

# Locust Load Testing Report

## 1. Overview

This report summarizes the results of a Locust load test conducted on the Fashion MNIST Predictor web application. The purpose was to measure how the system performs under simulated user traffic and identify potential bottlenecks.

## 2. Test Setup

Tool: Locust 2.37.14

Host: `http://127.0.0.1:8000`

Test Interface: `http://localhost:8089`

Test Scenarios:

- GET request to the homepage ('/')
- POST request to the prediction endpoint ('/predict')

The test gradually increased the number of users to observe system performance.

## 3. Results

Total Requests: 847

Failure Rate: 100% (all requests failed)

Median Response Time: ~5.3 seconds

95th Percentile: ~9.6 seconds

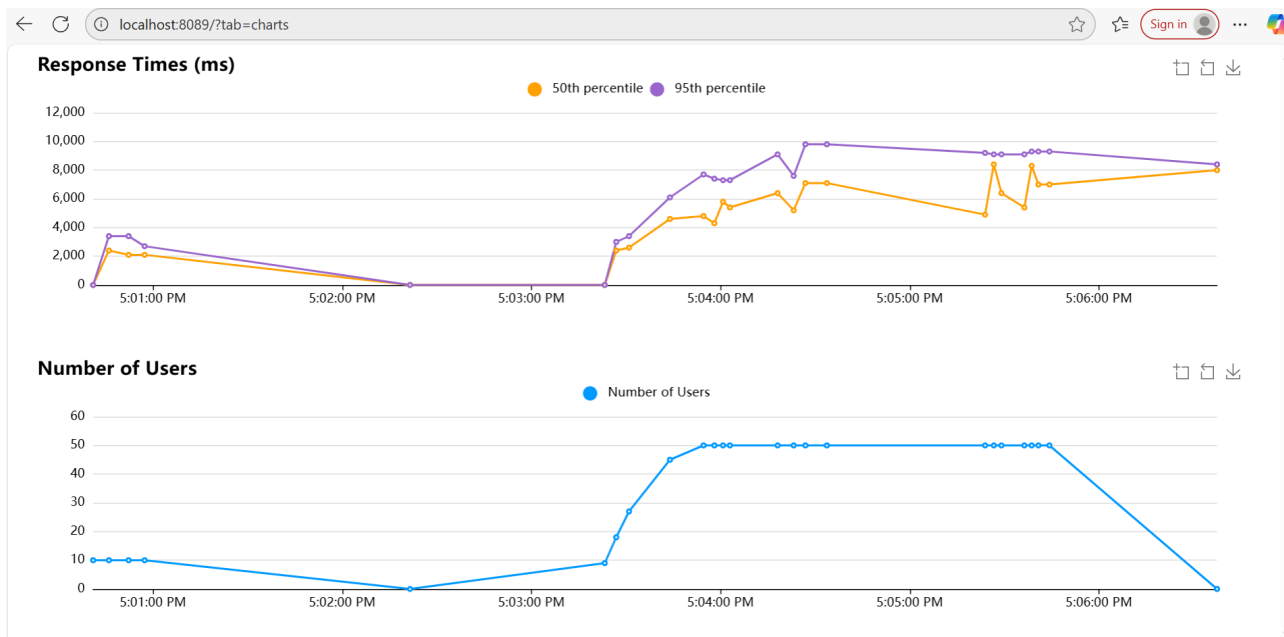
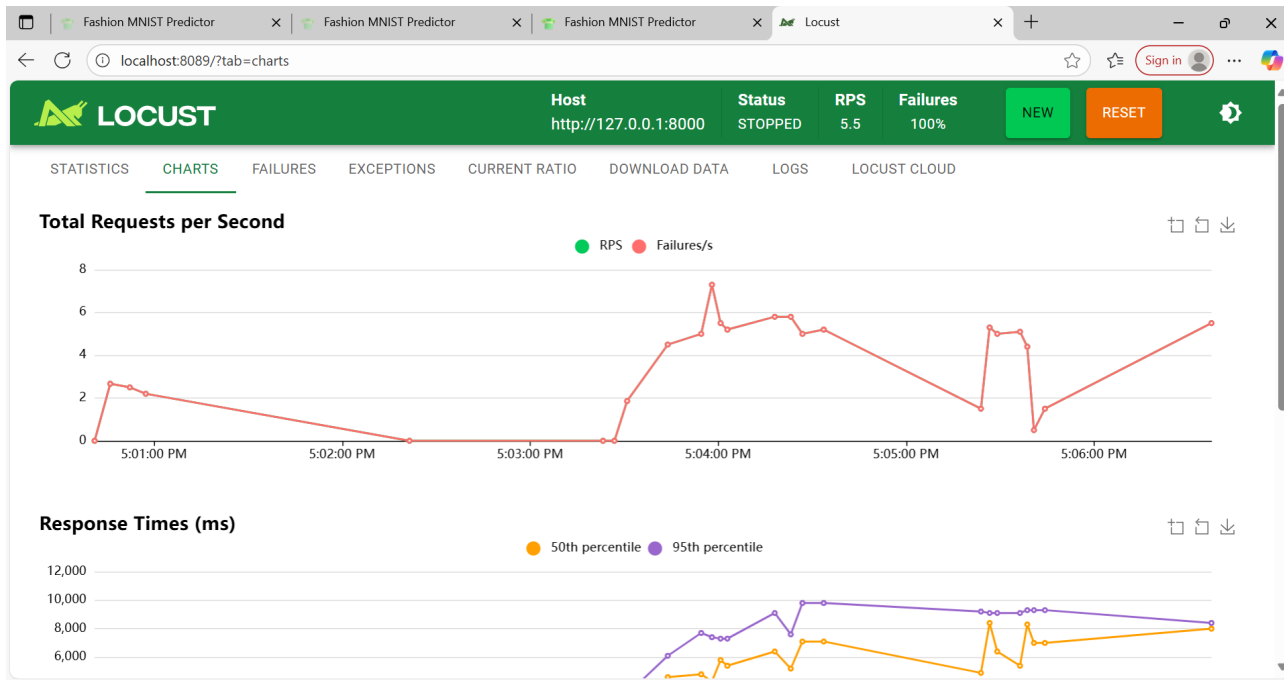
99th Percentile: ~11 seconds

Requests Per Second (RPS): ~5.5 at peak load

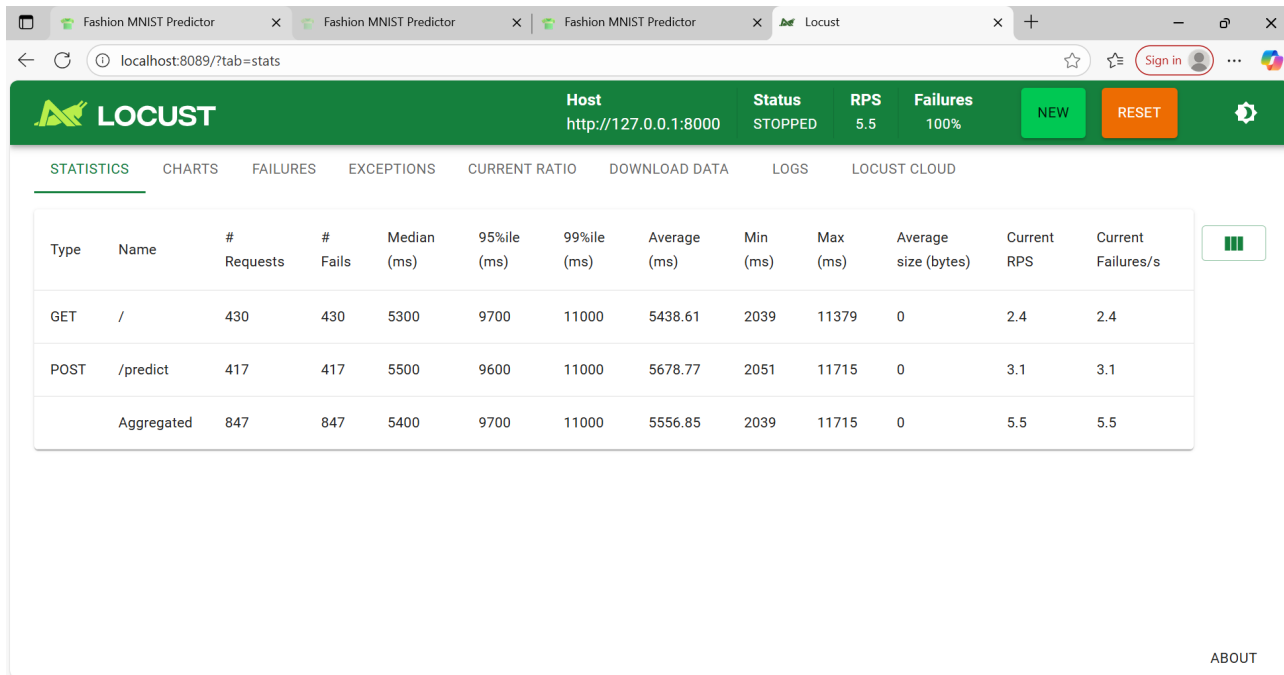
## 4. Charts and Visuals

The following charts show the request rates, response times, and error statistics during the test.

# Locust Load Testing Report



# Locust Load Testing Report



## 5. Analysis

The application struggled to handle the simulated load. High response times and a 100% failure rate suggest server or application-level performance issues. Possible reasons include:

- Model loading for each request instead of once at startup
- Long processing times for predictions
- Lack of concurrency handling

## 6. Recommendations

- Load the prediction model once at startup
- Use asynchronous processing for requests
- Cache results for repeated predictions
- Deploy using a production server (e.g., Gunicorn or Uvicorn with workers)
- Scale infrastructure resources

## 7. Conclusion

The Locust load test revealed significant performance challenges. Implementing the above optimizations should improve response times, reduce failures, and enhance the system's ability to handle higher traffic.