# Troubleshooting Scenarios for Production Engineering Interviews

## Scenario 1: "A service you support is no longer responding to requests"

### 1. Understand the Scope of the Problem

**Ask Clarifying Questions:**

- Is the issue affecting all users or only a subset?
- Which endpoints or services are impacted?
- When was the issue first noticed, and what triggered the alert?
- Are there any recent changes to the system (e.g., deployments, config updates, infrastructure changes)?

### 2. Check for Obvious Causes

**Verify Monitoring and Alerting Systems for:**

- **Service Health Metrics**: CPU, memory, disk usage, network bandwidth.
- **Error Rates**: HTTP 500s, 502s, or connection timeouts.
- **Latency**: Are requests delayed or timing out?

### 3. Narrow Down the Issue

**Application Layer:**

- Inspect application logs for errors or warnings.
- Check recent code deployments or configuration changes. Roll back if necessary.

**Network Layer:**

- Verify that the service is reachable using tools like `ping`, `curl`, or `telnet`.
- Check for DNS issues or incorrect IP mappings.

**Infrastructure Layer:**

- Validate that all dependent services (databases, APIs, caches) are healthy.
- Confirm the service is running on all intended instances.
- Look for any failed health checks in load balancers.

### 4. Deep Dive into Specific Areas

**Service Logs:**

- Are there stack traces or specific errors? Example: "Database connection refused" or "Out of memory."

**Dependencies:**

- Verify upstream and downstream dependencies.
- For example, if the service relies on a database, ensure that the database is reachable and performant.

**Configuration:**

- Confirm that environment variables, configs, and secrets are correct.

**Resource Limits:**

- Check if the service is hitting resource limits (e.g., container memory limits or thread pool exhaustion).

## 5. Quick Fixes for Immediate Recovery

- **Restart the Service**: Temporarily resolves issues caused by transient states.
- **Scale the Service**: Add more instances if the issue is related to traffic overload.
- **Redirect Traffic**: Route traffic to a healthy region or instance using a load balancer.

## 6. Long-Term Mitigation

**After the Service is Restored:**

- Conduct a **Postmortem** to identify root causes.
- Implement safeguards to prevent recurrence:
    - Improve monitoring to detect the issue earlier.
    - Add automated rollback mechanisms for failed deployments.
    - Optimize resource utilization to avoid limits.

---

## Key Points for the Interview

- Show your systematic approach to troubleshooting.
- Emphasize the importance of communication with stakeholders during outages.
- Highlight your ability to balance short-term fixes with long-term solutions.

---

# Scenario 2: "Your metrics indicate that there are sporadic issues that are impacting user experience"

This scenario focuses on identifying intermittent issues, which can be tricky. Here's how to approach it systematically:

## 1. Understand the Problem

Ask clarifying questions:

- What specific user experiences are impacted (e.g., latency, errors, incomplete functionality)?
- How frequent are the issues, and are there identifiable patterns (e.g., time of day, specific regions, devices, or user actions)?
- When did the problem first occur?
- Are there error reports or user feedback available?

## 2. Examine Metrics and Logs

**Metrics:**

- Look at key indicators such as request latency, error rates, CPU/memory usage, and traffic patterns.
- Identify anomalies: Spikes in 500 errors? Latency outliers? Traffic drops?

**Logs:**

- Look for timestamps or errors during affected times.
- Correlate logs across services for patterns.

## 3. Isolate Patterns

**User-Specific Patterns:**

- Are only certain users affected? (e.g., specific geographic locations, ISPs, devices, or versions of the app).

**Request-Specific Patterns:**

- Does the issue occur with specific endpoints or functionalities (e.g., searches, uploads)?

**Time-Specific Patterns:**

- Are issues tied to high traffic periods or specific scheduled tasks (e.g., batch jobs or data migrations)?

## 4. Narrow Down Possible Causes

**Application Issues:**

- Analyze application behavior:
    - Are there any buggy code paths that only execute under rare conditions?
    - Are retries or circuit breakers misconfigured, causing cascading failures?

**Network Issues:**

- Check for packet loss or increased latency.
- Inspect CDN or load balancer performance, especially for regions or routes.

**Dependency Issues:**

- Are backend services or third-party APIs intermittently failing?
- Check service health across regions or nodes.

**Infrastructure Issues:**

- Inspect resource usage (e.g., burst CPU or memory exhaustion).
- Check for failing containers, spot instances, or faulty hardware.

## 5. Reproduce and Validate

- Attempt to reproduce the issue in a controlled environment:
    - Use test accounts in different regions.
    - Simulate traffic patterns or rare scenarios (e.g., edge cases).
    - Use distributed tracing tools to track specific request flows across services.

## 6. Quick Mitigation Steps

**Redirect Traffic:**

- Shift traffic to healthier regions or nodes if issues are geographically isolated.

**Rate Limiting:**

- Apply rate limits or backpressure to prevent overload on downstream services.

**Temporary Fixes:**

- Restart impacted services or scale up resources temporarily.

## 7. Post-Incident Actions

**Root Cause Analysis:**

- Identify the exact trigger (e.g., resource contention, network issues, or a hidden application bug).

**Monitoring Improvements:**

- Add fine-grained metrics to detect sporadic patterns sooner.

**Resilience Improvements:**

- Optimize retry logic, circuit breakers, and timeout configurations.

**Testing Enhancements:**

- Expand test coverage for edge cases or traffic spikes.

---

## Key Points for the Interview

- Demonstrate your ability to analyze patterns and isolate the root cause.
- Highlight the importance of end-to-end visibility through metrics and tracing.
- Emphasize proactive actions, like improving monitoring and hardening systems to avoid recurrence.

---

# Scenario 3: "An application terminates unexpectedly after a period of time for unknown reasons"

This scenario focuses on diagnosing and resolving issues that cause an application to crash after running for a certain period. Here's a structured approach:

## 1. Understand the Symptoms

Ask clarifying questions:

- What are the application logs showing before termination? Any specific error messages or warnings?
- Does the termination happen consistently after a specific period, or is it random?
- Are there recent changes (code, dependencies, configurations, or environment)?
- Is the issue occurring in production, staging, or all environments?

## 2. Gather Evidence

**Application Logs:**

- Check for errors, warnings, or stack traces near the termination time.
- Look for patterns like resource exhaustion, segmentation faults, or exceptions.

**System Logs:**

- Check the host or container logs for events like `OOM Killer` (Out of Memory), kernel panics, or hardware failures.

**Monitoring Metrics:**

- Observe resource usage over time (e.g., memory, CPU, disk I/O, file descriptors).
- Identify gradual increases or spikes that align with the termination.

## 3. Analyze Potential Causes

**Memory Leaks:**

- Does memory usage grow over time without release? Tools: `valgrind`, `heap profiler`.

**Resource Limits:**

- Check container resource limits (CPU/memory) or system-level quotas (e.g., open file descriptors, threads).

**Uncaught Exceptions:**

- Inspect error handling in the application. Are exceptions bubbling up and terminating the process?

**External Dependencies:**

- Does the app rely on external services (e.g., databases, APIs) that might cause timeouts, retries, or cascading failures?

**Runtime Issues:**

- Is the app running into timeout configurations, thread pool exhaustion, or stale connections?

## 4. Reproduce the Issue

Replicate the workload or environment in a test system:

- Run the application under similar conditions (e.g., traffic, configuration, resources).
- Use stress-testing tools (e.g., `Apache JMeter`, `locust`) to simulate production-like behavior.

Enable debugging tools:

- Attach debuggers or profilers (e.g., `gdb`, `strace`, `perf`) to analyze runtime behavior.

## 5. Apply Immediate Fixes

**Restart the Service:**

- Temporary mitigation to restore availability while debugging.

**Scale Resources:**

- Increase resource limits (e.g., memory or CPU) to prevent immediate termination.

**Enable Verbose Logging:**

- Increase log detail to capture more context around the termination.

## 6. Long-Term Mitigations

**Debug and Resolve Root Causes:**

- **Fix Memory Leaks**: Identify and free unused memory.
- **Handle Exceptions**: Improve error handling to prevent unhandled exceptions from crashing the app.
- **Tune Resource Usage**: Optimize thread pools, connection pools, or caching mechanisms.

**Monitoring and Alerts:**

- Add proactive monitoring for memory, CPU, and other critical resources.
- Set up alerts for termination signals or crash-loop behavior.

**Resilience Improvements:**

- Implement retry mechanisms, graceful degradation, and circuit breakers to handle transient issues.
- Consider crash recovery mechanisms (e.g., restart policies, checkpoints).

## Key Points for the Interview

- Show a methodical approach to diagnosing and resolving runtime issues.
- Emphasize the use of tools and data to narrow down potential causes.
- Highlight the importance of balancing short-term fixes with long-term system improvements.

# Scenario 4: "An alert fires that shows a sustained latency regression"

When an alert indicates sustained latency regression, it's critical to quickly diagnose the issue to minimize user impact. Here's a structured approach:

## 1. Understand the Alert

Ask clarifying questions:

- What specific service, endpoint, or functionality is impacted?
- How significant is the regression (e.g., from 100ms to 500ms)?
- When was the latency first detected, and is it still ongoing?
- Are there specific thresholds or SLAs that are being violated?

## 2. Gather Initial Data

**Metrics:**

- Look at latency trends for key operations (e.g., P50, P90, P99 latency).
- Examine correlated metrics like error rates, request rates, or resource usage.

**Logs:**

- Analyze application logs for slow queries, timeouts, or other anomalies.

**Traces:**

- Use distributed tracing tools (e.g., Jaeger, Zipkin) to identify bottlenecks in request flows.

## 3. Narrow Down Potential Causes

**Application Layer:**

- Recent code deployments introducing inefficient algorithms or logic.
- Increased processing time for specific endpoints or user actions.

**Database Layer:**

- Slow queries, lock contention, or high replication lag.
- Increased query volume or unoptimized indexes.

**Infrastructure Layer:**

- Resource contention (e.g., CPU, memory, disk I/O).
- Scaling issues with application instances or backend services.

**Network Layer:**

- Increased latency due to packet loss, DNS resolution issues, or changes in routing.
- Overloaded load balancers or proxies.

## 4. Correlate Patterns

**Traffic Patterns:**

- Is the latency regression tied to increased traffic or specific times of day?

**Geographical Impact:**

- Are certain regions or CDNs more affected than others?

**Specific Requests:**

- Does the regression occur only for specific endpoints, users, or payload sizes?

## 5. Mitigate the Impact

**Temporary Scaling:**

- Scale up application instances or database replicas to handle increased load.

**Redirect Traffic:**

- Shift traffic to healthier regions or instances if possible.

**Bypass Bottlenecks:**

- Disable non-essential features or routes contributing to the regression.

## 6. Identify Root Cause

**Code Profiling:**

- Use profiling tools to identify slow functions or inefficient logic.

**Database Analysis:**

- Examine slow query logs and optimize queries or indexes.
- Add caching layers for frequently accessed data.

**Network Analysis:**

- Trace latency spikes in network communication, DNS lookups, or API calls.

**Resource Usage:**

- Check for resource exhaustion or limits at the infrastructure level.

## 7. Post-Incident Follow-Up

**Monitoring Improvements:**

- Add fine-grained monitoring to detect and alert on latency trends earlier.

**Performance Testing:**

- Conduct load and stress testing to simulate traffic patterns and identify bottlenecks.

**System Optimizations:**

- Optimize resource configurations, database queries, and application logic.

**Deploy Safeguards:**

- Implement circuit breakers, backpressure mechanisms, and rate limiting to prevent similar issues.

## Key Points for the Interview

- Highlight the importance of metrics correlation and end-to-end visibility.
- Emphasize the balance between immediate mitigation and long-term root cause analysis.
- Showcase familiarity with tools like tracing systems, profilers, and query optimizers.

---

# Scenario 5: "User reports suggest your service has an outage, but your team hasn't received any alerts"

This scenario focuses on investigating a potential outage reported by users that hasn't triggered any automated alerts. Here's a structured approach:

## 1. Gather Initial Information

Ask clarifying questions:

- What specific issues are users experiencing (e.g., errors, timeouts, incorrect behavior)?
- How widespread is the issue (e.g., all users, specific regions, or platforms)?
- When did users first notice the issue, and is it ongoing?
- Are there any patterns among the reports (e.g., device types, ISPs, regions)?

## 2. Validate the Reports

**Simulate the User Experience:**

- Test the service manually, especially for the affected functionality or region.
- Use tools like `curl`, browser developer tools, or service-specific APIs.

**Check External Monitoring:**

- Use third-party services (e.g., Pingdom, Uptrends, or UptimeRobot) to verify service availability and latency from multiple locations.

**Review User Reports:**

- Analyze user feedback for common symptoms or error messages.

## 3. Verify System Metrics and Logs

**Metrics:**

- Review key service metrics (e.g., request rates, error rates, latency, resource usage).
- Look for anomalies that might not have breached alert thresholds.

**Logs:**

- Inspect application logs for error spikes or warning messages.
- Look for unusual patterns, such as frequent retries or dropped connections.

**Tracing:**

- Use distributed tracing tools to identify bottlenecks or errors in specific requests.

## 4. Correlate Patterns

**Regional Impact:**

- Check for geographic, ISP, or CDN-specific issues that might explain isolated user reports.

**Recent Changes:**

- Investigate recent deployments, configuration changes, or infrastructure updates.

**Traffic Patterns:**

- Examine whether there's a traffic drop in specific regions or timeframes.

## 5. Check Monitoring and Alerting Gaps

**Alert Thresholds:**

- Verify if alerts are too lenient (e.g., thresholds set too high).

**Metrics Coverage:**

- Check if critical paths or components lack sufficient monitoring.

**Alert Configuration:**

- Ensure alerts are correctly routed and functional (e.g., check alerting integrations with Slack, PagerDuty, etc.).

## 6. Mitigate the Impact

**Redirect Traffic:**

- Shift traffic to healthy regions or services if the issue is isolated.

**Restart Services:**

- Restart impacted services or scale them to handle unexpected loads.

**Communicate:**

- Notify users via status pages or social media to acknowledge the issue and provide updates.

## 7. Root Cause Analysis

**Fix Monitoring Gaps:**

- Add or refine metrics and alerts for the affected components.
- Ensure comprehensive coverage for key service paths.

**Improve Observability:**

- Add distributed tracing, log aggregation, or real-time dashboards for better visibility.

**Resilience Enhancements:**

- Implement redundancy, auto-scaling, and self-healing mechanisms.

## Key Points for the Interview

- Emphasize the importance of validating reports and avoiding assumptions.
- Highlight the need to investigate monitoring gaps to prevent recurrence.
- Demonstrate proactive steps like communicating with users and rapid mitigation while performing root cause analysis.

---

# Scenario 6: "A production host gets shut down from time to time due to out of memory (OOM)"

When a production host shuts down intermittently due to OOM issues, it's critical to identify the root cause of memory exhaustion and implement both immediate and long-term fixes. Here's a structured approach:

## 1. Understand the Issue

Ask clarifying questions:

- How often does the OOM shutdown occur? Is it periodic or random?
- Are there specific workloads, users, or services running at the time of the shutdown?
- What are the consequences of the shutdown (e.g., degraded service, loss of data)?

Gather information:

- When was the issue first observed?
- Were there any recent changes in software, configuration, or traffic patterns?

## 2. Gather Evidence

**System Logs:**

- Check kernel logs (`/var/log/messages`, `/var/log/syslog`, or `dmesg`) for OOM-related messages. Look for details on which processes were terminated and their memory usage.

**Metrics:**

- Review memory usage trends for the host.
- Look at swap usage, memory allocation rates, and cache/buffer utilization.

**Application Logs:**

- Check for memory-related errors or warnings in application logs (e.g., out-of-memory exceptions).

**Traffic Patterns:**

- Analyze traffic data for unusual spikes that might correlate with increased memory usage.

## 3. Analyze Potential Causes

**Memory Leaks:**

- Are applications consuming memory without releasing it? Use tools like `valgrind` or heap profilers.

**Traffic Spikes:**

- Are certain traffic patterns causing excessive memory usage?

**Misconfigured Applications:**

- Check for improper memory limits (e.g., Java heap size, database connection pools, or caching layers).

## Resource Constraints:

- Is the host under-provisioned (insufficient RAM for the workload)?

**Faulty Processes:**

- Are runaway processes consuming excessive memory? Tools: `top`, `htop`, or `ps`.

## 4. Immediate Mitigation

**Kill Rogue Processes:**

- Identify and terminate high-memory processes before they trigger an OOM shutdown.

**Scale Resources:**

- Temporarily increase memory allocation for the host.

**Redistribute Workloads:**

- Offload some services or tasks to other hosts.

**Enable Swap:**

- If swap is not enabled, configure it to provide temporary relief (though this may degrade performance).

## 5. Deep Dive Diagnostics

**Memory Profiling:**

- Use tools like `perf`, `strace`, or application-specific profilers to identify memory-intensive operations.

**Heap Dump Analysis:**

- For applications like Java, analyze heap dumps to pinpoint memory leaks or over-allocation.

**Load Testing:**

- Simulate production workloads to reproduce the issue in a controlled environment.

## 6. Long-Term Solutions

**Fix Memory Leaks:**

- Patch or optimize applications that are not releasing memory correctly.

**Optimize Configurations:**

- Tune parameters like cache sizes, thread pools, or connection pools to use memory efficiently.

**Resource Limits:**

- Use cgroups or container memory limits to prevent runaway processes from exhausting memory.

**Scale Out:**

- Distribute workloads across multiple hosts to reduce memory pressure on a single instance.

**Capacity Planning:**

- Ensure hosts are adequately provisioned for peak workloads.

## 7. Monitoring and Alerts

**Add Memory Monitoring:**

- Track memory usage trends, cache utilization, and swap usage with tools like Prometheus, Datadog, or CloudWatch.

**Set Alerts:**

- Configure alerts for high memory usage thresholds or early OOM warnings.

## Key Points for the Interview

- Highlight your ability to triage and mitigate quickly to restore stability.
- Emphasize the importance of profiling and diagnostics to identify root causes.
- Discuss implementing resilient configurations and proactive monitoring to prevent recurrence.

---

# Scenario 7: "A small subset of users seems to get surprisingly slow responses"

## Steps to Troubleshoot:

1. **Understand the Problem**

   - Ask clarifying questions:
     - What functionality or endpoints are affected?
     - How are the impacted users distributed (e.g., geographically, by ISP, device type, or browser)?
     - What response times are they seeing compared to unaffected users?
     - Is this issue reproducible for the affected subset, and does it happen consistently?

2. **Gather Evidence**

   - **User Reports**:
     - Analyze logs, screenshots, or other data from affected users.
     - Identify patterns such as geographic regions, device types, or specific actions.
   - **Application Logs**:
     - Look for anomalies tied to affected users, such as specific user IDs, IPs, or session tokens.
     - Check for long-running queries or timeouts for their requests.
   - **Metrics**:
     - Review latency metrics by user segment (e.g., by region, device, API endpoint).
     - Compare P50, P90, and P99 latencies for affected users vs. the general population.
   - **Tracing**:
     - Use distributed tracing tools to pinpoint slow operations for requests from affected users.

3. **Narrow Down Potential Causes**

   - **Geographical Factors**:
     - Latency due to distance from data centers or suboptimal CDN edge nodes.
   - **Network Issues**:
     - ISP throttling, packet loss, or DNS resolution delays for the affected subset.
   - **Application Layer**:

- Misconfigurations causing delays for certain user segments (e.g., slow database queries or caching inefficiencies).
  - **Traffic Routing**:
    - Suboptimal load balancer or traffic routing decisions.
  - **Client-Specific Issues**:
    - Outdated app versions, incompatible browsers, or specific device constraints.

4. **Immediate Investigation**

  - **Reproduce the Issue**:
    - Simulate requests from affected regions, devices, or browsers using tools like `curl`, `Postman`, or browser dev tools.
    - Test using VPNs or proxies to mimic the affected users' network conditions.
  - **Check CDN Behavior**:
    - Verify if the CDN is routing traffic to the nearest edge servers for the affected users.
  - **Analyze Load Balancer Logs**:
    - Look for uneven traffic distribution or suboptimal routing to backend servers.
  - **Query Optimization**:
    - Review slow query logs or database performance for requests tied to the affected subset.

5. **Mitigate Impact**

  - **Route Traffic**:
    - Manually redirect affected users to a closer or healthier data center or CDN edge.
  - **Optimize Network Paths**:
    - Adjust DNS configurations or CDN settings to improve routing.
  - **Cache Responses**:
    - Add caching layers for endpoints that generate high latency for affected users.
  - **Provide Workarounds**:
    - Offer an alternative endpoint or simplified functionality for impacted users.

6. **Root Cause Analysis**

  - **Geography and Network**:
    - Diagnose network latency using tools like `mtr`, `ping`, or `traceroute` from affected regions.
  - **CDN Behavior**:
    - Check CDN logs and configurations for edge server performance or cache-hit rates.
  - **Database and Backend**:
    - Identify database queries or backend operations that take longer for specific user conditions.
  - **Client Configurations**:
    - Test client-side settings like app versions or browser compatibility.

7. **Prevent Recurrence**

  - **Improve Observability**:
    - Add segment-based metrics to monitor latency trends by region, device, or user group.
  - **Enhance Routing**:

- Use tools like geolocation-based routing or global traffic management for optimized performance.
  - **Optimize Caching**:
    - Ensure high cache-hit rates for static and frequently requested content.
  - **User Experience Safeguards**:
    - Implement client-side retries, fallback mechanisms, and loading spinners for slow responses.