

Key Performance Engineering Metrics Every Engineer Should Track

Performance engineering focuses on ensuring applications run efficiently under expected workloads. Here are **the most critical performance engineering metrics** to track, categorized into **response time, throughput, resource utilization, error rates, and scalability metrics**.

1. Response Time Metrics (Latency & Performance)

Definition: The time taken for a request to be processed end-to-end.

- **Average Response Time** → Measures typical request latency.

$$\text{Avg Response Time} = \frac{\sum \text{Response Times}}{\text{Total Requests}}$$

- **Percentile Response Times (90th, 95th, 99th)** → Identifies outliers impacting performance.
 - **90th percentile:** 90% of responses are faster than this value.
 - **95th percentile:** Indicates delays affecting a small percentage of users.
 - **99th percentile:** Helps track extreme worst-case response times.
- **Maximum Response Time** → The slowest response time recorded.

◆ Insights:

- ✓ High response times indicate slow API/database calls, thread contention, or GC issues.
 - ✓ A large gap between the **median** and **95th percentile** suggests inconsistent performance.
-

2. Throughput Metrics (System Capacity)

Definition: Number of requests processed per second.

- **Transactions per Second (TPS)** → Measures system workload handling capacity.

$$\text{TPS} = \frac{\text{Total Transactions}}{\text{Total Time (sec)}}$$

- **Requests per Second (RPS)** → Measures HTTP request handling capacity.
- **Data Throughput (MB/s)** → Measures data processed per second.

◆ Insights:

- ✓ If **TPS plateaus** while response time increases, a bottleneck exists.
 - ✓ **Comparing TPS and CPU Utilization** helps determine system efficiency.
-

3. Error Rate Metrics (System Stability)

Definition: The percentage of failed transactions.

- **Error Rate (%)**

$$\text{Error Rate} = \left(\frac{\text{Failed Requests}}{\text{Total Requests}} \right) \times 100$$

- **HTTP Status Code Analysis:**

- **4xx Errors:** Client-side issues (e.g., bad requests, authentication failures).
- **5xx Errors:** Server-side failures (e.g., database timeouts, crashes).

- **Timeout Errors** → High timeout rates indicate slow response issues.

- ◆ **Insights:**

✓ High **5xx errors** → Check server logs, DB health, and exception traces.

✓ High **4xx errors** → Look into user request handling logic.

4. CPU & Memory Utilization Metrics

Definition: Measures system resource usage.

- **CPU Utilization (%)**

$$\text{CPU Usage} = \left(\frac{\text{CPU Time}}{\text{Total Available CPU Time}} \right) \times 100$$

- High CPU (above 85%) can cause latency spikes.
- Low CPU but high response times → Possible **DB or disk bottlenecks**.

- **Memory Utilization (%)**

$$\text{Memory Usage} = \left(\frac{\text{Used Memory}}{\text{Total Available Memory}} \right) \times 100$$

- High memory usage (>85%) may cause **OutOfMemory (OOM) errors**.
- Continuous memory increase → Possible **memory leaks**.

- ◆ **Insights:**

✓ If **CPU spikes with low TPS**, optimize **thread pools**, **GC tuning**, or **DB connections**.

✓ If **high memory usage**, analyze **heap dumps for memory leaks**.

5. Garbage Collection (GC) Metrics

Definition: Measures JVM memory management efficiency.

- **GC Pause Time** → Time spent in garbage collection.
- **Frequency of Full GC** → Frequent full GCs can cause high latency.
- **Heap Usage (%)**

$$\text{Heap Usage} = \left(\frac{\text{Used Heap Memory}}{\text{Total Heap Memory}} \right) \times 100$$

◆ **Insights:**

✓ **High GC time** → Optimize heap size, GC tuning, or reduce object creation.

✓ **Frequent full GC events** → Leads to application stalls and slow response times.

6. Disk I/O Metrics

Definition: Measures how efficiently the system handles disk operations.

- **Disk Read/Write Latency (ms)**
 - Slow disk I/O can cause database and logging delays.
- **I/O Wait Time (%)**
 - High **I/O wait** → Indicates slow storage performance.

◆ **Insights:**

✓ High **disk I/O** → Optimize logging, use SSDs, or cache frequently accessed data.

✓ High **I/O wait** → May indicate **disk contention issues**.

7. Network Performance Metrics

Definition: Measures application network efficiency.

- **Network Latency (ms)**
 - Measures round-trip time for a request.
- **Packet Loss (%)**
 - High loss can cause degraded application performance.
- **Bandwidth Utilization (%)**
 - High usage can lead to network congestion.

◆ **Insights:**

- ✓ High **latency** → Optimize **CDN usage, load balancing, or API gateway response times.**
 - ✓ High **packet loss** → Check for **network congestion or connectivity issues.**
-

8. Scalability Metrics

Definition: Determines system ability to handle increased load.

- **Load vs. Response Time**
 - Measures how response time scales with user load.
- **Load vs. Throughput**
 - Identifies the maximum TPS the system can handle.
- **Autoscaling Efficiency**
 - Ensures cloud resources scale efficiently under load.

◆ **Insights:**

- ✓ If **response time increases with load**, system **scalability issues** exist.
 - ✓ If **autoscaling is slow**, consider **tuning scale-up thresholds.**
-

9. System Availability & Reliability Metrics

Definition: Measures system uptime and resilience.

- **Uptime (%)**

$$\text{Uptime} = \left(\frac{\text{Available Time}}{\text{Total Time}} \right) \times 100$$

- 99.9% uptime = ~8.76 hours downtime/year.
- **Mean Time to Recovery (MTTR)**
 - Measures how quickly the system recovers from failures.
- **Mean Time Between Failures (MTBF)**
 - Measures system reliability.

◆ **Insights:**

- ✓ High MTTR → Improve **incident response automation and monitoring.**
 - ✓ Low MTBF → Improve **fault tolerance and failover mechanisms.**
-

10. Database Performance Metrics

Definition: Tracks database query performance and efficiency.

- **Query Execution Time (ms)**
 - Measures time taken to execute SQL queries.
- **DB Connection Pool Utilization (%)**
 - High usage can cause delays.
- **Lock Wait Time (ms)**
 - Long wait times indicate contention issues.

♦ **Insights:**

✓ High **query times** → Optimize SQL queries, add **indexes**.

✓ High **connection pool usage** → Tune **DB connection settings**.

11. Request Queue Length (Backlog)

Definition: The number of requests waiting to be processed by the application or server.

Formula:

Queue Length = Total Incoming Requests - Processed Requests

♦ **Insights:**

✓ A **high queue length** indicates that the system is struggling to keep up with incoming requests.

✓ May suggest the need for **autoscaling, increased thread pools, or optimized query performance**.

12. Thread Pool Utilization

Definition: Measures how effectively an application manages worker threads.

Formula:

$$\text{Thread Pool Utilization (\%)} = \left(\frac{\text{Active Threads}}{\text{Max Threads}} \right) \times 100$$

♦ **Insights:**

✓ **High utilization (>80%)** → Increase thread pool size or optimize request handling.

✓ **Low utilization (<30%)** → Over-provisioning of resources, adjust thread allocation.

13. Connection Pool Utilization (Database or HTTP)

Definition: Measures the percentage of available connections in a pool being used.

Formula:

$$\text{Connection Pool Utilization (\%)} = \left(\frac{\text{Active Connections}}{\text{Max Connections}} \right) \times 100$$

◆ **Insights:**

✓ **Near 100% utilization** → Increase connection pool size or optimize connection reuse.

✓ **Low utilization** → Over-allocation of connections, leading to wasted resources.

14. Cache Hit vs. Miss Ratio

Definition: Tracks how often requested data is found in the cache versus being fetched from the backend.

Formula:

$$\text{Cache Hit Ratio} = \frac{\text{Cache Hits}}{\text{Total Cache Requests}}$$

$$\text{Cache Miss Ratio} = 1 - \text{Cache Hit Ratio}$$

◆ **Insights:**

✓ **High cache misses (>30%)** → Optimize caching strategy (e.g., increase cache size, use LRU/LFU eviction policies).

✓ **High cache hits (>90%)** → Effective caching, reducing backend load.

15. User Experience Score (Apdex Score)

Definition: Measures user satisfaction based on response time thresholds.

Formula:

$$\text{Apdex Score} = \frac{\text{Satisfied Requests} + \left(\frac{\text{Tolerating Requests}}{2} \right)}{\text{Total Requests}}$$

- **Satisfied Requests** → Requests with response times below **T (Threshold Time)**.
 - **Tolerating Requests** → Requests with response times **between T and 4T**.
 - **Frustrated Requests** → Requests **above 4T**.
-

◆ **Insights:**

✓ **Apdex Score near 1.0** → Excellent user experience.

✓ **Apdex Score below 0.5** → Performance issues affecting users significantly.

16. Time to First Byte (TTFB)

Definition: The time taken for a server to send the first byte of data in response to a request.

Formula:

TTFB=Time to Send Request+Server Processing Time+Network Latency

◆ **Insights:**

✓ **High TTFB (>500ms)** → Indicates slow backend processing, network latency, or database delays.

✓ **Optimize by** improving API/database performance, reducing network hops, and enabling HTTP/2.

17. Disk Swap Usage

Definition: Measures the amount of memory swapped in/out to disk when physical RAM is exhausted.

Formula:

$$\text{Swap Usage (\%)} = \left(\frac{\text{Swap Memory Used}}{\text{Total Swap Memory}} \right) \times 100$$

◆ **Insights:**

✓ **High swap usage (>30%)** → Can cause significant slowdowns due to disk-based memory access.

✓ **Optimize by** increasing RAM, tuning JVM heap sizes, or reducing memory-intensive processes.

18. Request Latency Distribution

Definition: Measures response time variation across different requests.

Formula:

Latency Percentiles (50th, 90th, 95th, 99th)

◆ **Insights:**

✓ **A high 99th percentile response time** indicates sporadic delays that impact user experience.

✓ **Optimize by** reducing API delays, optimizing database queries, or increasing worker threads.

19. Thread Contention (Lock Wait Time)

Definition: Measures time spent waiting for a resource lock in multi-threaded applications.

Formula:

Lock Wait Time=Total Time Threads Spend Waiting for Locks

◆ **Insights:**

✓ **High lock wait times (>500ms per request)** → Indicates synchronization issues or DB locking.

✓ **Optimize by** reducing contention, using **optimistic locking**, or increasing parallelism.

20. Concurrent Users

Definition: The number of users actively interacting with the application at a given time.

◆ **Insights:**

✓ Helps determine system limits and bottlenecks under load.

✓ If **response time degrades as users increase**, evaluate **scalability** strategies (autoscaling, DB optimizations).

Summary of Key Performance Metrics

Metric Category	Key Metrics
Response Time	Avg, 95th/99th percentile, Max Response Time
Throughput	TPS, RPS, Data Throughput (MB/s)
Error Rate	% Errors, HTTP 4xx/5xx, Timeout Errors
Resource Usage	CPU %, Memory %, GC Time
Disk & I/O	Read/Write Latency, I/O Wait
Network	Latency, Bandwidth Usage, Packet Loss
Scalability	Load vs Response Time, Autoscaling Efficiency
Availability	Uptime %, MTTR, MTBF
Database	Query Execution Time, DB Connection Pool, Lock Waits
Request Queue Length	Detects backlogs in request processing
Thread Pool Utilization	Measures efficiency of thread usage

Connection Pool Utilization	Tracks database/API connection limits
Cache Hit vs. Miss Ratio	Measures cache efficiency
Time to First Byte (TTFB)	Measures server response speed
Disk Swap Usage	Detects memory exhaustion issues
Latency Distribution	Identifies slow response outliers
Thread Contention (Lock Wait Time)	Measures multi-threading bottlenecks
Concurrent Users	Tracks system scalability under load