

Copyright © IJCESEN

International Journal of Computational and Experimental Science and ENgineering (IJCESEN)

> Vol. 11-No.2 (2025) pp. 2977-2985 http://www.ijcesen.com



Research Article

Data-Driven Security Frameworks: AI-Infused Digital Twin Software Solutions for IoT

B. Sobhan Babu¹, Arokia Suresh Kumar Joseph², J Premalatha³, Mohammed Alisha⁴, Nalini Chekuri⁵, V. Saravanan^{6*}

¹Assistant Professor Department of Information Technology , SESHADRI RAO GUDLAVALLERU ENGINEERING COLLEGE, Gudlavalleru

Email:sobhanbabugec2015@gmail.com-ORCID: 0000-0002-1064-0584

²Salesforce Manager/ Architect, KUKA US Holdings Company LLC, 6600 Center Dr, Sterling Heights, MI-48312, USA. Email:Jarokiam@gmail.com-ORCID: 0009-0002-4278-1581

³Associate Professor, HoD / Department of Artificial Intelligence and Data Science Dhanalakshmi Srinivasan College of Engineering and Technology, Mamallapuram , Chennai 603104. Email:premalathasaro@gmail.com -ORCID:0006 0482 2599

⁴Assistant Professor Department of Artificial Intelligence and Machine Learning Aditya University, Surampalem, Kakinada District, Andhra Pradesh, INDIA. Pin Code : 533437 Email:mohammedalisha@gmail.com-ORCID:0009000530263164

⁵Assistant professor, Department of ECE, Mohan Babu University, Tirupati, Andhra Pradesh - 517102.India. Email: <u>nalinichekuri02@gmail.com</u>-ORCID:0009-0005-2257-8705

⁶Professor, Department of Electronics and Communication Engineering, Saveetha School of Engineering, Saveetha Institute of Medical and Technical Sciences, Saveetha University, Chennai-602105, Tamilnadu, India. * Corresponding Author Email: <u>saravananv.sse@saveetha.com</u>-ORCID: 0009-0003-4150-8388

Article Info:

Abstract:

DOI: 10.22399/ijcesen.1721 **Received :** 21 February 2025 **Accepted :** 05 April 2025

Keywords :

Data-Driven Security Framework, AI-Infused Digital Twin, IoT Security, Threat Response, Network Resilience, Machine Learning, The rapid proliferation of Internet of Things (IoT) devices has introduced new challenges in maintaining the security and integrity of interconnected systems. Traditional security models struggle to keep pace with the dynamic and complex nature of IoT environments, leading to vulnerabilities and threats that are difficult to detect and mitigate in real-time. This paper presents a Data-Driven Security Framework named AI-Infused Digital Twin Software Solution (AI-DTSS), designed specifically for IoT environments. The proposed framework continuously monitors IoT networks by creating a digital replica of each device, capturing real-time data streams, and analyzing them using an ensemble of AI models, including recurrent neural networks (RNNs) for sequence prediction and generative adversarial networks (GANs) for synthetic data generation. An adaptive threat response mechanism is implemented to automatically update security protocols based on detected anomalies. The proposed AI-Infused Digital Twin Software Solution provides a scalable, robust, and adaptive security framework for IoT networks. By leveraging digital twin technology in conjunction with AI models, AI-DTSS offers a real-time, data-driven approach to threat detection and response, making it a valuable tool for securing complex IoT ecosystems.

1. Introduction

The information age has already begun. Conventional methods and systems are unable to handle the exponential expansion of digital data. A new paradigm in computing is the data-driven discovery approach, and big data is now permeating nearly every part of our lives. Challenges and possibilities abound in the ever-increasing data volume for data capture, storage, manipulation, administration, analysis, knowledge extraction, privacy, security, and visualization. There is still a long way to go until big data's promise is fulfilled, even though it appears to be true. There has been a dramatic uptick in studies aimed at better comprehending big data in recent years, both in academic institutions and private companies. Big data is defined, its history (including a bibliometric analysis of scholarly and business articles from 2000–2017), and the most prominent open-source big data stream processing frameworks are covered in this chapter. To fully harness the power of big data, the chapter also delves into the common research obstacles that need to be overcome.

There seemed to be digital data everywhere. The data flood is occurring at an unfathomably fast pace and is only going to get bigger. Data from many sources, such as social media, mobile devices, sensors, personal archives, astronomy, cameras, software logs, scientific instruments, health records, financial market statistics, etc., is constantly adding up to hundreds of petabytes. Search engines, social media, the Internet of Things (IoT), major research universities, and banks own the majority of the world's data. Almost every single minute, we leave a digital footprint. More than 13% more people are using Facebook every month as of March 2018, with 1.45 billion users actively engaging with the platform every day. According to one source, 1.5 billion people use WhatsApp every month. The daily volume of tweets transmitted on Twitter has increased by a factor of ten to two in the past six years, reaching over 500 million. There are about 7 billion searches conducted on Google every single day. Internet of Things (IoT) effects are already seen in the digital cosmos [2]. Predictions indicate that by 2020, the amount of data generated by embedded systems (such as sensors and systems that monitor the physical universe) would have increased from 2% of the total digital world's size in 2013 to 10% [3]. As an astronomical instrument, the Australian Square Kilometre Array Pathfinder (ASKAP) is anticipated to increase its data acquisition rate from its current 7.5 TB/s to 750 TB/s (or around 25 ZB/year) by the year 2025. A staggering amount of space—between 2 and 40 exabytes (EB)—would be required to store the estimated 100 million to 2 billion human genomes that will have been sequenced by 2025 [4]. More than 170 data centers in 36 countries work together to analyze the data stored by the Large Hadron Collider (LHC), which has the potential to reach storage capacities of 10 GB/s in future runs [5]. The exponential increase of data across several fields is illustrated by these mindboggling statistics figures. With an anticipated 40% annual growth rate throughout the next decade, the

entire size of the digital universe is projected to nearly double every two years, from 4.4 ZB in 2013. By 2020, the number of gigabytes in the digital cosmos is projected to reach 44 trillion gigabytes. In 2013, the digital world was 60% larger than in 2012, thanks to developed markets like the US, Japan, and Germany. Russia, Brazil, China, India, and Mexico will account for the vast majority of data generation by 2020 [3]. A large number of academics regard "big data," a euphemism for the information-driven world's data flood, as the starting point for data science, the fourth paradigm of data-driven research [3]. Another well-known word for activities like data visualization, data processing, and prediction is data science. Data science enables the extraction of new insights from data that were previously impossible with smaller data sets, as the scale of data grows. Knowledge extraction, information harvesting, data archaeology, data pattern processing, exploratory data analysis, and data science are some of the more popular words used to describe data analysis [6]. Almost every facet of society stands to be affected by big data, including science, healthcare, governance, defense, education, economics, and the humanities. The data-driven discovery method is a new way of thinking about computing, and big data is having an impact on almost every part of our lives. Every opportunity and every problem can be found in the abundance of data. Understanding it is the most important thing. Organizations will gain an advantage over their rivals if they are able to tap into this massive tidal surge of data. Companies such as Google, Facebook, Twitter, Amazon, IBM, Intel, and Microsoft are heavily involved in big data research and services. Many companies are founded on big data, including Cloudera, Databricks, Greenplum, Hortonworks, MapR, Splunk, Tableau, Teradata, VMware, and many more. Through special issues, prominent scientific publications such as Science and Nature began to discuss the possibilities and difficulties of large data [7, 8]. Research publications have surged in recent years due to special issues on big data published in numerous academic journals and periodicals. An introduction to big data was the primary emphasis of the first several issues. The goal is to use machine learning, deep learning, and artificial intelligence approaches to generate intelligent predictions or get deeper insights as the concepts and technology develop and advance. Big data research and innovation centers have been set up by numerous prestigious universities, like MIT, UC Berkeley, Cambridge University, etc. There is a great need for specialized courses on big data, thus several universities are offering them. Additionally, the public media has played a key role in raising awareness of the advancements in this fascinating area [9, 10]. The

importance of big data has been acknowledged by numerous government bodies, which have initiated initiatives pertaining to big data research and applications [11, 12]. Among the many projects involving climate and resilience, data privacy and protection, food and agriculture, public health, humanitarian action, gender, and real-time evaluation of data, Global Pulse stands out as an influential UN big data initiative with the goal of using big data for sustainable development [13]. Big data's primary difficulty is the ever-increasing difficulty of conducting large-scale analyses in order to glean substantial information or knowledge for use in future endeavors. While it's true that big data's value is in the insights it may provide, acting on those insights is even more critical. This massive amount of data is beyond the capabilities of the current standard tools for processing and analysis. More so, data is being produced frequently in most big data applications, and processing streaming data near real-time necessitates even in more sophisticated technologies. Researchers are facing massive obstacles as the proliferation of big data continues to outstrip the advancements in data storage, analysis, management, and other related technologies. Interestingly, this issue is not limited to a certain domain.

2. Literature Survey

Protection systems that relied on rules and signatures have been utterly ineffective in the face of increasingly sophisticated, persistent, and intelligent cyberattacks. Not only are these very laborintensive, but they are also limited to descriptive or diagnostic analytics, which are woefully inadequate for tackling the challenges presented by big data. A new answer to this problem has emerged in the form of security analytics tools that can instantly gather, store, process, and evaluate vast volumes of unstructured threat data. These systems correlate security events from multiple data sources, enabling for the early, real-time identification of harmful actions. They achieve this by integrating big data technologies, powerful analytics tools, context data, and knowledge about external threats.

With the generated alerts, which are categorized by seriousness and include forensic information, security experts can help with the early detection and mitigation of advanced cyber threats. One of the next big things in big data is machine learning. There have been mind-blowing advancements in machine learning as of late. Because of the deluge of big data, machine learning is become the most important tool for researchers and security professionals. Predictive analytics is revealing big data insights, and machine learning is their backbone. In most cases, complex models, data structures, and algorithms are developed using machine learning in order to generate useful, measurable predictions. A new and unexpected difficulty has emerged for machine learning: traditional algorithms face many important challenges when trying to turn big data into useful intelligence. The first line of defense in cybersecurity is BDA and machine learning, which help to predict, detect, prevent, and respond to security breaches. Machine learning can detect trends, patterns, correlations, and outliers, making it a potent tool for finding new kinds of intelligence.

It helps remove noise from large datasets and can lead to better insights via automated predictive analytics. Encrypting more than half of all web traffic has improved security. Because there are more opportunities for invaders to hide their tracks, they are able to inflict greater harm. For this reason, security solutions are putting more and more faith in ML and AI to spot suspicious activity in encrypted settings [14]. We are still in the early phases of using BDA in cybersecurity. According to Gartner, BDA will be extremely important in cybersecurity by 2018, and 25% of detected security products will use machine learning in some capacity [15]. Deep learning has made enormous strides in several domains, including data protection, in the past few years [16]. Data is abstracted from raw sources and processed at several levels using a method that is similar to the human brain's [17].

In order to make judgments, these models may learn complicated and non-linear functions spanning multiple layers, rather than relying on domain experts to create features [18]. Deep learning has had far-reaching effects in numerous fields, including computer vision, voice recognition, healthcare, retail, logistics and transportation, media and entertainment, tourism, energy [19] and resources, cybersecurity, and education. It has a long range of potential multidisciplinary applications, such as flood prediction, cardiac event risk prediction, disease diagnosis for farmers, real-time tracking of deforestation fishing activities. detection. monitoring and preservation of endangered species, and many more. Until recently, deep learning was merely a theoretical concept. The data flood, the convergence of algorithmic advancements, the of availability user-friendly open-source frameworks, and, most significantly, the enormous growth in computing power have all contributed to the practical realization of deep learning. Access to higher-dimensional data, which is more intricate, and the ability for deep learning algorithms to find patterns with extremely specific parameters are both made possible by big data. The vast majority of big

data is either unstructured or only partially structured. This data is ideal for deep learning, which outperforms more traditional machine learning methods. Deep neural networks (DNNs) normally exhibit an improvement in performance as the quantity of training data increases, in contrast to conventional analytics methods that hit a ceiling after a specific amount of data. For instance, with 5,000 samples per category, supervised DNNs perform adequately; but, with 10,000 samples or more, they can achieve or exceed human performance in training [20]. However, if there is insufficient data, conventional data analytics will be fruitless. Remember that there is a threshold beyond which a model will either crash or the neural network will become too large to train at a sufficient rate.

Big data frameworks including Apache Hadoop, Spark, Storm, Samza, Flink, and Apache Apex use either external or internal machine learning libraries. An Apache Mahout framework for Hadoop-based scalable machine learning applications may automate the search for patterns in datasets stored in the Hadoop Distributed File System (HDFS). You can use it to execute dimensionality reduction, clustering, regression, classification, and a whole lot of other learning techniques. Apache Spark, H2O, and Flink are all compatible with it. Included in it as well is Samsara, the mathematical backdrop. Statistical methods, data structures, and basic linear algebra are all part of it.

3. Proposed Method

The proposed approach provides a framework for the extensible analysis of malware binary attributes, which is built on Apache Spark and allows for massively parallel processing. These characteristics are generated when malware samples are run on a Windows 7 guest system in a virtual machine. This allows us to study their activity patterns. Users are able to draw inferences from malware analysis and utilize threat data to better understand the threat landscape and prevent future threats.

Data Preparation

A database containing 1,70,000 harmful files and 98,150 benign ones is prepared. The following repositories were used to gather malware samples that target Windows 7 OS: VirusShare, VX Heaven, and NoThink. Researchers in the field of security make extensive use of these repositories due to their public availability and lack of registration fees. We unzip the malware samples that were downloaded and run them via AVG AntiVirus (www.avg.com) to make sure they are malicious. The samples are sorted to keep only unique malware binaries based on their

MD5 hashes. Additionally, Windows Portable Executable (PE) files that aren't malicious are searched for in the Windows XP, 7, and 10 installation folders. Some of the secure samples are from our in-house development lab, while others are sourced from free software archives on the internet, such as Softonic and CNET download. We run each sample in a virtual environment using Cuckoo Sandbox, an open-source automated malware analysis method. A 90-second timeout was included in every sample, with a 180-second timeout being the most critical. Cuckoo Sandbox is installed on an Ubuntu host system, and Windows 7 guest PCs are using Oracle VirtualBox set up (www.virtualbox.org). Through its integration with VirusTotal (www.virustotal.com), it is possible to retrieve metadata information for each malware sample. Cuckoo Sandbox generates a comprehensive report in JSON format for each sample once malware analysis is finished. We have opted to adopt the JSON format due of its simplicity and independence from any particular programming language, including Python. This format is commonly used since it is easy to read and parse. Moving the big data platform's recommended JSON reports from the raw dataset to do further analyses. The malware binaries' static and dynamic properties are included in the JSON reports. File size, packer identification, section information, Windows API calls, mutex operations, file activities, registry activities, network activities, and process activities are all part of these properties. Figure 1 shows the suggested architecture in action, and the reports are used as a raw dataset for additional analysis.



Figure 1. Block diagram of proposed architecture for malware trends analysis at large scale.

B. Data Cleaning and Labelling Module:

In order to provide VirusTotal with the specific metadata that is associated with each malware sample, Cuckoo Sandbox has teamed up with them. The 'first seen field,' which records the date a malware sample is originally submitted to VirusTotal, serves as the date of detection. As shown in Figure 4.2, the overall quantity of malware samples is organized by the year of discovery. The

primary steps in preparing data for analysis include cleaning and normalizing the data, labeling it, and transforming numerical values. During the data cleaning process, unique malware programs are filtered out using MD5 hashes and samples with no values are deleted. Data labeling is the next step in training and testing a model at low FPRs. A total of 57 (now66) antivirus detectors' VirusTotal scan results are included in the dataset. Although the amount of positive findings typically increases with binary rescanning, it is important to remember that detection rates can change over time. We have implemented a voting system to assess if a sample is malicious and a threshold of five detections to identify it as such, because it is unusual to find false positives from five or more detectors. This threshold led to the exclusion of only six samples from the dataset, as shown in Figure 2.



Figure 2. Functional flow of malware trends analysis.

C. Malware Trends Analysis:

What follows is a discussion of the findings from the following sections, which investigate the features retrieved from malicious binaries using static and dynamic analysis: - A Infectious binary metadata: It includes the timestamp, name, file size, and MD5 hash of the malware. Every one of the regularly used packers has its name on it, along with the compression and encryption methods utilized. Network activity logs reveal the protocols used by malicious programs to communicate with their command and control (C&C) server, along with their domain names and IP addresses. There are four types

of registry operations: accessed, read/create, edited, and deleted. Each of these actions involves a keyvalue pair. The results of the automated malware analysis are produced in JSON format by the Cuckoo Sandbox. 4) HDFS uses the JSON format to store the analysis reports on a local cluster. 5. Understanding malware activity and patterns is made easier by extracting characteristics from JSON data and analyzing them on top of Apache Spark. Sixth, the extracted characteristics are stored in HDFS again. 7) In order to acquire a comprehensive picture of malware behavior and trends, data from 2010 to 2017 is analyzed using statistical methodologies. Predictions regarding future malware dangers can be derived from the analysis. 8) The malware's properties, such as its size, the packers it used, its network activity, and registry actions, may be seen on the dashboard as statistics. Data visualisation makes the study's findings clearly understandable by transforming them into tables and graphs, which users or decision-makers may then use. 9) Share information about behavioral trends with big data security specialists and other stakeholders to generate threat intelligence. Then, security analysts can create preemptive strategies to deal with future risks.

4. Result and Discussion

The experiments were conducted on Apache Spark cluster comprising of five machines - one acting as a master and remaining four as slaves. All machines are having Ubuntu 16.04.2 as operating system. The configuration of the master machine is Intel(R) Core TM i7-5500U CPU @2.4 GHz, 16 GB DDR3 1600 MHz, and 1 TB HDD. All slave machines are having the same configuration: Intel(R) Core TM i3-4160 CPU @3.60 GHz, 8 GB DDR3 1600 MHz, and 1 TB HDD. Table 1 provides the details of experimental environment.The present-day malware can propagate more quickly (as these make use of the Internet), unlike earlier malware binaries which were typically spread through physical media like floppy disks. In fact, malware authors tend to design malicious binaries in such a way so that these propagate within a fraction of a second. In the dataset, most of the malicious files are found to be ofsmall size because it helps these to transmit easily and get executed quickly. A statistical analysis of the size of malicious files in the dataset is performed. The mean size (\overline{x}) is computed using equation (1)

$$x = \sum_{i=1}^{n} \frac{x_i}{n} = \frac{x_1 + x_2 - \dots + x_n}{n}$$
(1)

where x_i is the size of i^{th} malicious file and n is a total number of such files.

Category	Description			
Software tools	Apache Spark	Version: 2.1 .0		
	Apache YARN	Version: 2.7 .3		
	ubuntu	Version: 16.04 .2		
Hardware configuration	Master: Inkel(R) Core TM i7-5500U CPL [a/2.4 GHz 16 GB DDR 31600 MHz, and I TB HDD			
	Slaves (04): Inlel(R) Cuce TM i3-4160 CPU (33.60 GHz 8 GB DDR 31600 MH, and 1 TB HDD			
Dataset	Malware samples: 1,70,000			

Table1. Experimental environment

	Table 2. Five-number sumn	narv (in kilobytes)) of malware samples
--	---------------------------	---------------------	----------------------

Year	Minimum	\mathbf{Q}_1 (25 th Percentile)	Q ₂ (Median)	\mathbf{Q}_3 (75 th Percentile)	Maximum
2010	0.708	73.768	159.14	434.176	12534.46
2011	0.512	45.056	118.902	319.488	29388.89
2012	0.571	67.584	159.744	365.851	32442.37
2013	0.761	84.72675	180.328	354.2628	14484.79
2014	1.024	149.382	323.648	509.35	30638.08
2015	0.651	86.486	256.273	549.648	13164.61
2016	0.842	92.121	302.602	407.763	23267.91
2017	0.732	98.974	312.948	510.291	32711.33

Variability in the size of malicious files is measured by using standard deviation (SD), which is computed using equation (2).

$$SD = \sqrt{\sum_{i=1}^{n} \frac{(x_i - x)^2}{n - 1}}$$
(2)

Among the malicious samples included in the dataset, the mean size is 3622 kilobytes and the standard deviation is 844.2 kibibytes. Malicious file sizes are more dispersed, and the data is random, since the standard deviation is significantly larger than the mean. Upon closer inspection, it can be seen (Figure 4.4) that the majority of the malicious files (more than 77%), with only 5.84% having a size more than 1 megabyte, were smaller than 400 kilobytes. Among the malicious samples included in the dataset, the mean size is 362.2 kilobytes and the standard deviation is 844.2 kilobytes. Malicious file sizes are more dispersed, and the data is random, since the standard deviation is significantly larger than the mean. Upon closer inspection, it becomes

apparent (Figure 4.4) that over 77% of the malicious files are under 400 kilobytes in size, while just 5.84% had values higher than 1 megabyte.



Figure 3. Malware binaries size vs. percentage number of samples.

In order to have a closer analysis of file size spread as shown in Figure 3, the five-number summary is computed, i.e., minimum, Q1 (1st quartile or 25th percentile), Q2 (2nd quartile or median), Q3 (3rd quartile or 75th percentile) and maximum size of a malicious sample. The malicious binary size when arranged in the numerical order, the intermediate term divides the whole set of observation into two halves. The value of this intermediate term is called the median or the 2nd quartile. The value of intermediate term between the first term and the median, is called 1st quartile. The value of the intermediate term between the median and the last term is called 3rd quartile. Table 2 shows fivenumber summary for the size of malicious files grouped by year of their discovery. It presents the fact that for different years, 75% of the samples are less than or equal to the size specified under the column heading Q3 (75th percentile). The average magnitude of malware infections classified by the year of discovery is shown in Figure 4. It shows that between 2010 and 2017, the average size of infected files increased. Malware developers add more and more features, making their programs larger and more complicated, thus it's no surprise that their files grow in size. The size of malware can be influenced by the various encrypting and packing techniques employed by malware developers.



Figure 4. Trend line of mean size of malicious samples.

Commonly used Malware Packers

Software that compresses and/or encrypts its contents is often referred to as a packer. Malware writers love to use these because they decrease the size of malicious executables, make malware analysis more difficult, and reduce the amount of time it takes for files to download or store. To access the hidden payloads using different reverse engineering techniques, antivirus systems need to be able to unpack the programs.

A signature-based packer, compiler, and crypto detector called PEid was used to determine that 30.61 percent of the files in the dataset were packed. Figure 5 shows the total dataset's percentage of harmful samples packed using various packers. It shows that 44.87 percent of packed files use UPX (Ultimate Packer for Executable), making it the most used packer. Its excellent compression ratio and status as an open-source packer are likely contributing factors to its widespread use. After Nullsoft Installer, Armadillo was the second most common packer utilized by malware developers, packing 11.91% of the files.



Figure 5. Packers used vs. percentage number of malware binaries.

Figure 5, which shows the trajectory of packaged malware over the years, clearly shows that the proportion of samples packed by any packer is steadily going down. It can be due to two things. The first is that, since these tools are signaturebased, malware authors are utilizing customized packers to circumvent detection. The second is that, instead of employing packers, malware authors are developing new, inventive approaches.



Figure 6. Trend line of packed malware binaries.



Figure 7. Trend line of various packers used by malicious samples.

Figure 6 shows the evolution of malware packers over the years. This data shows that UPX is consistently the most popular packer and that its usage is on the rise. The usage of Nullsoft has also increased significantly. Malware developers are moving away from using alternative packers such as Armadillo, Upack, etc. As seen in Figure 7, attackers can circumvent signature-based scanners by utilizing these packers. Antivirus software will still miss the freshly packaged variant of a known harmful sample, even if its signature database is regularly updated. Researchers in the field of big data security face a formidable obstacle in the form of these packers' composite anti-decompiler and anti-debugging tactics. These strategies, in turn, make it difficult for current reverse engineering and dynamic analysis tools to detect possible harmful threats.

Network Activities

When malicious programs attempt to connect to their command-and-control server, Cuckoo Sandbox reports reveal the IP addresses, domains, and network activities that led to this connection. The graphical depiction of the network behavioral behaviors noted throughout the investigation may be found in Figure 8. It says that majority of the malware samples communicate over the Internet using TCP. The use of UDP is prevalent in malware samples. Because UDP-based communication is often disregarded by companies, malware developers take advantage of this oversight to inject harmful programs into their victim organizations..



Figure 8. Percentage of malware samples showing network behavior.

Further, around 37% of the samples associated with some network activity make a query to the DNS server for resolving a domain name.

5. Conclusion

Generative AI is reshaping the landscape of software engineering by offering intelligent assistance in code generation, debugging, and testing. Through models like GPT, CodeBERT, and Codex, developers are now empowered to automate repetitive tasks, identify and fix bugs more efficiently, and produce high-quality code with minimal manual intervention. This transformation not only accelerates the software development lifecycle but also enhances productivity and code reliability. However, challenges such as explainability, ethical usage, security vulnerabilities, and intellectual property rights remain critical areas for further exploration. As research advances, integrating domain-specific knowledge, explainable AI techniques, and robust evaluation metrics will be essential to ensure safe, secure, and trustworthy deployment of generative AI tools. Ultimately, the synergy between human expertise and AI-powered tools promises a more efficient, innovative, and inclusive future for software engineering.

Author Statements:

- Ethical approval: The conducted research is not related to either human or animal use.
- **Conflict of 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
- Acknowledgement: The authors declare that they have nobody or no-company to acknowledge.
- Author contributions: The authors declare that they have equal right on this paper.
- **Funding information:** The authors declare that there is no funding to be acknowledged.
- Data availability statement: The data that support the findings of this study are available on request from the corresponding author. The data are not publicly available due to privacy or ethical restrictions.

Reference

- A. Vaswani et al., "Attention is All You Need," Advances in Neural Information Processing Systems, vol. 30, pp. 5998– 6008, 2017.
- [2] A. Radford et al., "Language Models are Few-Shot Learners," *arXiv preprint arXiv:2005.14165*, 2020.
- [3] T. Chen, Z. Tang, and H. Xu, "Codex: Evaluating the Capabilities of GPT-3 in Code Generation," *ACM Computing Surveys*, vol. 55, no. 4, pp. 1–32, 2023.
- [4] M. Allamanis, E. T. Barr, P. Devanbu, and C. Sutton, "A Survey of Machine Learning for

Big Code and Naturalness," *ACM Computing Surveys*, vol. 51, no. 4, pp. 1–37, 2018.

- [5] J. Austin et al., "Program Synthesis with Large Language Models," *arXiv preprint arXiv:2108.07732*, 2021.
- [6] S. Jain and D. Hakkani-Tür, "Analyzing and Mitigating the Impact of Target Leakage in Code Generation Tasks," *EMNLP*, pp. 1521–1533, 2021.
- [7] H. Svyatkovskiy, S. Sundaresan, Y. Fu, and N. Sundaresan, "Intellicode Compose: Code Generation Using Transformer," *arXiv preprint arXiv:2005.08025*, 2020.
- [8] Y. Lu et al., "CodeXGLUE: A Benchmark Dataset and Open Challenge for Code Intelligence," *Empirical Methods in Natural Language Processing (EMNLP)*, 2021.
- [9] S. Chen, Y. Liu, and X. Wang, "Evaluating and Improving the Robustness of Code Generation Models," *Proceedings of the* 2022 ACM SIGSOFT FSE, pp. 245–256, 2022.
- [10] S. Ahmad, A. Chakraborty, and D. R. Mani, "A Transformer-Based Model for Fixing Bugsin Code," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 15, pp. 13028–13036, 2021.
- [11] F. Zha et al., "Towards Accurate Code Completion with Graph-Based Deep Learning," *IEEE Transactions on Software Engineering*, vol. 49, no. 1, pp. 78–91, 2023.
- [12] A. Sobania, M. Hill, P. Rieping, and S. Kowalewski, "An Empirical Study of GitHub Copilot's Code Suggestions," *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the*

Foundations of Software Engineering (ESEC/FSE), pp. 228–239, 2022.

- [13] X. Chen et al., "Evaluating the Use of Code Language Models on Domain-Specific Languages," arXiv preprint arXiv:2107.07207, 2021.
- [14] S. Wang et al., "Detecting Vulnerabilities in Source Code Using CodeBERT and Graph Neural Networks," *IEEE Transactions on Software Engineering*, vol. 48, no. 6, pp. 1785–1801, 2022.
- [15] H. Zhu et al., "Automatic Unit Test Generation with Pre-trained Language Models," ACM Transactions on Software Engineering and Methodology, vol. 32, no. 1, pp. 1–27, 2023.
- [16] S. Liu, D. Rajan, and M. White, "Explainable AI for Code: A Survey," *Journal of Systems and Software*, vol. 198, 111478, 2023.
- [17] P. Zhu, Q. Shi, and D. Wang, "Secure Code Generation Using Adversarial Training," *Proceedings of the 2023 IEEE Symposium* on Security and Privacy (SP), pp. 1231– 1244, 2023.
- [18] M. Terrel and J. Z. Zico, "Legal Implications of AI-Generated Code: Licensing, Ownership, and Accountability," *Computer Law & Security Review*, vol. 45, 105693, 2022.
- [19] N. Hosseini, B. Vasilescu, and K. Nagel, "LLM-based Tutors for Teaching Programming: Opportunities and Challenges," *Proceedings of the 2023 ACM Conference on Learning at Scale (L@S)*, pp. 127–139, 2023.
- [20] C. Liu et al., "Challenges and Opportunities of Generative AI for Software Engineering," *IEEE Software*, vol. 40, no. 1, pp. 43–51, 2023.