# Explain Caching in API GATEWAY

Caching in an API Gateway refers to the process of storing responses from API calls and serving them directly from the cache instead of forwarding the requests to the backend servers. This caching mechanism helps improve the performance and scalability of an API by reducing the response time and the load on the backend infrastructure.

Caching in API Gateway can be configured and customized based on various parameters. Some common caching options include:

1. Cache Key: It is a unique identifier for each API request that determines the cache entry. The cache key is typically generated based on specific request parameters such as URL, headers, query parameters, or a combination of these.

2. Cache TTL (Time-to-Live): It represents the duration for which a cached response remains valid before it expires. After the TTL expires, subsequent requests for the same resource will be forwarded to the backend servers to fetch a fresh response.

3. Cache Invalidation: APIs often serve dynamic data that can change over time. Cache invalidation mechanisms allow you to define rules or triggers to invalidate and refresh the cache when the underlying data is modified. This ensures that clients always receive the latest data when necessary.

4. Cache-Control Headers: The API Gateway can honor the Cache-Control headers sent by the backend API or modify them before caching the response. These headers provide instructions on how the response should be cached, such as setting the maximum cache age, allowing caching only for authenticated users, or disabling caching altogether.

5. Cache Size and Storage: The API Gateway allocates a certain amount of memory or disk space to store the cached responses. The cache size can be configured based on the expected workload and available resources.