

DC/OS Day Two Operations

Day Two

Agenda:

Day Two

1. Monitoring and Metrics
2. Failure Handling
3. Production Checklist

Monitoring & Metrics

Gathering Health and Performance Metrics

Data-Driven Reporting

Monitoring and performance metrics help to answer the following questions:

- Are we close to cluster capacity?
- Are resources being used optimally?
- Is the system performing better or worse over time?
- Are there bottlenecks in the system?
- What is the response time of applications?

Metric Sources

Application Metrics

- QPS, latency, response time, hits, active users, errors

Container Metrics

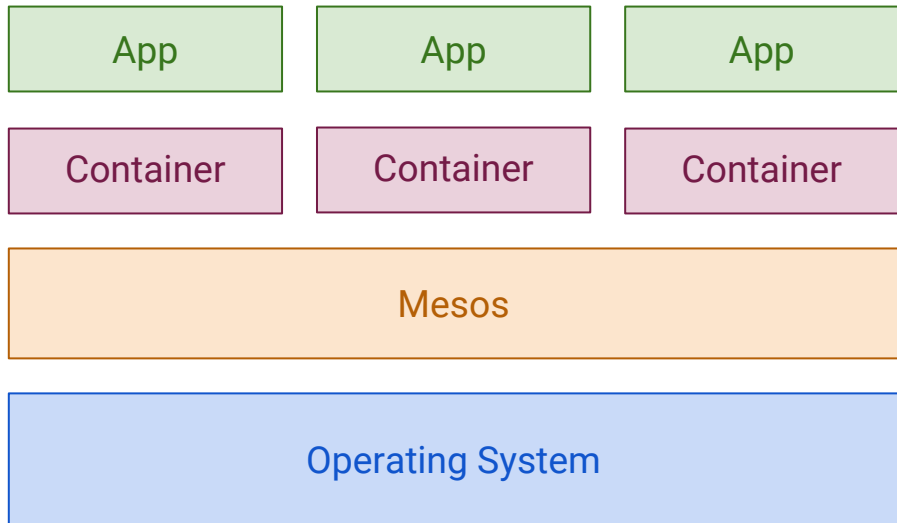
- CPU, memory, disk I/O, network I/O

Mesos Metrics

- Resources, frameworks, masters, agents, tasks, system, events

Host Metrics

- CPU, memory, disk I/O, network I/O



Data Sources

There are 4 components built-in to DC/OS that provide data that can be used for performance metrics and health monitoring:

1. DC/OS Metrics API (`dcos-metrics-master`, `dcos-metrics-agent`) exposes node, container, and application metrics
2. Marathon API (`dcos-marathon`) exposes metrics relating to Marathon's lifecycle, services, and connectivity with Mesos
3. Mesos APIs (`dcos-mesos-master`, `dcos-mesos-agent`) exposes performance and health metrics
4. DC/OS Diagnostics (`dcos-3dt`) aggregates and exposes health of DC/OS software components

Metrics API

- REST API that can retrieve metrics pertaining to nodes, containers, and applications running on the cluster
- Data can be sent to external services for visualization and alerting
- Runs on every node in the cluster
- Access to API is routed through `dcos-adminrouter` running on master node
- Master route:
 - `http(s)://<master_ip>/system/v1/metrics/v0/`
- Agent route:
 - `http(s)://<master_ip>/system/v1/agent/<agent_id>/metrics/v0/`
- [Documentation](#)

Metrics API: Benefits

Simplified configuration

- Automated collection of host and container level metrics
- Application level metrics integration via environment variables

Context injection

- Automated source tagging (agent IDs, container IDs, task IDs, etc.)
- Distributed aggregation

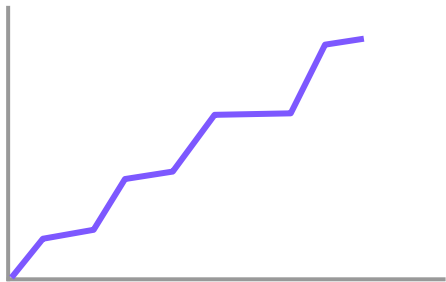
Collector per node

- Decoupled for faster upgrades and reconfigurations

Flexible output

- Kafka, logstash, Prometheus, etc.

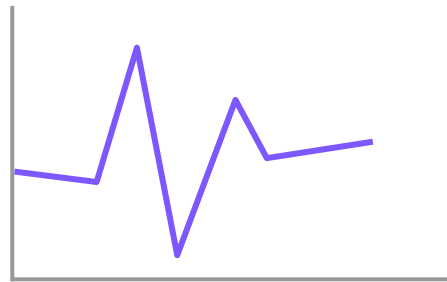
Metric Types



Counters

Discrete events that are monotonically increasing

- # of failed tasks
- # of agent registrations



Gauges

An instantaneous sample of some magnitude

- % of memory allocated in cluster
- # of connected slaves

Metrics API: Master Nodes

- Metrics for a master node are available at the following endpoint:
 - `http(s)://<master_IP>/system/v1/metrics/v0/node`
- Response is returned as a JSON document
- Metrics returned include the following measurements:
 - CPU
 - System load
 - Memory/Swap
 - Process count
 - Uptime
 - Filesystem usage
 - Network I/O

Metrics API: Agent Nodes

- Metrics for an agent node are available at the following endpoint:
 - `http(s)://<master_IP>/system/v1/agent/<agent_id>/metrics/v0/node`
- Response is returned as a JSON document
- Metrics returned include the following measurements:
 - CPU
 - System load
 - Memory/Swap
 - Process count
 - Uptime
 - Filesystem usage
 - Network I/O
 - Metadata

Metrics API: Node Level Measurements

Metric	Description
<code>cpu.{cores, idle, total}</code>	Percentage of cores used/idle/total
<code>load.{1, 5, 15}min</code>	Load average for the past 1, 5, and 15 minutes
<code>memory.{free, total}</code>	Amount of free/total memory in bytes
<code>processes</code>	Number of processes that are running
<code>swap.{free, used, total}</code>	Amount of swap free/used/total
<code>uptime</code>	Host uptime
<code>filesystem.{name}.capacity.{free, total, used}</code>	Amount of free/total/used storage capacity for filesystem
<code>network.{name}.in.{bytes, dropped, errors, packets}</code>	Incoming network metrics for a network interface
<code>network.{name}.out.{bytes, dropped, errors, packets}</code>	Outgoing network metrics for a network interface

Metrics API: Containers

- Metrics for a container on an agent are available at the following endpoint:
 - `http(s)://<master_IP>/system/v1/agent/<agent_id>/metrics/v0/node/container/<container_id>`
- Response is returned as a JSON document
- Metrics returned include the following measurements:
 - CPU
 - Memory
 - Disk
 - Network
 - Metadata

Metrics API: Container Level Measurements

Metric	Description
<code>cpus.limit</code>	Number of allocated CPUs
<code>cpus_system_time_secs</code>	Total CPU time spent in kernel space
<code>cpus_throttled_time_secs</code>	Total time container has been throttled by the host kernel
<code>cpus_user_time_secs</code>	Total CPU time spent in userspace
<code>disk_limit_bytes</code>	Hard limit for disk
<code>disk_used_bytes</code>	Disk consumption
<code>mem_limit_bytes</code>	Hard RAM limit
<code>mem_total_bytes</code>	RAM consumption
<code>net_rx_{bytes,dropped,errors,packets}</code>	Incoming network metrics
<code>net_tx_{bytes,dropped,errors,packets}</code>	Outgoing network metrics

Metrics API: Application Measurements

- Application metrics can be published to the dcos-metrics component through environment variables
 - `STATSD_UDP_HOST`: Initialized to the IP address of the agent where the container is running
 - `STATSD_UDP_PORT`: Initialized to the bind port for the statsd server on the agent
- Metrics are appended with metadata such as:
 - `agent_id`
 - `container_id`
 - `framework_id`
 - `hostname`
 - `labels`
- Only works with UCR containers

Metrics API: `dcos` CLI

- Node and container level metrics can be retrieved through the `dcos` CLI
- Node level metrics:
 - Retrieve the ID of an agent node: `dcos node`
 - Retrieve a summary of metrics for a node: `dcos node metrics summary <ID>`
 - Retrieve a detailed report of each metric for a node: `dcos node metrics details <ID>`
- Container level metrics:
 - Retrieve the ID of a task: `dcos task`
 - Retrieve a summary of metrics for a task: `dcos task metrics summary <ID>`
 - Retrieve a detailed report of each metric for a task: `dcos node task details <ID>`

Lab 6

Metrics API

Marathon Metrics

- Metrics for Marathon are available at the following endpoints:
 - `http://<master_IP>/marathon/metrics`
 - `http://<master_IP>/marathon/v2/apps`
- Metrics can be automatically exported to Graphite or Datadog via Marathon start up options:
 - `--reporter_graphite <graphite_server>`
 - `--reporter_datadog <datadog_server>`

Marathon Service Performance (Step 1 of 2)

```
curl http://<master_IP>/marathon/v2/apps/hello-world
```

```
{
  "appId": "/hello-world",
  "healthCheckResults": [],
  "host": "10.0.0.99",
  "id": "hello-world.57911772-47d6-11e8-bc94-821c2bfac18f",
  "ipAddresses": [
    {
      "ipAddress": "10.0.0.99",
      "protocol": "IPv4"
    }
  ],
  "ports": [
    3969
  ],
  "servicePorts": [],
  "slaveId": "729d2884-b78f-4aee-8ad4-3cec73d05a4b-S2",
  "state": "TASK_RUNNING",
  "stagedAt": "2018-04-24T15:44:15.251Z",
  "startedAt": "2018-04-24T15:44:16.458Z",
  "version": "2018-04-24T15:44:15.192Z",
  "localVolumes": []
}
```

Marathon Service Performance (Step 2 of 2)

```
curl http://<agent_IP>:5051/monitor/statistics
```

```
{
  "executor_id": "hello-world.57907b30",
  "executor_name": "Command Executor (Task: hello-world.57907b30) (Command: sh -c 'while [ true...')",
  "framework_id": "729d2884-b78f-4aee-8ad4-3cec73d05a4b-0001",
  "source": "hello-world.57907b30-47d6-11e8-bc94-821c2bfac18f",
  "cpus_limit": 0.2,
  "cpus_system_time_secs": 41.39,
  "cpus_throttled_time_secs": 0.386315377,
  "cpus_user_time_secs": 23.84,
  "mem_limit_bytes": 44040192,
  "mem_rss_bytes": 8445952,
  "mem_swap_bytes": 0,
  "mem_total_bytes": 8445952,
  "mem_unevictable_bytes": 0,
  "timestamp": 1524687903.8652
  ...
}
```

Alert Metrics for Mesos

Metric	Inference
<code>master/uptime_secs</code> is low	The master has restarted
<code>master/uptime_secs</code> < 60 for sustained periods of time	The cluster has a flapping master node
<code>master/tasks_lost</code> is increasing rapidly	Tasks in the cluster are disappearing. Possible causes include hardware failures, bugs in one of the frameworks, or bugs in Mesos
<code>master/slaves_active</code> is low	Slaves are having trouble connecting to the master
<code>master/cpus_percent</code> > 0.9 for sustained periods of time	Cluster CPU allocation is close to capacity
<code>master/mem_percent</code> > 0.9 for sustained periods of time	Cluster RAM allocation is close to capacity
<code>master/disk_percent</code> > 0.9 for sustained periods of time	Cluster disk allocation is close to capacity
<code>master/elected</code> is 0 for sustained period of time	No master is currently elected

Lab 7


Marathon & Mesos Metrics

Lab until :30

DC/OS Components Health

- DC/OS is packaged with a component called DC/OS Distributed Diagnostics Tool (3DT) which collects and aggregates the health of the DC/OS software components running on your cluster nodes
- 3DT runs on every node in the cluster, monitoring the state of `systemd` units
- Workflow:
 - 3DT running on leading master retrieves a list of all masters and agents
 - 3DT running on leading master queries 3DT on all other masters and agents and collects their respective states
 - 3DT aggregates the data and is exposed by dcos-adminrouter at the `/system/health/v1/` route
- Can also be queried via the DC/OS CLI to generate and download a diagnostics bundle that can be sent to the Mesosphere support team

DC/OS Components Health: GUI

 Components

52 Components

All 52

Healthy 52

Unhealthy 0

Download Snapshot

Name	Health ▾
Admin Router Agent	Healthy
Admin Router Master	Healthy
CockroachDB	Healthy
CockroachDB Disable Diagnostics	Healthy
CockroachDB Set Configuration Timer	Healthy
DC/OS Backup Master	Healthy
DC/OS Backup Master Socket	Healthy
DC/OS Certificate Authority	Healthy
DC/OS Checks Timer	Healthy
DC/OS Cluster Linker Service	Healthy
DC/OS Cluster Linker Socket	Healthy

DC/OS Components Health: GUI

 Components > Admin Router Agent (Healthy)

3 Health Checks

 Filter

All Health Checks ▼

Health ▼	Node	Role
Healthy	10.0.0.11	Agent_public
Healthy	10.0.0.242	Agent
Healthy	10.0.0.99	Agent

Lab 8

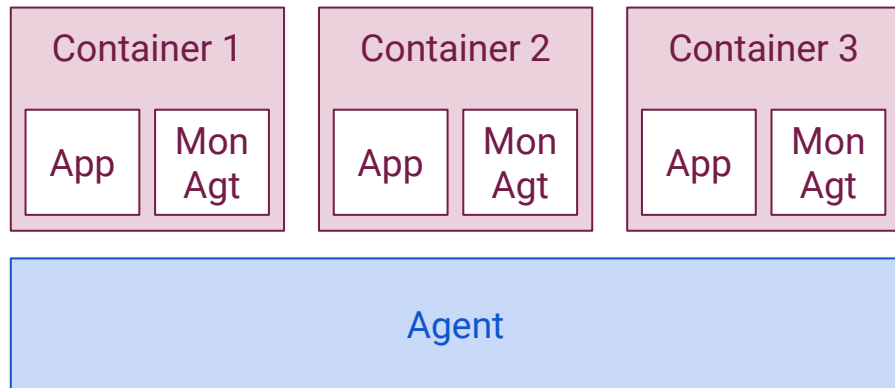
DC/OS Health Monitoring

Lab until :00

Container Level Metrics

Monitoring agent per container?

- Not scalable
- Increased footprint
- Higher likelihood of user error



High Level Workflow

DC/OS clusters subscribe to an event-driven monitoring model

Metric Collection

Host, container, and application level metrics are collected from the systems

Metric Routing

Metrics are fed to a central routing system which contains rules indicating where metrics are destined to go

Metric Storage

Metrics are persisted in a storage system

Metric Visualized

Metrics are plotted out in a graphing tool to make consumption of information more human-friendly

Step 1: Metric Collection

- A metric collector is a process that runs on a node as either a OS level daemon or as a container
- Responsible for collecting metrics related to:
 - The host itself
 - All the containers running on the host
 - The applications running within the containers
- Upon capturing metrics, the agent sends the data points to an event router
- Examples:
 - `collectd`
 - `cAdvisor`
 - `statsd`

Container Advisor (cAdvisor)

- Open source and created by Google
- Process that runs on a host that collects, aggregates, processes, and exports information about running containers
- Collects metrics such as isolation parameters, resource usage
- Supports both Docker and UCR containers
- Data is available through a REST API
- Data can automatically be exported to an external storage system such as InfluxDB, Elasticsearch, Kafka
- Can be installed through the DC/OS Package Catalog
- [Documentation](#)

Step 2: Event Router

- Receives metrics from the metric collector
- Central dispatcher which makes metrics consumable by other systems
- Examples:
 - fluentd
 - Kafka
 - Logstash
 - cAdvisor

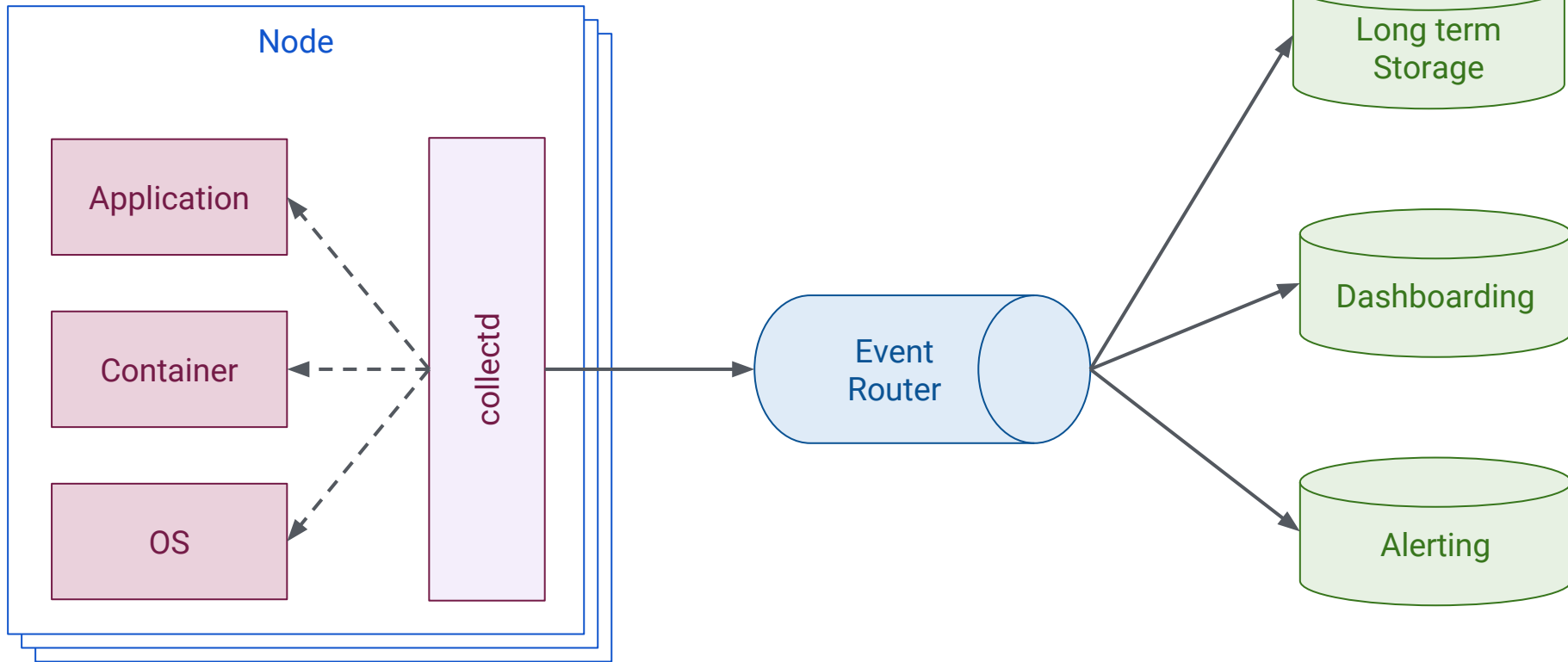
Step 3: Metric Storage

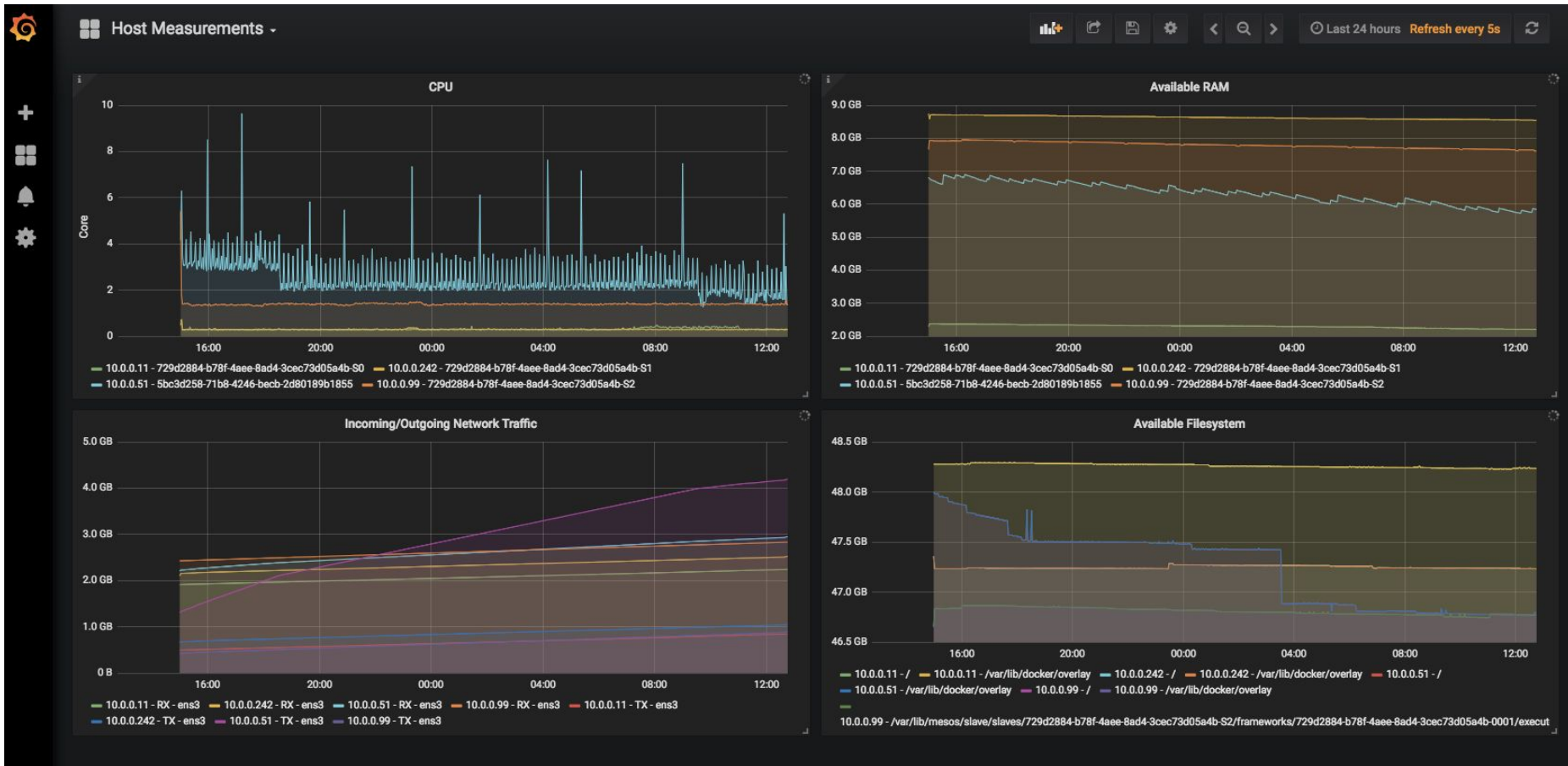
- Provides storage system to persist collected metrics
- Can have three forms from short-term/low-latency systems to long-term/persistent systems:
 - Alerting: If a metric has a certain unhealthy trend
 - Dashboarding: Plot the last hour in high resolution and a summary of the past 24hrs
 - Long Term: For legal and auditing obligations
- Examples:
 - Elasticsearch
 - Graphite
 - InfluxDB
 - Prometheus
 - Cassandra
 - Filesystem (local, CephFS, HDFS, etc.)

Step 4: Metric Visualization

- Provides customizable dashboarding
- Simplifies digestion of collected metrics
- Examples:
 - Grafana
 - D3
 - SignalFX

Putting it All Together





Lab 9

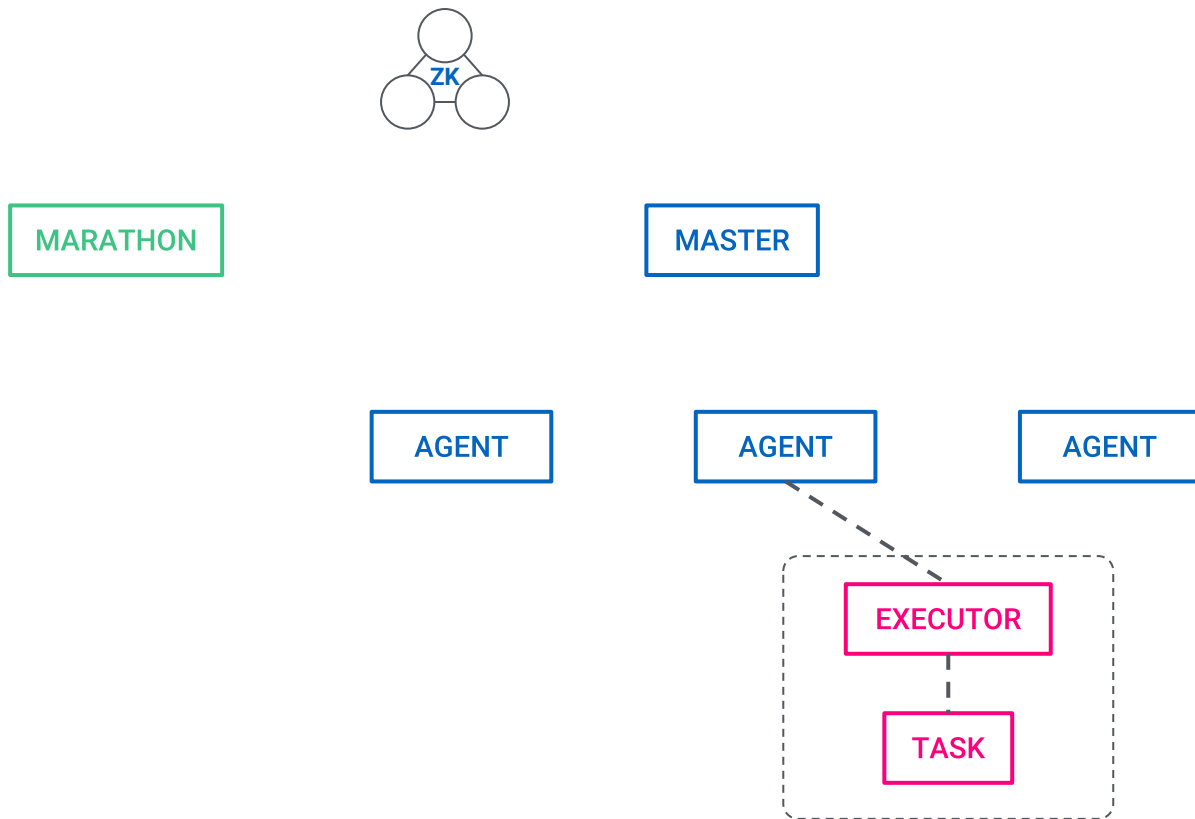
Monitoring DC/OS with cAdvisor,
InfluxDB, and Grafana

Lab until :35

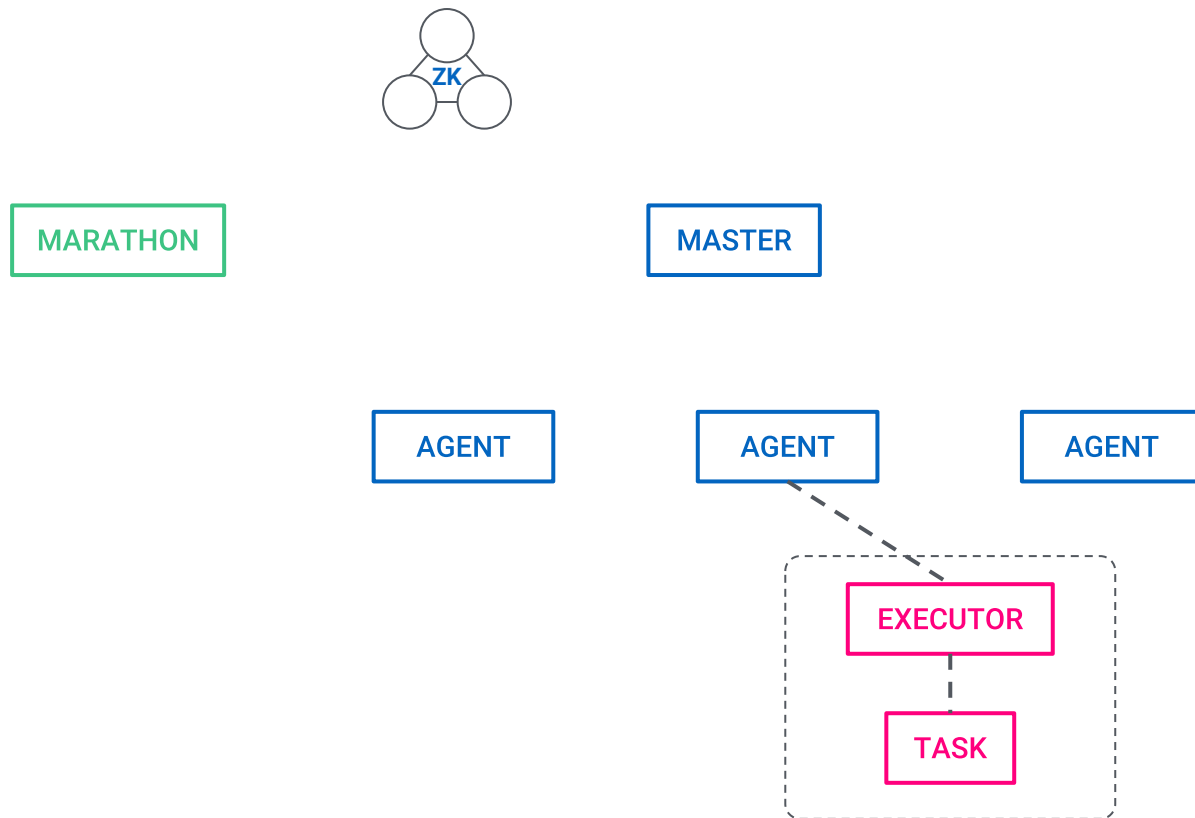
Cluster Failures

Task (Container) Failure

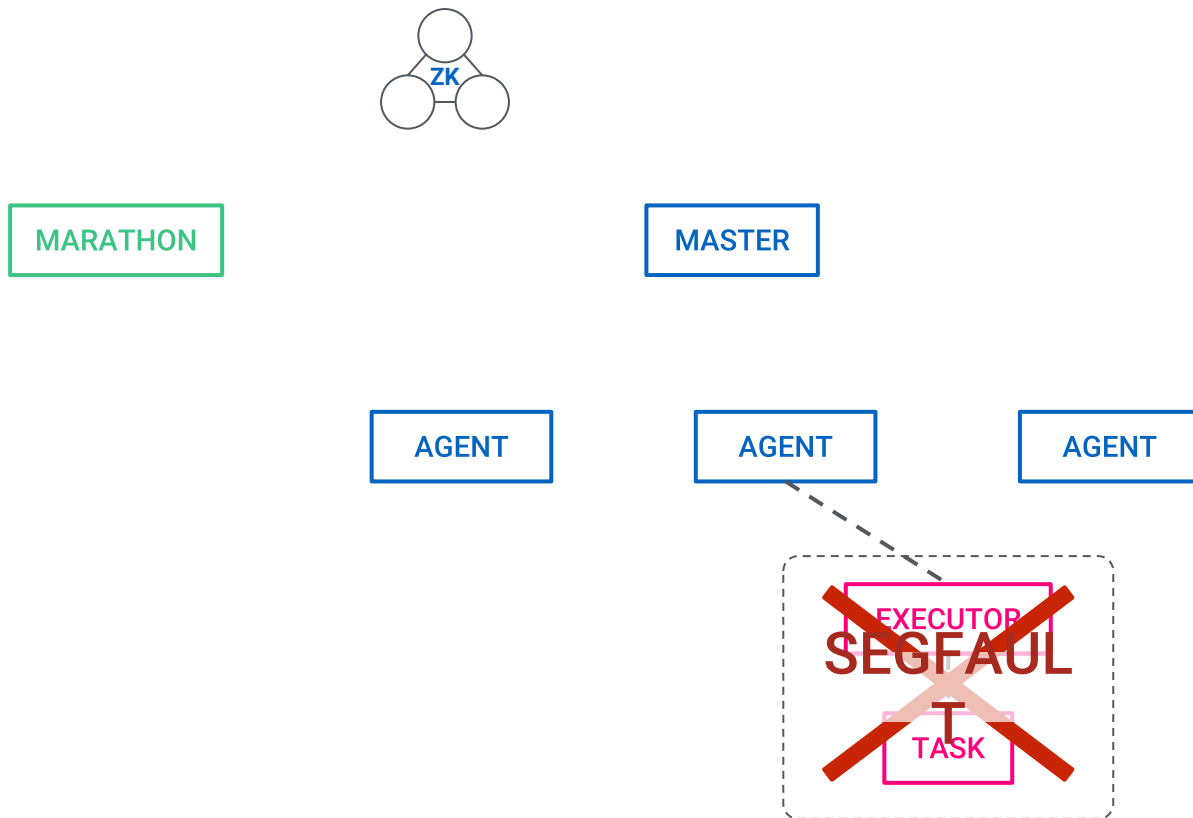
Task Failure



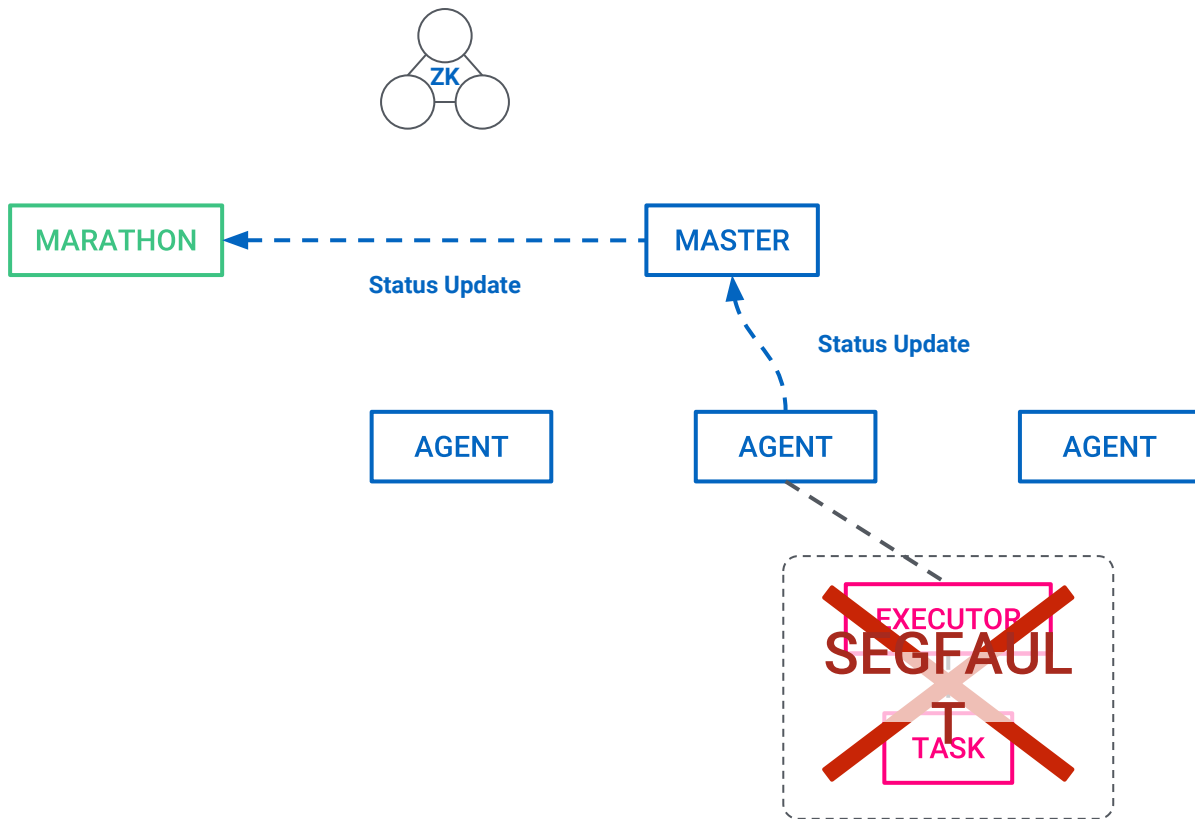
Task Failure



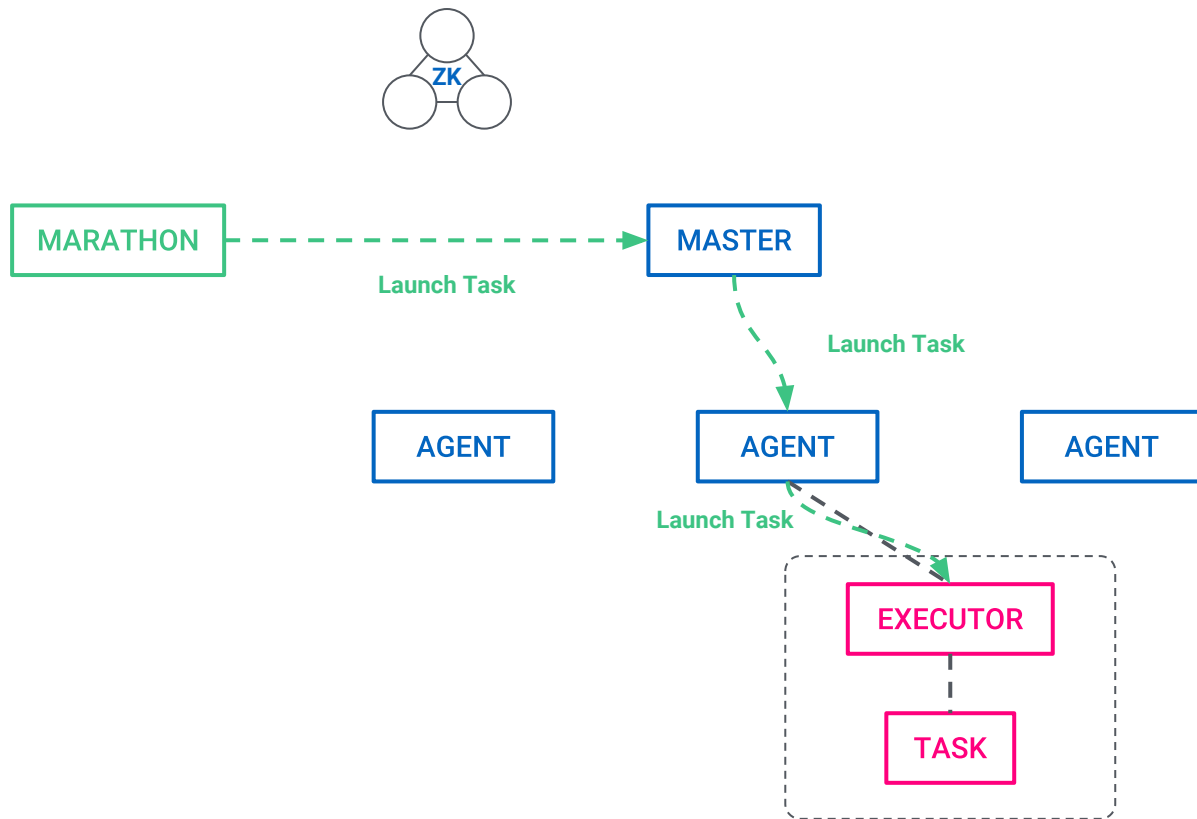
Task Failure



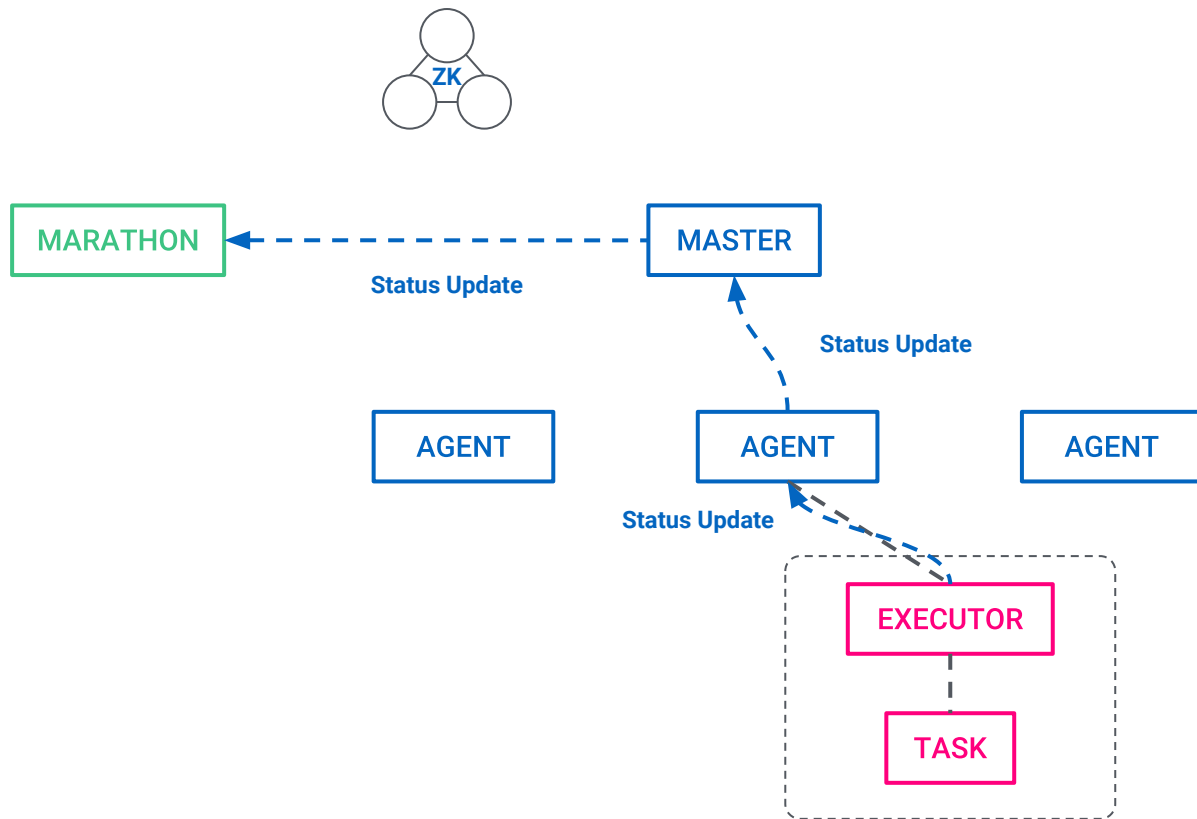
Task Failure



Task Failure

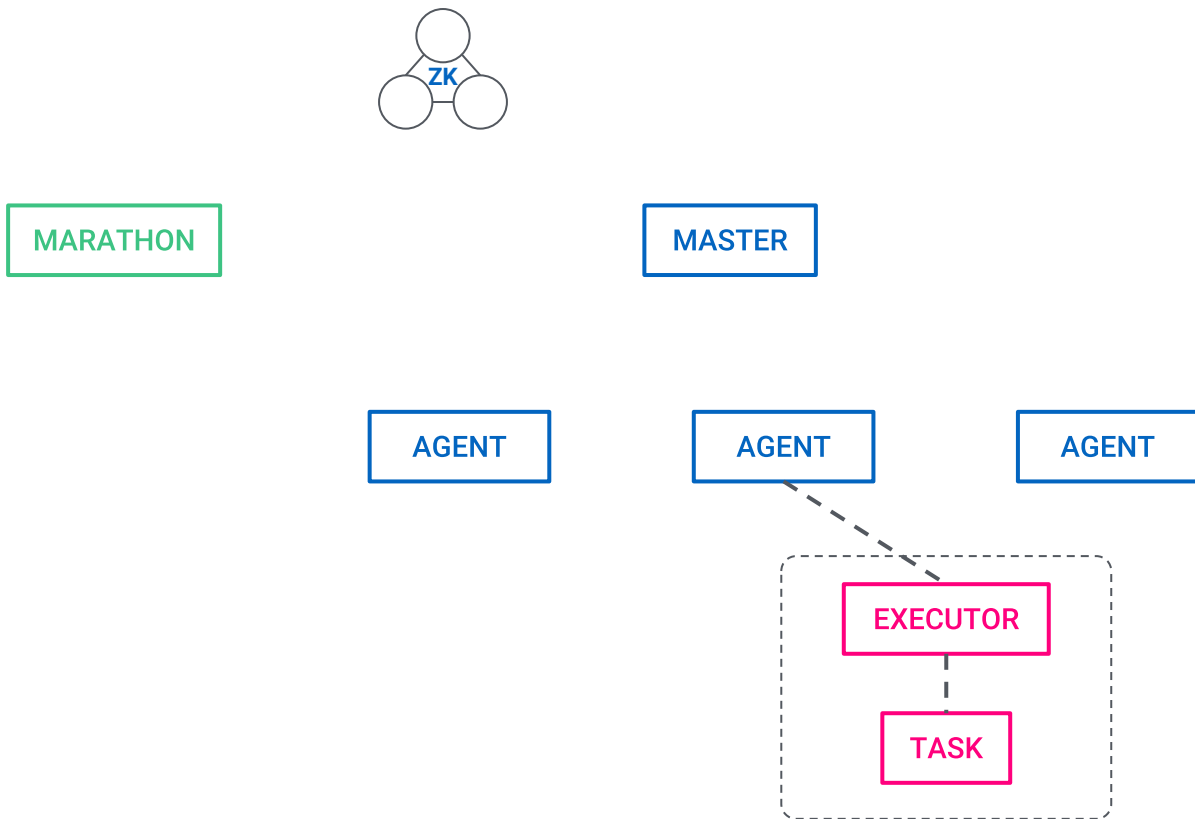


Task Failure

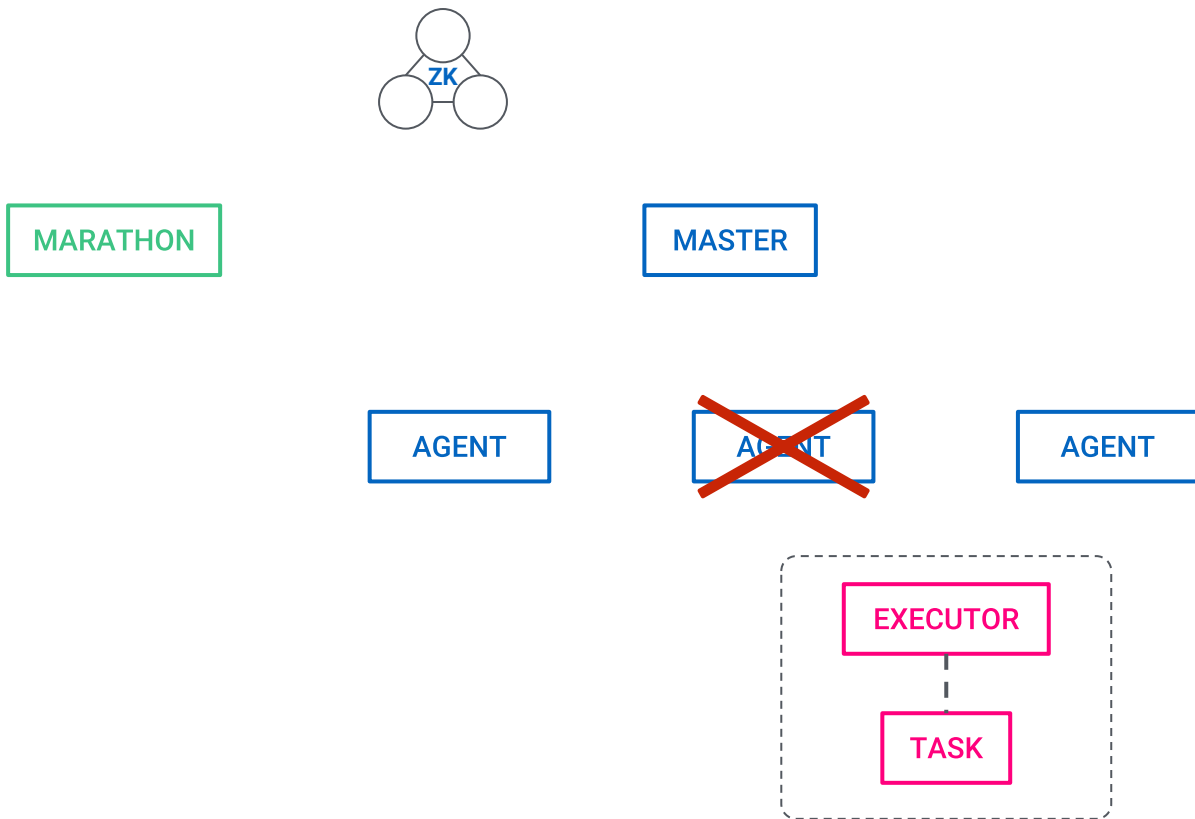


Mesos Agent Failure

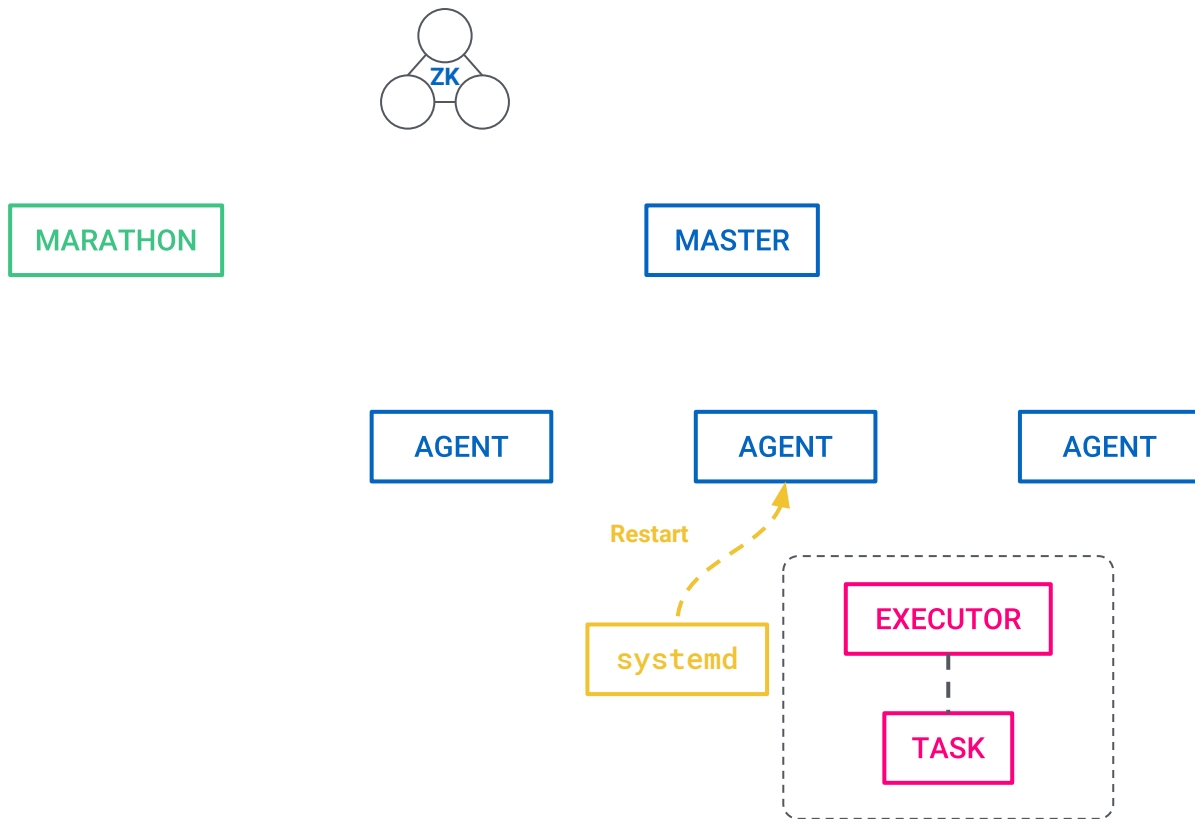
Mesos Agent Failure



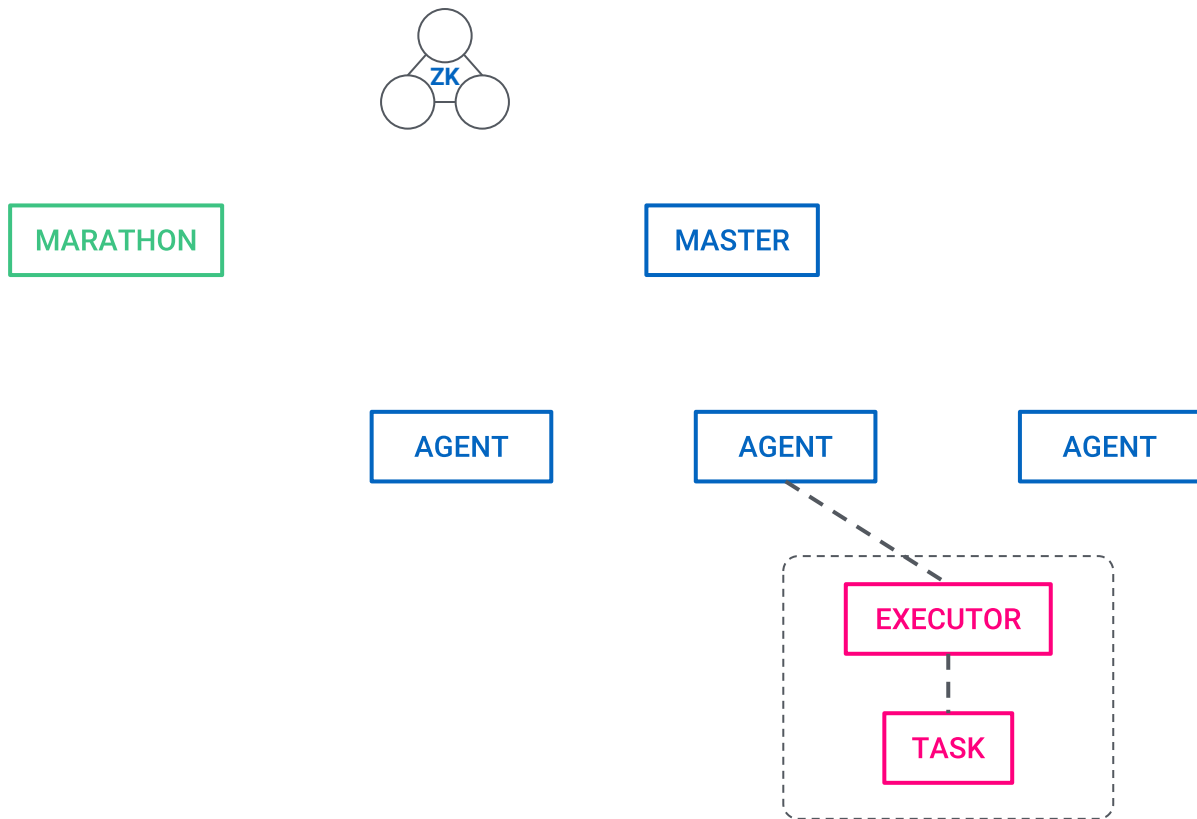
Mesos Agent Failure



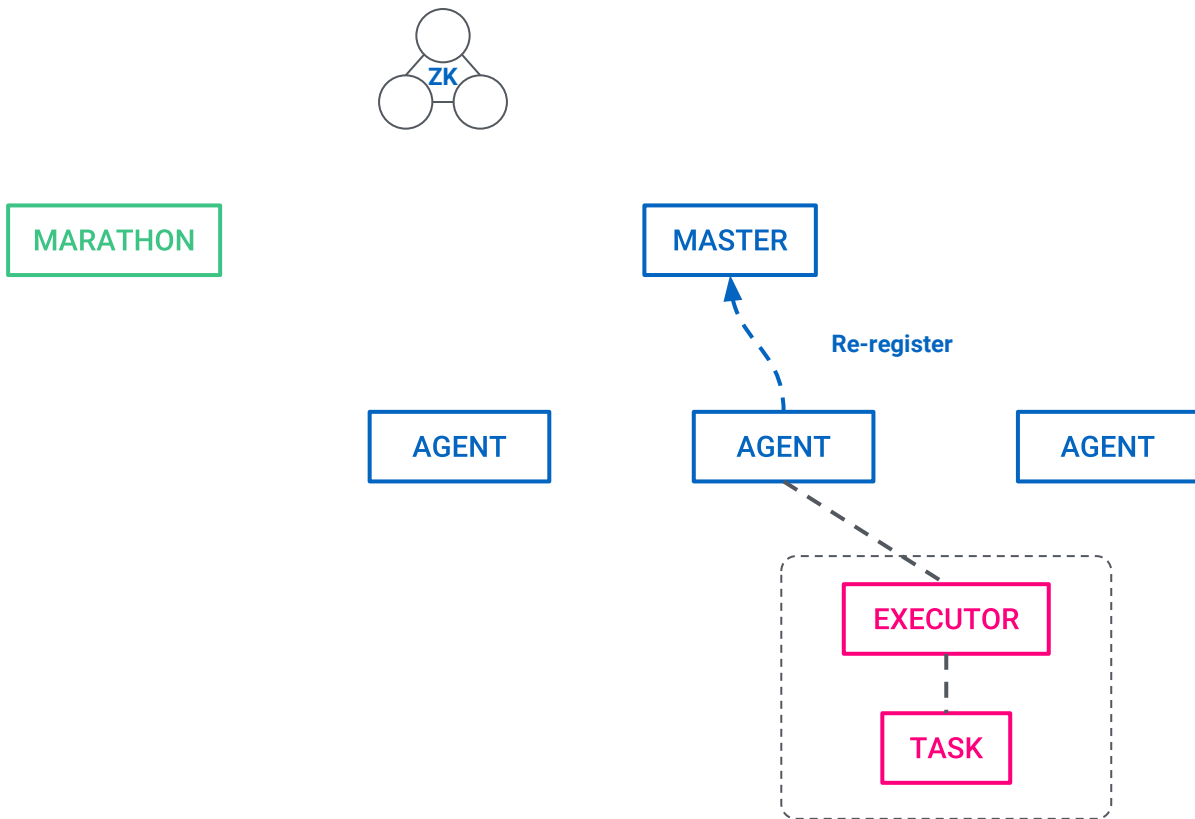
Mesos Agent Failure



Mesos Agent Failure

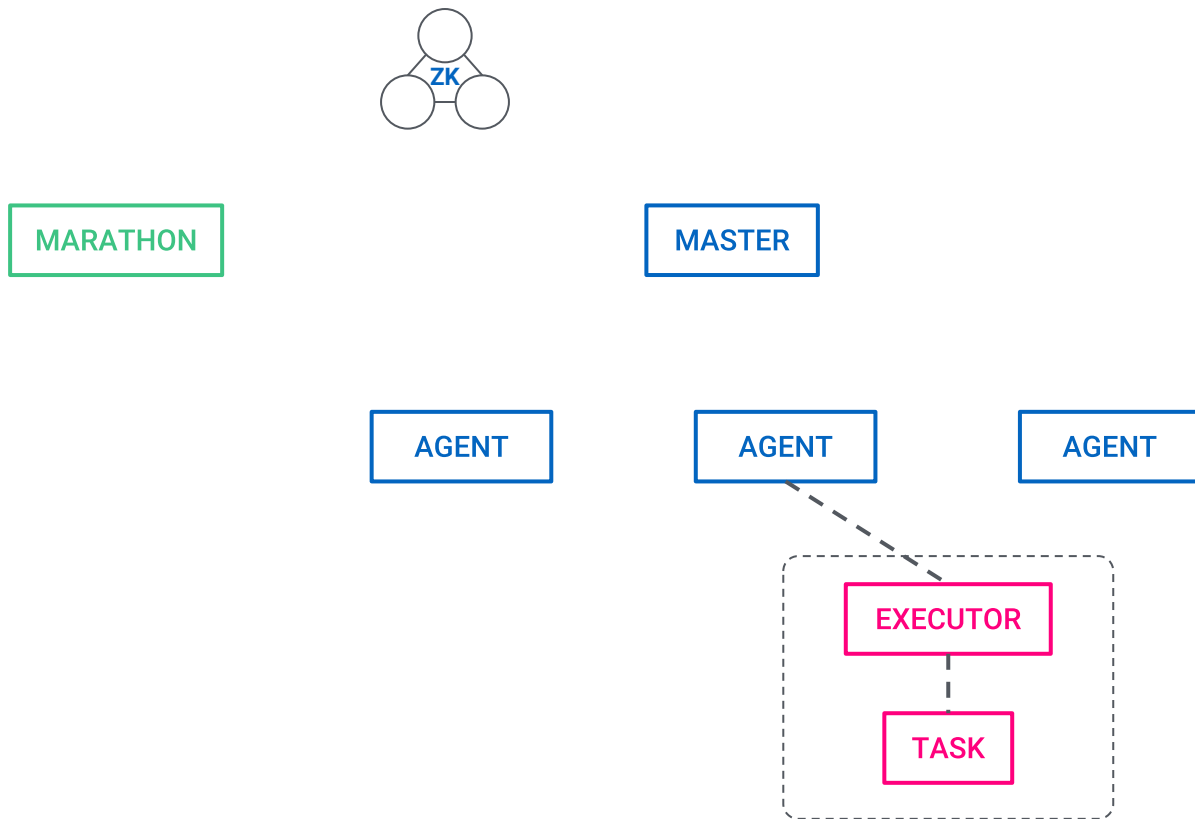


Mesos Agent Failure

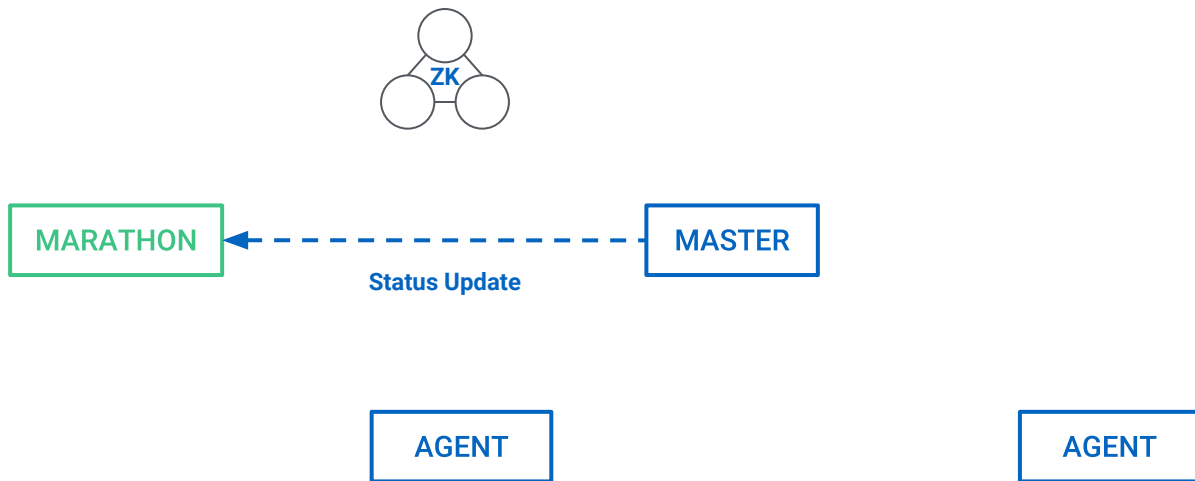


Agent Host Failure

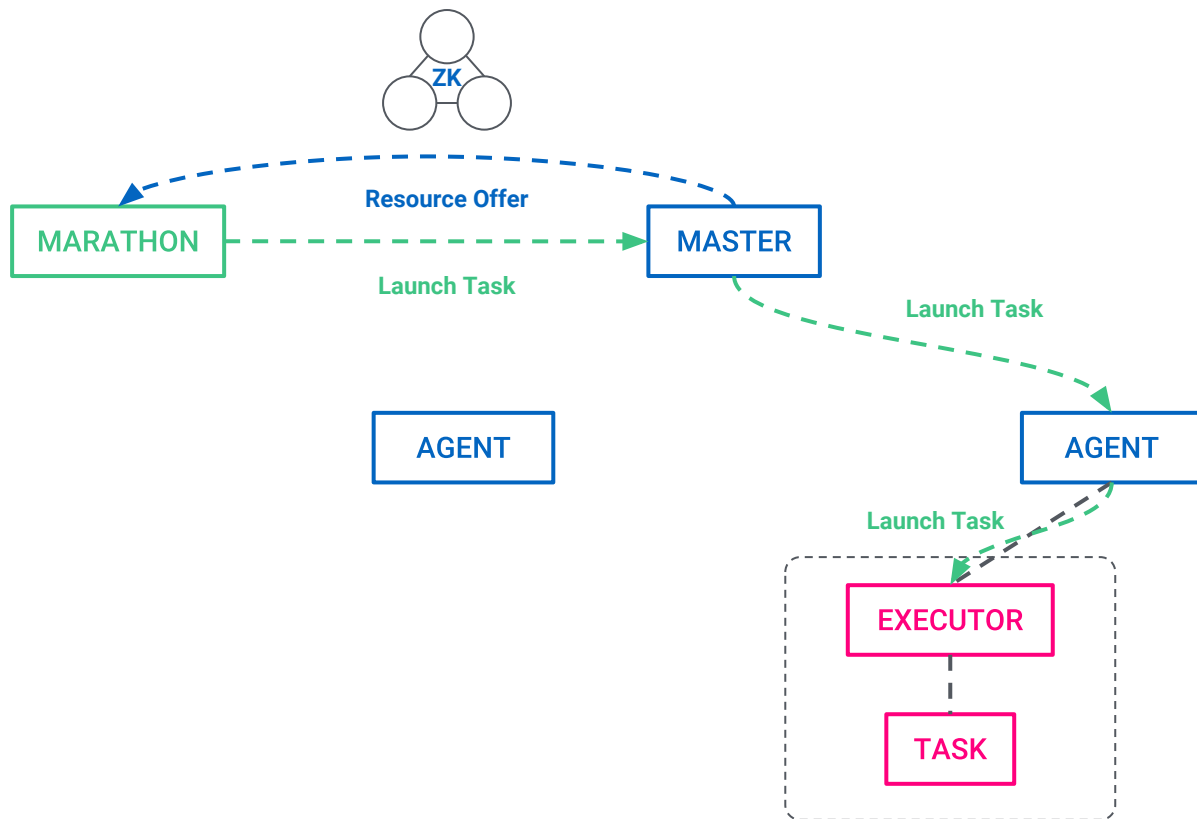
Agent Host Failure



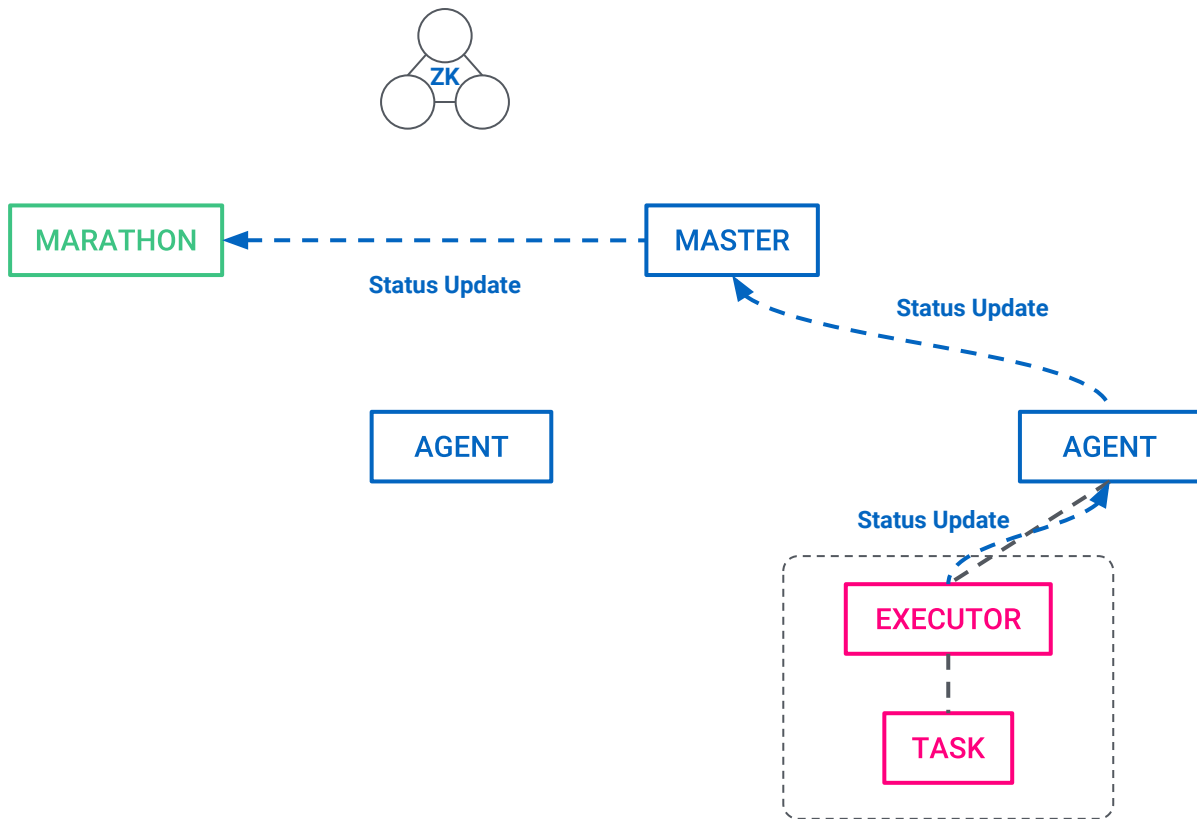
Agent Host Failure



Agent Host Failure

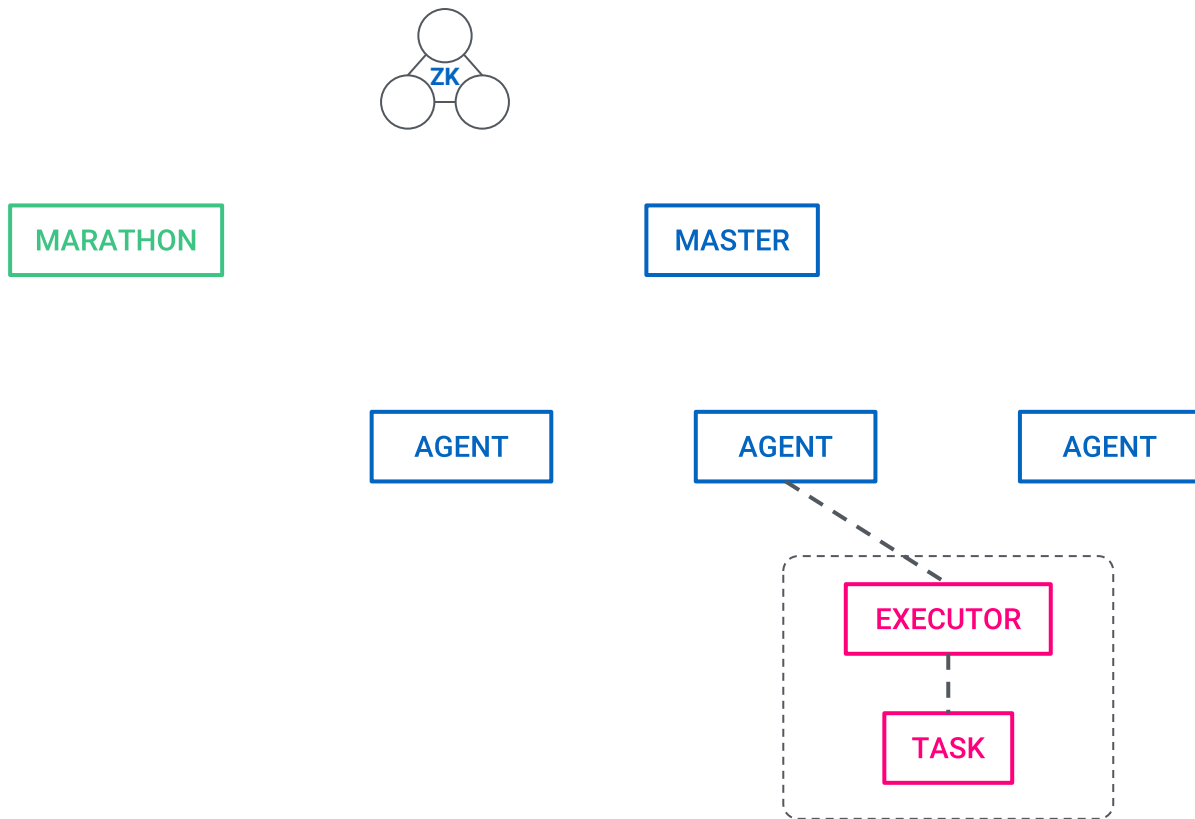


Agent Host Failure

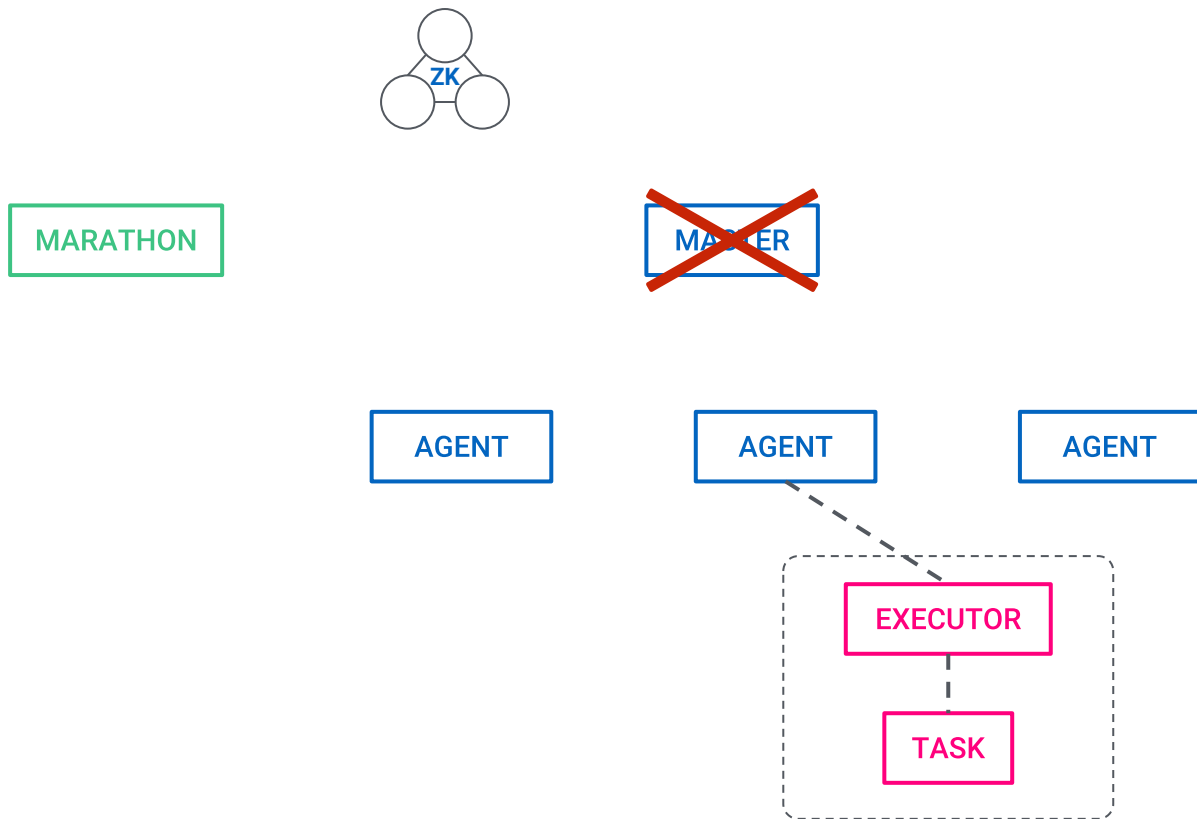


Mesos Master Failure

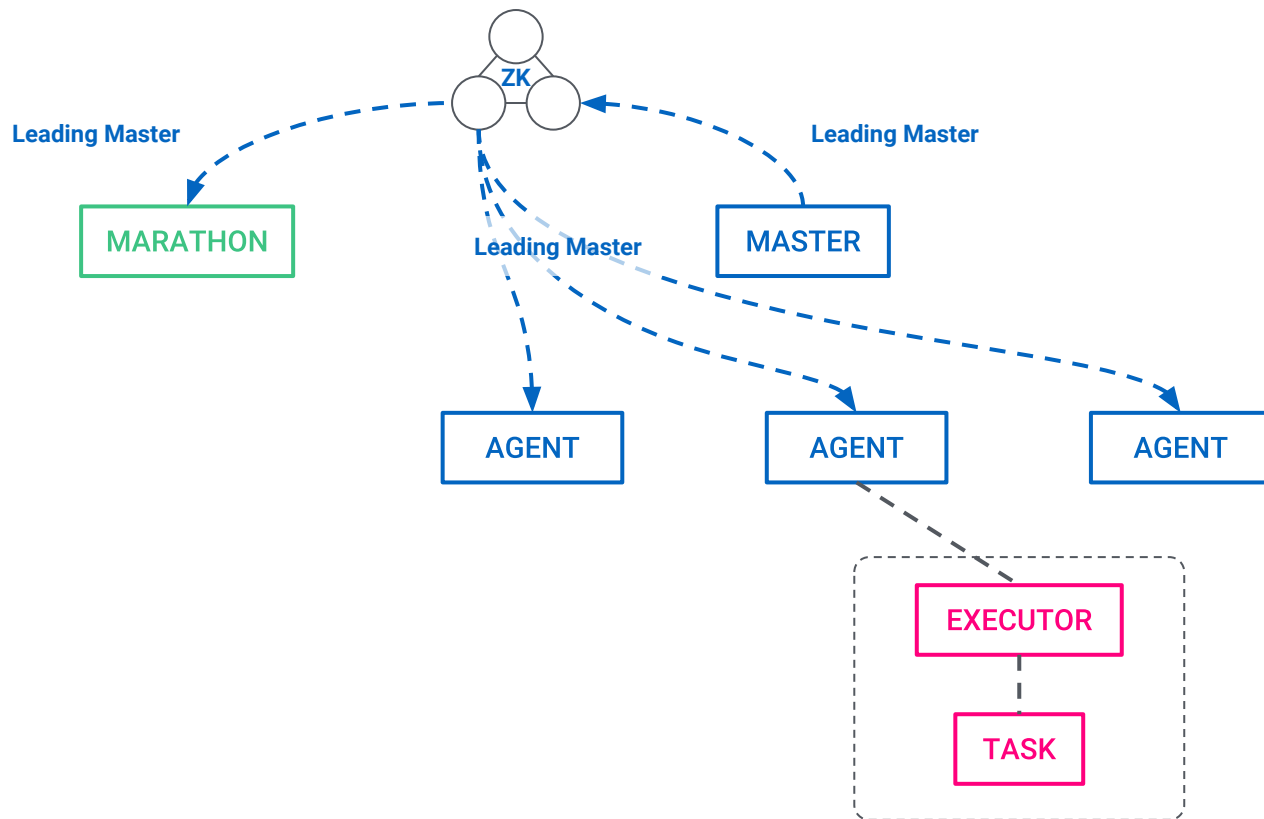
Mesos Master Failure



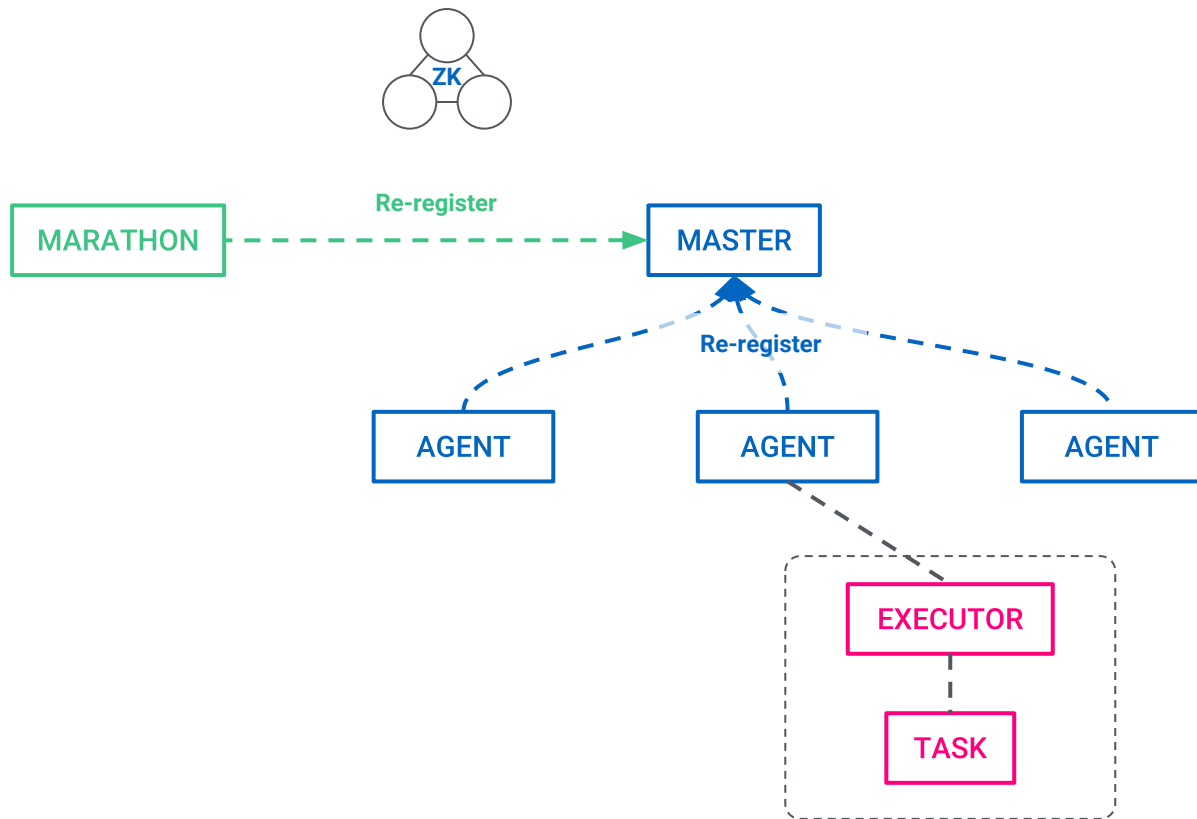
Mesos Master Failure



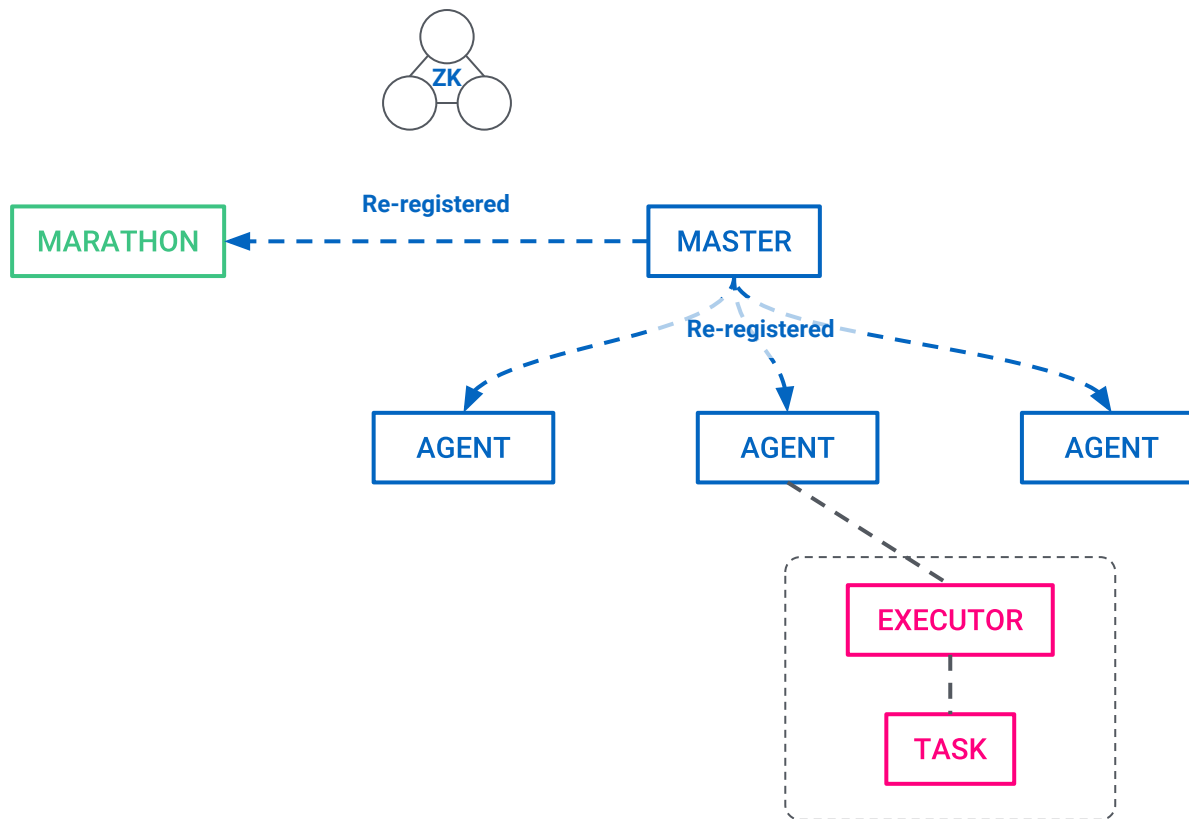
Mesos Master Failure



Mesos Master Failure

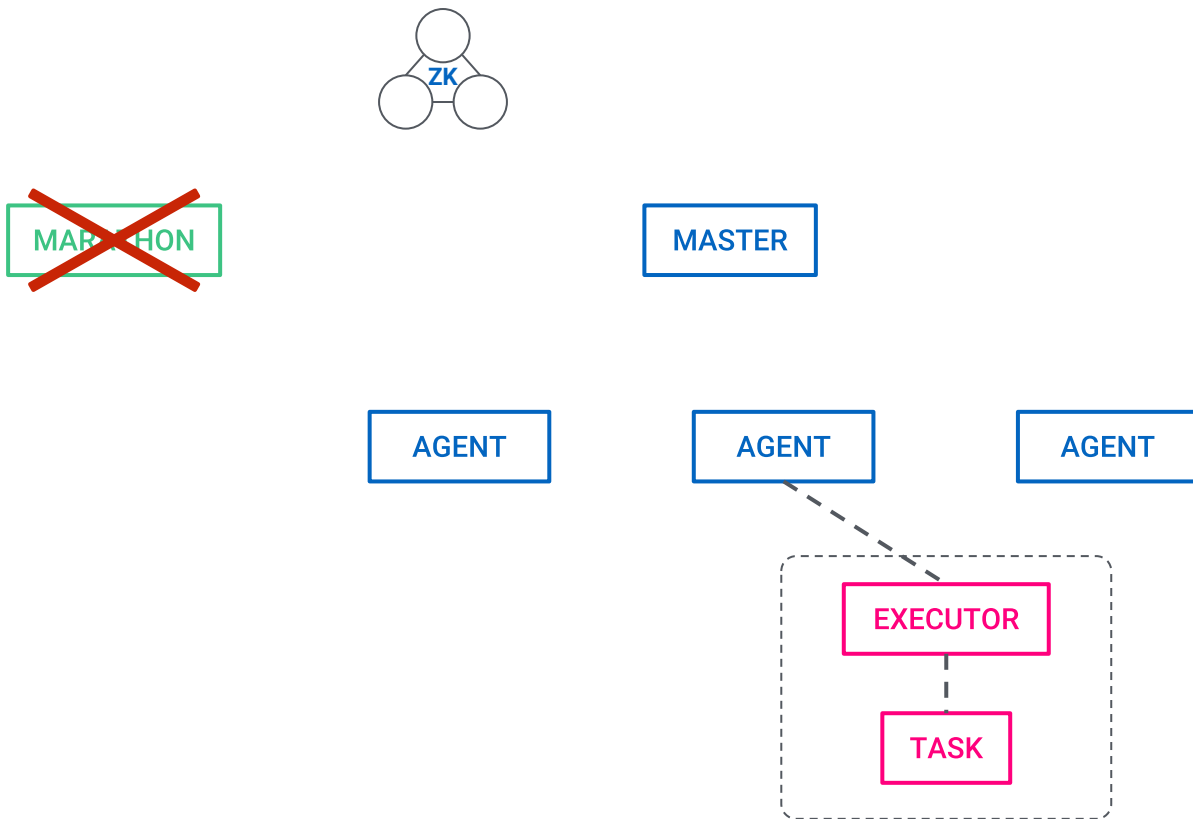


Mesos Master Failure

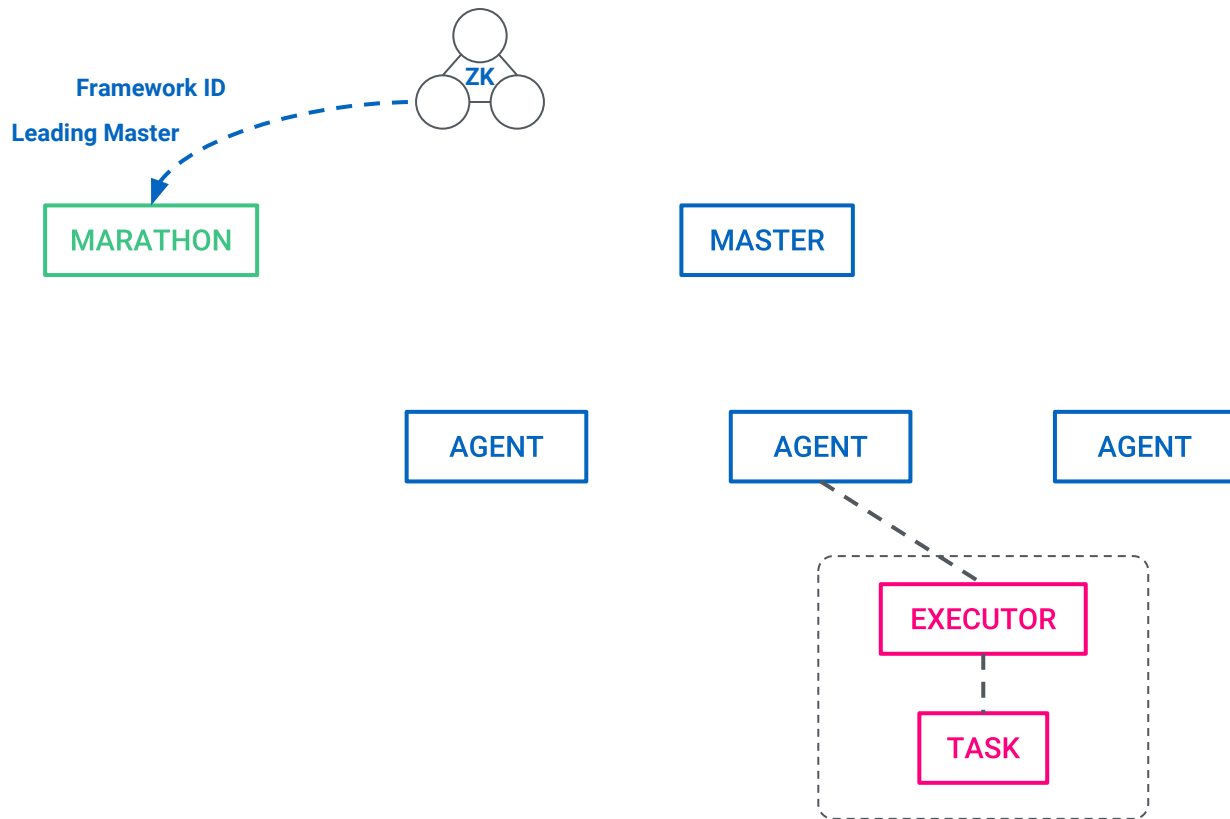


Scheduler Failure

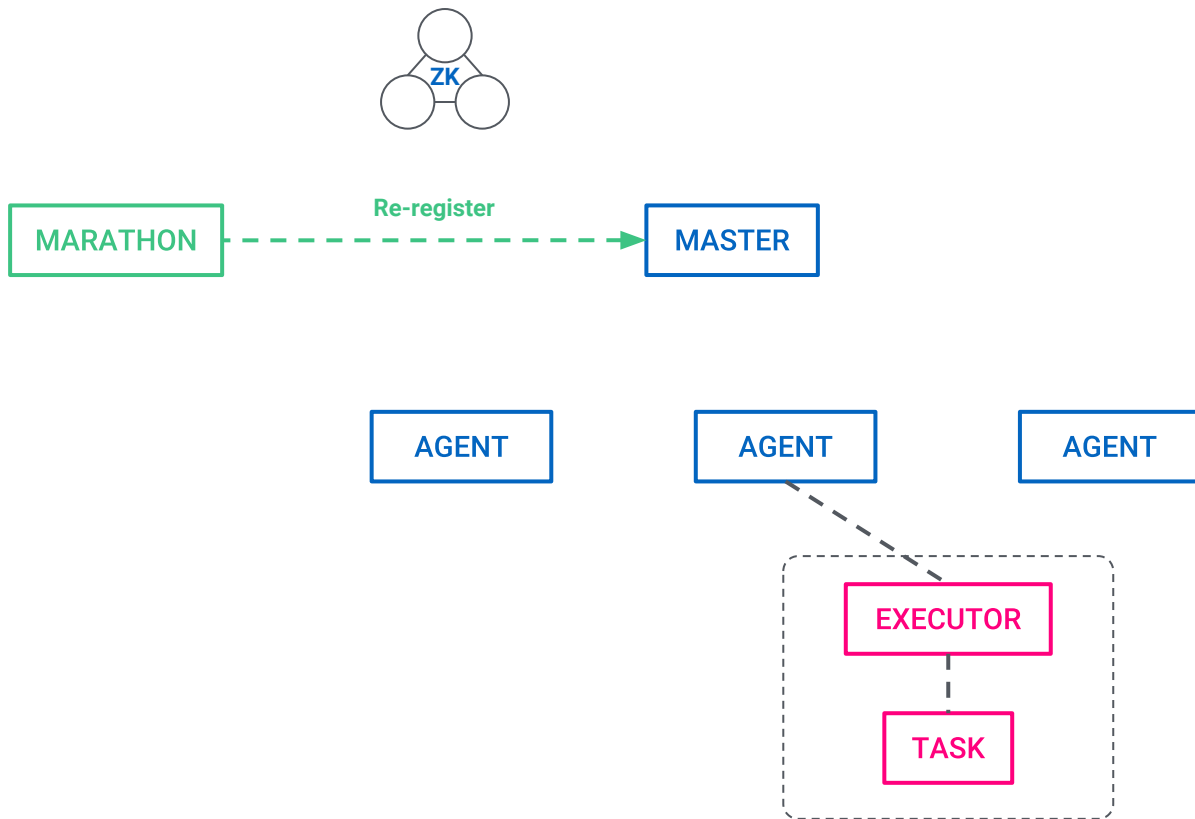
Scheduler Failure



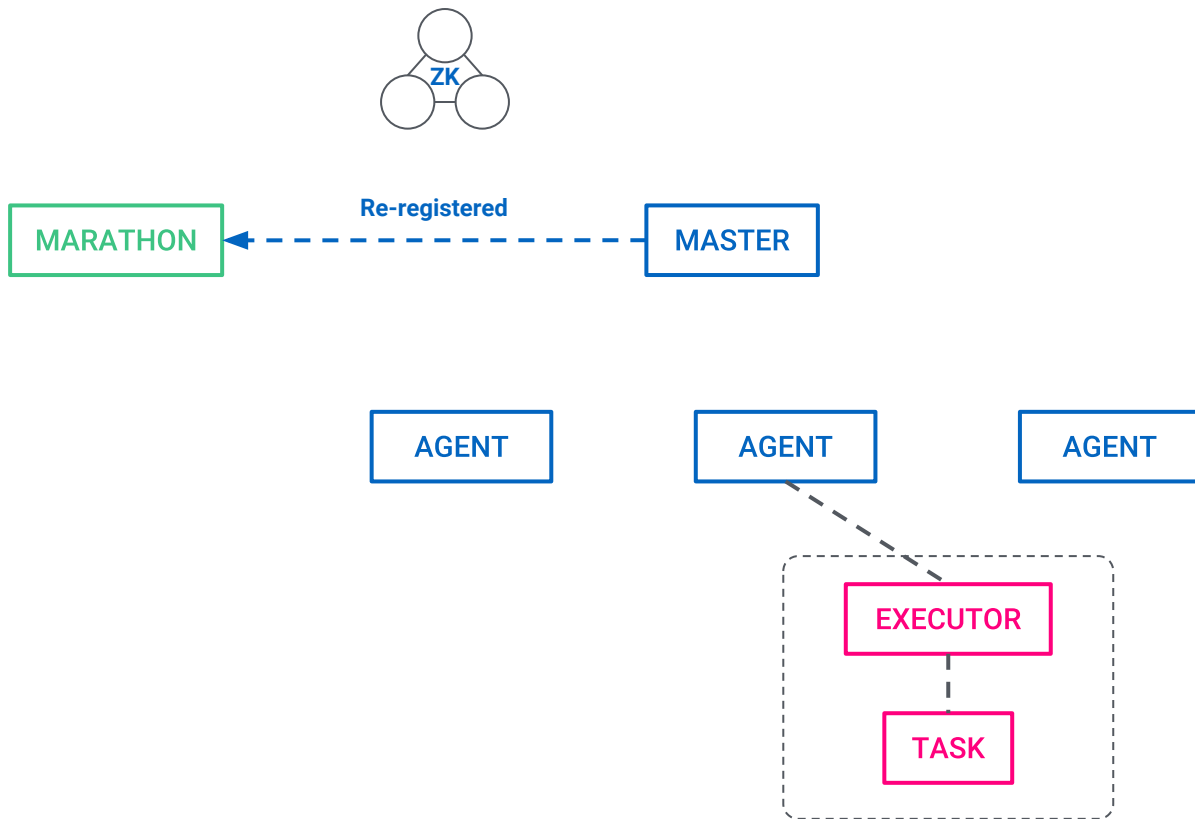
Scheduler Failure



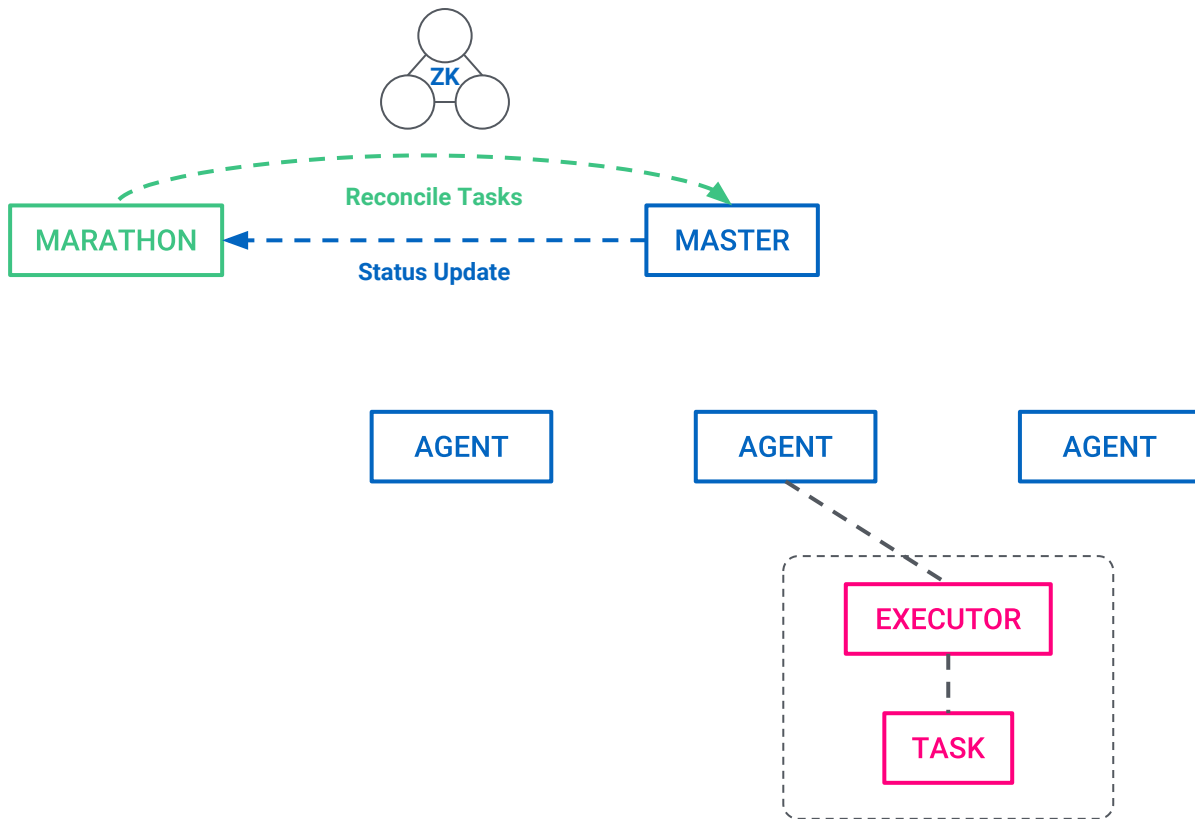
Scheduler Failure



Scheduler Failure



Scheduler Failure



Production Checklist

Mesos

- Configure ulimit settings appropriately
 - `open files`: A value of 32,000 (soft) and 262,144 (hard) have been used successfully in large production deployments
 - `max locked memory`: Increase to account for huge pages, if required
- Monitor both masters and agents for flapping (continuously restarting)
- Monitor the rate of changes in terminal task states, including `TASK_FAILED`, `TASK_LOST`, and `TASK_KILLED`
- Use five masters in production. Three is sufficient for HA in dev/stage/test
- Place master in separate racks or availability zones if possible

Mesos Agents

- Set agent attributes before you run anything on the cluster
 - Once an agent joins the cluster, change to it's attributes may break recovery of containers running on the node
- Explicitly set the resources on the nodes to leave capacity for other services running outside the control of Mesos
 - e.g. Storage daemon that is running on the host operating system

Zookeeper

- Backup ZooKeeper snapshots and logs at regular intervals
 - Guano and zkConfig.py
- Marathon, Metronome, and other frameworks store their states in ZK
- Services launched on the cluster should not store state in the ZK ensemble that is in use by the DC/OS control plane processes
- Monitor ZK's JVM metrics, such as heap usage, GC pause times, and full-collection frequency
- Monitor ZK for number of client connections, total number of znodes, size of znodes (min, max, avg, P99), and read/write performance metrics

Getting Help

Getting Help

- Enterprise and community DC/OS docs
 - <https://docs.mesosphere.com>
- Github repos
 - <https://github.com/mesosphere>
 - <https://github.com/dcos>
- Mesosphere enterprise support
 - <https://support.mesosphere.com> (requires login)
- Community
 - <https://dcos.io/community>



MESOSPHERE

Appendix

Unit Name	Description	Node(s) of Execution M = Master A = Agent
<code>dcos-adminrouter.service</code>	Exposes a unified control plane proxy for components and services using NGINX	M
<code>dcos-backup-master.service</code>	Backup & restore service	M
<code>dcos-bouncer.service</code>	Controls access to DC/OS components and services by managing users, user groups, service accounts, permissions, and identity providers	M
<code>dcos-ca.service</code>	Issues signed digital certificates for secure communication	M
<code>dcos-cluster-linker.service</code>	Service for DC/OS Cluster Linker	M
<code>dcos-cockroach.service</code>	Database for the DC/OS IAM	M
<code>dcos-cosmos.service</code>	Installs and manages DC/OS packages from DC/OS package repositories, such as the Mesosphere Universe	M
<code>dcos-diagnostics.service</code>	Aggregates and exposes component health	M, A
<code>dcos-exhibitor.service</code>	Zookeeper supervisor service	M
<code>dcos-history.service</code>	Caches and exposes historical system state	M

Unit Name	Description	Node(s) of Execution M = Master A = Agent
<code>dcos-licensing.service</code>	Licensing audit service	M
<code>dcos-log-master.service</code>	Exposes master node and component logs	M
<code>dcos-marathon.service</code>	Container orchestration engine	M
<code>dcos-mesos-dns.service</code>	Domain name based service discovery	M
<code>dcos-mesos-master.service</code>	Distributed systems kernel	M
<code>dcos-metrics-master.service</code>	Exposes node metrics	M
<code>dcos-metronome.service</code>	Job orchestration	M
<code>dcos-net-watchdog.service</code>	Restarts dcos-net when it is unhealthy	M, A
<code>dcos-net.service</code>	A distributed systems & network overlay orchestration engine	M, A
<code>dcos-pkgpanda-api.service</code>	Installs and manages DC/OS components	M, A

Unit Name	Description	Node(s) of Execution M = Master A = Agent
<code>dcos-secrets.service</code>	Provides a secure API for storing and retrieving secrets from Vault, a secret store	M
<code>dcos-vault.service</code>	DC/OS Default Secret Store Backend	M
<code>dcos-adminrouter-agent.service</code>	Exposes a unified control plane proxy for components and services using NGINX	A
<code>dcos-log-agent.service</code>	Exposes agent node, component, and container (task) logs	A
<code>dcos-mesos-slave-public.service</code>	Distributed systems kernel public agent	A
<code>dcos-mesos-slave.service</code>	Distributed systems kernel private agent	A
<code>dcos-metrics-agent.service</code>	Exposes node, container, and application metrics	A
<code>dcos-rexray.service</code>	A vendor agnostic storage orchestration engine	A