zendesk developers

API Basics > Best practices

# Best practices for avoiding rate limiting

**ON THIS PAGE**

If you make a lot of Zendesk API requests in a short amount of time, you may bump into the API rate limit. When you reach the limit, the API stops processing any more requests until a certain amount of time has passed.

The rate limits for the Zendesk APIs are outlined in Rate limits in the API reference. Get familiar with the limits before starting your project.

This article covers the following best practices for avoiding rate limiting.

## Monitoring API activity against your rate limit

You can use the API usage dashboard in Admin Center to monitor your API activity against your rate limit. See Managing API usage in your Zendesk account in Zendesk help.

## API usage

Monitor API usage to avoid failures and provide more efficient customer support.
Learn about API usage

⚠ You're over your API usage limit. Contact Zendesk Sales to increase your limit.

### 7-day summary

| ⊠ 429 errors | ⚠ Limit near-breaches | 👁 Limit approaches |
|---|---|---|
| Too many requests in a given timeframe | Times that 90–99% of limit was reached | Times that 80–89% of limit was reached |
| **27.2% of requests** ⓘ | **6** ⓘ | **2** ⓘ |

### API requests breakdown

| Last 30 days ⌄ | Daily ⌄ |
|---|---|

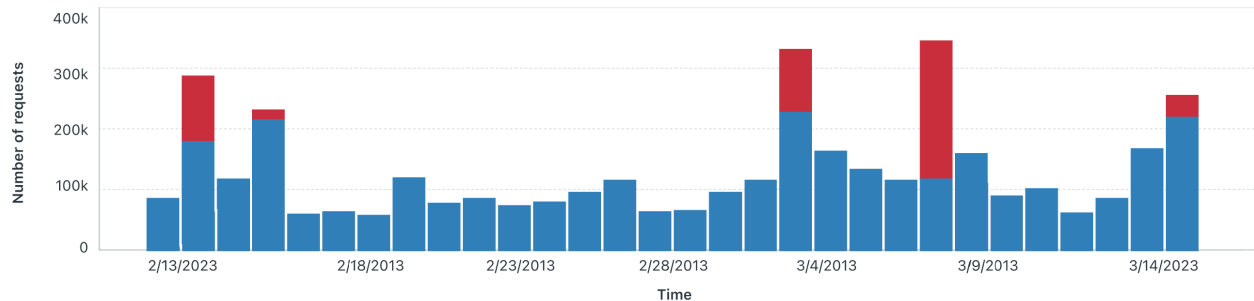| Total requests | 429 errors | Limit near-breaches | Limit approaches |
|---|---|---|---|
| **2,035k** | **521k (25.6%)** | **231** | **231** |

**Select data to show**

| Number of requests ⌄ |
|---|

● Successful requests   ● 429 errors



The following response headers contain the account's rate limit and the number of requests remaining for the current minute:

```
1    X-Rate-Limit: 700
2    X-Rate-Limit-Remaining: 699
```

The responses from Ticketing API list endpoints such as List Tickets or Search Users provide the following headers:

```
1    x-rate-limit: 700
2    ratelimit-limit: 700
3    x-rate-limit-remaining: 699
4    ratelimit-remaining: 699
5    ratelimit-reset: 41
6    zendesk-ratelimit-tickets-index: total=100; remaining=99; resets=41
```

The `zendesk-ratelimit-tickets-index` header is only available in list ticket endpoints. For more information, see List Tickets limits.

# Using rate limit headers in your application

For each Ticketing API call you make, Zendesk includes the above account-wide rate limit information in the response headers. To utilize this information, ensure that your application properly reads and interprets these headers. Always check the account wide limits header first. If you receive a 429 response, then look for endpoint-specific headers.

In the examples below, the headers are processed to check for an endpoint rate limit. **Note**: These examples are for reference only.

## Python

```python
1   import requests
2   import time
3
4   def call_zendesk_api():
5       url = "https://subdomain.zendesk.com/api/v2/tickets"
6       headers = {"Authorization": "Bearer YOUR_ACCESS_TOKEN"}
7       response = requests.get(url, headers=headers)
8       should_continue = handle_rate_limits(response)
9       if should_continue:
10          # Process the API response
11          # Your code here...
12
13  def handle_rate_limits(response):
14      account_limit = response.headers.get("ratelimit-remaining")
15      endpoint_limit = response.headers.get("Zendesk-RateLimit-Endpoint")
16      accountLimitResetSeconds = response.headers.get("ratelimit-reset")
17
18      if account_limit:
19          account_remaining = int(account_limit)
20          if account_remaining > 0:
21              if endpoint_limit:
22                  endpoint_remaining = int(endpoint_limit.split(";")[1].split("=")
    [1])
23                  if endpoint_remaining > 0:
24                      return True
25                  else:
26                      endpointLimitResetSeconds = int(endpoint_limit.split(";")[2]
    .split("=")[1])
27                      # Endpoint-specific limit exceeded; stop making more calls
28                      handle_limit_exceeded(endpointLimitResetSeconds)
29              else:
```

```python
30                    # No endpoint-specific limit
31                    return True
32            else:
33                # Account-wide limit exceeded
34                handle_limit_exceeded(accountLimitResetSeconds)
35        return False
36
37    def handle_limit_exceeded(limitHeaderResetTime):
38        reset_time = 60
    # default time if reset time is not available from the header
39        if limitHeaderResetTime:
40            reset_time = limitHeaderResetTime
41        wait_time = reset_time - time.time() + 1  # Add 1 second buffer
42        print(f"Rate limit exceeded for {limit_header}. Waiting for {wait_time}
     seconds...")
43        time.sleep(wait_time)
```

## Javascript

```javascript
1    const axios = require('axios');
2
3    async function callZendeskAPI() {
4        const url = "https://subdomain.zendesk.com/api/v2/tickets";
5        const headers = {"Authorization": "Bearer YOUR_ACCESS_TOKEN"};
6        try {
7            const response = await axios.get(url, { headers });
8            const shouldContinue = handleRateLimits(response);
9            if (shouldContinue) {
10               // Process the API response
11               // Your code here...
12           }
13       } catch (error) {
14           // Handle other errors
15           console.error(error);
16       }
17   }
18
19   function handleRateLimits(response) {
20       const accountLimit = response.headers["ratelimit-remaining"];
21       const endpointLimit = response.headers["Zendesk-RateLimit-endpoint"];
22       const accountLimitResetSeconds = response.headers["ratelimit-reset"]
23
24       if (accountLimit) {
25           const accountRemaining = parseInt(accountLimit);
26           if (accountRemaining > 0) {
27               if (endpointLimit) {
```

```
28                    const endpointRemaining = parseInt(endpointLimit.split(";")[1].
    split("=")[1]);
29                    if (endpointRemaining > 0) {
30                        return true;
31                    } else {
32                        const endpointLimitResetSeconds = parseInt(endpointLimit.
    split(";")[2].split("=")[1]);
33                        // Endpoint-specific limit exceeded
34                        handleLimitExceeded(endpointLimitResetSeconds);
35                    }
36                } else {
37                    // No endpoint-specific limit
38                    return true;
39                }
40            } else {
41                // Account-wide limit exceeded
42                handleLimitExceeded(accountLimitResetSeconds);
43            }
44        }
45        return false;
46    }
47
48    async function handleLimitExceeded(limitHeaderResetTime) {
49        const resetTime = limitHeaderResetTime || 60; // default to 60
50        const waitTime = resetTime - Math.floor(Date.now() / 1000) + 1;
    // Add 1 second buffer
51        console.log(`Rate limit exceeded for {limitHeader}. Waiting for ${waitTime}
     seconds...`);
52        await new Promise(resolve => setTimeout(resolve, waitTime * 1000));
53    }
54
55    callZendeskAPI();
```

Both examples verify that the limit is not exceeded for each response before proceeding to the next request. If the limit is exceeded, parse the reset time from the rate limit headers, calculate the wait time until the reset, and then pause until that time. After waiting, retry the API call.

Alternatively, you can choose to check only when the response status code is 429 (rate limit exceeded).

Here are some best practices for utilizing rate limit headers:

- Regularly monitor your API usage to prevent unexpected rate limit breaches.
- Design your application to gracefully handle rate limit headers, ensuring continuous service for your users.
- Implement exponential backoff strategies to effectively manage rate limit exceeded errors.

# Handling errors caused by rate limiting

If the rate limit is exceeded, the API responds with a 429 Too Many Requests status code.

It's best practice to include error handling for 429 responses in your code. If your code ignores 429 errors and keeps trying to make requests, you might start getting null errors. At that point, the error information won't be useful in diagnosing the problem.

For example, a request that bumps into the rate limit might return the following response:

```
1   < HTTP/1.1 429
2   < Server: nginx/1.4.2
3   < Date: Mon, 04 Nov 2013 00:18:27 GMT
4   < Content-Type: text/html; charset=utf-8
5   < Content-Length: 85
6   < Connection: keep-alive
7   < Status: 429
8   < Cache-Control: no-cache
9   < X-Zendesk-API-Version: v2
10  < Retry-After: 93
11  < X-Zendesk-Origin-Server: ****.****.***.*****.com
12  < X-Zendesk-User-Id: 338231444
13  < X-Zendesk-Request-Id: c773675d81590abad33i
14  <
15  * Connection #0 to host SUBDOMAIN.zendesk.com left intact
16  * Closing connection #0
17  * SSLv3, TLS alert, Client hello (1):
18  Rate limit for ticket updates exceeded, please wait before updating this
      ticket again
```

The response contains the following information:

- Status: 429
- Retry-After: 93

The 429 status code means too many requests. The Retry-After header specifies that you can retry the API call in 93 seconds. Your code should stop making additional API requests until enough time has passed to retry.

The following pseudo-code shows a simple way to handle rate-limit errors:

```
1   response = request.get(url)
2   if response.status equals 429:
3       alert('Rate limited. Waiting to retry…')
4       wait(response.headers['retry-after'])
5       retry(url)
```

## Node.js

The following snippet shows how you can handle rate-limit errors in JavaScript for Node.js.

```
 1  const axios = require("axios")
 2
 3  async function requestWithRateLimit(url, username, apiToken) {
 4    const tokenUsername = username + '/token';
 5    const password = apiToken;
 6    const response = await axios.get(url, {
 7        auth: {
 8          tokenUsername,
 9          password
10      }
11    })
12    if (response.status === 429) {
13      const secondsToWait = Number(response.headers["retry-after"])
14      await new Promise(resolve => setTimeout(resolve, secondsToWait * 1000))
15      return requestWithRateLimit(url, username, apiToken)
16    }
17    return response
18  }
```

## Browser JavaScript

The following snippet shows how you can handle rate-limit errors in client-side JavaScript for the browser.

```
 1  async function requestWithRateLimit(url, accessToken) {
 2    const options = {
 3      method: "GET",
 4      headers: {
 5        Authorization: `Bearer ${accessToken}`
 6      }
 7    }
 8    const response = await fetch(url, options)
 9    if (response.status === 429) {
10      const secondsToWait = Number(response.headers.get("retry-after"))
11      await new Promise(resolve => setTimeout(resolve, secondsToWait * 1000))
12      return requestWithRateLimit(url, accessToken)
13    }
14    return response
15  }
```

**Note:** To make an authenticated request from the browser to a Zendesk API, you must authenticate the request using an OAuth access token. For more information, see Making client-side CORS requests to the Ticketing API.

## Python

The following snippet shows how you can handle rate-limit errors in Python.

```python
import requests
import time
import os

# Store the API token in an environment variable for security reasons
ZENDESK_API_TOKEN = os.getenv('ZENDESK_API_TOKEN')

def request_with_rate_limit(url, username, token):
    auth = (f'{username}/token', token)
    response = requests.get(url, auth=auth)
    if response.status_code == 429:  # Check if the rate limit has been reached
        seconds_to_wait = int(response.headers["Retry-After"])
# Get the number of seconds to wait from the Retry-After header
        time.sleep(seconds_to_wait)  # Sleep for that duration
        return request_with_rate_limit(url, username, token)
# Recursively call the function again until the rate limit is lifted
    return response

# Usage example:
# Assuming the API endpoint you're trying to access is located at
    `api/some_endpoint`

url = 'https://your_subdomain.zendesk.com/api/some_endpoint'
response = request_with_rate_limit(url, ZENDESK_USER_NAME, ZENDESK_API_TOKEN)
    # Pass in the URL, username, and API token
if response.status_code == 200:
    # Success
    print(response.json())
else:
    # Error handling
    print(f'Error occurred: {response.status_code} - {response.text}')
```

## Reducing the number of API requests

Make sure you make only the requests that you need. Here are areas to explore for reducing the number of requests:

1. Optimize your code to eliminate any unnecessary API calls.

   For example, are some requests getting data items that aren't used in your application? Are retrieved data items being put back to your Zendesk product instance with no changes made to them?

2. Cache frequently used data.

   You can cache data on the server or on the client using DOM storage. You can also save relatively static information in a database or serialize it in a file.

3. Use bulk and batch endpoints, such as Update Many Tickets, that let you update up to 100 tickets with a single API request.

# Regulating the request rate

If you regularly exceed the rate limit, update your code to distribute requests more evenly over a period of time. This is known as a throttling process or a throttling controller. Regulating the request rate can be done statically or dynamically. For example, you can monitor your request rate and regulate requests when the rate approaches the rate limit.

To determine if you need to implement a throttling process, monitor your request errors. How often do you get 429 errors? Does the frequency warrant implementing a throttling process?

# Frequently asked questions

- **Is there an endpoint that returns all endpoints for a resource, including the different rate limits and times to retry?**

  No, we don't provide a rate-limit endpoint.

  However, you can use the following response headers to monitor your account's rate limit and the number of requests remaining for the current minute:

  ```
  1    X-Rate-Limit: 700
  2    X-Rate-Limit-Remaining: 699
  ```

  Tickets use the following response headers:

  ```
  1    x-rate-limit: 700
  2    ratelimit-limit: 700
  3    x-rate-limit-remaining: 699
  4    ratelimit-remaining: 699
  5    ratelimit-reset: 41
  6    zendesk-ratelimit-tickets-index: total=100; remaining=99; resets=41
  ```

- **What happens when the rate limit is reached?**

  The request isn't processed and a response is sent containing a 429 response code and a `Retry-After` header. The header specifies the time in seconds that you must wait before you can try the request again.

- **Will batching calls reduce the number of API calls? How about making parallel calls?**

  No, batching calls won't reduce the number of API calls.

  However, using a bulk or batch endpoint, such as Update Many Tickets or Ticket Bulk Import, will reduce calls.

---

### Join our developer community

🗨 Forum    📄 Blog    🍀 Slack

**Zendesk**  181 Fremont Street, 17th Floor, San Francisco, California 94105

Privacy Policy ⎤ Terms & Conditions ⎤ System Status ⎤ Cookie Settings