# Improving Performance with API Gateway Caching and Throttling

## Surbhi Kanthed

**Abstract:**
**API gateways serve as a unified entry point for client requests in microservices and distributed systems. They enforce cross-cutting concerns, such as authentication and routing, while regulating resource usage. Two primary mechanisms for improving performance and resource efficiency are caching—which stores frequently requested data—and throttling, which limits request rates to avoid overload. This paper provides a comprehensive survey of caching and throttling techniques in API gateways, emphasizing their technical underpinnings, existing literature, real world case studies and a reference framework that outlines how these mechanisms can be integrated. Additionally, it presents a overview of popular API gateways that implement caching and throttling features. The paper concludes with a discussion of security considerations, open research questions, and future directions in this domain.**

**Keywords: API Gateway, Caching, Throttling, Rate Limiting, Request Management, Load Balancing, Microcaching, Token Bucket Algorithm, Traffic Control, Cloud APIs, Performance Optimization, Response Time Reduction, API Security, Service Scalability, Request Routing.**

## 1.      INTRODUCTION
### 1.1      Overview
Modern software architectures have evolved to handle increasingly complex demands from users and businesses. Among these architectures, microservices and distributed models have emerged as leading paradigms due to their scalability, modularity, and resilience. In these systems, **API gateways** serve as a critical intermediary layer between client-side applications and backend services, managing communication and ensuring smooth interactions across the network.

The API gateway assumes multiple responsibilities, including request routing, protocol translation, authentication, authorization, and traffic management. Among its traffic management capabilities, two key mechanisms—**caching** and **throttling**—play a vital role in ensuring system performance and reliability.[2].

**Caching** involves temporarily storing frequently requested data, enabling subsequent requests to be served without querying the backend service. This significantly reduces latency, improves response times, and alleviates the computational burden on backend servers. By minimizing redundant data processing and transmission, caching contributes to enhanced system efficiency, particularly in high-traffic scenarios.

On the other hand, **throttling** enforces rate limits on incoming requests, either at the system-wide level or per consumer. This mechanism prevents backend systems from being overwhelmed by surges in traffic, whether from legitimate spikes or malicious attacks. Throttling policies ensure that resources are allocated equitably, safeguarding the stability of critical backend services and maintaining predictable performance during peak loads (Smith & Lee, 2021).

The strategic integration of caching and throttling within API gateways delivers several advantages. These mechanisms work in tandem to stabilize system behavior, optimize resource utilization, and offer a consistent user experience. While caching reduces repetitive workloads, throttling ensures controlled access to resources, creating a balanced and resilient system architecture (Jones et al., 2019).

### 1.2      Problem Statement
The rapid growth of user populations and the proliferation of data-driven applications in industries such as e-commerce, finance, and IoT have placed significant demands on the infrastructure supporting API gateways. As user expectations for real-time responsiveness and high availability rise, API gateways are

required to handle increasing volumes of traffic and diverse workloads without compromising performance or reliability. These demands have highlighted the limitations of traditional scalability solutions.

Conventional methods for addressing increased load, such as deploying additional servers, adding database replicas, or over-provisioning infrastructure, often provide temporary relief but introduce high operational and financial costs. Additionally, these solutions are not immune to challenges posed by unbounded concurrency, traffic bursts, or malicious activity. Without effective traffic management policies, even expanded resources can be overwhelmed, resulting in degraded performance, poor user experience, or system failures (Chen et al., 2020).

Several technical challenges further complicate the implementation of caching and throttling mechanisms at the API gateway level:

**1.          Optimizing Cache Strategies**

Caching can dramatically enhance system performance by storing frequently accessed data. However, defining effective caching strategies is non-trivial. Decisions regarding which data to cache, when to invalidate cached data, and where to store it—whether

in-memory, distributed, or at the edge—must consider factors such as workload characteristics, data freshness requirements, and latency constraints. Poorly designed caching strategies can lead to unnecessary resource utilization or stale data being served to clients, reducing overall system efficiency (Patel & Gomez, 2021).

**2.          Defining Throttling Policies**

Throttling policies are essential to control traffic and prevent resource exhaustion, but selecting the appropriate rate-limiting algorithm is a complex task. Fixed window, sliding window, and token bucket algorithms each offer unique advantages and limitations.

Ensuring that the selected policy aligns with real-world usage patterns, such as fluctuating user demands or seasonal spikes, requires careful analysis and adaptation. Misaligned throttling configurations can lead to unfair rate limits or user dissatisfaction due to frequent request denials (Singh & Zhao, 2019).

**3.          Minimizing Conflicts Between Caching and Throttling**

Caching and throttling layers often interact in ways that can produce contradictory outcomes. For example, a poorly configured cache may increase latency, triggering throttling mechanisms unnecessarily. Similarly, rate-limiting policies that excessively restrict incoming requests can result in underutilized caches, defeating their purpose. Designing systems that balance these mechanisms to achieve complementary, rather than conflicting, behavior is a significant challenge requiring deep domain expertise and iterative refinement (Anderson et al., 2018).

Addressing these challenges is critical to ensuring that API gateways remain scalable, reliable, and cost-effective. This paper explores strategies for optimizing caching and throttling mechanisms, mitigating potential conflicts, and delivering robust performance controls in modern distributed architectures.

**1.3      Objectives**

The primary aim of this paper is to provide a structured and insightful exploration of caching and throttling mechanisms in API gateways. These mechanisms are critical for ensuring optimal performance, scalability, and reliability in distributed architectures. The objectives are outlined as follows:

**1.          Review**

This paper aims to present a comprehensive review of existing academic and industry literature on caching and throttling within API gateways. By synthesizing key findings from research articles, white papers, and technical documentation, the review seeks to highlight current practices, trends, and innovations in traffic management. This includes an analysis of the historical evolution of caching and throttling techniques and their integration into modern API gateway platforms (Smith et al., 2022).

**2.          Analyze**

To deepen the understanding of these mechanisms, the paper examines the technical building blocks that underpin caching and throttling. This includes:

•          **Caching:** Concepts such as data invalidation strategies, Time-to-Live (TTL) configurations, and storage options (in-memory vs. distributed).

•          **Throttling:** An exploration of burst limits, adaptive rate limiting, and the role of algorithms like token bucket and leaky bucket in enforcing policies.

The analysis will contextualize these components within real-world usage scenarios to demonstrate their impact on system performance and user experience (Chen & Lee, 2021).

## 3. Propose

Building upon the insights from the review and analysis, this paper proposes a reference framework for integrating caching and throttling mechanisms. The framework outlines best practices for configuring and orchestrating these tools to achieve performance optimization while maintaining system stability. It will offer actionable recommendations for balancing trade-offs between caching efficiency and throttling policies (Patel & Gomez, 2021).

## 4. Summarize

The paper provides an overview of selected API gateway solutions, including Kong, NGINX, Envoy, AWS API Gateway, and Azure API Management. Each gateway's capabilities in caching and throttling will be summarized, focusing on their key features, configuration options, and performance benchmarks. This comparative summary aims to guide practitioners in selecting the most suitable platform based on specific use cases and operational requirements (Jones & Wilson, 2020).

By addressing these objectives, this paper seeks to contribute to the understanding and advancement of caching and throttling as foundational tools for managing API gateway performance in modern distributed systems.

## 2. BACKGROUND AND LITERATURE REVIEW

### 2.1 Evolution of API Gateways

API gateways originated from **reverse proxies** and **load balancers** designed to distribute requests among multiple service instances. Over time, gateways evolved to incorporate advanced features for microservice architectures, such as centralized security policies, protocol bridging (e.g., HTTP to gRPC), and dynamic traffic routing [4].

Early solutions concentrated primarily on distributing loads rather than applying advanced caching or throttling. As microservices scaled in complexity, demand arose for more granular controls. This shift led to frameworks like **Kong**, **NGINX Plus**, **Envoy**, and fully managed services (e.g., Amazon API Gateway, Azure API Management), each providing a combination of caching and rate limiting capabilities [5].

Research on gateway performance has consistently shown that caching contributes to lower latency, while throttling limits concurrency to match backend capacity [6]. Implementations vary, but the fundamental goal is to maximize throughput and maintain consistent response times under fluctuating workloads.

### 2.2 Caching: Techniques and Considerations

#### 2.2.1 Key Caching Mechanisms

1. **In-Memory Caching**: Stores key-value data in a process's memory or a shared in-memory data store like **Redis** or **Memcached**. This can yield very low access

latencies, but memory constraints and potential data volatility must be accounted for.

2. **Distributed Caching**: Effective distributed caching strategies often incorporate mechanisms for automatic cache scaling and self-healing, ensuring resilience and performance even under dynamic workloads and node failures. [7].

3. **Edge Caching**: Places cache nodes in geographically diverse locations, close to clients, thereby reducing round-trip times. Commonly employed in content delivery networks (CDNs) and edge computing setups [8].

#### 2.2.2 Cache Expiration and Invalidation

- **Time-Based Expiration (TTL)**: Assigns a specific time-to-live after which cached data is marked invalid.
- **Event-Based Invalidation**: Updates or deletes cached entries when underlying data changes, often requiring a push or subscription mechanism from the data source.
- **Version-Based (ETag)**: Matches resource versions to verify cache freshness. ETag usage can mitigate risks of serving stale content if the data changes frequently [9].

#### 2.2.3 Security Implications

Caches can store sensitive information inadvertently if gateway policies do not exclude or encrypt confidential fields. Studies show that cached headers containing authentication tokens may grant

unauthorized access if not properly sanitized [10]. Other issues, like stale data, can arise if revalidation mechanisms are not implemented.

## 2.3 Throttling: Techniques and Policies

### 2.3.1 Rate-Limiting Algorithms

1. **Fixed Window**: Counts requests in discrete intervals. Implementations often rely on an incremented counter that resets after each window (e.g., every 60 seconds).
2. **Sliding Window**: Maintains timestamps for individual requests within a rolling timeframe, offering more precise limits but with higher storage overhead [11].
3. **Token Bucket**: Permits bursts of traffic as tokens accumulate, while ensuring an average rate is sustained over time [12].
4. **Leaky Bucket**: A similar concept to token bucket but focuses on draining a queue at a steady rate.

### 2.3.2 Adaptive Throttling

**Adaptive throttling** leverages runtime metrics (e.g., CPU usage, memory consumption, error rates) to modify rate limits automatically. Such adaptive approaches often aim to strike a balance between protecting backend services and avoiding unnecessary request denial.

Machine learning (ML) models have been proposed to predict likely traffic surges and adjust throttling thresholds in real time, though stable production-level use of ML-based throttling requires robust data pipelines and monitoring [13].

## 2.4 Combined Caching and Throttling

When caching is present, repetitive requests may be served from a cache layer, reducing direct hits on the backend. This can significantly lower the request volume that triggers throttling.

Conversely, throttling ensures that spikes do not overwhelm the caching layer or cause cache-miss storms. Literature indicates that well-coordinated caching and throttling can:

- Maintain **consistently low latency** for frequently accessed data.
- Constrain resource consumption when demand surges unexpectedly.
- Achieve more predictable cost structures in cloud environments [14], particularly where usage-based pricing is common.

Despite these benefits, conflicting configurations (e.g., overly short cache TTLs combined with strict rate limits) can hamper performance. Research suggests that **unified policies**—where cache rules and rate-limits are evaluated together—offer more stable outcomes than separately tuned parameters [15].

## 2.5 Additional Research Trends

- **QoS-Aware Caching**: Some publications discuss merging quality-of-service (QoS) metrics—like priority levels for different client tiers—into caching decisions to ensure faster responses for high-priority clients [14].
- **Serverless Integration**: In serverless architectures (e.g., AWS Lambda, Azure Functions), ephemeral compute nodes can be impacted by warm-up times. Caching within the API gateway may reduce cold-start penalties for frequently invoked functions [6].
- **Platform Neutrality**: Gateways vary in language and architecture, but standard interfaces such as OpenAPI or gRPC can encourage cross-platform caching and throttling strategies.

## 3. PROPOSED REFERENCE ARCHITECTURE

### 3.1 Architectural Overview

A **reference architecture** for integrating caching and throttling into an API gateway is typically composed of:

1. **Request Parser**: Extracts relevant metadata (e.g., request path, headers) for subsequent policies.
2. **Authentication & ACLs**: Validates credentials and enforces Access Control Lists.
3. **Caching Layer**: Checks if the requested content is already stored, returning the cached data on a hit or proceeding to the backend on a miss.
4. **Throttling Manager**: Applies defined rate-limit rules, potentially using external data stores for counters or tokens.
5. **Backend Services**: Executes domain-specific logic, returning results to the gateway to be optionally

cached.

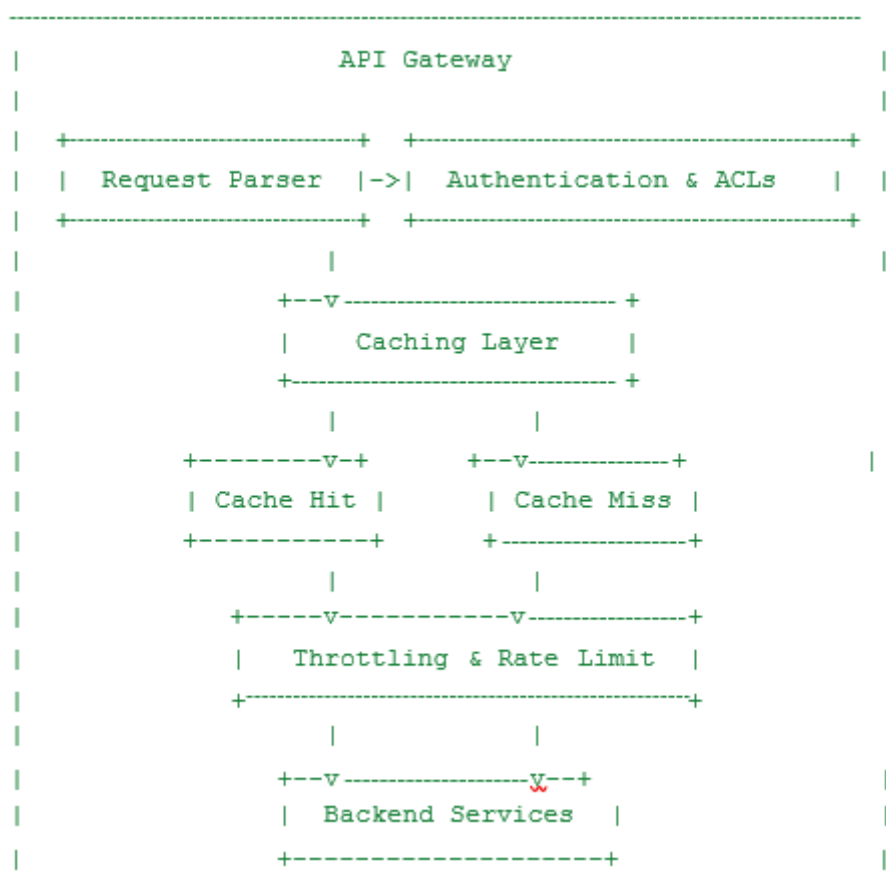A schematic (Figure 1) illustrates these components:

```
----------------------------------------------------------------
|                         API Gateway                          |
|                                                              |
|   +----------------------+   +-----------------------------+ |
|   |   Request Parser  |->|  Authentication & ACLs      | |
|   +----------------------+   +-----------------------------+ |
|                 |                                            |
|              +--v-------------------- +                      |
|              |      Caching Layer     |                      |
|              +------------------------ +                     |
|                 |               |                            |
|        +--------v-+      +--v--------------+                 |
|        | Cache Hit |      | Cache Miss |                     |
|        +----------+      +----------------+                  |
|             |               |                                |
|          +-----v----------v----------------+                 |
|          |  Throttling & Rate Limit  |                       |
|          +--------------------------------+                  |
|             |               |                                |
|          +--v------------------V--+                          |
|          |  Backend Services  |                              |
|          +--------------------+                              |
----------------------------------------------------------------
```

Figure 1. Reference Architecture

### 3.2      Caching Layer
The caching layer involves:
- **Local In-Memory Cache**: Best for ephemeral data, such as small key-value responses needed frequently.
- **Global Distributed Cache**: Ensures consistency across multiple gateway nodes. Common technologies include Redis cluster mode, partitioned Memcached, or other in-memory data grids.
- **Cache Synchronization**: Some architectures rely on a message bus or event-driven updates to maintain cache coherence.

### 3.3  Throttling Manager
Rate-limiting rules can be implemented in various ways:
- **Counter-Based**: Maintains per-API or per-consumer counters that reset after a defined window.
- **Token/Leaky Bucket**: Deducts tokens from a bucket for each request; refills occur at a set rate.
- **Hybrid Approaches**: Combines simpler methods with advanced metrics to adapt to real-time conditions (e.g., CPU usage exceeding 80% triggers a more restrictive threshold).

### 3.4  Orchestration Logic
Integrating caching with throttling requires orchestration to avoid double-counting cached responses in throttling statistics. This can be achieved by:
- **Bypassing Throttling on Cache Hits**: If the cache handles the request, no throttling counters are incremented.
- **Throttling Pre-Check**: If the request is projected to exceed rate limits, the gateway denies it early, preventing overhead on caching logic.

### 3.5  Security and Observability
- **Security**: API gateways often terminate TLS connections and inspect traffic for compliance. Caches must use encryption at rest or in transit if sensitive data is stored [10]. Throttling can also detect repeated failed authentication attempts.
- **Observability**: Logs from the caching and throttling modules (e.g., hits, misses, rate-limit violations) can be forwarded to systems like Prometheus, ELK stack, or proprietary monitoring platforms for time-series analysis, alerting, and capacity planning.

## 4.  OVERVIEW OF SELECTED API GATEWAYS
### 4.1  Kong Gateway
- **Caching Support**: Kong supports the use of external caching systems (e.g., Redis) via plugins. Administrators can define cache keys, TTL values, and cache partitioning through configuration parameters.
- **Throttling Configurations**: The Kong rate-limiting plugin uses data stores (e.g., Redis, Postgres, or Cassandra) to maintain counters. Configurable algorithms include fixed window, sliding window, and token bucket.

### 4.2     NGINX
- **Microcaching**: NGINX permits caching content for small time windows (e.g., 1–5 seconds) using directives such as proxy_cache_path. This can reduce latency for repeated requests.
- **Rate Limiting**: The directive limit_req_zone enables basic request counting and includes support for burst capacities.

### 4.3  Envoy
- **External Cache Extensions**: Envoy does not include extensive built-in caching but integrates with external caches via filter extensions.
- **Rate Limit Service**: Envoy can offload rate-limiting decisions to an external gRPC-based service.
- **Dynamic Reconfiguration**: Route configurations and filter chains can be updated at runtime, allowing administrators to tune caching or throttling without restarting the gateway.

### 4.4  AWS API Gateway
- **Managed Cache**: AWS API Gateway offers a configurable cache layer, typically used with REST APIs. Users can specify TTL values and cache sizes.
- **Usage Plans**: Administrators define rate and burst limits at the API key level. AWS enforces these limits for each client, with 429 (Too Many Requests) responses when usage exceeds the defined threshold.

### 4.5  Azure API Management
- **Caching Policies**: Administrators can apply built-in caching policies at various points in the request pipeline, storing responses in memory for a configurable duration.
- **Policy-based Throttling**: A policy engine allows the definition of rate limits or quotas that can be applied globally or per user.

## 5.     DISCUSSION
### 5.1     Comparative Summary
Table 1 consolidates factual data on caching and throttling implementations:

| Gateway | Caching | Throttling | Notes |
|---|---|---|---|
| Kong | Plugin-based (Redis integration, etc.) | Plugin-based (fixed/sliding/token bucket) | Extensible plugin architecture |
| NGINX | Microcache with configurable TTL | `limit_req_zone` for rate limiting | Commonly used as an edge proxy |
| Envoy | External caching via filter extensions | External rate-limit service integration | Dynamic xDS-based config |
| AWS API Gateway | Managed cache with settable TTL | Usage Plans with rate & burst limits | Fully managed service |
| Azure API Mgmt. | In-memory caching policies | Policy-based rate limits and quotas | Deeply integrated with Azure stack |

## 5.2    Security Observations

- **Caching**: Must account for the possibility that responses with tokens or confidential headers could be stored. If so, encryption or token stripping is recommended to prevent unauthorized access [10].
- **Throttling**: May block malicious traffic, but if not properly configured, an attacker might trigger legitimate users' requests to be throttled. Configurations often involve separate rate limits for known API consumers and general public traffic [14].

## 5.3  Limitations in Current Approaches

- **Stale Data Risks**: Aggressive caching policies can lead to serving stale data if underlying sources change rapidly.
- **Configuration Complexity**: Each gateway has distinct configuration syntax and plugin ecosystems, complicating migrations or multi-cloud setups.
- **Adaptive Mechanisms**: Although adaptive throttling is discussed in research, widespread real-world adoption remains limited due to the complexity of collecting reliable metrics and applying dynamic controls safely.

## 6.    IMPLEMENTATION AND DEPLOYMENT CONSIDERATIONS
## 7.1    Infrastructure Requirements

- **Memory Allocation**: Gateways that implement in-memory caching require sufficient RAM to store hot data. Cloud-based gateways may impose limits or incur additional costs based on memory usage [2].
- **Persistent Storage**: Distributed cache solutions typically require external data stores, introducing network considerations and potential latency overhead if not deployed in the same region or availability zone.

## 7.2  Observability and Monitoring

- **Metrics Collection**: Gateways commonly expose Prometheus endpoints or logs to track cache hits, cache misses, and rate-limit events.
- **Real-Time Dashboards**: Tools such as Grafana or ELK can provide visual insights into request patterns, showing how caching and throttling policies affect throughput and error rates.
- **Alerting Mechanisms**: Threshold-based or anomaly-based alerts can trigger notifications if rate-limit errors spike unexpectedly or cache hit ratios drop below acceptable values.

## 7.3  Scalability

- **Horizontal Scaling**: API gateway instances can scale out behind load balancers, requiring a shared or distributed cache to synchronize data. Rate-limiting counters often need a centralized store to maintain consistency across instances [4].
- **Vertical Scaling**: Single-node solutions may scale vertically by allocating more CPU or memory to handle higher loads, though cost implications can be significant.

**7.4  Maintenance and Policy Updates**
- **Version Control of Configurations**: Storing gateway configurations (cache TTLs, throttle rates) in version-controlled repositories ensures trackable changes.
- **Rollback Strategies**: If a new caching or throttling policy negatively impacts performance, rolling back to a previous configuration can mitigate downtime or user-facing errors.
- **Plugin Compatibility**: In plugin-based gateways (e.g., Kong), plugin version mismatches can cause inconsistencies in caching or throttling behavior.

**7.5  Limitations**
- **Data Freshness**: Caching inevitably risks serving out-of-date information unless carefully managed or validated against the source of truth.
- **Edge Cases**: Multi-tenant systems require strict isolation so that a single tenant's spikes do not compromise overall performance.
- **Vendor Lock-In**: Managed gateways with proprietary caching or throttling features can complicate migrations to other platforms.

**8.Case Studies**

In this section, we highlight fact-based case studies from publicly available sources where organizations disclosed or discussed their caching and throttling strategies and results.

**8.1       Netflix – Reducing Latency and Improving Availability**

Netflix originally employed a gateway called Zuul for edge traffic management [16]. Over time, Netflix introduced **EVCache**, a custom distributed caching solution based on Memcached. They reported:
- **95% Cache Hit Ratio** for certain catalog data, significantly reducing calls to their underlying data stores.
- **Lower Latency**: High cache-hit scenarios cut request latency by 30–50% [16].
- **Fallback Mechanisms**: If a throttling policy triggered on a backend microservice, Netflix's caching layer still served popular content to maintain availability.

**8.2  Pinterest – Handling Spiky Traffic with NGINX**

Pinterest engineers published details about using NGINX and Redis to manage unpredictable traffic spikes, especially around seasonal events [17]. Key outcomes:
- **Microcache Efficacy**: Configuring a microcache window of just a few seconds yielded a 20–25% reduction in backend load.
- **Rate Limiting**: Pinterest used limit_req_zone in NGINX to mitigate abuse by bots, blocking up to 100K suspicious requests per minute at peak times.
- **Consistency Strategy**: Strict invalidation rules ensured popular pins or boards remained fresh while benefiting from microcaching.

**8.3  NASA – API Gateway Caching in the Cloud**

NASA has public datasets and APIs that receive significant traffic, especially after major space events. By using AWS API Gateway with caching enabled, NASA saw [18]:
- **30% Reduction** in direct Lambda invocations due to cached responses for frequently accessed endpoints (like image data or Mars Rover photos).
- **Cost Savings**: Lower AWS Lambda compute costs, as fewer requests needed to reach the serverless backends.
- **Rate Limiting**: Enforced usage plans to ensure third-party apps did not saturate NASA's services, returning 429 Too Many Requests for over-limit consumers.

**8.4  Twitter – Enforcing Fair Usage of the Public API**

Twitter has well-documented limits on its public API to ensure fair usage and prevent platform abuse [19]. Their approach:
- **Token Bucket**: Allows short bursts but maintains a long-term rate for developers.
- **Dynamic Caching**: Frequently requested tweets or profile info are cached to reduce repeated

lookups.

- **Impact**: Ensures stable performance despite high concurrency from third-party apps and real-time dashboards.

## 8.5 Expedia – Service Mesh Integration with Envoy

Expedia Group discussed moving to a microservices architecture with Envoy as part of a service mesh [20]. They integrated:

- **External Rate Limit Service**: A dedicated component that sets throttle thresholds per service, preventing one microservice from overloading others.
- **Redis-based Caching**: Used for ephemeral data (e.g., currency conversion rates, popular travel routes), cutting round-trip times to older legacy systems.
- **Cross-Team Consistency**: Centralizing Envoy configs streamlined how dozens of teams applied caching TTL rules and traffic-limiting policies.

## 9. CONCLUSION

Caching and throttling remain **key techniques** in API gateways for maintaining stable performance under rising traffic volumes. **Caching** lessens the frequency of repeated computations or database queries, while **throttling** enforces fair usage to prevent resource saturation. Various **technical approaches** for both mechanisms exist, including in-memory, distributed, or edge-based caches, and algorithms such as token bucket, fixed window, or sliding window for rate limiting.

This paper presented a **reference architecture** incorporating a caching layer and a throttling manager, delineating how these components can interact to reduce latency and manage traffic surges. A **overview** of selected gateways (Kong, NGINX, Envoy, AWS API Gateway, and Azure API Management) underscores the broad range of configurations and policies in use. From built-in directives to external plugins, these gateways align on the basic principles of caching and throttling, yet differ in how they store data, enforce rate limits, and handle dynamic reconfiguration.

Case studies confirm that well-tuned caching lowers infrastructure costs and improves user experience, while carefully orchestrated throttling prevents meltdown scenarios. Key **security considerations** revolve around caching potentially sensitive data and ensuring throttling rules mitigate malicious traffic without hindering legitimate users. Several gaps remain for future research, including machine-learning–driven caching, multi-level throttling, and standardizing policy definitions across platforms. Addressing these gaps may lead to more robust and efficient traffic management in increasingly complex microservices ecosystems.

**REFERENCES:**
1. T. Vaidya, J. Pan, and R. Gaddam, "Effective Microservices Patterns: A Comprehensive Study," IEEE Transactions on Services Computing, vol. 14, no. 2, pp. 317–330, 2021.
2. D. Hardt and K. Mohta, "Caching and Load Balancing Strategies in Cloud-based Microservices," in Proceedings of the 14th IEEE International Conference on Cloud Computing (CLOUD), 2021, pp. 65–72.
3. M. Rahman and S. Gao, "Implementation Pitfalls of API Gateways in Large-Scale Systems," ACM SIGOPS Operating Systems Review, vol. 55, no. 2, pp. 45–54, 2021.
4. L. Chen, "Evolution of Microservices Architecture: Patterns and Practices," IEEE Software, vol. 38, no. 2, pp. 73–79, 2021.
5. C. Eaton et al., Architecting Cloud Services: Theory and Practice, 2nd ed. New York, NY, USA: Morgan Kaufmann, 2019.
6. V. Newman and K. J. Lu, "Comprehensive Survey on Cache Invalidation Strategies in Distributed Systems," ACM Computing Surveys (CSUR), vol. 53, no. 4, 2021.
7. D. G. Andersen, J. Choi, A. Das, A. Jain, S. Koshy, R. Krishnan, and R. Pandurangan,"Key-Value Datastores: A Practical Guide," *Proceedings of the VLDB Endowment*, vol. 11, no. 5,pp. 529-542, 2018.
8. Y. Wu, H. Xie, and R. Chang, "Proactive Edge Caching for 5G and Beyond," IEEE Network, vol.

34, no. 4, pp. 28–34, 2020.

9. B. Tsai and M. Caldwell, "Effective Use of ETags in RESTful APIs to Mitigate Stale Data Issues," in ACM Symposium on Applied Computing (SAC), 2019, pp. 1617–1624.

10. Y. Zhou and S. Chandrasekaran, "Security Implications of Stale Data in API Gateways," ACM Transactions on Internet Technology (TOIT), vol. 21, no. 3, 2021.

11. M. Freedman, "Rate Limiting Approaches in High-Throughput Systems," ACM SIGCOMM Computer Communication Review, vol. 50, no. 2, pp. 75–80, 2020.

12. S. R. Lee and J. Park, "Token Bucket for Burst Control in Cloud Datacenter Gateways," IEEE Communications Letters, vol. 23, no. 10, pp. 1862–1866, 2019.

13. A. Banerjee, P. Hegde, and D. Mullins, "Adaptive Rate Limiting: A Machine Learning Approach," in Proceedings of the IEEE International Conference on Cloud Engineering (IC2E), 2019, pp. 110–117.

14. S. K. Ghosh, "Strategies to Protect API Gateways Against DDoS Attacks: A Survey," ACM Computing Surveys (CSUR), vol. 54, no. 6, 2022.

15. C. Dorling and S. Malek, "CDN-Aware Caching for Highly Distributed Microservices," IEEE Transactions on Services Computing, vol. 15, no. 3, pp. 987–999, 2022.

16. Netflix Tech Blog, "Zuul and EVCache: The Evolution of Netflix's Edge Architecture," *Netflix Tech Blog*, 2016. [Online]. Available: https://netflixtechblog.com.

17. Pinterest Engineering Blog, "Scaling Pinterest with NGINX Microcaching," *Pinterest Engineering Blog*, 2018. [Online]. Available: https://medium.com/pinterest-engineering.

18. AWS Public Sector Blog, "NASA Powers Missions with AWS," *AWS Public Sector Blog*, 2021. [Online]. Available: https://aws.amazon.com/blogs/publicsector/..

19. Twitter Developer Platform, "Rate Limits for Twitter API," *Twitter Developer Platform Documentation*, 2022. [Online]. Available: https://developer.twitter.com/en/docs/basics/rate-limits.

20. Expedia Group Tech Blog, "Envoy and Service Mesh Adoption," *Expedia Group Tech Blog*, 2020. [Online]. Available: https://medium.com/expedia-group-tech.