

Python

# FastAPI Latency & Response Times

Understanding Performance Metrics and Optimizing API Speed in Modern Web Applications

April 15, 2025



Source Code

## 1. Introduction to API Performance

In modern web development, speed and efficiency are crucial. When developing APIs or web services with FastAPI, understanding fundamental concepts like latency and response time is essential to ensure our applications are scalable and provide an optimal user experience.

### 1.1. Latency vs Response Time: Understanding the Difference

Although often used interchangeably, **latency** and **response time** are distinct concepts:

- **Response Time:** What the client perceives—includes request processing time (service time), network delays, and queue wait times.
- **Latency:** The duration a request spends waiting to be processed—the time it remains latent before processing begins.

### 1.2. Measuring Performance: Beyond Averages

When evaluating API performance, average response time is commonly reported. However, averages are not ideal metrics for understanding "typical" system behavior because they don't reflect real user experience.

Using **percentiles** provides a more accurate picture:

- **50th Percentile (p50) or Median:** Half of requests complete in less time,



**Alejandro Sánchez Yalí**

Software Developer | AI & Blockchain Enthusiast

 [www.asanchezyali.com](http://www.asanchezyali.com)

half take longer

- **High Percentiles** like p95, p99, and p999 (95th, 99th, and 99.9th percentiles):  
Reveal how bad your outliers are and how they affect a small but significant percentage of users

High percentiles, also known as "tail latencies," are crucial because they directly affect user experience. For example, Amazon describes response time requirements for internal services in terms of the 99.9th percentile, even though it affects only 1 in 1,000 requests. Why? Because customers with the slowest requests often have the most data in their accounts—they're the most valuable customers.

## 2. Performance Monitoring with FastAPI

Let's build a small service with FastAPI and simulate load to analyze its behavior. Our goal will be to measure and visualize different response time percentiles.


### 2.1. Creating a FastAPI Application

First, we'll create a simple API with different endpoints that simulate various processing times:



**Alejandro Sánchez Yalí**


Software Developer | AI & Blockchain Enthusiast

 [www.asanchezyali.com](http://www.asanchezyali.com)

```
1 # api_server.py
2 import asyncio
3 import random
4 import time
5 from fastapi import FastAPI
6 from fastapi.middleware.cors import CORSMiddleware
7
8 app = FastAPI(title="Latency Demo API")
9
10 # Configure CORS
11 app.add_middleware(
12     CORSMiddleware,
13     allow_origins=["*"],
14     allow_credentials=True,
15     allow_methods=["*"],
16     allow_headers=["*"],
17 )
18
```

**Alejandro Sánchez Yalí**

Software Developer | AI &amp; Blockchain Enthusiast


 [www.asanchezyali.com](http://www.asanchezyali.com)

```
19 @app.get("/")
20 async def root():
21     """Fast endpoint with consistent response time"""
22     return {"message": "Latency demonstration API"}
23
24 @app.get("/fast")
25 async def fast_endpoint():
26     """Endpoint with fast response (10-30ms)"""
27     await asyncio.sleep(random.uniform(0.01, 0.03))
28     return {"response_type": "fast"}
29
30 @app.get("/medium")
31 async def medium_endpoint():
32     """Endpoint with medium response time (50-150ms)"""
33     await asyncio.sleep(random.uniform(0.05, 0.15))
34     return {"response_type": "medium"}
35
36 @app.get("/slow")
37 async def slow_endpoint():
38     """Endpoint with slow response time (200-500ms)"""
```



## Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast

 [www.asanchezyali.com](http://www.asanchezyali.com)

```
39     await asyncio.sleep(random.uniform(0.2, 0.5))
40     return {"response_type": "slow"}
41
42 @app.get("/variable")
43 async def variable_endpoint():
44     """
45     Endpoint with variable response time:
46     - 80% of requests: fast (10-50ms)
47     - 15% of requests: medium (100-300ms)
48     - 5% of requests: very slow (500-1500ms)
49     """
50     random_value = random.random()
51
52     if random_value < 0.8:
53         await asyncio.sleep(random.uniform(0.01, 0.05))
54         category = "fast (80%)"
55     elif random_value < 0.95:
56         await asyncio.sleep(random.uniform(0.1, 0.3))
57         category = "medium (15%)"
58     else:
```



## Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast

 [www.asanchezyali.com](http://www.asanchezyali.com)

```
59         await asyncio.sleep(random.uniform(0.5, 1.5))
60         category = "very slow (5%)"
61
62         return {"response_type": "variable", "category": category}
63
64 if __name__ == "__main__":
65     import uvicorn
66     uvicorn.run("api_server:app", host="0.0.0.0", port=8000,
67                reload=True)
```

## 2.2. Creating a Load Testing Agent

Now, let's create an agent that performs multiple requests to our API and collects performance data:

```
1 # load_tester.py
2 import asyncio
3 import time
4 import statistics
```



**Alejandro Sánchez Yalí**


Software Developer | AI & Blockchain Enthusiast

 [www.asanchezyali.com](http://www.asanchezyali.com)

```
5 import matplotlib.pyplot as plt
6 import numpy as np
7 import aiohttp
8 from collections import defaultdict
9
10 class LoadTester:
11     def __init__(self, base_url="http://localhost:8000"):
12         self.base_url = base_url
13         self.endpoints = {
14             "fast": "/fast",
15             "medium": "/medium",
16             "slow": "/slow",
17             "variable": "/variable"
18         }
19         self.results = defaultdict(list)
20
21     async def make_request(self, session, endpoint):
22         """Makes an HTTP request and measures response time"""
23         url = f"{self.base_url}{self.endpoints[endpoint]}"
24         start_time = time.time()
```

**Alejandro Sánchez Yalí**

Software Developer | AI &amp; Blockchain Enthusiast

 [www.asanchezyali.com](http://www.asanchezyali.com)



```
25
26     try:
27         async with session.get(url) as response:
28             await response.json()
29             response_time = (time.time() - start_time) * 1000 #
30             in milliseconds
31             self.results[endpoint].append(response_time)
32             return response_time
33     except Exception as e:
34         print(f"Error in request to {url}: {e}")
35         return None
36
37     async def generate_load(self, endpoint, num_requests,
38                             concurrency):
39         """Generates load for a specific endpoint"""
40         async with aiohttp.ClientSession() as session:
41             tasks = []
42             for _ in range(num_requests):
43                 tasks.append(self.make_request(session, endpoint))
44             if len(tasks) >= concurrency:
```



## Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast


 [www.asanchezyali.com](http://www.asanchezyali.com)

```
43         await asyncio.gather(*tasks)
44         tasks = []
45
46         if tasks:
47             await asyncio.gather(*tasks)
48
49     def calculate_percentiles(self, endpoint):
50         """Calculates percentiles for response times"""
51         if not self.results[endpoint]:
52             return {}
53
54         data = sorted(self.results[endpoint])
55         return {
56             "min": min(data),
57             "p50": statistics.median(data),
58             "p90": np.percentile(data, 90),
59             "p95": np.percentile(data, 95),
60             "p99": np.percentile(data, 99),
61             "p999": np.percentile(data, 99.9) if len(data) >= 1000
        else None,
```



## Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast


 [www.asanchezyali.com](http://www.asanchezyali.com)

```
62         "max": max(data),
63         "mean": statistics.mean(data),
64         "stdev": statistics.stdev(data) if len(data) > 1 else 0
65     }
66
67     def plot_results(self):
68         """Generates charts to visualize the results"""
69         fig, (ax1, ax2) = plt.subplots(2, 1, figsize=(12, 14))
70
71         # Percentile bar chart
72         endpoints = list(self.results.keys())
73         metrics = ["p50", "p90", "p95", "p99"]
74
75         x = np.arange(len(endpoints))
76         width = 0.2
77
78         for i, metric in enumerate(metrics):
79             values = [self.calculate_percentiles(ep)[metric] for ep
80 in endpoints]
81             ax1.bar(x + i*width, values, width, label=f'{metric}')
```



## Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast


 [www.asanchezyali.com](http://www.asanchezyali.com)

```
81
82     ax1.set_ylabel('Response Time (ms)')
83     ax1.set_title('Response Time Percentiles by Endpoint')
84     ax1.set_xticks(x + width * 1.5)
85     ax1.set_xticklabels(endpoints)
86     ax1.legend()
87     ax1.grid(axis='y', linestyle='--', alpha=0.7)
88
89     # Histogram for the variable endpoint
90     if "variable" in self.results and self.results["variable"]:
91         data = self.results["variable"]
92         # Use logarithmic scale for data
93         bins = np.logspace(np.log10(min(data)),
np.log10(max(data)), 50)
94         ax2.hist(data, bins=bins, alpha=0.7, color='green')
95         ax2.set_xscale('log')
96         ax2.set_title('Response Time Distribution (variable
endpoint)')
97         ax2.set_xlabel('Response Time (ms) - logarithmic scale')
98         ax2.set_ylabel('Number of Requests')
```



## Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast


 [www.asanchezyali.com](http://www.asanchezyali.com)

```
99
100     # Add vertical lines for percentiles
101     percentiles = self.calculate_percentiles("variable")
102     for metric, value in [(k, v) for k, v in
percentiles.items()
103                             if k in ["p50", "p90", "p95", "p99"]
and v is not None]:
104         ax2.axvline(x=value, color='red', linestyle='--',
alpha=0.6,
105                     label=f"{metric}: {value:.2f}ms")
106
107     ax2.legend()
108
109     plt.tight_layout()
110     plt.savefig('latency_results.png', dpi=300)
111     plt.close()
112
113     def print_summary(self):
114         """Prints a summary of the results"""
115         for endpoint in self.results:
```



**Alejandro Sánchez Yalí**

Software Developer | AI & Blockchain Enthusiast

 [www.asanchezyali.com](http://www.asanchezyali.com)

```
116         print(f"\n== Endpoint: {endpoint}  
      ({len(self.results[endpoint])} requests) ==")  
117         percentiles = self.calculate_percentiles(endpoint)  
118         for metric, value in percentiles.items():  
119             if value is not None:  
120                 print(f"{metric}: {value:.2f} ms")  
121  
122     async def main():  
123         tester = LoadTester()  
124  
125         # Generate load for each endpoint  
126         print("Starting load tests...")  
127  
128         for endpoint in tester.endpoints:  
129             requests = 1000 if endpoint == "variable" else 200  
130             print(f"Testing endpoint '{endpoint}' with {requests}  
requests...")  
131             await tester.generate_load(endpoint, requests, concurrency=50)  
132  
133         # Print results and generate charts
```



## Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast


 [www.asanchezyali.com](http://www.asanchezyali.com)

```
134     tester.print_summary()
135     tester.plot_results()
136     print("Tests completed. Results saved in 'latency_results.png'")
137
138 if __name__ == "__main__":
139     asyncio.run(main())
```

## 2.3. Running and Analyzing the Results


To run our experiment:

1. Start the FastAPI server:



```
1 python api_server.py
```

2. In another terminal, run the load testing agent:




```
1 python load_tester.py
```



**Alejandro Sánchez Yalí**

Software Developer | AI & Blockchain Enthusiast

 [www.asanchezyali.com](http://www.asanchezyali.com)

PYTHON

| CONCURRENCY & PARALLELISM

The agent will generate a file called `latency_results.png` with visualizations of the results, and will also print a summary to the console.



**Alejandro Sánchez Yalí**

Software Developer | AI & Blockchain Enthusiast



[www.asanchezyali.com](http://www.asanchezyali.com)




## 2.4. Visualization of Results

latency\_results.png



**Alejandro Sánchez Yalí**

Software Developer | AI & Blockchain Enthusiast

 [www.asanchezyali.com](http://www.asanchezyali.com)

### 3. Analyzing the Results

The results clearly show why percentiles are more informative than averages:

- For the **variable** endpoint:
  - The median (p50) is quite low, showing that most requests are fast
  - p95 and p99 reveal the slower requests affecting a small percentage of users
  - The average could be distorted by extreme outliers
- Comparing different endpoints:
  - We observe how response time distribution varies considerably
  - We see the difference between endpoints with predictable vs. variable behavior

### 4. The Tail Latency Amplification Problem


Imagine that your frontend application needs to make several calls to these endpoints to compose a single page. If a page requires 5 API calls:

- With a 5% probability that each call is slow, the probability that at least one call is slow is:



**Alejandro Sánchez Yalí**

Software Developer | AI & Blockchain Enthusiast

 [www.asanchezyali.com](http://www.asanchezyali.com)

$$1 - (0.95)^5 \approx 23\%$$

This means that although only 5% of individual calls are slow, approximately 23% of page loads will experience delays. This phenomenon is known as "tail latency amplification" or "tail at scale."

## 5. Conclusion

When developing and monitoring APIs, it's critical to:


- **Measure Percentiles, Not Just Averages:** High percentiles reveal problems affecting real users that might go unnoticed in means.
- **Understand the Complete Distribution:** Visualize response time histograms to better understand system behavior.
- **Consider Tail Latency Amplification:** If your frontend or client makes multiple API calls to complete a user action, optimizing high percentiles is essential.
- **Establish SLOs and SLAs Based on Percentiles:** For example, "p99 must be less than 300ms" is a more meaningful metric than "average time must be less than 100ms."

Modern applications must be designed and optimized with these principles in mind to provide consistent, high-quality experiences to all users, not just the average user.



**Alejandro Sánchez Yalí**

Software Developer | AI & Blockchain Enthusiast

 [www.asanchezyali.com](http://www.asanchezyali.com)

In a future post, we'll explore techniques for optimizing these high percentiles, including implementing priority queues, circuit breakers, and caching strategies to mitigate tail latency amplification.

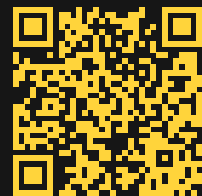
## 6. Explore My Other Posts

### Enjoyed This Content?

Don't miss my previous post about:


#### **Python Generators: Elegant, Memory-Efficient Iterations**

Discover how Python Generators can help you process large datasets efficiently, create elegant data pipelines, and write cleaner code with minimal memory footprint.



**Alejandro Sánchez Yalí**

Software Developer | AI & Blockchain Enthusiast

 [www.asanchezyali.com](http://www.asanchezyali.com)

# Ready to **Optimize** Your APIs?

You've learned how to measure and analyze response times.

Try implementing percentile monitoring in your next project

and watch your applications deliver a better experience for all users!