



#### OPEN ACCESS

SUBMITTED 21 February 2025

ACCEPTED 19 March 2025

PUBLISHED 30 April 2025

VOLUME Vol.07 Issue 04 2025

#### CITATION

Klimkov Ilia. (2025). Modeling Scaling Strategies for Shopify Platforms in International Expansion. The American Journal of Management and Economics Innovations, 7(04), 103–109.  
<https://doi.org/10.37547/tajmei/Volume07Issue04-13>

#### COPYRIGHT

© 2025 Original content from this work may be used under the terms of the creative commons attributes 4.0 License.

# Modeling Scaling Strategies for Shopify Platforms in International Expansion

**Klimkov Ilia**

E-commerce Founder & CEO VAOVAC

E-commerce Founder & CEO OXYFIT

San Diego, California, USA

**Abstract:** This article explores technological solutions aimed at scaling Shopify platforms as they expand into international markets. The relevance of the topic is driven by the growing number of online stores that require high service availability, even under rapidly increasing traffic. The novelty of the study lies in the comprehensive analysis of architectural strategies used to distribute load, ensure fault tolerance, and reduce latency during global operations. The paper outlines key principles of flexible scaling through client isolation (pods), geographic distribution of servers across regions, the use of content delivery networks (CDNs) for accelerated content delivery, and load testing approaches designed to simulate peak scenarios such as flash sales. Additionally, it examines methods for balancing and dynamically reallocating resources to protect the system from failure. The study aims to offer practical recommendations for those seeking to maintain platform stability amid international expansion. To achieve this, the article applies comparative analysis of architectural solutions and model's scalability potential. The research also draws on the experiences of Shopify engineers to deepen the understanding of practical challenges in global growth. The conclusion emphasizes the importance of constant monitoring of throughput and strategic distribution of key services across global data centers. This article will be useful to SaaS professionals, developers, and analysts responsible for planning large-scale e-commerce projects.

**Keywords:** Shopify, scaling, global traffic, pods, geo-distribution, international commerce, CDN, load testing, SaaS, architecture.

**Introduction:** Shopify is one of the largest e-commerce platforms in the world, serving millions of online stores

across numerous countries. Its success is largely attributed to its ability to scale effectively in response to growing demand and global expansion. From a technical perspective, Shopify originally operated as a monolithic application; however, to ensure system resilience and performance, the company was compelled to redesign its architecture. Today, Shopify is a sophisticated multi-tenant system distributed across numerous data centers.

This article investigates the strategies Shopify uses to scale its platform at the international level and explores methods for modeling these strategies. The objective is to analyze how Shopify isolates load across clients (via its pods architecture), ensures geographic distribution (via multi-region deployments), and performs scalability testing and monitoring (e.g., load testing and flash sale simulations). Special attention is given to how these strategies allow Shopify to maintain performance across different countries and adapt to the challenges of global expansion, including localization, time-zone-specific demand peaks, and resiliency requirements.

## MATERIALS AND METHODS

In preparing this article, a range of foundational works were reviewed. B. de Water [1; 2] provides in-depth insights into the design of payment systems and architectures built to handle multimillion-scale demand surges. X. Denis [3] describes the pods approach, which eliminates the risk of resource contention among stores within a shared environment. H. Khalid [4] analyzes

Shopify's migration to Vitess as a means of enabling horizontal scaling without disrupting workflows. P. Madan [5] discusses shard balancing and techniques for moving client data without downtime—even at terabyte scale. A. Rodukov [6] presents Shopify's global expansion strategy with a focus on distributed data centers.

The literature review revealed a need to synthesize and extend current findings through the modeling of a holistic scaling strategy. The novelty of this article lies in its unification and systematization of Shopify's international scaling practices—accounting for all aspects of the platform's infrastructure, from load isolation and geographic distribution to modular architecture and traffic simulation.

## RESULTS

One of Shopify's key scaling solutions is the implementation of so-called pods. A pod, in the context of Shopify, is a fully isolated segment of the platform that contains its own database cluster and other storage systems [3]. Initially, as load increased, Shopify turned to database sharding but encountered a situation where the failure of one shard affected the entire system. As a result, the company decided to separate resources more radically: each pod serves a specific group of stores and has its own database (MySQL), cache (Redis, Memcached), and background job queues (see Fig. 1).



Figure 1. Pod in the context of Shopify [3]

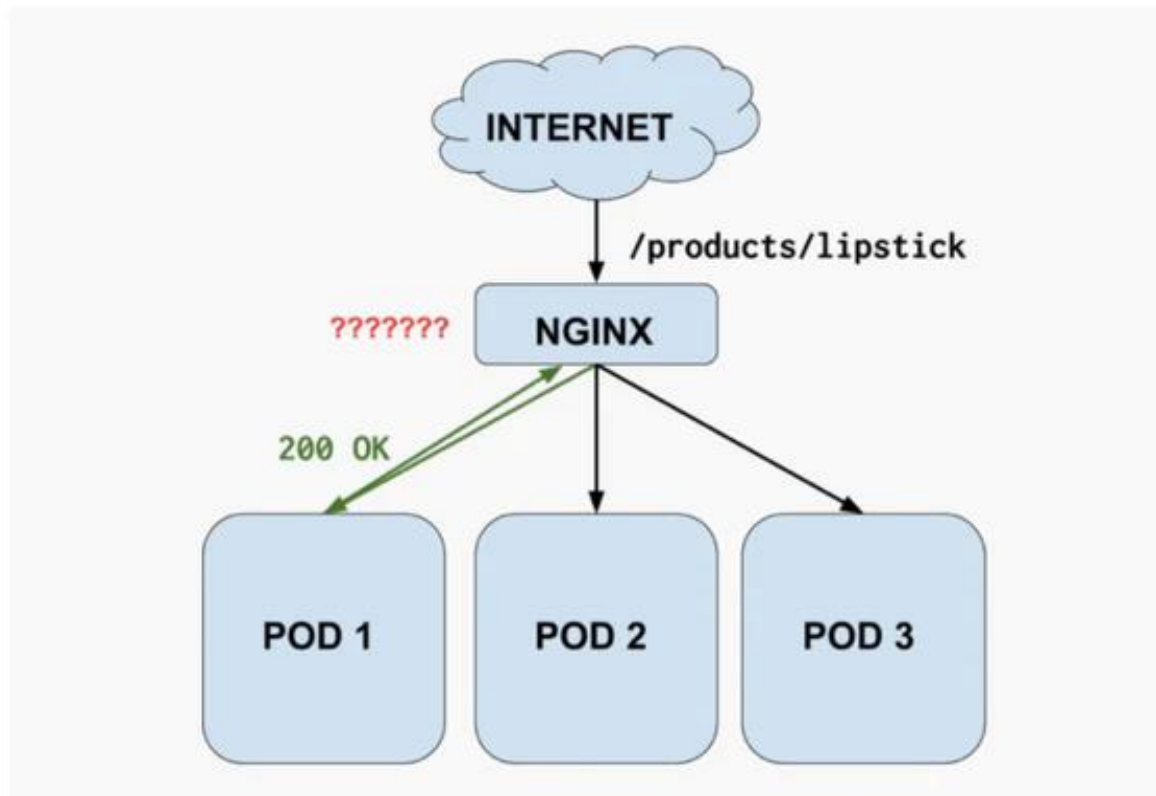
No requests cross between pods; each store is strictly tied to its own pod. This provides two advantages: first, horizontal scalability—the company can simply add new pods as the number of clients grows (similar to adding new nodes to a system); second, fault tolerance—issues in one pod do not affect the operation of stores in other pods. This architecture allowed Shopify to go beyond the limitations of a single monolithic database server. Xavier Denis, a Shopify engineer, noted that pods provide independence and eliminate mutual resource impact: “Adding a new pod does not cause unexpected load on existing ones” [3].

The pods approach means the platform remains multi-tenant, but tenants are grouped into clusters. This is important for international growth, as Shopify can optimally distribute stores across pods: for example, large stores generating heavy traffic can be allocated to separate pods so as not to interfere with smaller ones. In the case of a forecasted sharp traffic spike (e.g., Black Friday sales), Shopify can temporarily move an especially large store to a dedicated pod (or even several pods) [2]. This flexibility protects other clients from the “noisy neighbor” effect, when one store consumes all the resources. Shopify’s architecture includes a tool called Shop Mover, which allows stores to be moved between pods for load balancing without downtime [5]. For example, if several fast-growing merchants are concentrated on one pod, some of them are “moved” to another pod to balance the load. These solutions, modeled and tested within the company, ensure service stability as it scales.

For international expansion, it is crucial to place

infrastructure closer to users in different regions. Shopify initially operated from a single data center but eventually moved to a multi-data center deployment. Today, Shopify operates several major regions (North America, Europe, Asia, etc.), each of which runs a set of pods. As noted by Bart de Water (Shopify), each pod is tied to a specific region and actively runs in one data center, but also exists in at least two centers—the primary and the backup [2]. That is, each pod has a “partner” in another region where data is replicated in real time. If the primary data center fails, a pod failover occurs—a switch to the backup center (Shopify developed a special Pod Mover mechanism at the data center level for this purpose). This architecture ensures both geographic proximity (each merchant is served from the nearest center) and global resilience (if an entire region goes offline, stores continue to operate from the backup).

For international growth, it is also important that Shopify can open new regions by deploying the required number of pods and other components in a new data center. For example, when entering the Southeast Asian market, the company can launch pods in Singapore or another region, configure data replication. Shopify’s network system (including its own load balancer Sorting Hat) determines which region to route a client request to, based on the store domain and routing rules. “Sorting Hat” is the name of the internal routing system, which, at the incoming request stage, determines which pod (and region) the requested store belongs to, adds the appropriate header, and then sends the request directly to the required center (see Fig. 2) [3].



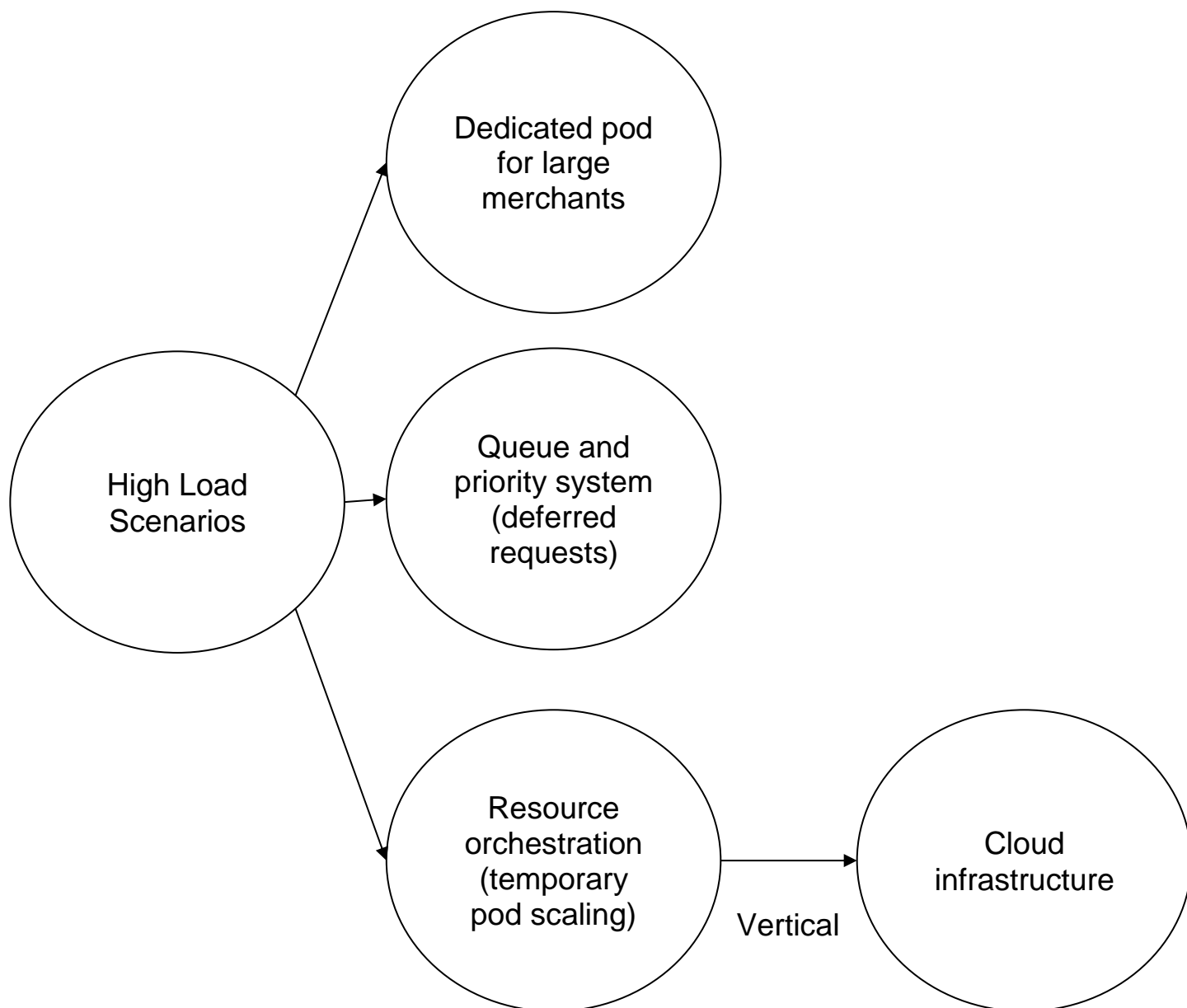
**Thanks to this, even global distribution remains transparent: customers simply visit `shopname.myshopify.com`, and the infrastructure automatically routes them to the appropriate nearby pod.**

In addition to server-side distribution, Shopify actively uses content delivery networks (CDNs) to accelerate the delivery of static resources (product images, theme files, etc.) worldwide. Shopify partners with CDN providers (such as Cloudflare), allowing content to be cached closer to end users. For example, a customer in Europe receives images and scripts from European CDN nodes, even if the store's main server is located in North America. This reduces response time and makes store performance faster for international visitors. As part of the scaling strategy, the CDN is a critical component: it reduces the load on core servers during global traffic peaks. According to the company, the global CDN network shortens page load times through file compression and geo-distribution [6]. For merchants, this means higher conversion rates in foreign markets; for Shopify, the ability to serve more concurrent users without expanding server capacity (a significant portion of traffic is served from cache). Thus, the combined solution—pods + multi-region + CDN—has become the foundation of Shopify's internationalization strategy.

Scaling modeling is an integral part of a platform like Shopify. The company places great emphasis on

automated load testing. For instance, large-scale load tests of key components are conducted weekly, simulating traffic volumes comparable to or exceeding Black Friday peaks [1]. During these tests, Service Level Objectives (SLOs)—acceptable response times and resource utilization—are monitored. If any component fails to meet the defined load threshold, the team receives an immediate alert and must improve it before the next testing round. This approach ("push the system until it breaks") allows potential bottlenecks to be identified in advance. For example, ahead of the next year, the team may increase the number of pods or optimize queries if the tests indicate proximity to the limit. Bart de Water explained that Shopify uses two types of load tests: regular verification of protective mechanisms (ensuring rate limiters, queues, etc. function properly), and breakpoint scale tests—intentionally pushing the system to failure to determine its limits and eliminate weak points [1]. Forecasting based on such tests supports capacity planning: determining how many new servers or pods will be needed if the number of stores grows by X, or if traffic increases by Y.

A particular challenge is when a single store suddenly receives a massive traffic spike (e.g., a well-known brand launches a limited product line). For Shopify, this poses a risk—even with pod isolation, a large surge within one pod can exhaust its resources. Several strategies are used for such scenarios (see Fig. 3):



**Figure 3. Load Management Strategies (compiled by the author based on [2])**

As mentioned, major merchants can be placed on a dedicated pod. Shopify has developed a queuing and prioritization system to manage extreme surges in a controlled way. For example, requests to an overloaded store can be placed in a queue instead of affecting the performance of other stores on the same pod. Resource orchestration: in a cloud environment (Shopify is gradually migrating parts of its infrastructure to the cloud), it is possible to temporarily scale a

specific pod up during the event.

However, the most reliable approach is forecasting and pre-scaling. Shopify works closely with its large merchants and tracks events via marketing (e.g., global sales). Knowing the promotion schedule, the team can run an unscheduled load test for that specific store or launch a copy of it in reserve infrastructure to confirm capacity. If uncertainty remains, a dedicated pod may be temporarily allocated to that store. Ultimately, the

goal is to prevent a spike from bringing the system down or slowing it. According to Shopify, during Black Friday Cyber Monday 2021, the platform processed a peak of 32 million requests per minute without major downtime [2]. This was made possible by proper scaling and proactive load modeling.

Scaling is not only about traffic—it is also about architectural flexibility across countries. Shopify scales the platform functionally: adding support for various payment gateways, tax systems, languages, and currencies. From an architectural standpoint, this is implemented through modularity: external integrations (e.g., local payment systems) are deployed as services that can be hosted closer to the respective regions. For example, for stores in Europe, the EU payment processing service can reside in a European data center, reducing latency when communicating with European banks.

## DISCUSSION

Shopify's scaling strategies demonstrate a balance between proven practices and innovation. The pods architecture—relatively uncommon in the industry due to the significant restructuring it demands from a monolithic system—has proven effective in Shopify's case. By isolating groups of stores, the company was able to scale performance linearly: adding a new pod adds new capacity. Modeling showed that without pods, the platform would have hit a ceiling at a certain database size (as operations on a single shard would block parts of the functionality). However, pods are not a panacea: they require administration (e.g., ensuring even load distribution), and they complicate some global functions (such as aggregate analytics across all stores, which now have to query every pod separately). Shopify likely addressed these issues by introducing higher-level components above pods (centralized search, indexes).

International scalability is a clear example of how technical solutions serve business objectives. Regional distribution enabled Shopify to provide low latency and high fault tolerance globally, which is critical for competitiveness. Notably, Shopify long relied on its own servers, but for international expansion, it partially transitioned to cloud partners—opening new regions within cloud platforms rather than building every data center from scratch. This is a typical trade-off: for the

sake of faster growth, using existing infrastructure (e.g., GCP or AWS) is sometimes more efficient—a path Shopify followed for some components (it is known, for example, that parts of Shopify's databases were moved to GCP using Vitess to simplify MySQL scaling) [4].

What stands out is how Shopify integrated continuous scalability testing into its development process. This reflects a mature approach: scalability is treated not as a one-time optimization (“before a major event”), but as an ongoing aspect of quality assurance. This mindset is worth emulating by other global platforms. Many companies focus solely on functional testing, overlooking a non-functional (but critical) requirement—withstanding future load. Shopify effectively implemented a “capacity planning in CI” model, treating infrastructure capacity checks as part of CI/CD.

Shopify's approach to flash sales is particularly notable from a strategy modeling perspective. The company developed a combination of technologies: from a virtual customer queuing system (to throttle simultaneous cart and payment activity) to excess compute capacity on standby. One could say Shopify models not only system-level load but also user flow behavior. When millions of users add items to carts at once, the system does not serve them all instantly but spreads the load over time (possibly showing waiting pages)—while ensuring that all are eventually served. The result is a compromise: users get a reliable, if not instantaneous, experience, and the platform maintains controlled load.

Despite Shopify's highly resilient and scalable architecture, it has its limitations. For instance, while pods can be added, the infrastructure that manages them (load balancers, routing tables) must support a growing number of records. This presents a challenge at scale: systems like Sorting Hat must continue to function reliably with hundreds or thousands of pods. A possible solution is a hierarchical model (grouping pods into clusters). In addition, expanding functionality (e.g., introducing large new services such as Shopify Functions) may increase load on existing databases. Shopify has likely already begun adopting microservices for specific modules—for example, separating payments, search, and notifications into standalone services (Bart de Water mentioned breaking the monolith into components) [2]. This means that scaling

considers not only the horizontal dimension of pods but also the vertical separation of services.

Analyzing Shopify's experience, similar principles can be recommended to other SaaS platforms aiming for global scale. First, logical tenant segmentation (whether through pods, shards, or multi-instance setups) is a fundamental mechanism for SaaS scalability—it enables horizontal growth rather than vertical. Second, geographic replication and CDN integration are de facto standards for global systems, without which users in distant regions would experience latency. And finally, continuous testing and refinement of the scalability model is what distinguishes a proactive engineering culture.

## CONCLUSION

The analysis confirms the effectiveness of Shopify's architectural solutions, including the pod structure, geographic distribution of servers, and the use of CDNs to ensure performance and fault tolerance during the platform's global expansion. The company's strategy of horizontal scaling combined with load modeling ensures stability under extreme traffic spikes and enables efficient resource management on an international scale.

However, the examined solutions have certain limitations, particularly regarding the management of a large number of pods and increased system complexity as functional expansion continues. Future research may focus on analyzing hierarchical models for pod management and assessing the impact of microservices architecture on the scalability of large-scale SaaS

platforms.

## REFERENCES

de Water, B. (2022). 10 tips for building resilient payment systems. Shopify Engineering. <https://shopify.engineering/building-resilient-payment-systems> (accessed April 11, 2025)

de Water, B. (n.d.). Shopify's architecture to handle the world's biggest flash sales [Conference presentation]. InfoQ. <https://www.infoq.com/presentations/shopify-architecture-flash-sale/> (accessed April 12, 2025)

Denis, X. (2018). A pods architecture to allow Shopify to scale. Shopify Engineering. <https://shopify.engineering/a-pods-architecture-to-allow-shopify-to-scale> (accessed April 10, 2025)

Khalid, H. (2024). Horizontally scaling the Rails backend of Shop app with Vitess. Shopify Engineering. <https://shopify.engineering/horizontally-scaling-the-rails-backend-of-shop-app-with-vitess> (accessed April 8, 2025)

Madan, P. (2021). Shard balancing: Moving shops confidently with zero-downtime at terabyte-scale. Shopify Engineering. <https://shopify.engineering/mysql-database-shard-balancing-terabyte-scale> (accessed April 9, 2025)

Rodukov, A. (n.d.). Shopify's strategy for seamless global expansion in eCommerce. LinkedIn. <https://www.linkedin.com/pulse/shopifys-strategy-seamless-global-expansion-ecommerce-alex-rodikov-z3exe> (accessed April 14, 2025)