



# Rate Limiting em Go com SDK próprio e suporte a Multi-Tenancy

## Contexto

Microserviços sem controle de requisições podem se tornar instáveis quando submetidos a cargas excessivas de tráfego. Clientes malcomportados ou scripts com bugs conseguem enviar *floods* de requisições, esgotando recursos do serviço (CPU, memória, conexões de banco de dados) e degradando a experiência para todos os usuários <sup>1</sup> <sup>2</sup>. Em ambientes **multi-tenant** (multi-inquilino), surge também o problema do "**noisy neighbor**", em que um cliente em um tier inferior pode consumir recursos desproporcionalmente e afetar a performance percebida pelos demais clientes, inclusive aqueles de tier premium <sup>3</sup>. Sem um controle adequado, um usuário gratuito pode monopolizar o serviço e impactar usuários pagos, violando expectativas de *fairness*. Nesse contexto, **rate limiting** desponta como uma ferramenta essencial para estabilizar microservices, garantindo que nenhuma entidade (usuário, IP ou chave de API) exceda limites razoáveis de uso e impedindo interferência indesejada entre clientes de diferentes planos de serviço <sup>4</sup>. A seguir, exploraremos estratégias clássicas de rate limit, sua aplicação em Go via um SDK customizado com Redis e memória local, e como tais técnicas permitem isolar e controlar o consumo por tenant/tier em sistemas distribuídos.

## Conceitos de Rate Limiting

**Rate limiting** é a técnica de restringir a taxa de eventos (requisições, por exemplo) que um sistema processa, geralmente definindo um número máximo de solicitações permitidas por intervalo de tempo. Diferente de *throttling* (que tenta apenas desacelerar ou enfileirar requisições sob alta carga), o rate limiting clássico atua como um “porteiro rigoroso”: ao atingir o limite, requisições adicionais são bloqueadas ou descartadas (tipicamente retornando HTTP 429 Too Many Requests) <sup>5</sup>. Isso garante previsibilidade de performance e protege recursos compartilhados, impondo uso justo mesmo em cenários de pico. Existem vários algoritmos e métodos para implementar rate limiting; os principais estão listados a seguir:

- **Janela Fixa (Fixed Window Counter)** – Essa estratégia divide o tempo em janelas discretas de tamanho fixo (por exemplo, intervalos de 1 minuto) e mantém um contador de requisições por janela <sup>6</sup>. Cada requisição incrementa o contador da janela atual; se o contador excede o limite definido, as requisições restantes dentro daquela janela são bloqueadas até a próxima janela temporal começar. A principal vantagem é a simplicidade de implementação (pode-se usar um contador com expiração pelo período da janela) <sup>7</sup>. No entanto, há uma desvantagem notória: **picos em bordas de janela**. Como os contadores são zerados no início de cada janela, um cliente pode efetuar até o limite no final de uma janela e imediatamente repetir o limite no início da janela seguinte, efetivamente dobrando o tráfego permitido em um curto intervalo <sup>8</sup>. Por exemplo, com limite de 100 requisições/min, um usuário pode disparar 100 requisições nos segundos finais de um minuto e mais 100 nos primeiros segundos do minuto seguinte, conseguindo ~200 requisições em ~1 segundo. Essa natureza discreta torna o controle menos uniforme e pode causar **bursts** abruptos de carga.

- **Janela Deslizante (Sliding Window)** – Para mitigar os bursts do método de janela fixa, utiliza-se a abordagem de janela móvel. Aqui, em vez de resetar contadores rigidamente a cada intervalo, considera-se um período contínuo “deslizante”. Existem duas implementações comuns:
  - *Sliding Log (Registro Deslizante)*: armazena-se um log temporal (timestamps) de cada requisição recente. A cada nova requisição, remove-se do log as timestamps fora da janela de tempo desejada (por exemplo, mais de 60 segundos atrás) e calcula-se o número de requisições dentro da janela atual móvel <sup>9</sup>. Se esse número for menor que o limite, a nova requisição é permitida e seu timestamp é registrado; se já atingiu o máximo, a requisição é bloqueada. Essa técnica é **precisa**, pois contabiliza exatamente as ações no intervalo móvel, eliminando completamente o problema de burst nas bordas <sup>10</sup>. Contudo, sua desvantagem é de performance e memória: em APIs de alto volume, manter e filtrar uma lista de timestamps por usuário pode ser custoso <sup>11</sup>.
  - *Sliding Window Counter (Contador Deslizante com Peso)*: é uma aproximação mais eficiente que combina o contador da janela atual com o da janela anterior para suavizar a transição <sup>12</sup>. Por exemplo, suponha que a janela seja 1 minuto e o limite 60 req/min. Se estamos 50% no meio da janela atual, podemos ponderar 50% das requisições da janela passada + as da atual para decidir se o limite foi alcançado. Em termos simples, calcula-se uma métrica como `total = req_atual + proporção_da_janela_anterior * req_anterior`. Se `total` exceder o limite, bloqueia-se <sup>13</sup>. Dessa forma, evita-se um salto brusco no reset da janela – quanto mais perto do fim da janela, mais peso se atribui ao histórico recente. Embora não seja tão exato quanto o log completo, esse método reduz muito o overhead de armazenamento (mantém basicamente dois contadores) e elimina os picos instantâneos nas bordas, resultando em um controle **mais uniforme** do fluxo de requisições.
- **Token Bucket (Balde de Fichas)** – É um dos algoritmos de rate limit mais difundidos, graças à sua simplicidade e flexibilidade <sup>14</sup>. Imagine um balde que acumula *tokens* (fichas) a uma taxa constante, até um máximo de capacidade. Cada token representa permissão para uma requisição. Quando uma requisição chega, ela “retira” um token do balde:
  - Se houver tokens disponíveis, a requisição é permitida e um token é consumido.
  - Se o balde está vazio (sem tokens), a requisição deve ser rejeitada (ou em alternativa, pode aguardar até haver token, dependendo da política).
  - O balde se enche no tempo ocioso – tokens são adicionados continuamente a uma taxa fixa (por exemplo, 5 tokens por segundo) até o limite de capacidade (que representa o **burst** máximo acumulável) <sup>15</sup>.

O token bucket permite absorver picos curtos de tráfego: um cliente pode acumular créditos durante a ociosidade e depois realizar várias requisições de uma vez até esvaziar o balde (até o tamanho do *burst*) <sup>16</sup>. Isso é útil para cenários em que desejamos tolerar rajadas esporádicas sem negar serviço (e.g., um usuário faz várias requisições rapidamente após um período de inatividade). Em média, porém, a taxa sustentada não excederá a reposição de tokens. As desvantagens potenciais incluem: (1) necessidade de manter um estado por cliente (o número de tokens), o que em sistemas com muitos usuários pode consumir memória proporcionalmente <sup>17</sup>; e (2) não impõe espaçamento perfeitamente uniforme – se o burst permitido for grande, ainda poderão ocorrer aglomerados de requisições (embora limitados pelo tamanho do balde).

- **Leaky Bucket (Balde Furado)** – Esse algoritmo é análogo a um balde com um furo no fundo por onde a água (requisições) escoa a uma taxa fixa <sup>18</sup>. Em termos de requisições:

- As requisições chegam e “entram” no balde (podemos imaginar enfileirar as solicitações).
- O processamento sai do balde a uma velocidade constante predefinida (por exemplo, 10 req/segundo), independentemente de quanto rápido as requisições entram.
- Se as requisições chegam mais rápido do que a taxa de saída, o balde enche. Quando o balde atinge sua capacidade máxima, requisições adicionais **transbordam** (são descartadas) <sup>19</sup> <sup>20</sup>.

O efeito do leaky bucket é **suavizar picos**: mesmo que um cliente faça muitas requisições de uma vez, o sistema libera elas de forma nivelada ao longo do tempo, prevenindo sobrecarga instantânea do backend <sup>21</sup>. Em outras palavras, garante um fluxo de saída constante (limite fixo de throughput) – “suave como manteiga”, conforme analogia em uma fonte <sup>22</sup> <sup>23</sup>. Uma vantagem é fornecer grande previsibilidade: o sistema processa requisições num ritmo estável, evitando rajadas. A contrapartida é que **não tolera bursts** do cliente além do balde: qualquer pico que excede a capacidade imediata resulta em perdas de requisições (o balde furado simplesmente deixa transbordar requisições extras) <sup>23</sup>. Implementar literalmente um balde furado às vezes envolve uma fila de tarefas com um trabalhador liberando em velocidade fixa, ou um contador que decresce periodicamente. Uma implementação ingênuo requer um processo de *drip* (gotejamento) para esvaziar contadores regularmente <sup>24</sup>. Isso traz desafios: se o processo de drip falhar ou não der vazão, o mecanismo pode se desequilibrar e bloquear indevidamente requisições novas <sup>25</sup>. Apesar disso, o conceito é poderoso e, com ajustes, bastante utilizado para garantir estabilidade do lado do servidor.

- **GCRA (Generic Cell Rate Algorithm)** – O GCRA é uma variante aprimorada do leaky bucket que evita a necessidade de um gotejamento contínuo, sendo muito apropriado para implementação distribuída. Originado nas redes de telecomunicações (ATM), o GCRA é essencialmente um algoritmo de escalonamento de células que verifica conformidade de taxa de maneira matemática em vez de simular fisicamente tokens ou vazão <sup>26</sup>. Em vez de decrementar um bucket ao longo do tempo, ele calcula o instante teórico em que a próxima requisição poderá ocorrer sem violar a taxa: mantém-se um **TAT (theoretical arrival time)**, que é o horário em que uma próxima ação seria permitida dado o histórico <sup>26</sup>. Quando uma nova requisição chega:
  - Compara-se o horário atual com o TAT ajustado por um fator de tolerância (basicamente **TAT - (burst\_interval)** onde *burst\_interval* corresponde à capacidade de burst em tempo) <sup>27</sup>.
  - Se o horário atual é posterior ou igual a esse limite (ou seja, a requisição está vindo “atrasada” o suficiente para estar dentro da taxa permitida), então ela é **permitida**. Nesse caso, recalcula-se um novo TAT adicionando o intervalo base da taxa (por exemplo, se a taxa é 5 req/s, adiciona 200ms ao TAT) <sup>28</sup>.
  - Se o horário atual ainda está antes do próximo tempo permitido (a requisição chegou cedo demais), então a requisição é **negada** por estar fora da conformidade <sup>28</sup>.

Em resumo, o GCRA impõe espaçamento mínimo entre eventos de forma eficiente, sem precisar manter contadores decrementando constantemente. Ele é conceitualmente equivalente ao leaky bucket e produz um resultado similar (taxa suave e bursts controlados) <sup>26</sup>, mas é implementado calculando tempos ao invés de remover tokens ou água gota a gota. Uma grande vantagem prática é em cenários distribuídos: pode-se manter apenas uma última marcação temporal e um “crédito” de tolerância, e usar operações atômicas para atualizar isso em um store central (como Redis) a cada requisição. De fato, bibliotecas populares de rate limit usam GCRA para facilitar a atomicidade no Redis <sup>29</sup>. O GCRA também é conhecido como a versão “sem balde” do leaky bucket (“leaky bucket as a meter”) – ao invés de simular vazão, ele mede diretamente se uma ação cabe na taxa. A implementação requer atenção à sincronização de tempo: em ambientes distribuídos, é importante usar uma referência consistente de horário (por exemplo, o comando **TIME** do Redis ou relógio NTP sincronizado) para evitar deriva de relógio entre nós <sup>30</sup>. Em

resumo, o GCRA fornece a robustez do leaky bucket com menos complexidade operacional, sendo usado em sistemas como o **Stripe** (via biblioteca *Throttled*) e em comandos nativos de alguns bancos chave-valor modernos <sup>31</sup>.

- **Limitador de Concorrência** – Diferente dos algoritmos acima (que limitam taxa no tempo), o limitador de concorrência foca em quantas operações *simultâneas* podem ocorrer de uma vez. É uma forma de evitar sobrecarga instantânea limitando o paralelismo. Por exemplo, um serviço pode estipular que no máximo 50 requisições podem ser processadas concorrentemente; se a 51<sup>a</sup> chegar enquanto 50 estão em andamento, ela será enfileirada ou rejeitada imediatamente. Técnicas simples implementam isso usando semáforos ou *tokens* de vaga: cada requisição em andamento ocupa um token, liberando-o ao terminar. Caso não haja token livre, a requisição nova não inicia até que um seja liberado <sup>32</sup>. Essa abordagem garante controle sobre recursos como threads, conexões de BD ou throughput imediato, complementando o rate limiting temporal. **Importante:** limitadores de concorrência não substituem o controle de taxa, eles agem em conjunto. Por exemplo, num pico repentino, o limitador de concorrência previne centenas de processos simultâneos, mas se a carga sustentada ficar alta (várias requisições em fila constantemente), o rate limiter temporal deve reduzir a taxa média de chegada. Em resumo, concorrência máxima impõe um teto de *trabalho paralelo*, enquanto rate limit impõe um teto de *trabalho por unidade de tempo*. Ambas técnicas evitam sobrecarga, mas em formas diferentes, e podem ser combinadas para maior eficácia.

## Arquiteturas de Rate Limiting Distribuído

Projetar um rate limiter para microsserviços distribuídos (vários nós atendendo requisições) requer decidir onde e como aplicar os limites para suportar **multi-tenancy** de forma consistente. Há diferentes padrões arquiteturais possíveis:

- **Rate Limiting Centralizado (Gateway/API Gateway):** Uma abordagem comum é implementar o controle de requisições em um ponto central de entrada – por exemplo, um API Gateway ou proxy reverso por onde todo o tráfego passa. Esse gateway mantém contadores globais e aplica as regras de rate limit antes de repassar as requisições aos microsserviços backend. A vantagem é a simplicidade de **visão única**: como todo tráfego passa por um único componente (ou camada), é fácil contabilizar e limitar globalmente por cliente/tenant. Isso evita que um usuário contorne limites alternando entre instâncias, pois o gateway vê tudo. Além disso, gateways frequentemente suportam rate limiting configurável out-of-the-box (por exemplo, Amazon API Gateway oferece tiers e throttling por API key nativamente <sup>33</sup>). A desvantagem é adicionar uma dependência central na arquitetura – o gateway se torna crítico para a disponibilidade – e menos flexibilidade dentro de cada serviço para regras específicas. Ainda assim, combinar um limite global no gateway com limites locais nos serviços pode fornecer camadas de proteção (defesa em profundidade) <sup>34</sup>.
- **Rate Limiting Distribuído via Armazenamento Compartilhado:** Aqui, cada instância do microsserviço aplica o rate limit consultando um armazenamento central (p.ex. Redis, Memcached, banco de dados) que guarda o estado dos contadores/buckets. Todas as instâncias referenciam os mesmos dados de limite, garantindo um comportamento consistente cluster-wide. Essa arquitetura é suportada por bibliotecas como o `go-redis/redis_rate`, que utiliza o Redis para armazenar e atualizar contadores atômicos usando GCRA <sup>29</sup>. O fluxo típico é: ao chegar uma requisição, a instância do serviço calcula a chave de rate limit (por exemplo, combina o tenant ou API key com o nome da API) e executa um comando/consulta no Redis para verificar e atualizar o contador/token.

Com comandos atômicos (p.ex. usando scripts Lua ou comandos com increment e expiração), assegura-se que mesmo requisições concorrentes em nós distintos obedecem a um limite comum. A grande vantagem é escalabilidade horizontal com consistência – independente de quantos servidores atendem, todos compartilham a mesma visão de uso do cliente. A latência adicional da chamada ao store externo (Redis) é geralmente baixa (milissegundos), mas deve ser levada em conta; otimizações como *pipelining* ou redução de chamadas podem ser necessárias em tráfego altíssimo. Outra preocupação é a **disponibilidade** do store: se o Redis falhar, o rate limiting pode ficar indisponível. Uma prática comum é projetar um *fail-open* (se o store estiver off, permitir requisições em vez de bloquear tudo), para não causar interrupção total do serviço – ainda que isso signifique operar temporariamente sem limites.

- **Rate Limiting Local em Memória (por Instância):** Nessa estratégia, cada instância do serviço aplica limites de forma isolada, usando armazenamento local (memória, variáveis no processo) para rastrear as requisições de cada cliente. Por exemplo, uma instância pode usar a biblioteca [golang.org/x/time/rate](https://golang.org/x/time/rate) para criar limitadores de token bucket para cada API key. Isso é **muito eficiente** em termos de latência (tudo em memória local, sem I/O remoto) e simplicidade – se seu serviço roda em apenas um nó, é provavelmente suficiente. Contudo, em cenários multi-nó, os limites tornam-se **aproximados**: cada instância só conhece as requisições que ela mesma atendeu. Assim, se você tem N instâncias e configura cada uma com limite L por cliente, um cliente poderia potencialmente fazer NL requisições totais (*espalhando uniformemente pelas instâncias*) sem que nenhuma instância isolada detecte a violação global. Em outras palavras, o cliente poderia “driblar” o limite global usando paralelismo contra diferentes nós do cluster. Uma mitigaçāo simples é dividir o limite global desejado pelo número de instâncias e aplicar por nó (por ex, se quer 100 req/s por cliente global e há 5 instâncias, cada instância limita 20 req/s daquele cliente). Porém, isso assume tráfego balanceado e um número estático de instâncias – em auto scaling dinâmico, essa divisão complica. Por essas razões, rate limiting puramente local é indicado quando: (a) garantidamente cada cliente sempre vai para o mesmo nó (sessões sticky – raramente garantido em escala); ou (b) aceitamos uma margem de erro no limite global em prol de eliminar a dependência externa. Na prática, muitas arquiteturas usam limitadores locais para proteção individual de instância\* (evitar que um nó se sobrecarregue, caso o平衡ador de carga não distribua perfeitamente), combinados com um limitador distribuído global para consistência forte.
- **Arquiteturas Híbridas:** Em sistemas de altíssimo desempenho, podem-se mesclar abordagens. Por exemplo, um modelo hierárquico onde cada instância aplica um limite local agressivo (para resposta rápida) e reporta contadores periodicamente a um serviço central que faz ajustes globais ou feedback. Ou então usar caches locais de tokens: o serviço pega um lote de permissões do Redis (digamos tokens para X requisições) e depois consome localmente sem ir ao Redis para cada requisição, renovando quando o lote esgota – isso amortiza a latência de acesso distribuído. Esses padrões aumentam a complexidade, mas podem melhorar a eficiência. Em geral, a escolha arquitetural depende do trade-off entre **precisão global vs latência/complexidade**: para limites estritos por tenant (como contratos de SLA), costuma-se preferir a consistência global via store central<sup>35</sup>; para proteção genérica de estabilidade, limitadores locais já contribuem bastante a evitar saturação de um nó.

Em todas essas arquiteturas, um ponto fundamental em multi-tenancy é o design das **chaves de identificação** dos limites. Normalmente utiliza-se identificadores únicos de cliente (ID do tenant, API key, IP do caller, etc.) para indexar os contadores. É importante concatenar ou separar por contexto – por exemplo,

chave no Redis como `limite:{tenant_id}:{serviço}` – para evitar colisões. Além disso, definir níveis de granularidade: podemos ter limite por API key, e um outro por IP de origem para evitar abuso de IPs diferentes com a mesma chave. Da mesma forma, limites diferentes por *endpoint*: e.g., talvez a API `POST /upload` tenha um limite mais restrito que o `GET /status`. Numa arquitetura robusta, podemos combinar várias dimensões de rate limit simultaneamente (por usuário, por IP, por serviço) para cobrir cenários de abuso mais complexos – embora isso aumente o número de verificações a cada requisição.

## Componentes de um SDK de Rate Limit em Go

Ao construir um SDK próprio de rate limiting em Go com suporte a multi-tenancy, podemos decompor a solução em diversos componentes e camadas de responsabilidade:

- **Definição de Políticas de Limite:** O SDK deve fornecer uma forma de definir os limites de cada categoria de cliente. Isso inclui parâmetros como *rate* (taxa permitida) e *burst* (tolerância a picos). Pode ser representado por uma estrutura – por exemplo, uma struct `Limit` contendo `Requests` por `Period` e talvez o *burst*. Bibliotecas existentes modelam isso; o `redis_rate`, por exemplo, define `Limit{Rate, Burst, Period}` permitindo criar facilmente limites “por segundo”, “por minuto” etc <sup>36</sup> <sup>37</sup>. O SDK deve permitir configurar limites diferentes para cenários diferentes, como planos distintos (gratuito, premium) ou diferentes chaves de API.
- **Identificação do Chave de Rate Limit:** Integrado ao SDK deve haver mecanismos para extrair a identidade relevante de cada requisição. Em serviços HTTP, isso pode ser:
  - Chave de API ou token do cliente (geralmente vindo em um header ou parâmetro) – útil para autenticar e aplicar limites por cliente autenticado.
  - Endereço IP de origem – para limitar usuários anônimos ou atuar como medida adicional contra abusos distribuídos. Importante considerar proxies: caso o serviço receba tráfego via proxy ou load balancer, usar `X-Forwarded-For` ou `X-Real-IP` para obter o IP real do cliente.
  - Identificador do tenant – em apps multi-tenant, pode haver um identificador do cliente no próprio caminho ou token que indica a qual cliente aquela requisição pertence. Por exemplo, um SaaS multi-tenant pode passar um header `X-Tenant-ID`.

O SDK pode oferecer uma interface para fornecer a chave de rate limit dada uma *request* (por ex, função callback configurável que extrai do contexto da requisição a chave certa). Em casos simples, a chave pode ser literalmente a API key fornecida.

- **Interface do Limitador (Rate Limiter):** O núcleo do SDK será uma abstração do limitador, exposta via métodos como `Allow()` ou `Take()`. Um possível design é uma interface `Limiter` com métodos:
  - `Allow(key string) bool` – verifica se a requisição identificada por `key` pode prosseguir (true = permitida, false = bloqueada). Pode também retornar informações extras, como contagem restante.
  - Alternativamente, um método `Wait(ctx)` se quiser implementar comportamento de *throttling* (aguardar até poder prosseguir), embora para APIs normalmente prefira-se negar imediatamente com erro 429 do que segurar a conexão.

Essa interface pode ter múltiplas implementações internas de acordo com o backend.

- **Backend de Armazenamento:** O SDK deve suportar pelo menos dois backends: **memória local** e **Redis** (ou outro store distribuído). Provavelmente com implementações que atendem a mesma interface `Limiter` mencionada:
- *MemoryLimiter*: mantém um mapa em memória de estado por chave (por ex., contadores, timestamps ou instâncias do algoritmo). Pode usar locks ou estruturas thread-safe para atualizações concorrentes. Em Go, por exemplo, poderíamos usar a biblioteca `golang.org/x/time/rate` para criar um `rate.Limiter` para cada chave de API e guardá-los em um `sync.Map`. Ou implementar manualmente token bucket/leaky bucket usando `time.Now` e variáveis atômicas.
- *RedisLimiter*: ao invés de armazenar local, em cada chamada de `Allow(key)` faz-se uma operação no Redis. O SDK aqui age como um *wrapper* das funcionalidades do Redis – pode usar scripts Lua para garantir atomicidade. Muitos implementam GCRA no Redis, mantendo para cada chave um timestamp de referência e um crédito restante. O pseudo-código pode ser: *pegar tempo atual do Redis, calcular novo teórico\_tempo\_chegada, comparar e setar se permitir*. A biblioteca `redis_rate` já implementa isso internamente: por exemplo, seu método `limiter.Allow(ctx, key, limit)` retorna `allowed` e `remaining` contando o novo estado <sup>38</sup>. O SDK pode simplesmente integrar essa biblioteca ou reutilizar sua lógica. Vale lembrar de configurar corretamente conexões Redis (pool de conexões, tratamento de erros, etc).

Abstrair o backend permite ao desenvolvedor usar o mesmo SDK em cenários diferentes – por exemplo, em um serviço simples, usar o modo memória; em um serviço distribuído crítico, plugar o modo Redis.

- **Módulos de Middlewares HTTP:** Para priorizar o uso em servidores HTTP Go (net/http puro ou frameworks como Gin e Chi), o SDK deve facilitar a integração como middleware. Isso significa fornecer funções ou handlers prontos que façam:
- Receber a requisição HTTP e identificar a chave do cliente (via o componente de identificação mencionado).
- Consultar o limitador (interface) para aquela chave. Exemplo em pseudocódigo:

```
if !limiter.Allow(apiKey) {
    // Excedeu limite
    w.WriteHeader(429)
    w.Write([]byte("Too Many Requests"))
    return // aborta a request
}
// senão, chama o próximo handler
next.ServeHTTP(w, r)
```

Em `net/http`, isso seria um wrapper do `http.Handler`. Em **Gin**, pode ser um `gin.HandlerFunc` que dá um `c.AbortWithStatus(429)` se estourou o limite, ou então `c.Next()` se permitido. De modo similar para **Chi** (que usa padrão de middlewares compatíveis com `net/http`). O SDK pode oferecer um middleware configurável, por exemplo: `RateLimitMiddleware := NewRateLimitMiddleware(limiter, extractKeyFunc)` para facilitar plugar no roteador. Ademais, é útil incluir nos handlers informações de header quando bloquear: ex. adicionar `Retry-`

`After` indicando em quantos segundos o cliente pode tentar novamente, ou cabeçalhos personalizados mostrando seu limite e quanto já usou – isso melhora a transparência para quem consome a API.

- **Diferenciação de Limites por Tier/Tenant:** Um requerimento chave é suportar limites variados conforme o plano do cliente. Na prática, isso significa que o SDK deve aplicar valores de limite diferentes dependendo de atributos do cliente. Há algumas formas de projetar isso:
  - **Instâncias distintas de Limitador:** Por exemplo, manter dois `Limiter` internos – um configurado com limite alto (p.ex. 1000 req/min) para tier *Premium*, outro com limite baixo (p.ex. 100 req/min) para tier *Free*. Na hora de processar a requisição, o middleware decide qual limiter usar baseado no tier do cliente (que poderia ser obtido do contexto de autenticação). Isso é simples, mas se há muitos tiers ou regras, pode escalar mal (múltiplos objetos para cada combinação).
  - **Parâmetro no Método:** Alternativamente, passar o limite desejado a cada chamada `Allow`. Por exemplo, a API do `redis_rate` permite especificar o limite desejado em cada chamada (`limiter.Allow(ctx, key, redis_rate.PerSecond(10))` para 10 rps)<sup>38</sup>. Assim, a lógica do SDK seria: determinar o perfil do cliente e então invocar `Allow` com os parâmetros adequados. Essa abordagem é flexível – as regras de tier podem vir de uma configuração externa ou banco, e basta mapear. Uma tabela de configuração poderia dizer: *tenant X (plano Gold): 5000/hora; tenant Y (plano Free): 1000/hora*, etc<sup>39</sup>.
  - **Chaveação das métricas:** Outra forma é embutir o tier na própria chave do rate limit e armazenar múltiplos limites no backend. Por exemplo, a chave no Redis poderia ser `premium:<idCliente>` vs `free:<idCliente>` e o código do limitador reconhecer pelo prefixo qual limite aplicar. Contudo, isso é menos direto e propenso a erro. É mais seguro manter a lógica de decisão de limite no código do SDK ou do serviço.

O importante é que o SDK facilite o desenvolvedor em **diferenciar clientes de tiers distintos**. Isso pode ser conseguido expondo uma interface para registrar políticas: e.g., `limiter.SetLimit(tenantID, LimitConfig)`. Em produção, essa informação de tier muitas vezes vem do sistema de cadastro/assinaturas – possivelmente o SDK poderia integrar uma chamada ou hook para buscar o limite do usuário (cacheando para performance). Em casos mais simples, define-se constantes para cada nível. O blog da DreamFactory dá um exemplo claro: “*tenants premium podem ter 10.000 req/hora, enquanto básicos 1.000*”<sup>39</sup> – o SDK deve suportar facilmente essa diferença.

- **Monitoramento e Telemetria:** Um SDK robusto também incorpora pontos de extensibilidade para registro de métricas e logs. Por exemplo, ao negar uma requisição por limite excedido, poderia emitir um log de aviso ou incrementar um contador de métricas (para sistemas como Prometheus). Métricas como quantidade de requisições bloqueadas por chave, taxa de utilização do limite, etc., são valiosas (ver seção de Métricas). O design poderia incluir callbacks ou simplesmente instruir o usuário do SDK a instrumentar essas partes (por exemplo, expondo um método para obter estatísticas atuais dos limiters, ou integrando com context para passar um struct onde se acumula dados).
- **Testabilidade:** Ao projetar o SDK, convém torná-lo fácil de testar. Isso significa, por exemplo, permitir injetar um relógio customizado (um *Clock interface*) para simular passagem do tempo – assim, testes unitários dos algoritmos podem avançar o tempo deterministicamente sem esperar segundos reais. A biblioteca do Uber (`ratelimit`) suporta injetar clock custom<sup>40</sup>; no

nossa SDK, poderíamos seguir prática similar. Igualmente, separar bem as camadas (lógica de cálculo vs I/O Redis vs middleware HTTP) facilita escrever testes isolados para cada parte.

Em resumo, um SDK de rate limiting em Go tipicamente se compõe desses módulos: *definição de limites*, *identificação de cliente*, *implementação de algoritmos local/distribuído*, *integração HTTP (middleware)*, *suporte a múltiplos tiers*, além de considerações de *observabilidade e teste*. A seguir discutiremos os trade-offs dessas escolhas e tecnologias disponíveis.

## Trade-offs e Decisões de Design

Ao implementar rate limiting, diversas decisões impactam o equilíbrio entre precisão, performance e complexidade. Abaixo discutimos alguns trade-offs importantes:

- **Precisão vs Simplicidade (Fixed Window vs Sliding Window):** A estratégia de **janela fixa** brilha pela simplicidade, usando operações atômicas simples (ex: `INCR` no Redis com TTL igual à janela)<sup>7</sup>. Porém, como vimos, ela pode permitir rajadas problemáticas nas bordas da janela. **Janelas deslizantes** evitam esses picos e fornecem controle mais justo e uniforme<sup>10</sup>, mas exigem computação extra – seja mantendo logs de timestamps (custoso em memória e CPU em altos volumes) ou calculando médias ponderadas entre janelas (exigindo lógica adicional e coordenação). O trade-off aqui é entre uma implementação trivial (que pode ser “boa o bastante” em muitos cenários) e uma implementação mais refinada que evita cenários adversos de burst. Muitas APIs públicas começaram com fixed windows por facilidade e eventualmente evoluíram para sliding windows para melhorar a experiência do cliente. Por exemplo, o GitHub API originalmente usa janela fixa de 1 hora (5000 requisições/hora) e sofreu do problema de resets abruptos<sup>41</sup>; suspeita-se que adicionaram depois um limite secundário de curto prazo para evitar saturação instantânea a cada reset<sup>42</sup>. Assim, se a carga do seu sistema apresenta picos precisamente nas trocas de janela ou se clientes maliciosos explorariam isso, vale adotar uma janela deslizante ou token bucket.
- **Permitir Bursts vs Suavizar Tráfego (Token Bucket vs Leaky Bucket/GCRA):** A escolha entre **Token Bucket** e **Leaky Bucket** reflete uma decisão de *política de burst*. Com token bucket, valorizamos permitir bursts esporádicos – a experiência do usuário tende a ser melhor, pois ele pode fazer um conjunto de requisições rápidas após inatividade (até o limite do bucket) sem ser bloqueado. Em compensação, o servidor pode enfrentar pequenas explosões de carga, que devem estar dentro do que consegue suportar. Já o leaky bucket e algoritmos equivalentes (GCRA) priorizam estabilidade do lado servidor: eles efetivamente espaçam as requisições, alisando picos. O trade-off está explicitado em uma comparação sucinta: *Token Bucket: “domínio dos bursts” (permite picos) mas requer ajuste fino da taxa e tamanho do bucket; Leaky Bucket: “saída suave” (nenhum burst significativo) mas não “ama” bursts – estes serão descartados*<sup>43</sup>. Em outras palavras, **Token Bucket** é cliente-friendly, **Leaky/GCRA** é server-friendly em termos de carga constante. A decisão deve considerar a natureza da aplicação: em APIs onde pequenas explosões de tráfego são aceitáveis e parte do uso normal (ex: chamadas em lote esporádicas), token bucket é indicado. Se for crucial evitar mesmo curtos sobrecarregos (ex: proteger um banco de dados sensível), leaky/GCRA podem ser mais apropriados. Vale notar que o GCRA, sendo equivalente ao leaky bucket, alcança suavidade similar porém é mais fácil de implementar corretamente em um sistema distribuído (devido à abordagem baseada em tempo).

• **Rejeitar vs Enfileirar (Limitar vs Throttle):** Ao atingir o limite, deve-se bloquear imediatamente as requisições extras, ou apenas retardá-las? O *rate limiting* estrito rejeita imediatamente com erro (HTTP 429) após o limite <sup>5</sup>, fornecendo um corte claro. Já um *throttling* adaptativo poderia, por exemplo, enfileirar algumas requisições ou inserir pequenas esperas antes de atender, ao invés de recusar sumariamente. Rejeitar é simples e alivia totalmente a carga extra (as requisições excedentes nem são processadas). Em contrapartida, causa uma experiência “dura” ao cliente, que recebe erros. Enfileirar ou atrasar melhora a experiência (tenta ainda servir todas as requisições sem erro, apenas mais lentamente), mas tem custos: requer buffer de requisições pendentes (consumindo memória) e pode elevar latências e até causar *timeouts*. Em sistemas multi-tenant, a escolha se inclina a depender do caso de uso: para **APIs públicas** com contrato de uso, costuma-se rejeitar com 429 – isso força o cliente a adaptar seu comportamento e não pune outros usuários nem congestiona servidores <sup>44</sup>. Já internamente, em integrações assíncronas ou filas de trabalho, um leve enfileiramento pode ser aceitável. Nosso SDK foca em rate limiting (rejeitar), mas vale mencionar que frameworks de *throttling* poderiam ser uma camada adicional caso se deseje garantir *throughput* sob carga ao custo de latência.

• **Local vs Distribuído:** Conforme discutido, usar um **backend distribuído (Redis)** vs **memória local** é um dos trade-offs centrais:

- *Local*: latência ínfima e sem dependências externas – ideal para performance máxima. Entretanto, não consegue, sozinho, uma visibilidade global numa arquitetura multi-nó. Aplica-se quando não se necessita de total rigor por cliente, ou quando cada cliente interage majoritariamente com uma única instância. Também é suscetível a perda de estado ao reiniciar instâncias (embora, para rate limiting de curto prazo, isso não seja catastrófico – resetar contadores ao reiniciar geralmente é ok).
- *Redis*: garante consistência em todo o cluster e persiste temporariamente os contadores, sobrevivendo a restarts de instâncias de aplicação. Por outro lado, cada requisição impõe um hit de rede e operação no Redis. Em cargas moderadas (p.ex. até alguns milhares de req/s), um cluster Redis robusto consegue lidar tranquilamente com essas operações; em escala gigantesca, esse pode virar o gargalo, exigindo particionamento ou otimizações. A robustez do Redis é crucial: implementar scripts Lua atômicos evita condições de corrida e garante exatidão mesmo sob clientes concorrentes <sup>29</sup>. Algumas plataformas (como a própria Redis através de extensões ou o DragonflyDB) já oferecem comandos dedicados de rate limit usando GCRA, otimizados internamente <sup>45</sup>, evidência de que essa abordagem distribuída é prática comum.

Em resumo, a decisão local vs distribuído é uma troca entre **velocidade** e **precisão global**. Muitos projetos adotam um híbrido: começam com limitadores locais (mais simples) e, conforme a necessidade de controle global cresce, migram para um Redis ou similar. Nosso SDK contempla ambos para maximizar flexibilidade.

- **Bibliotecas Prontas vs Implementação Própria:** O ecossistema Go oferece bibliotecas consagradas de rate limiting, cada uma com prós e contras:
- [golang.org/x/time/rate](https://golang.org/x/time/rate) – Implementa o algoritmo de token bucket no pacote oficial *x/time*. É altamente confiável, thread-safe e flexível (permite esperar por tokens via context ou simplesmente consultar disponibilidade) <sup>46</sup> <sup>47</sup>. Aceita configurar a taxa (como eventos/segundo) e o *burst* máximo <sup>48</sup>. Por ser parte do projeto Go, é bem otimizada. Trade-off: funciona somente em memória local (não tem suporte nativo a distribuição) e utiliza floats internamente para cálculo de tokens, o que adiciona leve complexidade. Mas para um único processo, é excelente – inclusive a Go Stdlib e Go Cloud SDKs frequentemente a usam. Se seu caso requer apenas limitar chamadas dentro

de um serviço individual, `x/time/rate` é indicado. Para usar com Redis, porém, seria necessário esforço manual de sincronização (p.ex. não trivial).

- `uber-go/ratelimit` – O pacote open-source da Uber fornece um limitador baseado em **leaky bucket** (balde furado) com uma API minimalista <sup>49</sup>. A ideia principal é que cada chamada a `Take()` retorna quando a requisição pode prosseguir, dormindo o tempo necessário se estiver adiantada <sup>50</sup>. Essa abordagem *blocking* é útil em contextos como limitar a taxa de operações de um loop ou de chamadas a um serviço externo, onde você prefere esperar ao invés de descartar. O foco desse lib é desempenho e simplicidade – por exemplo, ele recalcula a posição no bucket com base no tempo decorrido desde a última aquisição ao invés de usar ticks contínuos <sup>49</sup>. Na documentação ele é descrito como tendo API simples e overhead mínimo, em contraste ao `x/time/rate` que cobre mais casos complexos <sup>51</sup>. Trade-off: por ser bloqueante e não suportar contexto, pode não se encaixar direto em servidores HTTP onde desejamos não segurar a thread/goroutine muito tempo. Também é local apenas. Use-o se precisa de um *throttler* leve dentro de um serviço (ex: limitar o envio de logs ou chamadas a um cliente RPC).
- `go-redis/redis_rate` – Essa biblioteca integra com o cliente go-redis e implementa o algoritmo **GCRA** no Redis <sup>29</sup>. Ela lida com toda a complexidade de scripts Lua e compatibilidade de versão do Redis, expondo métodos simples como `Allow()` para verificar permissões. Pontos fortes: resolve o problema distribuído de forma atomicamente segura (usando recursos do Redis 3.2+), e suporta especificar limite por segundo, minuto, hora facilmente <sup>52</sup>. Internamente, segue uma abordagem robusta baseada no projeto `rwz/redis-gcra`, consagrado em outras linguagens. Trade-offs: adiciona a dependência do Redis (tanto no componente em si quanto no custo de manutenção de um Redis em infra) e implica que toda requisição precisa de chamada de rede. Além disso, se a aplicação já usa outro cliente Redis (por exemplo, um wrapper custom), pode precisar integrar. Porém, para muitos, esses são custos aceitáveis dada a confiabilidade de ter limites globais consistentes. Em benchmarks simples, a latência extra do Redis (geralmente <1ms LAN) é compensada pela garantia de proteção cluster-wide.
- Outras bibliotecas: Existe o antigo `juju/ratelimit` (token bucket em Go) e libs como `Throttled` (mencionada pela Stripe) que implementam GCRA e oferecem integração com stores (memória, Redis, Memcached) <sup>31</sup>. O Throttled foi bastante usado e testado em produção (Stripe, Heroku), embora hoje o go-redis/redis\_rate pareça ter assumido essa função com mais atualizações. Também há middlewares prontos para frameworks web (por ex., alguns projetos no GitHub oferecem *gin middleware rate limiters*), mas usar um SDK próprio permite customizar conforme necessidades específicas de multi-tenancy.

Em termos de decisão: se o objetivo é **não reinventar a roda**, pode-se usar `x/time/rate` para limites locais rápidos e `redis_rate` para limites distribuídos – nosso SDK inclusive poderia encapsular essas libs para o usuário final. Se optar por implementação própria dos algoritmos (por aprendizado ou requisitos especiais), tenha certeza de considerar as sutilezas já discutidas (ex.: atômicos no Redis, monotonic clock, etc.). Em geral, aproveitar bibliotecas maduras reduz risco de bugs nos algoritmos de rate limiting, mas pode-se abrir mão de um pouco de flexibilidade ou ter que adaptar o que já está disponível.

- **Integração com HTTP e UX do Cliente:** Uma decisão importante é *como informar e tratar o cliente* quando limites são atingidos. O padrão é retornar **HTTP 429** com talvez um cabeçalho `Retry-After`. Alguns provedores adicionam cabeçalhos como `X-RateLimit-Limit`, `X-RateLimit-Remaining` e `X-RateLimit-Reset` nas respostas, para permitir que os usuários monitorem seu consumo <sup>53</sup>. Existe um esforço de padronização IETF para definir cabeçalhos `RateLimit-*` padronizados (Limit, Remaining, Reset) <sup>54</sup>, o que indica uma tendência futura de uniformizar essa

comunicação. Decidir implementá-los agora seria um bônus de usabilidade. Entretanto, adicionar muitos cabeçalhos a cada resposta pode aumentar overhead (mínimo, mas existe) e nem sempre é trivial em frameworks – mas em Go puro ou Gin é bem possível configurar. Em qualquer caso, **documentar** os limites e **providenciar mensagens claras** é essencial para uma boa DX (developer experience). Boas práticas incluem também retornar um corpo JSON explicando o erro quando possível, e até endpoints para o cliente consultar sua cota restante <sup>55</sup>. Essas melhorias vão além do core técnico, mas fazem parte do design completo do sistema de rate limit.

- **Dimensões de Limite:** Como discutido, um design multitenant robusto pode envolver múltiplos critérios de limitação simultânea (por usuário, por IP, por região, etc). Cada critério extra adiciona custo (precisa manter contadores separados e checar todos). O trade-off aqui é **cobertura vs desempenho**. Por exemplo, se sua API é autenticada por chave sempre, talvez controle por chave seja suficiente e limitar por IP seja desnecessário (exceto para mitigar DDoS). Por outro lado, se há preocupação de um usuário distribuir suas chamadas entre muitos IPs (para burlar limites), pode ser válido implementar um limite secundário por IP. Um caso comum: limitar *IP anônimos* a, digamos, 60 req/min, mas se fornecer uma API key, permitir mais – combinando autenticação e IP. Na arquitetura do SDK, isso poderia ser feito em cascata: primeiro um middleware de IP (geral, sem considerar tenant) e depois um de API key. O risco é bloquear usuários legítimos atrás de NAT se o limite por IP for muito baixo. Portanto, deve-se ajustar esses valores e possivelmente permitir escapes (ex: IPs de provedores conhecidos podem precisar de limites maiores). Em resumo, mais dimensões = mais complexidade de configuração, mas maior robustez contra casos adversos.
- **Estado vs Stateless:** Rate limiting é inherentemente *stateful* (precisa lembrar do histórico de requisições). Em sistemas REST que idealmente são stateless, esse estado vai parar em algum lugar (memória ou Redis). Uma abordagem totalmente diferente são **tokens de cota (quota tokens)** emitidos ao cliente, que carregam informações de quantas requisições ele ainda pode fazer (talvez assinados digitalmente). Isso tornaria a verificação local (o cliente apresenta um token com sua cota restante). Porém, esse modelo é complexo e raro; a maioria das implementações opta por guardar estado no servidor mesmo. Mencionamos apenas como curiosidade arquitetural: sacrificar o stateless puro do REST é normal e necessário para implementar limites.

Em todas as decisões acima, o crucial é adequar a solução ao **perfil de tráfego e requisitos de fairness do seu sistema**. Não existe bala de prata – cada algoritmo e abordagem tem cenário ideal e trade-off. Muitas vezes, a implementação final é híbrida: ex., “*Fixed window de 1h para enforce de cota diária e token bucket por segundo para proteção em tempo real*”, ou “*Limiter global no gateway e limiter local nos serviços críticos*”. O importante é conhecer as opções (como fizemos) e combiná-las de forma coerente com os objetivos (proteger recursos, isolar tenants, cumprir contratos de SLA).

## Tecnologias e Integração em Serviços Go

**Linguagem Go** é bem servida de ferramentas para implementar rate limiting, e integrar isso a servidores HTTP é relativamente simples graças ao design de middleware da biblioteca padrão e frameworks:

- **net/http:** O servidor HTTP padrão do Go permite encadear middlewares manualmente. Nossa SDK pode expor, por exemplo, uma função `RateLimitHandler(next http.Handler)` que retorna um `http.Handler`. Esse handler internamente faz:

```

return http.HandlerFunc(func(w http.ResponseWriter, r *http.Request) {
    key := extractKey(r)           // identificar chave do cliente
    allowed := limiter.Allow(key)
    if !allowed {
        w.Header().Set("Content-Type", "text/plain")
        w.WriteHeader(http.StatusTooManyRequests)
        w.Write([]byte("Too Many Requests"))
        return
    }
    next.ServeHTTP(w, r)
})

```

O usuário do SDK então pode encadear isso na inicialização do servidor, e funcionará para qualquer roteamento. A biblioteca padrão não tem um middleware de rate limit embutido, mas pela flexibilidade do HandlerFunc, é trivial adicionar como acima.

- **Gin:** O Gin é um framework web popular que usa um conceito de *middleware* com o contexto `gin.Context`. O SDK pode oferecer um helper do tipo:

```

func RateLimitMiddleware(limiter Limiter, extractKey func(*gin.Context) string) gin.HandlerFunc {
    return func(c *gin.Context) {
        key := extractKey(c)
        if !limiter.Allow(key) {
            c.AbortWithStatusJSON(429, gin.H{"error": "rate limit exceeded"})
            return
        }
        c.Next()
    }
}

```

Isso aproveita o mecanismo do Gin de abortar requisições. Muitos desenvolvedores também optam por middlewares prontos (há projetos open-source para Gin rate limiters), mas escrever um próprio usando nosso SDK é simples. A vantagem é poder inserir lógica customizada – por exemplo, definir que certas rotas ou métodos HTTP tenham limitadores distintos.

- **Chi:** O Chi segue padrão parecido ao net/http para middlewares (funções que recebem e devolvem `http.Handler`). Podemos adaptar o mesmo `RateLimitHandler` para ser usado via `r.Use(RateLimitHandler)` em um roteador Chi. Alguns frameworks micro como **Fiber** ou **Echo** têm padrões similares – praticamente todos suportam a ideia de interceptar a requisição e abortá-la.
- **Uber-go/ratelimit** e **golang.org/x/time/rate**: Se quisermos usar essas libs dentro do nosso serviço, como mencionado antes, devemos adequá-las ao contexto HTTP. Por exemplo, se usarmos o `x/time/rate.Limiter`, podemos criar um mapa de `Limiter` por chave:

```

var limiters = sync.Map{} // chave -> *rate.Limiter
func getLimiter(key string) *rate.Limiter {
    actual, _ := limiters.LoadOrStore(key, rate.NewLimiter(10,
20)) // ex: 10 req/s, burst 20
    return actual.(*rate.Limiter)
}
// no handler:
lim := getLimiter(clientKey)
if !lim.Allow() {
    // negar 429
}

```

Esse código daria ao cliente identificado por `clientKey` um token bucket de 10 RPS com burst 20.

**Nota:** O `sync.Map` ajuda a criar limiters sob demanda por chave de forma thread-safe. Para políticas diferentes por tier, precisaríamos condicionar a criação com valores diferentes conforme o cliente.

Usando o `uber-go/ratelimit`, seria semelhante, exceto que a chamada retorna o instante permitido. O padrão de uso é:

```

rl := ratelimit.New(rate)      // cria limitador X operações/seg
now := rl.Take()              // bloqueia até poder emitir próximo

```

Integrar isso em HTTP significaria que no handler faríamos `rl.Take()`. Isso bloquearia a goroutine do request até o horário permitido. Na prática, isso atrasaria a resposta ao cliente ao invés de rejeitar – implementando *throttling* em vez de hard limit. Se o objetivo é não rejeitar mas sim *colar* requisições, poderia ser uma escolha. Porém, segurar a goroutine de atendimento do request (especialmente se for muitos segundos) não é ideal – arrisca esgotar *workers* se houver fila. Portanto, raramente se usa `ratelimit.New().Take()` em handlers web; ele é mais útil para controlar intervalo de chamadas de saída (e.g., chamando uma API externa com limite). Para rejeitar imediatamente com Uber's lib, não há um método pronto – teríamos que calcular se `Take` não espera (se retorna imediatamente) ou não, mas isso é indireto. Então preferimos x/time para esse caso.

- **Uso do Redis no código:** Ao optar por Redis, o código de integração via `redis_rate` se parece com:

```

rdb := redis.NewClient(&redis.Options{Addr: "localhost:6379"})
limiter := redis_rate.NewLimiter(rdb)
limit := redis_rate.PerSecond(10) // configura 10 req/s
res, err := limiter.Allow(ctx, "api:"+clientID, limit)
if err != nil { /* tratar erro Redis */ }
if res.Allowed == 0 {
    // bloquear 429
    // opcional: res.RetryAfter indica em segundos ou ms quando poderá de
}

```

```
novo  
}
```

No objeto `res` retornado há informações como `.Remaining` (quantos ainda pode no período atual) <sup>38</sup> e possivelmente um campo `RetryAfter`. Essa biblioteca facilita bastante – internamente ela já usa scripts Lua com GCRA atomicamente. Em termos de design, podemos usar uma **conexão Redis pool** compartilhada e reutilizar o `limiter` para todas as requisições. O uso de chaves no Redis requer padronização (prefixos por tipo de limite). É sábio também definir um TTL para as chaves de contadores para não ficar guardando para sempre usuários inativos – o GCRA normalmente atualiza a cada interação e mantém TTL de alguns períodos sem nova requisição.

- **Concurrency Limiting no código:** Em Go, limitar concorrência pode ser tão simples quanto usar um canal com buffer de tamanho = limite. Por exemplo:

```
semaphore := make(chan struct{}, 100) // permite 100 simultâneos
handler := func(w, r) {
    select {
    case semaphore <- struct{}{}:
        defer func(){ <- semaphore }()
        // processar request normalmente
    default:
        // semaphore cheio => rejeitar
        http.Error(w, "Too Many Concurrent Requests",
http.StatusTooManyRequests)
        return
    }
}
```

Esse snippet tenta inserir um token no canal (non-blocking select). Se o canal estiver cheio (100 em curso), cai no `default` e rejeita. Caso contrário, prossegue e garante com `defer` remover o token ao terminar, liberando vaga. Essa técnica implementa um limitador de concorrência eficaz com pouquíssimas linhas, aproveitando primitivos nativos do Go. O SDK poderia incorporar isso de forma configurável também. Alguns frameworks suportam nativamente limitar concurrency (p.ex. gin tem configuração de pool de workers), mas explicitamente codificar nos handlers oferece controle por escopo (talvez uma rota específica tenha limite de 5 concorrentes, enquanto global do servidor pode ser 100).

- **Tecnologias complementares:** Vale mencionar que em arquiteturas modernas, às vezes delega-se rate limiting a camadas como **Service Mesh** ou **Ingress Controller**. Por exemplo, o Envoy Proxy oferece um filtro de rate limit que se comunica com um serviço gRPC central de quotas. Istio e Linkerd podem implementar políticas de rate limit no mesh. Se a sua infraestrutura inclui essas peças, você poderia mover a responsabilidade para lá ao invés de em cada serviço (ou usar ambas para diferentes propósitos). Contudo, nosso foco aqui é um SDK em Go dentro do serviço. Apenas considere que as tecnologias de orquestração (Kubernetes + Istio etc.) têm recursos para isso, que às vezes simplificam do ponto de vista aplicativo. A escolha por um SDK próprio costuma vir da

necessidade de customização avançada (tiers, lógicas específicas de tenant) que soluções genéricas nem sempre atendem perfeitamente.

Em resumo, integrar rate limiting em serviços Go envolve usar as construções nativas de middleware para interceptar requests e aplicar a lógica do limitador (seja ele implementado manualmente ou via biblioteca). Tecnologias como Redis entram como apoio para compartilhamento de estado. Os frameworks Go são suficientemente flexíveis para adicionar esses controles sem esforço exagerado, garantindo que possamos, por exemplo, "controlar por API Key, IP e tenant" conforme solicitado – basta extrair essas informações e consultá-las em estruturas de dados apropriadas, conforme demonstrado nos exemplos acima.

## Boas Práticas na Implementação de Rate Limiting

Ao desenvolver e operar um sistema de rate limit multi-tenant, existem práticas recomendadas que melhoram a eficácia e evitam armadilhas comuns:

- **Definir Limites Apropriados e Documentá-los:** Escolha limites condizentes com a capacidade do sistema e o contrato com os clientes. Por exemplo, se usuários gratuitos devem ter 1.000 requisições/dia e pagos 10.000/dia, garanta que isso está alinhado à carga que sua infraestrutura suporta e às expectativas de uso. **Documente claramente** esses limites nas APIs públicas, para que desenvolvedores clientes possam programar retries ou saber quando reduzir a velocidade. A transparência evita frustrações.
- **Retornar Sinalização Adequada:** Utilize o código HTTP correto (429 Too Many Requests) para indicar bloqueio por limite <sup>5</sup>. Inclua cabeçalhos indicando o estado do limite – por exemplo, `X-RateLimit-Limit` (cota total), `X-RateLimit-Remaining` (restante na janela atual) e `X-RateLimit-Reset` (tempo de reset) como faz o GitHub <sup>53</sup>. Ou adote já os cabeçalhos propostos no padrão (`RateLimit-Limit`, etc.) para futuro compliance. Fornecer também um cabeçalho `Retry-After` com segundos até poder tentar novamente melhora a cooperação do cliente. Se retornar um corpo JSON, uma mensagem amigável do tipo `"error": "Rate limit exceeded, retry after 30 seconds"` é útil.
- **Evitar Limites Ocultos ou Arbitrários:** Todos os limites impostos deveriam ser intencionais e conhecidos. Evite ter, por exemplo, um limitador de concorrência interno que rejeita silenciosamente sem o cliente entender o motivo. Se implementar limites adicionais, trate-os similarmente (código 429, etc.). Consistência é importante: se a política é "X req por minuto", certifique-se de que seu sistema não bloqueia antes desse limiar (a não ser por outros motivos como autenticação).
- **Tiering e Isolamento:** Como discutido, separe claramente as políticas por tier. **Nunca** aplique um limite pensado para free users também aos premium inadvertidamente. Uma boa prática é ter estruturas ou configurações explicitamente segmentadas por plano, para não misturar. Por exemplo, mantenha constantes ou valores configuráveis: `FreeTierLimit = 1000/day`, `PremiumTierLimit = 10000/day`. Isso evita confusão e permite ajuste independente. Além disso, se for implementar *burst* maior para premium (por ex, podem fazer rajadas mais intensas), deixe isso refletido tanto no token bucket (capacidade maior) quanto talvez em limites de concorrência (permitir mais conexões simultâneas a premium se fizer sentido). Esse isolamento garante a **experiência diferenciada** que justifica tiers – e evita que um plano interfira no outro.

- **Noisy Neighbor e Fairness:** Além dos limites por cliente, avalie a necessidade de um limite **global** para proteger a infraestrutura. Por exemplo, você pode ter 100 tenants cada um com 100 req/s permitidas (teoricamente 10k req/s totais). Se todos atingirem o topo simultaneamente, talvez seu sistema não aguente 10k/sec. Nesse caso, pode-se introduzir um limite global (por cluster ou serviço) – e.g., aceitar no máximo 5k req/s somando todos. Esse limitador global atuaria como válvula de segurança (talvez retornando 503 Service Unavailable se ultrapassado, pois não é culpa de um cliente específico). Isso evita que mesmo todos seguindo seus limites individuais acabem saturando coletivamente o serviço. Tal mecanismo deve ser usado com cuidado para não anular os acordos com clientes, mas em ambientes multi-tenant onde a soma de picos pode exceder a capacidade, é prudente. Essa técnica relaciona-se ao conceito de *controle de admissão global*.
- **Monitoramento e Alertas:** Implemente métricas para acompanhar o funcionamento do rate limiting:
  - Conte o número de requisições bloqueadas por chave e total. Por exemplo, métrica Prometheus: `api_requests_blocked_total{tenant="X"}`.
  - Também meça o número de requisições atendidas. Assim, pode-se calcular a porcentagem de requisições bloqueadas por cliente – se um cliente está tendo, digamos, 30% de suas chamadas bloqueadas, talvez seja caso de analisar (pode ser um cliente que precisa de um plano maior ou está insatisfeito).
  - Monitore latência adicional: por ex, tempo gasto na chamada Redis para limitar, para garantir que não está introduzindo grande overhead.
  - Se usar filas (throttling), monitore o tamanho das filas e tempos de espera para detectar degradação.
  - Crie alertas: ex., se a taxa de 429s sobe repentinamente, pode indicar abuso ou um bug no lado do cliente. Também alerta internos se o Redis falhar ou taxa de erros de script Lua subir.
  - Registre logs de eventos de rate limit – p. ex., logue quando bloquear pela primeira vez um cliente em um período. Isso ajuda o suporte a identificar se um cliente foi bloqueado continuamente.

Ferramentas de APM/tracing também podem ser integradas: adicionar uma anotação no trace quando um request é atrasado ou bloqueado pode ajudar a debugar comportamentos em ambientes distribuídos.

- **Testes em Cenário de Carga:** Simular altas cargas e múltiplos tenants é crucial. Utilize ferramentas como `hey`, `ab` (ApacheBench) ou `JMeter` para enviar rajadas de requisições e verificar se o limitador responde conforme esperado. Teste casos limites:
  - Um único cliente excedendo seu limite: ele deve ser bloqueado mas outros clientes simultâneos não devem ser afetados (fairness).
  - Muitos clientes simultâneos no limite: o throughput total deve estabilizar na soma dos limites, e não quebrar o serviço.
  - Teste também períodos prolongados para ver se não há vazamentos de memória (especialmente se armazenando keys para cada cliente – veja se chaves expiradas estão sendo removidas).
  - Simule cenários de clock estranhos se possível (ex: se um servidor tem clock incorreto no caso de GCRA, ou viajar no tempo em teste).
- Se possível, escreva testes unitários para lógica do algoritmo com tempos simulados (p. ex., insira N requisições com intervalo X e veja se contagem bate com esperado).

- **Considerar Failures Gracefully:** Tenha uma estratégia para falhas do componente de rate limit. Por exemplo, se o Redis estiver indisponível, o que faz? Há duas filosofias: *fail-open* (permite todas requisições para não interromper serviço, mas arrisca sobrecarga) ou *fail-closed* (bloqueia tudo, protegendo backend mas negando serviço inclusive a bons clientes). Muitas implementações preferem *fail-open* – é melhor servir sem limites temporariamente do que derrubar o sistema ou retornar erro generalizado, especialmente se a janela de falha for curta. O SDK poderia retornar allowed=true em caso de erro de backend, possivelmente logando um erro. Também, use timeouts curtos nas chamadas ao Redis para não atrasar requests demais se ele estiver lento. Outra possibilidade: ter um *circuit breaker* – se detectar que o Redis está instável, automaticamente passar a permitir tráfego (e alertar a equipe para arrumar o Redis, claro).
- **Proteção contra Exploração:** Atacantes podem tentar explorar o próprio mecanismo de rate limit, por exemplo:
  - Enviar requisições com identificadores de API key aleatórios para preencher seu banco de dados/Redis com milhões de chaves (tentativa de DoS indireto consumindo memória). Para mitigar, só crie estado de rate limit após verificar que a API key é válida/autenticada. Ou use eviction no Redis (ex.: definir um TTL máximo razoável para chaves inativas).
  - Outra tática: usar um pool de IPs para escapar de limites por IP – aqui a combinação de limites por API key E por IP ajuda, como já dito. Mantenha também limites no próprio firewall ou CDN para casos extremos (Cloudflare, etc, podem bloquear IPs maliciosos antes de chegar no seu serviço).
  - Certifique-se que o rate limiter não possa ser bypassado por alguma rota alternativa no seu sistema (ex: se você aplica limitador no serviço A, mas o cliente pode sobrecarregar o sistema via serviço B que eventualmente chama A internamente sem check – se aplicável, coloque limites também nas chamadas internas ou no B).
- **Atualização Dinâmica de Configuração:** Em ambientes de produção, pode ser útil ajustar limites sem redeploy. Considere armazenar as configurações de rate (especialmente as de tiers) em um local centralizado (um arquivo de config, ou banco, ou variável de ambiente recarregável) e implementar capacidade de *hot reload*. Assim, se você perceber que o tier gratuito precisa de um ajuste (ex.: aumentar de 100 a 120 req/min), pode mudar e aplicar sem downtime. Mas cuidado: mudanças de limite afetam contratos – só aumente para beneficiar ou temporariamente, e reduções devem ser comunicadas com antecedência aos clientes se forem significativas.
- **Testes de Regressão e Cenários de Bordas:** Além de testes de carga, faça testes unitários para cenários pontuais:
  - Exatamente no limite: 100 requisições em 60s, a 101<sup>a</sup> deve bloquear.
  - Limite zero (caso bizarro, mas se você quiser suspender totalmente um cliente? O sistema lida com 0 req? Deve bloquear todas).
  - Mudança de relógio: se o servidor hora muda (daylight saving ou sync), o sistema continua consistente? (usar monotonic time da Go ajuda).
  - Vários nós simultâneos incrementando mesmo limite: teste com ambiente multi-thread no mesmo processo simulando (ou integração multi-processo) para ver se o Redis mantém atomicidade – muito provável que sim se bem feito, mas teste confirma.

- **Comunicação com o Cliente:** Fora do escopo técnico direto, mas relevante: oriente consumidores da API a implementar **exponencial backoff** ao receber 429. Clientes bem comportados, ao serem limitados, devem esperar um pouquinho mais a cada tentativa. Fornecer o `Retry-After` ajuda nesse backoff calculado. Se há um portal ou dashboard para usuários, mostrar consumo da API e limites restantes é uma ótima prática (por ex., o painel do GitHub mostra quantas requisições da API você já fez naquele hora).

Resumindo, adotar essas boas práticas garante que o rate limiting cumpra seu papel (proteger e isolar) sem pegar ninguém de surpresa ou introduzir novos problemas. A implementação técnica é metade do caminho; a outra metade é monitorar, ajustar e comunicar adequadamente.

## Armadilhas Comuns e Riscos

Mesmo com um bom design, há algumas armadilhas clássicas em sistemas de rate limit:

- **Bursts inesperados em limites mal configurados:** Se usar janela fixa, já citamos o risco de bursts duplos. Mesmo com token bucket, um bucket muito grande pode permitir um spike que sobrecarrega o sistema. Por exemplo, se você dá 1000 de burst acreditando ser inofensivo, um cliente inativo pode juntar 1000 tokens e despejar todos de uma vez – avalie se seu backend aguenta 1000 ops instantaneamente. Configure bursts conscientemente ou limite o tamanho do bucket.
- **Contagem duplicada ou perdida em distribuído:** Se implementar seu próprio algoritmo distribuído sem Lua/atomic, cuidado com race conditions. Ex: dois nós fazem GET do contador quase simultâneo e ambos veem valor X e incrementam para X+1 – resultando efetivamente em duas requisições permitidas a mais que o limite. O uso de incrementos atômicos ou scripts é fundamental para evitar essa condição de *lost update*. Além disso, atenção às bordas de expiração de chave no Redis: um padrão comum para janela fixa é usar `INCR` e se o retorno for 1 (era nova janela) então `EXPIRE key window`. Se ocorrerem múltiplos INCR paralelos bem no boundary, pode acontecer de configurar TTL múltiplas vezes ou não configurar. Scripts Lua que combinam INCR+EXPIRE de forma idempotente resolvem isso.
- **Uso de tempo não sincronizado:** Em algoritmos como GCRA, se cada nó usar seu próprio clock, divergências de milissegundos podem causar pequenas injustiças (um nó pode permitir uma requisição que outro teria negado, se o clock estiver defasado). Em casos extremos de drift grande, pode ocorrer false positives (bloquear sem motivo ou deixar passar demais). Por isso, se possível, baseie em um tempo central (p.ex. Redis `TIME` comando retorna hora única)<sup>30</sup>, ou ao menos sincronize NTP nos servidores. Em Go, preferir `time.Now()` (monotonic) ao invés de `time.Now().Unix()` para intervalos curtos, pois o monotonic clock não é afetado por ajustes de sistema.
- **Memory leak de chaves/tokens:** Se você tem milhões de usuários e mantém um limiter em memória para cada um, isso pode consumir muita RAM. Use estruturas com eviction: por ex, um LRU cache que descarta limiters não usados recentemente. No Redis, use TTL nas chaves de limite – se um usuário parar de fazer requisições, eventualmente a chave some, liberando espaço. Sem isso, você pode acabar com um Redis inchado contendo chaves antigas de usuários que nem utilizam mais o serviço. Monitorar a quantidade de chaves é importante.

- **Interferência com Retries:** Às vezes clientes implementam retries automáticos quando recebem certos erros. Se não souberem que 429 significa *slow down*, podem inadvertidamente piorar a situação. Imagine um cliente que ao receber 429 imediatamente tenta de novo várias vezes – podendo gerar um *storm*. Eduque clientes a não fazerem isso, ou implemente pequenos delays antes de responder 429 (não usual, mas Cloudflare por exemplo faz *token bucket* que de fato atrasa um pouquinho as requisições em vez de mandar tudo 429, para tentar acomodar). Essa nuance é mais de cliente, mas impacta o sistema.
- **Limitar componentes errados:** Garanta que o rate limiter está posicionado no local certo da arquitetura. Se colocado muito internamente (por ex, apenas em um serviço de backend, não no gateway), talvez ataques diretamente ao gateway ou a outro serviço passem sem controle. Por outro lado, se colocar no gateway mas ignorar que certos caminhos ficam fora (ex: upload de arquivos não passa pelo mesmo pipeline), pode ter brechas. Faça uma análise de onde um atacante ou um usuário intensivo poderia gerar carga e cubra esses pontos.
- **Rebaixamento do SLA sem perceber:** Um risco é definir limites arbitrariamente e afetar negativamente usuários legítimos. Por exemplo, se você tem um endpoint que normalmente requer 100 chamadas para concluir uma operação útil (talvez paginação, etc.), e coloca um limite de 50 req/min, estará forçando o usuário a demorar 2 minutos para conseguir completar a operação típica. Isso pode ser frustrante. Avalie padrões de uso comum antes de cravar limites. Uma boa prática é analisar logs de tráfego antes de ter rate limiting para descobrir quantas requisições típicas os top 95% dos usuários fazem, e definir limites um pouco acima disso para não impactar a maioria, focando apenas nos outliers.
- **Sobrecarga do próprio Rate Limiter:** Em cenários extremos, o rate limiting em si pode virar gargalo. Ex: se você usar um algoritmo de sliding log e milhões de usuários, a manipulação de listas enormes de timestamps pode comer CPU. Ou um Redis subdimensionado pode ficar CPU-bound processando incrementos. Assim, desenhe o limitador para ser leve – prefira algoritmos mais simples quando possível (ex: GCRA é O(1) por operação, enquanto sliding log é O(n) com n = reqs na janela). E dimensione a infraestrutura de suporte (um cluster Redis dedicado para isso, talvez replicado mas cuidado com replicação pois cada op de rate limit deve ir a todos nós se não usar *replicate\_commands* do Redis<sup>29</sup>). Ferramentas como *Pipeline* (enviar múltiplos comandos juntos) ou até bateladas de verificações podem ajudar se você precisa checar muitos limites diferentes por request.
- **Atualizações de configuração e sincronização:** Se você mudar um limite (digamos aumentar premium de 100 → 200 req/min) e tem muitos nós, todos devem estar cientes. Se for código, terá deploy. Se for config carregada em runtime, assegure-se que todos recarreguem (talvez um broadcast ou fetch periódico). Enquanto há dessintonia, uns nós poderiam aplicar antigo e outros novo – não é grave (pior que acontece é bloquear um pouco a mais até propagar), mas saiba dessa possibilidade.

Enfim, estar atento a essas armadilhas ajuda a evitar bugs ou comportamentos inesperados ao lançar seu SDK de rate limiting. Testes abrangentes e monitoração proativa pegam muitos desses pontos antes que afetem usuários finais.

## Métricas para Avaliar a Eficácia

Implementar o rate limiting é só o primeiro passo; é fundamental medir sua eficácia e impacto. Algumas métricas úteis a acompanhar incluem:

- **Requisições Totais vs. Bloqueadas:** A métrica básica é quantas requisições foram recebidas e quantas foram limitadas (bloqueadas). Idealmente expresse isso em porcentagem: *% de requisições bloqueadas*. Se este número for muito alto para um determinado cliente ou geral, indica uso abusivo (que está sendo corretamente restringido) ou que talvez o limite esteja aquém do necessário para uso legítimo (por ex., se 30% das requisições de clientes pagos são bloqueadas, talvez o limite premium esteja baixo demais para a demanda típica).
- **Distribuição de Consumo por Cliente:** Avaliar como diferentes tenants usam suas cotas. Por exemplo, calcular a média e o desvio padrão do uso de requisições por janela entre os clientes. Se poucos clientes estão sempre no limite e a maioria bem abaixo, isso pode estar ok. Mas se muitos estão no teto, talvez seja necessário rever limites ou reforçar infraestrutura (depende da estratégia de produto – às vezes induzir upgrade de plano é desejado). Ferramentas de monitoramento multi-tenant ou logs por chave podem ajudar aqui.
- **Taxa de 429s ao longo do tempo:** Visualizar em um dashboard a taxa de respostas 429 por minuto/hora. Deve correlacionar com eventos de pico ou comportamentos específicos. Um aumento súbito de 429s pode sinalizar um ataque ou um bug em cliente (ex: loop infinito tentando bater API e sendo bloqueado). Essa métrica ajuda a identificar comportamentos anômalos.
- **Throughput e Latência do Sistema:** Uma vez habilitado o rate limiting, monitore se o throughput do sistema se manteve dentro da expectativa e se a latência de processamento das requisições permitidas não degradou. Idealmente, sob cenários de abuso, o throughput útil se estabiliza próximo do limite configurado ao invés de colapsar – isso provará que o limiter está efetivamente protegendo o serviço. Métricas de latência P95/P99 também devem ficar estáveis (exceto possivelmente requisições bloqueadas que terminam mais rápido com 429).
- **Utilização de Redis (se aplicável):** Se usar Redis, monitore o **QPS (queries per second)** no Redis por operação de rate limiting, e uso de CPU/memória desse nó. Essa visibilidade garante que o Redis não está saturado. Por exemplo, se cada request gera 1 operação e você tem 5k req/s, terá 5k ops/s no Redis – o que é considerável mas ainda factível para um Redis moderno. Mas se houver scripts muito pesados ou vários incrementos, isso pode subir. Utilize monitoramento Redis (INFO command, etc) para ver se precisa escalá-lo.
- **Métricas de Espera/Queue (se throttling):** Caso tenha optado por enfileirar requisições excedentes ao invés de rejeitar, meça o tamanho médio e máximo das filas de espera e o tempo médio de espera. Por ex, *"quantas requisições estão aguardando atendimento agora"*. Isso indica o nível de degradação sob carga – uma fila crescendo muito sugere que o sistema está constantemente sobrecarregado. Um tamanho fixo de fila (ex: usar `channel` com buffer fixo) também implica que requests além daquele buffer serão descartadas ou sofrerão erros, então monitore *dropped due to queue overflow* se for o caso.

- **Métricas de Conformidade de SLA:** Se existem SLAs de disponibilidade/performance ligados a tiers (ex: premium deve ter 99.9% requests atendidas, free best effort), as métricas acima alimentam isso. Por exemplo, você pode calcular a disponibilidade percebida por um cliente = (requisições atendidas / requisições totais do cliente). Se premium estiver tendo porcentagem de bloqueio muito maior que zero (exceto por abuso claro), talvez seu tier não esteja cumprindo o prometido implicitamente.
- **Custos e Utilização:** Em plataformas cloud, cada requisição consome recursos (CPU, memória, etc) e potencialmente custo. Ao limitar, você espera reduzir uso excessivo. Monitorar se houve impacto em métricas de custo – e.g., redução de CPU por não processar cargas indevidas. Claro que isso depende do modelo de negócio (alguns oferecem extra throughput pago). Mas num multi-tenant onde um preço fixo dá direito a X, limitar pode evitar custo de infraestrutura de outliers. Logo, mensurar *CPU load, throughput vs. sem limiter* nos tests A/B ajuda a provar ROI do limiter.
- **Métricas de Experiência do Usuário:** Indiretamente, acompanhe métricas como número de erros relatados, satisfação, etc. Se o rate limiting estiver muito restritivo, clientes podem reclamar ou abandonar. Se for bem calibrado, deve proteger o serviço sem gerar muitos tickets de suporte. Por vezes, medir quantos usuários subiram de plano (upgrade) devido a atingirem limite também é interessante do ponto de vista de negócios.
- **Cross-metrics:** Combine dados de logins ou uso funcional com rate limit. Ex: se sempre que um certo endpoint intensivo é chamado ele precede um estouro de limite, talvez precise otimizar aquele endpoint ou dar um limite separado para ele. Ou se um tenant específico tem repetidos bloqueios, equipe de sucesso do cliente pode contatá-lo para entender se precisa de um plano maior.

Em suma, as métricas de rate limiting vão além de apenas contagens – elas demonstram se o sistema está **efetivamente protegendo** (ao limitar abusos) e **servindo de forma justa** (sem prejudicar indevidamente bons usuários). Essas medições também fornecem **insights para ajustes**: talvez um limite está muito folgado (nenhum usuário chega nem perto, poderia reduzir e economizar recurso) ou muito apertado (todos batem no teto, talvez aumentá-lo ou subdividir por funcionalidades). O acompanhamento contínuo permite evoluir as políticas de controle conforme o uso real do sistema.

## Estudos de Caso e Referências Públicas

Para ilustrar conceitos de rate limiting multi-tenant em prática, vejamos exemplos de plataformas conhecidas:

- **GitHub API:** O GitHub impõe limites tanto para usuários autenticados quanto não autenticados. Historicamente, a API REST v3 permite **5000 requisições por hora** para chamadas autenticadas por token (e apenas 60/hora para chamadas sem autenticação) <sup>53</sup>. Esse é um exemplo de **janela fixa** (1 hora) por usuário/token. Ao atingir 5000, qualquer chamada adicional retorna 403 ou 429 até a hora completar e os contadores resetarem. Eles informam os cabeçalhos `X-RateLimit-Limit`, `X-RateLimit-Remaining` e `X-RateLimit-Reset` em cada resposta, para o cliente saber sua situação. Como mencionado, um problema desse modelo é que um token pode consumir todas as 5000 de uma vez e só depois disso será bloqueado, causando um burst potencial no servidor. Brandur (Stripe) demonstrou exatamente isso: ele consumiu 5000 requisições rapidamente e então

ficou bloqueado ~51 minutos até o reset <sup>56</sup> <sup>57</sup>. Para mitigar abusos extremos, o GitHub implementou uma camada extra: há indícios de um limite mais curto para bursts – possivelmente um token bucket interno que não é explicitamente documentado <sup>42</sup>. Isso impede, por exemplo, que um cliente faça 5000 requisições *em um único minuto* e sobrecarregue o sistema, mesmo que a política horária permita. O GitHub assim exemplifica a combinação de estratégias: uma cota longa (diária/horária) e restrições de curto prazo para suavizar. Além disso, GitHub oferece diferentes limites para aplicações integradas via OAuth (que podem ter cotas separadas) e GraphQL API (que usa um sistema de pontos por consulta). Eles também têm tiers implícitos – usuários enterprise ou tokens de GitHub Apps podem ter limites expandidos. O caso do GitHub mostra a necessidade de adaptar e comunicar (a doc oficial detalha os limites e recomenda tratar código 403/429 com backoff).

- **Stripe API:** A Stripe, provedora de pagamentos, também lida com diferentes tenants (cada desenvolvedor da API Stripe é um cliente). Segundo um artigo do Brandur (que foi engenheiro na Stripe), eles implementaram um sofisticado rate limiting baseado em **GCRA** no serviço deles <sup>26</sup> <sup>31</sup>. O Stripe precisava de algo robusto porque múltiplos produtos (API de pagamentos, API de connect, etc.) compartilham infraestrutura. Eles optaram pelo GCRA para ter suavidade e usaram a biblioteca *Throttled* (open-source) para implementação em Go <sup>31</sup>. O fato de mencionarem que está “rodando em produção no Stripe e em breve na Heroku” <sup>31</sup> indica confiança no método. A Stripe não publica claramente suas cifras de limite (até por segurança), mas é sabido que existem restrições – por exemplo, limites de requisição por segundo nas APIs e possivelmente limites mensais por conta dependendo do contrato. Um interessante ponto é que, por lidarem com dinheiro, a estabilidade do serviço é crítica – então eles priorizam evitar cascatas de falha isolando abusos. O uso de GCRA no Stripe mostra a preferência por não permitir grandes bursts que poderiam afetar operações financeiras. Também é notável o foco em *distribuído* – clientes da Stripe podem ter múltiplos servidores mandando requisições concorrentes, então o controle global via Redis (ou internal store) é vital. Em suma, o Stripe exemplifica um caso de multi-tenant (cada integrador é um tenant) com tiers possivelmente implícitos (talvez contas maiores tenham limites negociados) e algoritmo avançado para fairness.
- **APIs de Cloud (AWS, Google Cloud):** Os provedores cloud geralmente têm *quotas* e *throttling* para cada conta (tenant). Por exemplo, o AWS Lambda impõe um limite simultâneo de X execuções para usuários básicos, e mais para enterprise – e aplica *reservas* para garantir que um usuário não consome tudo de um pool compartilhado. No contexto de APIs web, o Amazon API Gateway permite definir planos (usage plans) com rate limit e burst por API key. Esses gateways conseguem isolar tenants facilmente por chave, e oferecem dashboards de uso. Um blog da AWS cita: “*Multitenant REST APIs em escala usando API Gateway*”, detalhando criar chaves separadas por tier e associar políticas de throttling (por ex., Free = 10 req/s, Basic = 50 req/s) <sup>58</sup>. Isso mostra um caso prático de tiering: a infraestrutura trata cada tier de forma diferenciada, e também protege globalmente. Em clouds, muitas vezes o enfoque é mais em quotas diárias/mensais (para faturamento) e throttle global para estabilidade.
- **APIs públicas diversas:** É comum ver planos gratuitos vs pagos. O Twitter, por exemplo, na API antiga tinha limites por 15 minutos (janela fixa) – ex: 450 tweets de busca/15min autenticado, 150/15min não autenticado, e para cada endpoint família havia limites separados. O GitLab oferece 600 req/min por token de API. O Slack aplica limites cerca de ~1 req/s por token em média, com

bursts limitados (documentam taxa de *20 req/min* de slack API WebClient, por exemplo). Esses casos reforçam:

- *Distinção de contexto*: muitas APIs grandes diferenciam limites por método/endpoint (p. ex., ler dados pode ter um limite, escrever outro).
- *Tiers bem definidos*: geralmente gratuito vs pago. Alguns chegam a oferecer “*upgrade de limite via pagamento*” como produto.
- *Uso de códigos 429 e mensagens*: todas provêm um aviso quando excede. O GitHub retorna um corpo JSON com mensagem `You have exceeded your rate limit.` e informa quando poderá tentar de novo via Unix timestamp no header.
- **Estudos Acadêmicos e Padrões**: O problema de rate limit em multi-tenant também aparece em literatura de algoritmos de escalonamento e fairness. Existem teorias de *fair queuing* (enfileiramento justo) que são análogas a rate limiting por peso de usuário – às vezes mencionado em contextos de rede (roteadores distribuindo banda entre fluxos). Em sistemas distribuídos, há pesquisas em *adaptive rate limiting*, onde o sistema detecta automaticamente quando precisa apertar ou afrouxar limites conforme a saúde geral. Por enquanto, a maioria das implementações práticas usa limites estáticos configurados, mas podemos esperar tendências de mais automação.

Esses casos e referências mostram que controlar taxa por cliente é uma prática estabelecida em praticamente todos os grandes serviços web. Cada um ajusta as técnicas às suas necessidades: GitHub focou em simplicidade inicialmente (janela fixa) e complementou para bursts; Stripe investiu em algoritmo mais avançado desde cedo; Cloud providers entregam flexibilização via planos configuráveis; outros escolhem janelas curtas vs longas conforme convém. O denominador comum é tratar clientes diferentes de forma isolada e proteger a estabilidade do sistema como um todo – exatamente os objetivos que definimos no início.

## Tendências Futuras e Considerações Finais

O campo de rate limiting continua evoluindo, especialmente conforme crescem demandas de escala e complexidade de ambientes multi-tenant. Algumas tendências e aspectos emergentes:

- **Padronização de Comunicações de Rate Limit**: Como citado, está em andamento um padrão IETF para cabeçalhos de rate limit (RateLimit-Limit, RateLimit-Remaining, RateLimit-Reset)<sup>54</sup>. Espera-se adoção crescente, tornando mais fácil para clientes e ferramentas entenderem e reagirem a limites de forma uniforme (similar a como códigos HTTP padronizados permitem comportamento consistente). Nossa SDK no futuro poderia adotar automaticamente esses cabeçalhos ou formatos de resposta para estar em linha com a especificação.
- **Rate Limiting Dinâmico/Adaptativo**: Hoje, os limites geralmente são estáticos (configurados pelo operador). No futuro, poderemos ver sistemas que ajustam automaticamente as quotas conforme condições. Por exemplo, durante horários de baixa utilização global, o sistema poderia temporariamente elevar limites para todos (dando um “bônus” de throughput) e reduzi-los quando os recursos estiverem sob pressão<sup>35</sup>. Isso já é discutido em contextos de cluster Kubernetes (tornar limites dependentes de uso de CPU). Algoritmos adaptativos poderiam também aprender padrões de um cliente e alojar dinamicamente um pouco mais se ele nunca abusa, e restringir mais

cedo se detecta comportamento errático. Tudo isso, porém, aumenta a complexidade e potencial imprevisibilidade – é uma linha tênue entre ser adaptável e confundir os usuários com limites flutuantes.

- **Integração com Machine Learning e Detecção de Anomalias:** Relacionado ao ponto acima, futuras implementações podem incorporar ML para detectar *anomalies* no tráfego em tempo real e aplicar *rate limiting polimórfico*. Por exemplo, se um cliente de repente dobrar a taxa fora do seu perfil usual, um modelo ML poderia sinalizar e temporariamente reduzir seu limite para mitigar possível abuso, antes mesmo de atingir o limite fixo tradicional. Ou distinguir entre tráfego humano normal e rajadas automatizadas e aplicar políticas diferentes.
- **Limites granulares e políticas personalizadas:** À medida que plataformas crescem, oferecer **controle de rate limit customizado por cliente** pode virar um diferencial. Imagine permitir que um cliente enterprise defina ele próprio sub-limites para suas integrações, via nosso SDK. Isso requer expor configurações e talvez APIs de gerenciamento de limites. Estamos indo além de *um tamanho serve todos* – cada tenant poderia ter políticas próprias (dentro de um máximo permitido, claro).
- **Melhorias em Backends de Armazenamento:** O uso de Redis já é consolidado, mas novas tecnologias estão surgindo. Por exemplo, alguns bancos NoSQL oferecem *counters distribuídos* de alta performance nativamente. O Redis, por sua vez, continua aprimorando – a comando **CL.THROTTLE** no Redis (inspirado no algoritmo de leaky bucket/GCRA) facilita implementar limiters no lado do servidor <sup>45</sup>. Há também o **RedisBloom** com comandos *cf.reserve* que podem ser usados em gating, embora não especificamente para rate. Além disso, surgem caches distribuídos ultrarrápidos, e até soluções de rate limit in-memory replicadas (usar CRDTs para counters, por exemplo). Tudo isso pode, no futuro, reduzir a latência e overhead de manter limites globais.
- **Serverless e Multicloud:** Em arquiteturas *serverless*, onde não há servidor fixo para implementar logicamente o limiter, a tendência é usar camadas de API Gateway ou Functions pre/post hooks que aplicam limites. Ferramentas como Cloudflare Workers permitem escrever limiters em edge locations geograficamente distribuídas. Isso traz o conceito de rate limiting **no edge**, próximo do cliente, reduzindo carga no core. A contrapartida é sincronizar limites entre múltiplos edges – alguns usam soluções como replicação global ou redução da janela para garantir que mesmo se cada edge aplica local, no agregado se aproxima do limite pretendido.
- **Considerações de GDPR/Privacidade:** Um ponto sutil, mas: armazenar identificadores de usuários (como API keys ou IPs) e contadores configura um tipo de dado. Em jurisdições com leis de privacidade, talvez seja necessário expurgar dados de usuários que cancelam o serviço (o que incluiria suas chaves no Redis) ou anonimizar logs. É algo a se atentar em políticas internas – o rate limit em si não é foco da lei, mas os dados derivados não devem escapar das práticas de retenção. Oferecer um meio de limpar os dados de um tenant (por ex., remover suas keys do limiter quando ele sai) pode ser relevante no futuro compliance.
- **Combinação com Políticas de Segurança:** Rate limiting está se misturando cada vez mais com mitigação de ataques (DDoS, brute-force). Espera-se ver mais ferramentas unificando esses controles – por exemplo, um *WAF* (firewall de aplicação) pode trazer não só regras de bloqueio por padrão de ataque, mas também algoritmos de rate limit adaptativos para IPs suspeitos. Ou seja, a linha entre *security* e *fair usage* se mescla. Nosso SDK poderia futuramente integrar com sistemas de

deteção de intrusão para, digamos, aplicar limites muito mais baixos temporariamente a IPs sob suspeita de ataque, independentemente do tier.

Em conclusão, **Rate limiting** continuará sendo um alicerce para construir serviços escaláveis e seguros. Neste documento, exploramos em detalhes os problemas de estabilidade em microsserviços multi-tenant e as soluções clássicas (fixed window, sliding window, token bucket, leaky bucket, GCRA, limitador de concorrência) que fornecem um toolkit para controlar excesso de requisições. Discutimos como implementar essas estratégias em Go, seja via memória local ou Redis, aproveitando bibliotecas open-source de referência como `x/time/rate`, `uber-go/ratelimit` e `go-redis/redis_rate` e considerando integrações em frameworks web comuns. Também abordamos as nuances de suportar múltiplos tenants e tiers, garantindo que clientes de níveis diferentes sejam tratados de forma isolada e justa [3](#) [39](#). Analisamos decisões arquiteturais, componentes de um SDK customizado e uma série de boas práticas para evitar ciladas e medir sucesso – desde testes rigorosos até monitoramento contínuo.

Em última análise, um sistema de rate limiting bem projetado age como um **dreno de alívio de pressão**: evita que a voracidade de alguns degrade a estabilidade para todos, assegurando que recursos compartilhados sejam usados de forma civilizada. Para desenvolvedores, isso significa microsserviços resilientes mesmo sob rajadas imprevisíveis; para o negócio, significa poder oferecer diferentes níveis de serviço (free vs premium) de forma controlada; e para os usuários, significa uma experiência mais consistente e previsível, onde erros por sobrecarga se tornam raros. Com a contínua expansão de APIs e serviços multi-tenant, as técnicas e ferramentas discutidas aqui serão cada vez mais valiosas – e certamente continuarão a evoluir com novas ideias e padrões, mantendo o delicado equilíbrio entre **performance, equidade e segurança** no consumo de serviços.

## Referências

- Vyom Srivastava, “*Different Algorithms to Implement Rate Limiting in APIs*” – *Nordic APIs Blog*, 2021 [59](#) [60](#) . (Discussão de métodos de rate limit – throttling, GCRA, leaky bucket, etc.)
- Ashish Pratap Singh, “*Rate Limiting Algorithms Explained with Code*” – *AlgoMaster Newsletter*, 2024 [61](#) [8](#) . (Explica Token Bucket, Leaky Bucket, Fixed Window, Sliding Log e Counter com prós/cons e exemplos)
- Brandur Leach, “*Rate Limiting, Cells, and GCRA*” – *brandur.org Tech Blog*, 2015 [26](#) [30](#) . (Descreve limitações de algoritmos simples e introduz GCRA; cita uso na Stripe e melhora via Throttled)
- DreamFactory, “*Rate Limiting vs Throttling: Multi-Tenant API Use Cases*” – *DreamFactory Blog*, 2023 [62](#) [39](#) . (Aborda estratégias multi-tenant, tiers de serviço, noisy neighbor e melhores práticas de monitoração e adaptação de limites)
- Go pkg.dev Documentation – [go.uber.org/ratelimit](#) [49](#) [63](#) , [golang.org/x/time/rate](#) [48](#) , [go-redis/redis\\_rate](#) [29](#) . (Documentação das bibliotecas mencionadas, descrevendo algoritmos e uso)
- Jones Charles, “*Implementing Rate Limiters in Go: Token Bucket and Leaky Bucket Made Simple*” – *Dev.to*, 2025 [43](#) . (Tutorial prático em Go com analogias e código, comparando abordagens Token vs Leaky)
- DreamFactory, “*Rate Limiting in Multi-Tenant APIs: Key Strategies*” – *DreamFactory Blog*, 2023 [3](#) [35](#) . (Cobertura de isolamento por tenant, políticas por usuário/endpoint, e necessidade de distribuição em microsserviços)
- Stack Overflow – “*Token bucket vs Fixed window (Traffic Burst)*”, 2019 [8](#) . (Discussão das diferenças de burst em janela fixa vs token bucket)
- GitHub API Documentation – *Rate limiting* (referência de limites e cabeçalhos X-RateLimit do GitHub) [53](#) .

- IETF Draft – “*RateLimit Header Fields for HTTP*”, 2021 <sup>54</sup>. (Especificação emergente definindo cabeçalhos HTTP para comunicar limites de forma padrão)
- 

[1](#) [7](#) [24](#) [25](#) [26](#) [27](#) [28](#) [30](#) [31](#) [41](#) [42](#) [53](#) [56](#) [57](#) Rate Limiting, Cells, and GCRA — brandur.org  
<https://brandur.org/rate-limiting>

[2](#) [3](#) [4](#) [5](#) [34](#) [35](#) [39](#) [44](#) [55](#) [62](#) Rate Limiting vs Throttling: Multi-Tenant API Use Cases  
<https://blog.dreamfactory.com/rate-limiting-vs-throttling-multi-tenant-api-use-cases>

[6](#) [8](#) [9](#) [10](#) [11](#) [12](#) [13](#) [14](#) [15](#) [16](#) [17](#) [18](#) [19](#) [20](#) [21](#) [61](#) Rate Limiting Algorithms Explained with Code  
<https://blog.algomaster.io/p/rate-limiting-algorithms-explained-with-code>

[22](#) [23](#) [43](#) Implementing Rate Limiters in Go: Token Bucket and Leaky Bucket Made Simple - DEV Community  
[https://dev.to/jones\\_charles\\_ad50858dbc0/implementing-rate-limiters-in-go-token-bucket-and-leaky-bucket-made-simple-5162](https://dev.to/jones_charles_ad50858dbc0/implementing-rate-limiters-in-go-token-bucket-and-leaky-bucket-made-simple-5162)

[29](#) [36](#) [37](#) [38](#) [52](#) redis\_rate package - [github.com/go-redis/redis\\_rate/v9](https://github.com/go-redis/redis_rate/v9) - Go Packages  
[https://pkg.go.dev/github.com/go-redis/redis\\_rate/v9](https://pkg.go.dev/github.com/go-redis/redis_rate/v9)

[32](#) [59](#) [60](#) Different Algorithms to Implement Rate Limiting in APIs | Nordic APIs |  
<https://nordicapis.com/different-algorithms-to-implement-rate-limiting-in-apis/>

[33](#) Throttling a tiered, multi-tenant REST API at scale using API Gateway  
<https://aws.amazon.com/blogs/architecture/throttling-a-tiered-multi-tenant-rest-api-at-scale-using-api-gateway-part-1/>

[40](#) [49](#) [50](#) [51](#) [63](#) ratelimit package - [go.uber.org/ratelimit](https://go.uber.org/ratelimit) - Go Packages  
<https://pkg.go.dev/go.uber.org/ratelimit>

[45](#) Redis CL.THROTTLE Command (Documentation) - Dragonfly  
<https://www.dragonflydb.io/docs/command-reference/strings/cl.throttle>

[46](#) [47](#) [48](#) rate package - [golang.org/x/time/rate](https://golang.org/x/time/rate) - Go Packages  
<https://pkg.go.dev/golang.org/x/time/rate>

[54](#) RateLimit Header Fields for HTTP - IETF  
<https://www.ietf.org/archive/id/draft-polli-ratelimit-headers-02.html>

[58](#) Throttling a tiered, multi-tenant REST API at scale using API Gateway  
<https://aws.amazon.com/blogs/architecture/throttling-a-tiered-multi-tenant-rest-api-at-scale-using-api-gateway-part-2/>