

Domain 5. Optimizing Service Performance @

5.1. Introduction to Service Performance Optimization

- Goal: Improve speed, efficiency, and reliability of systems for better user experiences while minimizing cost and resource use.
- Key Focus Areas: Reducing latency, enhancing throughput, lowering error rates, ensuring scalability and efficiency.
- Definition: Techniques to improve responsiveness, speed, efficiency, and reliability of cloud applications by identifying bottlenecks and enhancing resource usage.
- Key Goals:
 - Maximize User Experience with faster response and higher availability.
 - Reduce Operational Costs by optimizing resource usage.
 - Enhance Scalability and Reliability to handle increased loads seamlessly.
- Importance:
 - User Expectations: Instant access expected; slow apps cause dissatisfaction.
 - Business Impact: 1-second delay can reduce conversions by 7%.
 - Resource Efficiency: Optimize cloud resource use to maximize ROI.

5.2. Identifying Performance Bottlenecks

- Definition: System points limiting speed and efficiency, e.g., slow queries, network latency, inefficient algorithms.
- Characteristics:
 - Performance degradation with slow response or reduced throughput.
 - Increased resource consumption causing inefficiency.
 - User frustration due to timeouts or errors.
- Common Bottleneck Types:
 - CPU: High utilization, slow responses.
 - Memory: Excessive paging, crashes.
 - I/O: Slow data retrieval, low throughput.
 - Network: High latency, dropped packets.
- Diagnosis Tools:
 - Cloud Profiler for CPU/memory hot paths.
 - Cloud Trace for latency across microservices.
 - Cloud Debugger for live code debugging.
 - Cloud Monitoring for system metrics.

5.3. Tuning Infrastructure

- Instance Sizing: Choose machine types (E2, N2, C2, M2), use autoscaling (Cloud Run, GKE, Managed Instance Groups).
- Storage Optimization: Select Standard vs. SSD Persistent Disks, Filestore for NFS, Cloud Storage with lifecycle rules.
- Network Optimization: Use Global Load Balancers, Cloud CDN for static content, Private Google Access for internal traffic.
- Goals:
 - Performance Optimization to handle load without degradation.
 - Cost Efficiency by optimizing resource use.
 - Scalability and Flexibility for demand-based scaling.

5.4. Application Optimization Techniques

- Techniques:
 - Code Profiling with Cloud Profiler.
 - Lazy Loading resources when needed.
 - Asynchronous Processing using Pub/Sub, Cloud Tasks.
 - Caching with Cloud Memorystore or in-memory cache.
 - Compression using gzip or Brotli.
 - Connection Pooling for efficient DB connections.
- Key Techniques:
 - Code Optimization: Profile and optimize algorithms.
 - Caching Strategies: In-memory and HTTP caching.
 - Database Optimization: Indexing, connection pooling.
 - Asynchronous Processing: Background jobs, microservices.
 - Frontend Optimization: Minification, bundling, responsive design.

5.5. Load Testing and Benchmarking

- Load Testing:
 - Simulate user traffic to evaluate behavior under load.
 - Tools: JMeter, Locust, k6, custom scripts with Cloud Monitoring.
 - Metrics: Request latency (P50, P90, P95, P99), throughput, error rate, resource usage.
 - Goals: Assess limits, identify bottlenecks, validate scaling.
- Benchmarking:
 - Measure performance against standards.
 - Goals: Set baselines, guide optimization, facilitate comparisons.
 - Tip: Run tests on staging environments similar to production.

5.6. Scaling for Performance

- Vertical Scaling: Increase CPU/RAM of instances.
- Horizontal Scaling: Add instances to handle load.
- GCP Services:
 - Cloud Run auto-scales per request.
 - GKE uses Horizontal Pod Autoscaler.
 - Compute Engine uses Managed Instance Groups with autoscaling.
 - Autoscaler Triggers: CPU utilization, request latency, custom metrics.
- Importance:
 - Maintain user experience with consistent performance.
 - Optimize cost by scaling based on demand.
 - Ensure business continuity during traffic spikes.

Exam Prep Checklist

- Identify performance bottlenecks using GCP tools.
- Understand scaling strategies for Compute Engine, Cloud Run, GKE.
- Use autoscaling based on load metrics.
- Apply best practices in caching, networking, and storage.
- Perform load testing and improve performance using results.
- Use profiler, trace, and debugger for optimization.
- Optimize cost without sacrificing performance.

Summary Cheat Sheet

- Code Optimization: Cloud Profiler, Debugger
- Latency Reduction: Trace, CDN, in-memory caching
- Scaling: Cloud Run, GKE, Managed Instance Groups
- Load Testing: k6, JMeter, Locust
- Database Tuning: Indexes, caching, replicas
- Network Optimization: Cloud Load Balancing, Cloud CDN
- Performance Monitoring: Cloud Monitoring, Alerting
- Autoscaling Config: CPU, latency, custom metrics
- Cost Optimization: Right-sizing, Preemptible VMs, Committed Discounts

5.10. Continuous Performance Monitoring

- Use Cloud Monitoring dashboards for latency, error rates, saturation, CPU/memory usage.
- Set SLO-based alerts to act before performance degrades.
- Key Metrics:
 - Application: Response time, throughput, error rates.
 - Infrastructure: CPU/memory usage, network latency.
 - User Experience: Real User Monitoring (RUM), Synthetic Monitoring.
- Tools and Techniques:
 - Google Cloud Monitoring for unified insights and alerts.
 - Structured logging via Google Cloud Logging.
 - Integrate monitoring into CI/CD pipelines for early detection of issues.

5.9. Optimizing Cost While Improving Performance

- Strategies:
 - Autoscaling: Scale down during low traffic.
 - Committed Use Discounts: Prepay for predictable workloads.
 - Instance Rightsizing: Match VM size to workload.
 - Preemptible VMs: Use for fault-tolerant batch jobs.
 - Efficient Data Storage: Lifecycle rules to move data to cheaper tiers.
- Additional Tips:
 - Choose the right instance types based on workload.
 - Use preemptible VMs for cost savings.
 - Optimize storage class selection.
 - Leverage serverless architectures for cost-effective scaling.
 - Monitor performance to identify inefficiencies and adjust resources.

5.8. Using Caching Effectively

- Definition: Storing copies of data temporarily to speed up future requests.
- Types:
 - Edge Caching: Cloud CDN.
 - In-memory Caching: Memorystore (Redis/Memcached).
 - Application-level: HTTP caching headers, lazy loading.
- Benefits:
 - Reduce backend load.
 - Improve response times.
 - Save costs.
- Caching Levels:
 - Client-Side: Browser storage.
 - Server-Side: Application or dedicated infrastructure.
 - CDN Caching: Static assets close to users.
- Key Caching Types:
 - In-Memory: For session data, frequent queries.
 - HTTP Caching: Controlled by cache headers.
 - Database Caching: Cache query results or materialized views.

5.7. Database Optimization

- Best Practices:
 - Query Performance: Use indexes, avoid N+1 queries.
 - Connection Management: Use Cloud SQL Auth proxy, connection pools.
 - Caching: Cache frequent reads with Memorystore.
 - Scaling: Use read replicas, sharding, partitioning.
 - Storage Tiers: Choose pricing/performance tiers appropriately.
- Techniques:
 - Indexing: Choose suitable index types for queries.
 - Query Optimization: Analyze execution plans, avoid SELECT *.
 - Schema Design: Balance normalization and denormalization; choose proper data types.
 - Connection Pooling: Reuse connections to reduce overhead.
 - For high-throughput, low-latency needs consider Spanner, Bigtable, or Firestore.