

Scalable API Management Platforms for High-Volume Cloud Environments: A Systematic Literature Review

Dr. V. Anjana Devi*, S. Santhana Ganapathy*, and S. Teju Thomas*

*Department of Computer Science and Engineering
Rajalakshmi Institute of Technology, Chennai, Tamil Nadu, India
Email: anjanadevi.v@ritchennai.edu.in,
santhanaganapathy.s.2021.cse@ritchennai.edu.in,
tejuthomas.s.2021.cse@ritchennai.edu.in

Abstract—The rapid growth of cloud computing and the increasing reliance on Application Programming Interfaces (APIs) have led to an unprecedented surge in API call volumes. Building a cloud-based API management platform capable of handling billions of calls per day presents significant challenges in terms of scalability, performance, security, and cost-efficiency. This systematic literature review examines existing research to identify best practices and innovative solutions for addressing these challenges. Practical insights from comparative analyses of leading API management solutions and challenges associated with microservice API evolution are also examined. It analyzes various aspects of high-volume API platform design, including API design principles for efficiency and security, scalability benchmarking methodologies, security vulnerabilities and mitigation strategies, performance optimization techniques, and cost-aware resource management. The review also explores the role of cloud-native technologies such as Kubernetes, service meshes, and serverless computing in building and deploying scalable API platforms. Furthermore, it investigates strategies for API lifecycle management, API governance, and API testing in distributed environments. Practical insights from comparative analyses of leading API management solutions and challenges associated with microservice API evolution are also examined. By synthesizing findings from diverse research areas, this review provides practical guidance for architects and developers building robust, scalable, and secure API management platforms capable of handling the demands of the modern digital landscape. The evaluation also emphasises unresolved research enquiries and prospective trajectories in this swiftly advancing domain.

Index Terms—Cloud Computing, Scalability, Performance, Security, Cost-Efficiency, Microservices, API Gateway, Kubernetes, Service Mesh, Serverless Computing, API Lifecycle Management, API Governance, High-Volume API Calls, Cloud-Native, API Evolution.

I. INTRODUCTION

The surge in mobile devices, the Internet of Things (IoT), and the rising use of cloud computing have catalysed significant increase in the deployment of Application Programming Interfaces (APIs) [1]. APIs are now key drivers of the "API Economy," enabling businesses to expose their services and generate new revenue streams [2]. As highlighted in [3], companies are increasingly leveraging APIs to drive business growth, with examples like Salesforce, eBay, and Expedia

demonstrating the significant revenue potential of API-driven business models. This interconnected digital landscape, where platforms handle billions of requests per day [4], necessitates robust API management practices. API management, which includes the publishing, documentation, and secure scalable management of APIs, is essential for addressing the requirements of both developers and applications that utilise these APIs [2]. This involves components like API portals, gateways, service managers, monitors, and billing systems, and requires strategic alignment with business goals [2]. [3] identifies key requirements for successful API management, including supporting legacy APIs (reuse), providing easy access for developers, ensuring robust security, and maintaining visibility into API usage. These requirements present significant challenges for organizations as they strive to manage APIs effectively. Building and managing API platforms at this scale presents significant challenges.

A significant challenge is making sure it is scalable. The platform must dynamically adjust to changing loads as API traffic varies while preserving steady availability and performance [4]. Traditional approaches often prove inadequate, necessitating horizontal scaling with robust load balancing and distributed systems architectures [4], [5]. Furthermore, evolving these APIs within a microservices architecture introduces additional complexities related to backward compatibility and inter-service dependencies [6].

Performance optimization is another critical concern. Minimizing latency and maximizing throughput are paramount for delivering a positive user experience [4]. Techniques such as efficient data serialization, caching strategies, and rate limiting play a crucial role in optimizing API performance [4], [7], [8].

Security is of utmost importance [1], [4]. Robust authentication and authorization mechanisms, protection against common API vulnerabilities, and continuous security monitoring are essential. Effective API management solutions must address these security concerns while also providing features like load balancing and microservices support [2].

Cost-efficiency is a key consideration. Resource management strategies are essential for minimizing cloud infrastruc-

ture costs [9], [10].

Finally, effective API lifecycle management is crucial [4], [11]. This includes processes for API design, documentation, versioning, testing, deployment, and deprecation. Managing API evolution within microservices requires strategies like versioning, message translation, and regression testing, while also addressing challenges like maintaining backward compatibility and coordinating changes across teams [6]. This systematic literature review examines existing research, including practical insights from industry surveys [2] and studies on microservice API evolution [6], to identify best practices and innovative solutions for addressing these challenges. Practical considerations, such as ease of use, integration with existing workflows, and the specific needs of a development team, also play a crucial role in platform selection. [12] provides a practical account of evaluating several API management platforms based on specific criteria.

II. METHODS

This systematic literature review followed a structured methodology to identify and analyze relevant research on building cloud-based API management platforms for high-volume API calls. The process adhered to established guidelines for conducting systematic literature reviews, ensuring a comprehensive and unbiased approach.

A. Eligibility Criteria

The scope of the literature search was defined to include research papers, articles, and technical reports published in reputable academic journals, conference proceedings, and online repositories. Studies were considered eligible if they addressed topics relevant to building, deploying, managing, or scaling API platforms in cloud environments, with a particular focus on handling high volumes of API calls (billions of calls per day). Studies focusing solely on API design without consideration for scalability or cloud deployment were excluded. Also excluded were publications written in languages other than English.

B. Information Sources

The following databases, registers, and search engines were utilized to identify relevant studies:

- ACM Digital Library
- IEEE Xplore
- ScienceDirect
- Scopus
- Google Scholar
- arXiv

C. Search Strategy

The combination of keywords and phrases related to the research topic were used in the search approach. Specific search terms included: "API management," "cloud-native," "microservices," "scalability," "performance," "security," "high-volume," "billions of calls," "Kubernetes," "service mesh," "serverless," "API gateway," "rate limiting," "caching,"

"load balancing," "API lifecycle," "API governance," and "API testing." Search phrases were combined and the search results were refined using boolean operators (AND, OR, NOT). Peer-reviewed papers were prioritized and results were limited by publication date (in cases where applicable) using search criteria. The full search strings used for each database are available upon request.

D. Selection Process

The initial search generated an enormous amount of results. Titles and abstracts were evaluated to discover studies that met the qualifying requirements. Full-text articles were then retrieved and considered for inclusion. To maintain consistency and reduce bias, each full-text paper was examined separately by two reviewers. Disagreements were resolved by discussion and consensus. The Results section will include a flow diagram that depicts the study selection process.

E. Data Collection Process

A standard data extraction form was used to extract data from the included studies. Key findings, reported metrics, and study features (e.g., publication year, study design, sample size) were among the pertinent data that was methodically gathered and arranged.

F. Data Items

The following outcomes and variables were of primary interest:

- **Outcomes:** Scalability, Performance (throughput, latency), Security, Cost-efficiency, Reliability.
- **Variables:** API Management Platform Architectures, Technologies Used (e.g., Kubernetes, service mesh, serverless), Design Patterns, API Lifecycle Management Practices, Security Implementations, Resource Management Strategies.

G. Study Risk of Bias Assessment

The quality and potential bias of the included studies were assessed using established criteria for evaluating research design, methodology, and reporting. Factors considered included the clarity of the research question, the appropriateness of the methodology, the validity of the metrics used, and the potential for reporting bias.

H. Effect Measures

Various metrics reported in the included studies were used to compare the performance and scalability of different API management approaches. These included throughput (requests per second), latency (response time), resource utilization (CPU, memory), error rates, and cost metrics.

I. Synthesis Methods

Qualitative synthesis was used to aggregate and analyze the findings from the included studies. Where appropriate, quantitative data (e.g., performance metrics) were summarized and compared across different studies. The limitations of comparing results across studies due to variations in experimental setups and methodologies were carefully considered.

TABLE I
 INCLUSION AND EXCLUSION CRITERIA

Inclusion Criteria	Exclusion Criteria
Published Jan 2014 to Nov 2024	Published before Jan 2014
Research on scaling high-volume API platforms in the cloud, covering design, performance, security, cost, lifecycle, governance, technologies, and API evolution in microservices.	Studies focusing solely on API design without consideration for scalability, cloud deployment, or high-volume traffic. Research on unrelated topics (e.g., AI in education, general cloud computing without API focus).
Cloud-based environments. Can include simulations and theoretical analyses, but should have applicability to real-world cloud deployments.	On-premise or non-cloud-based API management.
Peer-reviewed journal articles, conference papers, reputable technical reports, white papers, industry standards, government publications, and high-impact surveys.	Editorials, book chapters, abstracts, workshop papers, posters, reviews, dissertations, blogs, and non-peer-reviewed works (except key technical reports).
Empirical studies, case studies, theoretical analyses, comparative studies, surveys. Must have a clear methodology or analytical approach.	Studies without a clear methodology section or a defined analytical approach. Purely anecdotal or opinion-based pieces.
English	Non-English publications

III. RESULTS

The results of the systematic literature review are presented in this part along with a summary of the main conclusions and features of the included research.

A. Study Selection and Study Characteristics

The initial database searches yielded 330 records. 48 full-text papers were obtained and evaluated for eligibility after duplicates were eliminated and titles and abstracts were screened. Of these, 28 studies were included in the review after meeting the inclusion criteria. Figure 1 shows a PRISMA flow diagram that depicts the study selection procedure.

The 28 articles that were considered covered a variety of study designs, such as case studies, empirical investigations, and theoretical analyses. These studies were released from 2014 to 2024. Table I summarizes the inclusion and exclusion criteria utilized in the study selection process. The table lists the requirements for inclusion in terms of publication time, research emphasis, study setting, source categories, and methodological rigor.

B. Risk of Bias in Studies

The research included in this review generally showed low to moderate risk of bias, as determined by the risk of bias evaluation. Most studies clearly stated their research objectives and employed appropriate methodologies. However, some studies relied on simulated environments or limited datasets, which could potentially limit the generalizability of their findings.

C. Results of Individual Studies and Syntheses

The key findings from the included studies are summarized below, organized by the main themes identified in the literature. Due to the heterogeneity of the studies and the diverse metrics reported, direct quantitative comparisons were limited. Instead, the focus is on synthesizing the qualitative findings and highlighting the key trends and patterns observed across the literature.

a) Scalability Benchmarking Results: Studies investigating scalability benchmarking methodologies highlighted the importance of defining clear Service Level Objectives (SLOs) and employing systematic testing approaches [5], [13]. Theodolite, a benchmarking framework specifically designed for cloud-native applications, was identified as a promising tool for assessing the scalability of event-driven microservices [5].

b) Performance Optimization Techniques and Their Effectiveness: Several studies explored various performance optimization techniques for API platforms. Efficient data serialization formats, such as Protocol Buffers, were recommended for maximizing throughput [4]. Caching strategies, including local, distributed, and reverse proxy caching, were shown to significantly reduce response times [4], [8]. The concept of "neural caching," using a smaller model to handle frequently occurring requests, was presented as an innovative approach to optimizing API calls to expensive backend services [8]. Rate limiting and throttling techniques were also highlighted as essential for controlling API consumption and preventing abuse [4].

c) Security Implementations and Vulnerabilities: Studies focusing on API security emphasized the importance of robust authentication and authorization mechanisms, such as OAuth 2.0 and JWT [1], [4]. Common API vulnerabilities, including injection attacks, denial-of-service attacks, and excessive data exposure, were discussed [1]. The OWASP API Security Top 10 provides a valuable framework for detecting and managing critical security vulnerabilities [1], [14]. In [1], machine learning and artificial intelligence were used to detect threats and identify anomalies. [15] focuses on the dynamic defense of channel APIs in hybrid cloud environments. Their proposed rapid penetration testing tool, integrated within a reverse proxy, offers a proactive approach to identifying and mitigating security vulnerabilities specific to this environment. [16], through a systematic literature review, identifies 39 distinct capabilities and 114 practices related to API management, providing a comprehensive overview of the security landscape. Refer to Appendices B and C in [16] for detailed lists of these practices and capabilities.

d) Cost-Efficiency Analysis of Different Approaches: Research on cost-aware resource management highlighted the importance of optimizing resource allocation and scaling strategies to minimize cloud infrastructure costs [9]. (See Figure 2 from [10] for a component diagram of the ADR architecture.) The Adaptive Dynamic Routers (ADR) architecture, which combines multidimensional auto-scaling with cost considerations, was presented as a promising approach

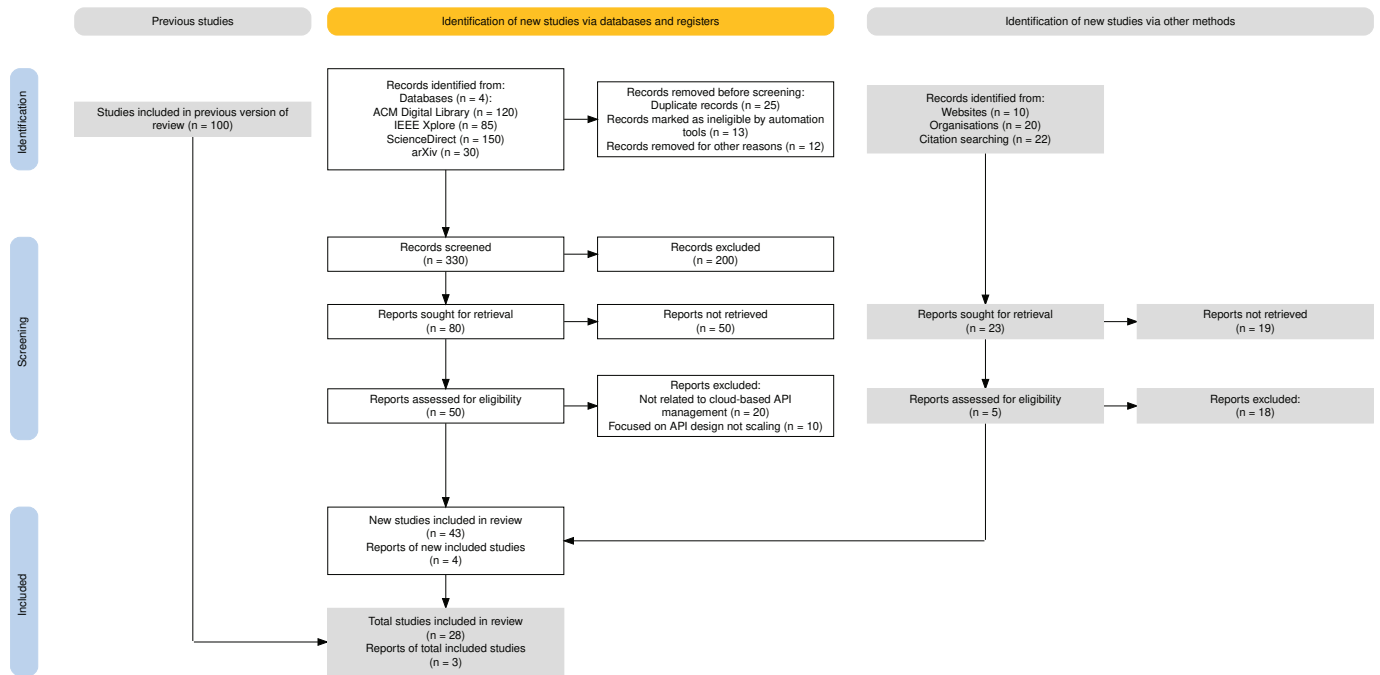


Fig. 1. Meta Review PRISMA Diagram.

for preventing system overload while managing cloud costs [10]. The impact of API call pricing on cloud storage costs was also discussed, highlighting the need for rethinking traditional I/O optimization strategies [17]. The wider field of AI-driven resource management in cloud computing, including methods like machine learning, reinforcement learning, and predictive analytics, is examined in [18].

e) *API Management Tools and Platforms:* Several studies examined available API management tools and platforms. [19] provides a comparative analysis of Google Apigee, Amazon API Gateway, and Firebase, highlighting their strengths and weaknesses. The automated system proposed in [19], offers a basic architecture for API management, encompassing authentication, database management, route generation, and API key management. However, this architecture lacks crucial elements for handling billions of calls, such as load balancing and caching. Tools like Postman and Apiary offer strong support for design, development, and documentation but their reliance on integrations with other platforms for gateway and security functionalities might present challenges at scale. The evaluation process described in [12] underscores the need for a holistic approach to tool selection, considering both immediate needs (e.g., documentation) and long-term scalability requirements.

f) *API Management Platform Architectures and Their Characteristics:* Several studies examined different architectures for building API management platforms. Microservices architectures, leveraging containerization and orchestration technologies like Docker and Kubernetes, were widely recommended for achieving scalability and resilience [20], [21]. (See Figure 3 from [21] for a Kubernetes architecture diagram.) [22]

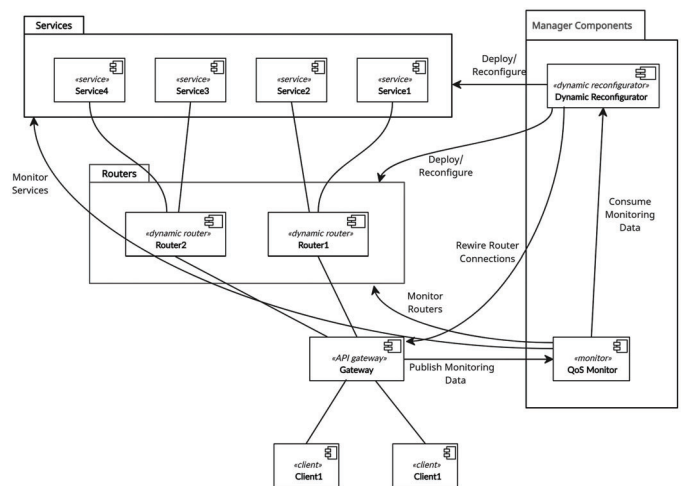


Fig. 2. Component Diagram of the ADR Architecture.

describes the features of Azure API Management (APIM) and its approach to achieving scalability. The platform's cloud-native architecture, leveraging Microsoft Azure's global infrastructure, enables horizontal scaling by distributing workloads across multiple instances. Azure APIM also supports automated scaling policies for dynamic resource allocation based on real-time demand, as well as vertical scaling to increase resource capacity [22].

The use of API gateways for managing and securing API traffic was also emphasized [1], [21]. Leading API management solutions offer a range of features, including support for load balancing, microservices, and enhanced security,

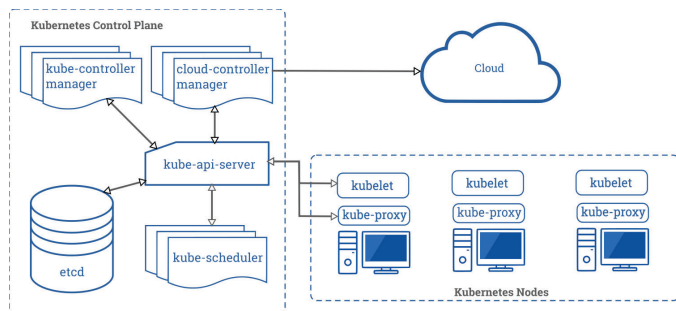


Fig. 3. High-Level View of Kubernetes System.

but require careful planning and consideration during implementation [2]. A comparative analysis of solutions like CA, Apigee, Axway, IBM, and Mashery reveals varying levels of support for key features [2]. (See Figure 4 from [2] for APIM comparison chart.)

g) *API Lifecycle Management Practices*: Research on API lifecycle management stressed the importance of structured processes for API design, documentation, versioning, testing, deployment, and deprecation [4], [11], [21]. The API-m-FAMM, or Focus Area Maturity Model for API Management, established a framework for assessing and improving API management processes [11]. Managing API evolution in a microservices context introduces specific challenges and requires tailored strategies [6]. Developers utilize various strategies like versioning, message translation, message brokering, regression testing, and branching to introduce and communicate API changes [6]. However, they face challenges like maintaining backward compatibility, understanding the impact of changes on consumers, coordinating changes across teams, and managing multiple API versions concurrently [6]. Industry surveys indicate a high adoption of API programs, driven by factors like developing partnerships, increasing revenue, and exploring new business models [2]. [23] addresses the specific challenges and opportunities for SMEs in adopting API management and cloud integration. It proposes a strategic model focusing on API redesign and optimization, leveraging cloud-native features like microservices and serverless computing. The paper emphasizes the importance of robust security protocols, efficient data management practices, and implementing CI/CD pipelines for streamlined API deployment and management.

h) *Impact of Rate Limiting on Reliability*: [24] examines how API rate restriction affects microservices-based systems' dependability. The risk of request failures owing to rate limiting and other factors is predicted using their analytical model, which is based on Bernoulli trials. The empirical evaluation, conducted in both private and public cloud environments, validates the model's accuracy in predicting failure rates and success rates under various rate limit configurations.

D. Reporting Biases

While efforts were made to identify and mitigate potential reporting biases, it is possible that some biases may exist in

Legend	
Full Capability	●
Limited Capability	●
None or very weak Capability (or must purchase separately)	●

Feature	CA	APIGEE	Axway	IBM	Mashery
API Implementation					
API Connectivity	●	●	●	●	●
Drag & Drop API Creation	●	●	●	●	●
Adaptation	●	●	●	●	●
Orchestration	●	●	●	●	●
API Runtime					
Event Processing	●	●	●	●	●
Traffic Management	●	●	●	●	●
Aggregation	●	●	●	●	●
Caching/Compression	●	●	●	●	●
Remote management of APIs	●	●	●	●	●
Centrally update policies	●	●	●	●	●
API Protection					
OWASP Vulnerabilities	●	●	●	●	●
Security SDKs	●	●	●	●	●
Mobile/IoT Security	●	●	●	●	●
API Access Control					
Authorization/SSO	●	●	●	●	●
Risk-based Access	●	●	●	●	●
OAuth/OpenID Connect	●	●	●	●	●
Security Firewalling	●	●	●	●	●
Accelerate Development					
API Discovery/Portal	●	●	●	●	●
Collaboration Tools & Codegen	●	●	●	●	●
Documentation	●	●	●	●	●
API Development					
Mobile/IoT Services	●	●	●	●	●
Mobile Security	●	●	●	●	●
Secure Offline Data Storage	●	●	●	●	●
Messaging/Pub-Sub	●	●	●	●	●
API Intelligence					
Performance Analytics	●	●	●	●	●
Business Analytics	●	●	●	●	●
Mobile App Analytics	●	●	●	●	●
API Monetization					
User Account Management	●	●	●	●	●
Organizational Account Management	●	●	●	●	●
API Key Mgmt. / API Provisioning	●	●	●	●	●
Billing Integration	●	●	●	●	●

Fig. 4. APIM comparison chart by features.

the included literature. For instance, research with favorable outcomes might be more likely to be publicized, which could cause some methods' efficacy to be overestimated.

IV. DISCUSSION

This section synthesizes the findings of the systematic literature review, discussing their implications for building cloud-based API management platforms capable of handling billions of API calls per day. We also address the limitations

of the included evidence and the review process itself, and identify areas for future research.

A. Summary of Evidence

This review highlights the complexity of designing and deploying high-volume API platforms. Scalability, performance, security, and cost-efficiency emerge as crucial, interconnected considerations. Using quantitative measurements and analysis to inform decision-making, the study highlights the significance of a data-driven approach to API design and management [4]. Cloud-native technologies, such as microservices architectures with Kubernetes orchestration [20], [21] and service meshes [7], are essential for achieving the required levels of scalability and resilience. Automated scaling and intelligent resource management are crucial for optimizing resource utilization and minimizing costs in a cloud environment [9], [10], [17]. Robust security implementations, informed by frameworks like the OWASP API Security Top 10 [1], [14], are paramount for protecting the platform and its users. Furthermore, structured API lifecycle management practices [4], [11], [21] and API governance [25] are essential for maintaining order and consistency as the API ecosystem grows. Finally, rigorous testing methodologies, including strategies for distributed environments [26], are vital for ensuring the reliability and performance of the platform. Practical insights from industry surveys highlight the growing adoption of API programs and the importance of key features in API management solutions [2]. The challenges and strategies associated with microservice API evolution, including versioning, message translation, and addressing backward compatibility, further emphasize the complexity of managing APIs in modern architectures [6].

B. Comparison with Academic Findings

While [19] provides a general comparison of API management tools, [12] offers a practical perspective on using Apigee, Postman, and Apiary, highlighting the trade-offs between features and ease of use. The observation in [12] of Apigee's complexity aligns with the challenges of implementing comprehensive API governance and lifecycle management discussed in [11]. The practical experience detailed in [12] emphasizes the importance of hands-on evaluation and considering factors beyond the theoretical capabilities of a platform.

C. Real-World Challenges and Considerations

[12] highlights practical challenges often encountered during API platform implementation, such as the complexity of setting up developer portals (in the case of Apigee) and the limitations of free tiers, factors not always addressed in academic studies. The need for creating 'API Products' in Apigee, as described in [12], adds another layer of complexity to the setup process, potentially increasing the overhead for smaller teams or projects. This practical perspective complements the theoretical discussions of API lifecycle management found in papers like [11] and [21], providing a more nuanced understanding of real-world implementation issues.

D. Limitations of the Evidence

While the included studies provide valuable insights, several limitations should be acknowledged. Some studies relied on simulated environments or limited datasets [27], [28], which may not fully represent the complexities of real-world, large-scale API platforms. The diversity of metrics and experimental setups across studies made direct quantitative comparisons challenging. Furthermore, the rapidly evolving nature of cloud technologies means that some findings may become outdated quickly. The limited focus on specific industry contexts within some research (e.g., [21] focusing on a specific enterprise use case) may also limit generalizability. Finally, many papers focused on individual aspects of API management rather than presenting a holistic view of building a complete platform, necessitating synthesis across multiple papers to draw comprehensive conclusions. Additionally, the comparison of API management solutions in [2] is limited to five vendors and could benefit from a broader scope.

E. Limitations of the Review Process

The review procedure itself contains inherent restrictions. Even with a thorough search technique, pertinent studies might have been overlooked. The overall results may have been impacted by publication bias, which favors research with favorable outcomes. Bias may also be introduced by the subjective character of some components of the quality assessment. Furthermore, resource constraints limited the ability to include and analyze all potentially relevant literature, especially grey literature.

F. Implications for Practice

This review offers practical guidance for practitioners building high-volume API platforms. It emphasizes the importance of:

- **Adopting a data-driven approach:** Using quantitative metrics and analysis to guide design and management decisions.
- **Embracing cloud-native technologies:** Leveraging microservices, containerization, and orchestration for scalability and resilience.
- **Prioritizing security:** Implementing robust authentication, authorization, and protection against common vulnerabilities.
- **Optimizing for cost-efficiency:** Employing automated scaling and intelligent resource management strategies.
- **Implementing structured API lifecycle management:** Adopting established best practices for API design, documentation, versioning, and testing.
- **Selecting appropriate API management solutions:** Considering factors like features, vendor support, and industry-specific requirements, informed by comparative analyses of available solutions [2].
- **Managing microservice API evolution:** Implementing effective versioning strategies, communication protocols, and testing procedures to address the challenges of backward compatibility and inter-service dependencies [6].

G. Implications for Policy

This review highlights the need for standardization and best practices in the API management domain. Industry bodies and standard organizations should collaborate to develop guidelines and frameworks for building secure, scalable, and cost-efficient API platforms.

H. Implications for Future Research

Several areas require further investigation:

- **Scalability benchmarking:** Developing standardized benchmarks and methodologies for evaluating the scalability of API platforms under realistic load conditions.
- **Performance optimization:** Exploring innovative techniques for optimizing API performance at scale, including the application of machine learning and AI.
- **Security:** Investigating new approaches for detecting and mitigating API security threats in dynamic cloud environments.
- **Cost optimization:** Developing more sophisticated cost models and resource management strategies for high-volume API platforms.
- **API lifecycle management:** Creating tools and frameworks for automating API lifecycle processes and improving collaboration among stakeholders.
- **Hybrid cloud and edge deployments:** Researching the challenges and best practices for deploying and managing API platforms in hybrid cloud and edge environments.
- **API evolution in microservices:** Further research is needed to explore automated tooling, standardized practices, and the impact of organizational structures on API evolution in microservices.

By addressing these research gaps, the API management community can further enhance the capabilities and robustness of high-volume API platforms, enabling the continued growth and evolution of the digital economy.

V. CONCLUSIONS

Building a cloud-based API management platform capable of handling billions of API calls per day requires a careful consideration of various interconnected factors. A cloud-native approach, leveraging microservices, containerization, and orchestration, is essential for achieving scalability and resilience. Automated resource management and intelligent scaling strategies are crucial for optimizing cost-efficiency. Robust security implementations, grounded in industry best practices, are paramount. Structured API lifecycle management and governance processes are vital for maintaining order and consistency. Finally, rigorous testing methodologies are essential for ensuring reliability and performance. By implementing the best practices and creative solutions found in this review, developers and architects may create scalable and reliable API platforms that can handle the demands of today's digital environment. Furthermore, successful API management requires consideration of practical aspects such as vendor selection, feature comparison, and addressing the challenges of

API evolution in microservices, as highlighted by the inclusion of [2], [6].

REFERENCES

- [1] F. Qazi, "Application programming interface (api) security in cloud applications," *EAI Endorsed Transactions on Cloud Systems*, vol. 7, no. 23, p. e1, Oct. 2023. [Online]. Available: <https://publications.eai.eu/index.php/cs/article/view/3011>
- [2] S. M. Ali and T. R. Soomro, "Comparative study of api management solutions," in *Proceedings of The 6th International Conference on Innovation in Science and Technology*, 2019.
- [3] Evolved Media and CITO Research, "Cloud-based api management."
- [4] N. Xie, "Strategic approaches to api design and management," in *Proceedings of the 6th International Conference on Computing and Data Science*, ser. Applied and Computational Engineering, A. Wang and R. Bauer, Eds., vol. 64. EWA Publishing, 2024, pp. 229–235. [Online]. Available: <https://www.confcds.org/>
- [5] S. Henning, "Scalability benchmarking of cloud-native applications applied to event-driven microservices," Ph.D. dissertation, Kiel, 2023. [Online]. Available: <https://doi.org/10.21941/kcss/2023/2>
- [6] A. Lercher, J. Glock, C. Macho, and M. Pinzger, "Microservice api evolution in practice: A study on strategies and challenges," *Journal of Systems and Software*, vol. 215, p. 112110, 2024. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0164121224001559>
- [7] S. Ambroszkiewicz and W. Bartyna, "A simple protocol to automate the executing, scaling, and reconfiguration of cloud-native apps," 2024. [Online]. Available: <https://arxiv.org/abs/2305.16329>
- [8] G. Ramirez, M. Lindemann, A. Birch, and I. Titov, "Cache & distil: Optimising api calls to large language models," 2023. [Online]. Available: <https://arxiv.org/abs/2310.13561>
- [9] P. Loncar and P. Loncar, "Scalable management of heterogeneous cloud resources based on evolution strategies algorithm," *IEEE Access*, vol. 10, pp. 68 778–68 791, 2022.
- [10] A. Amiri, U. Zdun, A. van Hoorn, and S. Dustdar, "Cost-aware multi-dimensional auto-scaling of service- and cloud-based dynamic routing to prevent system overload," in *2022 IEEE International Conference on Web Services (ICWS)*, 2022, pp. 379–384.
- [11] M. Overeem, M. Mathijssen, and S. Jansen, "Api-m-famm: A focus area maturity model for api management," *Information and Software Technology*, vol. 147, p. 106890, 2022. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0950584922000532>
- [12] S. Gupta, "A survey of api management platforms," May 2019. [Online]. Available: <https://medium.com/@shubhanshugupta/a-survey-of-api-management-platforms-348d80348a3f>
- [13] S. Henning and W. Hasselbring, "A configurable method for benchmarking scalability of cloud-native applications," *Empirical Software Engineering*, vol. 27, no. 6, p. 143, Aug. 2022. [Online]. Available: <https://doi.org/10.1007/s10664-022-10162-1>
- [14] M. Idris, I. Syarif, and I. Winarno, "Web application security education platform based on owasp api security project," *EMITTER International Journal of Engineering Technology*, vol. 10, no. 2, pp. 246–261, Dec. 2022. [Online]. Available: <https://emitter.pens.ac.id/index.php/emitter/article/view/705>
- [15] L. Nacson and A. Sandler, "Rapid penetration test for securing channel apis in hybrid cloud (dynamic defense of channel api)," Feb. 2023, pp. 95–104.
- [16] M. Mathijssen, M. Overeem, and S. Jansen, "Identification of practices and capabilities in api management: A systematic literature review," 2020. [Online]. Available: <https://arxiv.org/abs/2006.10481>
- [17] C. Tang, Y. Wang, B. Fan, B. Wang, S. Chen, Z. Qiu, C. Liang, J. Zhao, Y. Zhu, M. Chen, and Z. Hu, "Rethinking the cloudonomics of efficient i/o for data-intensive analytics applications," 2023. [Online]. Available: <https://arxiv.org/abs/2311.00156>
- [18] S. Kanungo, "Ai-driven resource management strategies for cloud computing systems, services, and applications," *World Journal of Advanced Engineering Technology and Sciences*, vol. 11, no. 2, pp. 559–566, 2024, last Modified: 2024-06-17T06:06+00:00, Publisher: World Journal of Advanced Engineering Technology and Sciences. [Online]. Available: <https://wjaets.com/content/ai-driven-resource-management-strategies-cloud-computing-systems-services-and-applications>

- [19] R. Gadia, R. Shah, S. Varshney, and V. Sawant, "A system on automated database and api (application programming interface) management," *International Journal for Research in Applied Science and Engineering Technology*, vol. 10, pp. 3226–3234, Apr. 2022.
- [20] M. Usman, D. Badampudi, C. Smith, and H. Nayak, "An ecosystem for the large-scale reuse of microservices in a cloud-native context," *IEEE Software*, vol. 39, no. 5, pp. 68–75, 2022.
- [21] I. Di Natali, "Deploying a scalable API management platform in an enterprise Kubernetes-based environment," laurea, Politecnico di Torino, Oct. 2020. [Online]. Available: <https://webthesis.biblio.polito.it/15945/>
- [22] B. Rob, "Api scaling effortlessly: Azure api management's seamless scalability," accessed: 2025-01-11. [Online]. Available: <https://www.dootrix.com/insights/azure-api-managements-seamless-scalability>
- [23] Z. Samira, Y. W. Weldegeorgise, O. S. Osundare, H. O. Ekpobimi, and R. C. Kandekere, "Api management and cloud integration model for smes," *Magna Scientia Advanced Research and Reviews*, vol. 12, no. 1, pp. 078–099, Sep. 2024. [Online]. Available: <https://magnascientiapub.com/journals/msarr/node/885>
- [24] A. El Malki, U. Zdun, and C. Pautasso, "Impact of api rate limit on reliability of microservices-based architectures," in *2022 IEEE International Conference on Service-Oriented System Engineering (SOSE)*, 2022, pp. 19–28.
- [25] C. Krintz, H. Jayathilaka, S. Dimopoulos, A. Pucher, R. Wolski, and T. Bultan, "Cloud platform support for api governance," in *2014 IEEE International Conference on Cloud Engineering*, 2014, pp. 615–618.
- [26] S. Ackerman, S. Choudhury, N. Desai, E. Farchi, D. Gisolfi, A. Hicks, S. Route, and D. Saha, "Towards api testing across cloud and edge," 2021. [Online]. Available: <https://arxiv.org/abs/2109.02540>
- [27] W. Funika, P. Koperek, and J. Kitowski, "Automated cloud resources provisioning with the use of the proximal policy optimization," *The Journal of Supercomputing*, vol. 79, no. 6, pp. 6674–6704, Apr. 2023. [Online]. Available: <https://link.springer.com/10.1007/s11227-022-04924-3>
- [28] W. Dawoud, I. Takouna, and C. Meinel, "Scalability and performance management of internet applications in the cloud," in *Communication Infrastructures for Cloud Computing*, H. T. Mouftah and B. Kantarci, Eds. Hershey, PA, USA: IGI Global, 2014, pp. 434–464. [Online]. Available: <https://services.igi-global.com/resolvedoi/resolve.aspx?doi=10.4018/978-1-4666-4522-6.ch019>