# Designing Observability for Fault Diagnosis

**Elton Stoneman**

Consultant & Trainer

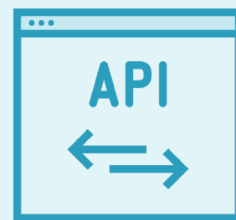@EltonStoneman    blog.sixeyed.com

# Incident Model

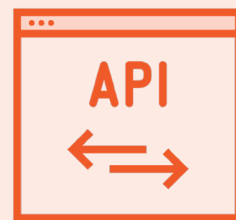**Triage**    **Examine**    Diagnose    Test    Cure
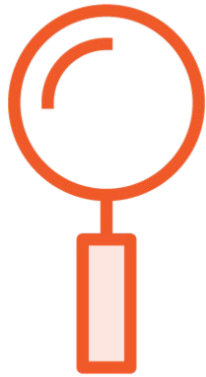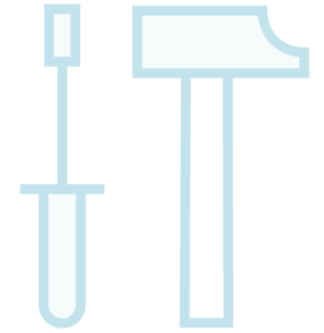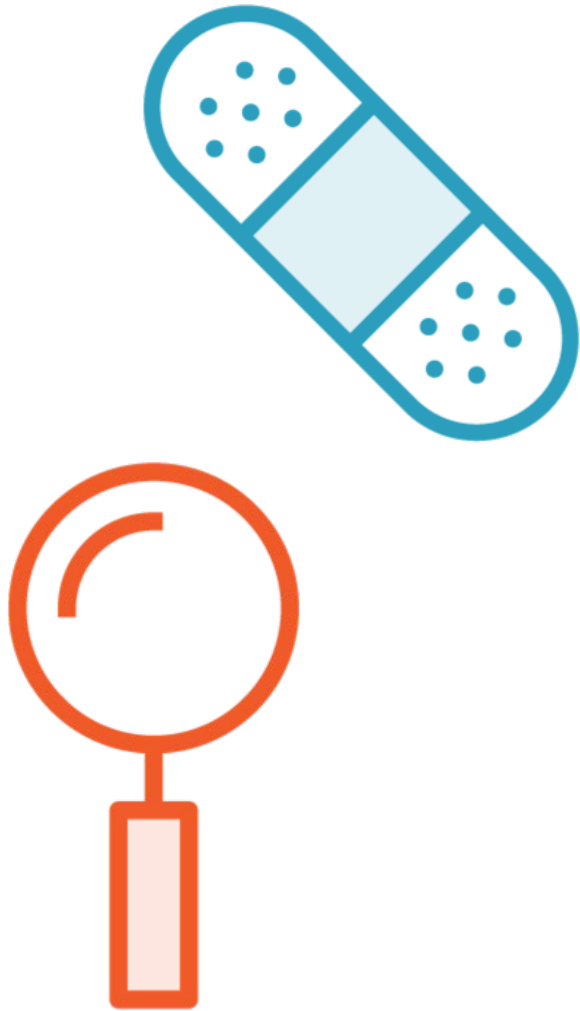
**Triage**
- Get back to "good enough"
- Application health dashboards
- High-level metrics overview

**Examine**
- Understand the problem & trigger
- Application logs & network traces
- Low-level details

# Golden Signals

**Latency**

Job processing time

Response generation time

**Traffic**

Length of message queue
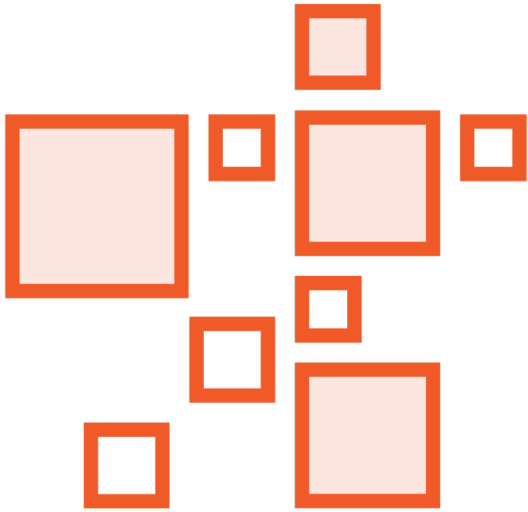
Requests per second

**Errors**

Request failures

Response correctness

**Saturation**
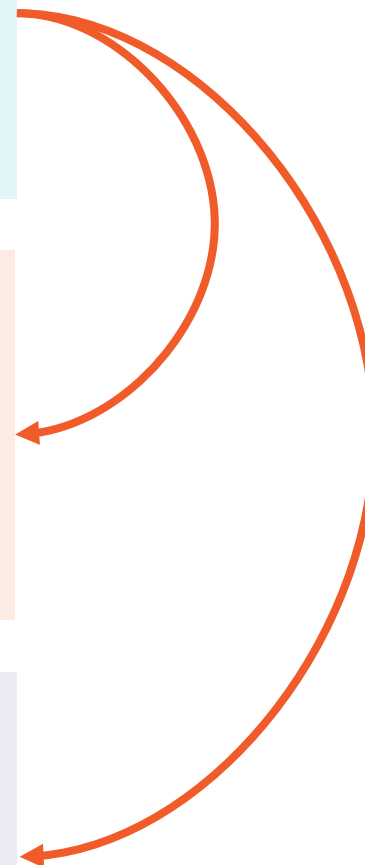
CPU & memory utilization

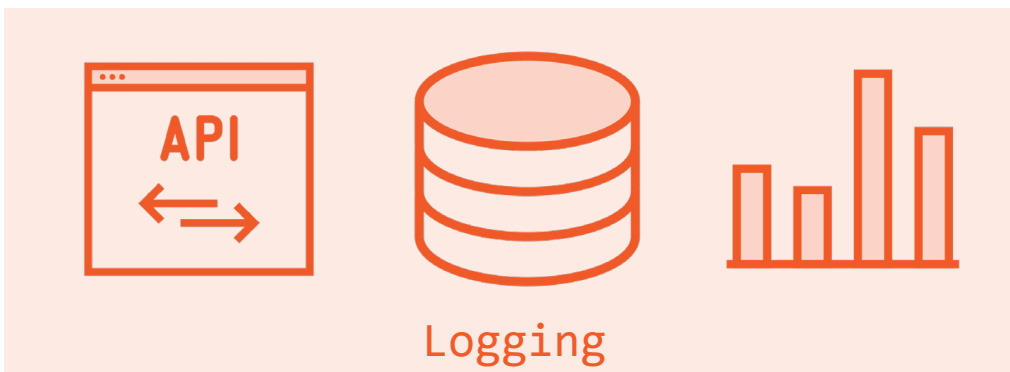Network bandwidth

# Observability Data

**Detail levels**

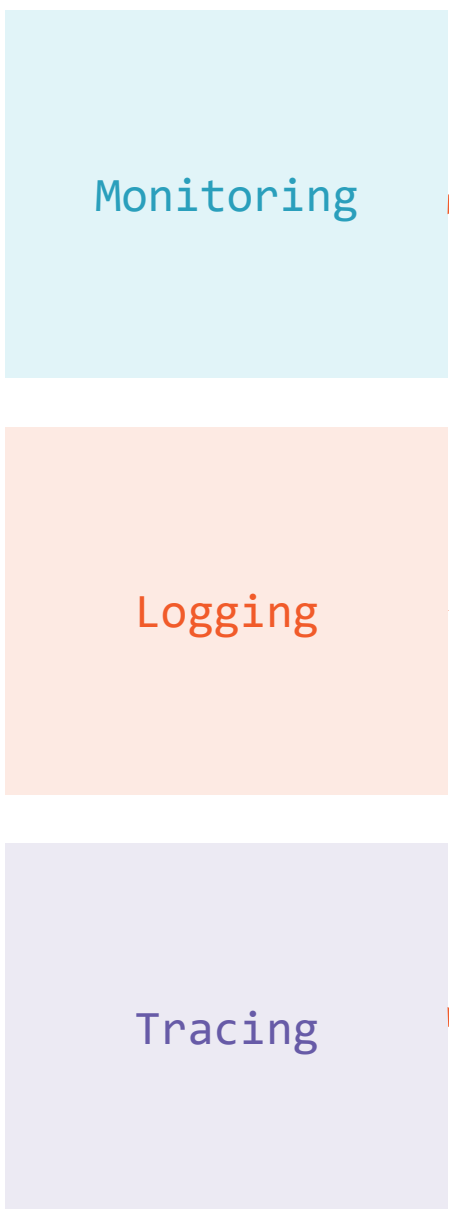**Amounts available**

**Investigative tools**

Monitoring

Logging

Tracing

# Exploring the Three Pillars of Observability
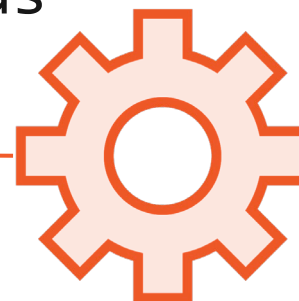
Collection

Storage

Visualization

Monitoring

Logging

Tracing

Monitoring

GET /metrics

- Latency
- Queue

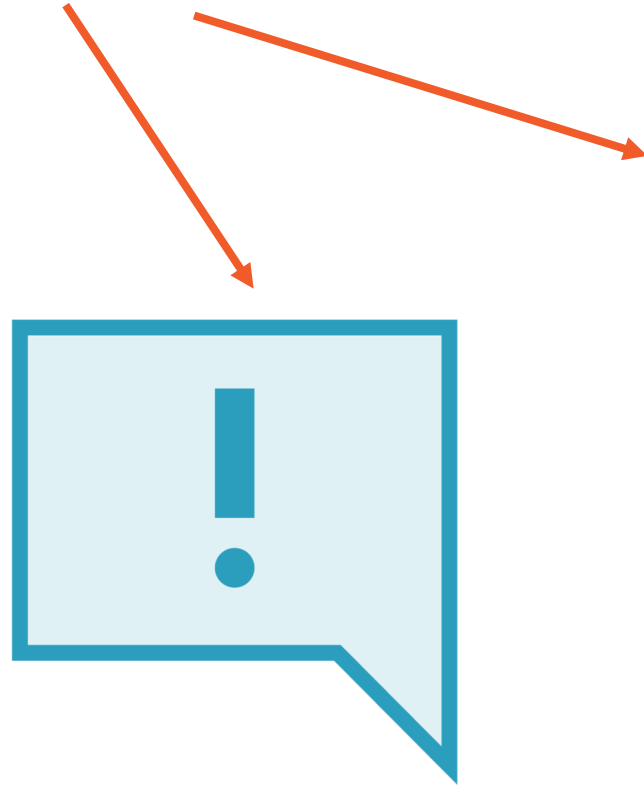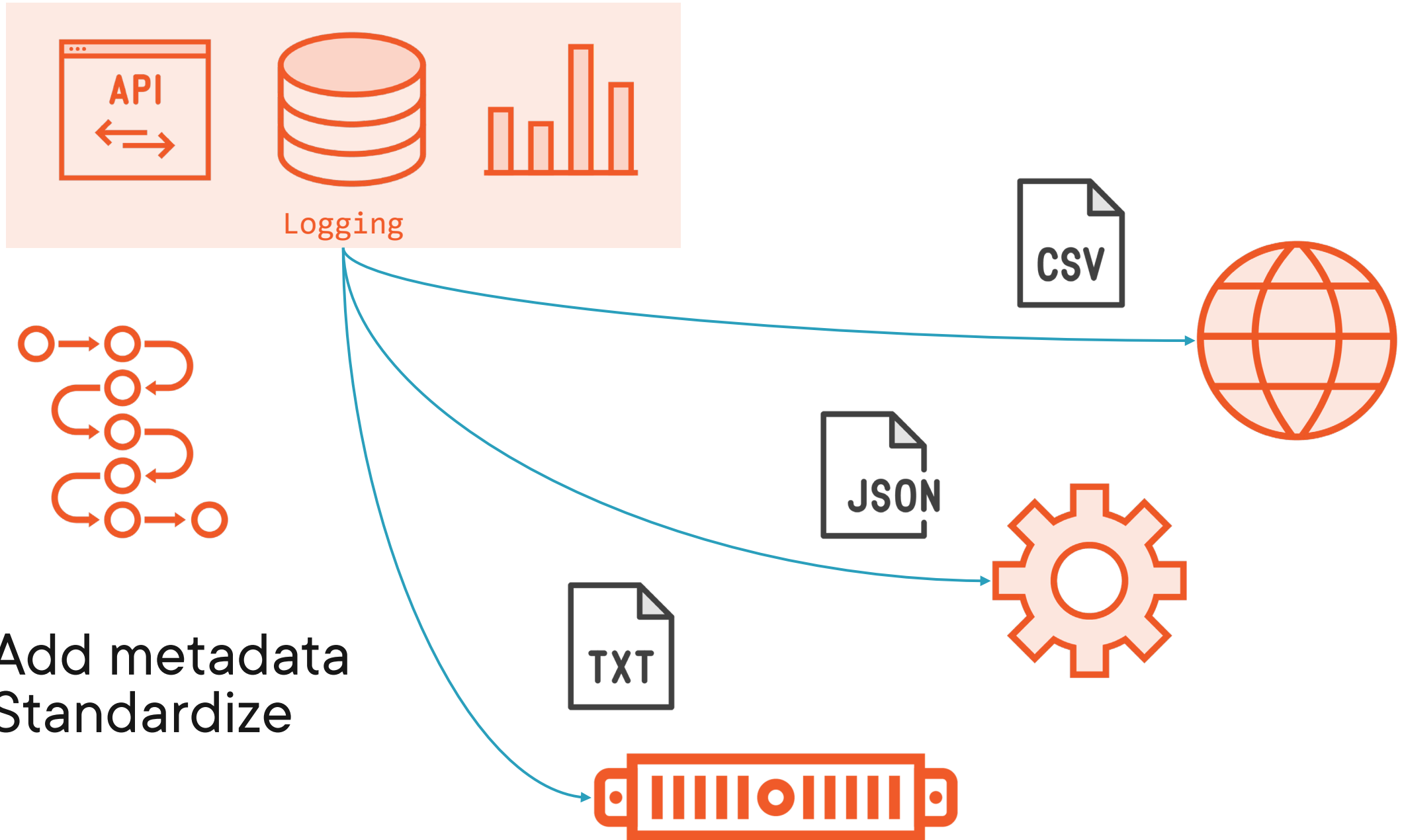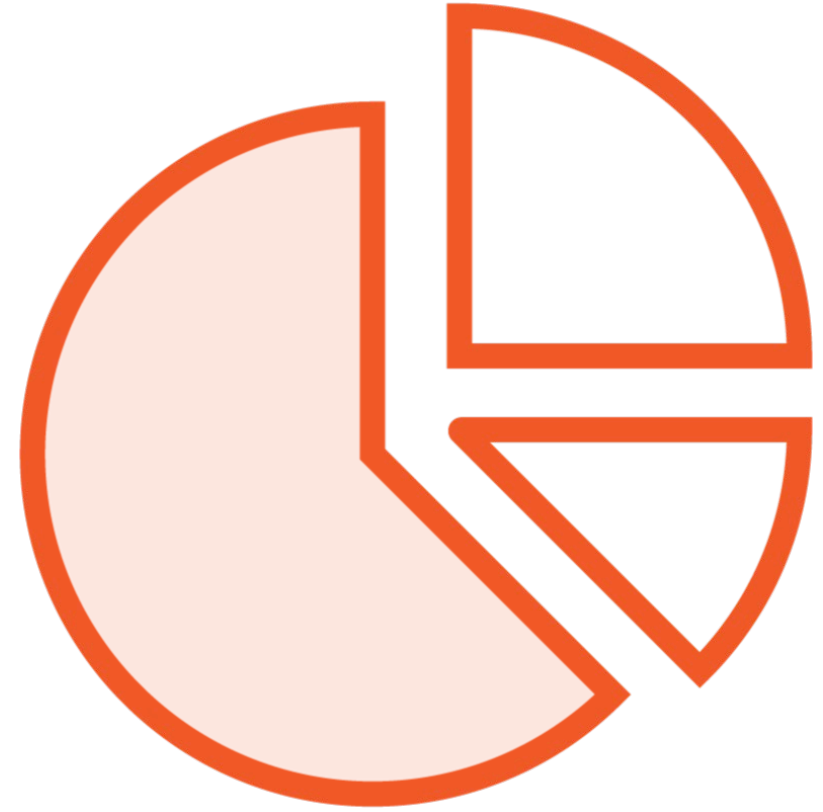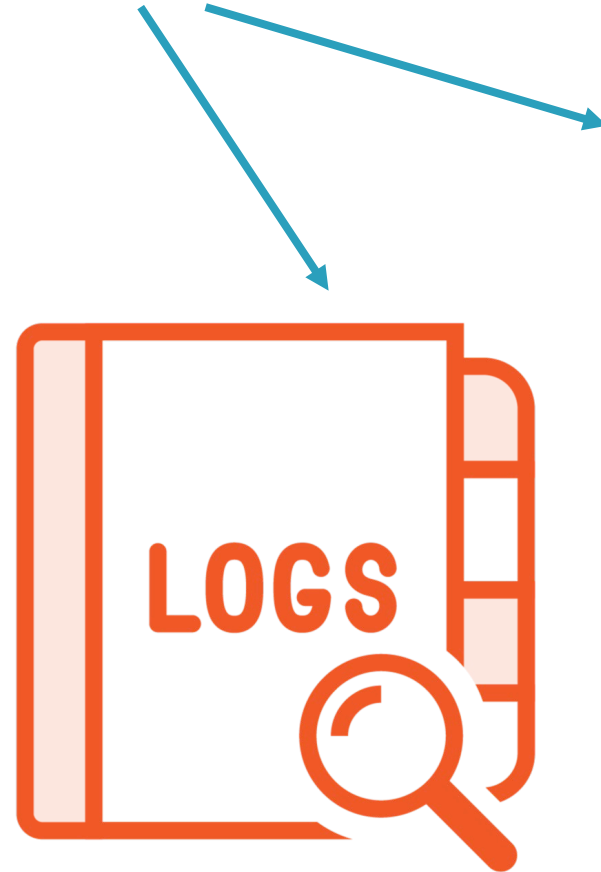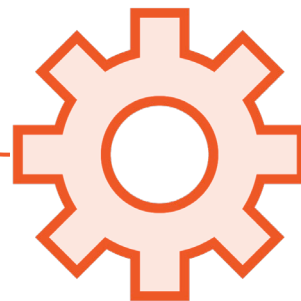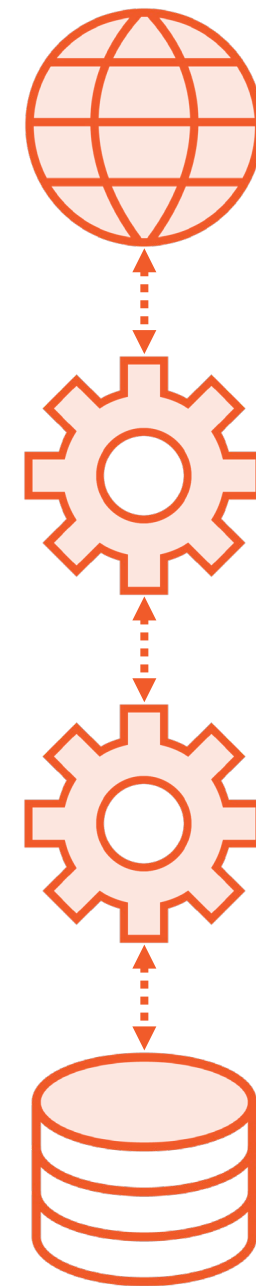- Threads
- GC

- CPU
- Memory

- Add metadata
- Standardize

API

Logging
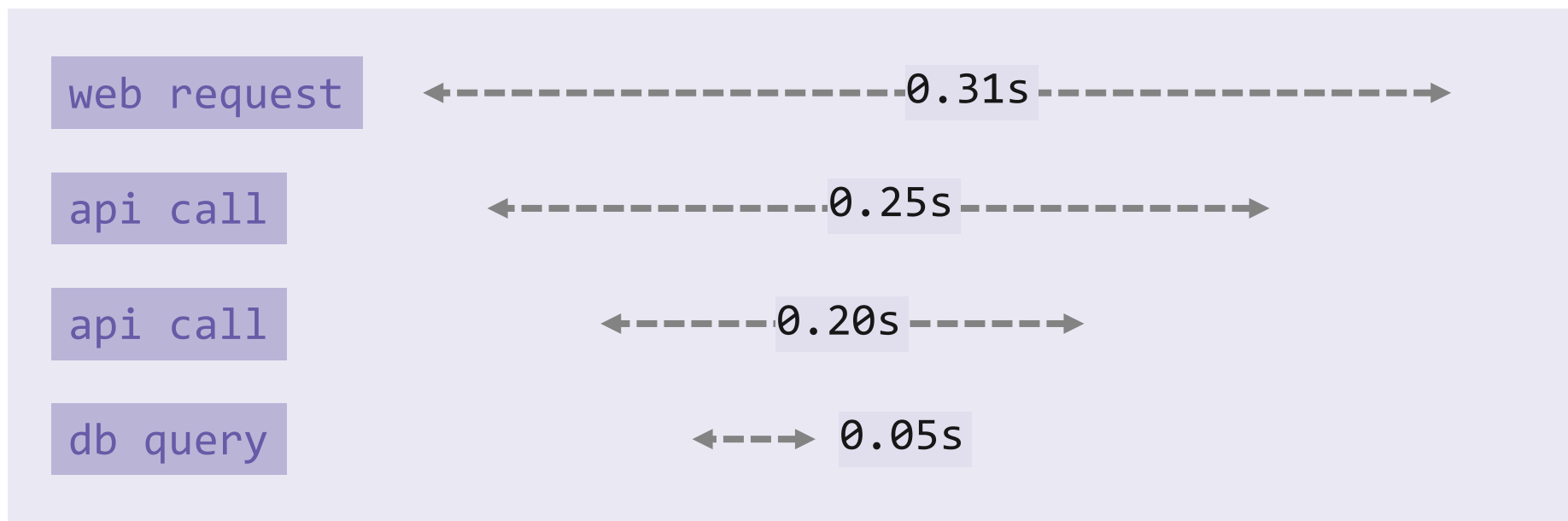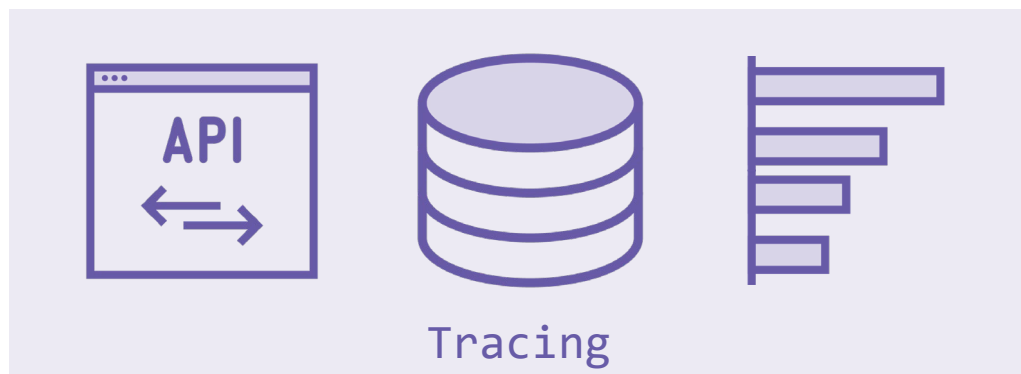
LOGS

**#1**

Monitoring
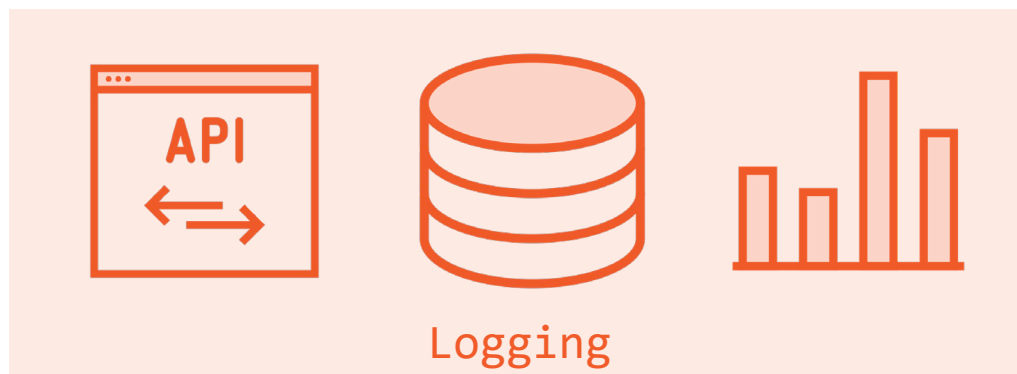
- Overall insight
- Uptime & usage
- App-level metrics

**#2**

Logging

- Low-level detail
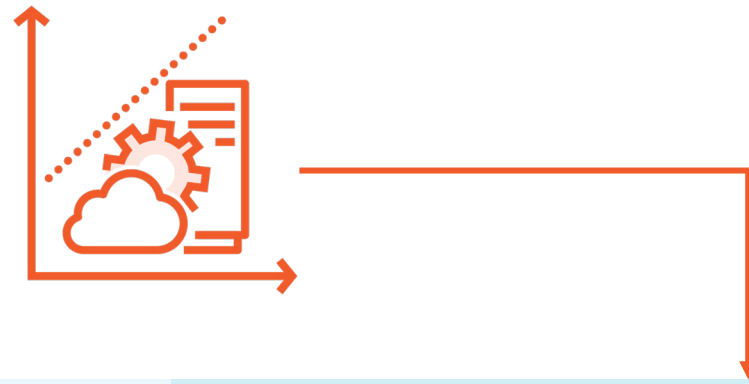- Collect existing logs
- Transform pipeline

**#3**

Tracing

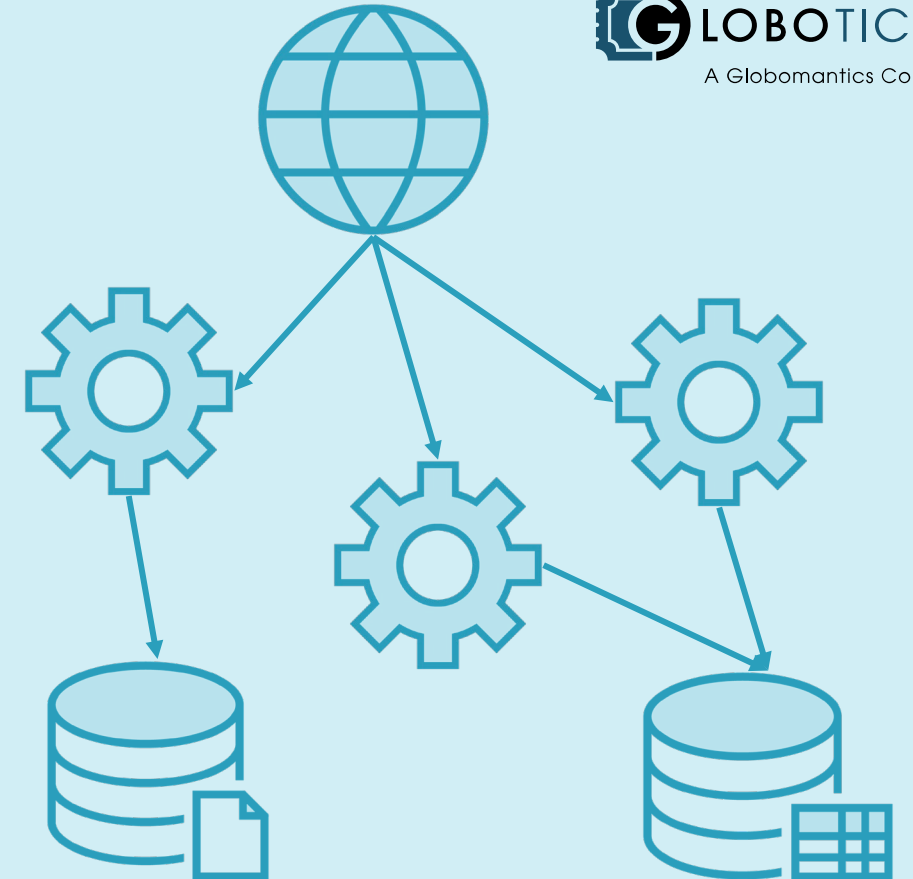- Communication detail
- Visualize interaction
- Needs tracing library
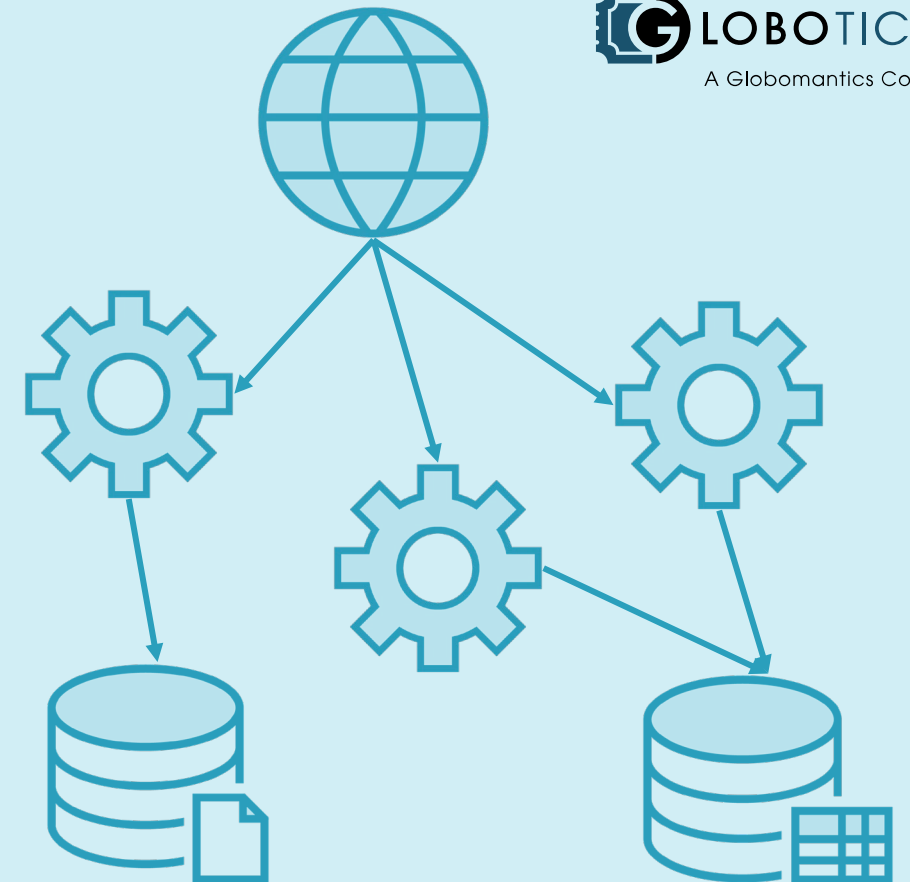
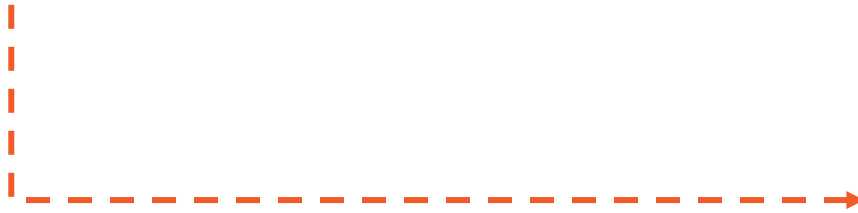# Scenario: Putting Metrics to Use
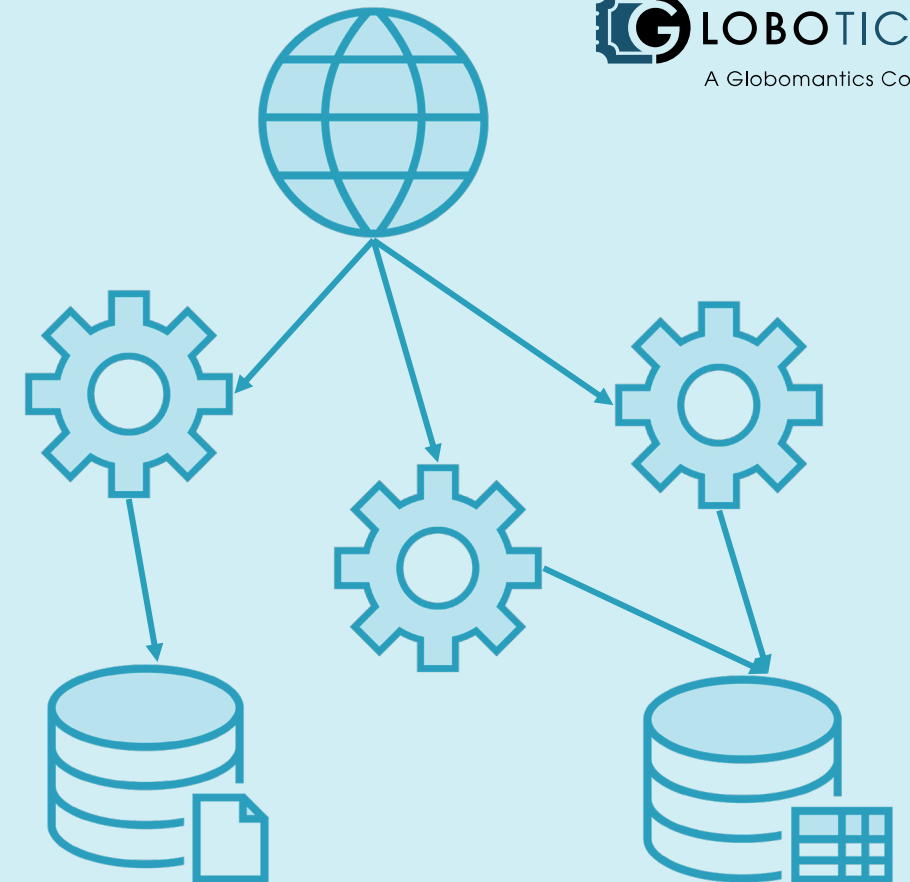
How do you add capacity when you need to?

Amila

- CPU & memory measured
- Tracked by component
- Alerts outside normal range

- Already auto-scaled
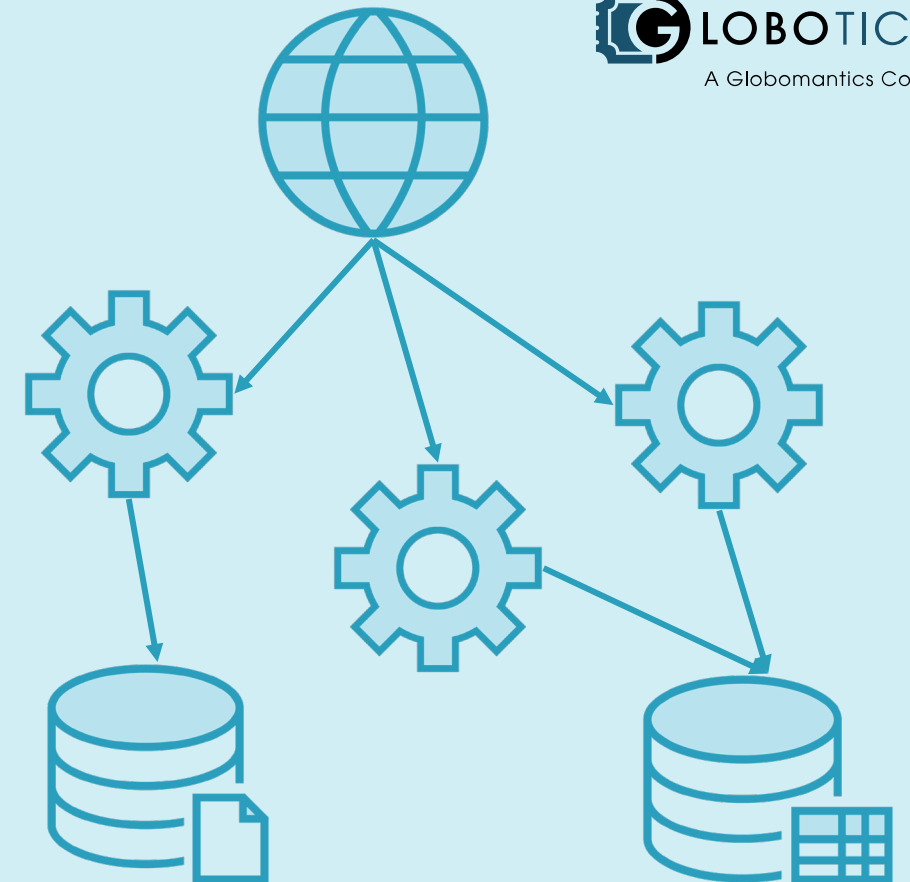- Engineer investigates
- Adds more scale if needed
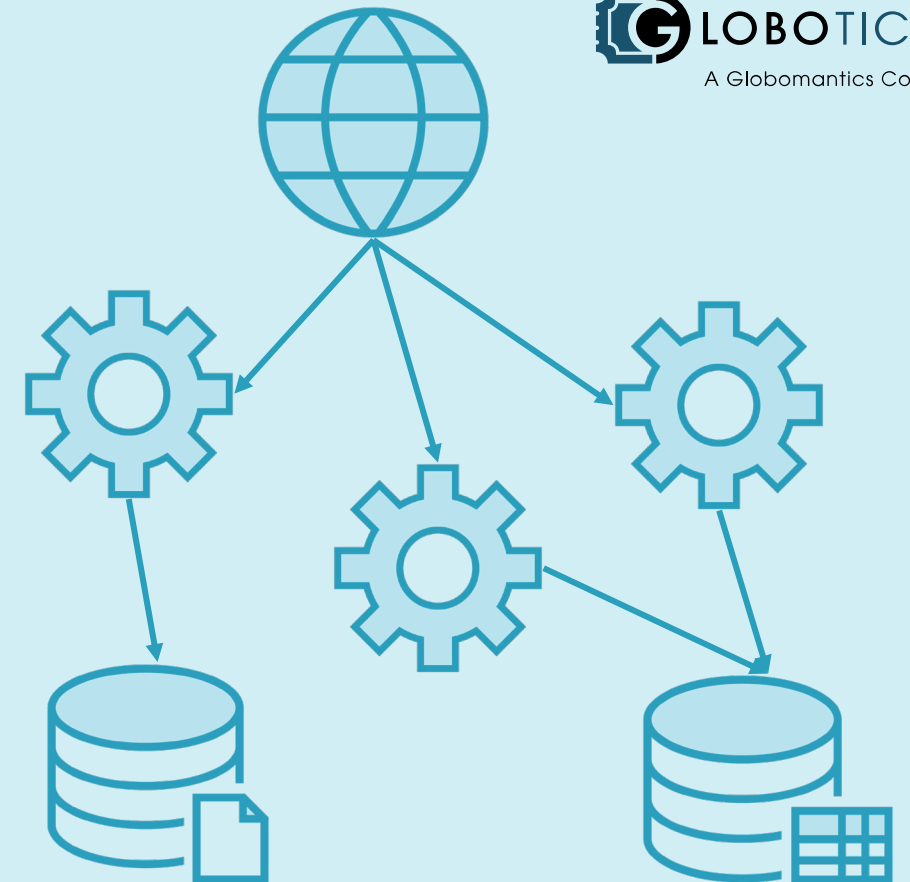
GLOBOTICKET
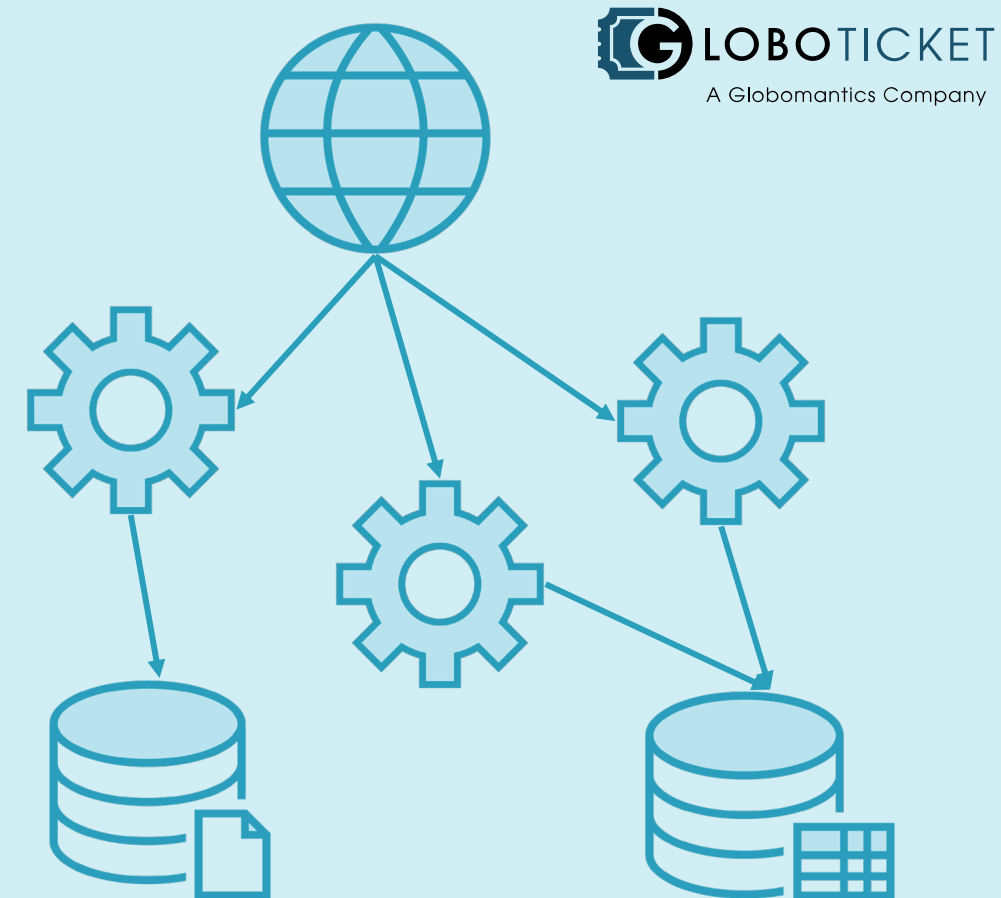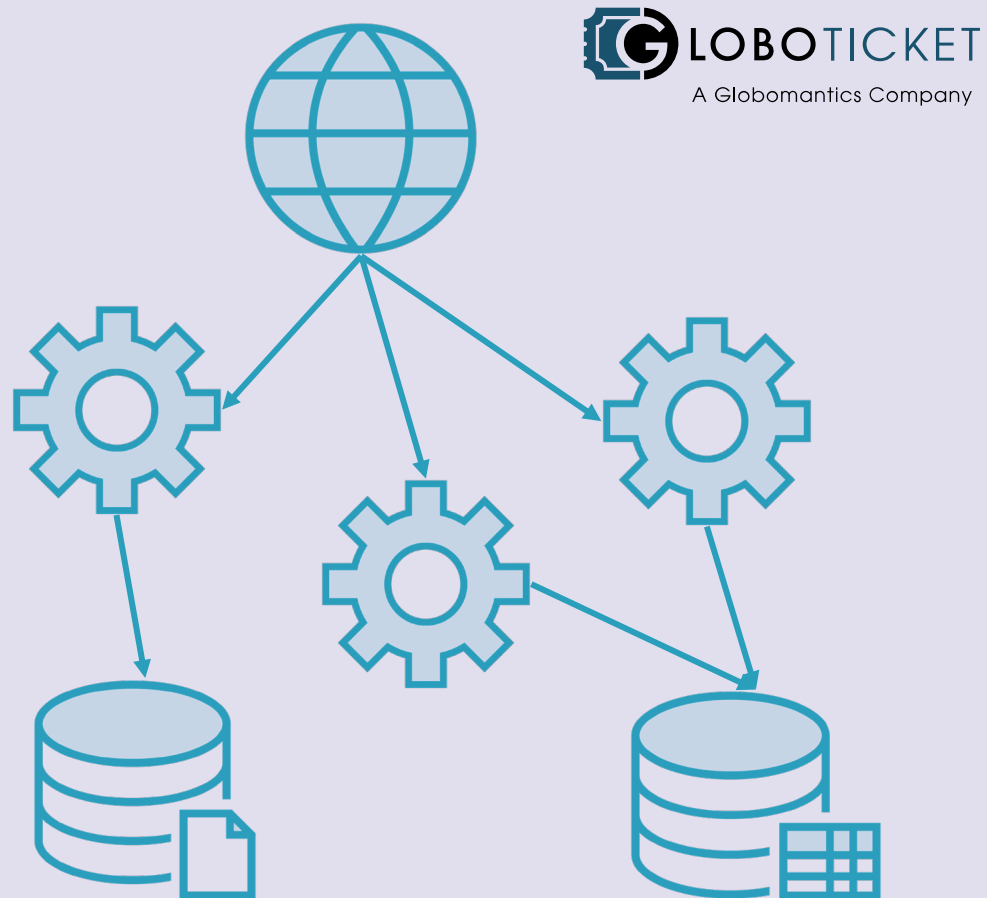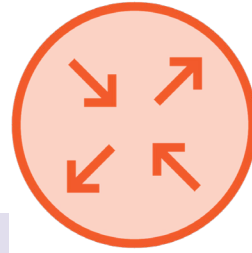A Globomantics Company

- Alert on budget
- Standard runbook
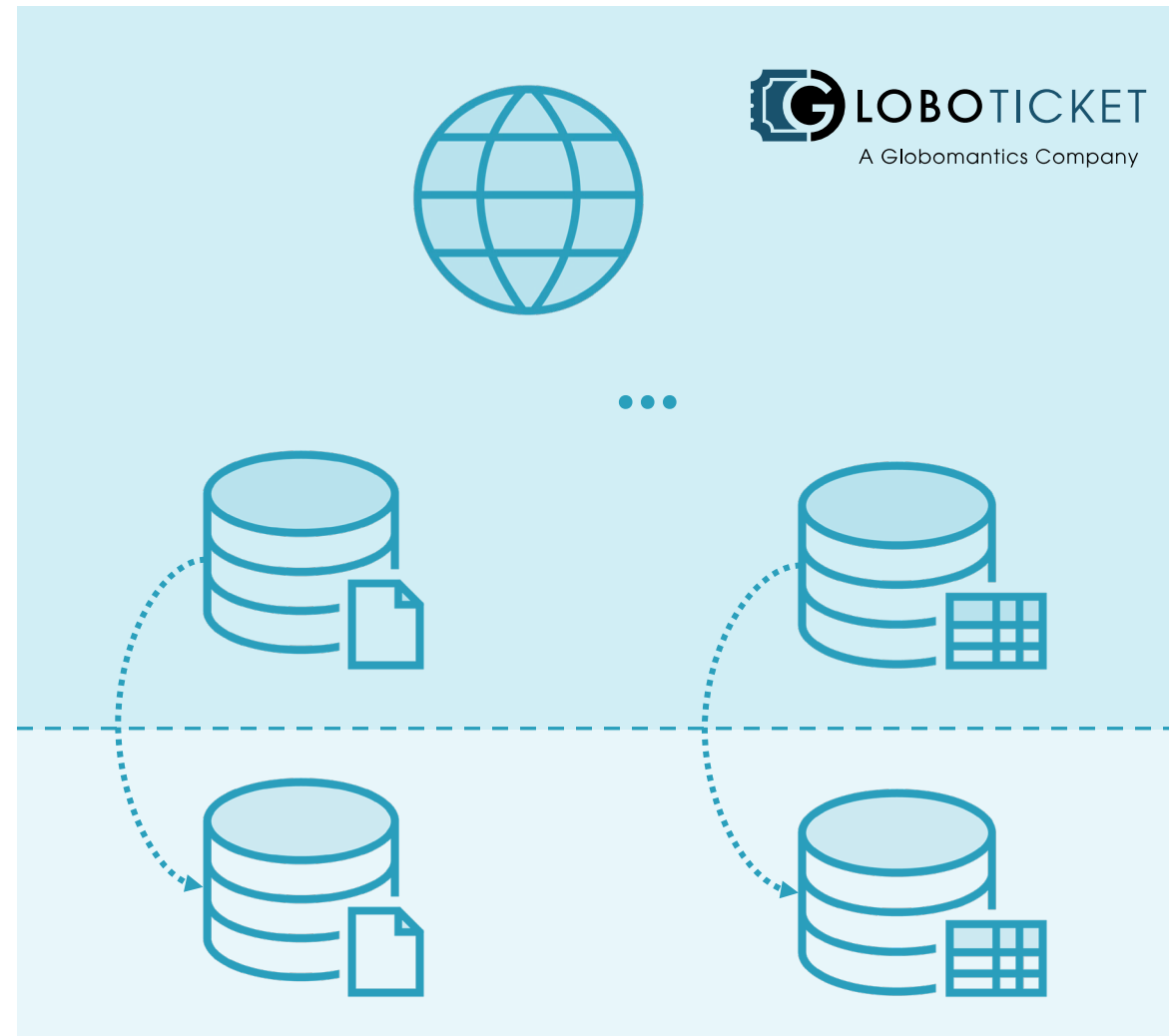- Scale guidance

Can you shift traffic between regions?
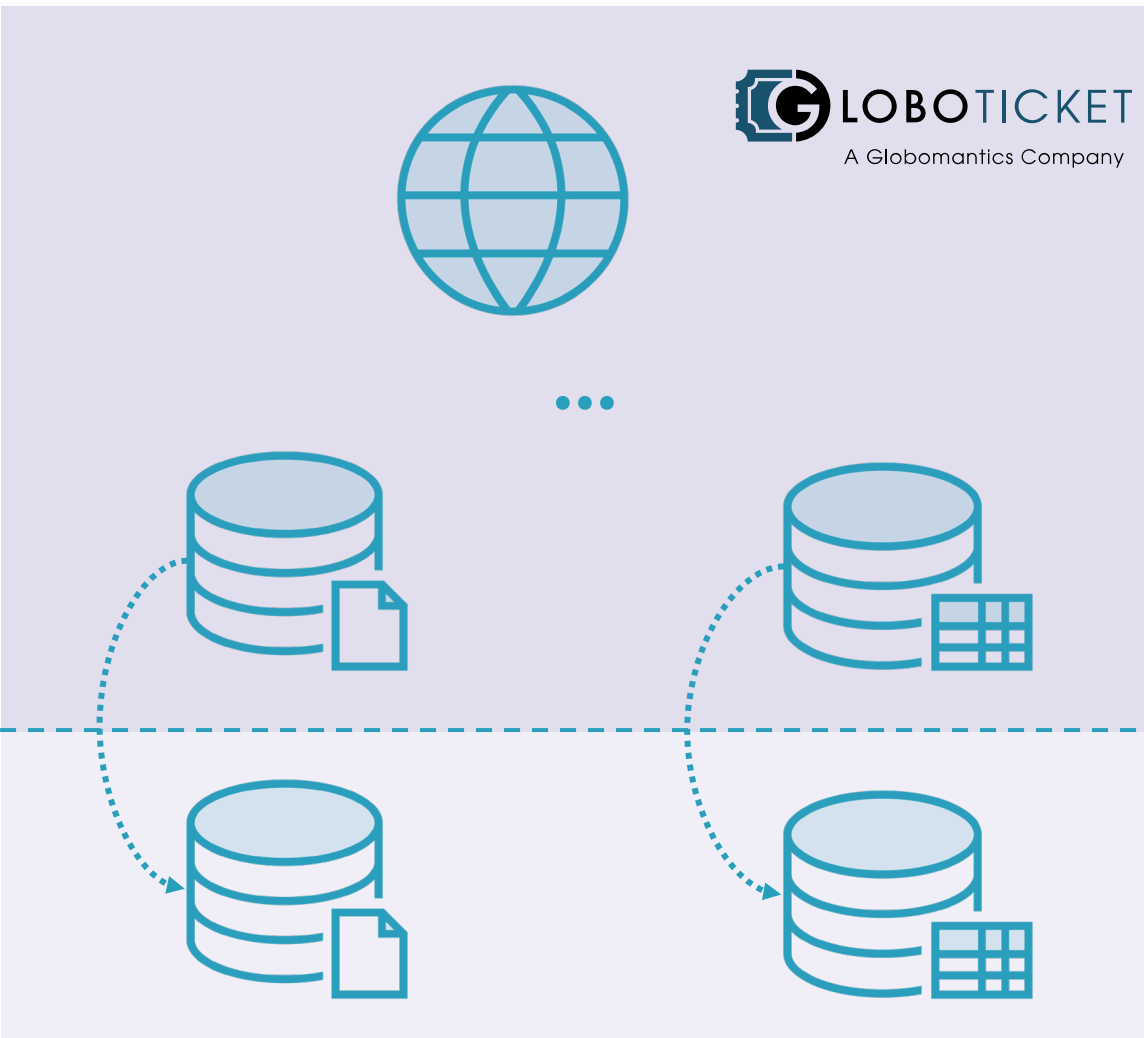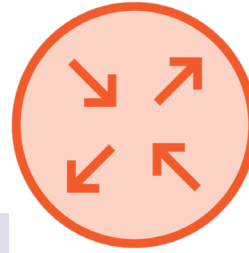
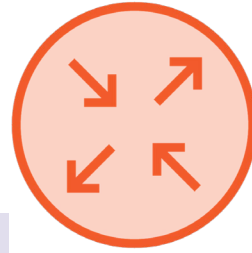GLOBOTICKET
A Globomantics Company

- DNS routing
- Latency + health
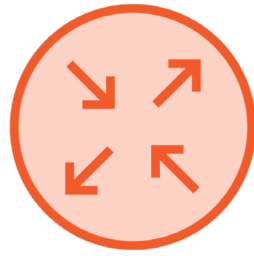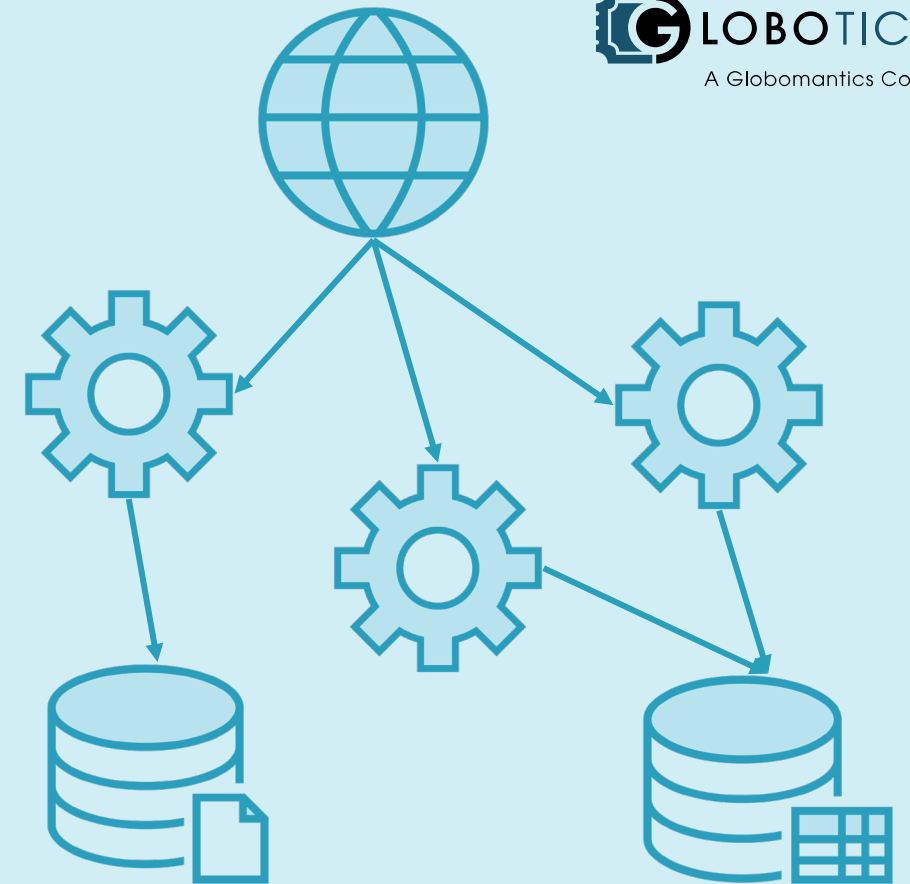
- Region-specific data
- Replicated in region

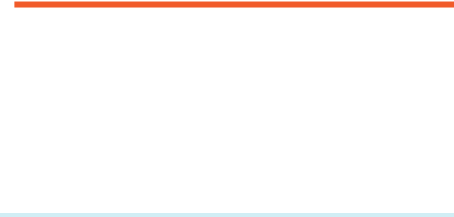- Global replication
- Platform feature

- Global data
- Routing metrics
- Manual override

GLOBOTICKET
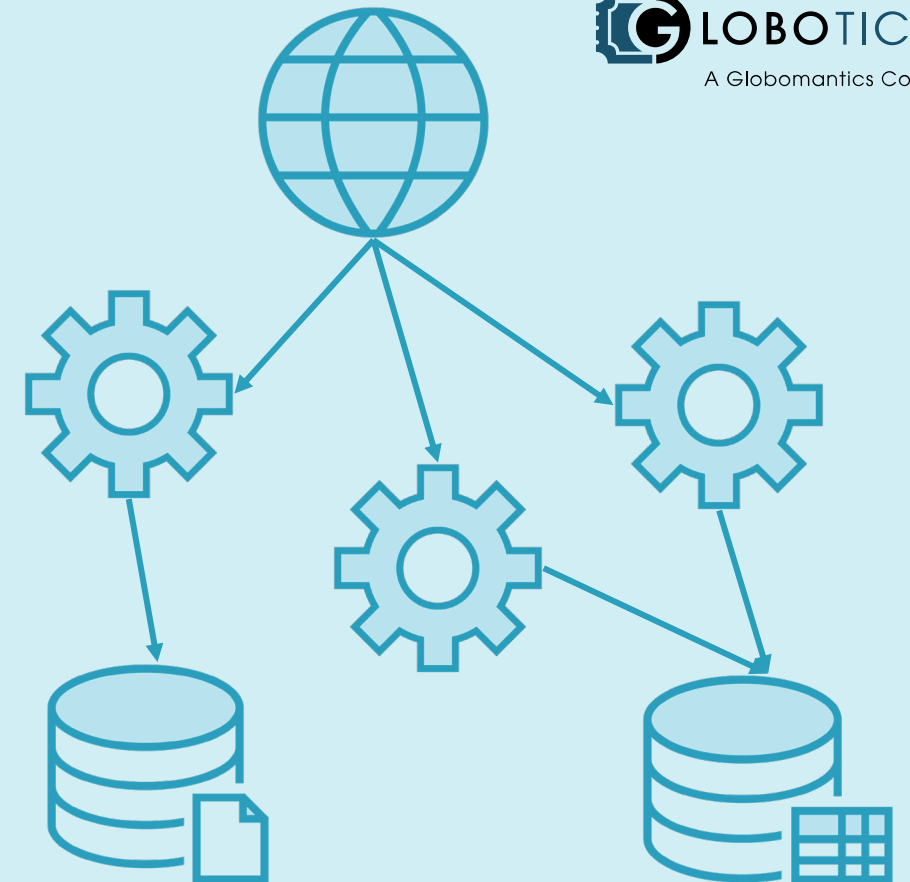A Globomantics Company

When do you downgrade service levels?
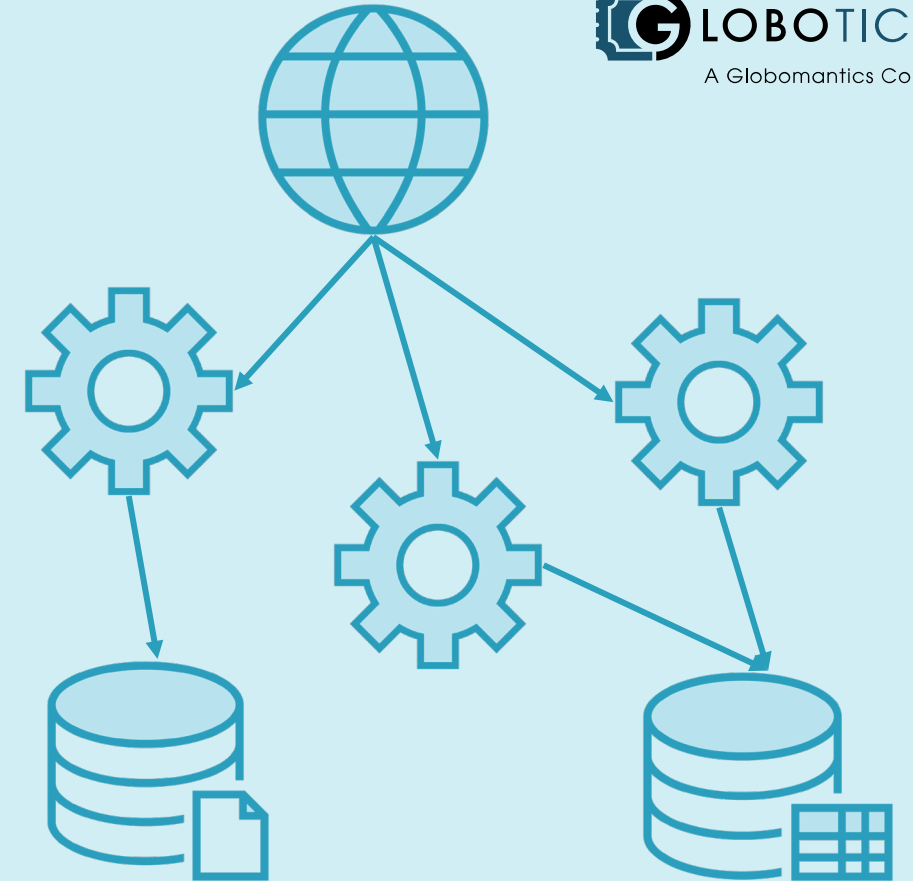
- Component saturation?
- Region saturation?
- Global saturation?

- Health dashboards
- Metric metadata
- Quick triage

LOBOTICKET
A Globomantics Company

# Supporting Triage with High-level Metadata
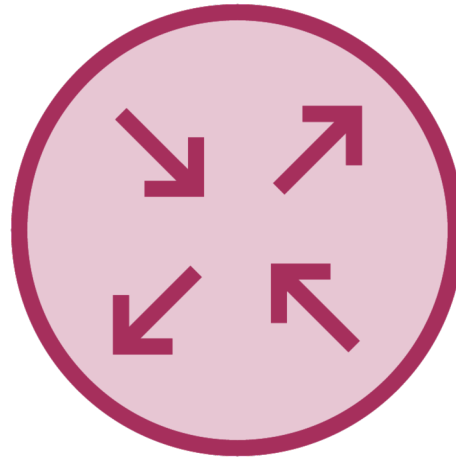
Monitoring

# Triage: First Responses



**Scale up**

*If compute-bound*

**Re-route traffic**

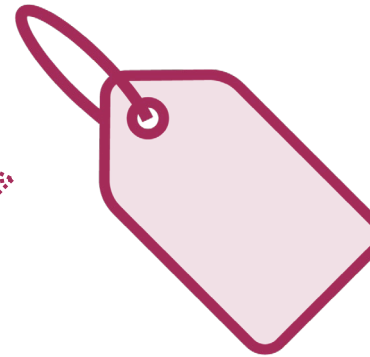*If regional*

**Downgrade service**

*If global*

# Core Metrics

**Traffic**
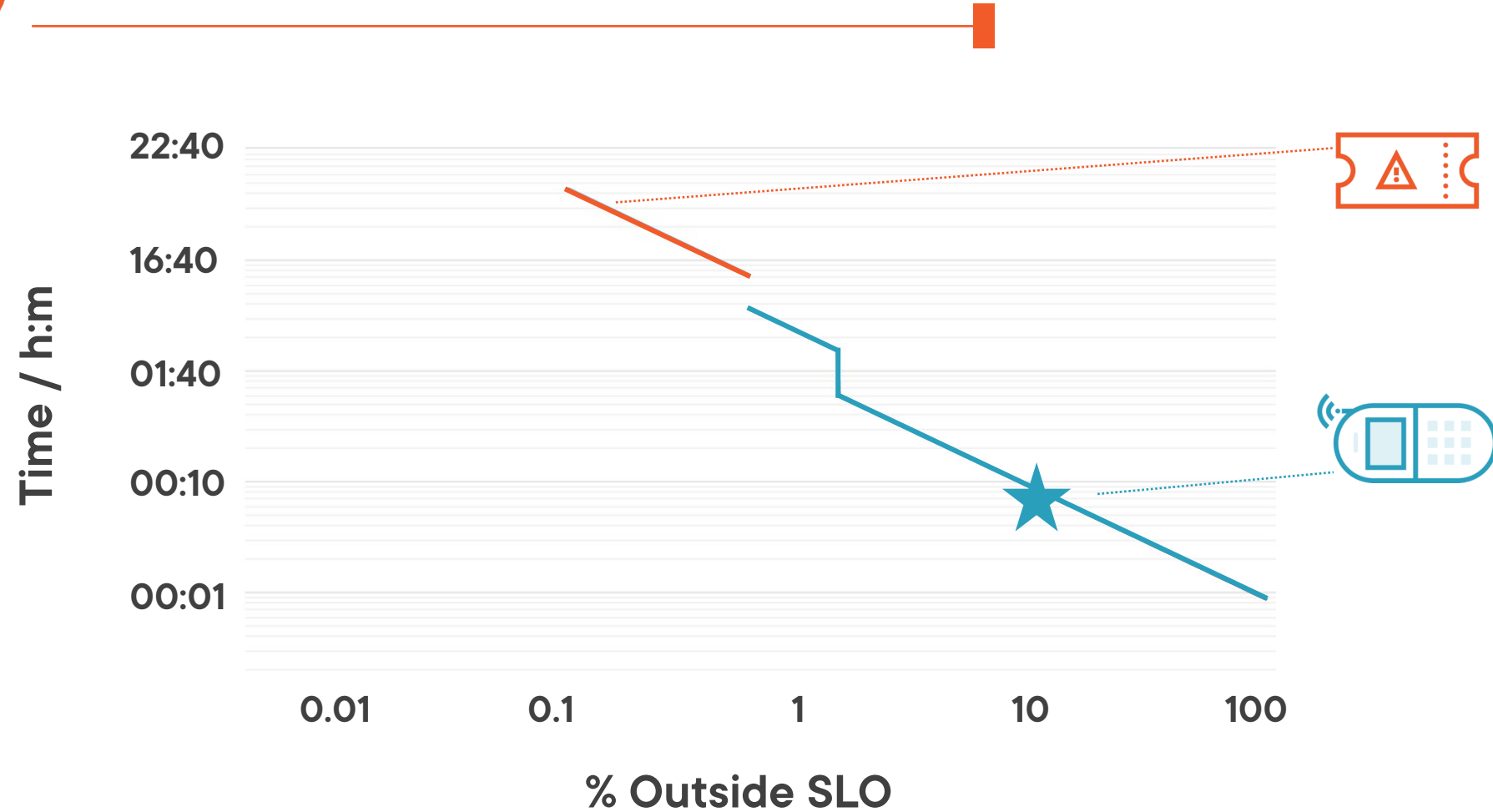
**Latency**

**Saturation**

- Region
- Component

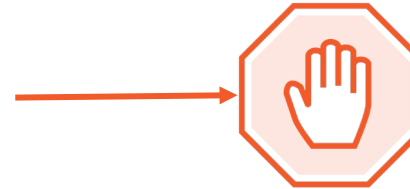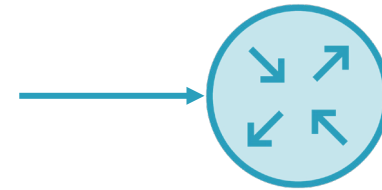# Alerting on Error Budget Burn

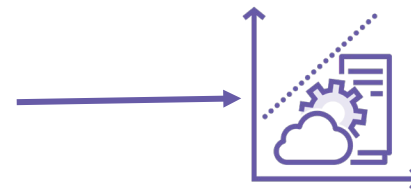**Latency**  **SLO: 99.9% within 2sec**

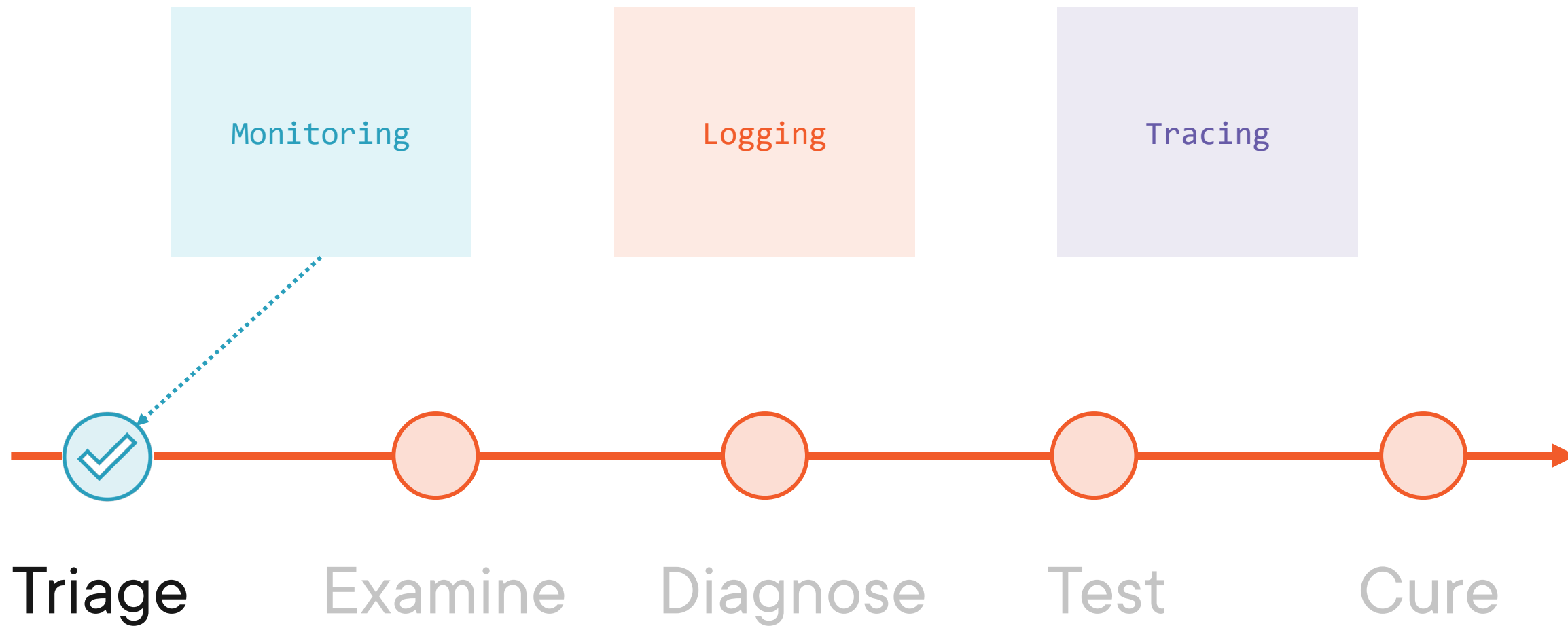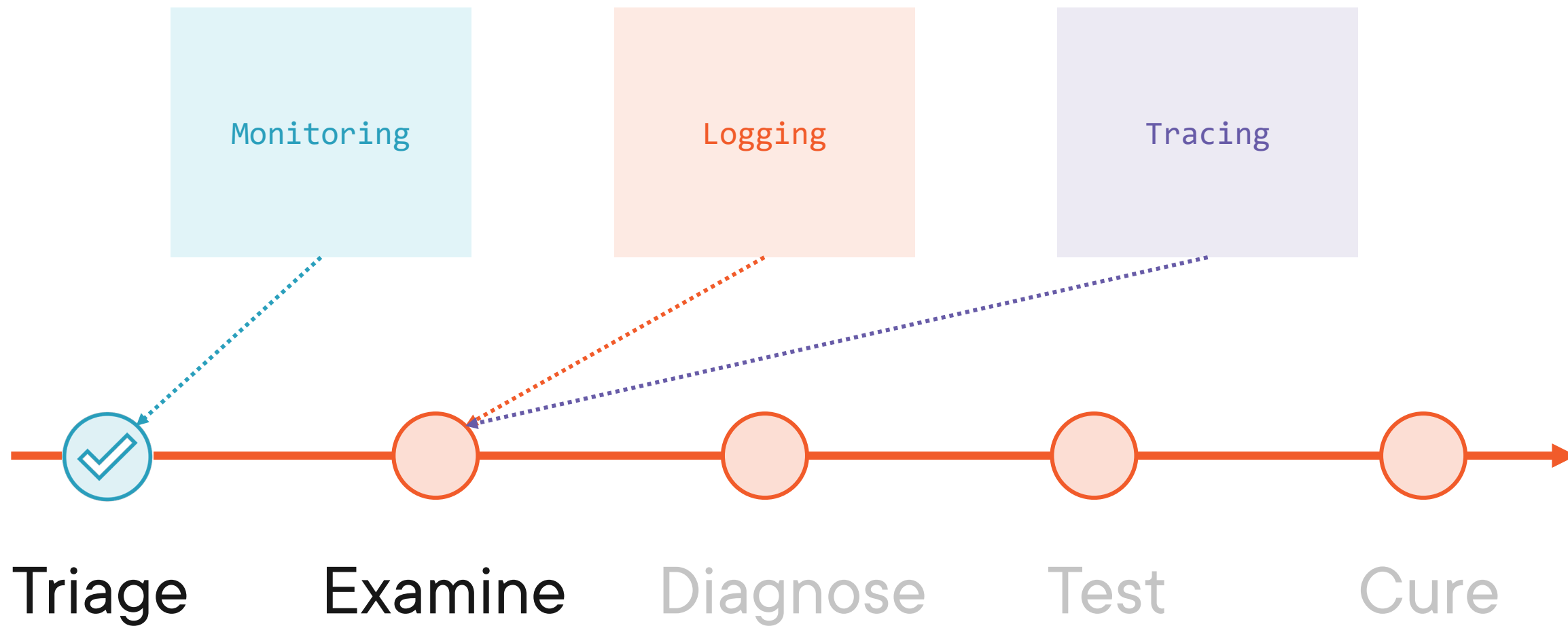# Identifying First Response


Traffic and saturation high globally


Traffic high for all components in one region


One component saturated in one region

Monitoring

Logging

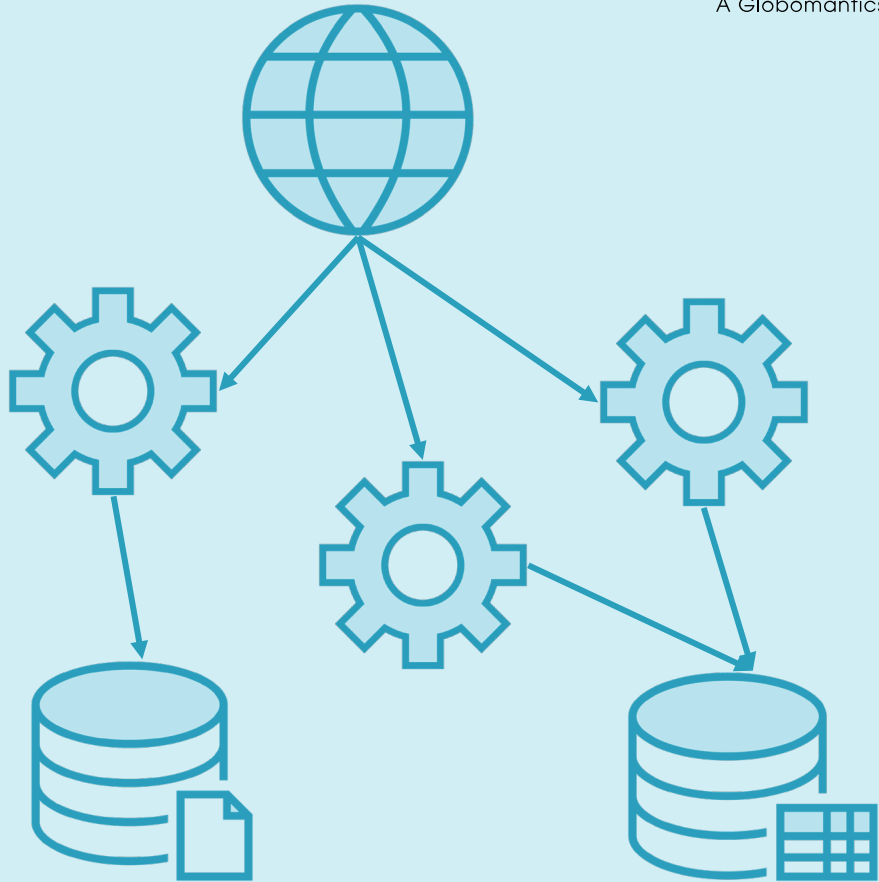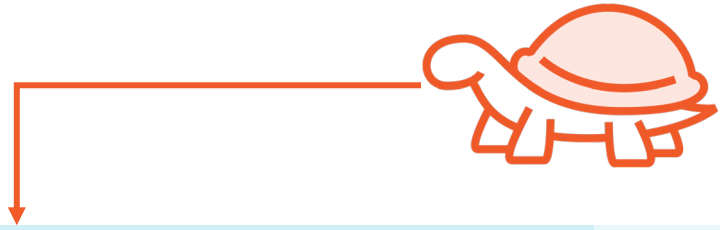Tracing

Triage    Examine    Diagnose    Test    Cure

Monitoring

Logging

Tracing

Triage  Examine  Diagnose  Test  Cure

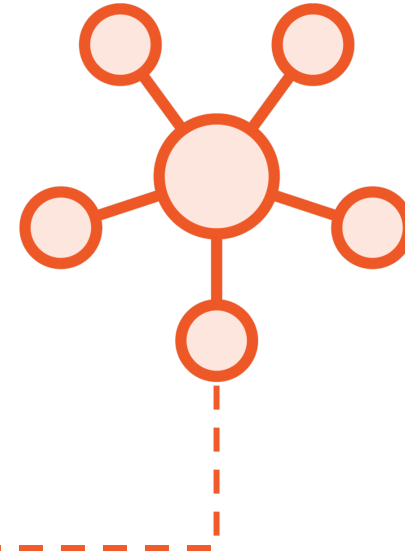# Scenario: Putting Traces to Use

LOBOTICKET

A Globomantics Company
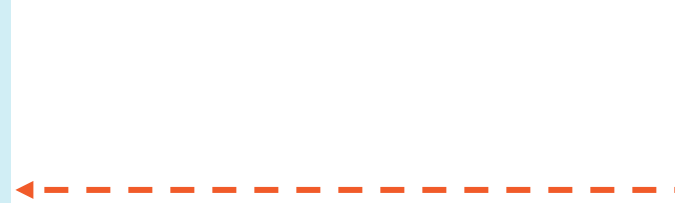
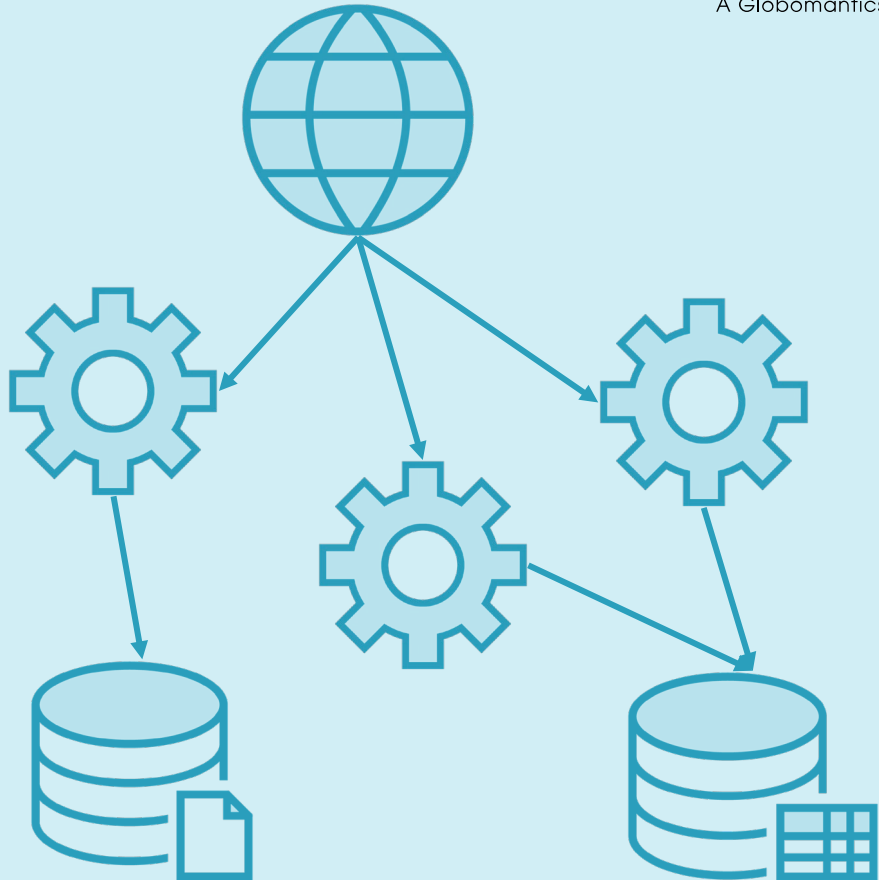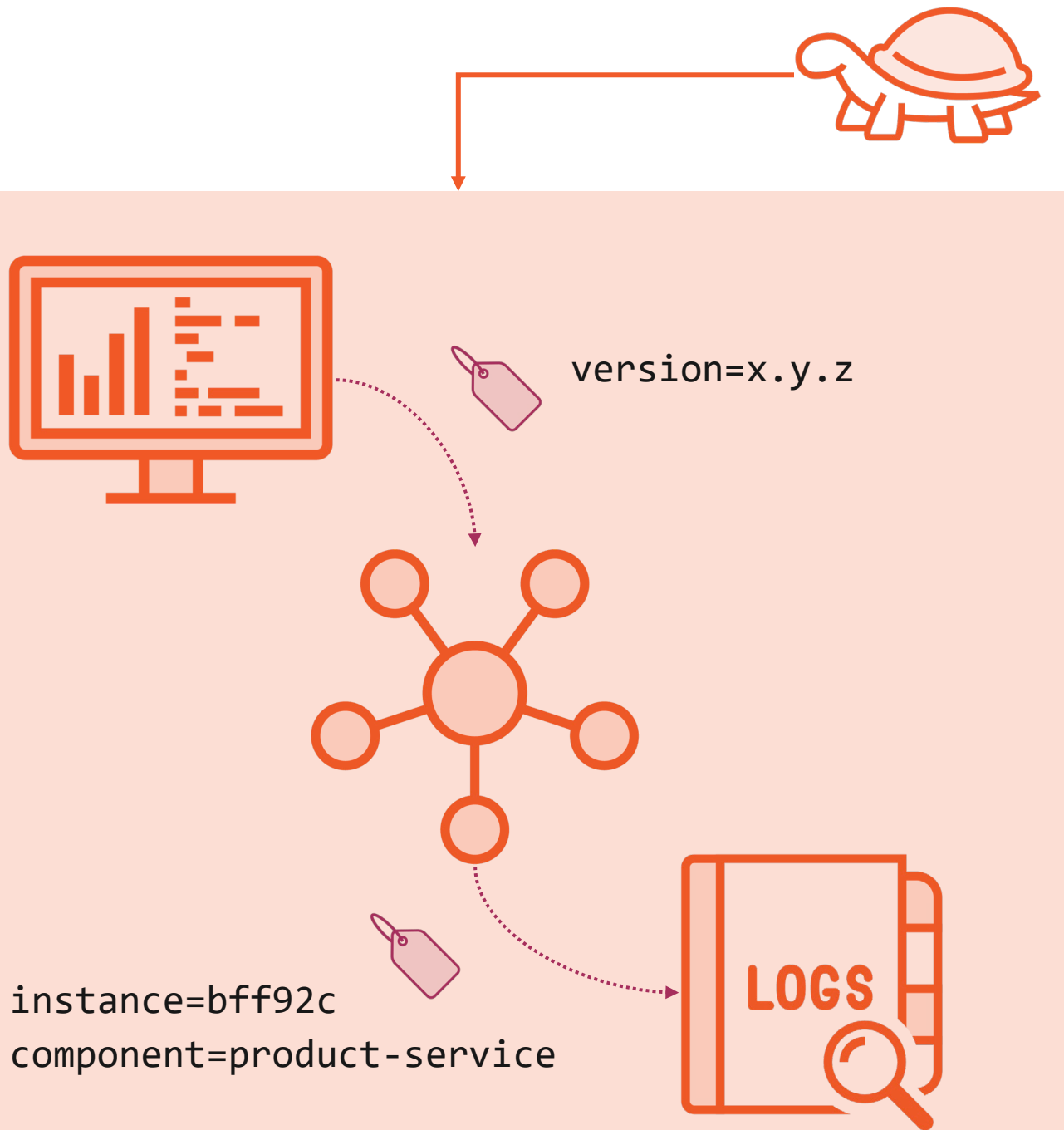Latency spikes - which is the slow service?

- Versioned metrics
- Filter before and after
- Confirm release trigger

- Filter traces by time
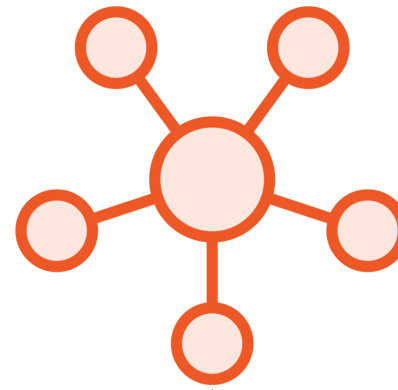- Sample before and after
- Identify slow instance

- Filter logs by instance
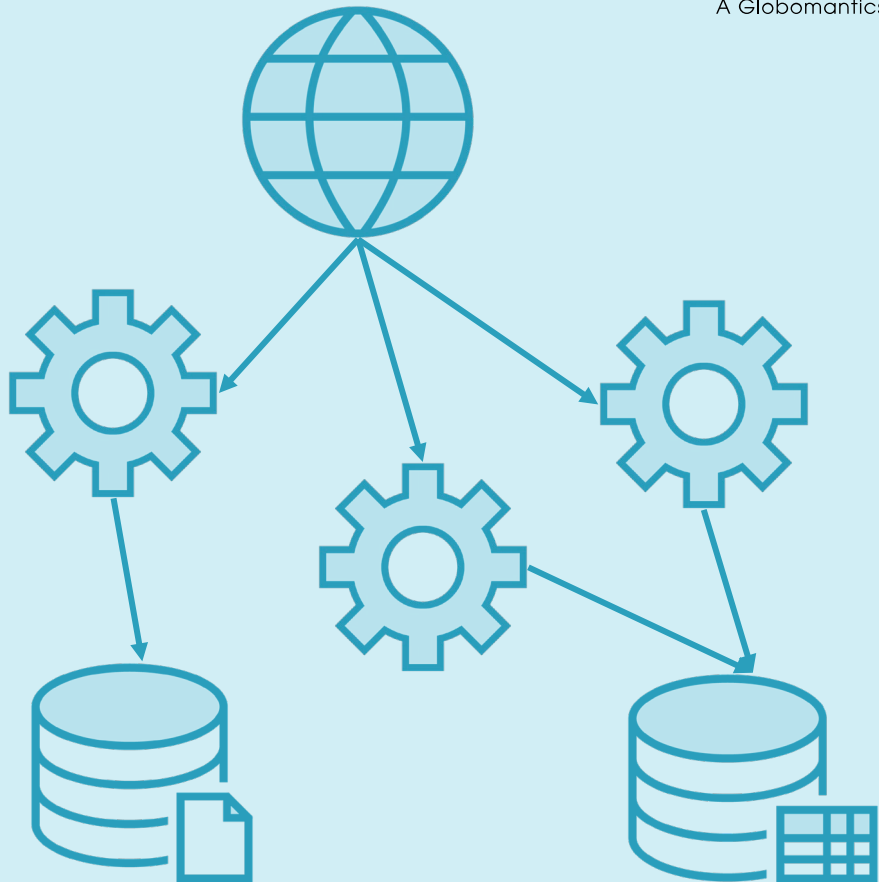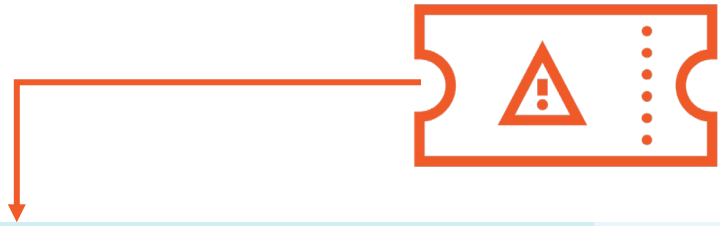- Filter by time
- Drill into detail

version=x.y.z

instance=bff92c
component=product-service

- Version in traces
- Trace ID in logs
- Is the data there?

- 1/1000 traces sampled
- Logs at WARN level

LOBOTICKET
A Globomantics Company

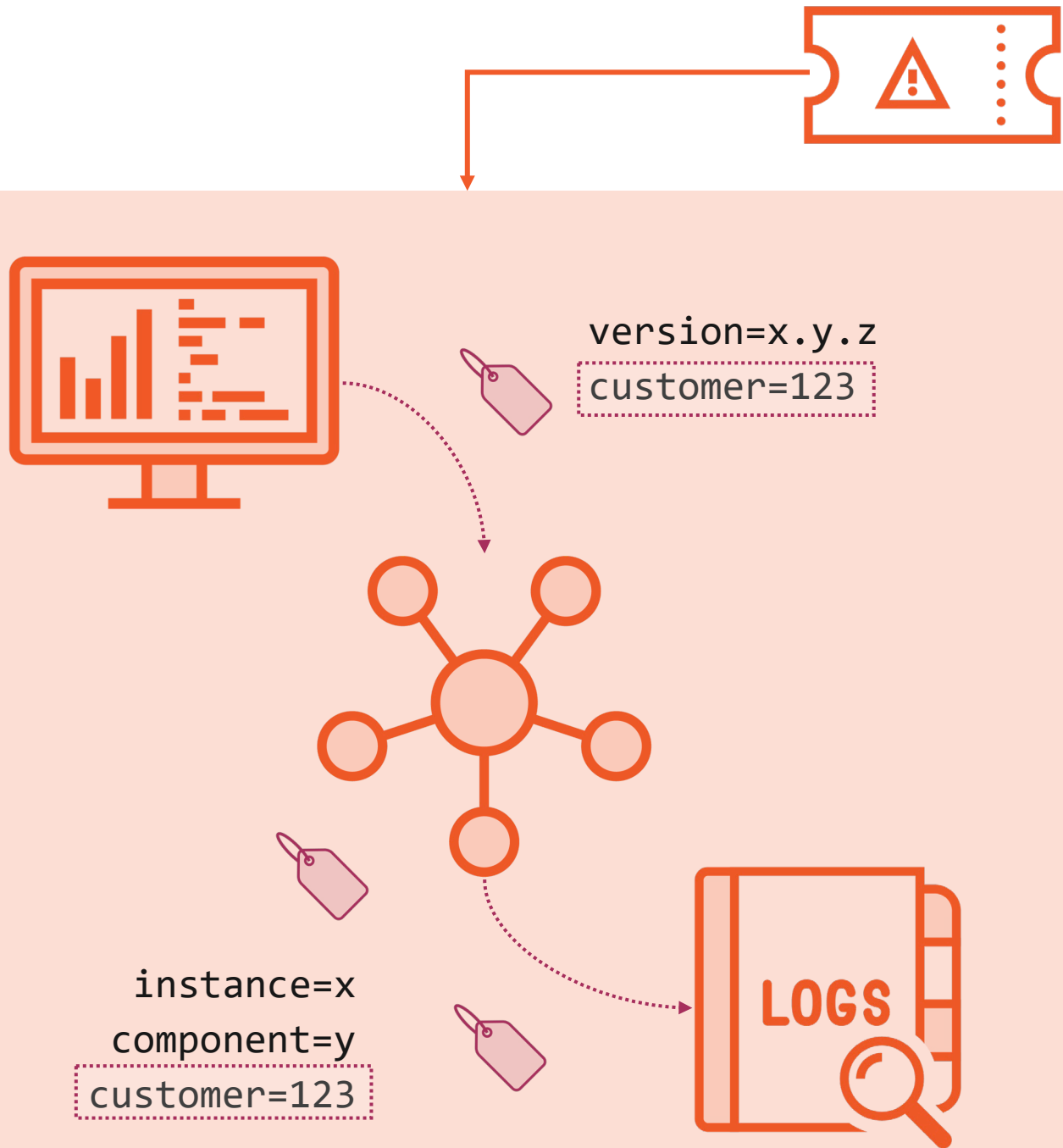Poor experience for one customer - why?

- No customer IDs
- Specific feature?
- Call product team

version=x.y.z
customer=123

instance=x
component=y

LOGS

- Low-level metadata
- Not for all metrics
- Customer SLOs

version=x.y.z
customer=123

instance=x
component=y
customer=123

LOGS

- Trace by customer
- Add ID to logs
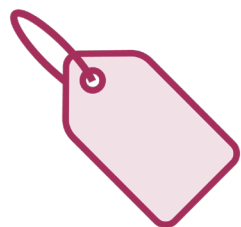- Work backwards

version=x.y.z
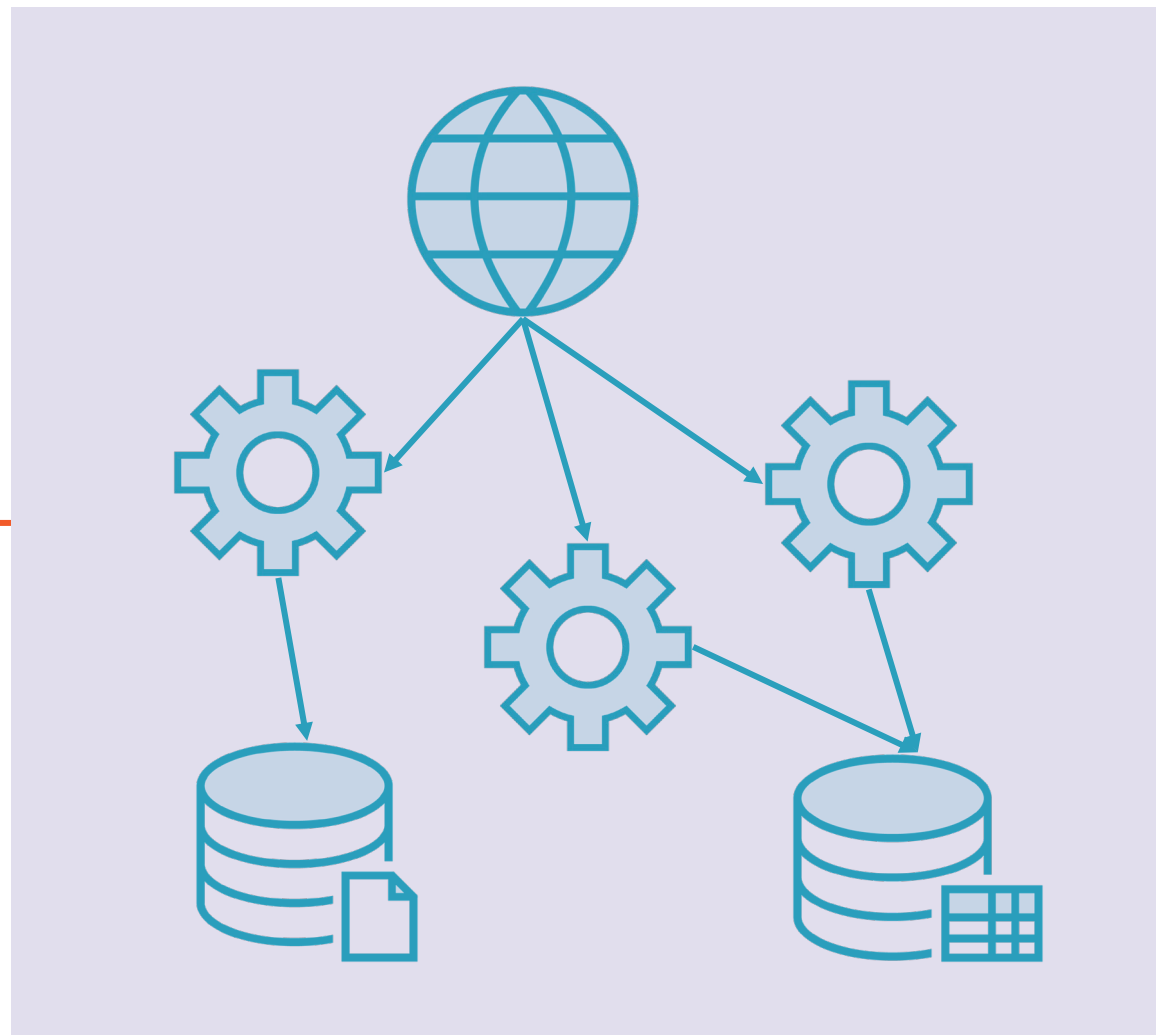customer=123

instance=x
component=y
customer=123

LOGS

- Is the data there?
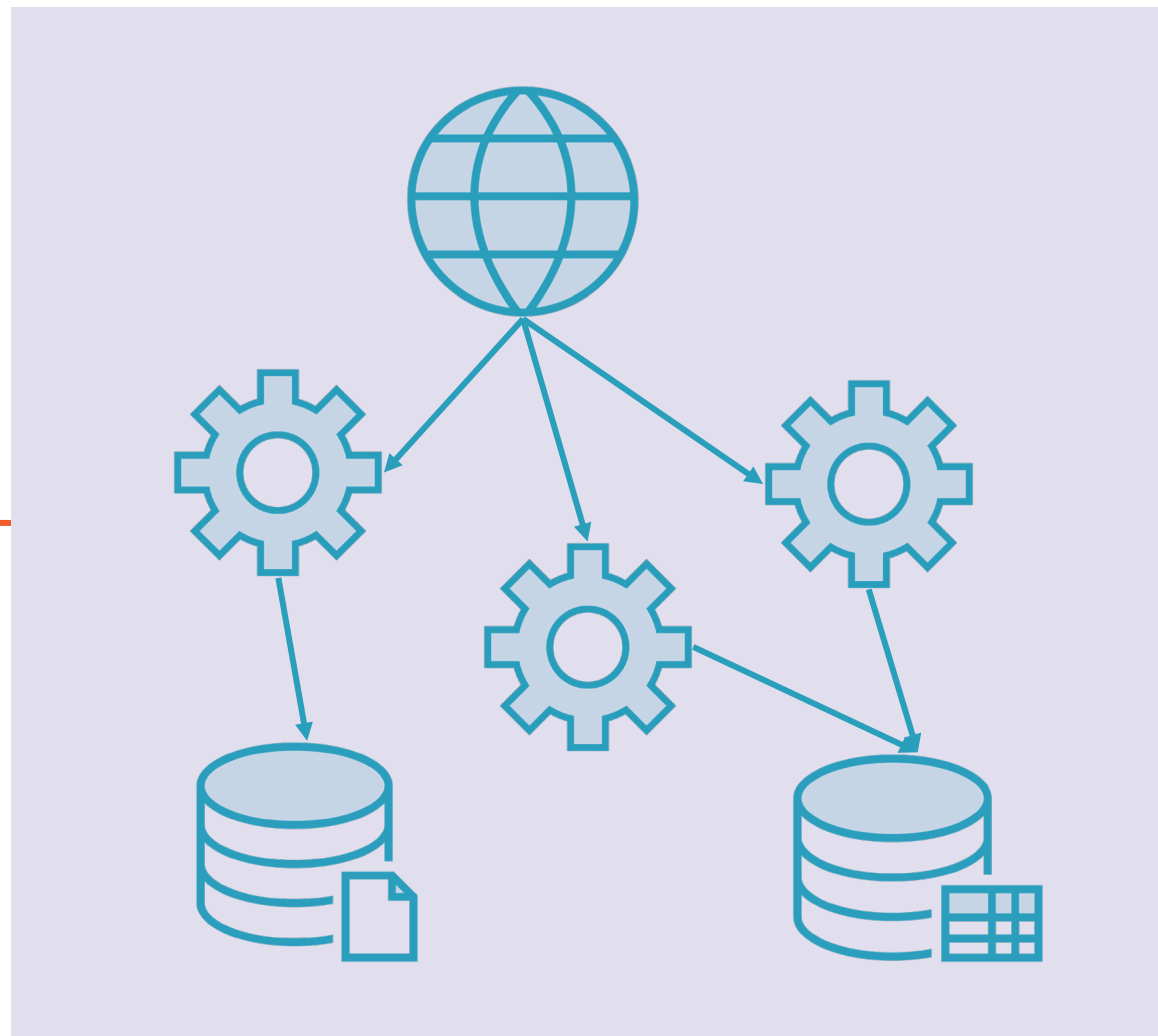
# Supporting Examination with Low-level Metadata

latency{component=1,region=1}: 0.9

latency{component=1,region=2}: 0.7

...

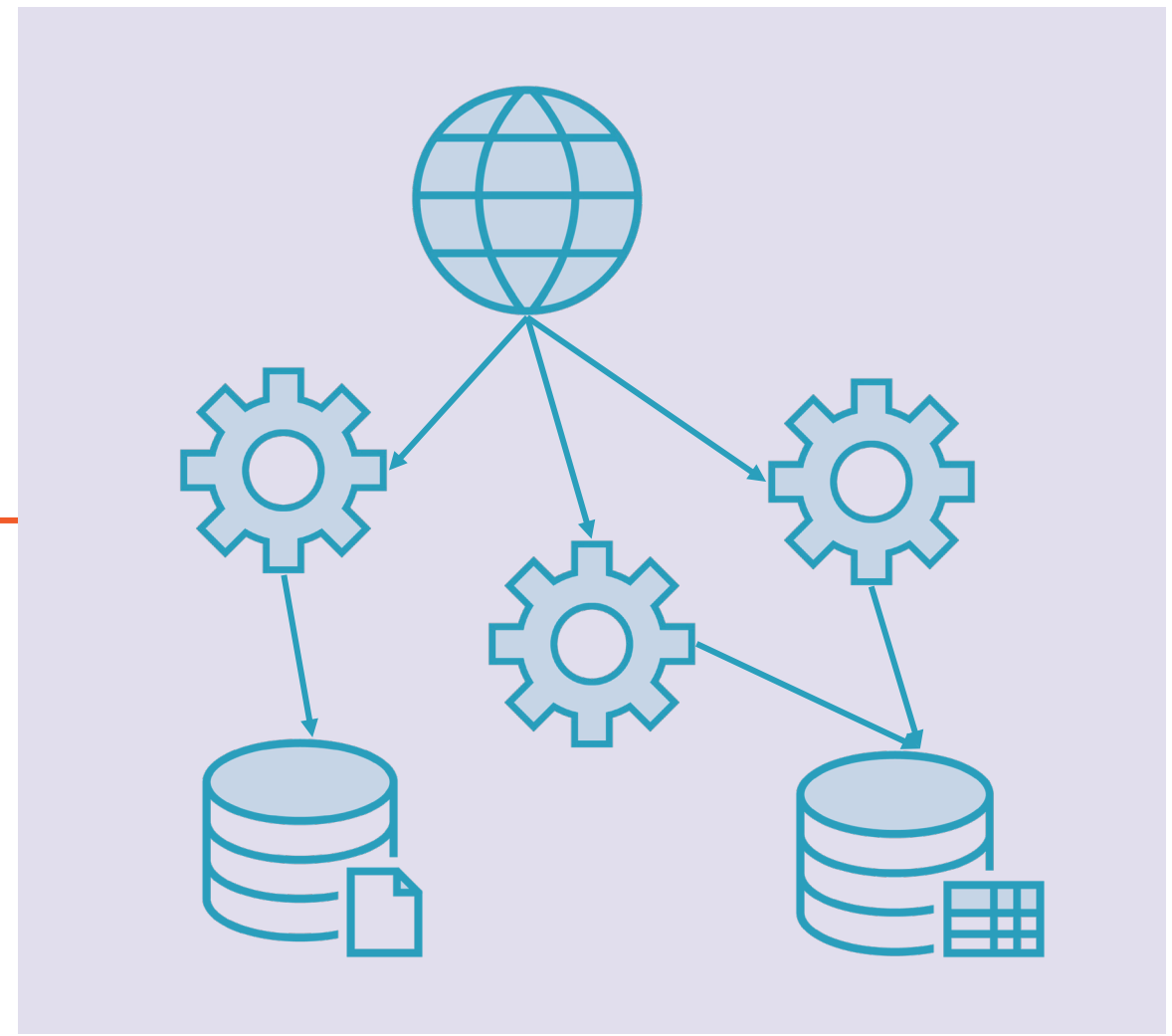latency{component=20,region=4}: 0.1
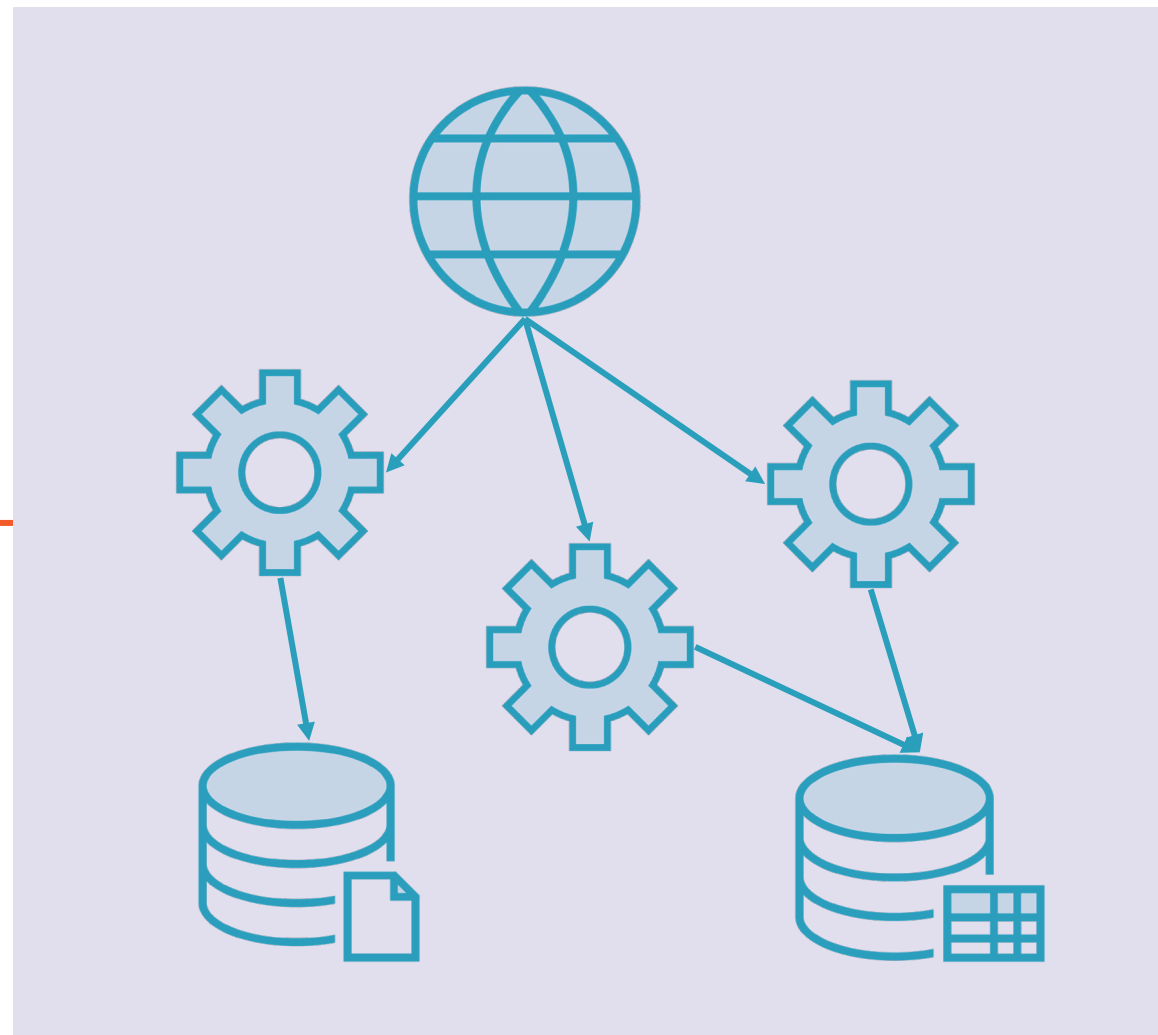
Monitoring

480/hr

Monitoring
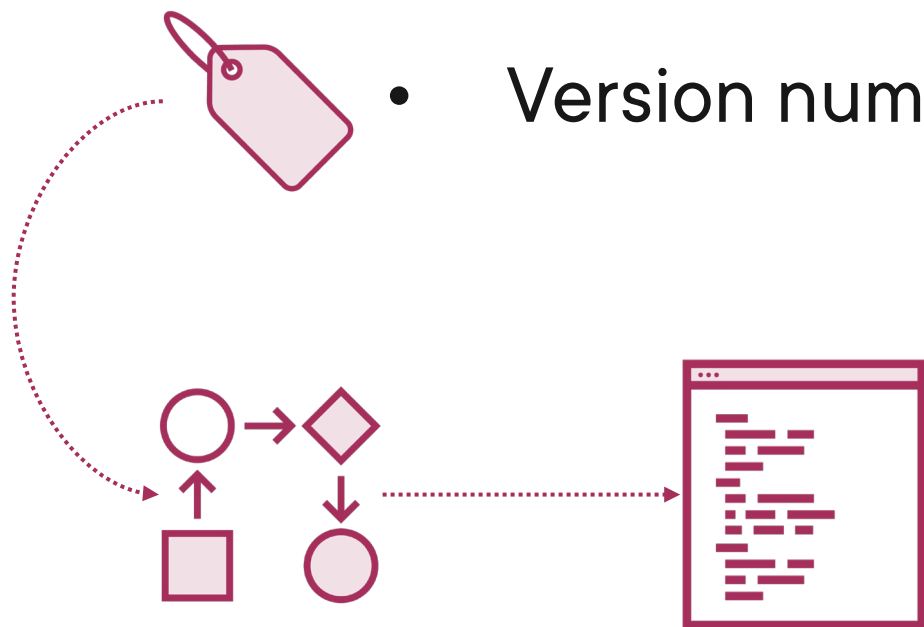
latency{customerId=1}: 0.8

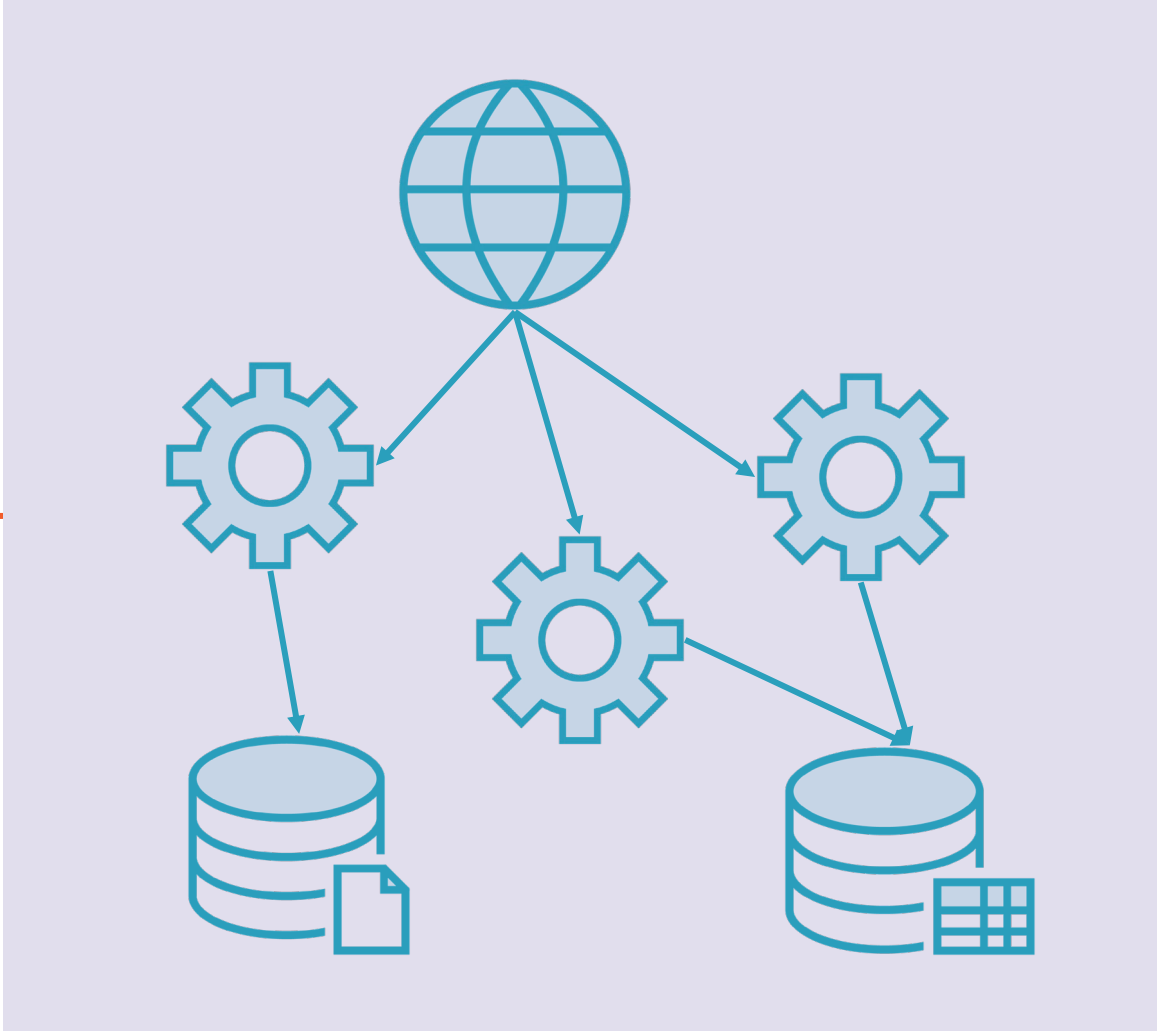latency{customerId=2}: 0.8
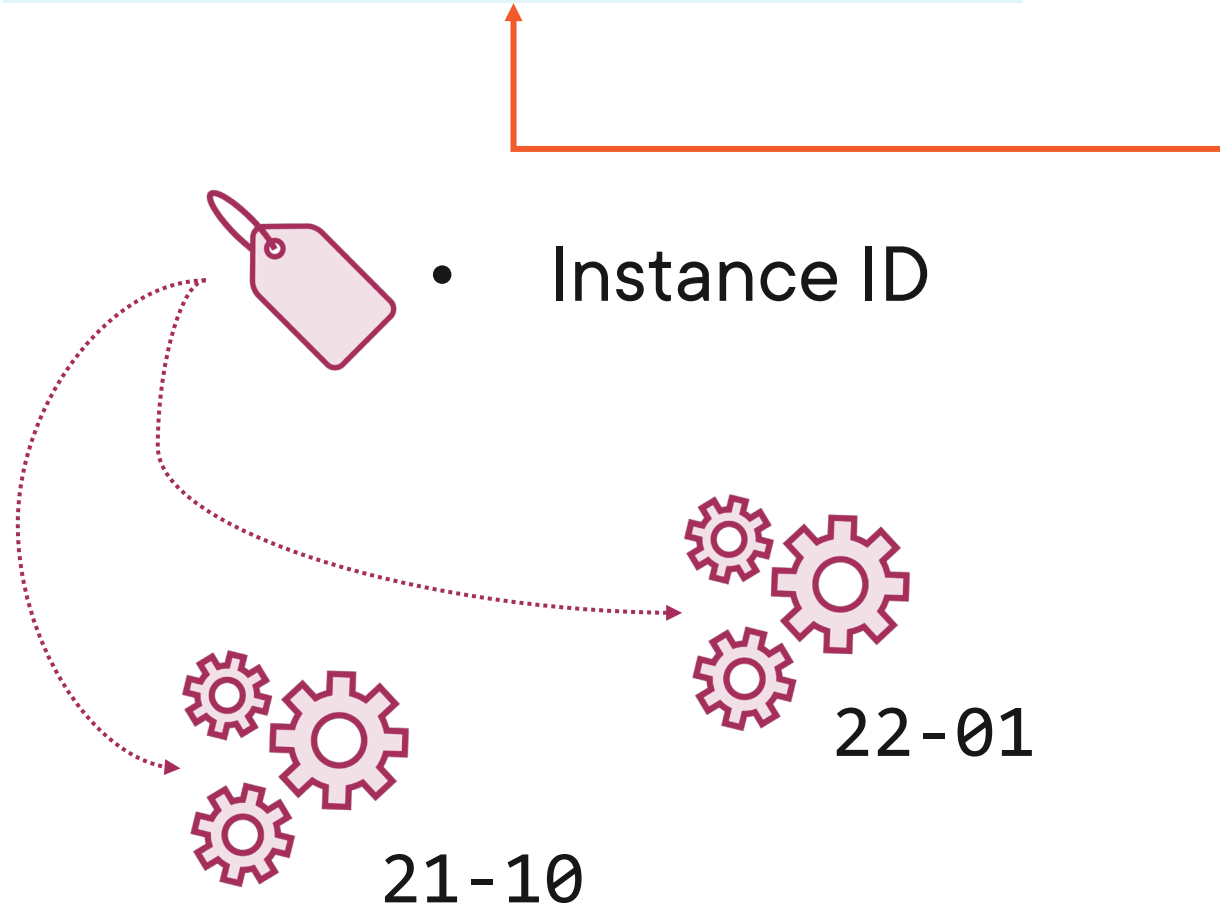
...

latency{customerId=10000}: 0.9

1.2M/hr

Monitoring

- Version number

API

Monitoring

Instance ID

22-01

21-10

API

Tracing
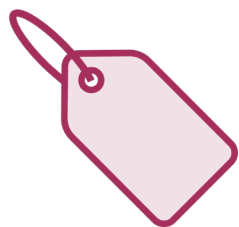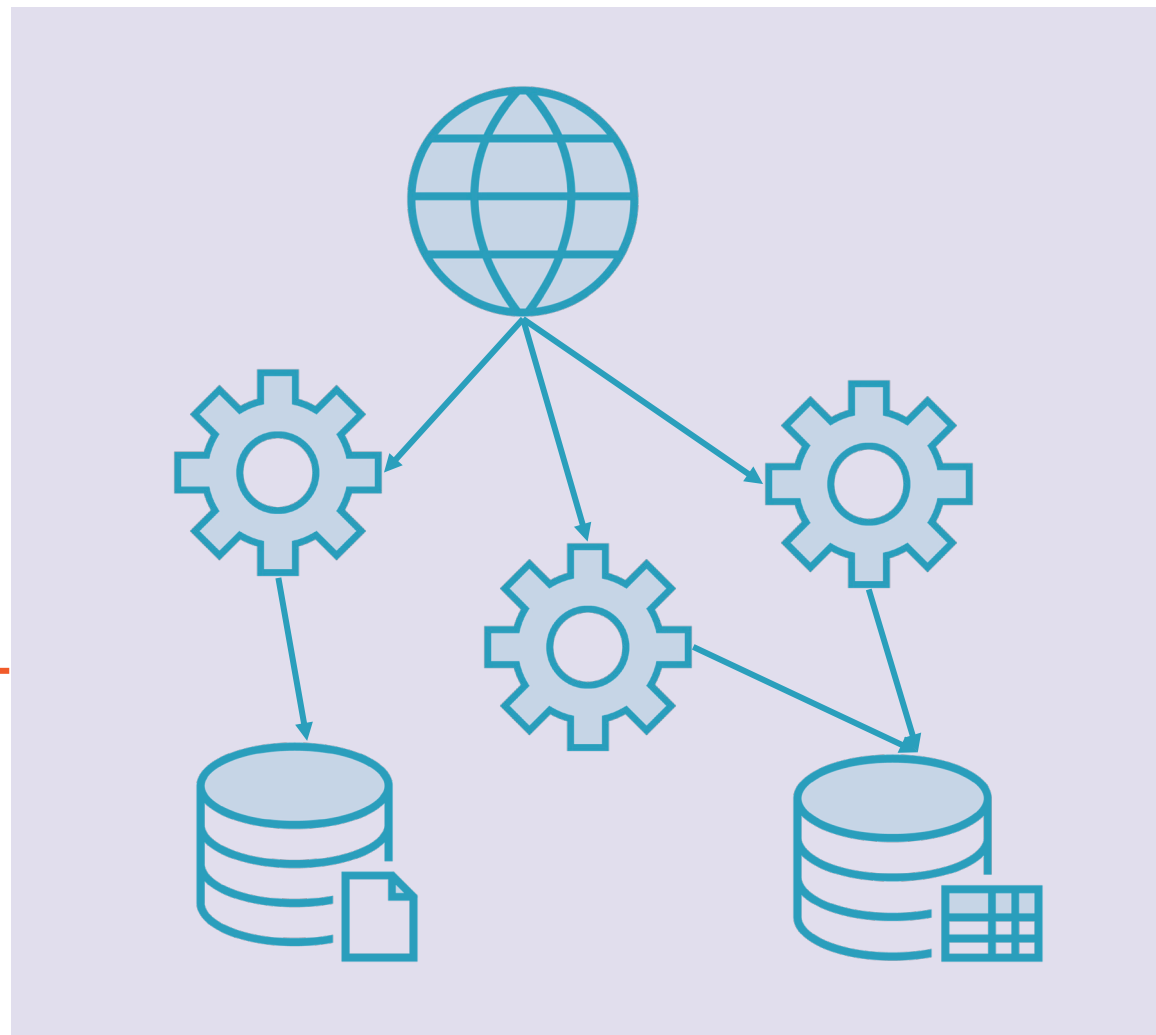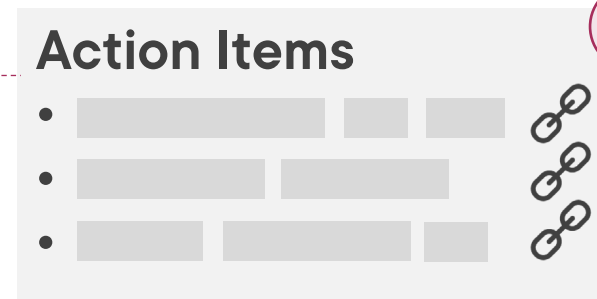
- Region
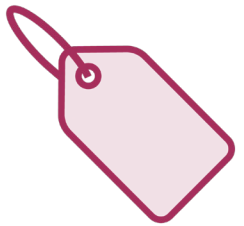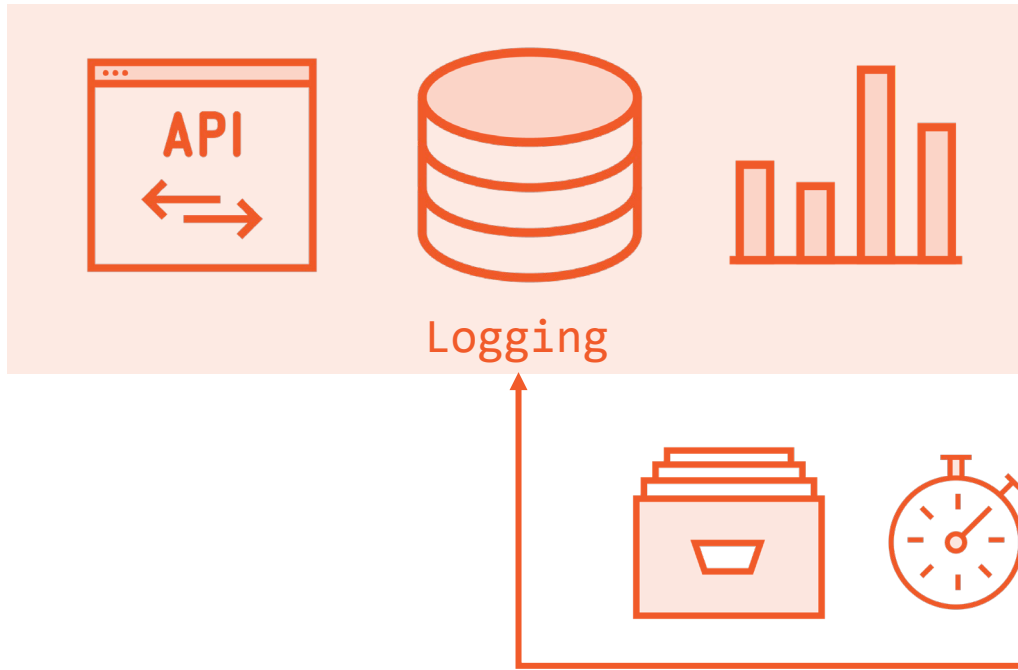- Instance ID
- Customer ID

- Identifying region took >1hr

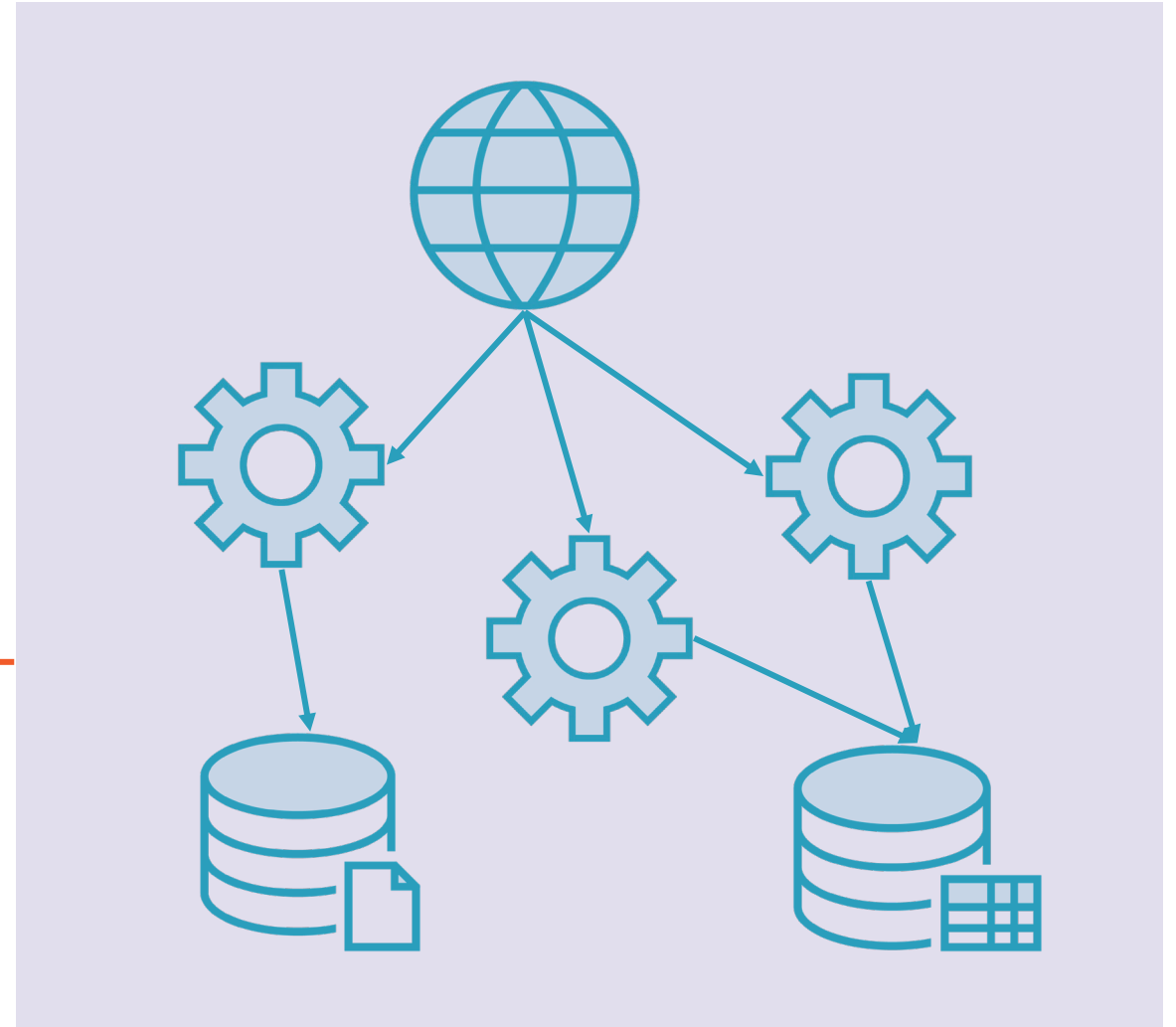- Need metadata to correlate metrics & traces

- Add region ID to Prometheus & Jaeger

**Timeline**

**Lessons Learned**

**Action Items**
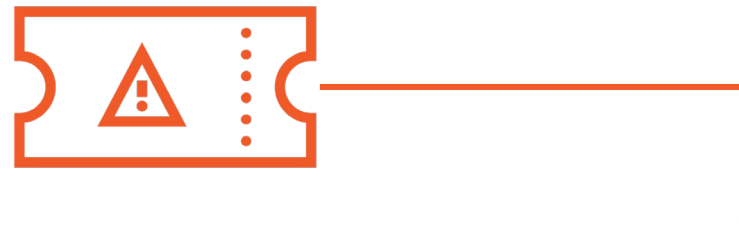
- Region
- Instance ID
- Version
- Customer ID
- Transaction ID

# Scenario: Putting Logs to Use

How do you find the log for a unique error ID?
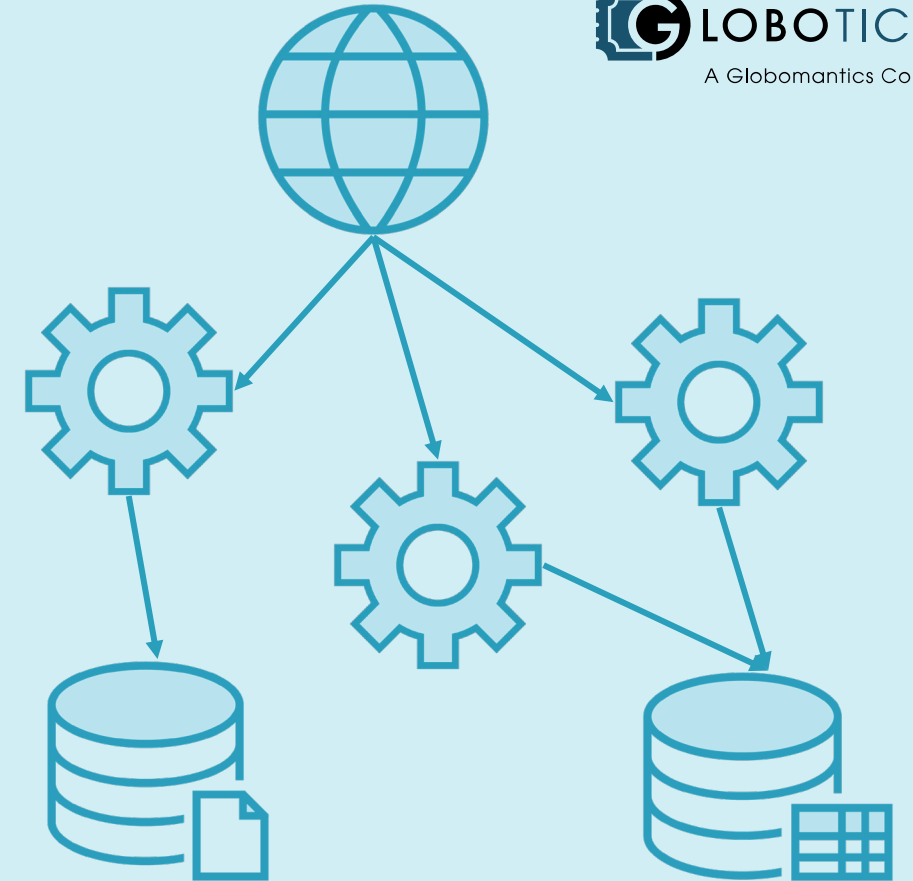
GLOBOTICKET
A Globomantics Company

LOGS

- All errors are logged
- UI for full-text search
- Component in metadata

GLOBOTICKET
A Globomantics Company

- Structured logs
- Key data in fields
- Efficient filtering
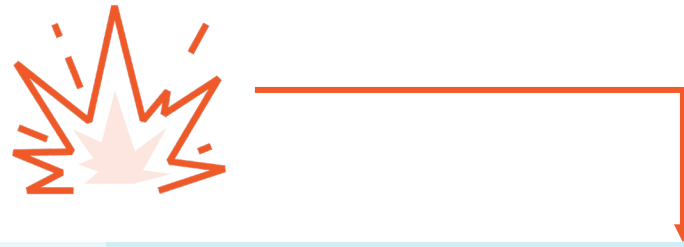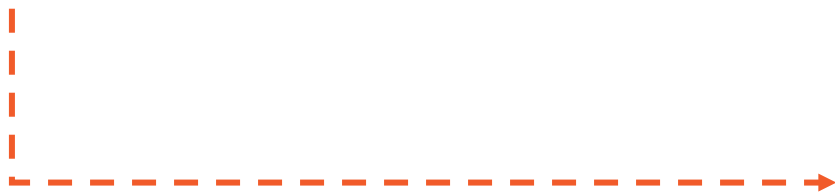
GLOBOTICKET
A Globomantics Company

- More investigation
- Transaction ID
- Rebuild user path

How do you track an increase in errors?

GLOBOTICKET
A Globomantics Company

- Standard log levels
- Logging dashboards
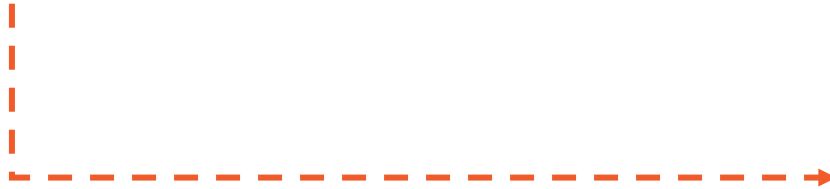- Alerts from monitoring

GLOBOTICKET
A Globomantics Company

- Metrics for logging
- Count by level
- Trigger alerts

Can you increase logging levels on the fly?

GLOBOTICKET
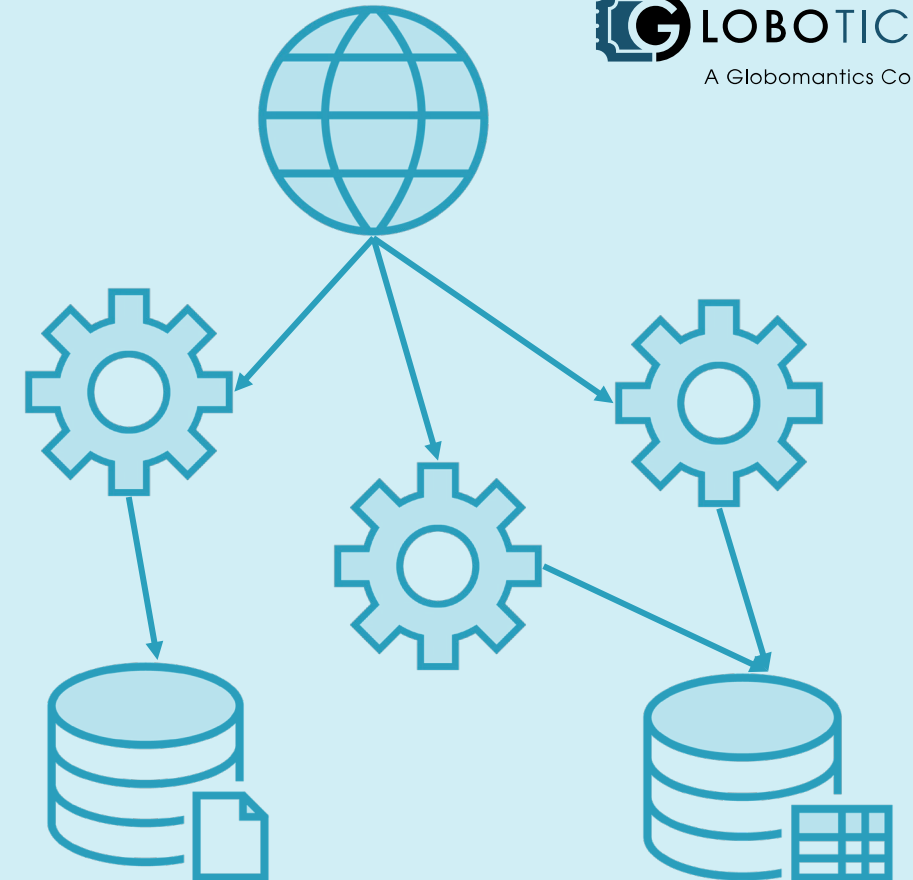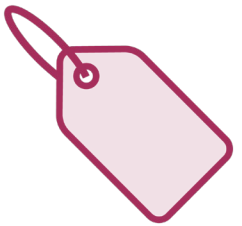A Globomantics Company

- Levels set by component
- Config rollout
- Instance restart

- Generate INFO logs
- Filter in pipeline
- Configure levels
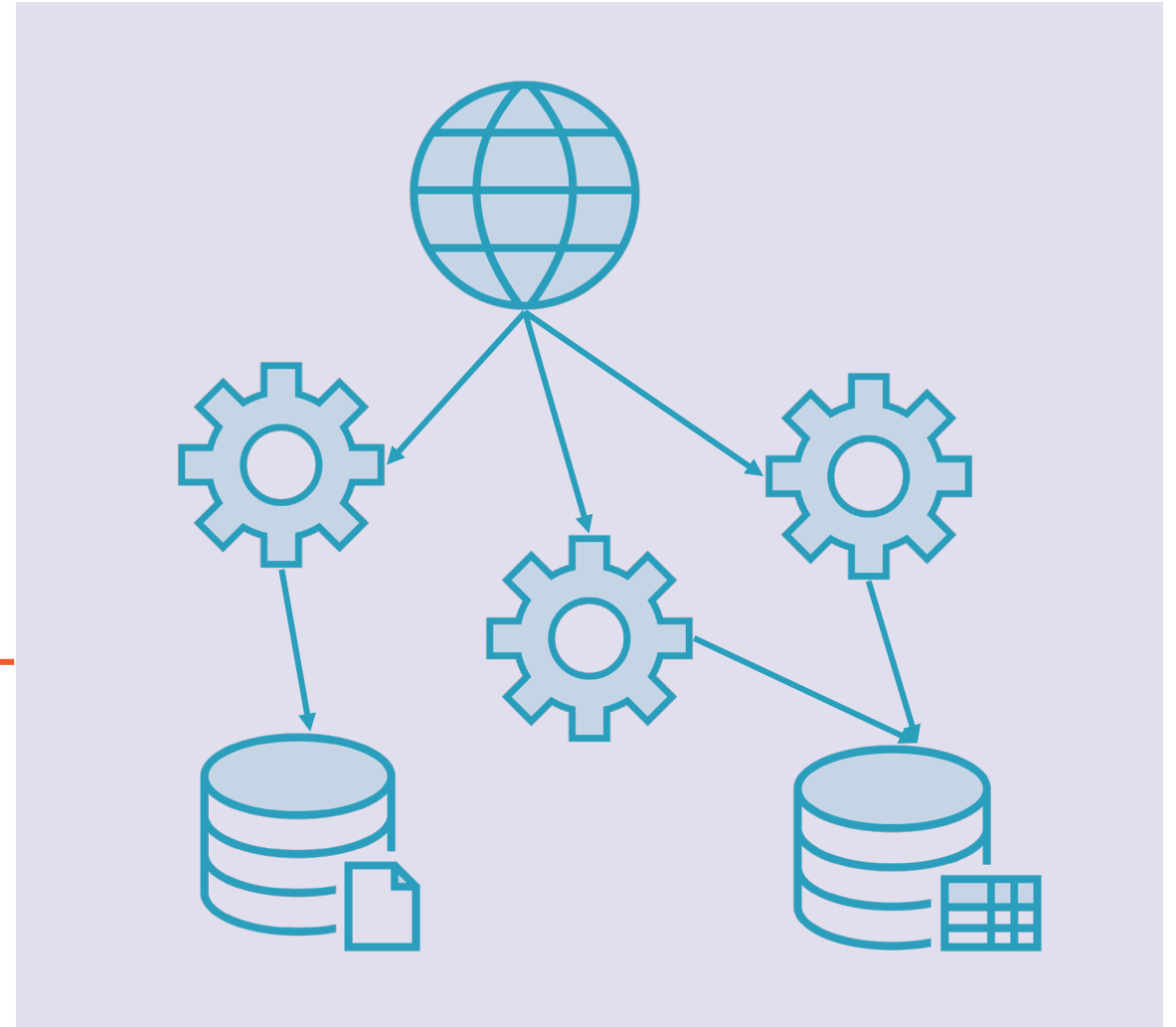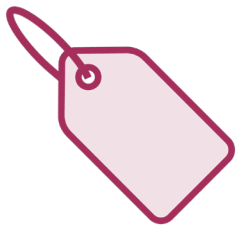
GLOBOTICKET
A Globomantics Company
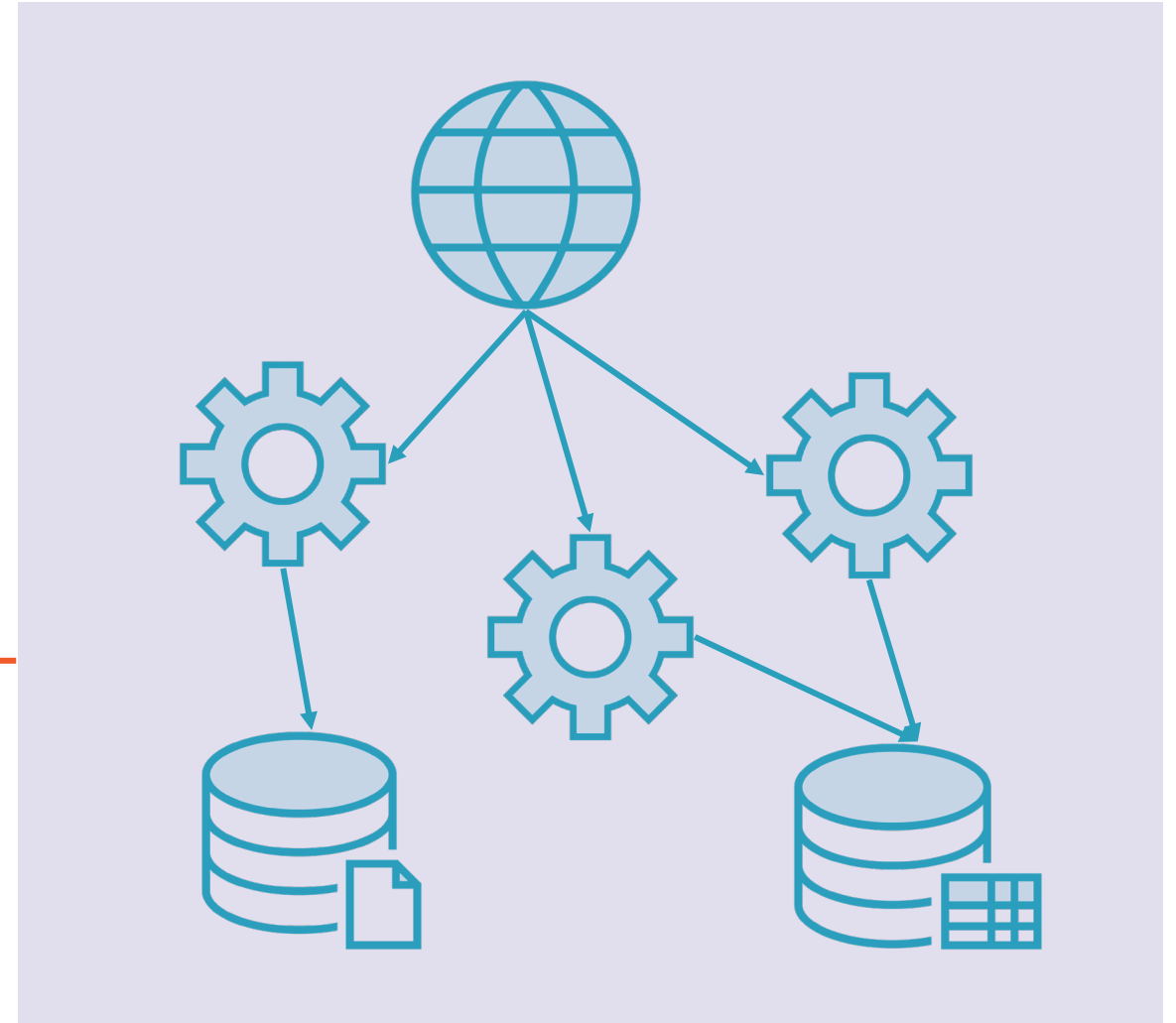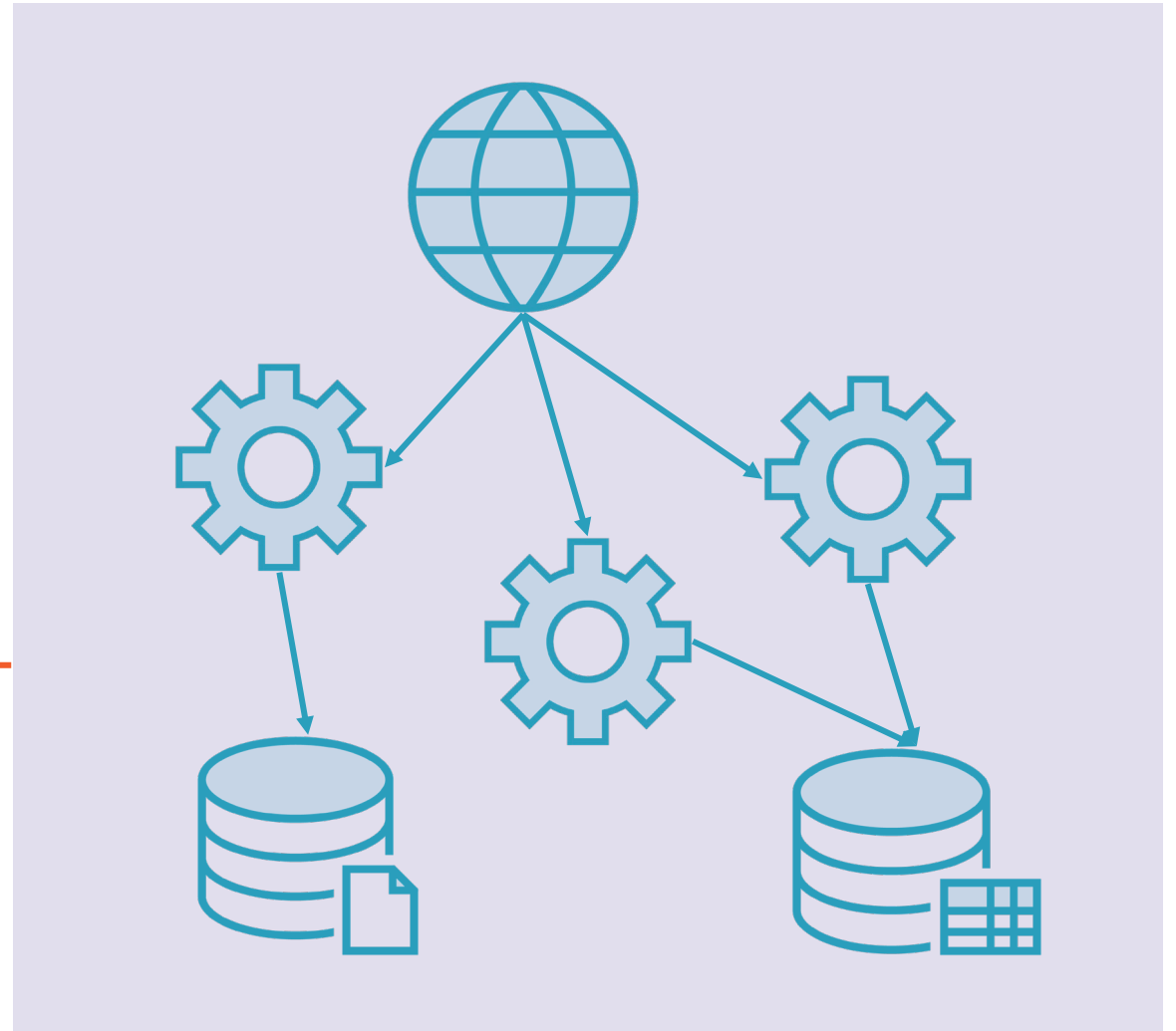
# Module Summary

Logging

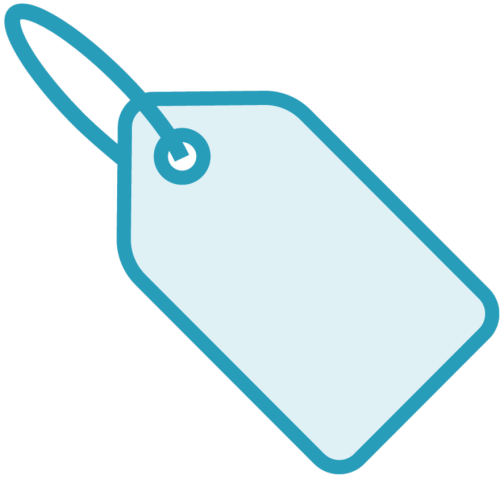- Region
- Instance ID
- Version
- Customer ID
- Transaction ID

Logging

- Timestamp
- Log level
- Component
- Log entry

# Designing Observability
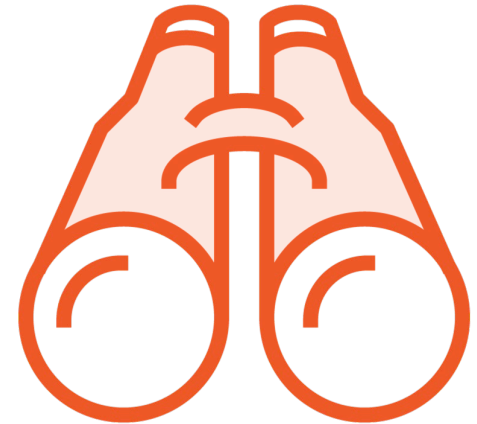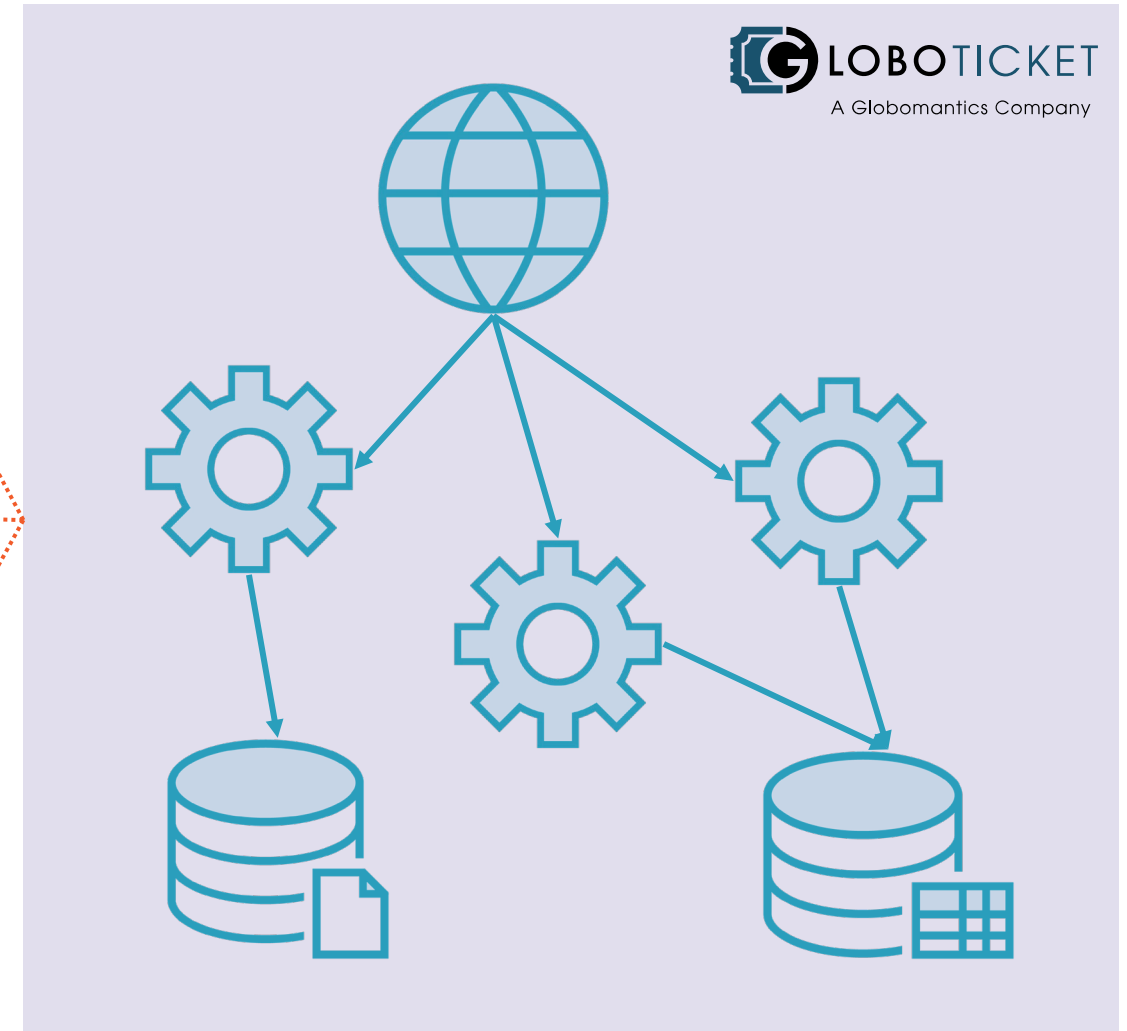
**Correlation**

**Real-time insight**

**Historical data**

Up Next:
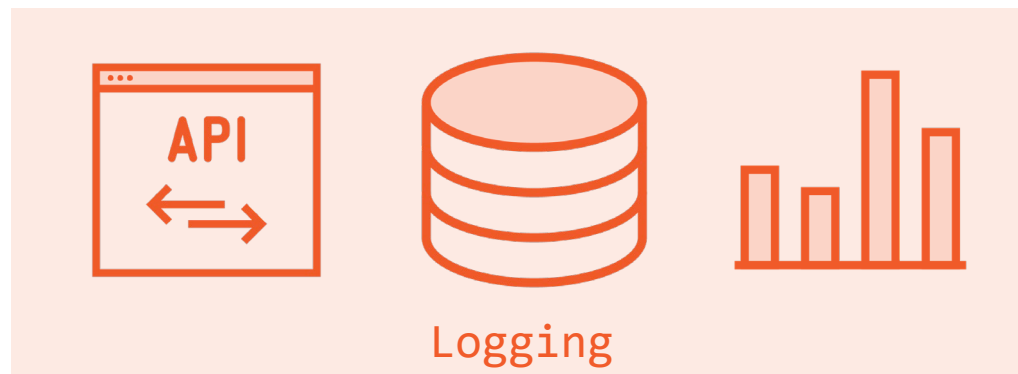Driving Continuous Improvement
with Service Levels