# ISS

## M3s: Apache Mesos and K3s

Resource-conscious Platform for ML and Data-Processing

Apachecon 2022 - Marcel Neuhausler

ISSGOVERNANCE.COM

# Agenda

- Motivation

- PaaS - M3s Developer Cloud

  - Apache Mesos

  - Frameworks - M3s

  - M3s Management

- SaaS - M3s Data-Science

  - Jupyterlab

  - Apache Spark

  - Data Pipeline - Apache NiFi. Apache Airflow

  - ML Platform, MLOps

- M3trics

- M3s Future

- Summary

#BrilliantTogether

# Motivation

ISS Tech Innovation Lab

Taking ML to Production - a Journey

Innovation: Data Inspired Continuous Exploration

**Guiding Principles**

- Strive for Simplicity
- Partition where and whenever you can .. and you have tools to support and automate it.
- Be Resource Conscious .. throughout each layer of the stack.
- Prioritize Independence .. "Don't be driven by technology, drive it" - BMW
- It is all about a healthy balance .. which to find is the hardest part.

You Build It You Run It

## con·scious
/ˈkän(t)SHəs/

*adjective*

aware of and responding to one's surroundings; awake.
*synonyms:* aware, awake, alert, responsive, sentient, compos mentis
"the patient was conscious"

- having knowledge of something; aware.
"we are **conscious of** the extent of the problem"
*synonyms:* aware, mindful, sensible;  More

- painfully aware of; sensitive to.
"he was very conscious of his appearance"

**Resource Conscious Infrastructure**

- Hardware Constraints
- Energy Consumption
- Head Count
- Operational Overhead
- Configuration, Deployment Complexity
- Effort to Debug, to Detect Anomalies
- Time to Develop versus Time to Operate

"Simplicity is a great virtue but it requires hard work to achieve it and education to appreciate it. And to make matters worse: complexity sells better."
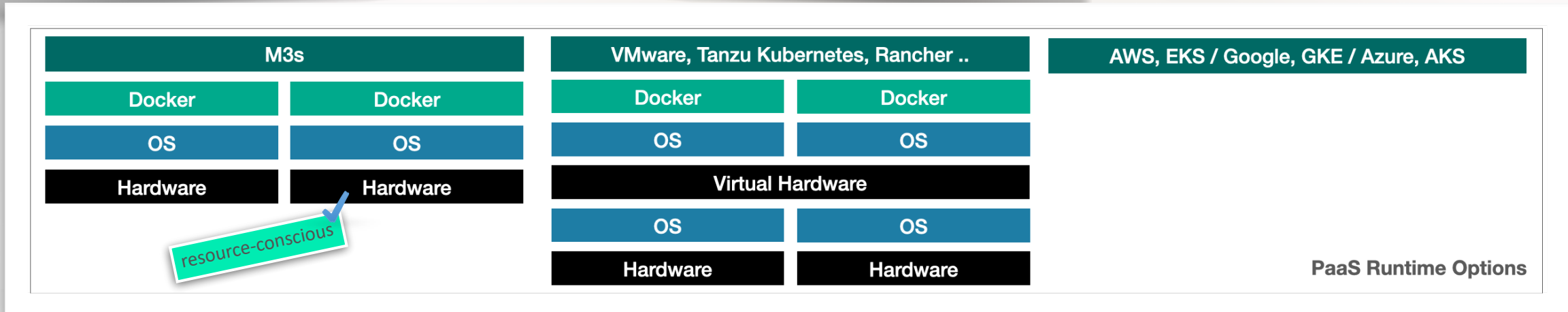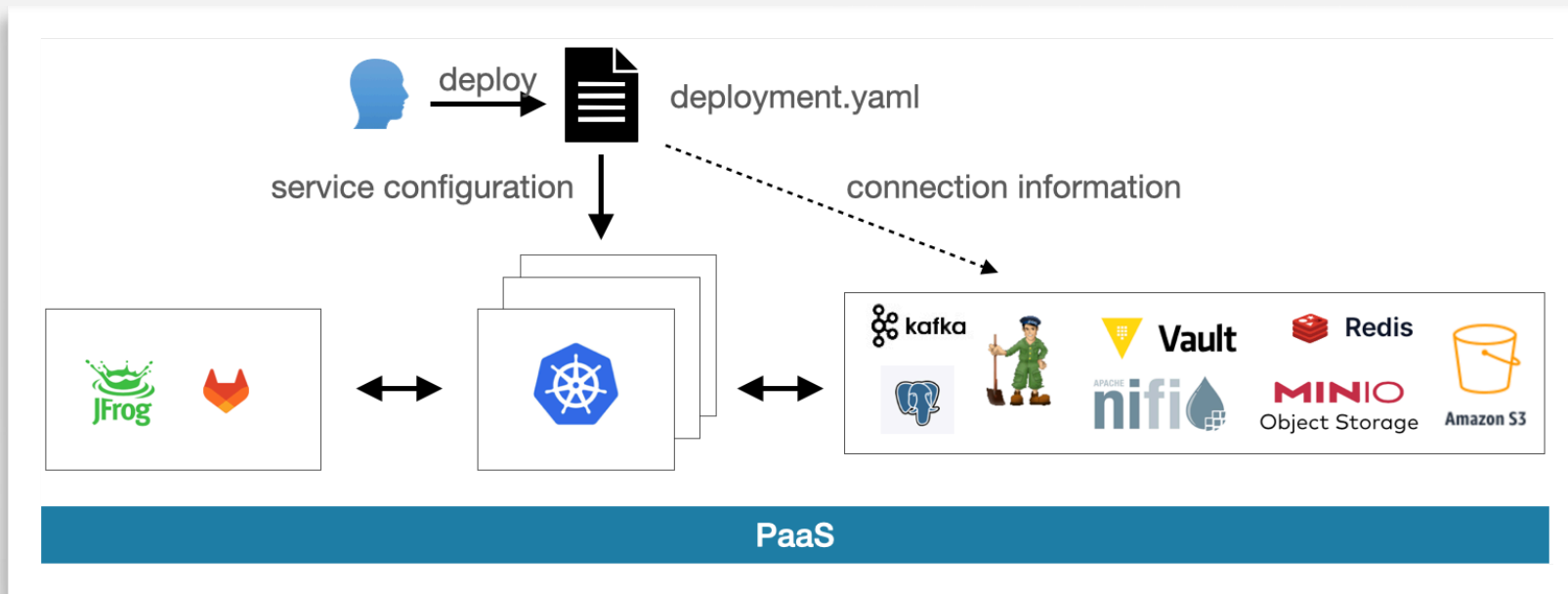- Edsger Dijkstra

Be careful ...

## The anomaly of cheap complexity

It is often more cost-effective to take a very complicated device, and make it simulate simplicity, than to make a simpler device.

https://freedom-to-tinker.com/2022/08/03/the-anomaly-of-cheap-complexity/

Apachecon 2022 - Marcel Neuhausler - Institutional Shareholder Services
This work is licensed under a Creative Commons Attribution 4.0 License.

#BrilliantTogether

# PaaS - Developer Cloud



Apachecon 2022 - Marcel Neuhausler - Institutional Shareholder Services
This work is licensed under a Creative Commons Attribution 4.0 License.

#BrilliantTogether

# Apache Mesos - Past & Present

Apache Mesos Narrowly Avoids a Move to the Attic (for Now)

That's why we have chosen to sunset DC/OS, with an end-of-life date of October 31, 2021. With DKP, our customers get the same benefits
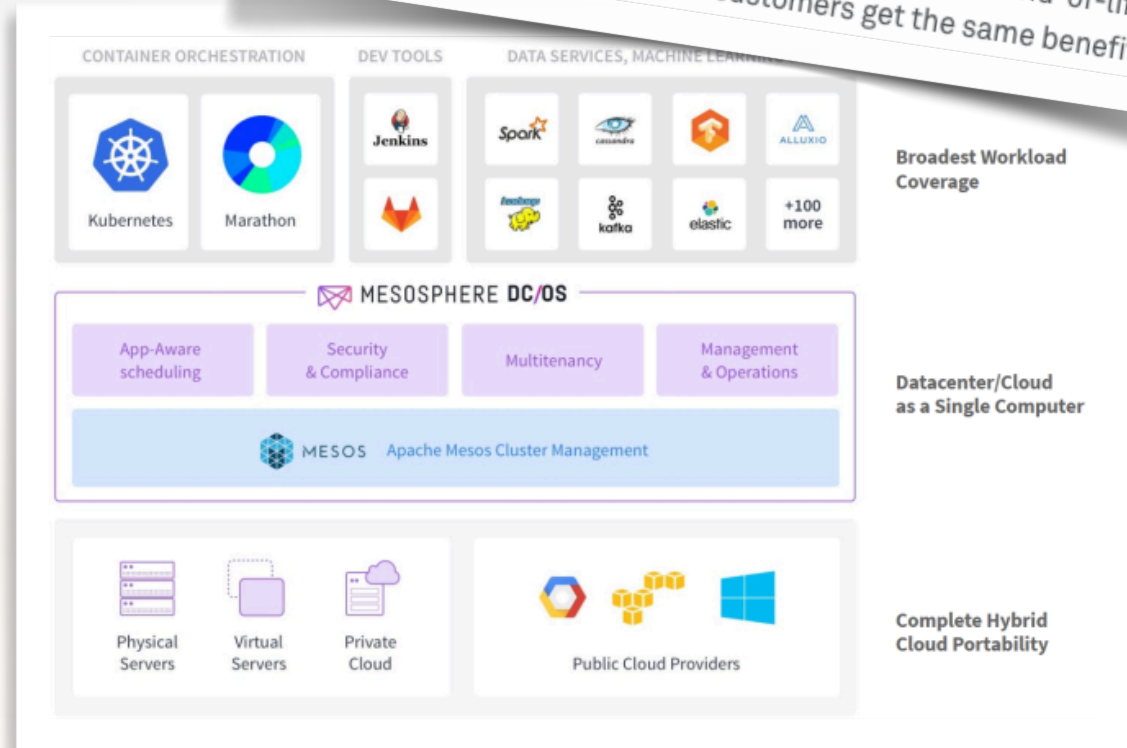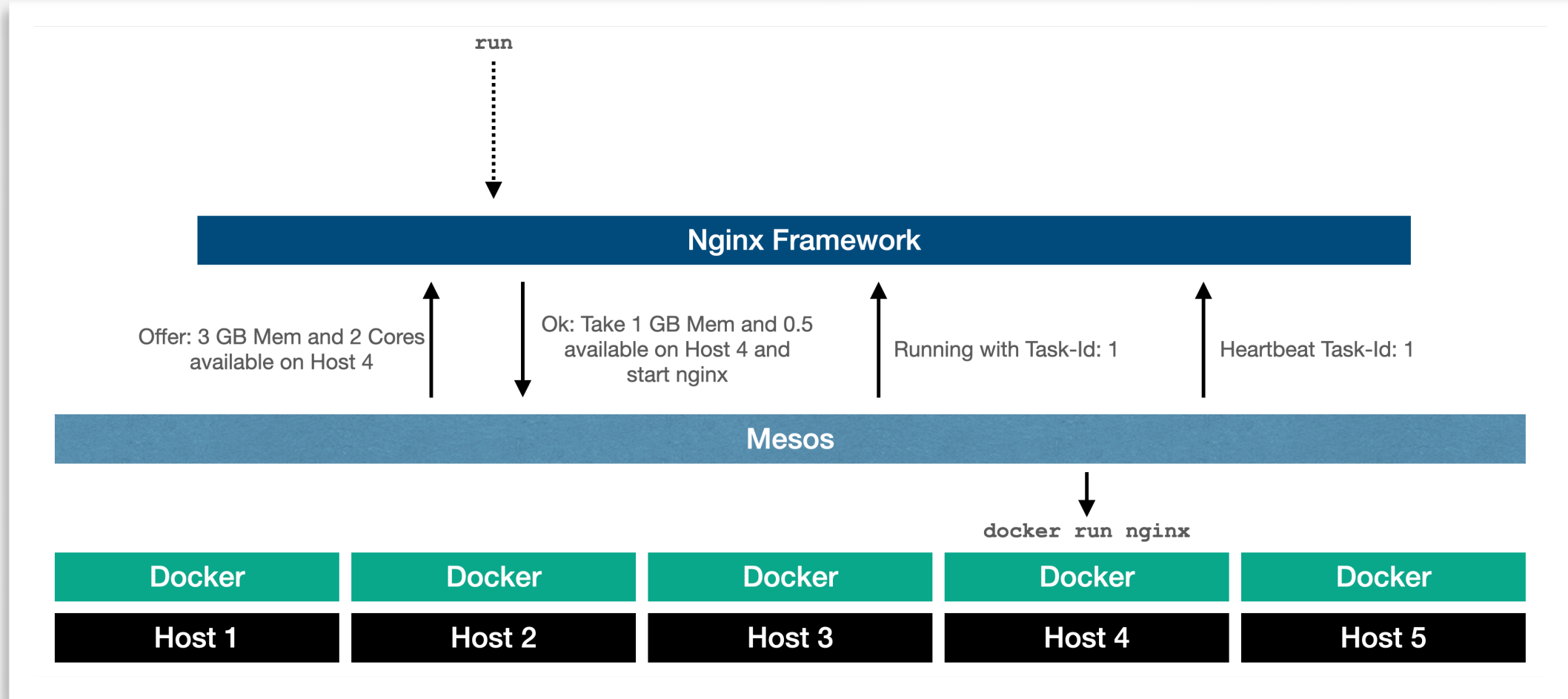


Soon after launch, Mesos was adopted by Twitter, Apple, Yelp, Uber, and Netflix.

**Mesos Is Dead. Long Live Mesos!**

Thanks to the current Apache Mesos main committers!
- Charles-Françoise Natali
- Qian Zhang
- Andreas Peters
- ...

https://mesos.apache.org/
https://thenewstack.io/apache-mesos-narrowly-avoids-a-move-to-the-attic-for-now/
https://d2iq.com/products/dcos

Apachecon 2022 - Marcel Neuhausler - Institutional Shareholder Services
This work is licensed under a Creative Commons Attribution 4.0 License.

#BrilliantTogether

# Apache Mesos

run

**Nginx Framework**

Offer: 3 GB Mem and 2 Cores available on Host 4

Ok: Take 1 GB Mem and 0.5 available on Host 4 and start nginx

Running with Task-Id: 1

Heartbeat Task-Id: 1

**Mesos**

`docker run nginx`

| Docker | Docker | Docker | Docker | Docker |
|--------|--------|--------|--------|--------|
| Host 1 | Host 2 | Host 3 | Host 4 | Host 5 |

Apachecon 2022 - Marcel Neuhausler - Institutional Shareholder Services
This work is licensed under a Creative Commons Attribution 4.0 License.

#BrilliantTogether

# Apache Mesos - Frameworks

## Long Running Services

- Aurora is a service scheduler that runs on top of Mesos, enabling you to run long-running services that take advantage of Mesos' scalability, fault-tolerance, and resource isolation.
- Marathon is a private PaaS built on Mesos. It automatically handles hardware or software failures and ensures that an app is "always on".
- Singularity is a scheduler (HTTP API and web interface) for running Mesos tasks: long running processes, one-off tasks, and scheduled jobs.
- SSSP is a simple web application that provides a white-label "Megaupload" for storing and sharing files in S3.

## Big Data Processing

- Cray Chapel is a productive parallel programming language. The Chapel Mesos scheduler lets you run Chapel programs on Mesos.
- Dpark is a Python clone of Spark, a MapReduce-like framework written in Python, running on Mesos.
- Exelixi is a distributed framework for running genetic algorithms at scale.
- Hadoop Running Hadoop on Mesos distributes MapReduce jobs efficiently across an entire cluster.
- Hama is a distributed computing framework based on Bulk Synchronous Parallel computing techniques for massive scientific computations e.g., matrix, graph and network algorithms.
- MPI is a message-passing system designed to function on a wide variety of parallel computers.
- Spark is a fast and general-purpose cluster computing system which makes parallel jobs easy to write.
- Storm is a distributed realtime computation system. Storm makes it easy to reliably process unbounded streams of data, doing for realtime processing what Hadoop did for batch processing.

## Batch Scheduling

- Chronos is a distributed job scheduler that supports complex job topologies. It can be used as a more fault-tolerant replacement for Cron.
- Cook is a job scheduler like Torque that not only supports individual tasks, but also Spark. Cook provides powerful automatic preemption and multitenancy features for shared clusters, in order to guarantee throughput to all users while allowing individuals to temporarily "burst" to additional resources as needed. Cook provides a simple REST API & Java client for interaction.
- Elastic-Job-Cloud is a distributed scheduled job cloud solution designed with HA and fault-tolerance in mind. It focuses on horizontal scaling, and provides transient and daemon jobs, event and schedule based job triggers, job dependencies, and job history.

https://mesos.readthedocs.io/en/latest/frameworks/
https://netflixtechblog.com/fenzo-oss-scheduler-for-apache-mesos-frameworks-5c340e77e543
https://github.com/att-innovate/torc_scheduler

Application-specific Open and Closed-Source Frameworks available under DC/OS



Apachecon 2022 - Marcel Neuhausler - Institutional Shareholder Services
This work is licensed under a Creative Commons Attribution 4.0 License.

#BrilliantTogether

# Apache Mesos - M3s Frameworks

Current

**Aventer UG - Andreas Peters**

## M3S - Apache Mesos Kubernetes Framework

### M3s - Apache Mesos Kubernetes Framework

#### Introduction

M3s is a Golang based Apache Mesos Framework to run and deploy Kubernetes through K3s from Rancher.

#### Requirements

- Apache Mesos min 1.6.0
- Mesos with SSL and Authentication is optional
- Persistent Storage to store Kubernetes data

#### How To use

Before we can run the following M3s example, we have to creat a docker network.

```
docker network create --subnet 10.32.0.0/24 mini
```

Afte that, we can create a docker-compose file to run the services we need.

Sidebar:
1. Introduction
2. Configuration
CLI
3. M3s CLI Installation for Mesos-CLI
4. M3s CLI Usage
5. Kubectl Usage with M3s
API
6. Scale up/down Kubernetes Service
7. Count of running Kubernetes Server/Agents
8. Get Kubeconfig
9. Get Kubernetes Version
10. Get M3s status information
Examples
11. Create NGINX ingress with traefik
12. Containerd sideload configuration

https://aventer-ug.github.io/mesos-m3s/index.html
https://github.com/AVENTER-UG/mesos-m3s
https://k3s.io/

K3S

resource-conscious

## mesos-m3s

Chat Support | Docs Support

resource-conscious

Mesos Framework to run Kubernetes (K3S)

## Requirements

- Apache Mesos min 1.6.0
- Mesos with SSL and Authentication is optional
- Persistent Storage to keep K3S data (not object storage)
- Redis DB

## Run Framework

The following environment parameters are only a example. All parameters and der default values are documented in the `init.go` file (real documentation will be coming later)

```
export FRAMEWORK_USER="root"
export FRAMEWORK_NAME="k3sframework"
export FRAMEWORK_PORT="10000"
export FRAMEWORK_ROLE="k3s"
export FRAMEWORK_STATEFILE_PATH="/tmp"
export MESOS_PRINCIPAL="<mesos_principal>"
export MESOS_USERNAME="<mesos_user>"
export MESOS_PASSWORD="<mesos_password>"
export MESOS_MASTER="<mesos_master_server>:5050"
export MESOS_CNI="weave"
export LOGLEVEL="DEBUG"
export DOMAIN="weave.local"
export K3S_SERVER_COUNT=1
export K3S_AGENT_COUNT=1
```
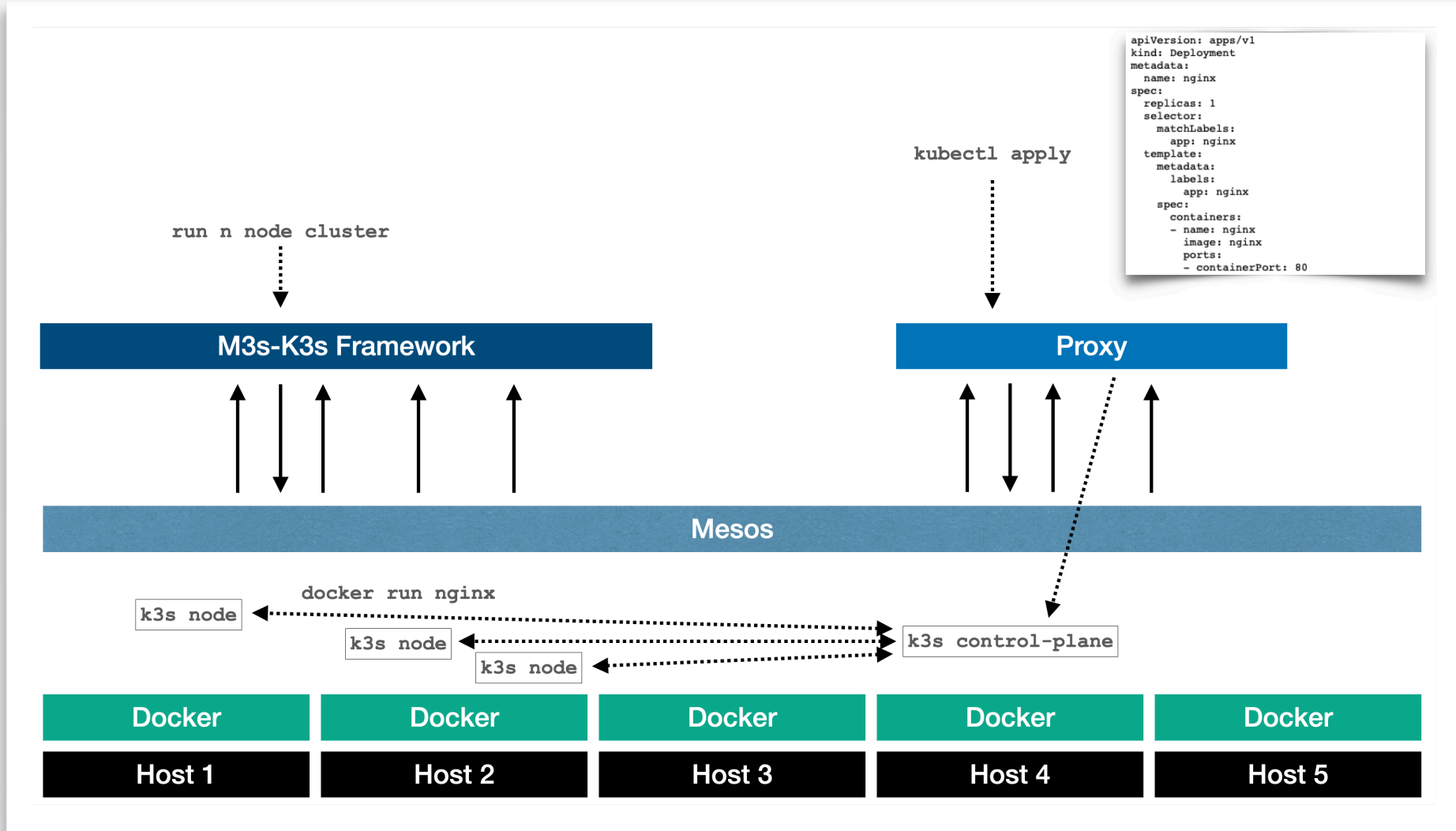
Apachecon 2022 - Marcel Neuhausler - Institutional Shareholder Services
This work is licensed under a Creative Commons Attribution 4.0 License.

#BrilliantTogether

# Apache Mesos - M3s



Baremetal-installs:
- Mesos-Master/Agent
- Mesos-DNS
- Apache Zookeeper
- Hashicorp Vault
- Docker

resource-conscious ✔

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: nginx
spec:
  replicas: 1
  selector:
    matchLabels:
      app: nginx
  template:
    metadata:
      labels:
        app: nginx
    spec:
      containers:
      - name: nginx
        image: nginx
        ports:
        - containerPort: 80
```

kubectl apply

run n node cluster

**M3s-K3s Framework**

**Proxy**

**Mesos**

docker run nginx

`k3s node`

`k3s node`

`k3s node`

`k3s control-plane`

| Docker | Docker | Docker | Docker | Docker |
|--------|--------|--------|--------|--------|
| Host 1 | Host 2 | Host 3 | Host 4 | Host 5 |

Apachecon 2022 - Marcel Neuhausler - Institutional Shareholder Services
This work is licensed under a Creative Commons Attribution 4.0 License.

#BrilliantTogether

# Apache Mesos - M3s Frameworks

**Current**

**Aventer UG - Andreas Peters**

**resource-conscious** ✓



Mesos Compose - The Docker-Compose Orchestrator for Mesos

1. Introduction
2. Custom Labels
3. Environment Variables

CLI

4. Installation for Mesos-CLI
5. Usage

Networking

6. Bridged Network
7. User Network
8. Host Network

Hashicorp Vault

9. How to use

## Introduction - mesos-compose, the docker-compose framework for Apache Mesos

### Requirements

- Apache Mesos min 1.6.0
- Mesos with SSL and Authentication is optional
- Redis Database
- Docker Compose Spec 3.9

### Example

Compose file with all supportet parameters:

```
version: '3.9'

services:
  app:
    image: alpine:latest
    command: ["sleep", "1000"]
    restart: always
```

# mesos-compose

Chat Support  Docs Support

Mesos Framework to use docker-compose files.

## Requirements

- Apache Mesos min 1.6.0
- Mesos with SSL and Authentication is optional
- Redis Database
- Docker Compose Spec 3.9

## Example

The compose file:

```
version: '3.9'

services:
  app:
    image: alpine:latest
    command: ["sleep", "1000"]
    restart: always
    volumes:
      - "12345test:/tmp"
    environment:
      - MYSQL_HOST=test
    hostname: test
```

https://github.com/AVENTER-UG/mesos-compose
https://aventer-ug.github.io/mesos-compose/

Apachecon 2022 - Marcel Neuhausler - Institutional Shareholder Services
This work is licensed under a Creative Commons Attribution 4.0 License.

#BrilliantTogether

# Apache Mesos UI

resource-conscious

In case you want to try it out:
https://hub.docker.com/extensions/avhost/docker-mesos-extension

Apachecon 2022 - Marcel Neuhausler - Institutional Shareholder Services

This work is licensed under a Creative Commons Attribution 4.0 License.

#BrilliantTogether

# M3s Management

Docker Swarm as replacement for systemd/puppet/ansible

resource-conscious

## Docker Swarm Services

pgAdmin

| M3s Dashboard | M3s Dashboard | M3s Dashboard | | |
| Swarm Visualizer | Swarm Visualizer | Swarm Visualizer | Traefik | Traefik |

| M3s-Manager | M3s Config | M3s-Manager | M3s Config | M3s-Manager | M3s Config | M3s-Manager | M3s Config | M3s-Manager | M3s Config |

| Docker | Docker | Docker | Docker | Docker |
| OS | OS | OS | OS | OS |
| Master 1 | Master 2 | Master 3 | Agent 1 | Agent n |

...

https://docs.docker.com/engine/swarm/
https://www.postgresql.org/
https://www.pgadmin.org/

https://traefik.io/
https://github.com/dockersamples/docker-swarm-visualizer
https://www.hashicorp.com/products/vault

Apachecon 2022 - Marcel Neuhausler - Institutional Shareholder Services
This work is licensed under a Creative Commons Attribution 4.0 License.

#BrilliantTogether

# M3s Management Service

**resource-conscious**

**Infrastructure as Code and APIs**
Provides access to Docker Engine, M3s Frameworks, Shell commands defined in M3s config
Will be open-sourced later this year.

```rust
use std::fs;
use std::io::prelude::*;
use std::net::TcpListener;
use std::net::TcpStream;
use std::thread;
use std::env::temp_dir;
use std::time::SystemTime;
use std::fs::File;

pub struct HttpServer {
    pub listening_port: u16,
    pub environment: String
}

impl HttpServer {
    pub fn run(&mut self) {
        let listener = TcpListener::bind(format!("0.0.0.0:{}", self.listening_port)).unwrap();

        for stream in listener.incoming() {
            let stream = stream.unwrap();
            let environment = self.environment.clone();
            thread::spawn(move || {
                handle_connection(stream, environment.clone());
            });
        }
    }
}
```

```rust
fn handle_connection(mut stream: TcpStream, environment: String) {
    let mut buffer = [0; 20000];
    stream.read(&mut buffer).unwrap();

    let get = b"GET / HTTP/1.1\r\n";
    let docker_version = b"GET /docker/version HTTP/1.1\r\n";
    let docker_containers = b"GET /docker/containers HTTP/1.1\r\n";
    let docker_swarm_deploy = b"POST /docker/swarm/";
    let docker_swarm_delete = b"DELETE /docker/swarm/";
    let git_clone_m3sconfig = b"POST /git/clone/m3sconfig HTTP/1.1\r\n";
    let git_pull_m3sconfig = b"POST /git/pull/m3sconfig HTTP/1.1\r\n";
    let executor_run_it = b"GET /executor/runit/";

    let (request_type, status_line, filename) = if buffer.starts_with(get) {
        (RequestType::Hello, "HTTP/1.1 200 OK", "./api/hello.json")
    } else if buffer.starts_with(docker_version) {
        (RequestType::DockerVersion, "HTTP/1.1 200 OK", "")
    } else if buffer.starts_with(docker_containers) {
        (RequestType::DockerContainers, "HTTP/1.1 200 OK", "")
    } else if buffer.starts_with(docker_swarm_deploy) {
        (RequestType::DockerSwarmDeploy, "HTTP/1.1 200 OK", "")
    } else if buffer.starts_with(docker_swarm_delete) {
        (RequestType::DockerSwarmDelete, "HTTP/1.1 200 OK", "")
    } else if buffer.starts_with(git_clone_m3sconfig) {
        (RequestType::GitCloneM3sConfig, "HTTP/1.1 200 OK", "")
    } else if buffer.starts_with(git_pull_m3sconfig) {
        (RequestType::GitPullM3sConfig, "HTTP/1.1 200 OK", "")
    } else if buffer.starts_with(executor_run_it) {
        (RequestType::ExecutorRunIt, "HTTP/1.1 200 OK", "")
    } else {
        (RequestType::PageNotFound, "HTTP/1.1 404 NOT FOUND", "./api/404.json")
    };

    let request = String::from_utf8(buffer.to_vec()).unwrap();
    let path = get_path_from_request(request.clone());
    let data = get_data_from_request(request);
```

```
measurementName,tagKey=tagValue fieldKey="fieldValue" 1465839830100400200
--------------- --------------- --------------------- -------------------
       |               |                 |                    |
  Measurement       Tag set          Field set            Timestamp
```

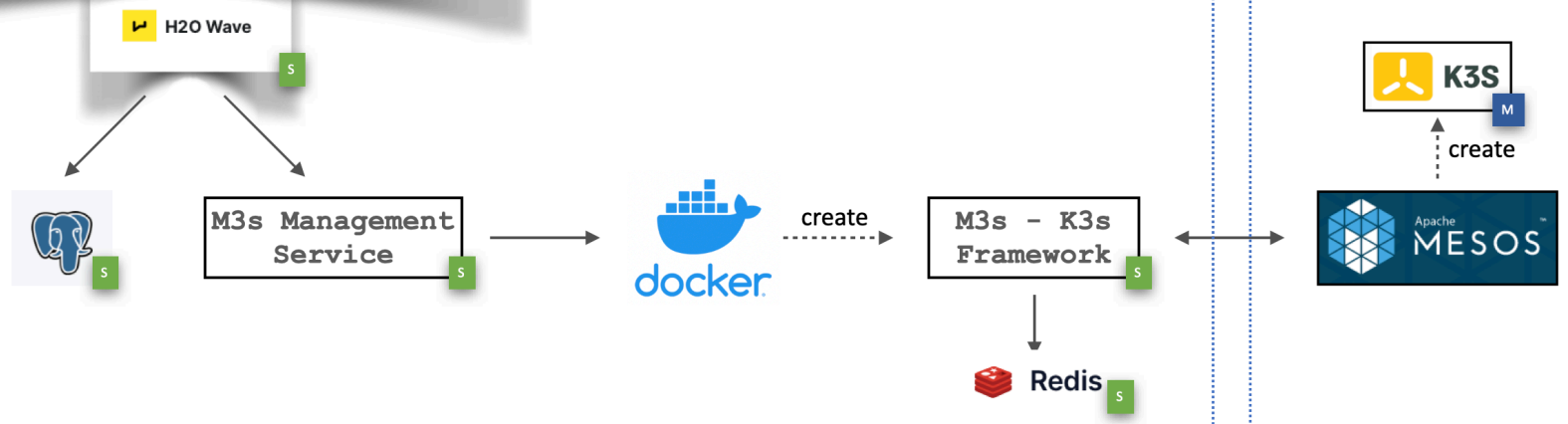Example: InflxuxDB Line Protocol
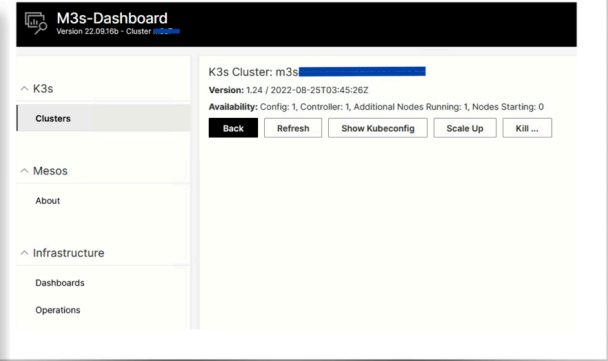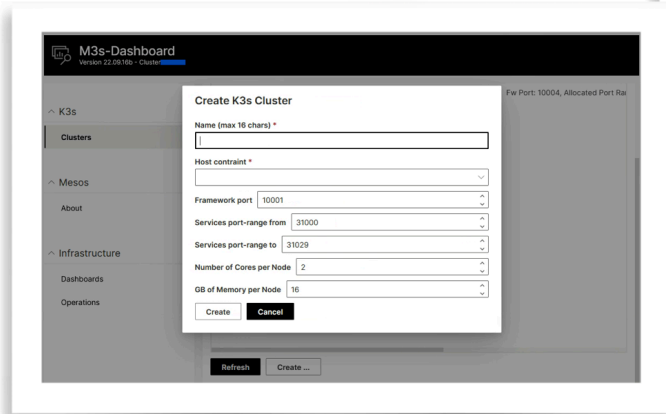https://docs.influxdata.com/influxdb/cloud/reference/syntax/line-protocol/

# M3s - Managing K3s Cluster

resource-conscious ✓

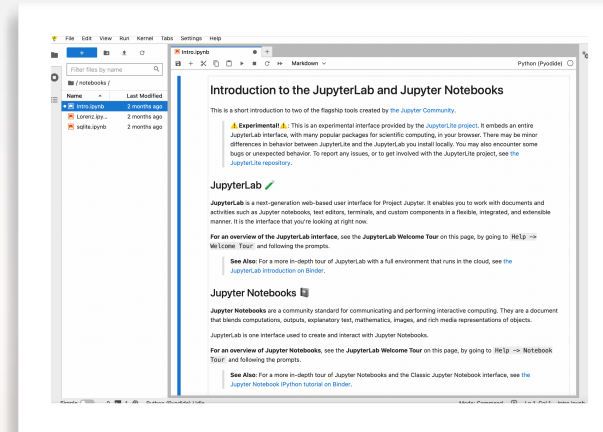**M3s-Dashboard**
Will be open-sourced later this year.



https://wave.h2o.ai/
https://docs.docker.com/engine/swarm/
https://www.postgresql.org/
https://redis.io/
https://k3s.io/

Apachecon 2022 - Marcel Neuhausler - Institutional Shareholder Services
This work is licensed under a Creative Commons Attribution 4.0 License.

#BrilliantTogether

# SaaS - M3s Data-Science



Apachecon 2022 - Marcel Neuhausler - Institutional Shareholder Services
This work is licensed under a Creative Commons Attribution 4.0 License.

#BrilliantTogether

# Jupyterlab

resource-conscious

Zero to JupyterHub with Kubernetes
https://zero-to-jupyterhub.readthedocs.io/en/stable

H2O Wave

M3s - Compose Framework

Apache MESOS

Apachecon 2022 - Marcel Neuhausler - Institutional Shareholder Services
This work is licensed under a Creative Commons Attribution 4.0 License.

#BrilliantTogether

# Apache Spark

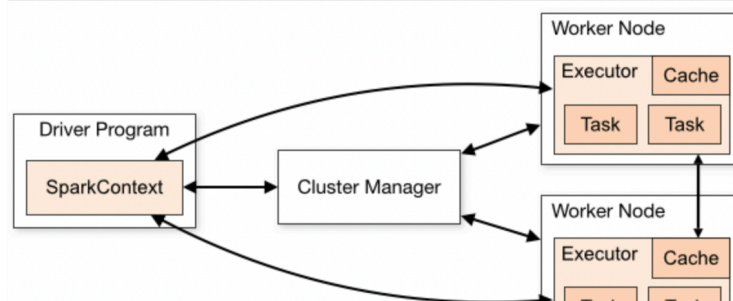https://spark.apache.org/docs/latest/index.html

## Mesos Run Modes

Spark can run over Mesos in two modes: "coarse-grained" (default) and "fine-grained" (deprecated).
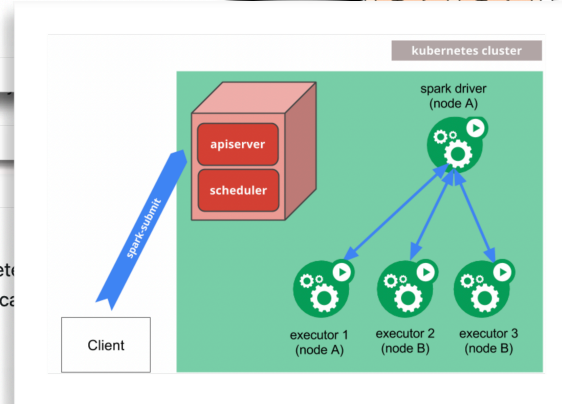
## Coarse-Grained

In "coarse-grained" mode, each Spark executor
configuration variables:

- Executor memory: `spark.executor.memory`
- Executor cores: `spark.executor.cores`
- Number of executors: `spark.cores.max/spa`

In "fine-grained" mode, each Spark task inside the S
(and other frameworks) to share cores at a very fine
but it comes with an additional overhead in launching
interactive queries or serving web requests.

Driver Program
SparkContext → Cluster Manager
Worker Node
Executor | Cache
Task | Task
Worker Node
Executor | Cache

kubernetes cluster

apiserver
scheduler

spark-submit

Client

spark driver (node A)

executor 1 (node A)
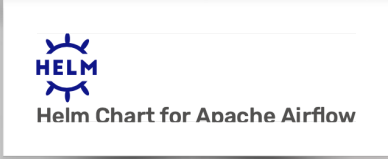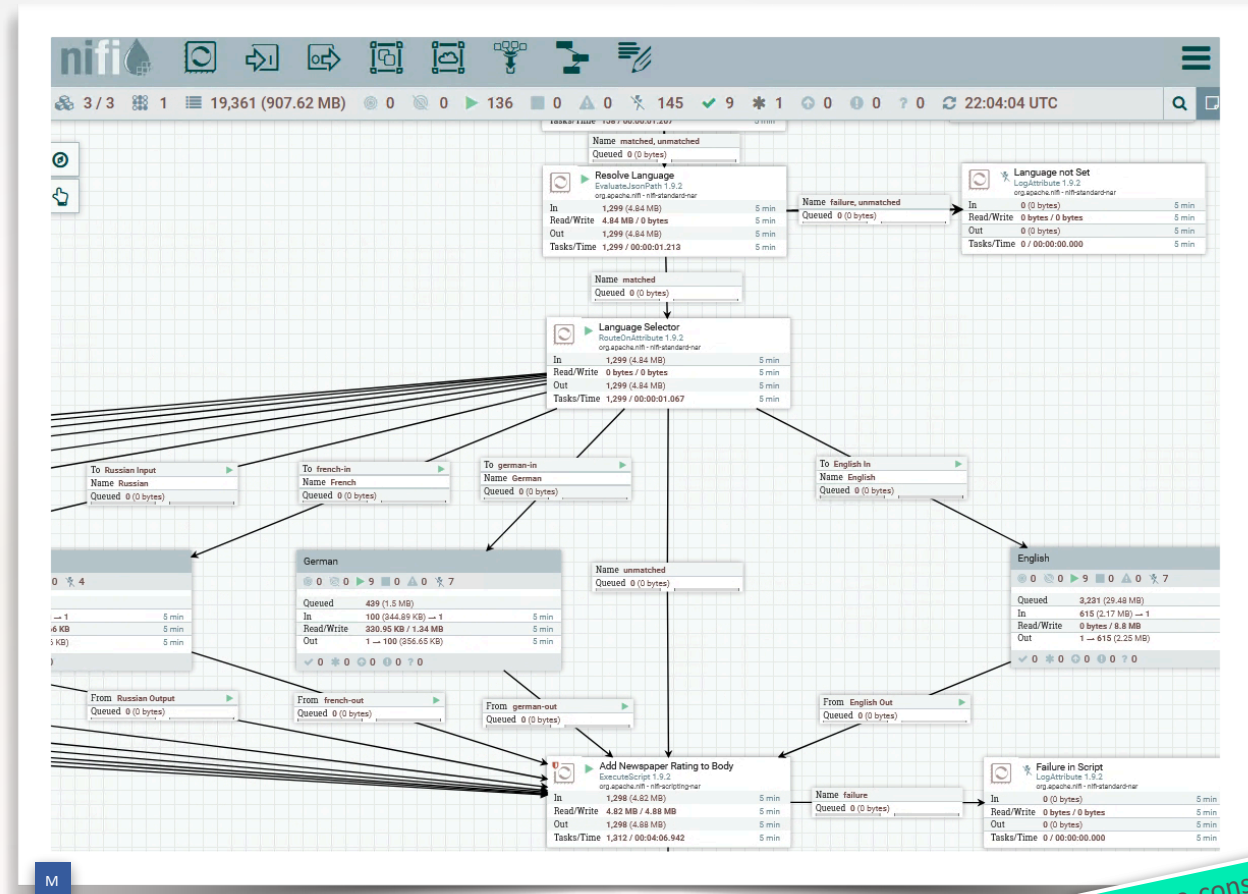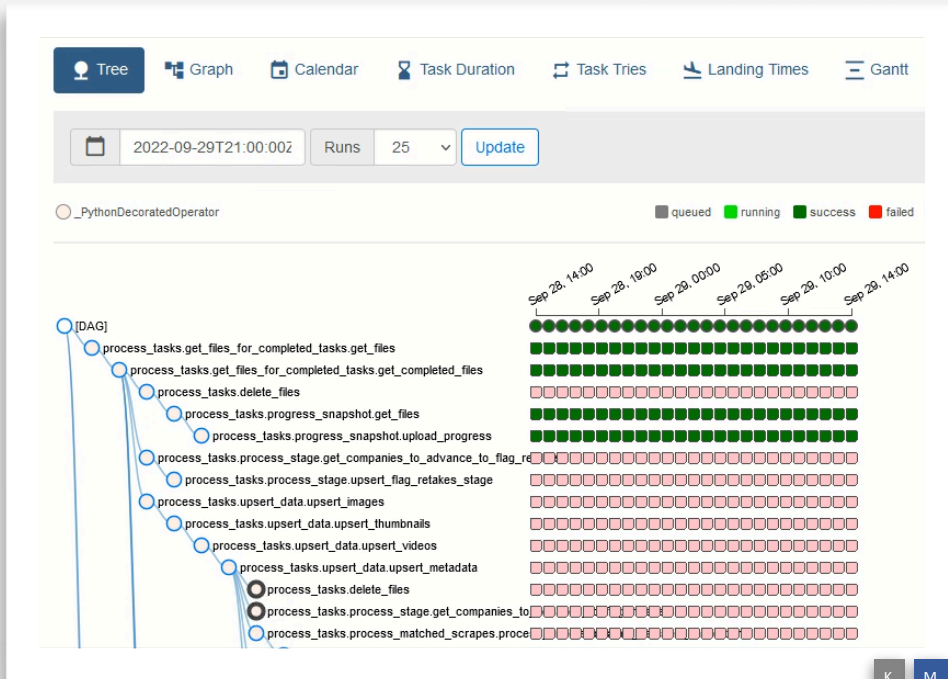executor 2 (node B)
executor 3 (node B)

## Client Mode

Starting with Spark 2.4.0, it is possible to run Spark applications on Kubernete
the driver can run inside a pod or on a physical host. When running an applica
following factors:

## Client Mode Networking

Spark executors must be able to connect to the Spark driver over a hostname and a port that is routable from the Spark executors. The specific network configuration that will be required for Spark to work in client mode will vary per setup. If you run your driver inside a Kubernetes pod, you can use a headless service to allow your driver pod to be routable from the executors by a stable hostname. When deploying your headless service, ensure that the service's label selector will only match the driver pod and no other pods; it is recommended to assign your driver pod a sufficiently unique label and to use that label in the label selector of the headless service. Specify the driver's hostname via `spark.driver.host` and your spark driver's port to `spark.driver.port`.

resource-conscious

Introduction to the JupyterLab and Jupyter Notebooks

JupyterLab

Jupyter Notebooks

M3s-DataScience

H2O Wave

M3s - Compose Framework

Apache MESOS

# Pipeline - Apache Airflow, Apache NiFi



https://nifi.apache.org/
https://airflow.apache.org/

resource-conscious

Apachecon 2022 - Marcel Neuhausler - Institutional Shareholder Services
This work is licensed under a Creative Commons Attribution 4.0 License.

#BrilliantTogether

# ML Platform, MLOps - Example News Processing

https://nifi.apache.org/
https://github.com/nteract/papermill

https://vespa.ai/
https://www.elastic.co/elasticsearch/

https://prestodb.io/
https://mlflow.org/

https://www.influxdata.com/
https://jupyter.org/

https://www.postgresql.org/
https://kafka.apache.org/

Apachecon 2022 - Marcel Neuhausler - Institutional Shareholder Services
This work is licensed under a Creative Commons Attribution 4.0 License.
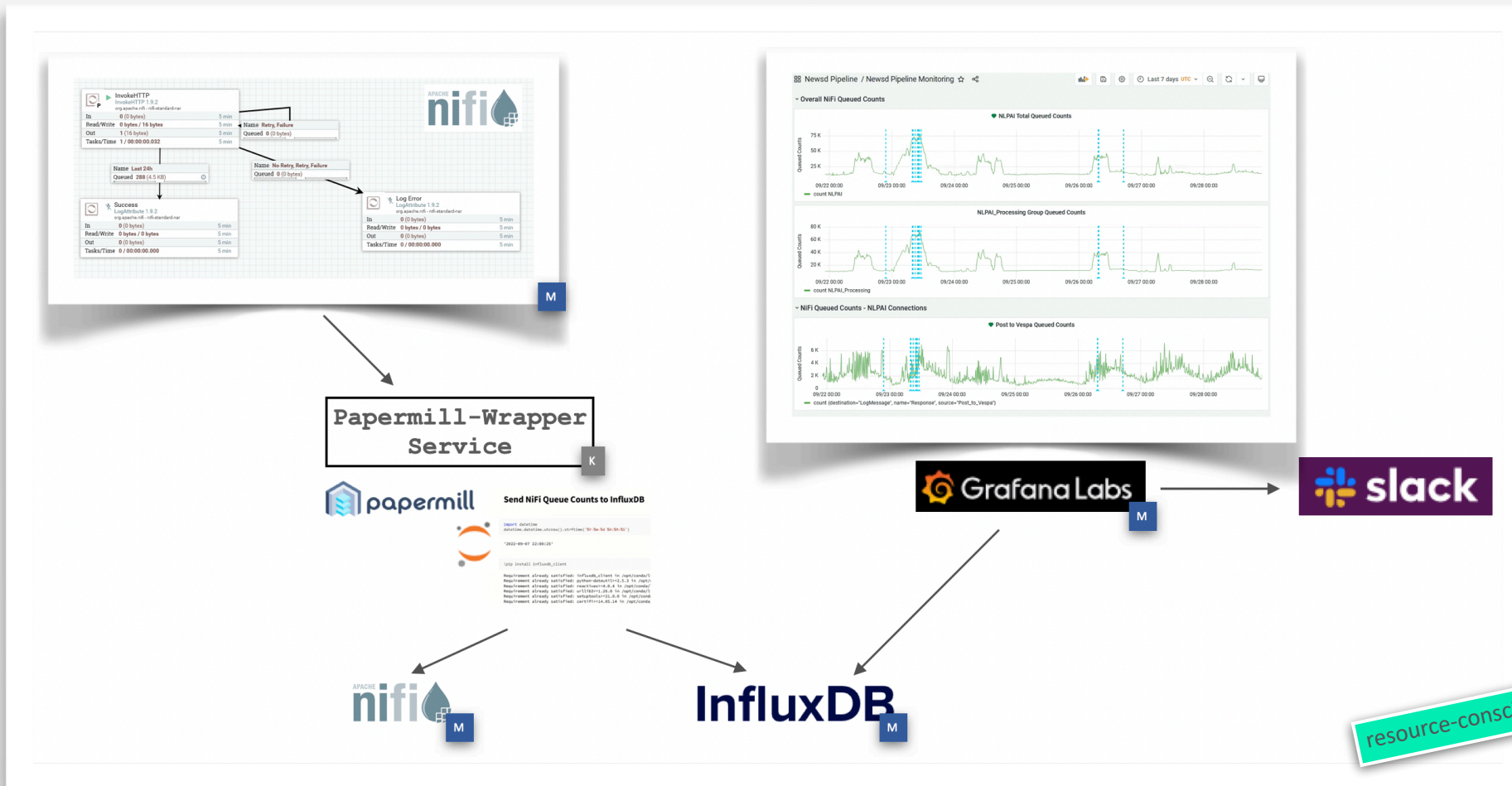
#BrilliantTogether

# M3trics

Apachecon 2022 - Marcel Neuhausler - Institutional Shareholder Services
This work is licensed under a Creative Commons Attribution 4.0 License.

#BrilliantTogether

# M3trics - Example Metrics Collection



Papermill-Wrapper Service

Send NiFi Queue Counts to InfluxDB

resource-conscious ✓

# M3s Future

### Data-Science / ML applied to Infrastructure

**Hybrid Cloud**

Communication

Proxy

AWS Management ↔ M3s Controller ↔ M3s Service Manager

EKS | EC2 | AWS

K3s | Marathon | Compose

Mesos

Host 1 | Host 1 | Host x

Apache **MESOS**  2.0 ?

| | Use Cases | Networking | Security | Observability & Tracing |

eBPF Projects — cilium — katran — Falco

eBPF SDKs

Application

eBPF Kernel Runtime — Verifier & JIT — Maps — Kernel Helper API — OS Runtime
- Tracing
- Profiling
- Monitoring
- Observability
- Security Controls
- Networking
- Network Security
- Load Balancing
- Behavioral Security

**Falco**, the cloud-native runtime security project, is the de facto **Kubernetes threat detection engine**

Detects threats at runtime by observing the behavior of your applications and containers.

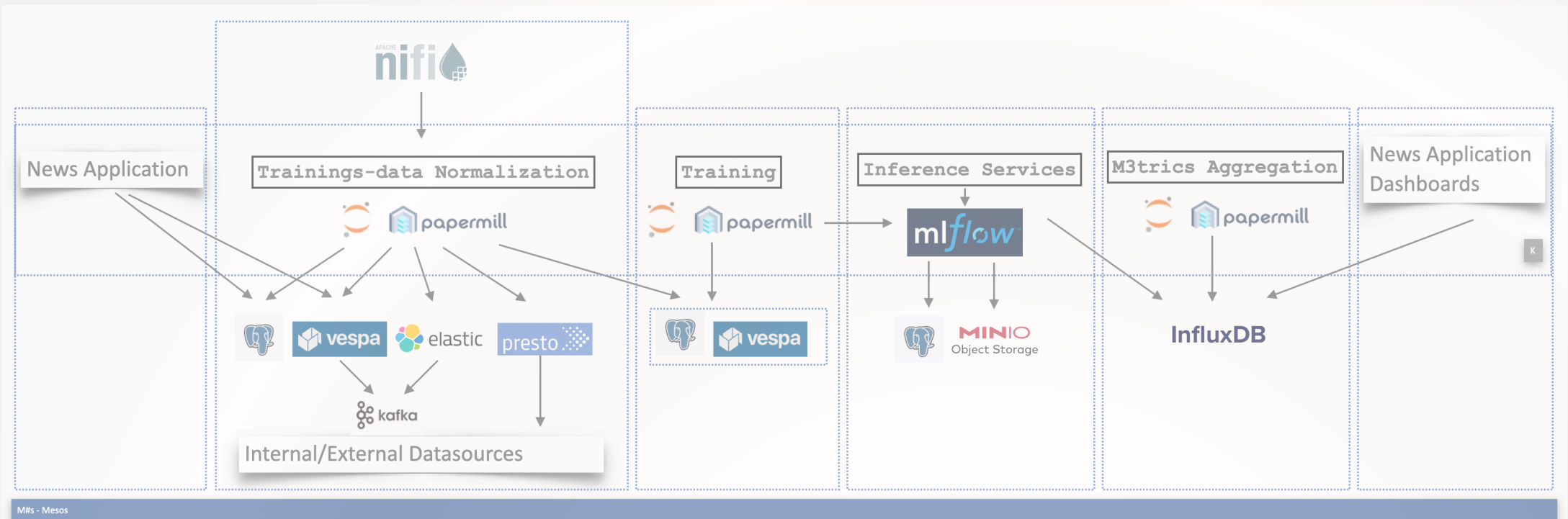Extends threat detection across cloud environments with Falco Plugins.

System Calls
Kubernetes Events
Cloud Activity

Metrics | Events
Logs | Traces

**Tetragon Agent**

Pod app.py — Func Calls
Pod app.go — Code Exec

**Linux Kernel**

**Kernel Runtime**
- Smart Collector
- Stack Traces | Ring Buffer
- Metrics | Hash Maps

Process Execution | Syscall Activity | System Calls
File Access | VFS | Seq Attack | TCP/IP
NS Escapes | Priv Escalations | Namespaces
Data Access | Storage | HTTP, DNS, TLS | Network

**Firecracker**

**WA WEBASSEMBLY**

https://mesos.apache.org/
https://ebpf.io/
https://github.com/cilium/tetragon
https://falco.org/
https://firecracker-microvm.github.io/
https://webassembly.org/

Apachecon 2022 - Marcel Neuhausler - Institutional Shareholder Services
This work is licensed under a Creative Commons Attribution 4.0 License.

#BrilliantTogether

# Summary

- ML to Production ✔.. it was an interesting journey .. and there is more to come.

- Brilliant SREs and Data Infrastructure Engineers can build exciting tools and platforms.

- "Data Inspired Continuous Exploration" - Innovation enabled by Apache Mesos, and many other resource-conscious tools and applications. Thanks to all the Committers and Contributors!

Apachecon 2022 - Marcel Neuhausler - Institutional Shareholder Services
This work is licensed under a Creative Commons Attribution 4.0 License.

#BrilliantTogether