

INTELLIGENT SCALING OF DISTRIBUTED PAYMENT SYSTEMS: APPROACHES TO REDUCING INFRASTRUCTURE COSTS IN HIGH-LOAD ECONOMIES

A.M. Kovalenko, *Bachelor's Degree*

Tomsk State University of Control Systems and Radioelectronics
(Russia, Tomsk)

DOI:10.24412/2411-0450-2025-7-74-79

Abstract. *The article analyzes architectural principles of intelligent scaling as implemented in distributed payment systems under conditions of high transactional load. It examines the use of micro-service and event-driven approaches, predictive analytics, and automated load adaptation mechanisms to enhance the resilience and flexibility of financial infrastructure. It also explores the impact of scaling strategies on infrastructure costs, evaluates resource optimization methods, including serverless models, network call management, and data lifecycle control, and assesses various risks associated with over-optimization, along with key metrics for measuring scaling efficiency.*

Keywords: *intelligent scaling, distributed payment systems, high-load systems, infrastructure optimization, load prediction, serverless computing, economic efficiency.*

Distributed payment systems play a pivotal role in the modern digital economy by ensuring the continuity and reliability of financial operations under high user traffic. As the volume of transactions grows and financial services become increasingly complex, the pressure on computational infrastructure intensifies. Its demands not only real-time system stability from developers but also the ability to scale efficiently without a proportional rise in operational costs. In high-load economic environments marked by market volatility, the efficiency of resource management emerges as a critical factor for maintaining infrastructure resilience.

One of the most promising directions for the evolution of such systems is the integration of intelligent scaling approaches, which rely on automated adaptation to dynamic workloads. The use of predictive models, real-time traffic management, and hybrid cloud architectures enables both improved fault tolerance and reduced expenses associated with maintaining excessive capacity. However, the implementation of these solutions requires a careful assessment of associated risks, cost-effectiveness, and technological constraints.

From a theoretical perspective, intelligent scaling encompasses a combination of algorithmic mechanisms, architectural strategies, and economic models aimed at balancing system performance with infrastructure expenditures. The goal of this study is to analyze current approaches to scaling distributed payment systems, with a

particular focus on reducing infrastructure costs while maintaining service reliability and quality. The paper explores architectural principles for building scalable payment infrastructures, emphasizing the application of intelligent load management algorithms. It also examines practical methods for optimizing infrastructure expenditures, outlines their use in high-load economic contexts, and assesses the potential risks and limitations of such strategies.

Architectural approaches to intelligent scaling in distributed payment systems

With the rapid growth of digital transactions, distributed payment systems have become the foundation of both retail and corporate financial infrastructure. Ensuring resilience under peak loads, typical in economies with high transaction density, makes scalability a critical design concern. Traditional scaling methods, such as over-provisioning or manual infrastructure management, are proving increasingly inadequate due to their rigidity and inefficiency. These approaches are being replaced by intelligent solutions that leverage automation, predictive analytics, and real-time adaptive behavior of system components.

Intelligent scaling comprises a set of architectural and algorithmic mechanisms that enable a system to dynamically adjust to fluctuating input traffic. In distributed payment environments, where workload patterns vary across hours, days, and seasons, such adaptability is essential. In practice, multiple architectural models are em-

ployed to address scalability at different layers of the payment infrastructure, from transaction processing to external API integration and user session management.

Among these models, the microservices architecture is the most widely adopted. It decomposes the system into loosely coupled, independent

components, such as authentication services, fee calculation modules, fraud detection filters, and gateways for transaction intake and dispatch [1]. The deployment of microservices typically relies on container orchestration platforms such as Kubernetes (fig. 1).

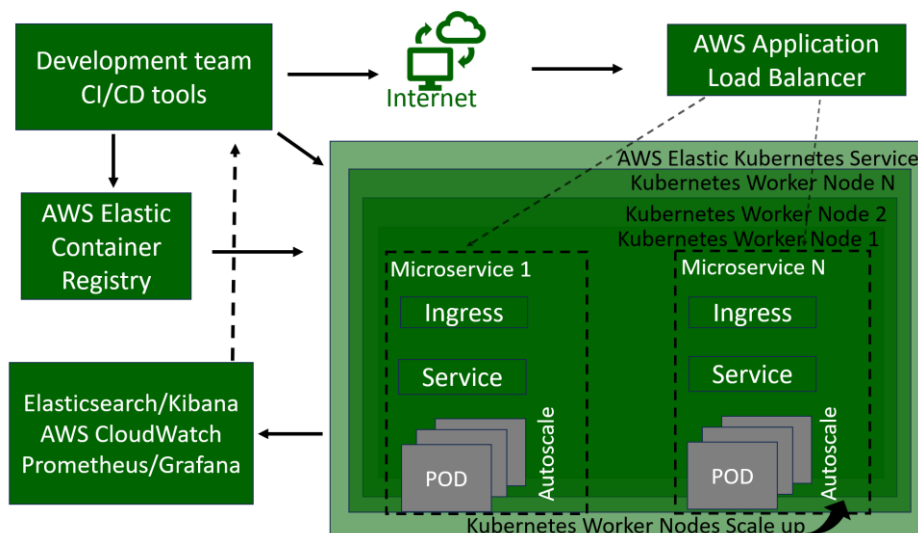


Fig. 1. Scalable microservices architecture using AWS EKS and Kubernetes

This modular design facilitates targeted scaling of specific services in response to demand and allows intelligent algorithms to be integrated precisely where they deliver the greatest benefit. Machine learning models embedded within anomaly detection services can forecast transaction spikes and trigger proactive scaling before bottlenecks arise.

In parallel with microservices architectures, event-driven approaches are increasingly employed to enable scaling in response to runtime occurrences within payment systems [2]. Events such as sudden surges in transaction volume, prolonged response times, or rising failure rates act as triggers for scaling scenarios. This facilitates proactive rather than reactive behavior, helping the system prevent performance degradation.

A critical component of intelligent scaling lies in its integration with external monitoring platforms and predictive analytics systems. Platforms

that collect real-time metrics serve as a foundation for constructing behavioral load profiles [3]. These profiles inform machine learning models capable of forecasting demand fluctuations, enabling the system to anticipate and prepare for load changes. Unlike traditional threshold-based scaling, intelligent strategies incorporate contextual factors such as transaction types, historical usage patterns, holidays, weekends, and market-specific events. Such predictive models, when deployed within modular services like fraud detection or dynamic pricing engines, benefit from algorithmic selection based on both accuracy and computational efficiency. As demonstrated earlier by comparative evaluations using Root Mean Square Error (RMSE) and Mean Absolute Error (MAE) metrics [4], algorithms such as Gradient Boosting exhibit superior predictive performance in load forecasting tasks (table 1).

Table 1. Algorithms performance evaluation

ML algorithm	RMSE, Mbit	MAE, Mbit	R ² , %	Wall time of learning, s	Advantages	Disadvantages
Bagging	2,59	1,66	50,8	58,1	Good quality of prediction. Minimizing model's overfit.	Long time of training.
Random Forest	3,38	2,19	34,2	172,0	Universality.	Low metrics compared to Gradient Boosting. Long time of training.
Gradient Boosting	2,18	1,43	60,2	24,6	High quality of prediction.	Long time of training.
Linear Regression	2,64	1,71	49,7	2,5	Simplicity of the algorithm and fast training.	Dependence on feature linearity and scaling.
Bayesian Regression	2,25	1,49	49,7	1,1	Adapts to the data at hand.	Inference of the model can be time consuming.
Huber Regression	2,43	1,62	26,8	4,0	More efficient to use on data with small number of samples.	Dependence on feature linearity.

Infrastructure orchestrators and management platforms, such as Kubernetes, Nomad, and OpenShift, form an integral part of such architectures. These platforms support container- and service-level autoscaling while allowing for the definition of custom scaling policies. A system can be configured to prioritize the scaling of international payment services over local transactions based on their business-criticality. This policy-driven adaptability ensures rational resource allocation and reduces the overhead of unnecessary deployments.

A key dimension of intelligent scaling in distributed systems is the localization of workload across geographically dispersed infrastructure. When serving a substantial number of users across multiple regions, load distribution can be optimized through geo-aware scaling strategies that automatically adjust the capacity of regional nodes. In this context, architectural components such as geo-aware load balancers and distributed databases with active replication play a crucial role.

In summary, intelligent scaling in distributed payment systems introduces a new level of adaptability and architectural efficiency. It integrates microservices modularity, event-driven responsiveness, predictive analytics, and automated resource orchestration. The success of such systems hinges on the coordinated interaction of architectural layers, orchestration tools, and monitoring platforms, as well as the system's ability to anticipate load patterns not only at the infrastructure level but also in terms of user behavior.

Reducing infrastructure costs: methods, risks, and efficiency

Despite its technological maturity, intelligent scaling in distributed payment systems remains highly sensitive to the cost of infrastructure maintenance. In high-load economic environments, where millions of transactions occur daily, any excess in computing resource consumption, redundant API calls, or inefficient data storage directly translates into increased operational expenses. As a result, beyond technological adaptability, economic efficiency becomes a critical consideration in scaling strategies, positioning intelligent scaling not only as a tool for resilience but also as a lever for cost optimization.

One of the most effective methods for reducing expenses in such systems is the optimization of scaling policies based on time-of-day usage patterns, transaction types, and anticipated demand. The application of demand-driven adaptation policies allows for dynamic adjustment of active instances, containers, or compute nodes in low-traffic windows, without compromising service-level agreements (SLA) [5].

A significant opportunity for cost reduction lies in the adoption of serverless computing (Function-as-a-Service) for handling irregular or event-driven tasks, such as notification delivery, report generation, KYC document verification, or one-off external API calls. A typical serverless architecture consists of multiple abstraction layers, each designed to isolate responsibilities and improve efficiency. The architecture also integrates with development and monitoring environments through a management API (fig. 2).

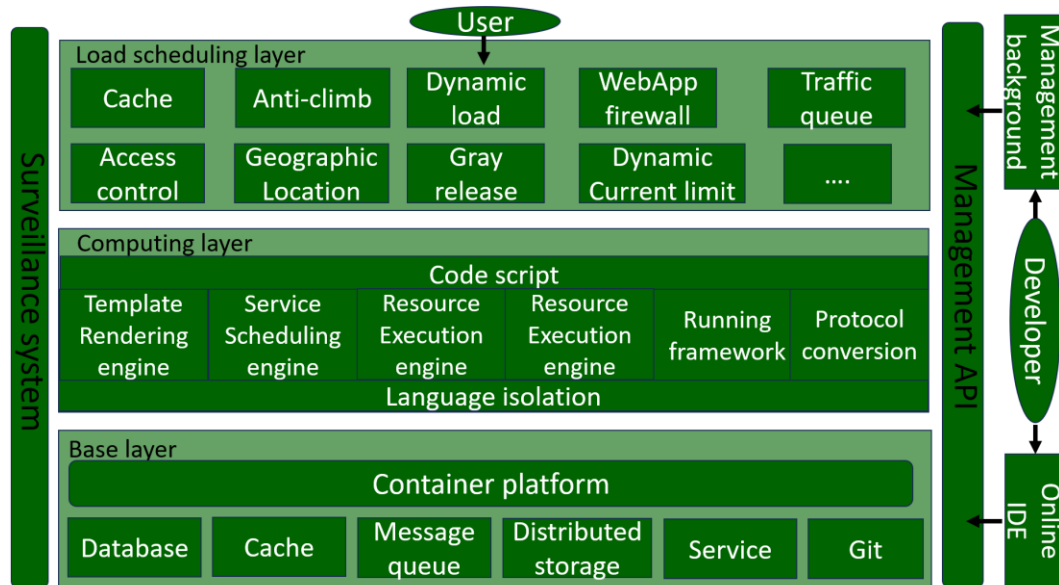


Fig. 2. Architecture design of the Serverless architecture [6]

Transitioning from continuously running services to a pay-per-execution model eliminates the costs associated with idle infrastructure. This approach is particularly beneficial for high-load systems that process a large volume of transient events, which do not require persistent memory or server-side presence.

Beyond compute-related savings, optimizing network interactions and external API usage also plays a critical role in cost efficiency. Commercial payment gateways, fraud detection systems,

currency conversion APIs, and banking interfaces often impose usage-based pricing models. Techniques such as transaction batching and caching of frequently accessed metadata (e.g., exchange rates or card limits) help reduce the number of outbound requests and lower third-party service costs. Asynchronous call aggregation also reduces network overhead and third-party billing costs. Multiple components send individual requests to a shared asynchronous aggregator, which groups them into a single outbound API call (fig. 3).

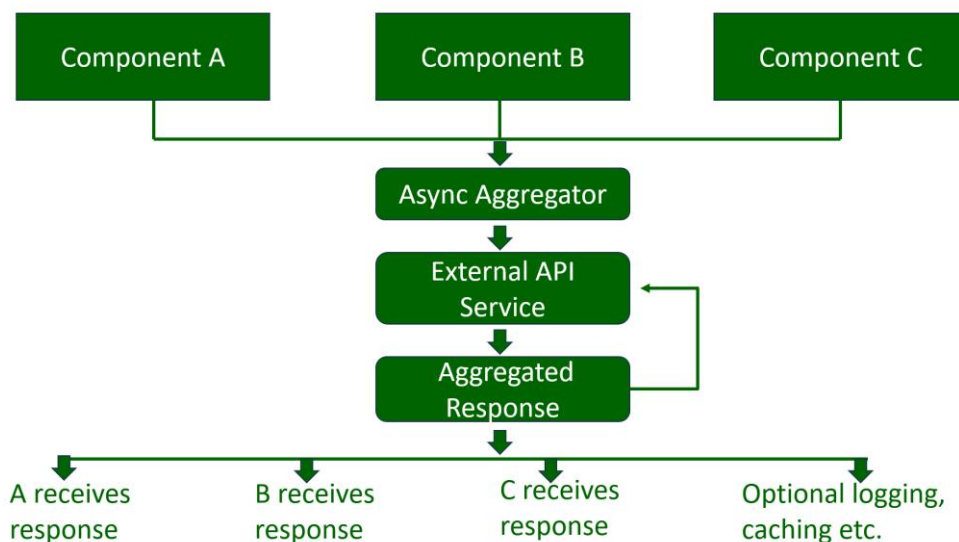


Fig. 3. Asynchronous API call aggregation architecture

Data storage and processing costs also represent a critical consideration. High-load payment systems generate vast volumes of logs, metrics, temporary tables, and redundant records. In this

context, implementing data lifecycle management policies becomes a highly effective strategy. These include automated deletion of outdated logs, storage compression, and the use of low-

cost archival storage tiers (cold storage) for analytical backups. Beyond direct cost savings, such practices reduce the load on indexing and analytics subsystems, enhance real-time responsiveness, and decrease the need for storage layer scaling.

The adoption of these optimization strategies requires systematic evaluation of their economic viability. To this end, several key performance indicators (KPI) are commonly employed (table 2).

Table 2. Common KPI for Scaling Efficiency [7, 8]

KPI	Description
Cost per transaction (CPT)	The average infrastructure cost associated with processing a single transaction.
Resource utilization efficiency (RUE)	The ratio of actual resource usage to allocated capacity.
Idle resource rate (IRR)	The proportion of provisioned but unused resources over time.
Scaling cost impact (SCI)	The incremental cost incurred with each additional unit of scaling.

These metrics help identify inefficiencies within the scalable architecture and serve as a basis for refining scaling strategies. If the SCI increases exponentially rather than linearly during horizontal scaling, it may indicate issues with load balancing or communication overhead between system modules.

Despite the apparent benefits of intelligent cost optimization, there are inherent risks associated with over-optimization. One such risk is the latent degradation of service quality when resources are aggressively reduced. Relying on the minimum number of instances can result in increased processing latency during unexpected traffic surges. Another risk involves overreliance on external analytics services for workload forecasting. Inaccurate predictions may lead to premature downscaling or suboptimal resource allocation.

Mitigating these risks requires the implementation of hybrid scaling strategies, where multiple rules and models are combined. A common approach involves maintaining a static baseline of resources while enabling dynamic scaling based on forecast-driven metrics. This balance between efficiency and reliability allows for more stable system behavior under variable load conditions. Continuous calibration of these strategies is essential and should be informed by retrospective analysis, comparing forecast models with actual

usage patterns, and revising thresholds and coefficients in light of seasonal or market-driven fluctuations.

Reducing infrastructure costs in distributed payment systems necessitates a comprehensive approach that integrates economic rationality, predictive analytics, and technological flexibility [9]. There is no one-size-fits-all solution, only an adaptive, continuously refined strategy that reflects the specific dynamics of the system, user behavior, and broader market environment. The effectiveness of such strategies directly impacts not only operational expenditures but also an organization's competitiveness in a high-load economy.

Conclusion

Intelligent scaling of distributed payment systems constitutes a multifaceted mechanism that integrates architectural, algorithmic, and managerial components. These strategies are designed to enhance adaptability and reduce operational overhead under conditions of high transactional load. Efficient resource allocation, the use of predictive modeling, and the implementation of automated scaling policies collectively enable resilient and cost-effective payment infrastructure. The successful deployment of such approaches requires systematic analysis, careful risk assessment, and ongoing refinement of strategies based on performance metrics and economic feasibility.

References

1. Li S., Zhang H., Jia Z., Zhong C., Zhang C., Shan Z., Shen J., Babar M.A. Understanding and addressing quality attributes of microservices architecture: A Systematic literature review // *Information and software technology*. – 2021. – Vol. 131. – № 106449.
2. Cabane H., Farias K. On the impact of event-driven architecture on performance: An exploratory study // *Future Generation Computer Systems*. – 2024. – Vol. 153. – P. 52-69.

3. Smirnov A. Monitoring and logging in distributed systems: application of OpenTelemetry and the ELK stack // Universum: technical sciences: electronic scientific journal. – 2025. – № 3(132). – P. 30-33.
4. Alekseeva D., Stepanov N., Veprev A., Sharapova A., Lohan E.S., Ometov A. Comparison of machine learning techniques applied to traffic prediction of real wireless network // IEEE Access. – 2021. – Vol. 9. – P. 159495-514.
5. Badshah A., Jalal A., Farooq U., Rehman G.U., Band S.S., Iwendi C. Service level agreement monitoring as a service: an independent monitoring service for service level agreements in clouds // Big Data. – 2023. – Vol. 11. – № 5. – P. 339-54.
6. Jiang L., Pei Y., Zhao J. Overview of serverless architecture research // InJournal of Physics: Conference Series. – 2020. – Vol. 1453. – № 1. – P. 012119.
7. Smoliarchuk V. Management strategies and models in conditions of high market volatility: application of Agile and Scrum in industrial business // The Eurasian Union of Scientists. Series: Economic and legal sciences. – 2025. – Vol. 1. – № 1(126). – P. 16-19.
8. Ramachandran K.K. Optimizing IT Performance: A Comprehensive analysis of Resource Efficiency // International Journal of Marketing and Human Resource Management (IJMHRM). – 2023. – Vol. 14. – № 3. – P. 12-29.
9. Arkhipov V. Analysis of strategic approaches to operational risk management in the context of digital transformation // Professional Bulletin: Economics and Management. – 2024. – № 1/2024. – P. 11-15.

ИНТЕЛЛЕКТУАЛЬНОЕ МАСШТАБИРОВАНИЕ РАСПРЕДЕЛЕННЫХ ПЛАТЕЖНЫХ СИСТЕМ: ПОДХОДЫ К СНИЖЕНИЮ ИНФРАСТРУКТУРНЫХ ЗАТРАТ В УСЛОВИЯХ ВЫСОКОНАГРУЖЕННОЙ ЭКОНОМИКИ

А.М. Коваленко, *бакалавр*

Томский государственный университет систем управления и радиоэлектроники
(Россия, г. Томск)

***Аннотация.** В статье анализируются архитектурные принципы интеллектуального масштабирования, реализуемые в распределенных платежных системах на фоне высокой транзакционной нагрузки. Изучается применение микросервисных и событийно-ориентированных подходов, предиктивной аналитики и автоматизированных механизмов адаптации нагрузки для обеспечения устойчивости и гибкости финансовой инфраструктуры. Также исследуется влияние стратегий масштабирования на инфраструктурные затраты, анализируются методы оптимизации ресурсов, включая бессерверные модели, управление сетевыми вызовами и жизненным циклом данных, а также оцениваются различные риски, связанные с чрезмерной оптимизацией, и метрики эффективности масштабируемых решений.*

***Ключевые слова:** интеллектуальное масштабирование, распределенные платежные системы, высоконагруженные системы, оптимизация инфраструктуры, прогнозирование нагрузки, бессерверные вычисления, экономическая эффективность.*