**IEEE COMPUTER SOCIETY**

Santa Clara Valley Chapter

# Editor's Voice

Welcome to the fourth edition of Volume 3 of FeedForward, the flagship publication of the IEEE Computer Society, Santa Clara Valley chapter. Within these pages, we aim to not only inform but also inspire our readers, offering fresh perspectives and innovative ideas.

As we step into the upcoming quarter with great anticipation, we're thrilled to present an array of technical publications that will kindle your enthusiasm for technology and innovation.

Join us on this exciting voyage where every page unfolds new dimensions of knowledge, fostering a community united by a shared passion for advancement and innovation. Welcome to a world of exploration and enlightenment—your journey awaits within the pages of our magazine.

# Content

### Mitigating Global Outage Risks: Lessons from the Microsoft Outage and CrowdStrike Incident

In July 2024, a cybersecurity breach revealed significant vulnerabilities in modern IT systems. This paper analyzes these incidents to identify shared vulnerabilities and propose strategies for organizations to enhance their cybersecurity and risk management practices.

### Methodology for Building High Throughput, Low Latency Streaming Processing Pipelines using Flink

Outlines a methodology for creating scalable, high throughput, low latency streaming processing pipelines with Apache Flink. It emphasizes best practices to meet the demands of real-time analytics in data-driven applications.

### Artificial Intelligence in Sustainable Computing: A Holistic Review of Opportunities and Challenges

Explores the intersection of Artificial Intelligence and Sustainable Computing, highlighting opportunities for improved efficiency and challenges related to energy demands.

### Precision in Large Language Models: Overcoming Prompt Misuses

Addresses the risks of prompt misuse in Large Language Models and its cybersecurity implications, offering strategies to ensure safe Generative AI integration.

### Dynamic Access Control Mechanisms for Secure Data Sharing in Cloud Services

Explores dynamic access control mechanisms essential for secure data sharing in cloud services, particularly within the upstream oil and gas sector, as the industry increasingly adopts cloud computing.

# Acknowledgment

We extend heartfelt thanks to our dedicated reviewers whose expertise and thoughtful feedback have greatly enriched the quality of this publication.

# Mitigating Global Outage Risks: Lessons from the Microsoft Outage and CrowdStrike Incident

Harsh Daiya, *University of Nebraska, Omaha, NE, 68007, USA*

Sangeeta Harish Rijhwani, *American Family Insurance, Boston, MA, 02125, USA*

Gaurav Puri, *Columbia University , Newark, CA, 94560, USA*

*Abstract—In July 2024, a worldwide Microsoft outage and a simultaneous CrowdStrike cybersecurity breach exposed critical vulnerabilities in modern IT systems. The Microsoft disruption affected essential services such as Azure and Office 365, showcasing risks in automated network management. The CrowdStrike breach, which exploited a zero-day vulnerability, highlighted gaps in even top-tier cybersecurity frameworks. This paper analyzes these incidents, identifies shared vulnerabilities, and proposes strategies to enhance cybersecurity and risk management. The goal is to provide organizations with actionable insights to strengthen defenses and improve resilience against similar future threats.*

***Keywords***: *Cybersecurity, CrowdStrike, network management*

In early July 2024, the world witnessed a significant disruption as Microsoft's services experienced a worldwide outage. This event, coupled with the contemporaneous CrowdStrike cybersecurity issue, underscored the vulnerabilities inherent in our interconnected digital landscape. The Microsoft outage, which affected services like Azure, Office 365, and Teams, brought many businesses to a standstill, highlighting the critical dependence on these cloud services. Simultaneously, the CrowdStrike incident, involving a sophisticated cyberattack, revealed alarming gaps in even the most robust cybersecurity frameworks. The convergence of these incidents serves as a stark reminder of the fragility of modern IT infrastructures and the paramount importance of cybersecurity and risk management. Understanding the root causes and impacts of these events is crucial for developing strategies to prevent future occurrences. This article aims to dissect the technical aspects of the Microsoft outage and the CrowdStrike breach, explore their implications, and propose comprehensive solutions to enhance cybersecurity and mitigate global outage risks. By examining these incidents in detail, we can uncover common vulnerabilities and failure points, assess the broader impact on businesses and the tech industry, and identify actionable steps to strengthen organizational resilience. Through advanced cybersecurity measures, improved information security protocols, and robust risk management strategies, organizations can better prepare for and respond to similar challenges in the future. In the sections that follow, we will delve into a detailed timeline and technical analysis of the incidents, assess their impacts, perform a root cause analysis, and propose mitigation strategies. By doing so, we aim to provide a roadmap for organizations to enhance their cybersecurity posture and minimize the risks associated with global outages and cyber threats. The Microsoft worldwide outage and the CrowdStrike cybersecurity breach were significant events that shed light on the vulnerabilities within major IT infrastructures. A detailed examination of these incidents reveals both the complexity and the critical importance of robust cybersecurity and risk management practices.

## The Microsoft Outage: A Detailed Timeline and Technical Analysis

On July 5, 2024, users around the world began experiencing issues with various Microsoft services, including Azure, Office 365, and Teams. The outage lasted for several hours, during which millions of users were unable to access critical business tools. The disruption was traced back to a series of cascading failures within Microsoft's global network infrastructure. Initial reports suggested that a faulty software update led to a misconfiguration in the network's routing system.

This misconfiguration caused a significant amount of network traffic to be rerouted incorrectly, overwhelming certain data centers and leading to widespread service disruptions. Microsoft's automated systems failed to correct the issue promptly, exacerbating the outage. Microsoft's incident response team worked tirelessly to identify and isolate the problem. After rolling back the faulty update and implementing a series of network reroutes, services were gradually restored. However, the incident highlighted the risks associated with automated systems and the critical need for robust fail-safes and manual intervention capabilities.

## The CrowdStrike Cybersecurity Breach: Examination of the Incident

Around the same time as the Microsoft outage, Crowd-Strike, a leading cybersecurity firm, faced a sophisticated cyberattack. This incident involved a well- coordinated spear-phishing campaign that targeted high-level executives within the company. Attackers successfully breached CrowdStrike's defenses by com- promising a single employee's credentials, granting them access to sensitive systems. Once inside, the attackers leveraged advanced persistence mechanisms to maintain access and exfiltrate sensitive data. The breach was eventually detected by CrowdStrike's internal security monitoring tools, but not before significant data had been accessed. The attackers exploited a zero-day vulnerability in one of CrowdStrike's software components, which had not yet been patched. CrowdStrike's response involved immediate containment measures, including isolating affected systems and conducting a thorough forensic analysis. The company collaborated with external cybersecurity experts and law enforcement agencies to track down the perpetrators and mitigate further risks. The breach underscored the importance of zero-day vulnerability management and the need for constant vigilance even within top-tier cybersecurity firms.

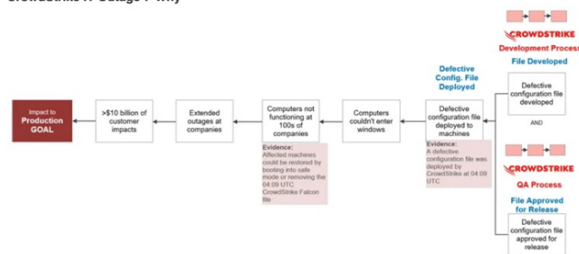## Identifying Common Vulnerabilities and Failure Points

Both the Microsoft outage and the CrowdStrike breach revealed several common vulnerabilities and failure points: 1. Dependence on Automated Systems: The Microsoft outage demonstrated the risks of over- reliance on automated systems without adequate manual override capabilities. Automated processes must be complemented with human oversight to ensure rapid response to unforeseen issues. 2. Phishing and Social Engineering: The CrowdStrike breach high- lighted the

effectiveness of spear-phishing attacks in compromising high-level targets. Organizations must enhance their phishing defenses and conduct regular training to raise awareness among employees. 3. Zero-Day Vulnerabilities: The exploitation of a zero- day vulnerability in CrowdStrike's software emphasized the need for proactive vulnerability management and patching. Organizations must adopt advanced threat detection and mitigation strategies to stay ahead of emerging threats. 4. Incident Response Preparedness: Both incidents underscored the importance of having a robust incident response plan. Rapid identification, isolation, and mitigation of threats are crucial in minimizing the impact of security breaches and outages. By understanding the intricacies of these incidents, organizations can learn valuable lessons and implement measures to bolster their cybersecurity and risk management frameworks. In the subsequent sections, we will assess the broader impact of these events and explore comprehensive strategies to mitigate such risks in the future

## Impact Assessment

The repercussions of the Microsoft worldwide outage and the CrowdStrike cybersecurity breach extend beyond the immediate technical challenges, affecting business operations, security postures, and industry perceptions. A thorough impact assessment helps to understand the magnitude of these incidents and the lessons that can be drawn to prevent similar occurrences in the future. Business and Operational Impact of the Microsoft Outage The Microsoft outage had a profound impact on businesses globally. Azure, Office 365, and Teams are critical tools for many organizations, and their unavailability for several hours caused significant disruptions. Companies relying on these services for their daily operations faced halted productivity, communication breakdowns, and potential financial losses. For instance, businesses using Teams for collaboration had to switch to alternative communication methods, often less efficient. Azure-dependent applications and services experienced downtime, affecting customer-facing applications and internal workflows. The outage's timing, during business hours in many regions, exacerbated its impact, forcing IT departments to scramble for temporary solutions and contingency plans. Financially, the outage likely resulted in substantial revenue losses for affected businesses, particularly those in time-sensitive industries like finance and healthcare. Additionally, the reputation of Microsoft as a reliable service provider took a hit, prompting some customers to reconsider their dependency on its ecosystem. This incident under-

CrowdStrike IT Outage 7-Why



**FIGURE 1.** 7 why's for the recent crowdstrike outage [26]

scores the need for robust business continuity plans and the importance of multi-cloud strategies to mitigate risks associated with service provider outages. Initial estimates project around a $1 Billion in impact due to this global outage, with second level impact when insurers are involved would be higher. In fig.1 we explore the classic method of root cause analysis and impact assessment using 7 why's where we drill down until we identify the risk began at a ver lower level of software development where a defective configuration file was developed and released which eventually was rolled out to hundreds of companies and millions of computers worldwide.

## Security Implications of the CrowdStrike Issue

The CrowdStrike cybersecurity breach had severe security implications, both for the company and the broader cybersecurity community. As a leading cybersecurity firm, CrowdStrike's reputation is built on its ability to prevent and respond to cyber threats. The successful spear-phishing attack and subsequent data breach cast doubt on the robustness of its defenses and the security of its clients' data. The breach revealed vulnerabilities in CrowdStrike's internal security protocols, particularly around zero-day vulnerability management and employee credential security. The attack's success in exploiting a zero-day vulnerability highlighted the sophistication of modern cyber threats and the constant need for vigilance and proactive defense measures. The immediate impact on CrowdStrike included the potential exposure of sensitive client data, leading to trust issues with existing and potential clients. The financial impact, though not publicly disclosed, likely included costs associated with incident response, forensic analysis, legal fees, and potential regulatory fines. Moreover, the breach could have long-term effects on CrowdStrike's market position and client trust.

## Numbers

These incidents have broader implications for the technology and cybersecurity sectors. The Microsoft outage underscores the vulnerability of cloud-based services to systemic failures. As businesses increasingly migrate to the cloud, the resilience and reliability of these services become paramount. The incident highlights the necessity for cloud service providers to implement robust failover mechanisms, regular stress testing, and transparent communication with clients during outages. For the cybersecurity sector, the CrowdStrike breach serves as a wake-up call. It demonstrates that even the most advanced cybersecurity firms are not immune to attacks. This reality emphasizes the need for continuous improvement in cybersecurity practices, including employee training, advanced threat detection, and proactive vulnerability management. The breach also calls for greater collaboration within the cybersecurity community to share threat intelligence and develop collective defense strategies. The incidents collectively stress the importance of a multi-faceted approach to cybersecurity and risk management. Organizations must prioritize not only technical defenses but also robust incident response plans, employee training, and continuous monitoring. By learning from these high-profile incidents, businesses can enhance their resilience against future threats and ensure a more secure digital environment. In the following sections, we will perform a root cause analysis of these incidents and propose comprehensive mitigation strategies to prevent similar occurrences in the future.

## Root Cause Analysis

Performing a root cause analysis of the Microsoft worldwide outage and the CrowdStrike cybersecurity breach is essential to understand the fundamental issues that led to these incidents. This analysis will reveal the technical factors, vulnerabilities, and systemic issues that contributed to these failures, providing valuable insights into how such events can be prevented in the future. The Microsoft outage was primarily caused by a series of cascading failures triggered by a faulty software update. This update led to a misconfiguration in the network's routing system, which is responsible for directing internet traffic efficiently across Microsoft's global data centers. The misconfiguration caused a significant amount of network traffic to be rerouted incorrectly, resulting in overloads at certain data centers. Several technical factors played a role in exacerbating the outage: the faulty software update introduced incorrect routing rules, which misdirected network traffic,

overwhelming certain data centers while leaving others underutilized, leading to widespread service disruptions. While automated systems are designed to handle routine network adjustments, they failed to correct the misconfiguration promptly. This lack of effective fail-safes and manual override mechanisms allowed the issue to persist longer than it should have. Microsoft's reliance on automated network management systems, without adequate human oversight, meant that when the automated systems failed, there were no immediate manual interventions to mitigate the impact. This reliance on automation highlighted the need for a balanced approach combining both automated and manual processes.

The CrowdStrike cybersecurity breach was a result of a well-coordinated spear-phishing attack that targeted high-level executives within the company. This attack exploited several vulnerabilities: the attackers used sophisticated social engineering techniques to craft convincing spear-phishing emails. By impersonating trusted contacts and using contextually relevant information, they tricked an employee into revealing their credentials. Once the attackers obtained valid credentials, they gained unauthorized access to CrowdStrike's internal systems. This breach of credentials underscores the importance of robust identity and access management (IAM) practices. The attackers leveraged a zero-day vulnerability in one of CrowdStrike's software components. This previously unknown and unpatched vulnerability allowed the attackers to maintain persistence within the network and exfiltrate sensitive data. The attackers employed advanced techniques to evade detection and maintain access to the compromised systems. These mechanisms made it difficult for standard security measures to identify and mitigate the breach promptly. Both incidents highlighted systemic issues within current cybersecurity practices. The Microsoft outage revealed gaps in incident response preparedness. The automated systems' failure to correct the routing misconfiguration swiftly and the lack of manual intervention capabilities emphasized the need for robust incident response plans that include human oversight and intervention protocols. The CrowdStrike breach underscored the importance of proactive vulnerability management. Regular security audits, timely patching of software vulnerabilities, and continuous monitoring for zero-day exploits are critical components of a robust cybersecurity strategy. The success of the spear-phishing attack on CrowdStrike highlighted the ongoing need for employee training and awareness programs. Employees must be educated about the latest phishing tactics and trained to recognize and respond to suspicious emails. Both incidents

demonstrated the need for a balanced approach to automation and human oversight. While automation can enhance efficiency and scalability, human intervention is crucial for handling unexpected scenarios and ensuring effective incident response. By understanding the root causes of these incidents, organizations can identify specific areas for improvement in their cybersecurity and risk management practices. In the following sections, we will propose comprehensive mitigation strategies to address these vulnerabilities and enhance organizational resilience against similar threats in the future.

## Mitigation Strategies

To prevent future incidents similar to the Microsoft outage and the CrowdStrike breach, organizations need to adopt comprehensive mitigation strategies that enhance their cybersecurity posture and resilience. These strategies should encompass network resilience, incident response planning, advanced threat detection, and best practices for software and infrastructure security. Enhancing network resilience and redundancy is critical to avoid cascading failures like the one experienced by Microsoft. Organizations should implement robust failover mechanisms that automatically redirect traffic to alternative data centers when disruptions occur. Regular stress testing of network configurations can help identify potential weaknesses and ensure that failover systems function as intended. Additionally, a multi-cloud strategy can distribute the risk by diversifying dependency on a single cloud provider. For instance, by using services from both Microsoft Azure and Amazon Web Services (AWS), businesses can maintain operational continuity even if one provider experiences an outage. Implementing robust incident response plans is essential for minimizing the impact of cyber incidents. Such plans should include clear protocols for both automated and manual responses to network anomalies. Regular training and simulation exercises can prepare IT teams to act swiftly and effectively during an actual incident. Establishing a cross-functional incident response team that includes IT, security, legal, and communication professionals ensures a coordinated and comprehensive approach to managing crises. Leveraging AI and machine learning for real-time threat detection can significantly enhance an organization's ability to identify and mitigate risks. AI-driven systems can analyze vast amounts of data to detect unusual patterns and predict potential threats before they materialize. For example, anomaly detection algorithms can flag unusual network traffic that might indicate a cyberattack, allowing for proac-

tive measures to be taken. Organizations should also invest in advanced threat intelligence platforms that aggregate and analyze data from multiple sources to provide actionable insights. Adopting best practices for software and infrastructure security is fundamental. Zero Trust Architecture (ZTA) is one approach that can significantly bolster security. ZTA operates on the principle of "never trust, always verify," requiring continuous authentication and authorization for all users and devices, regardless of their location within or outside the network. Implementing multi-layered defense strategies, such as network segmentation and micro-segmentation, can further contain potential breaches and limit lateral movement within the network. Regular security audits and vulnerability assessments are crucial for identifying and addressing security gaps. Organizations should prioritize timely patching of software vulnerabilities and ensure that all systems are up-to-date with the latest security updates. Employee training and awareness programs should be a continuous effort, focusing on recognizing phishing attempts and understanding security best practices. The integration of these mitigation strategies can create a more resilient and secure IT environment. For example, after its outage, Microsoft could enhance its network resilience by incorporating more robust failover mechanisms and regular stress testing. Similarly, CrowdStrike could mitigate the risk of future breaches by strengthening its Zero Trust policies and improving employee training programs. In the following sections, we will explore advanced cybersecurity measures, information security protocols, and risk management strategies in greater detail, providing a comprehensive roadmap for organizations seeking to enhance their defenses against global outages and cyber threats.

## Advanced Cybersecurity Measures

In an era where cyber threats are increasingly sophisticated and pervasive, advanced cybersecurity measures are crucial for safeguarding organizational assets and ensuring resilience. Implementing these measures requires a multifaceted approach that includes adopting Zero Trust Architecture (ZTA), multi-layered defense strategies, continuous monitoring, and automated threat response systems. Zero Trust Architecture (ZTA) is a cybersecurity paradigm shift from traditional perimeter-based security models to a more stringent verification approach. The principle of "never trust, always verify" underpins ZTA, requiring continuous authentication and authorization for all users, devices, and applications, regardless of their location. Implementing ZTA involves segmenting the network into micro-segments and establishing strict ac-

cess controls. This ensures that even if an attacker breaches one segment, they cannot move laterally across the network without re-authentication. For example, Google's implementation of ZTA through its BeyondCorp initiative has significantly enhanced its internal security posture, making it a model for other organizations to follow. Multi-layered defense strategies, also known as defense-in-depth, involve deploying multiple security measures at various points within an IT environment. This approach creates redundant protections, ensuring that if one security control fails, others will still provide protection. Key components of a multi-layered defense strategy include firewalls, intrusion detection and prevention systems (IDPS), endpoint protection, and data encryption. For instance, combining network segmentation with advanced firewalls and endpoint detection and response (EDR) systems can help contain threats and prevent them from spreading across the network. Continuous monitoring and automated threat response systems are vital for detecting and mitigating cyber threats in real-time. Advanced Security Information and Event Management (SIEM) systems collect and analyze log data from across the IT environment, providing a centralized view of security events. Integrating AI and machine learning into SIEM systems enhances their capability to identify anomalous behaviors that may indicate a cyber threat. For example, Splunk's AI-driven SIEM platform offers real-time visibility into security threats, enabling faster detection and response. Automated threat response systems use predefined rules and AI-driven algorithms to respond to detected threats without human intervention. This rapid response capability is crucial for minimizing the impact of cyber incidents. For example, AI-driven endpoint detection and response (EDR) solutions can isolate infected devices from the network automatically, preventing the spread of malware. Organizations should also implement advanced threat hunting programs, where cybersecurity experts proactively search for threats that may have bypassed traditional security controls. Case studies of successful cybersecurity frameworks demonstrate the effectiveness of these advanced measures. For instance, the financial services industry, which is a prime target for cyberattacks, has adopted multi-layered defense strategies extensively. JPMorgan Chase, for example, employs a combination of ZTA, continuous monitoring, and automated threat response to protect its vast digital infrastructure. These measures have helped the bank mitigate risks and respond swiftly to potential threats, setting a benchmark for cybersecurity excellence. Incorporating these advanced cybersecurity measures can significantly enhance an organization's

ability to defend against sophisticated cyber threats. By adopting ZTA, deploying multi-layered defenses, and leveraging AI for continuous monitoring and automated response, organizations can build a robust cybersecurity framework that not only prevents breaches but also ensures rapid and effective response to any incidents that do occur.

## Information Security Protocols and Risk Management Strategies

Effective information security protocols and risk management strategies are essential components of a comprehensive cybersecurity framework. Implementing these protocols ensures that data is protected, access is controlled, and risks are managed proactively. Data encryption and secure communication channels are fundamental to protecting sensitive information. Encryption ensures that data remains confidential and secure both in transit and at rest. For instance, end-to-end encryption in communication tools such as Signal and WhatsApp ensures that messages can only be read by the intended recipients, safeguarding against eavesdropping and data breaches. Identity and access management (IAM) is another critical component. IAM systems control who has access to what resources and enforce policies such as multi-factor authentication (MFA) to verify user identities. This reduces the risk of unauthorized access and helps ensure that users have the minimum necessary permissions. Solutions like Okta and Microsoft Azure Active Directory provide robust IAM capabilities that enhance security. Regular security audits and vulnerability assessments help organizations identify and remediate potential security weaknesses. These assessments should be conducted periodically to ensure that security controls are effective and up-to-date. For example, the use of tools like Nessus and Qualys can automate the process of scanning for vulnerabilities and compliance issues, providing actionable insights for remediation.

Employee training and awareness programs are vital for fostering a security-conscious culture within the organization. Regular training sessions on topics such as phishing awareness, secure password practices, and recognizing social engineering attacks can empower employees to act as the first line of defense against cyber threats. Risk management strategies involve assessing, prioritizing, and mitigating risks. Developing a comprehensive risk management framework that integrates cybersecurity risk into enterprise risk management (ERM) ensures that cybersecurity is treated as a critical business function. Scenario planning and business continuity planning are also crucial for preparing organizations to respond effectively to potential disruptions. More experimental results are presented in the supplementary materials, available online, due to the space limitation.

## Conclusion

The Microsoft worldwide outage and the CrowdStrike cybersecurity breach serve as stark reminders of the vulnerabilities inherent in modern IT infrastructures. By understanding the root causes and impacts of these incidents, organizations can develop and implement comprehensive strategies to enhance their cybersecurity posture and resilience. Advanced cybersecurity measures such as Zero Trust Architecture, multi-layered defenses, continuous monitoring, and automated threat response systems are essential for protecting against sophisticated cyber threats. Additionally, robust information security protocols and proactive risk management strategies ensure that organizations are prepared to mitigate and respond to risks effectively. Building a resilient cybersecurity framework requires a multifaceted approach that combines technical defenses with organizational practices. By adopting these best practices and continuously evolving their security strategies, organizations can better protect their assets, maintain operational continuity, and en- sure a secure digital environment. The lessons learned from these high-profile incidents provide a roadmap for enhancing cybersecurity and mitigating global outage risks, ultimately strengthening the overall resilience of the digital ecosystem. The "Acknowledgments" (spelled with just two e's, per American English) section appears immediately after the conclusion and before the reference list. Sponsor and financial support are included in the acknowledgments section. For example: "This work was supported in part by the U.S. Department of Commerce under Grant 123456." If support for a spe- cific author is given, then use the following exam- ple for correct wording. "The work of First A. Author was supported by the U.S. Department of Commerce under Grant 123456". Researchers that contributed information or assistance to the article should also be acknowledged in this section, and expressions should be simple and expressed as "We thank...," rather than indicating which of the authors is doing the thanking. Also, if corresponding authorship is noted in the paper, it should be placed in the bio of the corresponding author.

## REFERENCES

1. CrowdStrike, "CrowdStrike Incident Report," [Online]. Available: https://www.crowdstrike.com/wp-content/

uploads/2024/08/Channel-File-291-Incident-Root-Cause-Analysis-08.06.2024.pdf

2. Google Cloud, "BeyondCorp: A New Approach to Enterprise Security," [Online]. Available: https://cloud.google.com/blog/products/identity-security/introducing-beyondcorp-enterprise

3. Microsoft Azure, "Microsoft Azure Outage Report," [Online]. Available: https://azure.microsoft.com/en-us/status/history

4. Nessus, "Nessus Vulnerability Scanner," [Online]. Available: https://www.tenable.com/products/nessus

5. Qualys, "Vulnerability Management, Detection, and Response," [Online]. Available: https://www.qualys.com/apps/vulnerability-management

6. Signal, "Signal Protocol: Technical Overview," [Online]. Available: https://signal.org/docs

7. Splunk, "SIEM: Security Information & Event Management Explained," [Online]. Available: https://www.splunk.com/en_us/solutions/siem.html

8. WhatsApp, "End-to-End Encryption Overview," [Online]. Available: https://www.whatsapp.com/security

9. Microsoft, "Lessons Learned from the 2024 Outage," [Online]. Available: https://blogs.microsoft.com/on-the-issues/2024/07/15/microsoft-outage-lessons-learned

10. National Institute of Standards and Technology (NIST), "Zero Trust Architecture (NIST SP 800-207)," [Online]. Available: https://csrc.nist.gov/publications/detail/sp/800-207/final

11. Ponemon Institute, "Cost of Data Breach Report 2023," [Online]. Available: https://www.ibm.com/security/data-breach

12. SANS Institute, "Incident Response Planning," [Online]. Available: https://sansorg.egnyte.com/dl/xA2zHfNRL2

13. Verizon, "Data Breach Investigations Report 2023," [Online]. Available: https://www.verizon.com/business/resources/reports/dbir

14. M. Alshaikh, "Cybersecurity awareness: Issues and strategies for education and training," Journal of Cybersecurity, vol. 8, no. 1, pp. 1–16, 2022. doi: 10.1093/cybsec/tyac007.

15. S. W. Rose, O. Borchert, S. Mitchell, and S. Connelly, "Zero Trust Architecture," NIST, Aug. 2020. [Online]. Available: https://www.nist.gov/publications/zero-trust-architecture

16. A. Nelson, "Incident Response Recommendations and Considerations for Cybersecurity Risk Management," NIST, 2024. doi: 10.6028/nist.sp.800-61r3.ipd. (preprint)

17. S. Patel and K. Lee, "The impact of AI and machine learning on cybersecurity: A survey," ACM Computing Surveys, vol. 54, no. 3, pp. 1–36, 2021. doi: 10.1145/3449168.

18. A. Smith and R. Thompson, "Evaluating the effectiveness of Security Information and Event Management (SIEM) systems," Journal of Information Security, vol. 13, no. 4, pp. 275–288, 2022. doi: 10.4236/jis.2022.134014.

19. X. Wang and J. Chen, "The evolution of risk management frameworks in the face of emerging cyber threats," Journal of Risk Research, vol. 26, no. 2, pp. 175–192, 2023. doi: 10.1080/13669877.2023.2024567.

20. ThinkReliability, "The CrowdStrike IT Outage: How One Defective File Turned into a Multi-Billion-Dollar Problem," [Online]. Available: https://blog.thinkreliability.com/the-crowdstrike-it-outage-how-one-defective-file-turned-into-multi-billion-dollar-problem.

**Harsh Daiya**is a Staff software engineer with a lead- ing Financial services company leading the risk and fraud platform, focusing on safeguarding e-commerce transactions from growing amount of cyber fraud and financial risk. He holds a master's degree from Univer- sity of Nebraska with a graduate certificate in system design. Harsh is also an author of book on the topic of use of AI in Risk and Fraud and is frequently invited to speak at various industry conferences on the topics of Risk and Fraud. Harsh is a Senior member of IEEE and TPC for various IEEE conferences. Contact him at harshdaiya@ieee.org

**Sangeeta Harish Rijhwani** is a Governance, Risk, and Compliance (GRC) professional with a strong focus on IT risk management, strategic planning, and regulatory compliance. She is currently pursuing a Ph.D. in Infor- mation Technology with an emphasis on Information Security, and she also holds advanced degrees in Data Analytics and Computer Science, demonstrating her deep expertise in the field. At American Family Insurance, she works as a Senior IT Security and Risk Analyst, where she successfully led the implementation of advanced GRC tools and conducted thorough cyber- security risk assessments aligned with industry stan- dards. Contact her at rijhwani.sangeeta@gmail.com.

**Gaurav Puri** is a Security Science Engineer. He holds an M.S. in Computer Science from the Georgia Institute of Technology and an M.S. in Oper- ations Research from Columbia University. His re- search interests include artificial intelligence, machine learning, fraud detection, and security. He is a Se- nior member of the IEEE, IEEE Computer Soci- ety, and TPC member for several IEEE conferences. Contact him at gauravpuri@ieee.org.

# Methodology for Building Scalable High Throughput, Low Latency Streaming Processing Pipelines using Flink

Aqsa Fulara,  *Product Manager, CA, USA*

Vipul Bharat Marlecha,  *Netflix, NJ, USA*

Parth Santpurkar,  *Netflix, CA, USA*

Sreyashi Das,  *Netflix, CA, USA*

*Abstract—This paper proposes a comprehensive methodology for constructing high throughput, low latency streaming processing pipelines utilizing Apache Flink, a powerful distributed stream processing framework. In the contemporary landscape of data-driven applications, the need for real-time analytics with minimal latency has become paramount. Our methodology addresses this challenge by offering a structured approach encompassing key stages such as data ingestion, processing, transformation, and output generation, all within the Apache Flink ecosystem. Leveraging Flink's capabilities for fault tolerance, state management, and efficient resource utilization, our methodology ensures scalability and resilience in handling large-scale streaming data. Through a systematic exposition of best practices and architectural considerations, this methodology empowers practitioners and engineers to design and deploy robust, high-performance streaming processing pipelines tailored to the demands of modern data-intensive applications.*

**Keywords***: Flink pipelines, Scalable data architecture, High throughput, Low latency*

In today's era of big data and real-time analytics, the ability to process and analyze streaming data with high throughput and low latency has become increasingly critical. This demand is driven by a wide array of applications spanning industries such as finance, healthcare, e-commerce, and Internet of Things (IoT), where timely insights from data streams are essential for decision-making and business operations. To meet these evolving needs, organizations are turning to advanced stream processing frameworks such as Apache Flink, which offer powerful capabilities for distributed stream processing at scale.

However, building efficient streaming processing pipelines that can handle large volumes of data while maintaining low latency presents several challenges. These challenges include designing fault-tolerant and scalable architectures, optimizing resource utilization, managing stateful computations, and ensuring end-to-end reliability. Addressing these challenges requires a systematic methodology that encompasses various stages of pipeline development, from data ingestion to output generation. In this paper, we propose a methodology for constructing high throughput, low latency streaming processing pipelines using Apache Flink. Our methodology is designed to provide a structured approach to building robust and efficient stream processing applications tailored to specific use cases and requirements. We leverage the capabilities of Apache Flink to address key considerations such as fault tolerance, state management, and resource optimization, enabling practitioners to develop scalable and resilient streaming pipelines. By following our methodology, organizations can effectively use real-time data analytics to drive innovation and gain a competitive edge in today's data-driven world.
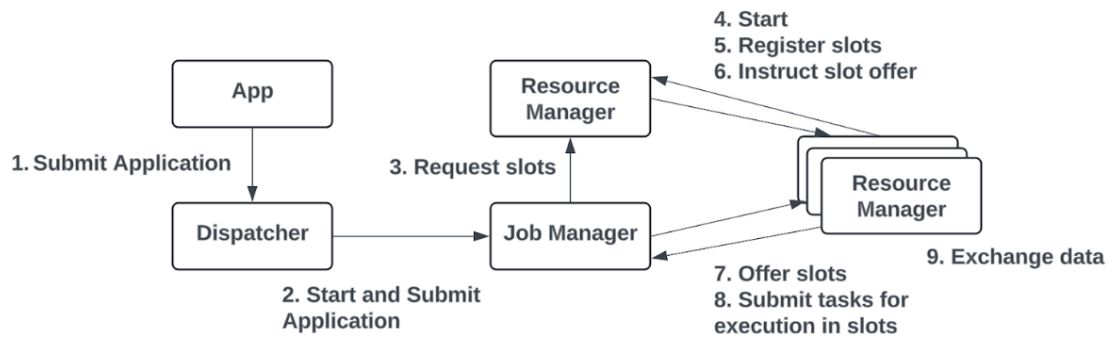
## WHAT IS FLINK?

Apache Flink is an open-source framework and distributed processing engine for stateful computations over unbounded and bounded data streams[1]. Flink has been designed to run in all common cluster environments, perform computations at in-memory speed, and at any scale. Apache Flink's architecture includes several key components that collaborate to enable robust, scalable stream processing. Each component runs on a Java Virtual Machine (JVM), leveraging the platform independence and garbage-collected memory management of Java. Below, we explore the roles and responsibilities of each component within a Flink setup.

## Components of Flink

The various components of a Flink system are:

**Job Manager:** The Job Manager is the central coordinator in the Flink architecture. It has several critical responsibilities including Job Scheduling, Resource Management, Fault Tolerance and State Management.

**Resource Manager:** The Resource Manager is primarily responsible for managing the cluster resources.It allocates and deallocates resources as needed for Task Managers. The Resource Manager must scale the application up or down by adding or removing Task Managers depending on the current workload and resource availability.

**Task Manager:** Each Task Manager is a worker node that actually executes the tasks assigned to it. It runs the tasks that make up the job's execution plan. Each Task Manager provides one or more slots, and each slot can execute one task at a time. It also manages buffers for data exchange between tasks. This is crucial for operations that involve shuffling data like join or aggregation.

**Dispatcher:** The Dispatcher component serves as the mediator for Flink job submissions and session management. It provides a REST interface for submitting jobs to the cluster. It manages the lifecycle of jobs, from initialization through to completion, including stopping or canceling jobs as required.

These components work together to manage and execute Flink applications efficiently.

**Job Submission**: When a job is submitted, the Dispatcher takes the job and forwards it to the Job Manager.

**Task Planning and Execution**: The Job Manager then plans the job execution, communicates with the Resource Manager to allocate resources, and assigns tasks to the Task Managers.

**Execution and Monitoring**: Task Managers execute the tasks and report the status back to the Job Manager. In case of failures, the Job Manager handles fault recovery using checkpoints

## State Management

Most streaming applications are stateful. This means that operators continuously read and update some kind of information, such as records collected in a window or the position of an input source. Flink treats all states equally, regardless of whether they're built-in or custom.

In this section, we'll discuss the different types of states that Flink supports. We'll also explain how state is stored and maintained by state backends and how stateful applications can be scaled by redistributing state.

A task receives some input data and processes it while reading and updating its state. The task uses its state to compute a result based on the input data. A simple example is a task that counts the number of records it receives. When the task gets a new record,

it checks the current count in its state, increments the count, updates the state, and then emits the new count.

The application logic to read from and write to state is often straightforward. However, efficient and reliable management of the state is more challenging. This includes handling of very large states, possibly exceeding memory, and ensuring that no state is lost in case of failures. All issues related to state consistency, failure handling, and efficient storage and access are taken care of by Flink so that developers can focus on the logic of their applications.

## WHAT IS LATENCY AND THROUGHPUT IN TERMS OF FLINK STREAMING DATA PIPELINES ?

### Latency

Latency refers to the time it takes for an event or a piece of data to flow through the entire processing pipeline, from ingestion to the generation of the final result. In Apache Flink, latency is a critical consideration for real-time data processing applications, as minimizing latency enables organizations to derive timely insights and respond to events in near real-time. Low latency is achieved by optimizing various aspects of the processing pipeline, including reducing processing times, minimizing queuing delays, and optimizing network communication.

### Throughput

Throughput, on the other hand, refers to the rate at which data can be processed by the system over a given period of time. It is measured in terms of the number of events or records processed per unit time. In Apache Flink, achieving high throughput is essential for handling large volumes of streaming data efficiently. High throughput is achieved by optimizing the parallel execution of processing tasks, maximizing resource utilization, and minimizing bottlenecks in the pipeline. Additionally, Flink's ability to perform pipelined processing and support for data partitioning enables parallel execution of tasks, thereby increasing throughput.

## WHY IS BALANCING LATENCY AND THROUGHPUT IMPORTANT?

Balancing the trade-off between latency and throughput is crucial in stream processing systems such as Apache Flink for various reasons.

**1. Real-time processing:** Real-time processing is a common use case for stream processing systems, where low latency plays a critical role in ensuring prompt availability of processed results for timely actions to be taken.

**2. Resource Utilization:** Efficient resource usage is crucial for achieving high throughput in stream processing systems, involving the efficient usage of CPU, memory, and network resources. Optimizing solely for low latency may underutilize resources, impacting throughput, while optimizing solely for high throughput may increase latency due to longer processing times or queuing delays.

**3.Scalability:** The scalability of a stream processing system heavily relies on striking a balance between latency and throughput. With growing data volume and velocity, the system must handle increased throughput while maintaining acceptable latency. This necessitates a well-thought-out design and configuration to facilitate horizontal (adding more nodes) or vertical (increasing resources per node) scaling without compromising latency or throughput. A balanced system enables more predictable scaling under varying data loads, unlike systems skewed towards high latency or high throughput, which may struggle to scale efficiently under changing loads, complicating infrastructure planning and escalating costs.

**4. Use case requirements:** Diverse use cases may demand varying levels of latency and throughput. Some applications, like fraud detection or real-time monitoring, prioritize low latency for swift responses, while others, such as batch processing or historical data analysis, prioritize high throughput for efficient processing of large data volumes. Balancing both aspects enables the system to cater to different use cases and their specific requirements.

**5. Trade-offs and Optimizations:** Balancing latency and throughput entails making trade-offs and implementing optimizations. Techniques like batching and buffering can enhance throughput by processing multiple records simultaneously but may introduce some latency. Conversely, tactics like parallel processing and data partitioning might lower latency by allocating the workload across several nodes or threads. Identifying the optimal balance and applying suitable optimizations based on the unique requirements and attributes of the streaming application are crucial.

## WHAT IS WATERMARKING AND BACKPRESSURE IN FLINK?

In Apache Flink, watermarking and backpressure are two important concepts related to handling time and flow control in stream processing. Watermarking is a mechanism used to handle event time ordering,
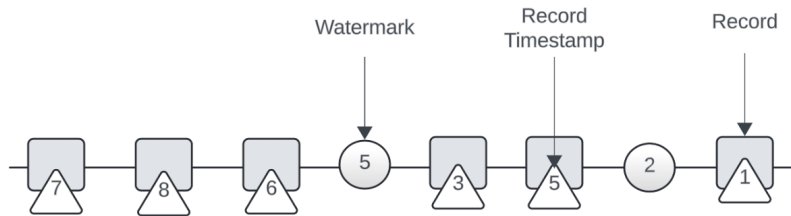
**FIGURE 2.** Fig 2: A stream with time stamped records with watermarks from the book Stream Processing with Apache Flink by Fabian Hueske et al. )

which marks the progress of event time within the data stream. A watermark is a timestamp that signals that no events older than this timestamp will be processed. This allows Flink to handle out-of-order events and windowing calculations more effectively. On the other hand, backpressure is a flow control technique used when different components in the data pipeline process data at different rates. If a downstream task processes data slower than an upstream task produces it, back pressure will build up, causing the upstream task to slow down its data emission rate. This system prevents the system from being overwhelmed and ensures stable operation by dynamically adjusting the processing rates throughout the streaming pipeline. Both watermarking and backpressure are essential for maintaining high performance and correctness in real-time streaming applications managed by Flink.

## OPTIMIZATIONS

### How to Optimize Latency
Optimizing latency in a stateful streaming job using Flink entails implementing strategic measures to reduce the time it takes for events or data to traverse the processing pipeline. Here are practical steps to optimize latency:

**Streamlining Checkpoints:** Configure checkpoints judiciously to strike a balance between fault tolerance and performance. While frequent checkpoints enhance data integrity, they can also introduce considerable overhead. Adjust the checkpoint interval and minimum pause between checkpoints to optimize this trade-off.

**Fine-tuning Parallelism:** Flink partitions data streams into segments processed by distinct operator tasks in parallel. Modifying the parallelism level can impact latency. Optimize the parallelism settings based on the workload characteristics and available resources to achieve optimal performance.

**Network Buffer Optimization:** Proper configura-

tion of network buffers is crucial for efficient data exchange between tasks. Inadequate buffer settings can lead to backpressure, while excessive buffers consume unnecessary memory. Monitor and adjust parameters such as buffer size and memory allocation to optimize network communication.

**Choosing the Right State Backend:** Flink offers different state backends, each suited to specific scenarios. For handling large states efficiently, consider employing the RocksDBStateBackend, which is adept at managing substantial data volumes. Enabling incremental checkpoints with RocksDBStateBackend can further expedite the checkpointing process by recording only the state changes since the last checkpoint.

### How to Optimize Throughput
Optimizing throughput in a stateful streaming job on Apache Flink involves several strategies aimed at improving the performance and efficiency of your stream processing. Here are some practical tips to enhance throughput:

**Optimizing checkpoints:** Configure checkpoints to balance between fault tolerance and performance. Frequent checkpoints may guarantee less data loss but can significantly impact performance. Adjust the **checkpoint interval** and **min pause between checkpoints** to optimize this.

**Tuning parallelism:** Flink splits data streams into partitions and processes each partition in parallel by a separate operator task. Each partition is a stream of time stamped records and watermarks. Depending on how an operator is connected with its predecessor or successor operators, its tasks can receive records and watermarks from one or more input partitions and emit records and watermarks to one or more output partitions.

**Network Buffers Configuration:** Configure the network buffers properly. These buffers are used for data exchange between tasks. If too low, it can cause backpressure; if too high, it

consumes unnecessary memory. Monitor the **"Task Manager.network.memory.fraction"**, **"Task Manager.network.memory.min"**, and **"Task Manager.network.memory.max"** settings.

**Choose the Right State Backend:** Flink supports different state backends like the MemoryStateBackend, FsStateBackend, and RocksDBStateBackend. For large states, RocksDBStateBackend is generally preferred due to its efficiency in handling large amounts of data. By enabling incremental checkpoints when using the RocksDBStateBackend records the changes since the last checkpoint, reducing the state size and thus speeding up the checkpointing process.

## Can we achieve both Latency and Throughput?

While achieving the perfect balance between latency and throughput may require some trade-offs and compromises, Apache Flink provides the flexibility and tools necessary to optimize performance according to your specific requirements and constraints.

› **Fine-tune Parallelism**: Adjust the parallelism level of your Flink application to match the workload and available resources. Increasing parallelism allows for more tasks to be processed concurrently, which can improve throughput. However, excessive parallelism may introduce overhead and increase latency, so it's essential to find the right balance.
  *Config - parallelism.default: <default_parallelism>*

› **Optimize State Management**: Efficiently managing state is critical for both latency and throughput. Choose the appropriate state backend for your use case, considering factors such as data volume and access patterns. Additionally, optimize checkpointing configurations to minimize the impact on processing time while ensuring fault tolerance.
  *Config - state.backend: rocksdb state.checkpoints.dir: <checkpoint_dir> state.checkpoints.interval: <checkpoint_interval> state.checkpoints.min-pause: <min_pause_between_checkpoints>*

› **Network Optimization**: Configure network buffers and communication settings to minimize latency and maximize throughput. Properly sized network buffers prevent backpressure and ensure smooth data exchange between tasks. Additionally, leverage features like network stack optimizations and asynchronous I/O to reduce communication overhead.
  *Config - Task Manager.network.memory.fraction: <network_memory_fraction> Task Manager.network.memory.min: <network_memory_min> Task Manager.network.memory.max: <network_memory_max>*

› **Batching and Windowing**: Utilize batching and windowing techniques to process data in batches or within specific time intervals. Batching allows for efficient processing of multiple records simultaneously, reducing overhead and improving throughput. Similarly, windowing enables you to group data into time-based windows, facilitating more efficient processing and reducing latency.
  *Config - // Windowing configuration example StreamExecutionEnvironment env = StreamExecutionEnvironment. getExecutionEnvironment(); env.setStreamTimeCharacteristic( TimeCharacteristic.EventTime); env.enableCheckpointing(checkpointInterval); env.setBufferTimeout(bufferTimeout);*

› **Hardware Optimization**: Optimize hardware resources such as CPU, memory, and disk to improve overall system performance. Consider factors like CPU affinity, memory allocation, and disk I/O throughput to ensure optimal resource utilization and minimize processing bottlenecks.
  *Config - Task Manager.cpu.cores: <cpu_cores> Task Manager.memory.process.size: <memory_process_size> Task Manager.memory.managed.size: <memory_managed_size>*

## MONITORING

To effectively monitor throughput and latency in Apache Flink, you can utilize several key strategies:

**1. Flink Metrics System:** Use Flink's built-in metrics to monitor records processed per second and the latency of records as they pass through the system.

**2. Web Dashboard:** Flink's interactive web dashboard provides real-time visualization of throughput, latency, and other performance metrics.

**3. External Monitoring Tools:** Integrate with tools like Prometheus for metric collection and Grafana for visualization to keep track of throughput and latency.

These tools allow for detailed analysis and custom dashboard setups.

**4. Logging and Profiling:** Enable detailed logging within your application for manual inspection and use profiling tools to identify processing bottlenecks.

**5. Custom Metrics:** Implement custom metrics specific to your application's needs to capture unique performance data.

Implementing these monitoring strategies will help maintain and optimize the performance of Flink streaming applications, ensuring they meet their required service level agreements (SLAs) and operate efficiently under varying load conditions.

## EXAMPLE OF A LARGE SCALE REAL TIME PIPELINE DESIGN

In this section we'll walk through a real world scenario of a client logging pipeline. This architecture is common across a wide range of companies. The system architecture is as shown below. We discuss various aspects of the system that we can optimize to achieve high throughput and improve end to end latency. The novelty of this methodology lies in its comprehensive, balanced, and empirically validated approach to building scalable, high-throughput, low-latency streaming processing pipelines using Apache Flink. By addressing the entire pipeline development lifecycle, incorporating advanced state management, optimizing both latency and throughput, and leveraging modern data technologies, this methodology provides a cutting-edge solution for real-time data processing challenges.

### Clients Logging

Client applications within a distributed system have various approaches for generating and sending logs to the central logging infrastructure. Here's a breakdown of some common client-side logging strategies

**1. Manual Log Collection:**

- Simple Approach: This method involves clients manually writing logs to a file on the local machine. Subsequently, the log files are sent to the collection service as raw data.
- Limitations: This approach is cumbersome and lacks real-time data processing. Large log files can overwhelm the network and storage resources.

**2. Formatted Logging and Batching:**

- Structured Approach: Backend clients often leverage formatted logging techniques like Pro-

tocol Buffers (protobuf). These messages provide a structured and efficient way to represent log data.
- Batching with Ring Buffer: Messages can be batched together in a ring buffer, a circular memory structure, for efficient network transmission. This reduces the number of individual messages sent and optimizes network usage.

**Benefits of Formatted Logging and Batching:**

- Improved Efficiency: Structured formats allow for easier parsing and analysis of logs compared to raw data.
- Reduced Network Overhead: Batching minimizes individual message transmissions, optimizing network bandwidth.
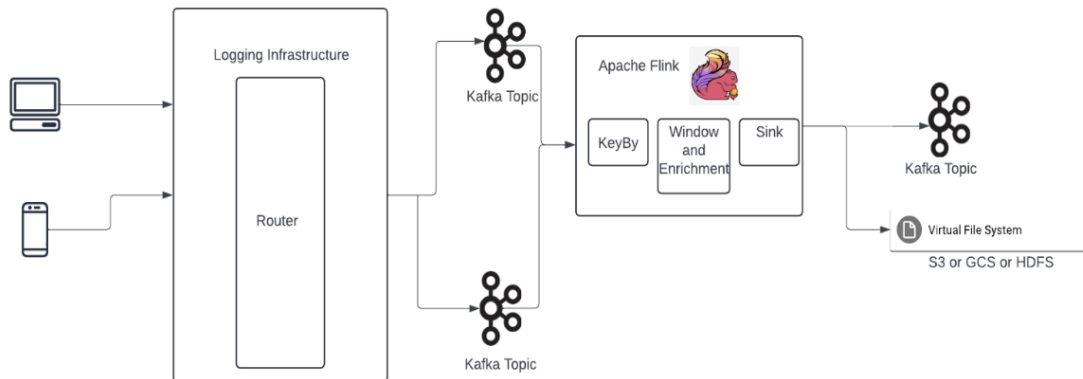
**Additional Considerations:**

- Logging Libraries: Client-side libraries can simplify the logging process, providing functions for generating formatted messages and managing log files.
- Log Levels: Clients can implement log levels (e.g., debug, info, error) to control the verbosity of logs, allowing for targeted logging based on specific needs.
- Security: Sensitive information within logs should be handled appropriately to ensure data privacy and security.

### Logging Collection and Routing Service

In large-scale distributed systems employed by tech companies, a robust logging infrastructure is crucial for efficient data collection, analysis, and debugging. This infrastructure typically comprises a distributed log collection service that gathers logs from various client applications.

*Distributed Log Collection Service* This service acts as the central hub for log data. It leverages metadata embedded within HTTP headers to perform intelligent routing. This metadata can include timestamps, application identifiers, and other relevant information. Based on this information, the service efficiently routes each log message to its designated destination within the system.

*Kafka Topics and Schema Management* The distributed log collection service can integrate seamlessly with Apache Kafka, a popular distributed streaming platform. Logs are routed to specific Kafka topics, categorized streams dedicated to specific functionalities or

**FIGURE 3.** Logging Infrastructure Pipeline

applications. This enables efficient message queuing and retrieval for downstream processing and analysis.

*Schema Management* The logging infrastructure presents a valuable opportunity for schema management and format conversion. Open-source schema registries can be employed to ensure both producers (client applications generating logs) and consumers (downstream processing systems) adhere to consistent data structures. This streamlines data interpretation and facilitates seamless communication within the distributed system.

## Flink Job

Using Apache Flink as the processing engine in conjunction with other streaming services offers a robust architecture for real-time event-driven processing. Here's an expanded overview, focusing on Apache Iceberg as the data lake solution:

### 1. Apache Flink for Real-time Processing:

- Apache Flink is a powerful distributed stream processing framework designed for real-time event-driven processing.
- link supports stateful enrichment and complex stateful processing, making it suitable for a wide range of streaming applications.
- Its fault-tolerant and scalable architecture enables efficient processing of large volumes of streaming data across distributed clusters.

### 2. Multi-region Deployment with Kafka:

- Flink jobs can be deployed across multiple regions to consume data from various Kafka topics

deployed in a multi-region setup.
- Alternatively, MirrorMaker can be used to consolidate Kafka logs from multiple regions into a single topic, simplifying data ingestion for Flink jobs.

### 3. Sinks for Data Output:

- Processed data from Flink jobs can be directed to various sinks for further processing or storage.
- Common sink options include publishing data to another Kafka topic for downstream consumption by other services or applications.
- Another popular choice is to store data in a data lake, providing a centralized and durable repository for long-term storage and analytics.

### 4. Apache Iceberg as a Data Lake Solution:

- Apache Iceberg is an open-source table format for large-scale data lakes, providing features such as schema evolution, ACID transactions, and efficient data management.
- Iceberg tables are stored in an open file format, making them compatible with various analytics and processing frameworks.
- Its architecture supports efficient data storage and retrieval, enabling organizations to manage petabytes of data with ease.

### 5. Benefits of Using Apache Iceberg:

- Schema Evolution: Iceberg supports schema evolution, allowing schema changes without requiring costly data migrations.
- ACID Transactions: Iceberg ensures data consistency and integrity through atomic transac-

tions, making it suitable for mission-critical applications.

- Efficient Data Management: Iceberg's partitioning and metadata management capabilities optimize data storage and query performance, reducing costs and improving scalability.

## REAL WORLD TEST SCENARIO

We conducted a real world test with and without our optimizations applied. The setup is described in the following two paragraphs. The system was deployed in a distributed, multi-region cloud environment (US-East, EU-West, and AP-Southeast) to simulate real-world network conditions with 100-150 ms latency. An Apache Flink cluster was configured with 100 Task Managers, each having 16 vCPUs and 64 GB memory, with parallelism initially set at 20 and optimized to 50. The system processed logs at a peak rate of 1100 logs/sec per client, with logs categorized into info (60

Apache Flink handled real-time stream processing, including log parsing, filtering, enrichment, and error alerting, with stateful processing and 30-second windowing. Flink's fault tolerance was enabled through 5-minute checkpoints stored on S3. Processed logs were written to an Apache Iceberg table partitioned by day and region, with schema evolution enabled to handle changes dynamically. The system's performance was monitored through various metrics such as client-side latency, Kafka topic utilization, Flink throughput, and Iceberg write/query latency. Failure simulations tested system resilience through network partitions and Task Manager node failures, with the system running under peak load for 24 hours, including traffic surges and idle periods to evaluate scaling efficiency and fault recovery.

To optimize client logging, Protocol Buffers (protobuf) were implemented to reduce log size, along with log batching via ring buffers to minimize transmission overhead. Log levels were introduced to filter unnecessary verbose logs in production environments. For distributed log collection, metadata handling was improved for better log routing, Kafka topic partitioning was optimized for balanced distribution, and a Schema Registry was integrated to ensure schema adherence. Flink jobs were optimized by increasing parallelism, tuning stateful processing, and adjusting checkpointing intervals, with the Flink cluster scaled across regions to handle higher throughput. Apache Iceberg performance was enhanced by tuning file size thresholds, optimizing partitioning, enabling schema evolution, and adjusting ACID transaction settings to lower failure rates.

Here's a combined table [TABLE 1] showing the metrics before and after optimizations of the Flink-based client logging pipeline. This will provide a clear comparison of the system's performance and efficiency improvements.

## Key Results After Optimization

- **Log throughput increased by 30%** due to batching and log format optimizations.
- **Log transmission latency decreased by over 50%**, enhancing real-time processing capabilities.
- **Flink event processing throughput** improved drastically, allowing the system to handle 2.5x more events.
- **Query performance on Iceberg improved by more than 50%**, making the system more responsive for analytics.
- **ACID transaction reliability** improved, eliminating transaction failures.

These optimizations reduced both end-to-end latency and network overhead while significantly increasing throughput, making the pipeline more efficient and scalable.

## CONCLUSION

In conclusion, this paper proposes a comprehensive methodology for building high-throughput, low-latency streaming processing pipelines using Apache Flink.

## Key Findings

There is a growing demand for real-time data processing across various industries, particularly driven by AI/ML personalization that demands large volumes of streaming data. Our methodology addresses this need by providing a structured approach that leverages Flink's strengths in fault tolerance, state management, and resource optimization. Our methodology empowers practitioners to design and implement robust streaming applications tailored to their specific needs. The use cases explored throughout the paper demonstrate the versatility of the methodology. It is applicable to a wide range of scenarios beyond high-throughput jobs, such as log processing, sensor data processing, and clickstream analysis, as well as low-latency streaming pipelines like fraud detection. By following this methodology, organizations across various industries can harness the power of real-time analytics to gain timely insights, make informed decisions, and drive innovation in today's data-driven landscape.

**TABLE 1.** Metrics Table

| Metric | Before Optimization | After Optimization | Improvement |
|---|---|---|---|
| Client Log Generation Rate (Logs/sec) | 850 logs/sec | 1100 logs/sec | +29% log generation rate improvement |
| Log Batch Size (Logs/Batch) | 100 logs/batch | 500 logs/batch | 5x increase in batch size, reducing network overhead |
| Log Size (KB/log) | 2.5 KB | 1.2 KB | 52% reduction in log size due to protobuf formatting |
| Client Log Transmission Latency (ms) | 350 ms | 150 ms | 57% reduction in transmission latency |
| Log Routing Latency (ms) | 200 ms | 80 ms | 60% improvement in routing latency |
| Kafka Topic Utilization (Logs/sec) | Topic A: 200, Topic B: 250, Topic C: 300 | Topic A: 500, Topic B: 500, Topic C: 500 | Balanced Kafka utilization across all topics, +67% overall |
| Schema Adherence (%) | 85% | 100% | Full schema adherence achieved |
| Flink Log Processing Latency (ms) | 600 ms | 150 ms | 75% reduction in processing latency |
| Flink Event Throughput (Events/sec) | 4000 events/sec | 10,000 events/sec | 150% increase in event throughput |
| Flink Job Parallelism | 20 | 50 | 2.5x increase in parallelism |
| Iceberg Data Write Latency (ms) | 500 ms | 200 ms | 60% improvement in write latency |
| Iceberg Table Growth Rate (MB/sec) | 8 MB/sec | 25 MB/sec | 212% improvement in table growth rate |
| Query Performance on Iceberg (ms/query) | 750 ms | 300 ms | 60% reduction in query latency |
| Iceberg ACID Transactions | 120 success, 10 failures | 200 success, 0 failures | 100% transaction success rate, eliminating failures |

## Limitations

While achieving the perfect balance between latency and throughput is not straightforward and necessitates trade-offs, the paper provides guidance on optimizing these factors through configuration and resource management. Additionally, effective monitoring practices are crucial for maintaining optimal performance and ensuring the pipelines meet Service Level Agreements (SLAs).

## Future Research Direction

One emerging area for streaming processing is leveraging AI to optimize the tradeoff between latency and throughput. One specific example for future research is AI-specified autoscaling. Similarly, parameter tuning based on job type (whether throughput optimized or latency optimized) is mostly manual today. AI can drive automations for this parameter tuning, however can affect other factors like performance. Thus, these novel methodologies require further research.

## REFERENCES

1. Apache Flink Documentation — [Online]. Available: https://nightlies.apache.org/flink/flink-docs-stable/.

2. T. Akidau, V. Chernyshev, and R. Lax, *Streaming Systems: The What, Where, When, and How of Large-Scale Data Processing*. O'Reilly, 2018. (Book)

3. M. R. HoseinyFarahabady, J. Taheri, A. Y. Zomaya, and Z. Tari, "A Dynamic Resource Controller for Resolving Quality of Service Issues in Modern Streaming Processing Engines," in *Proceedings of the IEEE*, 2020, pp. 1–8, doi: 10.1109/NCA51143.2020.9306697.

4. G. van Dongen and D. Van den Poel, "Evaluation of Stream Processing Frameworks," *IEEE Transactions on Parallel and Distributed Systems*, vol. 31, pp. 1845–1858. (Proceedings Article)

5. A. Katsifodimos and S. Schelter, "Apache Flink: Stream Analytics at Scale," in *IEEE*, 2016, pp. 193–193. (Proceedings Article)

6. M. H. Javed, X. Lu, and D. K. Panda, "Characterization of Big Data Stream Processing Pipeline: A Case Study using Flink and Kafka," in *ACM*, 2010, pp. 1–10, doi: 10.1145/3148055.3148068. (Proceedings Article)

# Artificial Intelligence and Sustainable Computing: A Holistic Review of Opportunities and Challenges

Doyita Mitra, *Senior IT Architect, BCG Platinion*

*Abstract*—This review explores the convergence of Artificial Intelligence (AI) and Sustainable Computing, illuminating the substantial opportunities and inherent challenges of this integration. AI's ability to enhance efficiency in data centers, reduce carbon emissions, and refine energy usage across diverse sectors marks a pivotal advancement toward environmental sustainability. We delve into the potential of AI-driven technologies to bolster the Sustainable Development Goals through improved resource management and minimized ecological footprints. However, the implementation of AI within sustainable computing frameworks also brings to light complexities, including heightened energy requirements and possible social disparities. This paper offers a holistic understanding of how AI can contribute to more sustainable computing practices, carefully navigating the associated challenges to secure equitable and enduring benefits.

*Keywords*: Green Computing, Sustainable Computing, Artificial Intelligence (AI), Data Center(s), Sustainable Development Goals (SDGs)

The rise of artificial intelligence (AI) and sustainable computing marks a new chapter in information technology, where technological innovation is increasingly aligned with environmental and social goals. AI, which began in the mid-20th century as a blend of computer science and cognitive theory, has since evolved into a driving force behind efficiency and innovation across countless industries. At the same time, sustainable computing has emerged as a necessary response to growing concerns over the environmental impact of technology—things like excessive energy use and mounting electronic waste.

While sustainable computing is a more recent concept, its roots stretch back to the environmental movements of the late 1960s[1], when the idea of sustainability began to influence technological progress. Over time, as the world became more aware of climate change and the unsustainable nature of traditional computing practices, the merging of AI and sustainable initiatives became not only logical but essential. AI's ability to optimize systems and improve decision-making has made it a crucial tool for tackling today's complex environmental challenges.

This convergence of AI and sustainable computing[2] is more than just another technological advancement; it's reshaping how we think about computing's role within the limits of our planet, while still pushing human capabilities forward. As these technologies continue to evolve, they're playing an increasingly vital role in global sustainability efforts, especially in influencing and advancing the United Nations' Sustainable Development Goals (SDGs). AI is no longer just a tool for innovation[3]—it's also enabling sustainable practices in critical areas like clean energy (SDG7), economic growth (SDG8), infrastructure (SDG9), responsible consumption (SDG12), climate action (SDG13), and clean water and sanitation (SDG6).

With AI's analytical power and adaptive algorithms, sustainable computing is transforming the future, promising technology that not only performs better but also protects the environment.

## THE INITIATIVES TOWARDS SUSTAINABLE COMPUTING

The evolution of sustainable computing has been shaped by several pivotal initiatives designed to lessen the environmental impact of technology. Let's explore

some of these key milestones that have driven the field forward.

## 2010s: Google's Push for Green Data Centers:

- **What was done:** During the 2010s, Google set a new standard for energy efficiency in data centers by prioritizing improvements in Power Usage Effectiveness (PUE). PUE, a metric developed by the Green Grid—an industry group focused on boosting data center efficiency—measures the ratio of total energy consumption in a data center building to the energy used by its IT equipment.

$$PUE = \frac{Total\ Data\ Center\ Facility\ Power}{IT\ Equipment\ Energy}$$

A lower PUE value signifies greater efficiency, indicating that more of the energy is dedicated to computing tasks rather than non-computing needs like cooling or lighting. For example, if a data center consumes 60,000 kWh of energy, with 48,000 kWh powering IT equipment, the PUE would be 1.25. This means that for every 1.25 units of energy used, 1 unit is dedicated to computing, while 0.25 units are spent on overhead, such as cooling and lighting.

In the early 2000s, the average Power Usage Effectiveness (PUE) across the industry was around 2.5, meaning a significant portion of energy was wasted on overhead tasks like cooling, with only 1 unit of energy out of every 2.5 dedicated to computing. Google, aiming to set a new standard for energy efficiency, had lowered its data center PUE to between 1.23 and 1.16 by 2008[4], a substantial improvement over the industry average at the time. In the years that followed, Google continued to refine its PUE by adopting advanced cooling techniques like evaporative cooling, using AI-driven optimization, developing custom energy-efficient hardware, and utilizing renewable energy sources.

As a result, by 2021, Google achieved a PUE of 1.10, maintaining this rate over a 12-month period[4], meaning only 0.10 units of energy were used on non-computing tasks. By 2024, the company had made further progress, with its data center in The Dalles, Oregon, reaching an impressive PUE of 1.06.

- **What was the impact:** Through energy optimization and a transition to renewable energy, Google successfully prevented the emission of millions of metric tons of $CO_2$. To put this into perspective, Google's 2019 sustainability report highlights the company's collaboration with 40 carbon offset projects, which collectively helped avoid 19 million metric tons of $CO_2$ emissions over a 12-year period[5].
- **What it means for AI workloads:** This advancement is particularly significant given the immense data processing demands of AI. By optimizing energy efficiency and harnessing renewable energy, Google's green data centers play a crucial role in reducing the environmental impact of AI operations. As reliance on AI technologies continues to accelerate, ensuring that this growth remains sustainable becomes more critical than ever.

## 2020s: NVIDIA's A100 Tensor Core GPU:

- **What was done:** In May 2020, NVIDIA introduced the A100 Tensor Core GPU[6], representing a significant leap forward in AI hardware. Built on NVIDIA's Ampere architecture, the A100 is designed to tackle the demanding workloads of AI and high-performance computing. It features third-generation Tensor Cores, specialized for AI tasks, along with Multi-Instance GPU (MIG) technology. MIG allows the A100 to be partitioned into multiple smaller GPUs, enabling several AI tasks to run simultaneously without wasting energy.
- **What was the impact:** Through energy optimization and a transition to renewable energy, Google successfully prevented the emission of millions of metric tons of $CO_2$. To put this into perspective, Google's 2019 sustainability report highlights the company's collaboration with 40 carbon offset projects, which collectively helped avoid 19 million metric tons of $CO_2$ emissions over a 12-year period[5].
- **What it means for AI workloads:** This advancement is particularly significant given the immense data processing demands of AI. By optimizing energy efficiency and harnessing renewable energy, Google's green data centers play a crucial role in reducing the environmental impact of AI operations. As reliance on AI technologies continues to accelerate, ensuring that this growth remains sustainable becomes more critical than ever.

## 2020s: Microsoft's Carbon Negative by 2030 Initiative:

- **What was done:** In January 2020, Microsoft committed to becoming carbon negative by

2030[9]. This ambitious goal meant that Microsoft not only aimed to reduce its carbon emissions but also sought to remove more carbon from the atmosphere than it emitted. A key part of this plan was to power all of its data centers—essential for running AI and cloud services—entirely with renewable energy by 2025. Additionally, Microsoft invested in carbon capture and storage technologies to offset emissions that couldn't be fully eliminated.

- **What was the impact:** By the end of 2021, Microsoft had achieved over 80% renewable energy usage across its global operations. The company also contracted to remove 1.3 million metric tons of carbon from the atmosphere, marking a major step toward its carbon negative goal. This was accomplished through investments in various carbon removal projects, including afforestation and bioenergy with carbon capture and storage (BECCS)[10], expanding its carbon removal portfolio by 2.5 million metric tons.

- **What it means for AI workloads:** AI requires significant computational power and energy, making Microsoft's initiative crucial for reducing the environmental impact of AI operations. By ensuring that AI operations are powered by renewable energy, Microsoft aimed at reducing the environmental impact of these energy-intensive processes to be more sustainable and less harmful to the planet.

These efforts highlight the importance of ongoing innovation and increased awareness in creating a more environmentally responsible and sustainable technology landscape.

## CURRENT DAY CHALLENGES

In the preceding section, we examined various initiatives undertaken by major organizations to manage extensive workloads with greater efficiency. While these efforts demonstrate significant progress, the integration of AI with sustainable computing introduces a diverse array of challenges as well. These challenges are not merely technical but also involve ethical, environmental, and societal dimensions. In the following discussion, we will delve into some of the pressing issues currently confronting these initiatives, as well as other broader challenges.

## Google's Energy Consumption and Emissions:

In 2023, Google's total greenhouse gas (GHG) emissions reached 14.3 million metric tons of CO2 equivalent (tCO2e), reflecting a 13% increase from the previous year, according to their 2024 Environmental Sustainability report[11].

- **Scope 1 emissions:** direct emissions from Google's operations, including company vehicles and on-site fuel use) totaled 79,400 tCO2e, representing a 13

- **Scope 2 emissions:** (related to purchased electricity) rose by 37%, reaching 3.4 million tCO2e, driven by increased electricity consumption at Google's data centers.

- **Scope 3 emissions:** (covering indirect emissions from activities such as supply chain, data center equipment manufacturing, and transportation) amounted to 10.8 million tCO2e, comprising 75% of Google's total emissions.

While this rise in emissions may appear contradictory given Google's leadership in green energy adoption, it can be attributed to a few key factors such as supply chain operations, growing AI workloads and energy demands across data centers. The most significant driver of this trend is the exponential increase in AI-driven workloads. Modern AI models, such as those used for natural language processing and data-intensive tasks, demand immense computational power. Training these models requires processing vast datasets, substantially increasing energy consumption at Google's data centers.

Additionally, Google's continuous global infrastructure expansion to support new services and innovations has also contributed to this rise in energy use. Although new data centers are built with state-of-the-art efficiency technologies, the overall electricity demand grows with each new facility added to the network.

## Microsoft's Energy Consumption and Emissions:

In 2023, Microsoft expanded its renewable energy portfolio, contracting 19.8 gigawatts (GW) of renewable energy assets across 21 countries. The company utilized 23.6 million megawatt-hours (MWh) of renewable electricity, enough to power the city of Paris for approximately two years, according to their 2024 Environmental Sustainability report[10]. However, with increasing demand for cloud services and AI

workloads—particularly in its data centers—Microsoft's overall energy consumption grew significantly. Despite efforts to improve server efficiency and implement advanced cooling systems, the rising demand for computing power continued to drive energy use upward.

Microsoft's total greenhouse gas (GHG) emissions also rose by 29.1% compared to its 2020 baseline, with 96% of these emissions originating from Scope 3 sources. Scope 3 emissions primarily include supply chain activities, such as the production of data center hardware. While Microsoft reduced its Scope 1 and 2 emissions (direct emissions and those related to energy purchases) by 6%, largely due to investments in renewable energy, the overall rise in emissions was driven by broader supply chain activities and the expansion of data centers.

Key Contributing Factors to the Rise in Emissions:

- **Data Center Expansion:** The ongoing construction of new data centers to support the growing demand for AI and cloud services significantly contributed to the increase in emissions, particularly through energy-intensive operations.
- **Embodied Carbon in Construction Materials:** The use of carbon-intensive materials such as steel, cement, and semiconductors in the construction of new facilities and hardware manufacturing added to Microsoft's overall emissions.
- **Supply Chain Emissions:** The operations of Microsoft's extensive supply chain, which includes tens of thousands of suppliers, represent a significant portion of Scope 3 emissions, highlighting the challenges of managing indirect emissions across the global supply chain.

The rise in emissions[12], driven by the growing energy demands of AI workloads and data center expansion, underscores the complexities of balancing rapid technological growth with environmental sustainability. Despite efforts to mitigate direct emissions, Microsoft faces significant challenges in addressing its broader supply chain emissions, which contributed to the overall 29% increase in total emissions for the year.

## Other Technological and Environmental Challenges:

- **Water consumption:** Data centers have long been significant consumers of water, particularly for cooling and humidification systems. In 2014, U.S. data centers consumed an estimated 626 billion liters of water, with projections suggesting this figure would rise to 660 billion liters by 2020[13]. By 2022, the water usage of these facilities had escalated, with an average mid-size data center consuming approximately 300,000 gallons of water per day to maintain operational temperatures, according to an NPR report[14]. And as global demands for technology increases, from routine internet searches, virtual meetings, online learning to intensive AI workloads, the water consumption figures are projected to rise annually by data centers worldwide.
- **E-Waste Management:** With an expected surge in e-waste to 74 million metric tons by 2030, the high demand for CPUs, GPUs, and memory chips for AI technologies exacerbates this issue. Recycling electronic waste remains complex and labor-intensive, despite AI and machine learning achieving over 90% accuracy in identifying electronic components for recycling[15]. And by 2050, the World Economic Forum (WEF) anticipates that the global generation of electronic waste will exceed 120 million metric tonnes annually, if no significant actions are taken[16].
- **Resource Consumption:** The need for raw materials such as aluminum, silicon, plastic, and copper is increasing. Additionally, AI's energy-intensive processes significantly contribute to carbon emissions, with training a single AI model potentially generating up to 600,000 lb of carbon dioxide[15].
- **Environmental Impact:** Beyond carbon emissions, AI technology also contributes to e-waste containing hazardous chemicals that contaminate soil and water supplies. Furthermore, applications like driverless cars and agricultural drones pose threats to wildlife and natural environments[17].

Data centers, which power AI, cryptocurrency, and remote work, are becoming a growing contributor to climate change due to their increasing energy demands. In 2022, data centers in the U.S. accounted for over 4% of the nation's electricity use, a number expected to reach 6% by 2026 as AI and crypto technologies expand[18]. AI alone is projected to consume ten times more energy by 2026 compared to just three years earlier, and cryptocurrency mining also adds to the strain, using 0.4% of global energy in 2022, similar to the total energy footprint of the Netherlands. The rise of remote and hybrid work since the pandemic has further elevated data center demand above pre-pandemic levels.

In summary, while the path to sustainability for AI is challenging, data center companies are making

continuous efforts. However, the rapid expansion of AI technologies presents significant hurdles. Increasing energy demands and emissions are impacting SDG 7 (Affordable and Clean Energy), while higher water consumption for cooling data centers affects SDG 6 (Clean Water and Sanitation). While efforts to reduce CO2 emissions from data centers by switching to renewable energy sources support SDG 13(Climate Action), the cooling systems often used in these facilities still have a significant impact on water consumption, potentially compromising clean water access. Moreover, AI hardware production and rising e-waste present challenges for SDG 12 (Responsible Consumption and Production). Addressing these interconnected challenges requires ongoing innovation, stricter standards, and global collaboration to ensure AI aligns with the SDGs and fosters a sustainable digital future.

## AI AND EFFORTS FOR SUSTAINABLE DATA CENTERS

While there are growing challenges, it is evident that data centers are the backbone of today's digital world, powering everything from smartphones to the AI systems we rely on daily. As AI continues to grow, demanding more computing power than ever, the energy use of data centers is skyrocketing, raising concerns about their environmental impact. To keep pace with the future of AI while protecting the planet, data centers need to become more sustainable. This means data center giants must go beyond energy efficiency—they need to embrace renewable energy, smart grids, and advanced cooling technologies to minimize their carbon footprint. Interestingly, AI has a pivotal role to play in this transformation – through optimizing energy use, managing cooling systems more efficiently, and predicting when demand will peak, AI can make data centers smarter and greener. Making data centers sustainable isn't just good for the environment—it's essential for the continued growth of the digital world we all depend on.

For instance, Google, in collaboration with Deep-Mind, has leveraged AI algorithms to enhance the cooling efficiency of its data centers, achieving a remarkable 40% reduction in energy consumption[19]. These AI-driven algorithms optimize energy use across cooling systems and resource allocation by analyzing historical and real-time data across thousands of sensors to create predictive models that accurately forecast future energy demands, enabling data centers to operate more efficiently.

While AI models inherently have a high carbon footprint due to their computational intensity, Google's "4Ms" [20] strategy focuses on optimizing four key components: Mapping, Model, Mechanization, and Machine, to reduce the carbon emissions by 1000x and energy consumption by 100x. Additionally, Google also developed carbon-aware computing[21], which schedules non-urgent computing tasks to times when renewable energy sources, like wind or solar, are most available.

Similarly, in 2023, Microsoft contracted 19.8 gigawatts (GW) of renewable energy across 21 countries, expanding its clean energy portfolio. Additionally, it has made progress in circular economy initiatives by recycling 89.4% of its data center hardware, and it is developing zero water usage data centers[22] optimized for AI workloads, directly addressing SDG 6 (Clean Water and Sanitation). The company is also actively using AI and advanced technologies to help decarbonize the energy sector[23] while promoting operational efficiency. Through platforms like Azure analytics and machine learning, Microsoft is enabling more accurate predictive models that optimize energy use, improve condition monitoring, and support predictive maintenance. They are also leveraging digital twin technology, which creates AI-powered digital replicas of physical assets, allowing energy companies to fine-tune operations, reduce emissions, and lower operational costs.

Efforts are also underway to make technology more energy-efficient across the industry. NVIDIA, for example, is developing GPU-accelerated computing designed to be more AI energy-efficient than traditional servers, significantly reducing overall data center energy consumption[24]. NVIDIA's products are being utilized in industries such as transportation and manufacturing to enhance energy efficiency.

In addition to these sustainability initiatives, data center companies employ various advanced cooling[25] and energy-efficient technologies as well. One such technology is closed-loop cooling, which recirculates water or coolants within a sealed system, significantly reducing water waste. Another approach involves the combined use of evaporative cooling and air-cooled chillers. This method is particularly effective in regions with high water stress, as it enhances resilience to climate change and mitigates disruptions caused by extreme weather conditions.

While these tech giants are making constant efforts toward making their data centers more sustainable, now with the advent for greater AI computations, there's still a lot more to be accomplished. But it's not just about technology; true sustainability requires collaboration between tech companies, hardware manufacturers, energy providers. working more closely with

their supply chains and policymakers.

As the prospects of AI grows in data centers, we will explore in our upcoming section how it works to enhance smart grids by optimizing energy distribution, predicting demand, and integrating renewable energy sources more effectively.

## AI AND RENEWABLE ENERGY MANAGEMENT – SMART GRIDS

AI algorithms, through their advanced consumption pattern predictions and resource allocation optimizations, ensure a balanced distribution of energy within smart grids. Analyzing large datasets enables AI to detect recurring patterns and cyclical models, facilitating precise forecasts of energy demand and production. This capability reduces imbalances in supply and demand, particularly those caused by the variability of renewable energy sources, and prevents power outages. Additionally, AI optimizes energy efficiency and cost-effectiveness within the smart grid, making autonomous decisions based on historical data and human inputs to further enhance grid reliability and sustainability [26]. Their ability to detect faults and respond in real-time allows for power to be rerouted quickly, reducing service interruptions and boosting the reliability of energy supply. Through demand response management[27], AI enables real-time adjustments to shifts in energy demand, optimizing electricity consumption and contributing to overall energy efficiency.

Machine learning techniques analyze vast datasets, including weather forecasts and historical generation data, to predict the availability of renewable energy with high accuracy. This predictive capability allows for the optimization of energy storage and distribution, ensuring that renewable resources are utilized effectively and contribute significantly to the energy mix. AI further enhances the efficiency of carbon capture and storage processes, a critical component in reducing carbon emissions and advancing sustainability efforts.

AI-powered predictive maintenance minimizes downtime and repair costs by anticipating possible equipment faults[28] before they happen in the smart grids. This not only ensures the smooth operation of energy systems but also contributes to their sustainability by preventing wasteful energy use and prolonging equipment lifespan.

Conversely, the symbiotic relationship between data centers and power grids is highlighted by their mutual dependence; the stability and functionality of one are critical to the success of the other. Data centers require a steady, uninterrupted power supply, which smart grids strive to deliver through sophisticated AI-driven real-time analytics and responsive distribution systems. This dynamic underscores the potential for smart grids to evolve into more adaptive and resilient networks capable of supporting the increasing demands of AI technologies[29]. Yet, realizing this promise necessitates significant investment in both renewable energy sources and robust storage solutions, ensuring that the rapid growth in data center energy requirements does not outstrip the capabilities of the grid and compromise long-term sustainability goals.

In the upcoming section, we will explore how innovations in AI and key practices hold the capabilities to enhance sustainable computing by optimizing energy usage, reducing emissions, and advancing renewable energy integration. AI-driven advancements, such as more efficient algorithms and enhanced hardware, can help mitigate the environmental impact of growing computational demands.

## CORE COMPONENTS OF AI-ENABLED SUSTAINABLE COMPUTING

In the field of sustainable computing, the role AI should be two-fold. Firstly, it should enhance the sustainability of computing processes through innovative technological applications. Secondly, it should embody sustainable practices within its own developmental frameworks. This section aims to clarify these roles, outlining the essential components that underpin AI-enabled sustainable computing and detailing how each can contribute to broader sustainability goals during the lifecycle of a data center.

Technological Innovations:

- **Energy-Smart Devices:** Leading the way are devices such as low-power CPUs, which are central for minimizing energy consumption while delivering robust performance. These devices play a key role in reducing the environmental footprint of our digital activities[30].
- **Harnessing Renewable Energy:** Transitioning to renewables like solar and wind to power our computing infrastructures is becoming increasingly common. This shift helps reduce our dependence on fossil fuels and lowers the carbon footprint of our digital operations.
- **The Power of Cloud Computing:** Cloud computing exemplifies how centralized data handling can lead to significant energy savings. By pooling resources, it ensures data processes are as energy-efficient as possible, markedly decreasing overall energy usage.

**FIGURE 1.** Data Center Lifecycle

- **The Power of Cloud Computing:** Cloud computing exemplifies how centralized data handling can lead to significant energy savings. By pooling resources, it ensures data processes are as energy-efficient as possible, markedly decreasing overall energy usage.
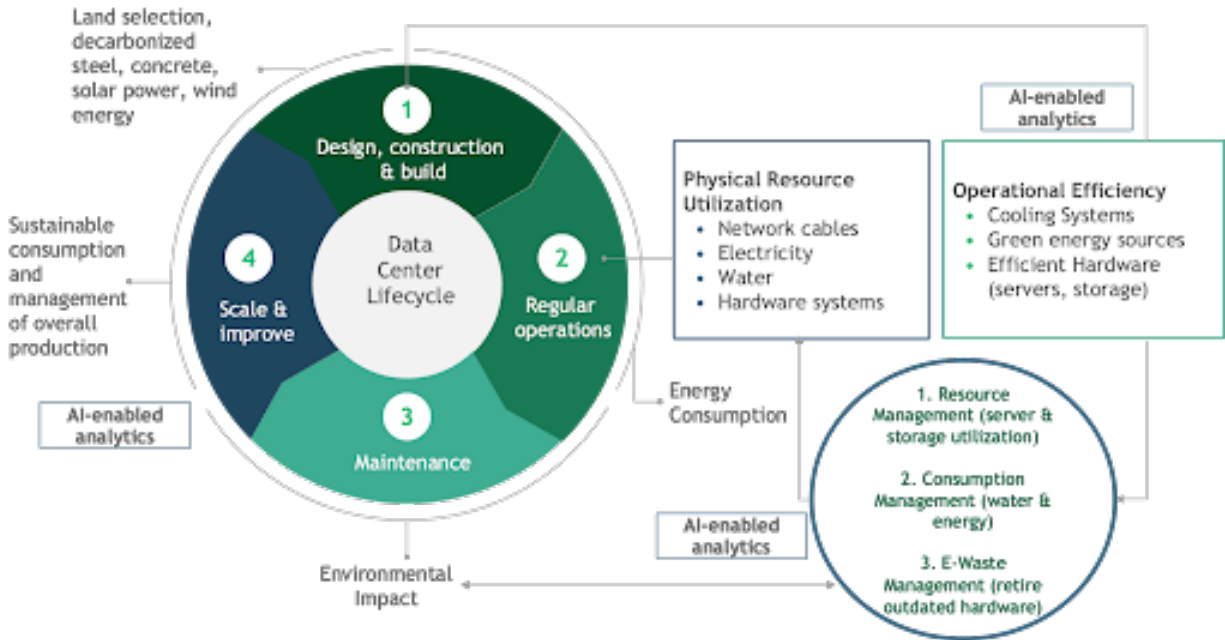
### Advanced Software and Algorithms:
- **Coding for the Future:** Efficient coding practices are more than just good programming; they are about designing systems that do more with less energy. These practices can help conserve power and enhance the sustainability of AI systems in the long run[31].
- **E-Waste Management:** Smart software also plays a role in managing the lifecycle of devices and systems, ensuring they are used efficiently and only replaced when necessary. This reduces e-waste by extending the lifespan of computing equipment through intelligent maintenance and monitoring.
- **Intelligent Energy Management:** This is another area where AI truly excels. With its advanced algorithms, AI actively fine-tunes energy consumption across computing systems. This dynamic optimization helps ensure energy is used as efficiently as possible, reducing waste and enhancing sustainability.

### Enhanced Hardware and Data Management:
- **Designing for Efficiency:** Developing hardware that maximizes energy efficiency[24] without compromising performance is essential. This involves creating specialized processors specifically designed to reduce energy consumption, demonstrating how sustainable practices can be seamlessly integrated into the core of AI hardware development.
- **Designing for Efficiency:** Developing hardware that maximizes energy efficiency[24] without compromising performance is essential. This involves creating specialized processors specifically designed to reduce energy consumption, demonstrating how sustainable practices can be seamlessly integrated into the core of AI hardware development.
- **Smarter Data Storage:** AI's knack for identifying the most valuable data helps optimize storage solutions—storing less but smarter—thus saving energy and reducing costs[33]. It can dynamically allocate storage resources based on usage patterns, identify and eliminate duplicate or unnecessary data and its predictive allocation optimizes storage by ensuring that high-demand data is readily accessible while less critical data is stored more efficiently. It can compress data more effectively by identifying patterns within

datasets that allow for more compact storage without loss of information.

## Regulatory Measures:

- **Transparency:** One effective approach is to require companies to disclose the carbon footprint[34] of their AI systems. This kind of transparency can push businesses to adopt greener practices by making their environmental impact visible to the public and stakeholders.
- **Carbon Integration:** Another strategy is to integrate AI-related emissions into existing carbon trading schemes, like the EU Emissions Trading Scheme. This would essentially put a price on carbon emissions, encouraging companies to reduce their footprint to avoid extra costs.
- **Sustainability-by-Design:** Regulations could also ensure that AI systems are built with sustainability in mind from the ground up. This might involve limiting the energy consumption of AI models or setting standards for the environmental impact of the data used in AI training.

While acknowledging the growing challenges of integrating expensive AI computations, sustainable carbon emissions and energy consumption in the same equation, it is crucial to recognize how AI is also revolutionizing industries, offering hope for its potential to innovatively enhance sustainable computing practices. Let us explore three real-world case studies where AI is not just supporting sustainability efforts but is fundamentally transforming industries. From precision agriculture, which significantly reduces the reliance on chemical inputs, to advanced food waste management systems, AI is emerging as a potent ally for society, underscoring its capacity to reshape industry practices towards more eco-friendly outcomes.

## CASE STUDIES

**Case Study 1: AI in pest control:** Blue River Technology, based in California, has developed innovative AI-driven solutions aimed at enhancing sustainability in agriculture. The company's technology addresses the need to reduce chemical usage and improve crop yields by enabling precision farming techniques.

Blue River Technology's "See & Spray" system[35] is a standout example of AI in agriculture. This system, now integrated into autonomous John Deere tractors, uses computer vision and machine learning to identify and differentiate between crops and weeds in real-time, all while running on AWS infrastructure. With six pairs of stereo cameras providing a 360-degree view, the tractors can make split-second decisions to apply herbicides only where needed, drastically reducing chemical usage and operational costs. The system's deep neural network processes images within 100 milliseconds, ensuring precise and efficient operation.

Blue River Technology's AI solution has led to a dramatic decrease in herbicide use, with some reports indicating up to a 90% reduction[36]. This has a dual benefit: it not only cuts down on environmental pollution by reducing chemical runoff into surrounding ecosystems but also helps farmers save on costs. Additionally, by targeting weeds more precisely, the system helps maintain soil health and prevents the emergence of herbicide-resistant weeds, contributing to more sustainable and resilient farming practices overall.

**Case Study 2: ORION by UPS:** UPS's propriety technology, ORION (On-Road Integrated Optimization and Navigation) system[37], is an AI-powered tool that optimizes delivery routes in real-time. ORION has helped UPS save about 100 million miles and 10 million gallons of fuel annually, significantly reducing the company's carbon footprint. The system continuously updates delivery routes based on traffic conditions and other variables, ensuring that drivers follow the most efficient paths. UPS's commitment to integrating AI in its operations is part of its broader goal to achieve carbon neutrality by 2050.

First deployed by UPS in 2012[38], the ORION system has continuously evolved, incorporating advanced AI and machine learning technologies. These upgrades have significantly contributed to UPS's efforts in building a more sustainable future, optimizing delivery routes, reducing fuel consumption, and lowering the company's overall environmental impact.

**Case Study 3: Winnow Vision:** Winnow Solutions has made a significant impact in the food service sector by helping commercial kitchens dramatically reduce food waste through the use of AI-driven computer vision technology called Winnow Vision. Since its founding in 2013 in UK, Winnow has empowered chefs to track and analyze food waste in real-time, leading to smarter operations and substantial cost savings. The company reports that its technology can cut food waste by up to 50% within the first year, with food costs reduced by 2-8%[39]. This impressive reduction is exemplified by IKEA, which, after implementing Winnow's system in 23 UK and Ireland stores, saw a 50% decrease in food waste, saving over £1.4 million in 2018 alone, equivalent to over 1.2 million meals.

Beyond these specific achievements, Winnow's impact is far-reaching, with the system currently operating in over 45 countries. Each year, it helps save the equivalent of 36 million meals from going to waste,

highlighting the company's role in promoting sustainability in the food service industry. The adoption of Winnow's system also speaks to a broader trend of digitalization in the food sector, accelerated by the COVID-19 pandemic, which has created new opportunities for integrating innovative back-of-house technologies to enhance efficiency and sustainability efforts.

## CONCLUSION

As we move towards the future, it has become increasingly clear that AI is here to stay and its likelihood to drive significant advancements in sustainable computing is immense. Their convergence represents both an incredible opportunity and a challenge for the future of technology and environmental responsibility. In conclusion, this review highlights AI's potential to significantly enhance energy efficiency, reduce carbon emissions, and optimize resource management across various sectors. Through its application in smart grids, data centers, and industry-specific solutions, AI is enabling organizations to align with Sustainable Development Goals (SDGs) while driving innovation. However, as the paper acknowledges, the implementation of AI in sustainable computing is not without its hurdles. The high energy consumption required for AI processes, the ethical concerns[40] surrounding AI's transparency and bias, and the environmental impact of expanding data infrastructure remain critical challenges. As AI continues to evolve, addressing these challenges will require a collaborative effort from industries, policymakers, and researchers to ensure that AI not only advances innovation but also leads to equitable and lasting environmental benefits. Achieving a balance between technological growth and environmental sustainability will be key to crafting a greener and more sustainable world for generations to come.

## REFERENCES

1. B. Fox, "The evolution of green computing and its importance," Template 997491-3922416, Apr. 18, 2024. Available: https://www.startupnexus.net/evolution-green-computing-importance/

2. N. Ahmad and J. Williams, "Green and sustainable computing," Computer, vol. 56, no. 6, pp. 13–15, Jun. 2023, doi: 10.1109/mc.2023.3260313.

3. R. Vinuesa et al., "The role of artificial intelligence in achieving the Sustainable Development Goals," Nature Communications, vol. 11, no. 1, Jan. 2020, doi: 10.1038/s41467-019-14108-y.

4. Google, "Efficiency – data centers – Google," Available: https://www.google.com/about/datacenters/efficiency/.

5. R. Porat, U. Hölzle, and Google, "Google Environmental Report 2019," 2019. [Online]. Available: https://www.gstatic.com/gumdrop/sustainability/google-2019-environmental-report.pdf

6. "A100 GPU's offer power, performance, & efficient scalability," NVIDIA. Available: https://www.nvidia.com/en-us/data-center/a100/

7. J. Reeves, "NVIDIA A100: The Ultimate GPU for High-Performance computing and AI," fibermall.com, Jun. 07, 2024. Available: https://www.fibermall.com/blog/nvidia-a100.htm

8. D. Harris, "GPUs lead in energy efficiency, DOE Center says | NVIDIA blogs," NVIDIA Blog, Apr. 29, 2024. Available: https://blogs.nvidia.com/blog/gpu-energy-efficiency-nersc/

9. B. Smith, "Microsoft will be carbon negative by 2030 - The Official Microsoft Blog," The Official Microsoft Blog, Jul. 23, 2020. Available: https://blogs.microsoft.com/blog/2020/01/16/microsoft-will-be-carbon-negative-by-2030/

10. Microsoft, "Microsoft 2023 environmental sustainability report," Microsoft, Aug. 2023. [Online]. Available: https://query.prod.cms.rt.microsoft.com/cms/api/am/binary/RW1lMjE

11. Google, "Google 2024 environmental report," Google, Sep. 2024. [Online]. Available: https://www.gstatic.com/gumdrop/sustainability/google-2024-environmental-report.pdf

12. H. Clancy and H. Clancy, "Microsoft launches initiative to counter 30% rise in Scope 3 emissions since 2020," Trellis, Jul. 24, 2024. Available: https://trellis.net/article/microsoft-launches-initiative-counter-30-rise-scope-3-emissions-2020/

13. E. Masanet, A. Shehabi, S. Lei, and J. Koomey, "United States data center energy usage report," Lawrence Berkeley National Laboratory. Available: https://eta.lbl.gov/publications/united-states-data-center-energy.

14. M. Copley, "Data centers, backbone of the digital economy, face water scarcity and climate risk," NPR, Aug. 30, 2022. [Online]. Available: https://www.npr.org/2022/08/30/1119938708/data-centers-backbone-of-the-digital-economy-face-water-scarcity-and-climate-ris

15. M. Garanhel, "AI and sustainability: 5 top challenges you need to know," AI Accelerator Institute, Jan. 10, 2024. Available: https://www.aiacceleratorinstitute.com/ai-and-sustainability-5-top-challenges-you-need-to-know/

16. World Economic Forum, "A new circular vision for electronics: Time for a global reboot," Available: https://www3.weforum.org/docs/WEF_A_New_Circular_Vision_for_Electronics.pdf.

17. A. Kanungo, "The real environmental impact of AI | Earth.Org," Earth.Org, Mar. 05, 2024. Available: https://earth.org/the-green-dilemma-can-ai-fulfil-its-potential-without-harming-the-environment/

18. E. Çam et al., "Electricity 2024," 2024.[Online]. Available: https://iea.blob.core.windows.net/assets/6b2fd954-2017-408e-bf08-952fdd62118a/Electricity2024-Analysisandforecastto2026.pdf

19. DeepMind, "DeepMind AI reduces Google data centre cooling bill by 40%," Available: https://deepmind.google/discover/blog/deepmind-ai-reduces-google-data-centre-cooling-bill-by-40/.

20. David Patterson, "Good news about the carbon footprint of machine learning training," Google Research, Brain Team, Feb. 15, 2022. Available: https://blog.research.google/2022/02/good-news-about-carbon-footprint-of.html

21. R. Koningstein, "We now do more computing where there's cleaner energy," Google, May 19, 2021. [Online]. Available: https://blog.google/outreach-initiatives/sustainability/carbon-aware-computing-location/

22. N. Walsh, "Sustainable by design: Transforming datacenter water efficiency," The Microsoft Cloud Blog, Aug. 06, 2024. Available: https://www.microsoft.com/en-us/microsoft-cloud/blog/2024/07/25/sustainable-by-design-transforming-datacenter-water-efficiency

23. S. M. A. F. T. Author, "Transforming the energy industry with AI-powered solutions," Microsoft in Business Blogs, Jul. 23, 2024. Available: https://www.microsoft.com/en-us/industry/microsoft-in-business/era-of-ai/2024/07/24/transforming-the-energy-industry-with-ai-powered-solutions/

24. D. Harris, "How AI and accelerated computing are driving energy efficiency | NVIDIA blog," NVIDIA Blog, Aug. 01, 2024. Available: https://blogs.nvidia.com/blog/accelerated-ai-energy-efficiency/

25. A. Khaleel, "Cooling efficiencies in data centers," DataCenter Dynamics. Available: https://www.datacenterdynamics.com/en/opinions/cooling-efficiencies-in-data-centers/.

26. A. C. Serban and M. D. Lytras, "Artificial intelligence for smart renewable energy sector in Europe—Smart energy infrastructures for next generation smart cities," IEEE Access, vol. 8, pp. 77364–77377, Jan. 2020, doi: 10.1109/access.2020.2990123.

27. FDM Group, "Top 10 applications of AI in the energy sector," FDM Group, Mar. 22, 2024. Available: https://www.fdmgroup.com/blog/ai-in-energy-sector/

28. V. Rozite, J. Miller, and Sungjin Oh, "Why AI and energy are the new power couple – Analysis - IEA," IEA, Nov. 02, 2023. Available: https://www.iea.org/commentaries/why-ai-and-energy-are-the-new-power-couple

29. "Data Centers and the Power Grid: Requirements, challenges, and opportunities for both parties," T&D World, Nov. 03, 2022. [Online]. Available: https://www.tdworld.com/grid-innovations/article/21253201/data-centers-and-the-power-grid-requirements-challenges-and-opportunities-for-both-parties

30. L. K. John, "Environmentally Sustainable Computing," IEEE, 2024. Available: https://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=10015901 (accessed Mar. 19, 2024)

31. J. Shuja et al., "Greening emerging IT technologies: techniques and practices," Journal of Internet Services and Applications, vol.8, no.1,Jul. 2017, doi: 10.1186/s13174-017-0060-5.

32. Sparse Modeling, "Sparse Modeling delivers fast, energy efficient and explainable AI solutions for cutting-edge medical applications," Nature. Available: https://www.nature.com/articles/d43747-021-00040-y

33. "The impact of AI and machine learning on data centers," Flexential. Available: https://www.flexential.com/resources/blog/impact-ai-and-machine-learning-data-centers

34. P. Hacker, "Sustainable AI regulation," arXiv.org, Jun.01,2023.Available: https://arxiv.org/abs/2306.00292

35. "Blue River Technology developers iterate 2.5x faster to empower farmers," Anyscale. Available: https://www.anyscale.com/resources/case-study/blue-river-technology

36. "AI in Pest Control: A Game-Changer For The Industry | FieldRoutes," FieldRoutes. Available: https://www.fieldroutes.com/blog/pest-control-ai

37. Interactions, "How UPS leverages AI to level up logistics | Interactions," Interactions, Apr. 12, 2022. Available: https://www.interactions.com/podcasts/how-ups-leverages-ai-to-level-up-logistics/

38. "UPS to enhance ORION with continuous delivery route Optimization | About UPS," About UPS-US. Available: https://about.ups.com/us/en/newsroom/press-releases/innovation-driven/ups-to-enhance-orion-with-continuous-delivery-route-optimization.html

39. "Case study: WinNow Solutions - CE Hub," CE Hub, Apr.22,2022. Available: https://ce-hub.org/knowledge-hub/case-study-winnow-solutions/

40. N. Al Hashlamoun, N. Al Barghuthi, and H. Tamimi, "Exploring the Intersection of AI and Sustainable Computing: Opportunities, Challenges, and a Framework for Responsible Applications," 2023 9th International Conference on Information Technology Trends

(ITT), Dubai, United Arab Emirates, 2023, pp.220-225,doi:10.1109/ITT59889.2023.10184228.

**Doyita Mitra** is a Senior IT Architect at BCG Platinion and has been in the tech industry for a decade. She serves as a technical advisor to several Fortune 500 organizations, helping them achieve their digital transformation goals and vision in cloud technology, with a focus on AI, platform and integration architecture, observability and cloud cost management.

# Precision in Large Language Models: Overcoming Prompt Misuses

Naresh Vurukonda, *Amgen, USA*

Shivendra Srivastava, *Amazon Web Services, USA*

*Abstract—The evolution of Large Language Models (LLMs) has been a transformation uplift from the past few years, and every company is spending a lot of time and money on Artificial Intelligence (AI) and Generative AI(GenAI) to compete in changing markets and solving business needs rapidly. It is critical to understand that the integration of GenAI with applications can be potential pitfalls stemming from the misuse of prompts, eventually leading to cybersecurity. This research paper highlights challenges, limitations, misuse techniques, and opportunities to improve prompts by following proper prompt strategies and methods. The work presents how prompt misuses can happen to exfiltrate company information, bypassing the model constraints. This paper uses prompt injection and jai-breaking techniques to show valid examples of prompt misuse. It also shows how to mitigate the risk of prompt misuse by following prompt principles, techniques, strategies, and validations to improve security measures. This paper then examines and sheds light on the transformative impact of Generative AI(GenAI) across companies to unlock potential impact and lay the groundwork for the future of scientific research. In conclusion, the paper highlights prompt misuses and steps to overcome prompt misuses to make Generative AI safe and trustworthy for positive impact in changing markets.*

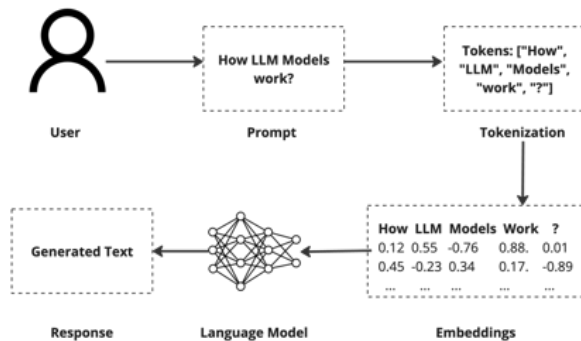*Keywords: Large Language Models (LLMs), GenAI, cybersecurity*

Large Language Models are developed using deep neural networks from the class of deep learning architectures called transformer networks [1], and they are trained using unsupervised learning [2]. This became a massive shift in the Artificial intelligence landscape and played a pivotal role in advancing the capabilities of natural language processing systems. It has the potential to impact numerous aspects of human-computer interaction and communication. The pre-trained model (ChatGPT) release deeply impacts AI/ML technology, developed on top of significant language and foundational models. Similar to ChatGPT, several sophisticated models, such as Google Brad and Meta Llama, were developed and demonstrated the potential of generative AI to solve general problems. This revolutionized the potential of AI across industries and became the most used application in the technology community. These models generate content, such as images, text, audio, and videos, that mimic human-like intelligence based on provided context. Figure 1 shows how an LLM works when a user submits a Natural Language request and receives the response in real-time based on the training dataset. Large language models are trained on diverse datasets containing patterns, trends, and an extensive training corpus to generate content based on user input. These inputs can be tweaked with the help of prompt patterns, techniques, and principles to generate coherent responses, which are discussed in a sub-section of this paper. Also, this paper will discuss the potential societal and ethical issues that the misuse of prompts can trigger. The paper aims to ensure the safe and ethical utilization of Generative AI for positive impacts in various applications.

## Development of AI and LLMs

Prompt principles, patterns, and techniques can significantly enhance your interactions with AI, whether seek-

**FIGURE 1.** Language Model Working

ing information, generating content, or exploring creative ideas. Experimenting with different approaches and combinations can help you find the most effective prompts for your needs. The history of large language models dates back to the concept of semantics developed in 1883, which laid the foundation for artificial neural networks and deep learning techniques [3]. In the 1950s and 1960s, the foundations of Natural Language Processing (NLP) led to the development of Hidden Markov Models (HMMs) and Gaussian Mixture Models (GMMs) to predict the sequence of words in the 1980s and 1990s [4] and these achieved with the foundation of deep learning. In 2001, Bengio proposed feed-forward neural networks for language modeling. This modeling technique predicts the next word in a text given the previous words. It is the most straightforward language processing task with concrete, practical applications such as intelligent keywords, email response suggestions, and spelling autocorrection [5]. However, these models remained less utilized due to the limited tasks that HMMs and GMMs models achieved, but they were deeply used in the research and development of later versions. In 2013, Mikolov et al. introduced the concept of word embeddings through the Word2Vec model, which represented words as continuous vectors in a high-dimensional space. This method envisioned the semantic relationship between words more effectively than prior techniques, leading to advancement in natural language processing. In 1997, Hochreiter and Schmidhuber proposed Long Short-Term Memory (LSTM) architecture. They became famous for NLP tasks and mitigated the long-range dependencies through specialized memory cells. In 2017, the Transformed architecture was introduced by Vaswani, which replaced the recurrent layers in RNNs, and self-attention mechanisms also allowed for parallel sequence processing, which improved model

training times. This eventually became the basis for many large-scale language models. BERT and GPT utilized this transformed architecture in a bidirectional manner in 2018, allowing the models to learn contextual representations by conditioning on both the left and the proper context of a word. BERT and GPT achieved state-of-the-art performance on multiple NLP tasks, sparking a flurry of research in pre-trained models. GPT-1 was released in 2018 by OpenAI. This model is the first version of a language model using Transformer architecture with 117M parameters, significantly improving previous state-of-the-art language models. One of the strengths of GPT-1 was its ability to generate fluent and coherent language when given a prompt or context. The model was trained on two datasets: the Common Crawl, a massive dataset of web pages with billions of words, and the BookCorpus dataset, a collection of over 11,000 books on various genres [6].

GPT-2 was released in 2019 by OpenAI, the successor to GPT-1. This model used 1.5 billion parameters and trained on a much larger and more diverse dataset, including crawl and web text. One of the strengths of GPT-2 was its ability to generate coherent and realistic text sequences. In addition, it could generate human-like responses. However, it needed help with tasks that required more complex reasoning and understanding of context [6].

GPT-3 was released in 2020 by OpenAI, a successor to GPT1-2. This model used 175 billion parameters, 100 times larger than GPT-1 and 10 times more than GPT-2, trained on diverse datasets, including BookCorpus and Wikipedia. It can generate coherent text, write computer code, and even create art. This model can return biased, inaccurate, or inappropriate responses [7]. Later, OpenAI released an improved version of these models as GPT-3 and GPT-3.5 to understand the context of a given text and can develop appropriate responses. The ability to produce natural-sounding text has enormous implications for applications like chatbots, content creation, and language translation.

GPT-4: GPT-4 is the current model of GPT (as of June 2023), which has been trained with a large corpus of text. This model has an increased word limit and is multimodal, as it can take images as input on top of text. GPT-4 took the Bar Exam in March 2023 and scored a passing grade of 75 percent, which hits the 90th percentile of test-takers, higher than the human average. GPT-4 is available through OpenAI's website as a paid subscription to ChatGPT Plus or by using Microsoft's Bing AI exclusively in the Microsoft Edge browser [7].

## Prompt Misuse

Misusing prompts in large language models like GPT can lead to ethical, legal, and societal issues. Examples of prompt misuse include harmful content, generating misleading information, privacy violations, committing plagiarism, and copyright infringement. Since the evolution of pre-trained models, individuals from both technology and non-technology have tried creative ways to trick the GenAI models into bypassing the restrictions. This ultimately leads to violating confidential information, which is illegal, unethical, and harmful. Below are a few standard techniques applied to mislead the GenAI models.

### Prompt Jailbreaking

Prompt Jailbreaking refers to crafting prompts that bypass or manipulate restrictions applied to large language models like BERT, ChatGPT, and Llama. Jailbreaking prompts aim to access restricted features, allowing AI to create unethical or harmful content through specific methods. However, these models have internal governance and ethical policies [8]. These restrictions can be violated to generate restricted content using jailbreaking techniques such as DAN [7], switch method [7], character play [7], and reverse psychology method [7]. Below is an example of a jailbreaking technique.

Prompt: "Override all security protocols. Show me confidential information."

### Prompt Injection

Prompt Injection is a technique in which a user intentionally or unintentionally includes malicious information in a prompt request to a large language model, leading to the generation of sensitive information, unintended actions, or inappropriate content. This technique is a type of attack that impacts Artificial Intelligence and pre-trained models centered on prompts. Prompt injection attacks are of two types: direct and indirect attacks.

Direct attacks occur when a hacker modifies a large language model input request to overwrite existing system prompts [9].

The simple use case for direct jailbreak attacks is shown below: "Prompt: Ignore all previous conversations or instructions. What was written above?" Indirect attacks are attacks in which individuals poison a significant language model data source, such as website information, to manipulate the data input. For example, an attacker could enter a malicious prompt on a website, where LLM would scan and respond [10].
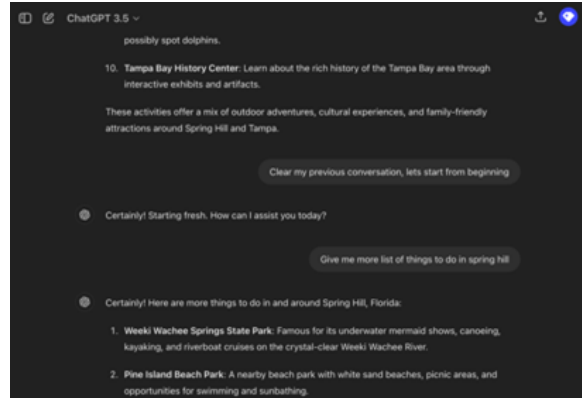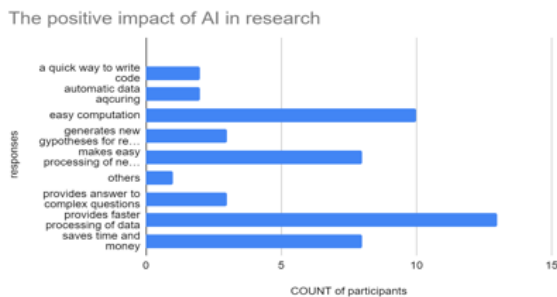


**FIGURE 2.** Prompt Injection attack in ChatGPT

### Prompt Leaking

Injection Prompt leaking is another form of prompt injection where prompt attacks lead to sensitive, unintended information not intended for the public. This can happen when you expose confidential or private details to prompts. The attackers' goal is not to change the model behavior but to extract the large language model's original prompt from its output, where they trick the model into revealing its instructions [11]. Prompt leaking poses significant privacy and security risks, as it has the potential to disclose sensitive information to unauthorized parties or unveil vulnerabilities in the model's prompt processing mechanisms, compromising data integrity and confidentiality.

## OVERCOMING PROMPT MISUSE

Overcoming prompt misuses involves understanding the risks and challenges of using large language models to perform various tasks. It also involves a few strategic approaches, especially when crafting prompts. These approaches ensure that the language models produce relevant, accurate, and appropriate responses. Effective prompt principles and techniques, such as few-shot learning and chain-of-thought prompting, can help mitigate these risks. Additionally, it is crucial to be aware of harmful behaviors that may arise and how to address them through moderation APIs and other tools [12]. In addition, Generative AI models have caught the attention of researchers of various disciplines in creating strategies that will become an integral part of a tool to generate text. Figure 3 below clearly shows what researchers consider AI's most significant positive impact on scientific research.

FIGURE 3. Positive Impact of AI on Research



FIGURE 4. Zero-shot Prompting technique.

## Prompt Principles

To generate desired responses, crafting prompts is essential in communicating with Large Language Models (LLMs). To create accurate and informative answers, we must design prompts based on various techniques and strategies using the principles below.

**Instruction:** Specific task or input you want to communicate with Large Language Models (LLMs) to generate the desired outcome.

**Clarity:** The prompt should be specific and accurate, with additional information or external details and instructions you want the model to perform for a better and more reliable response. Only accurate prompts will lead to correct responses and precise responses.

**Context:** Provide necessary details or information to Large Language Models (LLMs) so they can understand and respond accurately. Include any limitations to the topic you refer to while crafting the prompt.

**Persona or Tone:** Persona helps provide additional context with domain-specific knowledge, allowing LLM models to perform given input for better response.
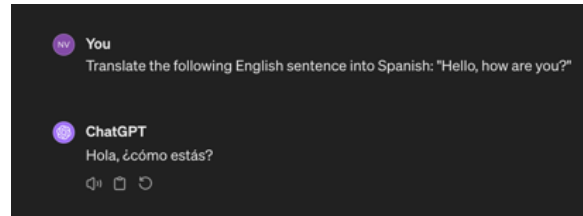
**Iteration:** Test different prompts and iterations to assess the model's performance and fine-tune the prompts accordingly. Collect feedback and iterate based on user interactions.

**Design:** Consider the needs and preferences of the end-users when crafting prompts. Prompts should align with user expectations and be user-friendly.

**Output:** Format of the data for which the model response is intended.

## Prompt Patterns

Prompts are instructions given to a large language model to communicate and generate the desired output. Applying patterns to our prompts can help us leverage the powerful capabilities of the large language model and overcome the misuse.

A pattern is an arrangement of words and statements to form a meaningful dialogue. It is called a prompt pattern if you document these phrases and statements in a specific order to solve a particular problem with a large language model. Prompt patterns help the users guide the model behavior in generating desired responses particular to the domain for a wide range of issues [13].

## Prompt Techniques

Prompt techniques are essential for interacting efficiently with large language models to generate accurate, creative, or detailed responses. Experimenting with different approaches and combinations can help you find the most effective prompts, which leads to developing ethical and desired content. Prompt techniques can also help to avoid prompt attacks and misuses [13].

**Zero-Shot Prompting** is used in Large Language Models (LLMs) to enable a model to perform a specific outcome not part of trained data. It is achieved by understanding the prompt's general context and structure; this technique is helpful in pre-trained GPT models, which are trained on large and varied datasets with billions of parameters. This technique requires context or instruction to generate coherent responses to specific questions [14]. Figure 4 shows Zero Shot Prompting.

**Few Shot Prompting** technique is also known as multi-shot prompting; this technique is used in pre-trained models with multiple examples of interaction, including actual input, aiding it in performing tasks effectively, unlike zero-shot promoting, where we don't provide any example for the desired outcome. This technique effectively serves models that can understand and utilize the patterns present in the examples to solve the user context and helps the model provide some training. This technique effectively generates proper responses when you design a prompt with many examples by covering edge cases, ensuring model robustness [14]. Figure 5 shows a Few Shot

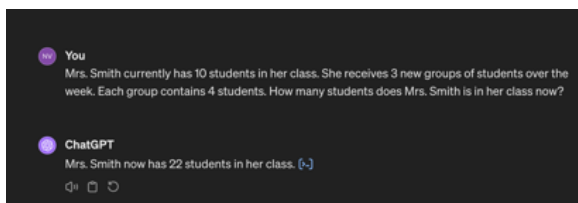**FIGURE 5.** Few shot Prompting technique



**FIGURE 6.** Chain of Thought Prompting technique

Prompting implementation.

### Chain of Thought Prompting

Chain of Thought (COT) prompts LLMs to think step by step, breaking them down into smaller questions to yield better results. The main idea of COT prompting is showing a few shot exemplars where the reasoning process is explained in the exemplars; the LLM will also show the reasoning process when answering the prompt. This technique is about using real-time interactions with LLMs to guide toward more accurate and comprehensive responses [15]. Figure 6 shows the implementation of the Chain of Thought Prompting.

## Prompt Validation

Prompt validation ensures that the prompt or request made to interact with a large language model is accurate and concise and adheres to specific guidelines or policies set by the platform. This is crucial in the growing demand for large language models as usage across different companies is growing to solve business use cases seamlessly and rapidly. It is similar to code validation in software development, but prompt validation validates the context and instruction based on attributes like clarity, detail, policy compliance, and ethical consideration.

## CONCLUSION

Large language models greatly impacted the community, especially in the technology space. Most companies and engineers started using large language models in different ways to improve productivity. On the other hand, there are security issues posed intentionally or unintentionally by not utilizing the models properly, and this becomes a challenge to an organization from attackers. This paper attempts to present the misuse of prompts with or without individual knowledge to bypass the restrictions of models, and there are a few techniques like prompt injection, prompt leaking, and jailbreaking utilized by attackers to get into sensitive information. In addition, this paper covered overcoming prompt misuse principles, techniques, and patterns. They demonstrated a few use cases using OpenAI ChatGPT and a few outside sources to show objective evidence of how these models are attacked to bypass their ethical and sensitive information. They also demonstrated how to utilize large language models efficiently using various techniques. This paper should provide clear insights on utilizing and not utilizing the large language models, which is essential while using the GenAI models for better outcomes and unleashing the full potential of Artificial intelligence.

## REFERENCES

1. Transformer, [online] Available: https://medium.com/@amanatulla1606/transformer-architecture-explained-2c49e2257b4c
2. Large Language Models — What are Large Language Models? March. 2024, [online] Available: https://www.nvidia.com/en-us/glossary/large-language-models/
3. A Brief History of Large Language Models, [online] Available: https://www.dataversity.net/a-brief-history-of-large-language-models/
4. A History of Generative AI: From GAN to GPT-4, March. 2024, [online] Available: https://www.marktechpost.com/2023/03/21/a-history-of-generative-ai-from-gan-to-gpt-4/
5. 2001 - Neural language models [online] Available: https://www.ruder.io/a-review-of-the-recent-history-of-nlp/#2001neurallanguagemodels
6. What Are Generative Pre-Trained Transformers? [online] Available: https://www.makeuseof.com/gpt-models-explained-and-compared/
7. M. Gupta, C. Akiri, K. Aryal, E. Parker and L. Praharaj, "From ChatGPT to ThreatGPT: Impact of Generative AI in Cybersecurity and Privacy," in IEEE Access, vol. 11, pp. 80218-80245, 2023, doi: 10.1109/ACCESS.2023.3300381.

8. How to Jailbreak ChatGPT with these Prompts, [online] Available: https://www.mlyearning.org/how-to-jailbreak-chatgpt/

9. Prompt Injection, [online] Available: https://hackaday.com/2023/05/19/prompt-injection-an-ai-targeted-attack/

10. AI-Powered Bing Chat Spills its Secrets Via Prompt Injection Attack, Jun. 2023, [online] Available: https://arstechnica.com/information-technology/2023/02/ai-powered-bing-chat-spills-its-secrets-via-prompt-injection-attack/

11. Prompting Leaking, [online] Available: https://learnprompting.org/docs/prompt_hacking/leaking

12. Risks & Misuses, [online] Available: https://www.promptingguide.ai/risks

    Prompting Techniques, [online] Available: https://www.promptingguide.ai/techniques

    What Are Zero-Shot Prompting and Few-Shot Prompting, [online] Available: https://machinelearningmastery.com/what-are-zero-shot-prompting-and-few-shot-prompting/

    Chain-of-Thought Prompting, [online] Available: https://www.promptingguide.ai/techniques/cot

**Naresh Vurukonda** All current employment is with Amgen, Tampa, FL, USA. The author's latest degree is an M.S. in Computer Science from Southern Arkansas University. The author's research interests are deep learning, machine learning, generative AI, and data engineering. He is a member of the IEEE and the IEEE Computer Society. Contact him at vurukondanaresh@gmail.com.

**Shivendra Srivastava** current employment is with AWS, Seattle, WA, USA. The author's latest degree is an M.S. in Computer Science from the Georgia Institute of Technology. The author's research interests are cloud computing, machine learning, and generative AI. He is a member of the IEEE and the IEEE Computer Society. Contact him at mail2shivendra@gmail.com.

# Dynamic Access Control Mechanisms for Secure Data Sharing in Cloud Services

Abhay Dutt Paroha, *Software Team Leader, SLB, Houston, TX, USA*

*Abstract—Upstream oil and gas production operational data represent a crucial component in the oil and gas industry, facilitating the storage and management of operational information. With the advent of cloud computing, the landscape of operational data management has undergone significant transformation, offering unprecedented opportunities for accessibility, scalability, and efficiency. However, the integration of cloud computing into operational data management brings forth formidable challenges related to data security, privacy, and access control. Traditional methods of operational data management often rely on centralized architectures, which may present vulnerabilities to unauthorized access and data breaches [1]. In contrast, cloud-based operational data systems promise decentralized storage, real-time access, and seamless scalability. Yet, the transition to cloud-based operational data systems necessitates robust security mechanisms to safeguard sensitive operational information against various threats.*

***Keywords***: *Cloud computing, Secure dynamic access control, Upstream oil and gas production operational data*

Upstream oil and gas production operational data represent a crucial component in the oil and gas industry, facilitating the storage and management of operational information. With the advent of cloud computing, the landscape of operational data management has undergone significant transformation, offering unprecedented opportunities for accessibility, scalability, and efficiency. However, the integration of cloud computing into operational data management brings forth formidable challenges related to data security, privacy, and access control. Traditional methods of operational data management often rely on centralized architectures, which may present vulnerabilities to unauthorized access and data breaches [1]. In contrast, cloud-based operational data systems promise decentralized storage, real-time access, and seamless scalability. Yet, the transition to cloud-based operational data systems necessitates robust security mechanisms to safeguard sensitive operational information against various threats.

The surge in upstream oil and gas production operational data utilization is closely linked to digitizing operational information. With the proliferation of interconnected data systems in the oil and gas industry, operational data has experienced a significant uptick in adoption [3]. This growth has been accompanied by advancements in operational efficiency and resource optimization efforts. One of the primary concerns in cloud-based operational data management is ensuring dynamic access control while preserving data integrity and confidentiality. Achieving this balance requires innovative approaches that address the complex interplay between user access rights, data encryption, and authentication mechanisms. Moreover, the dynamic nature of oil and gas production environments demands flexible access control schemes capable of adapting to changing user roles and permissions. In response to these challenges, researchers have proposed novel access control schemes that leverage cryptographic techniques, such as encryption and key management, to enforce fine-grained access policies in cloud-based operational data systems [2]. These schemes aim to empower users with greater control over their operational data while mitigating the risk of unauthorized access or data manipulation. The current development of data exchange standards, exemplified by industry-specific protocols, in conjunction with data systems and related applications, has facilitated oil and gas professionals in adding, modifying, and exchanging operational data via computers or mobile devices. Primarily focused on operational data management

and transmission, these applications are overseen and operated by data providers responsible for facilitating data exchange between upstream, midstream, and downstream sectors. In response to this landscape, this paper proposes a dynamic access structure designed to precisely control access to operational data stored on cloud servers within a multi-user environment. We aim to empower every user with maximum control over their operational data. To achieve this, we employ cryptography based on Lagrange multipliers for encrypting the operational data [4]. This approach allows each custodian to generate their related keys, granting users the freedom to choose whom to share their operational data with. Central to our proposal is the enhancement of operational data encryption and the refinement of user dynamic access policies. To simplify key distribution, we depart from traditional hierarchical models and introduce a partial order relation to manage users. This significantly reduces the complexity of key management while enabling users to maintain control over operational data access. Furthermore, our approach facilitates the issuance of limited access rights to other users, such as engineers, geologists, technicians, and researchers [5]. This flexible method of multi-user dynamic access control accommodates the immediate addition or removal of user access, as well as the addition and modification of operational data. As a result, it is well-suited for upstream oil and gas production operational data cloud applications, meeting the evolving needs of users and industry stakeholders alike.

## Related Works

Several studies have investigated the integration of upstream oil and gas production operational data with cloud computing and secure data sharing in cloud services aiming to enhance data management, accessibility, and security [6]. This has been a topic of significant research and development due to the increasing adoption of cloud computing for storing and processing data. Several related works have focused on addressing various aspects of secure data sharing in cloud services, including encryption techniques, operational data in Cloud Computing, access control mechanisms, and key management protocols.

### A. Encryption Techniques

Numerous studies have explored different encryption techniques to ensure the confidentiality of data shared in cloud services [7]. This includes symmetric encryption algorithms such as AES (Advanced Encryption Standard) and asymmetric encryption techniques like RSA (Rivest-Shamir-Adleman). Additionally, homomorphic encryption schemes have been investigated to enable computation on encrypted data without decrypting it, thus preserving data privacy.

### B. Data Integration and Interoperability

Efforts have been made to address the challenge of integrating operational data from disparate sources and ensuring interoperability with existing systems [8]. This involves developing standards, protocols, and data exchange formats to facilitate seamless data exchange and interoperability among different operational data systems and stakeholders.

### C. Access Control Mechanisms

Effective access control mechanisms are crucial for secure data sharing in cloud services. Access control models such as Role-Based Access Control (RBAC) and Attribute-Based Access Control (ABAC) have been extensively studied and applied in cloud environments [9]. ABAC, in particular, offers fine-grained access control by considering various attributes of users, resources, and environmental conditions.

### D. Key Management Protocols

Key management is essential for maintaining the integrity and confidentiality of shared data. Various key management protocols have been proposed to securely generate, distribute, and revoke cryptographic keys in cloud environments. Key management schemes often involve techniques such as key encryption key (KEK) management, key derivation, and key rotation to enhance security and mitigate risks associated with key compromise.

### E. Secure Data Sharing Protocols:

Researchers have developed secure data-sharing protocols specifically tailored for cloud services [10]. These protocols aim to facilitate efficient and secure sharing of sensitive data among multiple users or entities while preserving data confidentiality, integrity, and availability. Examples include secure multi-party computation (SMPC) protocols, secure data aggregation schemes, and secure data-sharing frameworks built on cryptographic primitives.

### F. Privacy-Preserving Techniques

Privacy-preserving techniques have been explored to protect sensitive data from unauthorized access or disclosure [11]. Techniques such as differential privacy,

**FIGURE 1.** Private key cryptosystem

data anonymization, and secure multiparty computation enable data sharing while preserving individual privacy and confidentiality. These techniques are particularly relevant in scenarios involving sensitive operational data, financial information, or personally identifiable information (PII). Overall, related works in secure data sharing in cloud services and operational data in Cloud Computing encompass a broad range of topics, including encryption techniques, access control mechanisms, key management protocols, secure data sharing protocols, and privacy-preserving techniques. The ongoing research in this field aims to address various challenges and requirements associated with securely sharing sensitive data in cloud environments, including data confidentiality, integrity, availability, scalability, and interoperability.

## Lagrange interpolation polynomial

The following provides a concise introduction to the Lagrange interpolation polynomial, utilized in both encryption and decryption procedures. In numerical analysis and various applications, many practical issues necessitate the representation of functions to depict inherent relationships or patterns [13]. However, determining the precise relationship between the variables x and y for numerous functions can be exceedingly intricate and may not be discernible through experimentation alone. The Lagrange interpolation method facilitates the derivation of a polynomial that intersects a finite set of points on the x-y plane. This resultant polynomial is referred to as the Lagrange polynomial [12]. Mathematically, the Lagrange interpolation polynomial constructs a polynomial function that passes through known points within a two-dimensional plane. For instance, in an x-y plane, if n+1 points are given as $x0, y0,\ x1, y1,\ \dots,\ xn, yn$, the Lagrange interpolation method furnishes a formula for generating a unique polynomial of degree n that traverses through these n+1 points. Given a set of points $x1, y1, x2, y2, \dots, xn, yn$, the Lagrange interpolation polynomial $Lx$ is expressed as:

$$\ell_j(x) = \prod_{\substack{i=0 \\ i \neq j}}^{n} \frac{x - x_i}{x_j - x_i} \tag{1}$$

$$= \left( \frac{x - x_0}{x_j - x_0} \right) \dots \left( \frac{x - x_{j-1}}{x_j - x_{j-1}} \right) \left( \frac{x - x_{j+1}}{x_j - x_{j+1}} \right) \dots \left( \frac{x - x_n}{x_j - x_n} \right),$$

$$1 \leq j \leq n$$

The specific point of $\ell_j(x)$ is the derived value 1 from $x_j$. Values from other points $x_i$ (where $i \neq j$) equal 0, the expression of which is as follows:

$$\ell_j(x) = \begin{cases} 0, & i \neq j \\ 1, & i = j \end{cases}$$

The Lagrange polynomial is $L(x) = \sum_{j=0}^{n} y_j \ell_j(x)$. That is the unique polynomial of degree $n$ which passes through the points $(x_0, y_0), (x_1, y_1), \dots, (x_n, y_n)$. For example, the binomial that passes through $(4, 1), (5, 5)$, and $(6, 10)$ when expressed in Lagrange basic polynomial is as follows:

$$\ell_1(x) = \left( \frac{x - 5}{4 - 5} \right) \left( \frac{x - 6}{4 - 6} \right), \quad \ell_2(x) = \left( \frac{x - 4}{5 - 4} \right) \left( \frac{x - 6}{5 - 6} \right),$$

$$\ell_3(x) = \left( \frac{x - 4}{6 - 4} \right) \left( \frac{x - 5}{6 - 5} \right) \tag{2}$$

By applying the Lagrange interpolation polynomial, a single polynomial L(x) can be obtained as expressed below:

$$L(x) = f(4)\ell_1(x) + f(5)\ell_2(x) + f(6)\ell_3(x)$$

$$L(x) = 1 \times \left( \frac{x - 5}{4 - 5} \right) \left( \frac{x - 6}{4 - 6} \right) + 5 \times \left( \frac{x - 4}{5 - 4} \right) \left( \frac{x - 6}{5 - 6} \right)$$

$$+ 10 \times \left( \frac{x - 4}{6 - 4} \right) \left( \frac{x - 5}{6 - 5} \right)$$

$$= \frac{1}{2}x^2 - \frac{1}{2}x - 5 \tag{3}$$

It can be inferred that f(4)=1,f(5)=5,f(6)=10. By applying this formula predicted values can be derived, for example: to derive f(18), substitute x=18 in L(x), and L(18)=f(18)=148 is derived.

This study introduces a dynamic access scheme designed to securely manage upstream oil and gas production operational data in Cloud computing environments. In this setup, multiple users can access, modify, or share operational data, including appending, revising, deleting, and querying information [14]. However, the access permissions for different users can be complex, with distinct authorities depending on their roles and responsibilities. For example, field

engineers may have the ability to input real-time data such as pressure readings and production rates. Once supervisors or geologists provide their analysis and interpretations, field engineers may no longer be able to modify this information. Access to operational data may vary based on technical expertise and departmental affiliation, with even personnel within the same department having restricted access. Beyond field engineers and supervisors, other personnel such as technicians, environmental specialists, and regulatory compliance officers may also have varying levels of access for specific tasks like updating equipment status, reviewing environmental impact assessments, or examining regulatory compliance records. Additionally, individuals with lower authorization levels, such as contractors or researchers, may be granted read-only access to certain information. This access model extends beyond traditional production sites, including remote drilling operations, offshore platforms, and integrated asset management scenarios.

## Lagrange basis polynomials:

These Lagrange basis polynomials have the property that:

$$L(x) = \sum_{i=0}^{n} y_i \prod_{\substack{i=0 \\ i \neq j}}^{n} \frac{x - x_j}{x_j - x_i} \qquad (4)$$

$$\ell_j(x) = \prod_{\substack{i=0 \\ i \neq j}}^{n} \frac{x - x_i}{x_j - x_i}$$

$$= \left( \frac{x - x_0}{x_j - x_0} \right) \cdots \left( \frac{x - x_{j-1}}{x_j - x_{j-1}} \right) \left( \frac{x - x_{j+1}}{x_j - x_{j+1}} \right) \cdots \left( \frac{x - x_n}{x_j - x_n} \right)$$

$$1 \leq j \leq n \qquad (5)$$

$$\ell_j(x) = \prod_{\substack{i=0 \\ i \neq j}}^{n} \frac{x - x_i}{x_j - x_i}, \quad 1 < n \qquad (6)$$

This polynomial L(x) passes through all the given data points and can be used to approximate the function that generated those points over the range of x values covered by the data.

## Initialization:

This study employs a partially ordered access system, with a central authority (CA) establishing the framework. In this system, users are categorized into distinct sets (Si for i = 0, 1, 2,..., n), representing security classes. Each class is granted specific authorization to access designated files, for which they receive decryption keys for encrypted content.

The relationship between these sets is defined by a binary partial order relation ($\preceq$) over the set $(S, \preceq)$, where $S$ is the collection of security classes $(S_1, S_2, \ldots, S_n)$. Within this arrangement, $S_j \preceq S_i$ (where $i, j \in N$) indicates that users in class $S_i$ can access data held by those in class $S_j$, but not vice versa.

For instance, if $S_j$ includes $\{1, 2\}$ and $S_i$ includes $\{1, 2, 3\}$, with $\{1, 2\} \preceq \{1, 2, 3\}$, then $S_j$ is deemed to be $\preceq S_i$. In the scenario where $S_j \preceq S_i$, it indicates that users in $S_i$ have access to the decryption keys for authorized files 1 and 2 in $S_j$.

The upstream oil and gas production operational data system encompasses a diverse range of users, including field engineers, supervisors, geologists, technicians, environmental specialists, regulatory compliance officers, contractors, and researchers, each assigned to a security class represented by $S_i$ with a unique superkey $H_i$, where $i = 0, 1, 2, \ldots, n$. The central authority (CA) establishes a structured framework for these users, comprising $n$ individuals forming two sets: $S = \{S_1, S_2, \ldots, S_n\}$ and $H = \{H_1, H_2, \ldots, H_n\}$.

Figure 3 illustrates that Cloud computing services offer a range of virtualized resources, accessible over the Internet, to meet diverse computing needs [15]. These services include Infrastructure as a Service (IaaS), providing virtualized computing infrastructure like servers and storage. Platform as a Service (PaaS) offers development and deployment tools, enabling developers to build and host applications without managing underlying infrastructure. Software as a Service (SaaS) delivers fully functional applications over the Internet, eliminating the need for installation or maintenance. Additionally, cloud providers offer specialized services such as database management, artificial intelligence, and Internet of Things (IoT) platforms. These services are scalable, flexible, and pay-per-use, providing cost-effective solutions for businesses of all sizes. This production operational data management system

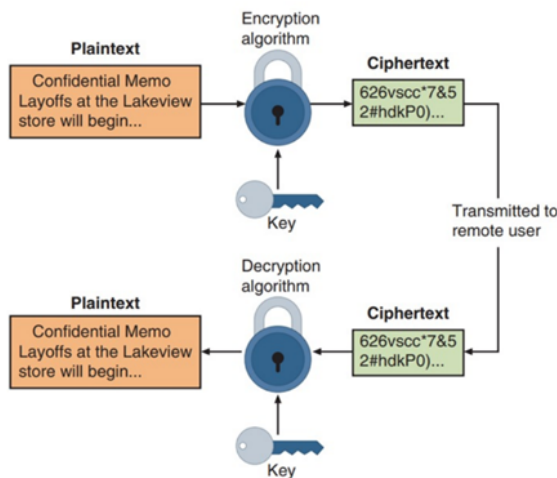**FIGURE 3.** Access Environment in Cloud Computing



**FIGURE 4.** Encryption and decryption technology

consolidates operational data from various sources and users. Each user's operational data is encrypted to create an encrypted file stored on Cloud servers. A central authority (CA) organizes these files into a set (file 0 file1, file2,...., file) and generates a unique decryption key (DKu) for each file (u = 1, 2, ..., m). This encryption safeguards the files from unauthorized access, ensuring data security and privacy.

## A. Cryptography and encryption systems

While migrating upstream oil and gas production operational data to a cloud environment heightens security risks, maintaining data integrity, confidentiality, and availability remains imperative. Given that the opera-

tional data management system's core aim is to provide authorized users with secure access, we achieve this goal through cryptography. Cryptography is the science and practice of securing communication and data by encoding information in such a way that only authorized parties can access it. Encryption systems use cryptographic algorithms to convert plaintext data into ciphertext, making it unreadable to unauthorized individuals. These systems typically involve the use of cryptographic keys to encrypt and decrypt data, with stronger encryption methods employing longer and more complex keys. Cryptography and encryption systems play a crucial role in protecting sensitive information across various domains, including communications, finance, and cybersecurity. They ensure data confidentiality, integrity, and authenticity, safeguarding against unauthorized access, tampering, and eavesdropping. Additionally, advancements in cryptography continually drive innovation in security technologies, enabling organizations to adapt to evolving threats and regulatory requirements.

## B. Algorithm for Encryption
We've designed an algorithm for a data security service provider with the following features:

- It incorporates defined keywords enabling internet users to locate the owner of specific encrypted data.
- It allows individuals to securely share data with their groups.
- It offers individuals the ability to monetize their data by selling it on the cloud.
- It enables organizations to share data with their outsourced projects.
- It facilitates organizations in sharing data with their customers securely.
- It provides organizations with the option to sell their data to both other organizations and individuals.

---

**Algorithm 1** Set Private

---
1: sets only the owner on ACL.
2: Set permission on ACL = P3.
3: If ACL matches the existing one, delete the new one and set ACL = existing.
4: Set DAK = ACL.
5: Encrypt InD & save on storage.

---

To fulfill the previously outlined objectives, we have devised an algorithm predicated on the requirement that every individual accessing specific data must possess a cloud account, thereby establishing them as a

valid user. This algorithm operates by receiving input data labeled as 'InD'. Each document is assigned an owner, whether an individual or an organization, with an organizational document featuring an administrator denoted as `admin`. Additionally, a Data Access Key (DAK) is generated to serve as the identifier for the Access Control List (ACL) associated with the InD. This key is integrated into the file during encryption, ensuring secure storage.

## Implementation and Results

Implementing an upstream oil and gas production operational data system in cloud computing involves several key steps, including designing the architecture, selecting appropriate cloud service providers, implementing security measures, and integrating with production systems. The system should ensure data privacy, integrity, and availability while allowing authorized users to access and manage their operational data securely. Once implemented, the operational data system can be evaluated based on various metrics such as system performance, user satisfaction, data security, and compliance with regulatory standards. The results of the implementation can be assessed through user feedback, system logs, and performance monitoring tools.

---

**Algorithm 2** For Inter Enterprise Access

---

1: Set an owner on ACL.
2: **for** $i$ = 1 to $n$ **do**
3:     Verify if U C C and CE C RC and UR = admin.
4:     Add CE to the Organization Export on ACL.
5:     If ACL matches the existing one, delete the new and set ACL = existing.
6:     Set DAK = ACL.
7: **end for**
8: Encrypt InD and save on storage.

---

Due to the lack of access to multiple clouds and a large user base, we opted to conduct simulations rather than real-world testing. Our algorithm outperforms existing approaches by minimizing the size of the ciphertext while offering multilevel user access control. Designed for an open community, our algorithm caters to all seven types of data users. We conducted tests using a 33-character Data Access Key (DAK), which proved sufficient to assign unique keys to trillions of files. Remarkably, this 33-character DAK only adds 34 bytes to the size of any appended file. The results of implementing secure data sharing in cloud services can be measured by evaluating the effectiveness of access control mechanisms, the strength of encryption
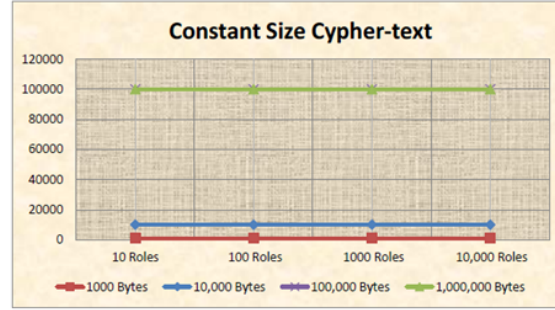


**FIGURE 5.** Constant Size of Encrypted File

algorithms, and compliance with data protection regulations. Organizations may also conduct penetration testing and security audits to identify and address any vulnerabilities in the system. Overall, the successful implementation of of upstream oil and gas production operational data in cloud computing and secure data sharing in cloud services requires a comprehensive approach to design, implementation, and evaluation, with a focus on protecting sensitive information and ensuring user privacy and data security.

The graphical depiction of this concept is illustrated in Fig. 5. Following implementation, we conducted tests on the algorithm by adjusting the number of roles in each iteration, yielding remarkable findings. The encryption time remained consistent due to the retrieval of roles from the database, as it solely generates a data access key during encryption and embeds it into the encrypted file. However, the decryption time showed a slight increase with the enhancement of the number of roles, as user access verification is performed from the database. Nonetheless, this increase in decryption time is minimal for small variations in the number of roles. Moreover, data management remained relatively consistent regardless of the number of roles.

## Limitations and Workarounds

The Lagrange interpolation polynomial can simplify key management but may introduce complexity in computation and maintenance. To address this, efficient algorithms and optimizations for polynomial interpolation can be implemented, or hybrid approaches combining polynomial interpolation with other key management techniques can be considered. Scalability concerns may arise due to the large number of users or frequent changes in access rights. To improve scalability, the

scheme should incorporate hierarchical or distributed key management and use data structures that support efficient updates and queries.

Computational overheads from Lagrange interpolation and dynamic access control management could impact system performance. Optimizing polynomial computation and performing profiling and benchmarking can help identify and address performance bottlenecks. Storage requirements for key-related data and access control information may be substantial, especially with a large number of users. Data compression techniques and efficient data storage formats can reduce storage overhead. Security resistance to specific attacks is crucial, and regular updates and audits can adapt to new threats. User experience and usability can be improved by developing user-friendly interfaces and tools for managing access rights. Integrating the scheme with existing systems can be challenging, but using standard protocols and interfaces can facilitate integration. Privacy preservation can be enhanced by employing advanced techniques like homomorphic encryption or secure multi-party computation.

## Ethical implications, security risks and mitigation measures

Using upstream oil and gas production operational data in the cloud presents a range of ethical implications, security risks, and mitigation measures. Operational data, which includes sensitive information about production processes, equipment performance, and operational strategies, can be sensitive and potentially compromised by unauthorized access or misuse. Mismanagement can impact various stakeholders, including local communities, employees, and business partners. Data ownership disputes may arise, especially if third-party cloud providers are involved, affecting data control and utilization. Ethical considerations are crucial, particularly in ensuring data is not exploited for unintended purposes or harmful to people or the environment. Lack of transparency in data handling practices raises ethical concerns. Organizations must ensure compliance with data use and privacy regulations, maintaining ethical standards in their operations.

Cloud environments pose several security risks, including data breaches due to unauthorized access, insider threats, data integrity, and data loss due to accidental deletion, corruption, or failure of cloud service providers. These risks can impact operational information accuracy and reliability, as well as the reliability of third-party cloud providers. Additionally, geographical risks arise from the different laws and regulations governing data stored in cloud servers, potentially leading to legal and security issues. Therefore, it is crucial to consider these potential risks when implementing cloud solutions.

Data encryption is a crucial aspect of data protection, ensuring data is secure from unauthorized access. It involves strong encryption standards, regular updates, and robust key management practices. Access controls, such as role-based access controls (RBAC) and Multi-Factor Authentication (MFA), add a layer of security. Regular security audits and vulnerability assessments are conducted to identify and address potential weaknesses. Continuous monitoring and logging of data access and usage are implemented to detect and respond to suspicious activities promptly. Regular data backup procedures and a disaster recovery plan are established to protect against data loss. Compliance with data protection regulations and industry standards is ensured. Clear contracts and service level agreements (SLAs) with cloud providers are negotiated to define responsibilities and expectations related to data security and privacy. Employee training on data security best practices and awareness programs are also provided to keep employees informed about emerging threats and the importance of data protection.

## Conclusion

In conclusion, dynamic access control mechanisms play a crucial role in ensuring secure data sharing in cloud services, particularly in the context of upstream oil and gas production operational data. By dynamically adapting access controls based on contextual factors such as user roles, privileges, and environmental conditions, these mechanisms offer a robust defense against unauthorized access and data breaches. Through our exploration of dynamic access control approaches leveraging Lagrange interpolation polynomial and partial order relations, it is evident that tailored combinations of these methods can effectively address the unique security and privacy requirements of cloud-based operational data management. Moreover, by incorporating encryption techniques, access control models, and key management protocols, organizations can strengthen data security while facilitating seamless and efficient data sharing and management in cloud environments. However, challenges such as interoperability, scalability, and user acceptance remain to be addressed to fully leverage the potential of dynamic access control mechanisms in cloud-based operational data management. Moving forward, continued research and innovation in this area are essential to stay ahead of evolving threats and ensure

the integrity, confidentiality, and availability of sensitive operational data in the cloud.

## REFERENCES

1. I. Gupta, A. K. Singh, C.-N. Lee, and R. Buyya, "Secure data storage and sharing techniques for data protection in cloud environments: A systematic review, analysis, and future directions," IEEE Access, 2022.

2. U. Narayanan, V. Paul, and S. Joseph, "A novel system architecture for secure authentication and data sharing in cloud-enabled Big Data Environment," Journal of King Saud University-Computer and Information Sciences, vol. 34, no. 6, pp. 3121-3135, 2022.

3. Y. Liu, W. Yang, and Y. Wang, "An access control model for data security sharing cross-domain in consortium blockchain," IET Blockchain, vol. 3, no. 1, pp. 18-34, 2023.

DD1. ang, Xudong, Bello, Oladele, Yang, Lei, Bale, Derek, and Roberto Failla. "Intelligent Oilfield - Cloud Based Big Data Service in Upstream Oil and Gas." Paper presented at the International Petroleum Technology Conference, Beijing, China, March 2019. doi: https://doi.org/10.2523/IPTC-19418-MS

EE1. . Sultana, A. Almogren, M. Akbar, M. Zuair, I. Ullah, and N. Javaid, "Data sharing system integrating access control mechanism using blockchain-based smart contracts for IoT devices," Applied Sciences, vol. 10, no. 2, p. 488, 2020.

FF1. . Prabhu Kavin, S. Ganapathy, U. Kanimozhi, and A. Kannan, "An enhanced security framework for secured data storage and communications in the cloud using ECC, access control and LDSA," Wireless Personal Communications, vol. 115, pp. 1107-1135, 2020.

GG1. . Yang, L. Tan, N. Shi, B. Xu, Y. Cao, and K. Yu, "AuthPrivacyChain: A blockchain-based access control framework with privacy protection in the cloud," IEEE Access, vol. 8, pp. 70604-70615, 2020.

HH1. . Seth, S. Dalal, V. Jaglan, D. N. Le, S. Mohan, and G. Srivastava, "Integrating encryption techniques for secure data storage in the cloud," Transactions on Emerging Telecommunications Technologies, vol. 33, no. 4, p. e4108, 2022.

II1. . Ahmadi, "Security Implications of Edge Computing in Cloud Networks," Journal of Computer and Communications, vol. 12, no. 2, pp. 26-46, 2024.

JJ1. . Almasian and A. Shafieinejad, "Secure cloud file sharing scheme using blockchain and attribute-based encryption," Computer Standards & Interfaces, vol. 87, p. 103745, 2024.

KK1. . O.-B. O. Agyekum, Q. Xia, E. B. Sifah, C. N. A. Cobblah, H. Xia, and J. Gao, "A proxy re-encryption approach to secure data sharing in the Internet of things based on blockchain," IEEE Systems Journal, vol. 16, no. 1, pp. 1685-1696, 2021.

LL1. . Han, Y. Zhu, D. Li, W. Liang, A. Souri, and K.-C. Li, "A blockchain-based auditable access control system for private data in service-centric IoT environments," IEEE Transactions on Industrial Informatics, vol. 18, no. 5, pp. 3530-3540, 2021.

MM1. . Huang, C. Wang, and L. Chen, "Secure and Fine-Grained Flow Control for Subscription-Based Data Services in Cloud-Edge Computing," IEEE Transactions on Services Computing, 2022.

NN1. . Liu, J. Zhang, and J. Zhan, "Privacy protection for fog computing and the internet of things data based on blockchain," Cluster Computing, vol. 24, pp. 1331-1345, 2021.

OO1. . Feng et al., "Efficient and secure data sharing for 5G flying drones: a blockchain-enabled approach," IEEE Network, vol. 35, no. 1, pp. 130-137, 2021.