

# Rate Limiting

Rate limiting for APIs is a technique used to control the amount of incoming and outgoing traffic to or from a web service. It is implemented to ensure that the API can handle the load and to protect against abuse, such as denial-of-service (DoS) attacks.

## Purpose of Rate Limiting

- **Preventing Overload:** It helps to manage the load on the server by limiting the number of requests a user can make in a given time period. This ensures that the server remains responsive and stable.
- **Fair Usage:** Ensures fair usage by preventing a single user or a small group from consuming all the available resources.
- **Security:** Protects the API from abusive behaviors such as brute-force attacks, scraping, or any other automated processes that could exploit the service.
- **Cost Management:** Helps manage operational costs associated with bandwidth, compute power, and other resources.

## How Rate Limiting Works

Rate limiting can be implemented in various ways, commonly using the following methods:

- **Fixed Window:** Limits requests within a fixed time window (e.g., 100 requests per hour). If the limit is reached, subsequent requests within that period are denied.
- **Sliding Window:** A more dynamic approach where the time window "slides" over time. It calculates the allowed requests over a continuously moving window, providing smoother rate limiting.
- **Token Bucket:** Each request consumes a token from a bucket. Tokens are added to the bucket at a constant rate, and if the bucket is empty, further requests are denied until more tokens are added.
- **Leaky Bucket:** Similar to the token bucket but with a fixed rate of processing requests. Excess requests are queued and processed at a fixed rate, ensuring a constant flow.

## Implementation Considerations

- Headers: APIs often include headers in responses to inform clients of their rate limit status (e.g., X-RateLimit-Limit, X-RateLimit-Remaining, X-RateLimit-Reset).
- Penalties: Once the limit is exceeded, APIs might return a specific HTTP status code like 429 (Too Many Requests) to indicate that the user needs to wait before making more requests.
- Granularity: Rate limits can be applied per user, per API key, per IP address, or any other identifier.

## Adding Rate limiting to a Spring Boot project using “Bucket4j” Library

To add rate limiting to a Spring Boot project, you can use several approaches. One popular method is to use the [Bucket4j](#) library,

### 1. Add Dependencies:

Add the Bucket4j and Redis dependencies to your pom.xml if you're using Redis for distributed rate limiting. Otherwise, Bucket4j can work with in-memory storage.

```
<dependencies>
  <dependency>
    <groupId>com.github.vladimir-bukhtoyarov</groupId>
    <artifactId>bucket4j-core</artifactId>
    <version>7.0.0</version>
  </dependency>
  <dependency>
    <groupId>org.springframework.boot</groupId>
    <artifactId>spring-boot-starter-redis</artifactId>
  </dependency>
  <!-- Other dependencies -->
</dependencies>
```

## 2. Configure Rate Limiting:

Create a rate limiting configuration class using Bucket4j.

```
import io.github.bucket4j.Bucket;
import io.github.bucket4j.Bucket4j;
import io.github.bucket4j.Refill;
import io.github.bucket4j.Bandwidth;
import org.springframework.context.annotation.Bean;
import org.springframework.context.annotation.Configuration;
import org.springframework.web.filter.OncePerRequestFilter;

import javax.servlet.FilterChain;
import javax.servlet.http.HttpServletRequest;
import javax.servlet.http.HttpServletResponse;
import java.time.Duration;

@Configuration
public class RateLimitingConfig {

    @Bean
    public Bucket bucket() {
        Refill refill = Refill.greedy(10, Duration.ofMinutes(1));
        Bandwidth limit = Bandwidth.classic(10, refill);
        return Bucket4j.builder().addLimit(limit).build();
    }

    @Bean
    public OncePerRequestFilter rateLimitingFilter(Bucket bucket) {
        return new OncePerRequestFilter() {

            @Override
            protected void doFilterInternal(HttpServletRequest request,
            HttpServletResponse response, FilterChain filterChain) {
                if (bucket.tryConsume(1)) {
                    filterChain.doFilter(request, response);
                } else {

                    response.setStatus(HttpServletResponse.SC_TOO_MANY_REQUESTS);
                    response.getWriter().write("Too many requests");
                }
            }
        }
    }
}
```

```
};  
}  
}
```

### 3. Apply the Filter:

Ensure that the filter is registered with your Spring Boot application.

```
import org.springframework.boot.SpringApplication;  
import org.springframework.boot.autoconfigure.SpringBootApplication;  
import org.springframework.boot.web.servlet.FilterRegistrationBean;  
import org.springframework.context.annotation.Bean;  
  
@SpringBootApplication  
public class RateLimitingApplication {  
  
    public static void main(String[] args) {  
        SpringApplication.run(RateLimitingApplication.class, args);  
    }  
  
    @Bean  
    public FilterRegistrationBean<OncePerRequestFilter>  
loggingFilter(OncePerRequestFilter rateLimitingFilter) {  
        FilterRegistrationBean<OncePerRequestFilter> registrationBean = new  
FilterRegistrationBean<>();  
        registrationBean.setFilter(rateLimitingFilter);  
        registrationBean.addUrlPatterns("/");  
        return registrationBean;  
    }  
}
```

## Conclusion

Rate limiting is a crucial aspect of API management that helps maintain service quality, ensures fair use of resources, and protects against misuse. By understanding and implementing appropriate rate limiting strategies, API providers can enhance the reliability and security of their services.