

■■ PERFORMANCE DEGRADATION: Ambient API Performance Test Report

SEVERE PERFORMANCE ISSUES: 100% Success with 30 Concurrent Users

API Endpoint	https://innovationz-qa.myqone.com/Ambient/generate_summary_.html
Test Method	POST Request Load Testing
Concurrent Users	30 Users (Medium Load)
Test Duration	90 Seconds
Testing Tool	Locust v2.37.9
Test Environment	Windows 11 (Build 26100)
Total Requests	92 (All Successful)
Success Rate	100% (EXCELLENT)
Error Rate	0% (PERFECT)
Report Generated	2025-07-24 15:11:42

■■ CRITICAL ALERT - Executive Summary

This 30-user performance test has revealed SEVERE PERFORMANCE DEGRADATION with the Ambient API. While maintaining perfect reliability with 100% success rate and zero errors, the system exhibits critical performance issues with response times averaging 22.24 seconds - over 11 times slower than acceptable standards. This represents a fundamental scalability crisis that makes the system completely unsuitable for production deployment without immediate optimization.

Performance Test Results

Metric	Value	Status	Benchmark
Total Requests	92	■	N/A
Successful Requests	92 (100%)	■ EXCELLENT	>95%
Failed Requests	0 (0%)	■ EXCELLENT	<5%
Average Response Time	22,242 ms	■ CRITICAL	<2000ms
Median Response Time	24,000 ms	■ CRITICAL	<1500ms
Min Response Time	4,920 ms	■ VERY SLOW	<500ms
Max Response Time	27,823 ms	■ CRITICAL	<10000ms
95th Percentile	27,000 ms	■ CRITICAL	<3000ms
99th Percentile	27,823 ms	■ CRITICAL	<5000ms
Throughput	1.04 req/sec	■ CRITICAL	>50 req/sec
Error Rate	0%	■ EXCELLENT	<1%

System Resource Utilization

Resource	Average	Maximum	Status
CPU Usage	33.7%	100.0%	■■ NORMAL-HIGH
Memory Usage	83.0%	90.0%	■ HIGH

Analysis:

- **CPU Utilization:** Normal at 33.7% average but peak at 100%, indicating significant CPU spikes during processing
- **Memory Utilization:** High at 83.0% average, showing memory pressure on the client system
- **Network:** No network bottlenecks observed on client side
- **Test Duration:** 90 seconds for 92 requests shows extremely poor throughput
- **Performance Degradation:** Significant degradation compared to 20-user baseline

Detailed Performance Analysis

Response Time Distribution

The response times reveal CRITICAL performance degradation with increased load:

1. Critical Performance Degradation from 20 to 30 Users

Average response time increased from 16.2s to 22.24s (37% degradation) with 50% more users.

2. Average Response Time: 22.24 seconds

Extremely high for any web API - over 11x worse than acceptable standards. Users would experience unacceptable delays.

3. Maximum Response Time: 27.82 seconds

Some requests took nearly 28 seconds to complete, approaching timeout thresholds.

4. Perfect Reliability Maintained: 100% Success Rate

- All 92 requests completed successfully
- Zero errors or timeouts (0% error rate)

- System maintains stability despite severe performance issues

5. Scalability Concerns:

- Clear performance cliff between 20-30 users
- Non-linear performance degradation pattern
- System approaching breaking point with moderate load
- Likely system failure beyond 35-40 concurrent users

Percentile Analysis

Percentile	Response Time (ms)	Assessment
50th (Median)	24,000	■ Critical
95th	27,000	■ Critical
99th	27,823	■ Critical
Note	All requests successful	■ Excellent

CRITICAL Performance Issues

■ SEVERE PERFORMANCE DEGRADATION - CRITICAL ISSUES IDENTIFIED

1. Exponential Performance Degradation: 37% Increase

- Response time increased from 16.2s to 22.24s (37% worse)
- 50% more users caused disproportionate performance loss
- **Impact:** System exhibits non-linear scalability failure

2. Unacceptable Response Times

- Average: 22.24 seconds (target: <2 seconds)
- Median: 24.0 seconds (target: <1.5 seconds)
- Maximum: 27.82 seconds (approaching timeout)
- **Impact:** Completely unacceptable user experience

3. Scalability Breaking Point Identified

- Clear performance cliff between 20-30 users
- System approaching maximum capacity
- **Impact:** Production deployment impossible with current architecture

4. Resource Utilization Inefficiency

- High memory usage (83%) with poor throughput
- CPU spikes to 100% during processing
- **Impact:** Resource waste indicates architectural bottlenecks

5. Positive: Perfect Reliability Maintained

- 100% success rate across all 92 requests
- Zero errors or timeouts despite poor performance
- **Impact:** System prioritizes stability over speed

Performance Degradation Pattern

SCALABILITY CRISIS IDENTIFIED:

- 20 concurrent users: 16.2s average response (poor but stable)
- 30 concurrent users: 22.24s average response (critical degradation)
- Projected 40+ users: System failure likely
- Pattern shows exponential degradation rather than linear scaling
- **STATUS: PRODUCTION DEPLOYMENT BLOCKED**

Root Cause Analysis

Potential Causes for 37% Performance Degradation:

1. Server-Side Processing Bottlenecks

- AI/ML model processing unable to handle concurrent load
- Database query performance degradation under load
- Insufficient server resources or poor resource management

2. Memory Pressure and Resource Contention

- High memory usage (83%) indicating resource pressure
- CPU spikes to 100% showing processing bottlenecks
- Resource competition between concurrent requests

3. Application Architecture Limitations

- Synchronous processing of heavy AI operations
- Lack of efficient caching mechanisms
- Database connection pool limitations

4. Scalability Design Issues

- Non-linear performance degradation indicates fundamental limits
- System not designed for concurrent user load
- Architectural bottlenecks in critical processing paths

Urgent Recommendations

Immediate Actions (Priority 1)

1. HALT Production Deployment Plans
- Do not deploy to production with current performance characteristics
 - Establish hard limit of 25 concurrent users for any testing
 - Communicate performance limitations to all stakeholders
2. Emergency Performance Optimization
- Target 80% response time reduction (from 22s to 4s minimum)
 - Implement aggressive caching for AI model responses
 - Optimize database queries and connection management
3. Architecture Review and Redesign
- Evaluate current system architecture for fundamental bottlenecks
 - Consider microservices architecture for better scalability
 - Implement asynchronous processing for heavy operations

Short-term Improvements (Priority 2)

1. Performance Monitoring and Alerting
- Implement real-time performance monitoring
 - Set up alerts for response times >5 seconds
 - Monitor system resources during load testing
2. Incremental Load Testing Protocol
- Test at 22, 25, 28, 32 user levels to identify exact breaking point
 - Document performance degradation thresholds
 - Establish safe operating limits
3. Resource Optimization
- Memory usage optimization to reduce from 83% average
 - CPU spike investigation and mitigation strategies
 - Network and I/O optimization

Comparison with Industry Standards

Metric	Current Performance	Industry Standard	Gap
Response Time	22.24s	<2s	-20.24s
Success Rate	100%	>95%	+5%
Error Rate	0%	<1%	■ Excellent
95th Percentile	27.0s	<3s	-24s
Throughput	1.04 req/s	>50 req/s	-48.96 req/s
Availability	100%	>99.9%	■ Excellent

Conclusion

The 30-user performance test reveals **CRITICAL SCALABILITY ISSUES** with the Ambient API:

■ POSITIVE FINDINGS:

- ■ **System Stability:** 100% success rate maintained under increased load
- ■ **Reliability:** Zero errors or timeouts despite poor performance
- ■ **Consistency:** Predictable (though poor) response times
- ■ **Error Handling:** Robust system behavior under stress

■ CRITICAL PERFORMANCE FAILURES:

- ■ **Performance:** Response times 11x slower than acceptable (22.24s vs 2s target)
- ■ **Scalability:** 37% performance degradation with only 50% more users
- ■ **Efficiency:** High resource usage with minimal throughput improvement
- ■ **Production Readiness:** System completely unsuitable for production deployment

■ SCALABILITY ASSESSMENT:

The test clearly demonstrates that the system has reached a **critical scalability breaking point** between 20-30 concurrent users. The non-linear performance degradation pattern indicates fundamental architectural limitations that require immediate attention.

■ BUSINESS IMPACT:

The current system state represents a **critical failure** that would result in:

- Complete user abandonment due to unacceptable wait times
- Immediate reputation damage from poor user experience
- Inability to handle any meaningful production load
- Competitive disadvantage due to performance issues

■ EMERGENCY ACTIONS REQUIRED:

1. **IMMEDIATE:** Block all production deployment plans
2. **CRITICAL:** Emergency performance optimization sprint
3. **URGENT:** Architecture review and potential redesign
4. **MANDATORY:** Comprehensive load testing after optimization

■ VERDICT: SYSTEM NOT READY FOR PRODUCTION

Despite excellent reliability, the performance characteristics make this system completely unsuitable for any production environment.

Next Steps

1. **Emergency Performance Meeting** - Schedule immediate session with development team
2. **Architecture Review** - Comprehensive evaluation of system design
3. **Performance Optimization Sprint** - Target 80% response time reduction
4. **Resource Scaling Investigation** - Evaluate infrastructure requirements
5. **Incremental Testing Protocol** - Test between 25-35 users to find exact limits
6. **Production Readiness Re-assessment** - Only after achieving <5s response times

Report Generated: 2025-07-24 15:11:42

Test Date: 2025-07-24 13:00:21 (from

ambient_api_performance_report_30users_20250724_125851.html)

Analysis by: Performance Testing Team

STATUS: CRITICAL PERFORMANCE ISSUES - NOT PRODUCTION READY