enforcement authorities to improve the identification of cybercriminals and malicious entities.

The Digital Forensic Cybercrime Language as a Service (DFClaaS) is an innovative system developed to address the complexities of text-based cybercrime (Al Mahdi and Baror, 2024). Using natural NLP techniques, including LLMs, sentiment analysis, and lexicon analysis, DFClaaS aims to improve capabilities in incident reporting, analysis, and investigation. The primary objectives of DFClaaS include implementing microservices to address specific challenges, proposing an advanced system to improve incident handling, and providing valuable tools for DF investigators. Designed to serve individual users, organisations, and forensic professionals, DFClaaS is a versatile and effective resource in the ongoing fight against cybercrime.

LLMs can be specifically fine-tuned for the analysis of various data types, including log files, email contents, chat transcripts, call records, file metadata, hex dumps, memory dumps, and registry hives. Incorporating contents such as event logs, timestamps, and network traffic captures further enables the effective recreation of incidents by correlating each data set with the assistance of LLMs. In addition, MLLMs that are audio and video specific can assist in analysing content within these formats. This specialised capability can significantly reduce the time investigators spend analysing audio and video data during investigations.

Table 3
DF functionalities by CFTT highlighting the usability of LLMs and example prompts.

CFTT Functionality	DF Phase(s)	Usable LLMs/Agent Frameworks	Example Prompt
Cloud Data Extraction	Acquisition	LLaMA (Fine-Tuned), Code Llama, StarCoder, AutoGen or CrewAI	Retrieve all the data inside given S3 bucket by using given credentials
Deleted File Recovery Specs	Acquisition, Examination	LLaMA (Fine-Tuned), Code Llama, StarCoder, AutoGen or CrewAI	Find all the deleted files from the X disk image and recover them to Y location
Disk Imaging	Acquisition	LLaMA (Fine-Tuned), Code Llama, StarCoder, AutoGen or CrewAI	Get a full disk image from this computer and save it in Z location
Forensic File Carving	Acquisition, Examination	LLaMA (Fine-Tuned), Code Llama, StarCoder, AutoGen or CrewAI	Find all the deleted PDF files from X disk image
Forensic Media Preparation	Acquisition	LLaMA (Fine-Tuned), Code Llama, StarCoder, AutoGen or CrewAI	Prepare the given X device for new investigation
String Search	Examination	LLaMA (Fine-Tuned), Code Llama, StarCoder, AutoGen or CrewAI	Search all the files containing the email of mail@test.com
Mobile Devices	Examination	LLaMA (Fine-Tuned), LLaVa (Fine-Tuned), Code Llama, StarCoder, AutoGen or CrewAI	Find all the photos taken with a computer within the last 3 months
MS Windows Registry	Examination	LLaMA (Fine-Tuned), Code Llama, StarCoder, AutoGen or CrewAI	Find information about the users from a given Windows disk image
SQLite Databases	Examination	LLaMA (Fine-Tuned), Code Llama, StarCoder, AutoGen or CrewAI	Find the access time of user Y to application X using given SQLite databases

The use of automated agents can effectively distribute the analysis workload. Moreover, leveraging Augmented Large Language Models (ALLMs) and RAG techniques can improve knowledge retrieval in real time, thus improving the accuracy of analysis and decision making processes (Lewis et al., 2020b). For example, integrating a source of intelligence with an RAG system can assist investigators in connecting the dots during a DF investigation.

Other than these applications, LLMs can increase productivity through enhanced information correlation during the analysis phase. Shafee et al. (2025) suggest that LLMs hold significant potential for data correlation and cybersecurity applications. The referenced study evaluated the performance of various LLM-based chatbots, including ChatGPT, GPT4all, Dolly, Stanford Alpaca, Alpaca-LoRA, Falcon, and Vicuna, specifically for text classification and Named Entity Recognition (NER) tasks using OSINT data. The findings indicated that, although the commercial chatbot GPT-4 and the open-source GPT4all performed well in text classification, all tested LLM-based chatbots showed limitations and were less effective for cybersecurity entity recognition compared to specialised models. The study concludes that there remains room for improvement.

4.6. Reporting phase

The quality and validity of the evidence, along with the thoroughness of the analysis, are encapsulated in the final report. The reporting phase holds significant weight, as the entire judgement may hinge on this crucial stage. Notably, DF is experiencing heightened scrutiny about the quality of the reports, emphasising the importance of precision and clarity in this phase (Karie et al., 2019). As pointed out by Champod et al. (2016), there is no standard framework for evaluating and reporting scientific findings to authorities and stakeholders. To provide assistance and alleviate scrutiny, incorporating LLMs for report creation is a viable solution. While LLMs are inherently non-deterministic, adhering to investigation standards such as ISO/IEC 27043:2015 can establish robust processes around data integrity and evidence handling, even though these standards do not directly address the randomness or variability in LLM outputs. The ISO/IEC 27043:2015 standard provides guidelines for a consistent DF investigation framework, focusing on maintaining procedural rigour rather than modifying model behaviour. Although it does not directly resolve issues of LLM determinism, it can serve as a protocol to ensure that procedures involving LLMs uphold investigative standards and maintain integrity throughout the process (Valjarević et al., 2016).

A feasibility study by Michelet and Breiteringer (2024) highlighted the potential of LLMs to assist in automating forensic report generation. These models can facilitate the creation of structured sections, such as methodologies, data analysis, and summaries, by generating coherent, case-specific insights from forensic data. Additionally, LLMs could automate the production of reports in alternative formats, such as HTML or \LaTeX , which are frequently used for dynamic, web-based, or highly technical documentation.

4.7. Other possibilities

Scanlon et al. (2023a) highlights that LLMs can play an important role in teaching scenarios. This involvement extends to activities such as storyboarding, creation of synthetic content, and synthetic character profiling. Fine-tuned models could further enhance training by generating more complex, realistic case examples that challenge trainees with nuanced scenarios, providing a robust foundation for practical skills development. These models may also help translate technical findings into accessible language, facilitating communication of insights to non-specialists, such as judges or other stakeholders.

4.8. Discussion on potential for LLMs in DF

To provide a comprehensive understanding of the potential use of LLMs, Table 3 clarifies the sample functionalities within the framework of the National Institute of Standards and Technology (NIST) Computer Forensics Tool Testing Program (CFTT), highlighting the usability of LLMs and example prompts. The CFTT project establishes overarching specifications to assess the capabilities of tools, a framework adopted by numerous prominent free and commercial tools.⁶

The potential for having a positive impact on the typical phases of the investigation increases as one progresses through the typical order of the phases. For example, there is little improvement that can be made by an LLM or automated scripting during the identification or acquisition phases, but significant potential for aiding investigators in the reporting phase (Michelet and Breiteringer, 2024; Wickramasekara and Scanlon, 2024) – these are first and foremost large *language* models. The low/medium/high potential outlined below evaluates each DF phase based on three key requirements: reliance on human expertise,

⁶ <https://www.nist.gov/itl/ssd/software-quality-group/computer-forensics-tool-testing-program-cftt>.

physical versus digital evidence handling, and scope for automation, as explained below.

- **Low Potential for the Identification and Collection Phases**
 - High dependency on human involvement, expertise, and/or specialised knowledge.
 - Involves extensive handling of physical evidence.
 - Limited or no feasibility for automation.
- **Medium Potential for the Preservation Phase**
 - Requires some level of human involvement or expertise, but is not critical to the process.
 - Primarily deals with digital evidence, with minimal physical evidence handling.
 - Feasible for automation to a significant extent.
- **High Potential for the Examination, Analysis and Reporting Phases**
 - Human involvement is needed for expert verification of the conducted analysis.
 - Exclusively focused on digital evidence.
 - Many common tasks are suitable for significant support from LLMs.

With these possibilities, the scope for research in DF is vast. Future research could be extended to the generation of digital forensic reports, as well as the summarisation of these reports for non-technical users. This would save time, but can also lead to more consistent documentation compared to manual documentation. Given the capacity of LLMs to manage large textual datasets, exploring pattern recognition holds significant value, particularly for investigations requiring the detection of anomalies or outliers in chats, log events, or emails.

In addition, LLMs' ability to interpret the tone of messages or chats enables their application in the sentiment analysis of text-based evidence. There is also potential in fine-tuning LLMs for domain-specific tasks, such as network forensics, where LLMs could analyse log files and application data related to specific activities. Automating LLM-based DF tools could further enable investigators to generate customised reports using natural language queries.

A critical future research direction lies in the ethical and legal considerations of LLM-generated content. As the application of LLMs is still emerging, future studies should focus on developing appropriate benchmarks, standardisation protocols, and addressing legal aspects to ensure responsible use of this technology.

5. Challenges and risks

This section discusses the challenges and risks of using LLMs in DF. Despite their promising potential, there are significant risk factors to consider. These risks can have severe consequences for DF if not adequately identified and considered in the DF process.

5.1. Challenges for LLMs in digital forensics

To optimise the results, the LLMs will likely need to be trained with specific forensic data (i.e., previous case data) to achieve the best results. Given the complexity and variation of the cases, it is questionable how good the training data are and whether there are sufficient data (Breitinger et al., 2024). Any bias in training data can lead to skewed interpretations and unjust outcomes. This problem of bias can be mitigated by using diverse and representative datasets during the training phase, e.g., datasets that come from diverse sources, different case types and geographic regions. Furthermore, techniques such as data filtering, distribution reconstruction, rebalancing, regularisation, and prompting can be implemented to actively identify and correct biased patterns in the base data sets of the model and its outputs (Dai et al., 2024; Zhou, 2024). These techniques involve adjusting model weights or incorporating fairness constraints during training to reduce the likelihood of

biased predictions. Regular audit of the model's performance against fairness benchmarks is also crucial to ensure that it remains unbiased over time (Mökander et al., 2023).

The experience level of investigators and the practical strategies employed in conducting investigations are challenging to replicate with LLMs. Initially, LLMs can excel in assisting with certain subtasks, such as parsing and data conversion, tasks in which output can be easily verified. However, when it comes to more interpretative or inferential tasks, LLMs' lack of inherent transparency introduces explainability challenges. Unlike deterministic software, whose logic can be easily traced, LLMs often act as black boxes, making it difficult to validate and understand the rationale behind their conclusions, particularly when these outputs extend beyond straightforward parsing into areas requiring judgement and reasoning. This underscores the importance of explainability in the application of LLMs to forensics, where understanding the basis of each result is crucial for accuracy and accountability (Michelet and Breitinger, 2024).

Publicly hosted and maintained LLMs are generally unsuitable for casework due to the sensitivity of the evidence and information involved, which require strict privacy and security controls that cannot be reliably ensured on public platforms. Furthermore, managing the substantial infrastructure needed for LLM training and deployment is both energy and resource intensive, presenting a financial hurdle, especially for smaller forensic laboratories with limited budgets. Although methods like retrieval-augmented generation (RAG) or prompt engineering can reduce some of the computational load by tailoring responses with existing models, they still require powerful GPU resources to effectively run these models, adding to the cost and accessibility barriers. Centralised systems could be a viable option, but they require well-defined guidelines for data sharing and stringent security standards to safeguard sensitive information.

Although LLMs can serve as valuable tools to support forensic investigations, it should be recognised that they currently function best as an aid, not a substitute for human expertise (Scanlon et al., 2023b). There is a risk that people may place too much trust in the results generated by LLMs (over-reliance), which could lead to complacency and overlook the need for detailed human expert analysis and validation.

To mitigate the potential misuse of LLMs, many LLMs are subjected to censorship (Yao et al., 2024; Brown et al., 2022). Although this censorship may serve as a preventive measure against unethical use, it can pose challenges in the field of DF. For example, if an investigator seeks evidence related to 'drugs' or evidence of other illegal material, censorship of the LLM may restrict access to accurate information related to the investigator's query. This limitation underscores the need for a nuanced approach to censorship in LLMs, balancing ethical considerations with the imperative of facilitating effective forensic investigations. In addition, the censorship of generic, publicly accessible LLMs further supports the argument for a discipline-specific DF LLM.

Finally, ethical and legal considerations must also be discussed. Determining accountability in cases where LLMs produce false information or are compromised by hacking. Clarifying responsibilities between developers, users, and regulators is crucial to establish a framework for accountability. If LLM generated DF results lead to incorrect information, the responsibility may lie with both the developers, for ensuring the model's accuracy, and the users, for appropriately interpreting and validating the results.

5.2. Risks of integration

The integration of LLMs within the DF process comes with inherent risks, in addition to the general LLM limitations outlined in Section 3.7. In particular, in the examination, analysis, and reporting phases, the use of LLMs introduces the risk of producing inaccurate information, primarily due to the phenomenon of inheritance hallucinations associated with these models (Michelet and Breitinger, 2024; Scanlon et al., 2023b). Additionally, the biases and obscurities present in an inheri-

tance model can significantly impact the performance of a DF-focused LLM – potentially leading to the unacceptable generation of biased or inaccurate information within the DF process.

Hallucinations in LLMs present a considerable risk, as they can produce information that appears credible but is incorrect. This can lead law enforcement authorities to form invalid assumptions and make flawed decisions based on unreliable results. Additionally, inherent biases in LLMs can influence investigative outcomes, which could affect the fairness and integrity of legal procedures. Data privacy concerns are also prominent, as sensitive information confidentiality may be compromised when using LLMs in DF processes. Together, these factors present substantial challenges to the reliable and ethical application of LLMs.

It is also crucial to acknowledge that DF LLMs, like any complex model, are susceptible to adversarial manipulation (Zou et al., 2024). This vulnerability poses a substantial risk in the context of sensitive domains such as DF, where the integrity of the information obtained is paramount. Adversarial attacks can compromise the reliability of LLM-generated outputs, potentially influencing the outcomes of various phases within the DF process.

Indeed, despite incorporating human verification, outputs and reports generated by LLMs within DF applications may encounter challenges regarding acceptance within the legal systems of different countries. This highlights a significant usability risk associated with LLM-based DF applications, but one that can be carefully mitigated by limiting the technology's deployment as a human-in-the-loop investigative aid as opposed to directly feeding into any investigative/judicial decision making processes.

Mitigating these changes and risks can be challenging, particularly in scenarios that involve adopting country-specific legal systems. However, there are potential strategies to address technical challenges such as hallucinations, censorship, and substantial infrastructure costs.

One solution to mitigate hallucination was suggested by Ji et al. (2023), who proposed an interactive self-reflection method for generated knowledge and answers, an approach that has shown promise. Another method of reducing hallucinations is the use of RAG, which provides a larger knowledge base for LLMs to minimise unknown information (Lewis et al., 2020b). Other methods such as knowledge graphs, bias detection mechanisms, active learning methods for LLMs, supervised fine-tuning strategies, hallucination mitigation frameworks, and new decoding strategies can also help mitigate hallucinations to some extent (Perković et al., 2024).

Censorship issues can be addressed by fine-tuning the model with uncensored information, a technique already applied to the LLaMA and Mistral models, leading to the development of the Dolphin models. An example is the Dolphin-2.0-mistral-7b, which is an uncensored version of the Mistral 7B model (Xu et al., 2024b).

The high infrastructure costs associated with these models can be mitigated by employing Data Forensics as a Service (DFaaS) platforms such as Hansken. With DFaaS, investigators only need to input queries related to their investigations using personal computers, while the platform manages the model maintenance and computational demands (van Beek et al., 2015).

Despite these promising integration risks, the use of LLMs may face limitations in adaptability. The performance of an LLM is inherently tied to the dataset on which it was trained, which means that its ability to respond to new or emerging information is constrained. For example, if an LLM is tasked with identifying possible malware in a system, it may struggle to detect newer malware variants that were not part of its training data (Yu et al., 2024a). To mitigate such issues, LLMs need to be fine-tuned frequently, which poses its own challenges due to the significant computational power required for such operations.

6. Conclusion

The convergence of LLMs with an array of technologies represents exciting synergy. Although the utilisation of LLMs in the realm of DF is

still in its nascent stages, there is evidence of their substantial potential to significantly increase the efficiency of investigations. The exploration of investments for LLMs across the entire DF process is considered, with the aim of improving the productivity and efficiency of investigations. Additionally, the integration of LLMs into current DF tools is posited to reduce user training times, as these models comprehend natural language input and provide output accordingly. In the dynamic landscape of LLM applications for DF, promising avenues for further exploration and advancement unfold.

Although the surge in LLM research is promising, it is crucial to balance enthusiasm with awareness of existing challenges. The propensity of LLMs to produce hallucinations highlights the need for human oversight in critical decision making processes, underscoring the irreplaceable value of human judgement, intuition, and expertise. A notable limitation is the language dependency issue, as most LLMs are predominantly trained on English data, reducing their effectiveness with non-English content. Furthermore, the deployment of LLMs in DF involves significant costs related to the infrastructure to process evidence. Questions also arise about the validation of task correctness and quality when automated by LLMs, as well as the legal and professional acceptance of results obtained with limited human intervention.

The trustworthiness of LLMs remains a debatable issue that requires careful attention. It is crucial to establish clear boundaries and measures to define LLM trustworthiness. Addressing this will be a key aspect in the field of DF, ensuring that LLMs can be trusted for accurate and secure analysis, with the explainability of their operations being paramount.

Integrating LLMs with automated agents offers a promising path to automating DF processes, potentially allowing multiple cases to be handled concurrently for more timely and precise outcomes. This integration could significantly streamline investigations. Future research should explore the role of LLMs and AI in the decision making of DF. It is essential to focus on validating LLM generated outputs to ensure their scope, accuracy, reliability, and trustworthiness in investigations. More studies comparing DF outcomes with and without LLM integration are critical, as they could highlight the benefits of LLMs and the controlled applicability of LLMs in DF and similar fields.

A future use case involves developing forensic-specific LLMs fine-tuned for automated examinations. These models could be optimised for script generation to support investigations where no existing tools are available, allowing forensic analysts to create customised solutions on demand. Integrating AI agents with these models could streamline evidence handling by allowing investigators to perform complex queries more intuitively, such as retrieving all messages from a specific date without the need to craft regular expressions.

In essence, while LLMs offer exciting prospects for the future of digital forensics, a balanced approach that integrates their strengths with human oversight is essential for harnessing their full potential. Inevitably, LLM-facilitated DF processes themselves will become the focus of future investigation.

CRedit authorship contribution statement

Akila Wickramasekara: Writing – review & editing, Writing – original draft, Validation, Methodology, Investigation, Formal analysis, Data curation, Conceptualization. **Frank Breitinger:** Writing – review & editing, Writing – original draft, Validation, Methodology, Investigation. **Mark Scanlon:** Writing – review & editing, Writing – original draft, Validation, Supervision, Resources, Project administration, Methodology, Investigation, Data curation, Conceptualization.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

No data was used for the research described in the article.

References

- Acheampong, F.A., Nunoo-Mensah, H., Chen, W., 2021. Transformer models for text-based emotion detection: a review of BERT-based approaches. *Artif. Intell. Rev.* 54 (8), 5789–5829. <https://doi.org/10.1007/s10462-021-09958-2>.
- Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., et al., 2023. GPT-4 technical report. *CoRR arXiv:2303.08774 [abs]*. <https://doi.org/10.48550/ARXIV.2303.08774>.
- Al Mahdi, M.M., Baror, S., 2024. Proof of Concept of a Digital Forensic Readiness Cybercrime Language as a Service. In: *International Conference on Cyber Warfare and Security*, vol. 19, pp. 191–199. <https://doi.org/10.34190/iccws.19.1.2059>.
- Al-Dhaq, A., Razak, S.A., Ikuesan, R.A., Kbande, V.R., Siddique, K., 2020. A review of mobile forensic investigation process models. *IEEE Access* 8, 173359–173375. <https://doi.org/10.1109/ACCESS.2020.3014615>.
- Al-Rfou, R., Choe, D., Constant, N., Guo, M., Jones, L., 2019. Character-level language modeling with deeper self-attention. In: *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, pp. 3159–3166. <https://doi.org/10.1609/aaai.v33i01.33013159>.
- Alawida, M., Mejri, S., Mehmood, A., Chikhaoui, B., Isaac Abiodun, O., 2023. A Comprehensive Study of ChatGPT: Advancements, Limitations, and Ethical Considerations in Natural Language Processing and Cybersecurity. *Informatics* 14 (8). <https://doi.org/10.3390/info14080462>.
- Alayrac, J.-B., Donahue, J., Luc, P., Miech, A., Barr, I., et al., 2022. Flamingo: a Visual Language Model for Few-Shot Learning. In: Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., Oh, A. (Eds.), *Advances in Neural Information Processing Systems*, vol. 35. Curran Associates, Inc., pp. 23716–23736. https://proceedings.neurips.cc/paper_files/paper/2022/file/960a172bc7fb0177ccccbb411a7d800-Paper-Conference.pdf.
- Ali, L., 2019. Cyber crimes-A constant threat for the business sectors and its growth (A study of the online banking sectors in GCC). *J. Dev. Areas* 53.
- Allal, L.B., Li, R., Kocetkov, D., Mou, C., Akiki, C., et al., 2023. SantaCoder: don't reach for the stars! Deep Learning for Code (DL4C) Workshop. <https://par.nsf.gov/biblio/10416454>.
- Alon, U., Sadaka, R., Levy, O., Yahav, E., 2020. Structural language models of code. In: *Proceedings of the 37th International Conference on Machine Learning*, pp. 245–256. <https://doi.org/10.5555/3524938.3524962>.
- Arawjo, I., Swoopes, C., Vaithilingam, P., Wattenberg, M., Glassman, E.L., 2024. ChainForge: a visual toolkit for prompt engineering and LLM hypothesis testing. In: *Proceedings of the CHI Conference on Human Factors in Computing Systems*. CHI'24. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3613904.3642016>.
- Bajwa, A., Farooq, S., Malik, O., Khalique, S., Suguri, H., Farooq Ahmad, H., Ali, A., 2007. Persistent architecture for context aware lightweight multi-agent system. In: Bordini, R.H., Dastani, M., Dix, J., Seghrouchni, A.E.F. (Eds.), *Programming Multi-Agent Systems*. Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 57–69. https://doi.org/10.1007/978-3-540-71956-4_4.
- Bakker, M., Chadwick, M., Sheahan, H., Tessler, M., Campbell-Gillingham, L., Balaguer, J., McAleese, N., Glaese, A., Aslanides, J., Botvinick, M., Summerfield, C., 2022. Fine-tuning language models to find agreement among humans with diverse preferences. In: Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., Oh, A. (Eds.), *Advances in Neural Information Processing Systems*, vol. 35. Curran Associates, Inc., pp. 38176–38189. https://proceedings.neurips.cc/paper_files/paper/2022/file/f978c8f3b5f399cae464e85f72e28503-Paper-Conference.pdf.
- Barbaroxxa, R., Gomes, L., Vale, Z., 2025. Benchmarking Large Language Models for Multi-agent Systems: A Comparative Analysis of AutoGen, CrewAI, and TaskWeaver. In: Mathieu, P., De la Prieta, F. (Eds.), *Advances in Practical Applications of Agents, Multi-Agent Systems, and Digital Twins: The PAAMS Collection*. Springer Nature, Switzerland, Cham, pp. 39–48. https://doi.org/10.1007/978-3-031-70415-4_4.
- Baryamureeba, V., Tushabe, F., 2004. The enhanced digital investigation process model. In: *Proceedings of the Digital Forensic Research Conference (DFRWS)*.
- Bender, E.M., Geburu, T., McMillan-Major, A., Shmitchell, S., 2021. On the dangers of stochastic parrots: can language models be too big? In: *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*. FAccT'21. Association for Computing Machinery, New York, NY, USA, pp. 610–623. <https://doi.org/10.1145/3442188.3445922>.
- Bonner, E., Lege, R., Frazier, E., 2023. Large language model-based artificial intelligence in the language classroom: practical ideas for teaching. *Teaching English with Technology* 23. <https://doi.org/10.56297/bkam1691/wieo1749>.
- Bran, A.M., Cox, S., Schilter, O., Baldassari, C., White, A.D., Schwaller, P., 2024. Augmenting large language models with chemistry tools. *Nat. Mach. Intell.* 6, 525–535. <https://doi.org/10.1038/s42256-024-00832-8>.
- Breitering, F., Hilgert, J.-N., Hargreaves, C., Sheppard, J., Overdorf, R., Scanlon, M., 2024. DFRWS EU 10-year review and future directions in digital forensic research. *Forensic Sci. Int. Digit. Investig.* 48, 301685. <https://doi.org/10.1016/j.fsi.2023.301685>.
- Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., et al., 2020. Language models are few-shot learners. In: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., Lin, H. (Eds.), *Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc., pp. 1877–1901. https://papers.nips.cc/paper_files/paper/2020/file/1457c0d6bfb4967418bfb8ac142f64a-Paper.pdf.
- Brown, H., Lee, K., Miresghallah, F., Shokri, R., Tramèr, F., 2022. What does it mean for a language model to preserve privacy? In: *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. FAccT'22. Association for Computing Machinery, New York, NY, USA, pp. 2280–2292. <https://doi.org/10.1145/3531146.3534642>.
- Caines, A., Benedetto, L., Taslimipour, S., Davis, C., Gao, Y., Andersen, Ø.E., Yuan, Z., Elliott, M., Moore, R., Bryant, C., Rei, M., Yannakoudakis, H., Mullooly, A., Nicholls, D., Buttery, P., 2023. On the application of large language models for language teaching and assessment technology. In: *Proceedings of the Workshop on Empowering Education with LLMs - the Next-Gen Interface and Content Generation 2023*, co-located with 24th International Conference on Artificial Intelligence in Education (AIED 2023), pp. 173–197. <https://ceur-ws.org/Vol-3487/paper12.pdf>.
- Carvalho, I., Ivanov, S., 2024. ChatGPT for tourism: applications, benefits and risks. *Tour. Rev.* 79 (2). <https://doi.org/10.1108/TR-02-2023-0088>.
- Casey, E., 2011. *Digital Evidence and Computer Crime: Forensic Science, Computers, and the Internet*. Academic Press.
- Champod, C., Biedermann, A., Vuille, J., Willis, S., De Kinder, J., 2016. ENFSI guideline for evaluative reporting in forensic science: a primer for legal practitioners. *Crim. Law Justice Wkly.* 180, 189–193.
- Chang, Y., Wang, X., Wang, J., Wu, Y., Yang, L., Zhu, K., Chen, H., Yi, X., Wang, C., Wang, Y., Ye, W., Zhang, Y., Chang, Y., Yu, P.S., Yang, Q., Xie, X., 2024. A survey on evaluation of large language models. *ACM Trans. Intell. Syst. Technol.* <https://doi.org/10.1145/3641289>.
- Chen, B., Zhang, F., Nguyen, A., Zan, D., Lin, Z., Lou, J.-G., Chen, W., 2023. CodeT: code generation with generated tests. In: *The Eleventh International Conference on Learning Representations*.
- Chen, M., Tworek, J., Jun, H., Yuan, Q., de Oliveira Pinto, H.P., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., Ray, A., Puri, R., Krueger, G., Petrov, M., Khlaaf, H., Sastry, G., Mishkin, P., Chan, B., Gray, S., Ryder, N., Pavlov, M., Power, A., Kaiser, L., Bavarian, M., Winter, C., Tillet, P., Such, F.P., Cummings, D., Plappert, M., Chantzis, F., Barnes, E., Herbert-Voss, A., Guss, W.H., Nichol, A., Paino, A., Tezak, N., Tang, J., Babuschkin, I., Balaji, S., Jain, S., Saunders, W., Hesse, C., Carr, A.N., Leike, J., Achiam, J., Misra, V., Morikawa, E., Radford, A., Knight, M., Brundage, M., Murati, M., Mayer, K., Welinder, P., McGrew, B., Amodei, D., McCandlish, S., Sutskever, I., Zaremba, W., 2021. Evaluating large language models trained on code. *CoRR. arXiv:2107.03374 [abs]*. <https://doi.org/10.48550/arXiv.2107.03374>.
- Chiang, D.C., Lee, H., 2023. Can large language models be an alternative to human evaluations? <https://doi.org/10.18653/v1/2023.acl-long-870>.
- Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., et al., 2023. PaLM: scaling language modeling with pathways. *J. Mach. Learn. Res.* 24, 1–113.
- Dai, S., Xu, C., Xu, S., Pang, L., Dong, Z., Xu, J., 2024. Bias and unfairness in information retrieval systems: new challenges in the LLM era. In: *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. KDD'24. Association for Computing Machinery, New York, NY, USA, pp. 6437–6447. <https://doi.org/10.1145/3637528.3671458>.
- Dale, R., 2021. GPT-3: what's it good for? *Nat. Lang. Eng.* 27, 113–118.
- DeepSeek-AI, Aixin, L., Feng, B., Bing, X., Bingxuan, W., et al., 2024. DeepSeek-V3 Technical Report. *CoRR arXiv.2412.19437 [abs]*. <https://doi.org/10.48550/arXiv.2412.19437>.
- Dettmers, T., Pagnoni, A., Holtzman, A., Zettlemoyer, L., 2023. QLoRA: efficient finetuning of quantized LLMs. In: Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., Levine, S. (Eds.), *Advances in Neural Information Processing Systems*, vol. 36. https://proceedings.neurips.cc/paper_files/paper/2023/file/1feb87871436031bdc0f2beaa62a049b-Paper-Conference.pdf.
- Devlin, J., Chang, M., Lee, K., Toutanova, K., 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In: Burstein, J., Doran, C., Solorio, T. (Eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. NAACL-HLT 2019, Minneapolis, MN, USA, June 2–7, 2019. In: Volume 1 (Long And Short Papers), Association for Computational Linguistics, pp. 4171–4186. <https://doi.org/10.18653/v1/n19-1423>.
- Di, P., Li, J., Yu, H., Jiang, W., Cai, W., Cao, Y., Chen, C., Chen, D., Chen, H., Chen, L., Fan, G., Gong, Z., Hu, W., Guo, T., Lei, Z., Li, T., Li, Z., Liang, M., Liao, C., Liu, B., Liu, J., Liu, Z., Lu, S., Shen, M., Wang, G., Wang, H., Wang, Z., Xu, Z., Yang, J., Ye, Q., Zhang, G., Zhang, Y., Zhao, Z., Zheng, X., Zhou, H., Zhu, L., Zhu, X., 2024. CodeFuse-13B: a pretrained multi-lingual code large language model. In: *Proceedings of the 46th International Conference on Software Engineering: Software Engineering in Practice*. ICSE-SEIP'24. Association for Computing Machinery, New York, NY, USA, pp. 418–429. <https://doi.org/10.1145/3639477.3639719>.
- Dong, S., Wang, P., Abbas, K., 2021. A survey on deep learning and its applications. *Comput. Sci. Rev.* 40, 100379. <https://doi.org/10.1016/j.cosrev.2021.100379>. <https://www.sciencedirect.com/science/article/pii/S1574013721000198>.
- Du, X., Scanlon, M., 2019. Methodology for the automated metadata-based classification of incriminating digital forensic artifacts. In: *Proceedings of the 14th International Conference on Availability, Reliability and Security*. ARES'19. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3339252.3340517>.
- Du, X., Le-Khac, N.-A., Scanlon, M., 2017. Evaluation of digital forensic process models with respect to Digital Forensics as a Service. In: *Proceedings of the 16th European*

- Conference on Cyber Warfare and Security. (ECCWS 2017), ACPI, Dublin, Ireland, pp. 573–581.
- Dubey, H., Bhatt, S., Negi, L., 2023. Digital forensics techniques and trends: a review, the international Arab. J. Inf. Technol. 20, 644–654.
- Eggmann, F., Weiger, R., Zitzmann, N.U., Blatz, M.B., 2023. Implications of large language models such as ChatGPT for dental medicine. J. Aesthet. Restor. Dent. <https://doi.org/10.1111/jerd.13046>.
- Fan, W., Ding, Y., Ning, L., Wang, S., Li, H., Yin, D., Chua, T.-S., Li, Q., 2024. A survey on RAG meeting LLMs: towards retrieval-augmented large language models. In: Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining. KDD'24. Association for Computing Machinery, New York, NY, USA, pp. 6491–6501. <https://doi.org/10.1145/3637528.3671470>.
- Fried, D., Aghajanyan, A., Lin, J., Wang, S., Wallace, E., Shi, F., Zhong, R., Yih, S., Zettlemoyer, L., Lewis, M., 2023. InCoder: a generative model for code infilling and synthesis. In: The Eleventh International Conference on Learning Representations. ICLR 2023, Kigali, Rwanda, May 1–5, 2023.
- Fröhling, L., Zubiaga, A., 2021. Feature-based detection of automated language models: tackling GPT-2, GPT-3, and Grover. PeerJ Comput. Sci. 7, e443. <https://doi.org/10.7717/peerj-cs.443>.
- Goel, D., Husain, F., Singh, A., Ghosh, S., Parayil, A., Bansal, C., Zhang, X., Rajmohan, S., 2024. X-lifecycle learning for cloud incident management using LLMs. In: Companion Proceedings of the 32nd ACM International Conference on the Foundations of Software Engineering. FSE 2024. Association for Computing Machinery, New York, NY, USA, pp. 417–428. <https://doi.org/10.1145/3663529.3663861>.
- Guo, Q., De Mello, S., Yin, H., Byeon, W., Cheung, K.C., Yu, Y., Luo, P., Liu, S., 2024a. RegionGPT: towards region understanding vision language model. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 13796–13806. <https://doi.org/10.1109/CVPR52733.2024.01309>.
- Guo, T., Chen, X., Wang, Y., Chang, R., Pei, S., Chawla, N.V., Wiest, O., Zhang, X., 2024b. Large language model based multi-agents: a survey of progress and challenges. In: Larson, K. (Ed.), Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence. IJCAI-24, pp. 8048–8057. <https://doi.org/10.24963/ijcai.2024/890>.
- Henseler, H., van Beek, H., 2023. ChatGPT as a copilot for investigating digital evidence. In: Conrad, J.G., Jr., D.W.L., Baron, J.R., Henseler, H., Bhattacharya, P., Nielsen, A., Vinjumur, J.K., Pickens, J., Jones, A. (Eds.), Proceedings of the Third International Workshop on Artificial Intelligence and Intelligent Assistance for Legal Professionals in the Digital Workplace (LegalAIIA 2023) Co-Located with the 19th International Conference on Artificial Intelligence and Law. (ICAIL 2023), Braga, Portugal. In: CEUR Workshop Proceedings, CEUR-WS.org, vol. 3423, pp. 58–69. <https://ceur-ws.org/Vol-3423/paper6.pdf>, 2023.
- Hou, X., Zhao, Y., Liu, Y., Yang, Z., Wang, K., Li, L., Luo, X., Lo, D., Grundy, J., Wang, H., 2024. Large language models for software engineering: a systematic literature review. ACM Trans. Softw. Eng. Methodol. <https://doi.org/10.1145/3695988>. just Accepted.
- Jarrett, F., Choo, K.-K., 2021. The impact of automation and artificial intelligence on digital forensics. WIREs Forensic Sci. 3 (6), e1418. <https://doi.org/10.1002/wfs2.1418>.
- Ji, Z., Yu, T., Xu, Y., Lee, N., Ishii, E., Fung, P., 2023. Towards mitigating LLM hallucination via self reflection. In: Bouamor, H., Pino, J., Bali, K. (Eds.), Findings of the Association for Computational Linguistics: EMNLP 2023, Association for Computational Linguistics. Singapore, pp. 1827–1843. <https://doi.org/10.18653/v1/2023.findings-emnlp.123>.
- Jiang, Z., Xu, F.F., Gao, L., Sun, Z., Liu, Q., Dwivedi-Yu, J., Yang, Y., Callan, J., Neubig, G., 2023. Active retrieval augmented generation. In: Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing. <https://doi.org/10.18653/v1/2023.emnlp-main.495>.
- Kalaimannan, E., Gupta, J.N., Yoo, S.-M., 2013. Maximizing investigation effectiveness in digital forensic cases. In: 2013 International Conference on Social Computing, pp. 618–623. <https://doi.org/10.1109/SocialCom.2013.93>.
- Kalyan, K.S., 2024. A survey of GPT-3 family large language models including ChatGPT and GPT-4. Nat. Lang. Process. J. 6, 100048. <https://doi.org/10.1016/j.nlp.2023.100048>. <https://www.sciencedirect.com/science/article/pii/S2949719123000456>.
- Kan, K.B., Mun, H., Cao, G., Lee, Y., 2024. Mobile-LLaMA: instruction fine-tuning open-source LLM for network analysis in 5G networks. IEEE Netw. 38, 76–83. <https://doi.org/10.1109/MNET.2024.3421306>.
- Karabacak, M., Margetis, K., 2023. Embracing large language models for medical applications: opportunities and challenges. Cureus 15. <https://doi.org/10.7759/cureus.39305>.
- Karie, N.M., Kebande, V.R., Venter, H., Choo, K.-K.R., 2019. On the importance of standardizing the process of generating digital forensic reports. Forensic Sci. Int. Rep. 1, 100008. <https://doi.org/10.1016/j.fsir.2019.100008>.
- Karlsen, E., Luo, X., Zincir-Heywood, N., Heywood, M., 2024. Benchmarking large language models for log analysis, security, and interpretation. J. Netw. Syst. Manag. 32 (59). <https://doi.org/10.1007/s10922-024-09831-x>.
- Ke, G., He, D., Liu, T., 2021. Rethinking positional encoding in language pre-training. In: 9th International Conference on Learning Representations. ICLR 2021, Virtual Event, Austria, May 3–7, 2021.
- Kiros, R., Salakhutdinov, R., Zemel, R., 2014. Multimodal neural language models. In: Proceedings of the 31st International Conference on Machine Learning. In: Proceedings of Machine Learning Research, vol. 32. PMLR, Beijing, China, pp. 595–603. <https://proceedings.mlr.press/v32/kiros14.html>.
- Koh, J.Y., Fried, D., Salakhutdinov, R., 2023. Generating images with multimodal language models. In: Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., Levine, S. (Eds.), Advances in Neural Information Processing Systems, vol. 36. https://proceedings.neurips.cc/paper_files/paper/2023/file/43a69d143273bd8215578bde887bb552-Paper-Conference.pdf.
- Koper, C.S., Lum, C., Willis, J.J., 2014. Optimizing the use of technology in policing: results and implications from a multi-site study of the social, organizational, and behavioural aspects of implementing police technologies. Policing J. Policy Pract. 8, 212–221. <https://doi.org/10.1093/policing/pau015>.
- Lai, S., Liu, K., He, S., Zhao, J., 2016. How to generate a good word embedding. IEEE Intell. Syst. 31, 5–14. <https://doi.org/10.1109/MIS.2016.45>.
- LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. Nature 521, 436–444. <https://doi.org/10.1038/nature14539>.
- Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.-t., Rocktäschel, T., Riedel, S., Kiela, D., 2020b. Retrieval-augmented generation for knowledge-intensive NLP tasks. In: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., Lin, H. (Eds.), Advances in Neural Information Processing Systems, vol. 33. Curran Associates, Inc., pp. 9459–9474. <https://proceedings.neurips.cc/paper/2020/file/6b493230205f780e1bc26945df7481e5-Paper.pdf>.
- Li, R., Allal, L.B., Zi, Y., Muennighoff, N., Kocetkov, D., et al., 2023a. StarCoder: may the source be with you! Trans. Mach. Learn. Res.
- Li, J., Li, D., Savarese, S., Hoi, S., 2023b. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In: Proceedings of the 40th International Conference on Machine Learning, pp. 19730–19742.
- Li, G., Hammoud, H.A.A.K., Itani, H., Khizbullin, D., Ghanem, B., 2023c. CAMEL: communicative agents for “mind” exploration of large language model society. In: Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., Levine, S. (Eds.), Advances in Neural Information Processing Systems, vol. 36. https://proceedings.neurips.cc/paper_files/paper/2023/file/a3621ee907de47c1b952ade25c67698-Paper-Conference.pdf.
- Li, H., Zhang, J., Liu, H., Fan, J., Zhang, X., Zhu, J., Wei, R., Pan, H., Li, C., Chen, H., 2024a. CodeS: towards building open-source language models for text-to-SQL. Proc. ACM Manag. Data 2. <https://doi.org/10.1145/3654930>.
- Li, N., Gao, C., Li, M., Li, Y., Liao, Q., 2024b. EconAgent: large language model-empowered agents for simulating macroeconomic activities. In: Duh, K., Gomez, H., Bethard, S. (Eds.), Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Association for Computational Linguistics, Bangkok, Thailand, pp. 15523–15536. <https://doi.org/10.18653/v1/2024.acl-long.829>.
- Li, Y., Choi, D., Chung, J., Kushman, N., Schrittwieser, J., Leblond, R., Eccles, T., Keeling, J., Gimeno, F., Lago, A.D., Hubert, T., Choy, P., de Masson d’Autume, C., Babuschkin, I., Chen, X., Huang, P.-S., Welbl, J., Goyal, S., Cherepanov, A., Molloy, J., Mankowitz, D.J., Robson, E.S., Kohli, P., de Freitas, N., Kavukcuoglu, K., Vinyals, O., 2022. Competition-level code generation with AlphaCode. Science 378, 1092–1097. <https://www.science.org/doi/abs/10.1126/science.abq1158>.
- Lira, O.G., Marroquin, A., To, M.A., 2024. Harnessing the advanced capabilities of LLM for adaptive intrusion detection systems. In: Barolli, L. (Ed.), Advanced Information Networking and Applications. Springer Nature, Switzerland, Cham, pp. 453–464. https://doi.org/10.1007/978-3-031-57942-4_44.
- Liu, J., Xia, C.S., Wang, Y., Zhang, L., 2023a. Is your code generated by ChatGPT really correct? Rigorous evaluation of large language models for code generation. In: Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., Levine, S. (Eds.), Advances in Neural Information Processing Systems, vol. 36. https://proceedings.neurips.cc/paper_files/paper/2023/file/43e9d647ccd3e4b7b5baab53f0368686-Paper-Conference.pdf, 2023.
- Liu, H., Li, C., Wu, Q., Lee, Y.-J., 2023b. Visual instruction tuning. In: Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., Levine, S. (Eds.), Advances in Neural Information Processing Systems, vol. 36. https://proceedings.neurips.cc/paper_files/paper/2023/file/6dcf277ea32ce3288914fa369fe6de0-Paper-Conference.pdf.
- Liu, Z., Zhou, Y., Zhu, Y., Lian, J., Li, C., Dou, Z., Lian, D., Nie, J.-Y., 2024a. Information retrieval meets large language models. In: Companion Proceedings of the ACM Web Conference 2024, WWW'24. Association for Computing Machinery, New York, NY, USA, pp. 1586–1589. <https://doi.org/10.1145/3589335.3641299>.
- Liu, H., Li, C., Li, Y., Lee, Y., 2024b. Improved baselines with visual instruction tuning. In: 2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 26286–26296. <https://doi.org/10.1109/CVPR52733.2024.02484>.
- Lu, G., Ju, X., Chen, X., Pei, W., Cai, Z., 2024. GRACE: empowering LLM-based software vulnerability detection with graph structure and in-context learning. J. Syst. Softw. 212, 112031. <https://doi.org/10.1016/j.jss.2024.112031>. <https://www.sciencedirect.com/science/article/pii/S0164121224000748>.
- Lund, B.D., Wang, T., 2023. Chatting about ChatGPT: how AI and GPT may impact academia and libraries? Libr. Hi Tech News 40 (3), 26–29. <https://doi.org/10.1108/LHTN-01-2023-0009>.
- Luo, Z., Xu, C., Zhao, P., Sun, Q., Geng, X., Hu, W., Tao, C., Ma, J., Lin, Q., Jiang, D., 2024. WizardCoder: empowering code large language models with Evol-Instruct. In: The Twelfth International Conference on Learning Representations.
- Lutui, R., 2016. A multidisciplinary digital forensic investigation process model. Bus. Horiz. 59, 593–604.
- Ma, K., Zang, X., Feng, Z., Fang, H., Ban, C., Wei, Y., He, Z., Li, Y., Sun, H., 2023. LLaV-iLo: boosting video moment retrieval via adapter-based multimodal modeling. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops, pp. 2798–2803. <https://doi.org/10.1109/ICCVW60793.2023.00297>.

- Marvin, G., Hellen, N., Jjingo, D., Nakatumba-Nabende, J., 2024. Prompt engineering in large language models. In: Jacob, I.J., Piramuthu, S., Falkowski-Gilski, P. (Eds.), *Data Intelligence and Cognitive Informatics*. Springer Nature, Singapore, Singapore, pp. 387–402. https://doi.org/10.1007/978-981-99-7962-2_30.
- Mesko, B., 2023. Prompt engineering as an important emerging skill for medical professionals: tutorial. *J. Med. Internet Res.* 25, e50638. <https://www.jmir.org/2023/1/e50638>.
- Michelet, G., Breiting, 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. Investig.* 48, 301683. <https://www.sciencedirect.com/science/article/pii/S2666281723002020>. <https://doi.org/10.1016/j.fsidi.2023.301683>. DFRWS EU 2024 - Selected Papers from the 11th Annual Digital Forensics Research Conference Europe.
- Michelet, G., Breiting, F., Horsman, G., 2023. Automation for digital forensics: towards a definition for the community. *Forensic Sci. Int.* 349, 111769. <https://www.sciencedirect.com/science/article/pii/S03790738230021900>. <https://doi.org/10.1016/j.fsidi.2023.111769>.
- Mikolov, T., Chen, K., Corrado, G., Dean, J., 2013. Efficient estimation of word representations in vector space. In: Bengio, Y., LeCun, Y. (Eds.), *1st International Conference on Learning Representations, Workshop Track Proceedings*. ICLR 2013, Scottsdale, Arizona, USA, May 2–4, 2013.
- Mir, S.S., Shoaib, U., Sarfraz, M.S., 2016. Analysis of digital forensic investigation models. *Int. J. Comput. Sci. Inf. Secur.* 14.
- Mökander, J., Schuett, J., Kirk, H.R., Floridi, L., 2023. Auditing large language models: a three-layered approach. *AI Ethics* 4, 1085–1115. <https://doi.org/10.1007/s43681-023-00289-2>.
- Moore, S., Tong, R., Singh, A., Liu, Z., Hu, X., et al., 2023. Empowering education with LLMs: the next-gen interface and content generation. In: *International Conference on Artificial Intelligence in Education*. Springer, pp. 32–37.
- Mukherjee, S., Haque, S., 2018. Review paper on digital forensics practices: a road map for building digital forensics capability. *Iconic Res. Eng. J.* 1, 96–99.
- Neuwirth, R.J., 2023. Prohibited artificial intelligence practices in the proposed EU artificial intelligence act (AIA). *Comput. Law Secur. Rev.* 48, 105798. <https://doi.org/10.1016/j.clsr.2023.105798>. <https://www.sciencedirect.com/science/article/pii/S0267364923000092>.
- Nijkamp, E., Pang, B., Hayashi, H., Tu, L., Wang, H., et al., 2022. CodeGen: an open large language model for code with multi-turn program synthesis. In: *The Eleventh International Conference on Learning Representations*.
- Nozza, D., Bianchi, F., Lauscher, A., Hovy, D., et al., 2022. Measuring harmful sentence completion in language models for LGBTQIA+ individuals. In: *Proceedings of the Second Workshop on Language Technology for Equality, Diversity and Inclusion*, Association for Computational Linguistics, pp. 26–34. <https://doi.org/10.18653/v1/2022.ltedi-1.4>.
- O'Leary, D.E., 2023. An analysis of three chatbots: BlenderBot, ChatGPT and LaMDA, intelligent systems in accounting. *Finance Manag.* 30, 41–54. <https://doi.org/10.1002/isaf.1531>. <https://onlinelibrary.wiley.com/doi/abs/10.1002/isaf.1531>.
- Ou, J., Lu, J., Liu, C., Tang, Y., Zhang, F., Zhang, D., Gai, K., 2024. DialogBench: evaluating LLMs as human-like dialogue systems. In: Duh, K., Gomez, H., Bethard, S. (Eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*. Association for Computational Linguistics, Mexico City, Mexico, pp. 6137–6170. <https://doi.org/10.18653/v1/2024.naacl-long.341>.
- Park, S.-M., Kim, Y.-G., 2023. Visual language integration: a survey and open challenges. *Comput. Sci. Rev.* 48, 100548. <https://doi.org/10.1016/j.cosrev.2023.100548>. <https://www.sciencedirect.com/science/article/pii/S1574013723000151>.
- Penedo, G., Maltart, Q., Hesslow, R., Cojocaru, R., Alobaidli, H., Cappelli, A., Pannier, B., Almazrouei, E., Launay, J., 2023. The RefinedWeb dataset for falcon LLM: outperforming curated corpora with web data only. In: Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., Levine, S. (Eds.), *Advances in Neural Information Processing Systems*, vol. 36. https://proceedings.neurips.cc/paper_files/paper/2023/file/fa3ed726cc5073b9c31e3e49a807789c-Paper-Datasets_and_Benchmarks.pdf.
- Perković, G., Drobnjak, A., Botički, I., 2024. Hallucinations in LLMs: understanding and addressing challenges. In: *2024 47th MIPRO ICT and Electronics Convention (MIPRO)*, pp. 2084–2088. <https://doi.org/10.1109/MIPRO60963.2024.10569238>.
- Polak, M.P., Morgan, D., 2024. Extracting accurate materials data from research papers with conversational language models and prompt engineering. *Nat. Commun.* 15, 1569. <https://doi.org/10.1038/s41467-024-45914-8>.
- Prayudi, Y., Ashari, A., Priyambodo, T.K., 2020. The framework to support the digital evidence handling: a case study of procedures for the management of evidence in Indonesia. *J. Cases Inf. Technol.* 22, 51–71. <https://doi.org/10.4018/JCIT.2020070104>.
- Qi, S., Cao, Z., Rao, J., Wang, L., Xiao, J., Wang, X., 2023. What is the limitation of multimodal LLMs? A deeper look into multimodal LLMs through prompt probing. *Inf. Process. Manag.* 60, 103510. <https://doi.org/10.1016/j.ipm.2023.103510>. <https://www.sciencedirect.com/science/article/pii/S0306457323002479>.
- Qian, C., Liu, W., Liu, H., Chen, N., Dang, Y., Li, J., Yang, C., Chen, W., Su, Y., Cong, X., Xu, J., Li, D., Liu, Z., Sun, M., 2024. ChatDev: communicative agents for software development. In: Ku, L.-W., Martins, A., Srikumar, V. (Eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Bangkok, Thailand, pp. 15174–15186. <https://doi.org/10.18653/v1/2024.acl-long.810>. <https://doi.org/10.18653/v1/2024.acl-long.810>.
- Qin, Y., Zhou, E., Liu, Q., Yin, Z., Sheng, L., Zhang, R., Qiao, Y., Shao, J., 2024a. MP5: a multi-modal open-ended embodied system in minecraft via active perception. In: *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 16307–16316. <https://doi.org/10.1109/CVPR52733.2024.01543>.
- Qin, Y., Liang, S., Ye, Y., Zhu, K., Yan, L., Lu, Y., Lin, Y., Cong, X., Tang, X., Qian, B., Zhao, S., Tian, R., Xie, R., Zhou, J., Gerstein, M., Li, D., Liu, Z., Sun, M., 2024b. ToolLLM: facilitating large language models to master 16000+ real-world APIs. In: *The Twelfth International Conference on Learning Representations*.
- Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., et al., 2021b. Learning transferable visual models from natural language supervision. In: *Proceedings of the 38th International Conference on Machine Learning*. In: *Proceedings of Machine Learning Research*, vol. 139, PMLR, pp. 8748–8763.
- Rahman, N., Santacana, E., 2023. Beyond fair use: legal risk evaluation for training LLMs on copyrighted text. In: *ICML Workshop on Generative AI and Law*.
- Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., et al., 2021. Zero-shot text-to-image generation. In: *Proceedings of the 38th International Conference on Machine Learning*. In: *Proceedings of Machine Learning Research*, vol. 139, PMLR, pp. 8821–8831.
- Rao, H., 2023. Ethical and legal considerations behind the prevalence of ChatGPT: risks and regulations. *Front. Comput. Intell. Syst.* 4, 23–29. <https://doi.org/10.54097/fcis.v4i1.9418>.
- Ray, P.P., 2023. ChatGPT: a comprehensive review on background, applications, key challenges, bias, ethics, limitations and future scope. *Internet Things Cyber-Phys. Syst.* 3, 121–154. <https://doi.org/10.1016/j.iotcps.2023.04.003>.
- Reshamwala, A., Mishra, D., Pawar, P., 2013. Review on natural language processing. *IRACS Energy Sci. Technol. Int. J.* 3, 113–116.
- Rillig, M.C., Ågerstrand, M., Bi, M., Gould, K.A., Sauerland, U., 2023. Risks and benefits of large language models for the environment. *Environ. Sci. Technol. Lett.* 57, 3464–3466. <https://doi.org/10.1021/acs.est.3c01106>. PMID: 36821477.
- Rizzo, M.G., Cai, N., Constantinescu, D., 2024. The performance of ChatGPT on orthopaedic in-service training exams: a comparative study of the GPT-3.5 turbo and GPT-4 models in orthopaedic education. *J. Orthop.* 50, 70–75. <https://doi.org/10.1016/j.jor.2023.11.056>. <https://www.sciencedirect.com/science/article/pii/S0972978X2300332X>.
- Rozière, B., Gehring, J., Gloeckle, F., Sootla, S., Gat, I., Tan, X.E., Adi, Y., Liu, J., Remez, T., Rapin, J., Kozhevnikov, A., Evtimov, I., Bitton, J., Bhatt, M., Canton-Ferrer, C., Grattafiori, A., Xiong, W., Défossez, A., Copet, J., Azhar, F., Touvron, H., Martin, L., Usunier, N., Scialom, T., Synnaeve, G., 2023. Code Llama: open foundation models for code. *CoRR*. <https://doi.org/10.48550/arXiv.2308.12950>.
- Scanlon, M., Breiting, 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. Investig.* 46, 301609. <https://doi.org/10.1016/j.fsidi.2023.301609>. <https://www.sciencedirect.com/science/article/pii/S266628172300121X>.
- Scanlon, M., Nikkel, B., Geradts, Z., 2023b. Digital forensic investigation in the age of ChatGPT. *Forensic Sci. Int. Digit. Investig.* 44, 301543. <https://doi.org/10.1016/j.fsidi.2023.301543>.
- Shafee, S., Bessani, A., Ferreira, P.M., 2025. Evaluation of LLM-based chatbots for OSINT-based cyber threat awareness. *Expert Syst. Appl.* 261, 125509. <https://doi.org/10.1016/j.eswa.2024.125509>. <https://www.sciencedirect.com/science/article/pii/S095717424023765>.
- Shah, M.S.M.B., Saleem, S., Zulqarnain, R., 2017. Protecting digital evidence integrity and preserving chain of custody. *J. Digit. Forensics Secur. Law* 12, 12.
- Shen, Y., Heacock, L., Elias, J., Hentel, K.D., Reig, B., Shih, G., Moy, L., 2023. ChatGPT and other large language models are double-edged swords. *Radiology* 307, e230163. <https://doi.org/10.1148/radiol.230163>. PMID: 36700838.
- Shinn, N., Cassano, F., Gopinath, A., Narasimhan, K.R., Yao, S., 2023. Reflexion: language agents with verbal reinforcement learning. In: Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., Levine, S. (Eds.), *Advances in Neural Information Processing Systems*, vol. 36. https://proceedings.neurips.cc/paper_files/paper/2023/file/1b44b878bb782e6954cd888628510e90-Paper-Conference.pdf.
- Silalahi, S., Ahmad, T., Studiawan, H., 2023. Transformer-based sentiment analysis for anomaly detection on drone forensic timeline. In: *2023 11th International Symposium on Digital Forensics and Security (ISDFS)*, pp. 1–6. <https://doi.org/10.1109/ISDFS58141.2023.10131749>.
- Srivastava, A., Rastogi, A., Rao, A., Shob, A.A.M., Abid, A., et al., 2023. Beyond the imitation game: quantifying and extrapolating the capabilities of language models. *Trans. Mach. Learn. Res.*
- Subramonyam, H., Pea, R., Pondoc, C., Agrawala, M., Seifert, C., 2024. Bridging the gulf of envisioning: cognitive challenges in prompt based interactions with LLMs. In: *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. CHI'24. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3613904.3642754>.
- Sweetser, P., 2024. Large language models and video games: a preliminary scoping review. In: *Proceedings of the 6th ACM Conference on Conversational User Interfaces*. CHI'24. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3640794.3665582>.
- Tan, W., Ding, C., Jiang, J., Wang, F., Zhan, Y., Tao, D., 2024. Harnessing the power of MLLMs for transferable text-to-image person ReID. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 17127–17137. <https://doi.ieeecomputersociety.org/10.1109/CVPR52733.2024.01621>.
- Tang, R., Chuang, Y.-N., Hu, X., 2024. The science of detecting LLM-generated text. *Commun. ACM* 67, 50–59. <https://doi.org/10.1145/3624725>.

- Tewel, Y., Shalev, Y., Schwartz, I., Wolf, L., 2022. ZeroCap: zero-shot image-to-text generation for visual-semantic arithmetic. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 17918–17928. <https://doi.ieeecomputersociety.org/10.1109/CVPR52688.2022.01739>.
- Thapa, A., Patil, R., 2024. ChatGPT based ChatBot application. In: *IEEE SoutheastCon*, pp. 157–164. <https://doi.org/10.1109/SoutheastCon52093.2024.10500264>.
- Thapa, S., Naseem, U., Nasim, M., 2023. From humans to machines: can ChatGPT-like LLMs effectively replace human annotators in NLP tasks. In: *Workshop Proceedings of the 17th International AAAI Conference on Web and Social Media*. <https://doi.org/10.36190/2023.15>.
- Thirunavukarasu, A.J., Ting, D.S.J., Elangovan, K., Gutierrez, L., Tan, T.F., Ting, D.S.W., 2023. Large language models in medicine. *Nat. Med.* 29, 1930–1940. <https://doi.org/10.1038/s41591-023-02448-8>.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., et al., 2023. LLaMA: open and efficient foundation language models abs/2302.13971. <https://doi.org/10.48550/arXiv.2302.13971>.
- Tsai, M.-L., Ong, C.W., Chen, C.-L., 2023. Exploring the use of large language models (LLMs) in chemical engineering education: building core course problem models with Chat-GPT. *Educ. Chem. Eng.* 44, 71–95. <https://doi.org/10.1016/j.ece.2023.05.001>. <https://www.sciencedirect.com/science/article/pii/S1749772823000180>.
- Uppal, S., Bhagat, S., Hazarika, D., Majumder, N., Poria, S., et al., 2022. Multimodal research in vision and language: a review of current and emerging trends. *Inf. Fusion* 77, 149–171. <https://doi.org/10.1016/j.inffus.2021.07.009>. <https://www.sciencedirect.com/science/article/pii/S1566253521001512>.
- Valjarević, A., Venter, H., Petrović, R., 2016. ISO/IEC 27043:2015 – Role and application. In: *2016 24th Telecommunications Forum (TELFOR)*. IEEE, pp. 1–4. <https://doi.org/10.1109/TELFOR.2016.7818718>.
- van Baar, R., van Beek, H., van Eijk, E., 2014. Digital Forensics as a Service: A game changer. *Digit. Investig.* 11, S54–S62. <https://doi.org/10.1016/j.diin.2014.03.007>. <https://www.sciencedirect.com/science/article/pii/S1742287614000127>. *Proceedings of the First Annual DFRWS Europe*.
- van Beek, H., van Eijk, E., van Baar, R., Ugen, M., Bodde, J., Siemelink, A., 2015. Digital Forensics as a Service: Game on. *Digit. Investig.* 15, 20–38. <https://doi.org/10.1016/j.diin.2015.07.004>. <https://www.sciencedirect.com/science/article/pii/S1742287615000857>. *Special Issue: Big Data and Intelligent Data Analysis*.
- van Beek, H.M., van den Bos, J., Boztas, A., Van Eijk, E., Schram, R., Ugen, M., 2020. Digital Forensics as a Service: Stepping up the game. *Forensic Sci. Int. Digit. Investig.* 35, 301021.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I., 2017. Attention is all you need. In: *Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (Eds.), Advances in Neural Information Processing Systems, vol. 30*. Curran Associates, Inc.
- Vidgof, M., Bachhofner, S., Mendling, J., 2023. Large language models for business process management: opportunities and challenges. In: *Di Francescomarino, C., Burattin, A., Janiesch, C., Sadiq, S. (Eds.), Business Process Management Forum*. Springer Nature, Switzerland, Cham, pp. 107–123. https://doi.org/10.1007/978-3-031-41623-1_7.
- Vincze, E.A., 2016. Challenges in digital forensics. *Policy Pract. Res.* 17, 183–194. <https://doi.org/10.1080/15614263.2015.1128163>.
- Wang, Y., Wang, W., Joty, S., Hoi, S.C.H., 2021. CodeT5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. In: *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, pp. 8696–8708. <https://doi.org/10.18653/v1/2021.emnlp-main.685>.
- Wang, Y., Le, H., Gotmare, A., Bui, N., Li, J., Hoi, S., 2023a. CodeT5+: open code large language models for code understanding and generation. In: *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 1069–1088. <https://doi.org/10.18653/v1/2023.emnlp-main.68>.
- Wang, W., Chen, Z., Chen, X., Wu, J., Zhu, X., Zeng, G., Luo, P., Lu, T., Zhou, J., Qiao, Y., Dai, J., 2023b. VisionLLM: large language model is also an open-ended decoder for vision-centric tasks. In: *Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., Levine, S. (Eds.), Advances in Neural Information Processing Systems*, vol. 36. https://proceedings.neurips.cc/paper_files/paper/2023/file/c1f7b1ed763e9c75e4db74b49b76db5f-Paper-Conference.pdf.
- Wang, L., Ma, C., Feng, X., Zhang, Z., Yang, H., Zhang, J., Chen, Z., Tang, J., Chen, X., Lin, Y., Zhao, W.X., Wei, Z., Wen, J., 2024a. A survey on large language model based autonomous agents. *Front. Comput. Sci.* 18, 186345. <https://doi.org/10.1007/s11704-024-40231-1>.
- Wang, Y., Jiang, Z., Chen, Z., Yang, F., Zhou, Y., Cho, E., Fan, X., Lu, Y., Huang, X., Yang, Y., 2024b. RecMind: large language model powered agent for recommendation. In: *Findings of the Association for Computational Linguistics: NAACL 2024*. <https://doi.org/10.18653/v1/2024.findings-nacl.271>.
- Wang, Z., Li, M., Xu, R., Zhou, L., Lei, J., et al., 2022. Language Models with Image Descriptors Are Strong Few-Shot Video-Language Learners. In: *Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., Oh, A. (Eds.), Advances in Neural Information Processing Systems*, vol. 35. Curran Associates, Inc., pp. 8483–8497. https://proceedings.neurips.cc/paper_files/paper/2022/file/381ceae4a1feb1abc59c773f7e61839-Paper-Conference.pdf.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., et al., 2022. Chain-of-thought prompting elicits reasoning in large language models. In: *Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., Oh, A. (Eds.), Advances in Neural Information Processing Systems*, vol. 35. Curran Associates, Inc., pp. 24824–24837. https://proceedings.neurips.cc/paper_files/paper/2022/file/9d5609613524ecf4f15af0f7b31abca4-Paper-Conference.pdf.
- Wickramasekara, A., Scanlon, M., 2024. A framework for integrated digital forensic investigation employing AutoGen AI agents. In: *2024 12th International Symposium on Digital Forensics and Security (ISDFS)*, pp. 01–06. <https://doi.org/10.1109/ISDFS60797.2024.10527235>.
- Wickramasekara, A., Densmore, A., Breiting, F., Studiawan, H., Scanlon, M., 2025. AutoDFBench: A Framework for AI Generated Digital Forensic Code and Tool Testing and Evaluation. In: *Digital Forensics Doctoral Symposium. DFDS 2025*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3712716.3712718>.
- Wiggins, W.F., Tejani, A.S., 2022. On the opportunities and risks of foundation models for natural language processing in radiology, radiology. *Artif. Intell.* 4, e220119. <https://doi.org/10.1148/ryai.220119>.
- Wu, J., Gan, W., Chen, Z., Wan, S., Yu, P.S., 2023a. Multimodal large language models: a survey. In: *2023 IEEE International Conference on Big Data (BigData)*, pp. 2247–2256. <https://doi.org/10.1109/BigData59044.2023.10386743>.
- Wu, Q., Bansal, G., Zhang, J., Wu, Y., Zhang, S., Zhu, E.E., Li, B., Jiang, L., Zhang, X., Wang, C., 2023b. AutoGen: Enabling Next-Gen LLM Applications via Multi-Agent Conversation. Technical Report MSR-TR-2023-33, Microsoft. <https://www.microsoft.com/en-us/research/publication/autogen-enabling-next-gen-llm-applications-via-multi-agent-conversation-framework/>.
- Wu, T., Breiting, F., O'Shaughnessy, S., 2020. Digital forensic tools: recent advances and enhancing the status quo. *Forensic Sci. Int. Digit. Investig.* 34, 300999. <https://doi.org/10.1016/j.fsdi.2020.300999>.
- Wu, W., Mao, S., Zhang, Y., Xia, Y., Dong, L., Cui, L., Wei, F., 2024. Mind's eye of LLMs: visualization-of-thought elicits spatial reasoning in large language models. In: *The Thirty-eighth Annual Conference on Neural Information Processing Systems. NeurIPS*. <https://neurips.cc/virtual/2024/poster/96156>, 2024.
- Xu, F.F., Alon, U., Neubig, G., Hellendoorn, V.J., 2022. A systematic evaluation of large language models of code. In: *Proceedings of the 6th ACM SIGPLAN International Symposium on Machine Programming. MAPS 2022*. Association for Computing Machinery, New York, NY, USA, pp. 1–10. <https://doi.org/10.1145/3520312.3534862>.
- Xu, E., Zhang, W., Xu, W., 2024a. Transforming digital forensics with large language models: unlocking automation, insights, and justice. In: *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management. CIKM'24*. Association for Computing Machinery, New York, NY, USA, pp. 5543–5546. <https://doi.org/10.1145/3627673.3679091>.
- Xu, Z., Jiang, F., Niu, L., Jia, J., Lin, B.Y., Poovendran, R., 2024b. SafeDecoding: defending against jailbreak attacks via safety-aware decoding. In: *Ku, L.-W., Martins, A., Srikumar, V. (Eds.), Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Bangkok, Thailand, pp. 5587–5605. <https://doi.org/10.18653/v1/2024.acl-long.303>.
- Yang, J., Jin, H., Tang, R., Han, X., Feng, Q., Jiang, H., Zhong, S., Yin, B., Hu, X., 2024a. Harnessing the power of LLMs in practice: a survey on ChatGPT and beyond. *ACM Trans. Knowl. Discov. Data* 18 (6), 160. <https://doi.org/10.1145/3649506>.
- Yang, Y., Tian, B., Yu, F., He, Y., 2024b. An anomaly detection model training method based on LLM knowledge distillation. In: *2024 International Conference on Networking and Network Applications (NaNA)*, pp. 472–477. <https://doi.org/10.1109/NaNA63151.2024.00084>.
- Yao, Y., Duan, J., Xu, K., Cai, Y., Sun, Z., Zhang, Y., 2024. A survey on large language model (LLM) security and privacy: the good, the bad, and the ugly. *HighConfidence Comput.* 4, 100211. <https://doi.org/10.1016/j.hcc.2024.100211>. <https://www.sciencedirect.com/science/article/pii/S266729522400014X>.
- Yu, Z., Wen, M., Guo, X., Jin, H., 2024a. Maltracker: a fine-grained NPM malware tracker copiled by LLM-enhanced dataset. In: *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis. ISSTA 2024*. Association for Computing Machinery, New York, NY, USA, pp. 1759–1771. <https://doi.org/10.1145/3650212.3680397>.
- Yu, H., Shen, B., Ran, D., Zhang, J., Zhang, Q., Ma, Y., Liang, G., Li, Y., Wang, Q., Xie, T., 2024b. CoderEval: a benchmark of pragmatic code generation with generative pre-trained models. In: *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering. ICSE'24*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3597503.3623316>.
- Yu, Z., Zhang, X., Shang, N., Huang, Y., Xu, C., Zhao, Y., Hu, W., Yin, Q., 2024c. Wave-Coder: widespread and versatile enhancement for code large language models by instruction tuning. In: *Ku, L.-W., Martins, A., Srikumar, V. (Eds.), Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Bangkok, Thailand, pp. 5140–5153. <https://doi.org/10.18653/v1/2024.acl-long.280>.
- Yu, Z., Zhao, Y., Cohan, A., Zhang, X.-P., 2024d. HumanEval Pro and MBPP Pro: Evaluating Large Language Models on Self-invoking Code Generation. *CoRR arXiv.2412.21199* [abs]. <https://doi.org/10.48550/arXiv.2412.21199>.
- Zhang, H., Du, W., Shan, J., Zhou, Q., Du, Y., Tenenbaum, J., Shu, T., Gan, C., 2023. Building cooperative embodied agents modularly with large language models. In: *NeurIPS 2023 Foundation Models for Decision Making Workshop*.
- Zhang, J., Huang, J., Jin, S., Lu, S., 2024. Vision-language models for vision tasks: a survey. <https://doi.org/10.1109/TPAMI.2024.3369699>.

- Zhang, Y., Sun, S., Galley, M., Chen, Y.-C., Brockett, C., et al., 2020. DIALOGPT: large-scale generative pre-training for conversational response generation. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations. Association for Computational Linguistics, pp. 270–278. <https://doi.org/10.18653/v1/2020.acl-demos.30>.
- Zhao, H., Chen, H., Yang, F., Liu, N., Deng, H., Cai, H., Wang, S., Yin, D., Du, M., 2024. Explainability for large language models: a survey. *ACM Trans. Intell. Syst. Technol.* 15 (2). <https://doi.org/10.1145/3639372>.
- Zhao, Y., Misra, I., Krähenbühl, P., Girdhar, R., 2023. Learning video representations from large language models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 6586–6597. <https://doi.org/10.1109/CVPR52729.2023.00637>.
- Zheng, Q., Xia, X., Zou, X., Dong, Y., Wang, S., et al., 2023. CodeGeeX: a pre-trained model for code generation with multilingual benchmarking on HumanEval-X. 5673–5684 <https://doi.org/10.1145/3580305.3599790>.
- Zhou, D., Schärli, N., Hou, L., Wei, J., Scales, N., Wang, X., Schuurmans, D., Cui, C., Bousquet, O., Le, Q.V., et al., 2022. Least-to-most prompting enables complex reasoning in large language models. In: *The Eleventh International Conference on Learning Representations*.
- Zhou, R., 2024. Empirical study and mitigation methods of bias in LLM-based robots. *Acad. J. Sci. Technol.* 12, 86–93. <https://doi.org/10.54097/re9qp070>.
- Zhou, Y., Muresanu, A.I., Han, Z., Paster, K., Pitis, S., Chan, H., Ba, J., 2023. Large language models are human-level prompt engineers. In: *The Eleventh International Conference on Learning Representations*.
- Zhu, D., Chen, J., Shen, X., Li, X., Elhoseiny, M., 2024. MiniGPT-4: enhancing vision-language understanding with advanced large language models. In: *The Twelfth International Conference on Learning Representations*.
- Zou, J., Zhang, S., Qiu, M., 2024. Adversarial attacks on large language models. In: Cao, C., Chen, H., Zhao, L., Arshad, J., Asyhari, T., Wang, Y. (Eds.), *Knowledge Science, Engineering and Management*. Springer Nature Singapore, Singapore, pp. 85–96. https://10.1007/978-981-97-5501-1_7.