# Step 1: The Contract (Requirements)

## Functional Requirements

✓ Identify Clients: User ID, IP Address, or API Key

✓ Configurable Rules: e.g., 100 requests / minute

✓ Feedback: Return HTTP 429 + Headers (Remaining, Reset Time)

## Non-Functional Requirements

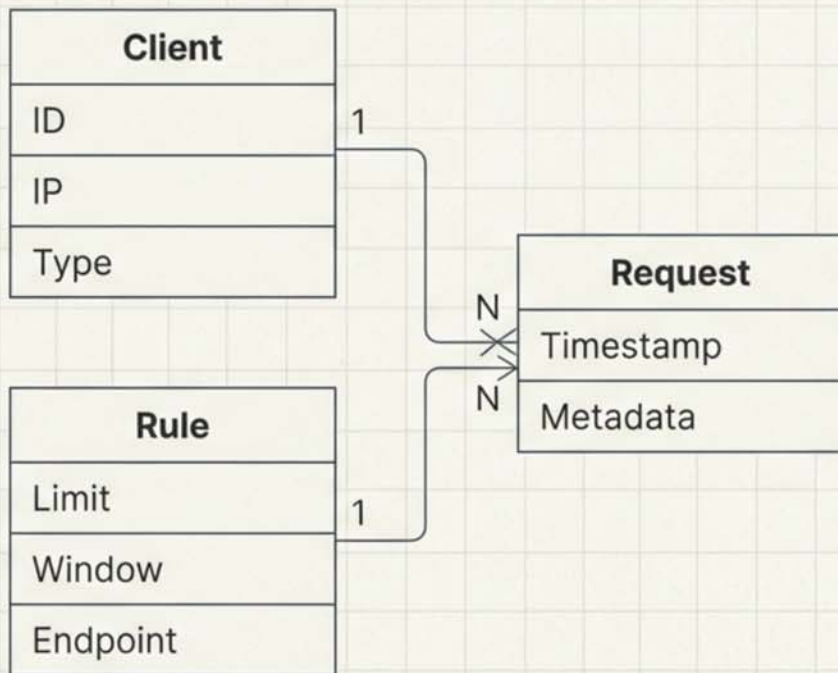⚖ Scale: **1M Requests Per Second (RPS)** / 10M DAU
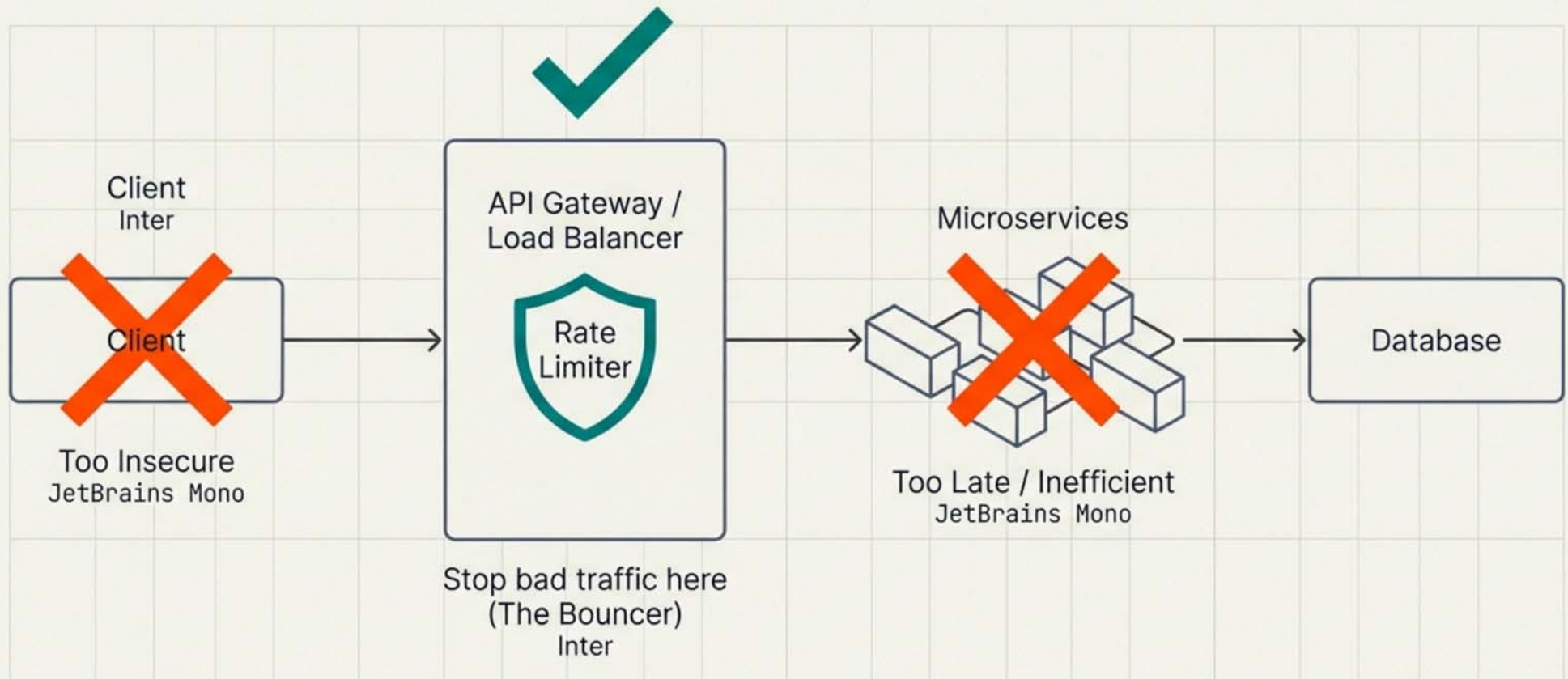
⏱ Latency: **< 5ms** per check

⚙ Availability: High availability > Strong consistency

# Step 2: The Blueprint

| Client | |
|---|---|
| ID | 1 |
| IP | |
| Type | |

| Rule | |
|---|---|
| Limit | 1 |
| Window | |
| Endpoint | |

| Request | |
|---|---|
| Timestamp | N |
| Metadata | N |

```
1   interface RateLimiter {
2       boolean isRequestAllowed(String clientId, String ruleId);
3
4       // Returns:
5       // - allowed: true/false
6       // - remaining: int
7       // - resetTime: timestamp
8   }
```

# Strategic Placement: The 'Bouncer' Analogy



Client
Inter

Client

Too Insecure
JetBrains Mono

API Gateway /
Load Balancer

Rate
Limiter

Stop bad traffic here
(The Bouncer)
Inter

Microservices

Too Late / Inefficient
JetBrains Mono

Database

# Identification Strategy: Who is Knocking?
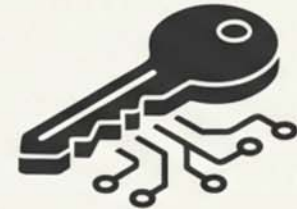
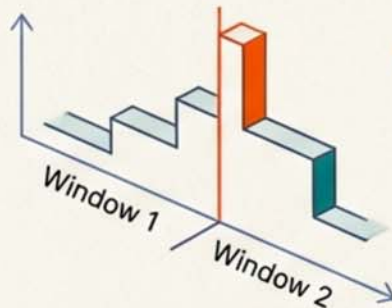| User ID | IP Address | API Key |
|---------|------------|---------|
| Best for Authenticated Users. Enables specific user limits. | Fallback for Anonymous/DDoS. **Risk:** Shared IPs (NAT). | Standard for B2B / Developer Tools. |

**The Senior Answer**

Use a Hybrid Approach. Authenticated users get higher limits; fallback to IP limits for unauthenticated traffic.
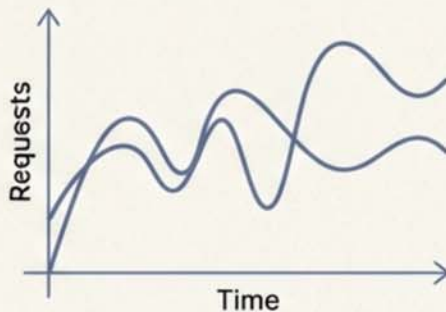
# The Chosen Algorithm: Token Bucket



Refill Rate

Tokens

Overflow

Empty

Request

Processing

Blocked

**State per User**

1. CurrentTokens (Integer)
2. LastRefillTimestamp (Long)

# High-Level Architecture (The MVP)



Client

1. Request →

API Gateway

3. Calculate Refill

2. Fetch (HMGET) →

4. Update (HSET) →

5. Allow/Deny ←

Redis
Single Instance

Externalized State
(In-Memory, Fast)

# Critical Deep Dive: The Race Condition

## The Problem



**T1**

Gateway A — Read: Tokens: 1 → Redis → Decrement to 0 → **T2**

Gateway B — Read: Tokens: 1 → Redis → Decrement to 0 → **T2**

**Result: Both requests pass. Limit Exceeded.**

## The Solution

**Redis Lua Script**

Get
→ Calculate
→ Set

Atomic Execution.
No other operation can
interrupt the script.

NotebookLM

# Scaling to 1 Million RPS: Sharding

User ID: Alice

Hash(Alice) % 16384
(Slots)

Shard #3

Determinism is key. Alice must always hit Shard #3.

50k OPS

Redis Cluster (1M+ OPS)

# Optimizing for < 5ms Latency
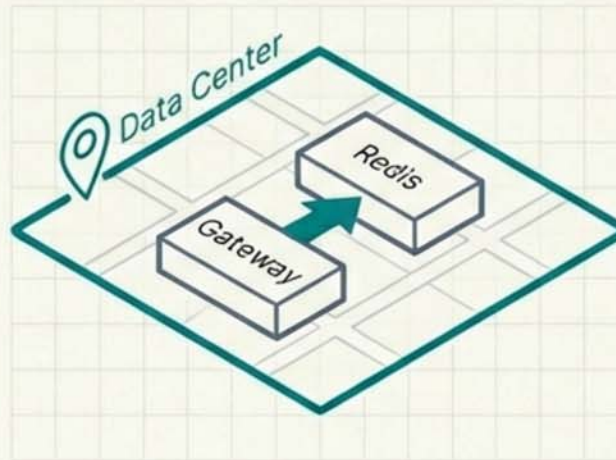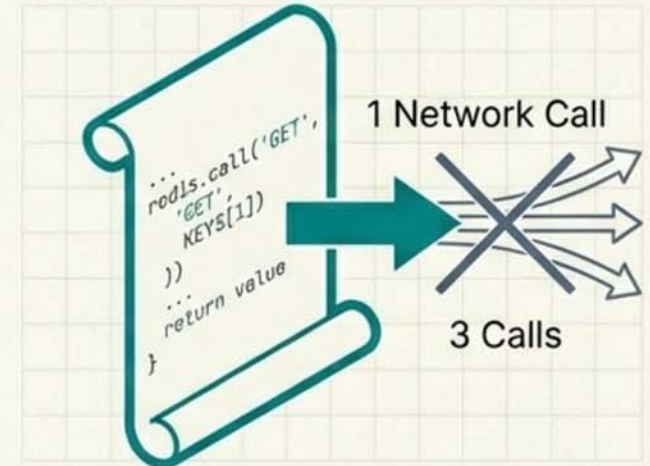
**Connection Pooling**

Reuse TCP connections.
Avoid handshake overhead.

**Geo-Proximity**

Colocate Cache
& Gateway.

**Lua Scripting**

1 Network Call
instead of 3.

# Dynamic Configuration