

### IMPROVING SERVICE RELIABILITY IN PYTHON UWSGI STACKS VIA PROACTIVE 500 ERROR HANDLING

Nikhita Kataria Independent Researcher Manager, Software Engineering nikhitakataria@gmail.com

#### Abstract

Python-based Flask applications are widely adopted in the software industry today due to their simplicity, flexibility, and ease of integration with other technologies. At peak usage, these applications may face an unsustainable volume of HTTP 500 errors, severely impacting end users by degrading service reliability, increasing latency, or causing complete request failures. These issues are typically triggered by resource exhaustion, unhandled edge cases in application logic, or limitations in the underlying uWSGI server configuration. This paper discusses a systematic approach to scaling and optimizing such applications by improving the reliability, scalability, and availability of the overall stack consisting of NGINX, Gunicorn, uWSGI, and MySQL. We present targeted optimizations and tuning strategies that can be generalized to other similar architectures.

Keywords – Nginx, uWSGI, Flask, MySQL, Caching, Rate Limiting, Connection tuning.

### I. INTRODUCTION

Imagine a scenario in which a Python-based application, built using the Flask framework, is experiencing a high frequency of 5xx server errors. These errors are indicative of deeper systemic issues related to the application's ability to reliably handle incoming requests. Such instability not only undermines user experience but also signals potential architectural or operational weaknesses. This paper aims to provide a comprehensive and methodical approach to identifying, diagnosing, and mitigating these server-side failures. It outlines a set of practical, field-tested strategies designed to improve service resilience and reliability. The proposed solutions are targeted toward a diverse range of stakeholders, including site reliability engineers (SREs), backend developers, and infrastructure teams. By implementing these recommendations, teams can significantly reduce the incidence of server errors, address root causes of recurring failures, and strengthen the overall fault tolerance of their applications.

### II. TARGET ARCHITECTURE

The target architecture for this research is whereclient requests are first handled by NGINX, which serves as a reverse proxy. NGINX forwards requests over a Unix domain socket – a low-overhead alternative to TCP when processes run on the same host – to the application server. The application server, Gunicorn, is a uWSGI-compatible server that manages multiple worker processes and routes incoming requests. Gunicorn passes these requests to a Flask application, which implements the core business logic and returns responses.



This layered setup – NGINX  $\rightarrow$  Unix socket  $\rightarrow$  Gunicorn  $\rightarrow$  Flask – improves performance, isolates failures, and allows each component to be scaled or tuned independently. On the backend, the Flask application interfaces with a data persistence layer – in this case, a MySQL database – responsible for managing application state. Each component in this stack introduces distinct performance bottlenecks and potential points of failure, which are analyzed and mitigated individually to enhance overall system reliability. Fig. 1. illustrates high level architecture experimented in this paper and subsequent sections cover optimizations for different layers.



Fig. 1. Architecture Overview

## III. CLIENT TO NGINX OPTIMIZATION

Often rate bursts from clients are a leading cause of service degradation. Two main techniques to solve this are rate limiting and connection tuning.

### A. Rate Limiting

Rate limiting can be controlled viathelimit\_req module provided by NGINX to control the rate of incoming requests by enforcing per-client rate limits. Defining a default burst size for all clients is crucial, as it ensures fair resource distribution and prevents any single client from adversely affecting the accessibility and performance experienced by others accessing the system concurrently.

One plausible configuration is in Fig. 2. Which applies the mylimit zone (which defines the rate limit) and allows up to 20 excess requests to be queued before applying the rate limit.



Fig. 2. Nginx Configuration for rate limiting.

### **B.** Connecting Tuning

In addition to rate limiting, consider further connection tuning to handle incoming traffic efficiently via the configurations outlined in Table 1. There are multiple configuration knobs in Nginx and these optimizations were experimented with to streamline request flows and reduce overhead.



limit_conn,	Restricts connections per client to			
limit_conn_zone	handle bursty traffic for a client			
	with high traffic.			
limit_rate	Limits the rate of response			
	transmission to the clients.			
	Setting a limit can prevent overall			
	system from being overloaded			
	ensuring QoS for all clients.			
max_conns	Controls the maximum number			
	of connections an upstream			
	accepts. Notably a 0 here means			
	no limit.			
queue(NGINXPlus)	Paired with max_conns to queue			
,	overflow requests.			

Table I. Nginx Configuration for client optimization

### IV. NGINX TO WSGI OPTIMIZATION

To scale the connection between NGINX and uWSGI effectively, it is crucial to understand how NGINX handles upstream communication with uWSGI through the ngx\_http\_uwsgi\_module. This module provides a wide range of configurable parameters – such as buffer sizes, queue limits, and rate limits – that allow fine-grained control over how requests are forwarded from NGINX to the uWSGI server. Proper tuning of these parameters can significantly impact the overall responsiveness, stability, and scalability of a web application.

Prior to delving into the intricacies and specifics of these configuration and tuning options, let's take a step-by-step journey through a clear and simplified illustrative example that effectively showcases the fundamental performance dynamics and behaviors inherent in a typical Flask web application when it is deployed using uWSGI as the application server and positioned behind a NGINX web server. Assume the following scenario:

- 1. Queries Per Second (QPS): 500 The application is receiving 500 incoming requests per second.
- 2. Average Response Time: 2 milliseconds Each request, on average, takes 2 ms to process.
- 3. Total Processing Window: 1 second We analyze the system over a one-second interval.
- 4. Total uWSGI Workers: 10 The application is served by 10 uWSGI worker processes.

This example helps highlight a key concern in scaling backend services: even small deviations from average request latency can cause significant ripple effects. For instance, if just 2% of incoming queries—around 10 requests in this case—become slow due to downstream issues (e.g., I/O delays or DB lock contention), they can block worker threads longer than expected. Since uWSGI does not support true concurrency for most Python workloads, those slow queries effectively reduce worker availability and can lead to queuing delays or even HTTP 502 errors if the queue overflows or times out.





Fig. 3.Request times in milliseconds with 10 and 20 workers, qps: 10, total queries: 100

Fig. 3.presents a scenario involving a slow query with a latency of approximately 5 seconds. With 10 uWSGI workers and a query-per-second (QPS) rate of 10, nearly 50% of the 100 total queries experienced latencies approaching 10 seconds. Increasing the number of uWSGI workers to 20 resulted in a noticeable improvement in response times. To address performance bottlenecks and improve the reliability of NGINX-to-uWSGI communication, several practical strategies can be implemented. These configurations help ensure the application remains responsive, especially under high load, and can prevent common issues such as 502 errors and worker saturation. To solve such scenarios, employ techniques outlined in next set of points.

#### A. Enable Caching

One of the most effective ways to reduce application load and accelerate response times is to implement caching at the NGINX level. By using the uwsgi\_cache directive, NGINX can store the results of frequently accessed uWSGI responses. This means that for repeat requests, NGINX can serve the cached content directly without forwarding the request to the backend application.Caching is particularly beneficial for:

- 1. Static or rarely-changing dynamic content.
- 2. Expensive computations that don't need to be re-evaluated for each request.
- 3. Reducing latency for end users by cutting down on application response time.



Fig. 4.Request times in milliseconds with and without caching enabled.



Fig. 4.depicts the distribution of request latencies for a simple GET operation against a Flask application that retrieves a record set from a MySQL database. The -axis shows latency percentiles, while the y-axis indicates response times in milliseconds. Under a workload of 100 concurrent connections at 1000 QPS, serving a 17,324-byte payload, observed latencies span from roughly 2,500 ms at the highest percentiles down to approximately 6 ms at the lower percentiles.

Proper cache invalidation policies and cache zones must be defined to ensure consistency and memory efficiency. When configured appropriately, caching can drastically lower backend CPU usage and increase system throughput.

B. Increase Buffer Sizes

NGINX communicates with uWSGI over a socket and uses buffers to temporarily store response data. When a response from uWSGI is larger than the default buffer size and no appropriate buffer expansion is configured, NGINX may fail to process the response and return a 502 Bad Gateway error.

Nginx offers key directives to handle such errors by increasing the buffer size and increasing these values ensures that NGINX has enough memory to process larger responses without dropping the connection or returning an error to the user. This is particularly important for applications returning large payloads, such as JSON APIs or file downloads.

- 1. uwsgi\_buffers: The number and size of the buffers.
- 2. uwsgi\_buffer\_size: The size of the initial buffer used before additional buffers are allocated.
- 3. uwsgi\_busy\_buffers\_size: The size limit for buffers that can be busy serving a response before client finishes reading it.



Fig. 5. illustrates the performance degradation that occurs when buffering limits are set too low. In this experiment, we enabled buffering and configured

- 1. uwsgi\_buffers: 8 buffers of size 128 bytes each
- 2. uwsgi\_buffer\_size: 128 bytes
- 3. uwsgi\_busy\_buffers\_size: 256 bytes.
- 4. The system was subjected to a query-per-second (QPS) rate of 10, processing a total of 100 requests.



#### C. Throttle Response Read Speed from uWSGI

In high-throughput environments, it's important to control how quickly NGINX reads responses from uWSGI. Reading responses too quickly can lead to resource contention and downstream system overload, especially if the backend or database cannot keep up with the traffic.

The uwsgi\_limit\_rate directive can be used to limit the data read rate from uWSGI, effectively throttling the response stream. Fig. 6.illustrates the outcomes of experiments performed under a concurrent query rate of 10 QPS, comprising 100 total requests and response size of 40 KB with uwsgi\_limit\_rate set to 10K, 20K and 40K. Note with these settings the response time does not reduce proportional to the decrease in limit.



Fig. 6. Request time in ms with 10K, 20K and 40K limit

Throttling responses can:

- 1. Prevent backend saturation by smoothing spikes in traffic.
- 2. Reduce the likelihood of worker thread starvation or socket queuing delays.
- 3. Contribute to more stable and predictable system performance under heavy load.

By applying the optimizations captured in above points, the system can better handle high traffic, avoid bottlenecks, and improve overall performance and reliability.

### V. UWSGI TO FLASK APPLICATION OPTIMIZATION

To minimize worker unavailability and ensure graceful degradation, consider the following approaches to ensure that the system remains available and responsive, providing users with meaningful feedback even during failure conditions.

#### A. Exception Handling

Uncaught exceptions can cause worker processes to crash or restart, leading to degraded performance, increased latency, and service instability. Implementing structured and context-aware exception handling helps ensure that your application can gracefully recover from common failure scenarios without triggering a full worker restart. Instead of crashing or returning a generic 500 Internal Server Error, your application should return meaningful and standardized HTTP



status codes that reflect the nature of the issue. This enables better observability, easier debugging, and more resilient client-side behavior.

Here are some examples of common failure scenarios and the recommended HTTP status codes to return:

- 1. Deadlocks: If a database operation results in a deadlock, catch the exception and return a 409 Conflict. This indicates that the request could not be completed due to a conflict with the current state of the resource. Clients can be designed to retry such requests intelligently.
- 2. Too Many Requests: When rate limiting or throttling mechanisms are triggered, return a 429 Too Many Requests. This status code informs clients that they have exceeded their allowed request quota and should back off or retry after a specified period.
- 3. Connection Timeout: If a dependency such as an upstream API or database becomes unresponsive and a timeout occurs, respond with a 502 Bad Gateway. This indicates that your service is acting as a proxy or gateway and received an invalid response (or no response) from the upstream server.

By mapping exceptions to meaningful status codes, you can avoid abrupt restarts, maintain service continuity, and provide better diagnostics. Consider logging stack traces and context for unexpected exceptions while using structured error responses to aid downstream services and client applications in decision-making.

### **B.** Queue Scaling

The uwsgi.listen parameter controls the size of the socket's listen queue—that is, how many incoming connections can be queued before the application starts rejecting them. By default, this value is often set conservatively (e.g., 100), which may be insufficient for high-traffic applications or during sudden bursts in load.

Increasing the uwsgi.listenvalue allows your application to handle more simultaneous incoming connections by queuing them instead of dropping them immediately when all workers are busy. This is especially important in scenarios where response times might spike temporarily, such as during cache misses, upstream latency, or CPU contention.However, simply increasing uwsgi.listen is not enough. You must also ensure that the underlying operating system's networking stack is configured to support a backlog of at least the same size.

### VI. FLASK $\leftrightarrow$ MYSQL OPTIMIZATION

Optimizing MySQL queries is dependent on the nature of traffic an application is serving. In certain cases, with high read throughput query optimization by setting up correct indexes and tuning the schema should result in enhanced performance. Following techniques are well tested to ensure performance gains.

### A. Partitioning

When designing your database, ask yourself whether it's necessary to store all data in a single database. If you decide that it is, consider partitioning the data. A common approach in MySQL is partitioning by date, which works well for time-series data and is easy to implement.



### **B.** Pagination

Efficient pagination is key to managing large datasets. In this paper, we initially used offset pagination but later switched to keyset pagination. This shift provided several advantages, including improved query performance and more consistent results.

### C. Revisiting Foreign Keys

While foreign keys enforce referential integrity, they can introduce performance overhead, particularly in write-heavy

applications. In this paper, removing unnecessary foreign keys reduced the strain on CRUD operations, as it eliminated the need to constantly evaluate integrity constraints.

### **D.** Revisiting Indexes

Indexes enhance query performance, but they can negatively impact write operations – such as updates, inserts, and deletes – because each change must also update the associated indexes. Regularly reviewing and cleaning up redundant or unnecessary indexes is essential for maintaining efficiency. Keep in mind that primary keys automatically create indexes, so adding extra indexes on the same columns may be redundant.

### **E.** Timeouts

Long-running queries (those taking more than 10 seconds) can significantly impact database performance. Explore solutions for automatically timing out long-running queries, helping maintain system efficiency. Putting such solutions via simple scripts can enhance client experience multi-fold by ensuring certain clients running heavy queries do not halt the system for clients with quicker queries.

## F. Stress Testing

It is essential to ensure that the data layer—often the most critical and performance-sensitive component of a web application—can sustain the anticipated query load under both normal and peak operating conditions. An inadequately optimized database can become a single point of failure, leading to degraded performance. As an example, Fig. 7. Illustrates the use of mysqlslap for adhoc load testing.

mysqlslap	user= <username></username>		create-		
schema= <data< td=""><td>base&gt;</td><td>query=«</td><td>your_</td><td>query&gt;</td><td>-h</td></data<>	base>	query=«	your_	query>	-h
<hostname> -p</hostname>	oco:	ncurrency	100		

Fig. 7. Sample usage of mysqlslap.

## VII. CONCLUSION

Improving reliability, scalability, and availability across each layer of a Python-based web stack delivers compounding and synergistic benefits that extend beyond isolated optimizations. When tackled holistically, enhancements at the application, web server, and infrastructure levels reinforce one another, creating a robust system capable of withstanding varied load patterns and failure scenarios. This paper outlines a comprehensive blueprint for engineering teams working on



similar Python web applications that encounter performance bottlenecks, service degradation, or availability issues.

The techniques and architectural adjustments described herein are grounded in pragmatic, realworld scenarios, ensuring they are not only technically sound but also operationally feasible. They span improvements in request handling, process management, resource utilization, and observability—each selected to strike a balance between implementation effort and measurable impact.

Prior to implementing the optimizations detailed in this paper, the application under studyexhibited significant instability, characterized by a recurringspike of HTTP 5xx errorsaveraging several per minute—especially during peak traffic hours. This level of unreliability posed both user experience risks and operational burdens.

Following a methodical application of the strategies presented – ranging from tuning the NGINX and Gunicorn/uWSGI layers to refactoring Python code paths and improving database access patterns – the application achieved notable stability. Post-optimization monitoring revealed that 5xx reduced by approximately 98%, accompanied by faster response times and reduced infrastructure strain. These improvements not only stabilized the platform but also enabled it to scale gracefully under high concurrency, demonstrating high availability and consistent performance even during traffic surges.

### VIII. LIMITATIONS

While this work provides useful insights and improvements, there are a few important limitations to keep in mind:

- 1. Our findings focus mainly on Python web apps running with uWSGI (or Gunicorn), NGINX, and MySQL. If you're using different languages, frameworks, or databases, the results might not apply directly.
- 2. The testing and recommendations assume everything runs on a single server or closely connected machines using Unix sockets. If your app runs on multiple servers, containers, or uses complex load balancing, some suggestions may need adjustment.
- 3. We mostly look at 5xx errors caused by app crashes or overloaded servers. Other problems like network failures or database outages aren't covered here in depth.
- 4. To catch and fix issues effectively, your system should already have decent logging and monitoring. Without those, it'll be harder to benefit fully from these strategies.
- 5. Our tuning assumes moderate to high traffic (around 500 requests per second) and reasonable hardware. If your app faces much heavier loads or limited resources, you might need extra tweaks beyond what's described.

### IX. ASSUMPTIONS

In the research outline in this paper, we make following assumptions:

- 1. The architecture outlined in this paper is assumed to be followed as is of a conventional web application backed via Nginx and Mysql as the backend data store with unix domain sockets being assumed to be used as communication mechanism
- 2. It is assumed that the system experiences spikes upto 500 QPS for the study however the proposed solutions are applicable for QPS higher than this. Average response time are also



assumed to be well under 5ms.

- 3. It is assumed that the application otherwise has enough resources in terms of cpu and memory available and the optimization points are between service intersection points.
- 4. It is assumed that sufficient observability is in place to identify operational issues between layers before proceeding to applying techniques outlined.

#### REFERENCES

- 1. NGINX. (n.d.). ngx\_http\_limit\_req\_module. NGINX Documentation. Retrieved May 13, 2025, from http://nginx.org/en/docs/http/ngx\_http\_limit\_req\_module.html
- 2. Oracle. (n.d.). mysqlslap. MySQL 8.4 Reference Manual. Retrieved May 13, 2025, from https://dev.mysql.com/doc/refman/8.4/en/mysqlslap.html
- 3. NGINX. (n.d.). Tuning NGINX for performance. NGINX Blog. Retrieved May 13, 2025, from https://www.nginx.com/blog/tuning-nginx/
- 4. Pallets Projects. (n.d.). Flask documentation. Retrieved May 13, 2025, from https://flask.palletsprojects.com
- 5. uWSGI Project. (n.d.). uWSGI documentation. Retrieved May 13, 2025, from https://uwsgidocs.readthedocs.io
- 6. Oracle. (n.d.). MySQL performance tuning. MySQL Reference Manual. Retrieved May 13, 2025, from https://dev.mysql.com/doc/refman/
- 7. Hadi, P. (n.d.). Gist: Flask uWSGI NGINX configuration. GitHub Gist. Retrieved May 13, 2025, from https://gist.github.com/prasetiyohadi/c24112871943aa21d1bc