



# DEVCOR

## 350-901

# Study Guide

**Constantin Mohorea**

CCIE® No. 16223, CCDE® No. 20170054, DevNet 500

**Cisco Press**

# Cisco DEVCOR 350-901 Study Guide

Edition: 2021.3

Constantin Mohorea

Cisco Press

- Find out API restrictions and limit the rate of your requests, thus preventing an HTTP 429 error.

The following partial code provides an example of a function that interacts with an API (for example, Webex) while handling network timeouts and rate-limiting errors:

```
def API_Read_Safe (url, timeout = 5, attempts = 3):
    """ Read API URL, handle HTTP 429 and timeout errors """

    tries = 0
    while (tries < attempts):
        tries += 1
        try:
            # read target URL
            response = requests.get(url, headers=headers, timeout = timeout)

            # if successful, return the response
            if response.ok:
                print ("Completed successfully")
                return response

            # If response is not "ok", handle error code 429:
            # - wait "Retry-After" seconds or 1 sec, if not present
            # - retry
            if response.status_code == 429:
                try:
                    retry_after = int(response.headers.get('Retry-After'))
                except Exception:
                    retry_after = 1

                print (f"Retry after {retry_after} seconds")
                sleep(retry_after)
                continue

            # if not OK and not 429, then it's something unrecoverable
            print(f"HTTP error: {response.status_code}")
            response.raise_for_status()

            ## if timeout error, retry. If any other error, leave it unhandled
            except requests.exceptions.ConnectTimeout:
                print ("Timeout, retry")
                continue

        # return "None" if unsuccessful after all attempts
        print ("Unable to get any response")
        return None
```

## 2.2 Implement control flow of consumer code for unrecoverable REST API errors

Network and REST API error handling in your application is one of the important development tasks. Good REST API error flow control should start with separation of errors into the following categories:

- **Unrecoverable errors:** These typically require user interaction for continuation.
- **Recoverable errors:** These require an application to wait and retry the process again without user interaction.

When an application receives a REST API response that falls into the unrecoverable error category, it should collect and display/store diagnostic information (response status code, headers, and body) for the troubleshooting process. No programmatic workaround is available.

However, if the response status code indicates a recoverable error, program logic should be implemented to act accordingly.

Here are some examples of unrecoverable errors:

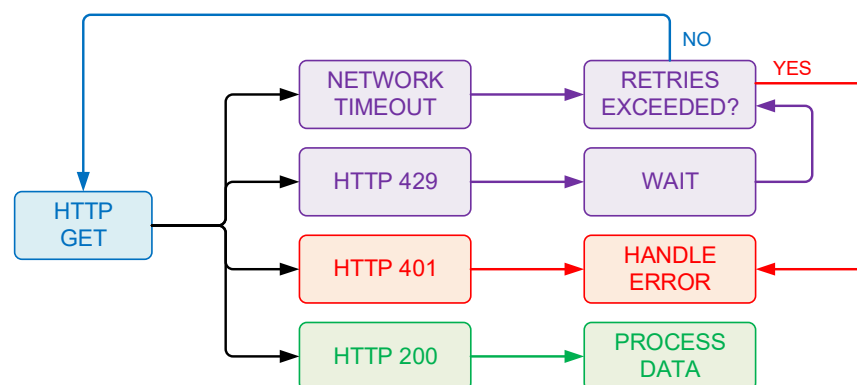
- **400 Bad Request:** You did not send a valid request, and the remote side is refusing to process further. You can get more information from HTTP headers or the body of the response.
- **401 Unauthorized:** Authentication is required (API key, API token, or username and password). If the API was previously accessible, be aware that keys and tokens can expire, and usernames and passwords can be changed.
- **403 Forbidden:** You are authenticated and authorized but do not have the permission to access the resource.

Here are some examples of recoverable errors:

- **405 Method Not Allowed:** Indicates an incorrect method for accessing a resource. The response must include the HTTP "Allow" header that lists valid methods.
- **408 Timeout:** Indicates a timeout sending a request to the server (for example, when POSTing large file). Your requests should be faster if you want the server to accept it.
- **429 Too Many Requests:** Indicates too many requests sent to the server. You should wait for some time before you issue a new one. The Cisco Webex REST API responds with the additional "Retry-After" HTTP header, which holds a value (in seconds) that the API consumer must wait before the rate-limit counter resets.

Figure 25 shows an example of the control flow in a sample application.

*Figure 25: Control Flow for REST API Error Handling*



## 2.3 Identify ways to optimize API usage through HTTP cache controls

When API load increases due to growth of user requests and application data, hardware requirements increase as well. There are a few ways to address this: you can scale hardware (horizontally or vertically) or optimize your application, or both. The goals of the application optimizations are as follows:

- Improve response time in user-interactive applications.
- Decrease the required bandwidth.
- Decrease processing times.
- Lower resource utilization.

These goals may be achieved with pagination (see Section 2.4), caching, and compression.

### HTTP Caching

HTTP caching is a technique that stores copies of web resources and serves them back when requested. Typically limited to caching responses to "GET" requests, HTTP caches intercept these requests and return local copies instead of redownloading resources from originating servers.

Caching may be implemented at different places: at the client (for example, in a web browser), locally as an organization's proxy server, on the server side as a standalone reverse proxy, or as a part of an API gateway.

Caching reduces the amount of network traffic and the load of web servers. It also improves the user experience since HTTP requests may be completed in less time. However, resources may change, and it is important to cache them properly: only until they change, not longer. To address this, caches use the concept of "freshness."

A stored resource is considered "fresh" if it may be served directly from the cache. As web servers cannot contact caches and clients when a resource changes, they simply indicate an expiration time for the provided content. Before this expiration time, the resource is "fresh," and after the expiration time, the resource is "stale."

"Stale" content is not ignored and might still be reused, but it needs to be validated first. During validation, a check is performed on the original resource to verify that the cached copy is still current. Such checks usually involve only HTTP headers (HTTP "HEAD" request) and no data, so they are "cheap" compared to retrieving the whole resource. Validation checks increase the response time; however, they are more reliable than a time-based mechanism because they ensure data is up to date.

Website content may change when the cache is still "fresh," so some servers may want to force validation on every request to avoid inconsistencies. One way to achieve that is to set the explicit expiration time in the past and thus indicate that the response is already stale when it's received.

Validation can only happen if the origin server previously provided a *validator* for a resource, which is some value that describes a specific version of a resource. There are two types of validators:

- The date of last modification of the document, provided in the "Last-Modified" HTTP header.
- Some string that uniquely identifies a version of a resource, called the Entity Tag (ETag) and provided in the "ETag" HTTP header. It is not specified how ETag values are generated, so it may be a hash of the content, a hash of the last modification timestamp, or just a revision number.

Validators may be *strong* when they guarantee that two resources are completely byte-to-byte identical to each other. With the *weak* validators, two versions of the document are considered as identical if the content is equivalent (for example, same web page but with the different ads on it). ETags may be used as strong validators (by default) or as weak ones (prefixed with "W/"; for example, W/"5bfb1b822b0f7").

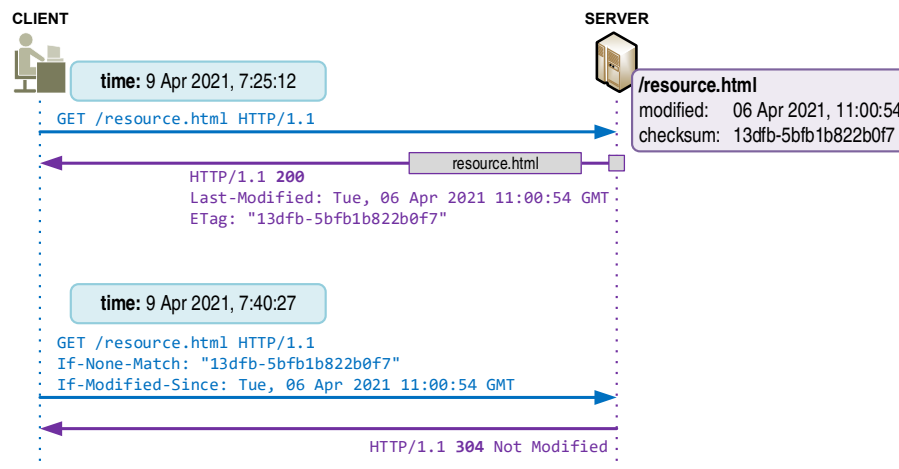
## HTTP Conditional Requests

With HTTP *conditional requests*, servers' responses are based on the result of a comparison between the current version of a resource with validators included in a request. Conditional requests enable an efficient validation process but have other uses as well: to verify the integrity of a document (for example, when resuming a download) or, when clients act in parallel, to prevent accidental deletion or overwriting of one version of a document with another.

Special headers are used to specify preconditions in requests, and their logic depends on the HTTP request type: safe (read) or unsafe (write), as presented in Figure 26. In these examples, the "Last-Modified" and "ETag" headers are already known from previous requests.

With safe methods (GET and HEAD), conditions are used to fetch a resource *only when needed*.

Figure 26: HTTP Conditional Requests: Fetch the Resource



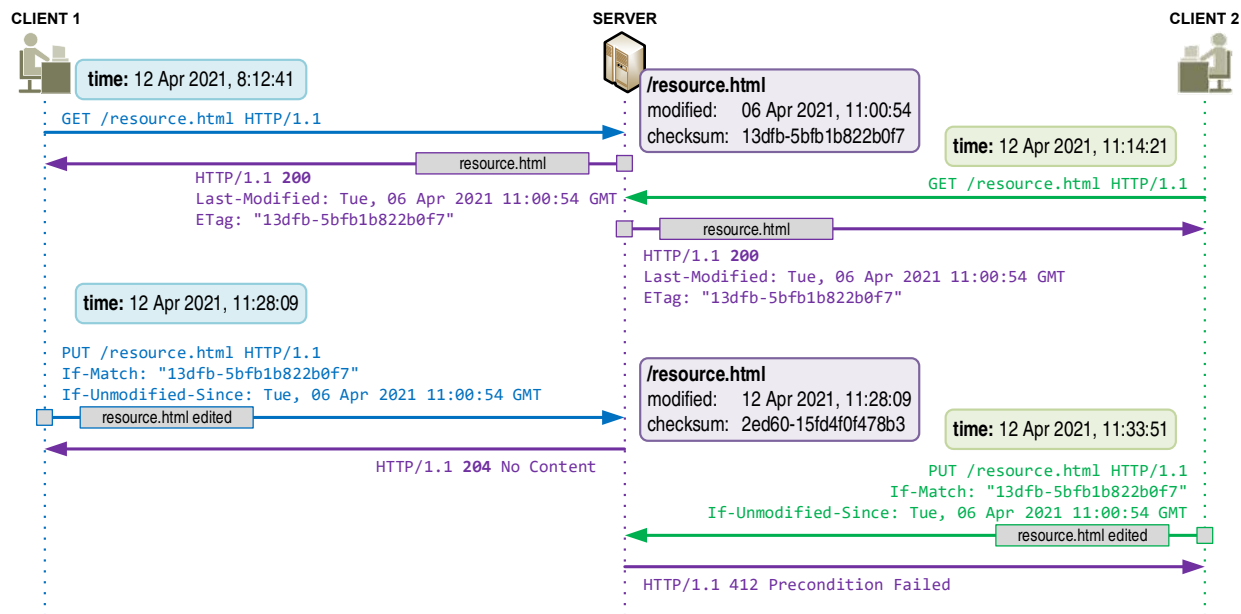
The following headers specify a condition for the request:

- **"If-None-Match"**: Set to previously known "ETag" value. The condition is true if the ETag of the resource on a server is different from the one given in this header. Both strong and weak ETags may be used.
- **"If-Modified-Since"**: Set to previously known "Last-Modified" value. The condition is true if the date of the resource on a server is more recent than the one listed in this header.

If the condition is true (that is, the resource has changed), then the server responds with the HTTP 200 code and the complete resource (for the GET method). If the resource hasn't changed, the server must respond with the HTTP status code 304 "Not Modified" and an empty body.

With unsafe methods (PUT and DELETE), conditions are used to allow modifications to a resource *only when it's still the same and hasn't changed since the last interaction* (see Figure 27).

**Figure 27: HTTP Conditional Requests: Modify the Resource**



The following headers specify a condition for the request:

- **"If-Match"**: Set to previously known "ETag" value. The condition is true if the ETag of the resource on a server is equal to the one supplied in this header. Only a strong ETag may be used, as weak tags never match under this comparison.
- **"If-Unmodified-Since"**: Set to previously known "Last-Modified" value. The condition is true if the date of the resource on a server is the same or older than the one listed in this header.

If the condition is true (that is, the resource has not changed), then the server proceeds with the update as normal. If the condition is not met, then the change is rejected with a 412 "Precondition Failed" error.

## Caching Controls

Sometimes, caching may interfere with the application behavior, and the server administrator may wish to provide explicit directions to cache engines. The *Cache-Control* HTTP header is used for this purpose. Some of the most often used options are:

- **no-store**: The response may not be stored in any cache; disable caching.
- **no-cache**: Allow caching, but stored response *must* be validated first before using it.
- **max-age**: Time in seconds a cached copy is valid. `max-age=0` is similar to `no-cache`.
- **public**: Can be cached by anyone.

- **private:** Only browser/client may cache.

Example: `Cache-Control: no-cache, max-age=0`

## HTTP Data Compression

By using HTTP compression, less data needs to be transferred over the network, which results in faster response times and lower bandwidth utilization.

When making an HTTP request, a client may include an HTTP "Accept-Encoding" header to indicate a list of compression algorithms that the client supports. Common examples are compress, deflate, gzip, and bzip2.

Compression is not always accepted by the server, as it may cause issues with antivirus software, next-gen firewalls, or proxies, and/or it may violate company policies.

When the server receives a request that indicates compression support, it can honor that request or ignore it and send an uncompressed response. If one of the client's compression methods is supported, and the server decides to use it to compress the response, the server must include it in the "Content-Encoding" HTTP header.

## 2.4 Construct an application that consumes a REST API that supports pagination

When a client makes a REST API request, the response can be very large (for example, event log). Quite often, the client is not interested in all the data. For example, it might only need log entries for yesterday or the 50 latest chat messages. Pagination is a method to split resulting data into manageable chunks to allow better handling of large data sets and faster request processing.

The main reasons for using pagination in REST APIs are as follows:

- **To improve response times and end-user experience:** Because paginated responses are much smaller, they can be handled by the server and the network faster, and faster responses provide better user experiences.
- **To save resources:** Providing small responses demands much less compute and network resources.

One of the simplest types of pagination is the "Offset/Page"-based pagination. With this method, the client supplies the data offset and data size as query parameters. Here are two examples:

- **`https://example.com/events?offset=100&limit=20`:** Returns 20 entries, starting with 100th
- **`https://example.com/events?page=6&per_page=20`:** Same logic, different syntax

Typically, if offset or limit parameters are omitted, the server will use defaults (for example, offset=0 and limit=50).

Offset-based pagination is simple to implement and to use; however, it has its downsides:

- Adding or deleting entries (known as *page drift*) between calls may cause confusion. Some results may be skipped over (for example, some item is deleted and now the first item on a new page moved to the previous) and some may be included more than once (when items are inserted).