

Example of available profile types based on the filtered profile data.

Each profile type represents a type of resource consumption, such as CPU, wall time, or memory. The profile types available for you to select may differ based on the language being profiled. The following are the most common profile types:

- **CPU** profiling measures which methods consume the most CPU on an application.
- **Allocation** profiling measures the amount (both in terms of count and size in Bytes) of memory allocated by a given method. Note: This isn't retained memory, which means that the allocated memory that is measured may or may not be garbage collected.
- **Heap** profiling measures the amount of heap memory allocated by each function that hasn't been garbage collected (yet). This is useful for investigating the overall memory usage of your service and identifying potential memory leaks.
- **Lock** profiling measures the amount of time a thread is waiting to acquire a lock and is hence doing nothing.
- **Wall Time** profiling measures the effective time spent by methods (regardless of whether those methods were running on CPU, waiting for network I/O, blocked by another thread, or just idle). It can be useful to debug latency at first glance and then dig into the other profiling types to find out what was causing the latency. The wall time profile can be considered to be the most similar to the associated APM flame graph.
- **File I/O** and **Socket I/O** measure the number of time spent by methods on disk (for example, reading a file from disk) and network I/O operations (for example, waiting for an API call to return).
- **Exceptions** measures the amount of exceptions thrown. The profiler doesn't catch/handle exceptions, but it tracks their creation. (For Java, exceptions and errors aren't synonymous, and errors are more likely to be unrecoverable, such as in the case of an `OutOfMemoryError`).

[Endpoint profiling](#) allows you to scope profiler flame graphs by any API endpoint of a service, so you can find endpoints that are slow and causing poor end-user experience.

Debugging and understanding why an endpoint has high latency can be tricky. For example, high latency could be caused by a method that is CPU heavy and unknowingly on the critical path of a request process where latency is important.

You can do the following with endpoint profiling:

- Identify the bottleneck methods that are slowing down the endpoint's overall response time.
- Isolate the top endpoints that are responsible for consuming resources like CPU and memory. This is particularly helpful when you're trying to optimize your service for performance gains.
- Understand if third-party code or runtime libraries are the reason for endpoints being slow or heavy on resource consumption.

<input checked="" type="checkbox"/> By	Minute by Endpoint ▾	
<input checked="" type="checkbox"/>	GET /credits	30.14s 
<input checked="" type="checkbox"/>	GET /old-movies	7.2s 
<input checked="" type="checkbox"/>	GET /movies	4.79s 
<input checked="" type="checkbox"/>	GET /stats	3.53s 

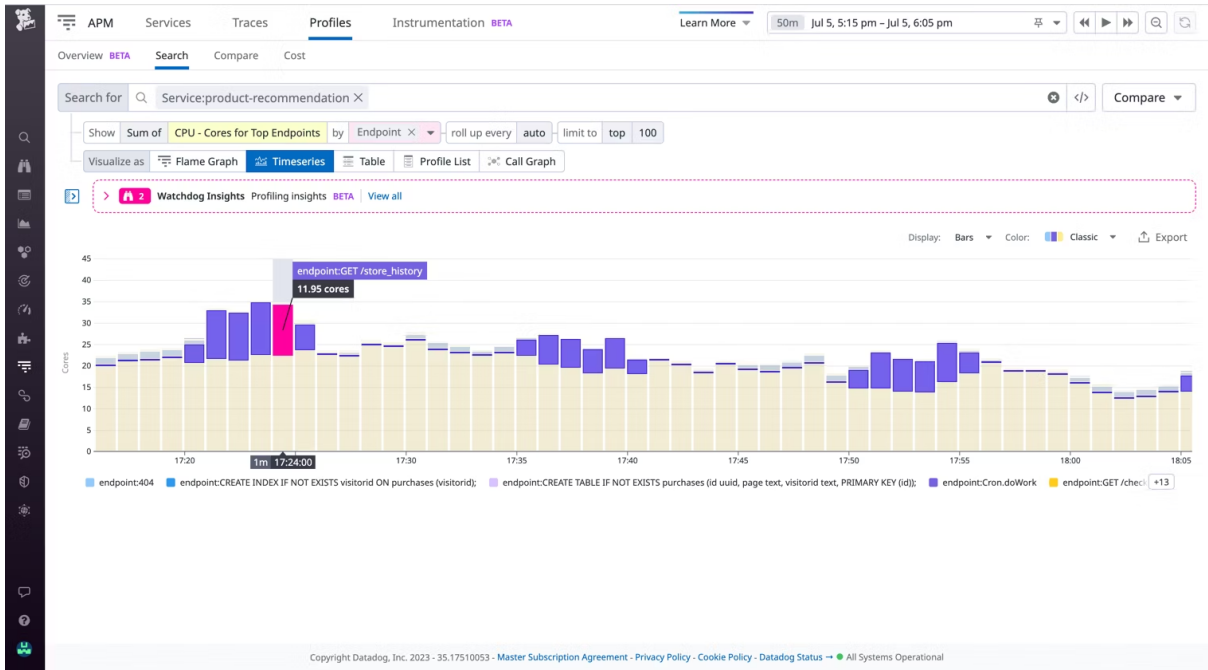
Example of list of endpoints that appears in the summary list for a profile.

In general, it's valuable to track which endpoints are consuming the most valuable resources, such as CPU and memory. The list can help you identify if your endpoints have regressed or if you have newly introduced endpoints that are drastically consuming more resources than expected and slowing down your overall service.

Track the endpoints that consume the most resources

It is valuable to track top endpoints that are consuming valuable resources such as CPU and wall time. The list can help you identify if your endpoints have regressed or if you have newly introduced endpoints that are consuming drastically more resources, slowing down your overall service.

The following image shows that `GET /store_history` is periodically impacting this service by consuming 20% of its CPU:



Track average resource consumption per request

Select Per endpoint call to see behavior changes even as traffic shifts over time. This is useful for progressive rollout sanity checks or analyzing daily traffic patterns.

The following video shows that CPU per request doubled for /GET train.