

Python

FastAPI Latency

Understanding and Optimizing Response Times

April 29, 2025



Source Code

1. Introduction to API Performance

In modern web development, speed and efficiency are crucial. Understanding fundamental concepts like latency and response time is essential to ensure APIs are scalable and provide an optimal user experience. This article uses FastAPI as a practical example to demonstrate these concepts, but the principles of latency, response time, and percentiles are applicable to all API systems, regardless of the technology used.

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:



Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast

 www.asanchezyali.com

- **50th Percentile (p50) or Median:** Half of requests complete in less time, half take longer.
- **High Percentiles** like p95, p99, and p999: Reveal how bad your outliers are and how they affect a small but significant percentage of users.

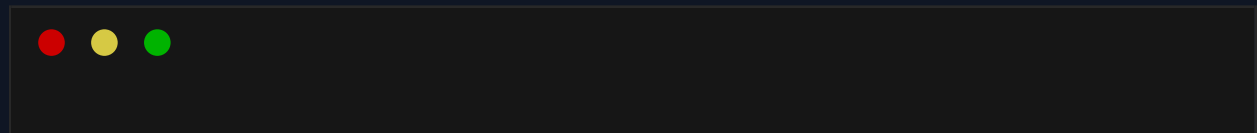
High percentiles, also known as "tail latencies," are crucial because they directly affect user experience. For example, Amazon describes response time requirements for internal services in terms of the 99.9th percentile, even though it affects only 1 in 1,000 requests, as these users often have the most valuable accounts.

2. Performance Monitoring with FastAPI

Let's build a small service with FastAPI and simulate load to analyze its behavior. Our goal is 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

```
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 app.add_middleware(
11     CORSMiddleware,
12     allow_origins=["*"],
13     allow_credentials=True,
14     allow_methods=["*"],
15     allow_headers=["*"],
16 )
17
18 @app.get("/")
19 async def root():
20     """Fast endpoint with consistent response time"""
```

**Alejandro Sánchez Yalí**

Software Developer | AI & Blockchain Enthusiast


 www.asanchezyali.com

```
21     return {"message": "Latency demonstration API"}
22
23 @app.get("/fast")
24 async def fast_endpoint():
25     """Endpoint with fast response (10-30ms)"""
26     await asyncio.sleep(random.uniform(0.01, 0.03))
27     return {"response_type": "fast"}
28
29 @app.get("/medium")
30 async def medium_endpoint():
31     """Endpoint with medium response time (50-150ms)"""
32     await asyncio.sleep(random.uniform(0.05, 0.15))
33     return {"response_type": "medium"}
34
35 @app.get("/slow")
36 async def slow_endpoint():
37     """Endpoint with slow response time (200-500ms)"""
38     await asyncio.sleep(random.uniform(0.2, 0.5))
39     return {"response_type": "slow"}
40
```



Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast

 www.asanchezyali.com

```
41 @app.get("/variable")
42 async def variable_endpoint():
43     """
44     Endpoint with variable response time:
45     - 80% of requests: fast (10-50ms)
46     - 15% of requests: medium (100-300ms)
47     - 5% of requests: very slow (500-1500ms)
48     """
49     random_value = random.random()
50
51     if random_value < 0.8:
52         await asyncio.sleep(random.uniform(0.01, 0.05))
53         category = "fast (80%)"
54     elif random_value < 0.95:
55         await asyncio.sleep(random.uniform(0.1, 0.3))
56         category = "medium (15%)"
57     else:
58         await asyncio.sleep(random.uniform(0.5, 1.5))
59         category = "very slow (5%)"
60
```



Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast

 www.asanchezyali.com

```
61     return {"response_type": "variable", "category": category}
62
63 if __name__ == "__main__":
64     import uvicorn
65     uvicorn.run("api_server:app", host="0.0.0.0", port=8000,
        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
5 import matplotlib.pyplot as plt
6 import numpy as np
7 from collections import defaultdict
```



Alejandro Sánchez Yalí


Software Developer | AI & Blockchain Enthusiast

 www.asanchezyali.com

```
8
9 class LoadTester:
10     def __init__(self, base_url="http://localhost:8000"):
11         self.base_url = base_url
12         self.endpoints = {
13             "fast": "/fast",
14             "medium": "/medium",
15             "slow": "/slow",
16             "variable": "/variable"
17         }
18         self.results = defaultdict(list)
19
20     async def make_request(self, session, endpoint):
21         """Makes an HTTP request and measures response time"""
22         url = f"{self.base_url}{self.endpoints[endpoint]}"
23         start_time = time.time()
24
25         try:
26             async with session.get(url) as response:
27                 await response.json()
```

**Alejandro Sánchez Yalí**

Software Developer | AI & Blockchain Enthusiast

 www.asanchezyali.com


```
28         response_time = (time.time() - start_time) * 1000
29         self.results[endpoint].append(response_time)
30         return response_time
31     except Exception as e:
32         print(f"Error in request to {url}: {e}")
33         return None
34
35     async def generate_load(self, endpoint, num_requests,
36                             concurrency):
37         """Generates load for a specific endpoint"""
38         async with aiohttp.ClientSession() as session:
39             tasks = []
40             for _ in range(num_requests):
41                 tasks.append(self.make_request(session, endpoint))
42                 if len(tasks) >= concurrency:
43                     await asyncio.gather(*tasks)
44                     tasks = []
45
46             if tasks:
47                 await asyncio.gather(*tasks)
```



Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast


 www.asanchezyali.com

```
47
48     def calculate_percentiles(self, endpoint):
49         """Calculates percentiles for response times"""
50         if not self.results[endpoint]:
51             return {}
52
53         data = sorted(self.results[endpoint])
54         return {
55             "min": min(data),
56             "p50": statistics.median(data),
57             "p90": np.percentile(data, 90),
58             "p95": np.percentile(data, 95),
59             "p99": np.percentile(data, 99),
60             "p999": np.percentile(data, 99.9) if len(data) >= 1000
        else None,
61             "max": max(data),
62             "mean": statistics.mean(data),
63             "stdev": statistics.stdev(data) if len(data) > 1 else 0
64         }
65
```



Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast


 www.asanchezyali.com

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



Alejandro Sánchez Yalí


Software Developer | AI & Blockchain Enthusiast

 www.asanchezyali.com

```
85         ax1.grid(axis='y', linestyle='--', alpha=0.7)
86
87         if "variable" in self.results and self.results["variable"]:
88             data = self.results["variable"]
89             bins = np.logspace(np.log10(min(data)),
np.log10(max(data)), 50)
90             ax2.hist(data, bins=bins, alpha=0.7, color='green')
91             ax2.set_xscale('log')
92             ax2.set_title('Response Time Distribution (variable
endpoint)')
93             ax2.set_xlabel('Response Time (ms) - logarithmic scale')
94             ax2.set_ylabel('Number of Requests')
95
96             percentiles = self.calculate_percentiles("variable")
97             for metric, value in [(k, v) for k, v in
percentiles.items()
98                                     if k in ["p50", "p90", "p95", "p99"]
and v is not None]:
99                 ax2.axvline(x=value, color='red', linestyle='--',
alpha=0.6,
```

**Alejandro Sánchez Yalí**

Software Developer | AI & Blockchain Enthusiast


 www.asanchezyali.com

```
100             label=f"{metric}: {value:.2f}ms")
101
102         ax2.legend()
103
104         plt.tight_layout()
105         plt.savefig('latency_results.png', dpi=300)
106         plt.close()
107
108     def print_summary(self):
109         """Prints a summary of the results"""
110         for endpoint in self.results:
111             print(f"\n=== Endpoint: {endpoint}
112                   ({len(self.results[endpoint])} requests) ===")
113             percentiles = self.calculate_percentiles(endpoint)
114             for metric, value in percentiles.items():
115                 if value is not None:
116                     print(f"{metric}: {value:.2f} ms")
117
118     async def main():
119         tester = LoadTester()
```



Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast

 www.asanchezyali.com

```
119
120     print("Starting load tests...")
121
122     for endpoint in tester.endpoints:
123         requests = 1000 if endpoint == "variable" else 200
124         print(f"Testing endpoint '{endpoint}' with {requests}
requests...")
125         await tester.generate_load(endpoint, requests, concurrency=50)
126
127     tester.print_summary()
128     tester.plot_results()
129     print("Tests completed. Results saved in 'latency_results.png'")
130
131 if __name__ == "__main__":
132     asyncio.run(main())
133
134 # Output
135 # Starting load tests...
136 # Testing endpoint 'fast' with 200 requests...
137 # Testing endpoint 'medium' with 200 requests...
```



Alejandro Sánchez Yalí


Software Developer | AI & Blockchain Enthusiast

 www.asanchezyali.com

```
138 # Testing endpoint 'slow' with 200 requests...
139 # Testing endpoint 'variable' with 1000 requests...
140
141 # === Endpoint: fast (200 requests) ===
142 # min: 10.22 ms
143 # p50: 22.67 ms
144 # p90: 31.31 ms
145 # p95: 34.34 ms
146 # p99: 37.17 ms
147 # max: 38.25 ms
148 # mean: 22.79 ms
149 # stdev: 6.84 ms
150
151 # === Endpoint: medium (200 requests) ===
152 # min: 53.42 ms
153 # p50: 102.49 ms
154 # p90: 143.29 ms
155 # p95: 150.27 ms
156 # p99: 161.89 ms
157 # max: 163.58 ms
```

**Alejandro Sánchez Yalí**


Software Developer | AI & Blockchain Enthusiast

 www.asanchezyali.com

```
158 # mean: 103.06 ms
159 # stdev: 30.29 ms
160
161 # === Endpoint: slow (200 requests) ===
162 # min: 206.44 ms
163 # p50: 361.78 ms
164 # p90: 476.99 ms
165 # p95: 489.22 ms
166 # p99: 501.95 ms
167 # max: 505.02 ms
168 # mean: 360.42 ms
169 # stdev: 84.46 ms
170
171 # === Endpoint: variable (1000 requests) ===
172 # min: 11.65 ms
173 # p50: 38.14 ms
174 # p90: 239.78 ms
175 # p95: 297.38 ms
176 # p99: 1240.78 ms
177 # p999: 1452.90 ms
```

**Alejandro Sánchez Yalí**

Software Developer | AI & Blockchain Enthusiast

 www.asanchezyali.com


```
178 # max: 1491.94 ms
179 # mean: 103.07 ms
180 # stdev: 213.06 ms
181 # Tests completed. Results saved in 'latency_results.png'
```

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
```

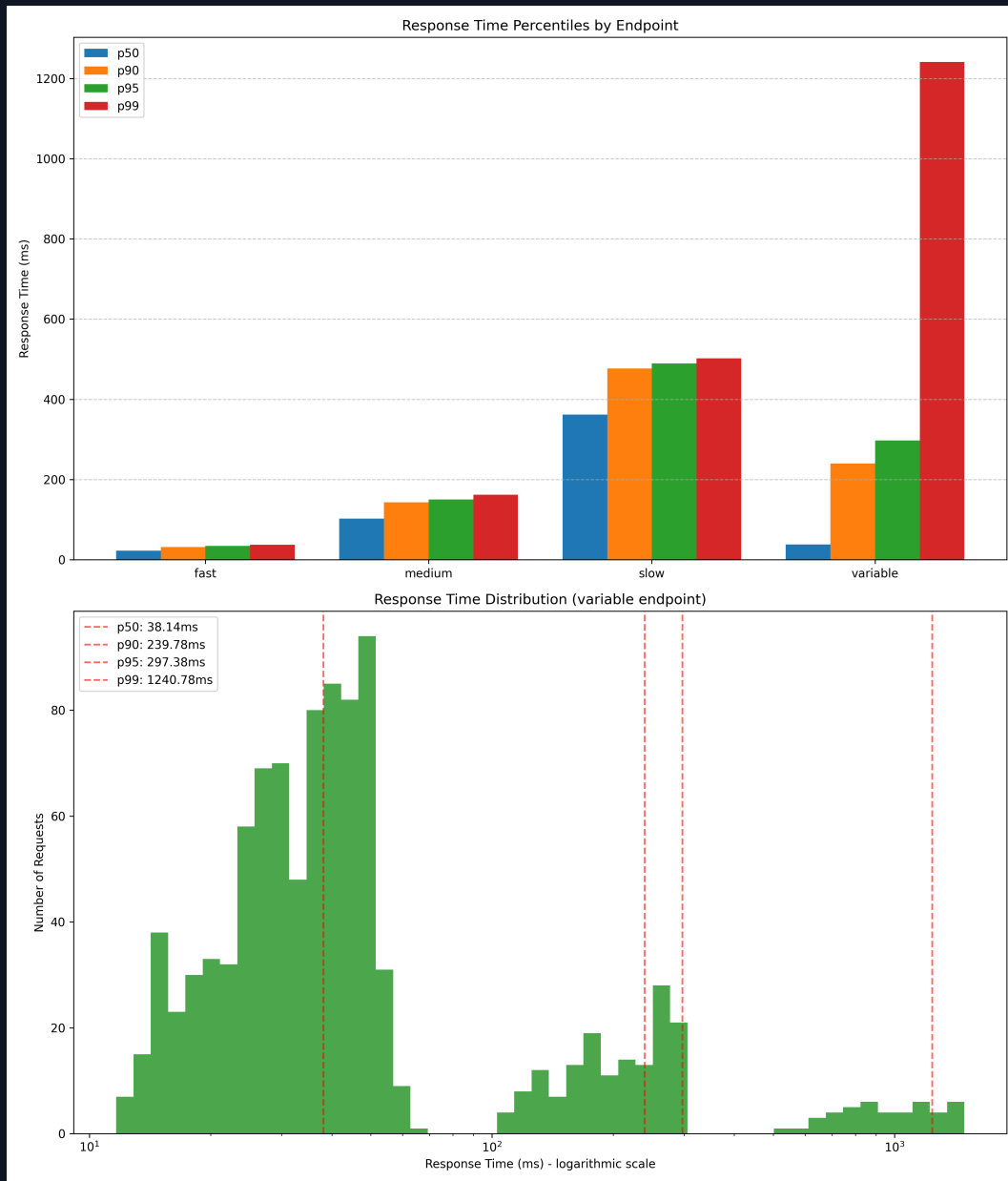
The agent will generate a file called `latency_results.png` with visualizations of the results, and will also print a summary to the console.




Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast

 www.asanchezyali.com

**Alejandro Sánchez Yalí**

Software Developer | AI & Blockchain Enthusiast

 www.asanchezyali.com

3. Analyzing the Results

The visualizations demonstrate why percentiles provide superior insights for API performance compared to averages:

3.1. Understanding the Variable Endpoint

The most revealing insights come from examining the **variable** endpoint:


- The **median (p50)** is only 38.14ms, indicating most users receive fast responses.
- The **p90** jumps to 239.78ms, revealing that 10% of requests experience significantly slower performance.
- The **p95** at 297.38ms shows further degradation for 5% of requests.
- Most critically, the **p99** at 1240.78ms demonstrates that 1% of users experience response times over 32 times slower than the median.
- The mean (103.07ms) obscures this reality, appearing deceptively moderate despite the extreme outliers.

The histogram's logarithmic distribution confirms these observations, showing three distinct clusters corresponding to the programmed response time categories (80% fast, 15% medium, 5% slow).



Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast

 www.asanchezyali.com

3.2. Comparative Analysis Across Endpoints

Comparing the four endpoints reveals important performance patterns:

- **Predictable endpoints** (fast, medium, slow) show relatively consistent behavior where p99 is only 1.3-1.4 times greater than p50.
- The **variable endpoint** exhibits dramatic tail latency, with p99 being 32.5 times higher than p50.
- While the **mean response time** of the variable endpoint (103.07ms) is nearly identical to the medium endpoint (103.06ms), their performance profiles are entirely different—a fact that would be missed by relying solely on averages.

3.3. Practical Benefits of Percentile-Based Monitoring


Using percentiles for performance monitoring offers several concrete advantages:

- **Detecting Hidden Issues:** The variable endpoint's mean suggests acceptable performance, while percentiles reveal severe degradation affecting a minority of requests.
- **User-Centric Metrics:** Percentiles directly correspond to user experience—p95 represents what 5% of your users are experiencing or worse.



Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast

 www.asanchezyali.com

- **Early Warning System:** Changes in high percentiles often precede system-wide degradation and can signal emerging problems before they affect most users.
- **SLA Alignment:** Service Level Agreements based on percentiles protect all users, while average-based SLAs may hide systematic failures affecting a subset of users.
- **Infrastructure Sizing:** Understanding tail latency helps properly size infrastructure for peak demands rather than average conditions.

The logarithmic histogram of the variable endpoint particularly emphasizes how response times are not normally distributed—they follow a multi-modal distribution with long tails. This pattern is common in real-world systems due to factors like cache misses, garbage collection pauses, or resource contention.

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:

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



Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast

 www.asanchezyali.com

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."

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:** Optimizing high percentiles is essential for multi-call user actions.
- **Establish SLOs and SLAs Based on Percentiles:** For example, "p99 must be less than 300ms" is more meaningful than "average time must be less than 100ms."

Modern applications must be designed with these principles to provide consistent, high-quality experiences to all users. Future posts will explore techniques like priority queues, circuit breakers, and caching to mitigate tail latency amplification.



Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast

 www.asanchezyali.com


6. References

- Sentry. (2024). *What's the difference between API Latency and API Response Time?*. [Link](#)
- Catchpoint. *API Performance Monitoring—Key Metrics and Best Practices*. [Link](#)
- FastAPI Documentation. *Benchmarks*. [Link](#)
- DeCandia, G., et al. (2007). *Dynamo: Amazon's Highly Available Key-value Store*. ACM SIGOPS Operating Systems Review, 41(6), 205-220.
- Kleppmann, M. (2017). *Designing Data-Intensive Applications: The Big Ideas Behind Reliable, Scalable, and Maintainable Systems*. O'Reilly Media, Inc.
- This article was translated, edited and written in collaboration with AI. If you find any inconsistencies or have suggestions for improvement, please don't hesitate to open an issue in our GitHub repository at [github](#) or reach out directly.



Alejandro Sánchez Yalí

Software Developer | AI & Blockchain Enthusiast

 www.asanchezyali.com

7. 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



Feedback

Found this helpful?

Save, comment and share

April 29, 2025