# A few words about me!

Hichem Kenniche,

OSS Product Architect @Instaclustr (part of NetApp)

Previously at:

- Databricks
- Capgemini Invent
- Sony PlayStation

# Disclaimer

- This is not a contribution to any OSS project !

- My vision of things is necessarily biased !

- Most of this is work is based on the principles of OSS, open data, and a culture of knowledge sharing ❤️

- (Human) Learning is a Lifelong WIP …

# Our agenda for today

- The fundamentals of Real-Time ML.

- The biggest challenges facing Data teams.

- The motivation behind running Spark on Kubernetes

- Some challenges of Running Spark on Kubernetes and solutions

- Conclusion and takeaways

# The fundamentals of Real-Time ML.

# Finding the shortest, fastest Cycling Route

# Finding the shortest, fastest, least traffic?

# Finding the shortest, fastest, least traffic!

# Finding the shortest, fastest, least traffic!

# Finding the Least Air Pollution Exposure Cycling Routes

# Finding the Least Air Pollution Exposure Cycling Routes

## Air quality adjusted routing for cyclists and pedestrians

Authors: Sebastian Müller, Agnès Voisard  Authors Info & Claims

Check for updates

**ABSTRACT**

Air quality adjusted routing

environmental friendly ve

air quality levels. We use

provides the amount of F

air pollution. The data is a

Open Source Routing Ma

at geofabrik.de. Our app

transportation mode prof

88, Rue de Paris

Esplanade de La

À vélo (GraphHo

Inverser les dire

## Itinéraire

Distance: 14 km. Tem
Croissant: 164 m. Dé

↑ 1. Continuez sur
↑ 2. Restez sur la Paris
↑ 3. Restez sur la la Porte des Lila
↗ 4. Tournez légè Avenue de la Porte des Lilas
↰ 5. Tournez fort à droite
↗ 6. Tournez légèrement à droite sur Avenue de la Porte du Pré Saint-Gervais
↑ 7. Tournez à gauche
↰ 8. Tournez fort à gauche sur Avenue de la Porte du Pré Saint-Gervais

## EAI Endorsed Transactions
on Scalable Information Systems          Research Article  EAI.EU

### Predicting the least air polluted path using the neural network approach

K.Krishna Rani Samal[1,*], Korra Sathya Babu[1], Santos Kumar Das[1]

[1]National Institute of Technology, Rourkela, India

Abstract

Air pollution exposure during daily transportation is becoming a critical issue worldwide due to its adverse effect on human health. Predicting the least air polluted healthier path is the best alternative way to mitigate personal air pollution exposure risk. Computing the least polluted path for the current time might not be helpful for real-time applications. Therefore, we develop a routing algorithm based on a neural network-based CNN-LSTM-EBK (CLE), a temporal-spatial interpolation model. The proposed model predicts pollution levels at high temporal granularity. This paper introduces a weight function to compute air pollution concentration at the road network. It also predicts the least air polluted path among all possible paths from a source to a destination at different time granularity. The results show that the predicted path may be longer than the shortest route but minimize pollution exposure risk all the time, which proves its effectiveness during daily transportation.

DRIEAT - Direction des routes Île-de-France (DIRIF) - www.sytadin.fr

## Health-oriented routes for active mobility

Paulo J.G. Ribeiro, Gabriel J.C. Dias, José F.G. Mendes

Show more ▾

Add t  Conferences  >  2021 17th International Confe...

### Real-time Route Planning using Mobile Air Pollution Detectors and Citizen Scientists

**Publisher: IEEE**   Cite This   PDF

Richard O. Sinnott ; Yuan Wang ; Yiqun Wang  All Authors

**79**
Full
Text Views

Abstract

Document Sections

I. Introduction (Heading 1)

II. Related Work

III. Data Collection

IV. Implementation of the Routing Algorithm

V. Discussion

Show Full Outline ▾

Abstract:
The increasing urbanization of society is resulting in numerous challenges. One of these challenges is transport congestion and the associated increase in pollution that is widely accepted as driving global warming. For many individuals and especially those with respiratory issues, e.g., those with asthma or chronic obstructive pulmonary disease, high levels of pollution can cause direct health events. The ability to measure pollution accurately in real time at disaggregated levels and subsequently avoid pollution hotspots is thus highly desirable. This paper describes a Cloud-based infrastructure and associated mobile application that utilizes real time, mobile pollution measurement technology to help individuals avoid pollution hotspots through real-time pollution-aware routing algorithms.

## Couches de carte

Standard

CyclOSM

de transport

de Tracestack

nitaire

les superpositions pour
la carte

es de carte
ées de carte
les GPS publiques

Bonne  Moyenne  Dégradée  Mauvaise  Très mauvaise  Extrêmement mauvaise

# Finding the Least Air Pollution Exposure Routes (in real time)

# Finding the Least Air Pollution Exposure Routes (in real time)

Real-time machine learning: the application of machine learning models to generate predictions or decisions in real-time and adapt to the changing environment.

Data Sources

Real time AIQ Route

# Real-time Machine Learning Platform

The biggest challenges facing Data teams.

# Real-time Machine Learning Challenges



Real-time machine learning: challenges and solutions
Jan 2, 2022 • Chip Huyen

O'REILLY
Designing Machine Learning Systems
An Iterative Process for Production-Ready Applications
Chip Huyen

O'REILLY
Practical MLOps
Operationalizing Machine Learning Models
Noah Gift & Alfredo Deza

| | |
|---|---|
| Feature Engineering | Stream Processing |
| Incremental Learning (online learning) | Scalability |
| Model Updating | Latency |
| Model / Data Drift | Monitoring |
| Performance Evaluation | Distributed Training & Inference |
| MLOps | Resource Management / Cost |

# Real-time Machine Learning Challenges

Real-time machine learning challenges (our experience) are largely an infrastructure problem.

Stream Processing

Scalability

Latency

Monitoring

Distributed Training & Inference

Resource Management / Cost

# Solving some Real-time Machine Learning Challenges

Addressing these challenges requires a significant investment in advanced (OSS) technologies.

Spark on k8s:

- Stream processing

- Training

- Scalability & Latency

- Resource Efficiency

# The motivation behind running Spark on kubernetes

# Apache Spark is the #1 analytics engine for Big Data & AI

**Fast\***
Massively parallelizable, efficient read and write

**Easy**
Interfaces with well-known programming languages

**Versatile**
Across multiple use cases

| Object stores | Data warehouses | Streams | SQL/NoSQL databases |
|---|---|---|---|

**Spark** *APACHE*

| Python | Scala/Java | SQL |
|---|---|---|

| ETL/ELT | Real-time | ML | BI |
|---|---|---|---|

# The role of resource manager in a Spark cluster

Spark depends on cluster manager for orchestration of a job on a cluster

# Kubernetes is the latest cluster manager for Spark

Standalone: built-in, limited functionalities

Apache Mesos: deprecated as of Spark 3.2.0

Hadoop YARN: most widely used today

**Kubernetes: most popular among new deployments**

# The Spark on Kubernetes Journey

**Feb 2018 - Spark 2.3**
Initial support released for Spark on Kubernetes

**June 2020 - Spark 3.0**
Dynamic Allocation, Local code upload, Kerberos Support

**Oct 2021 - Spark 3.2**
Dynamic PVC mounting and reuse, Faster S3 Writes (Magic Committer enabled)

**Apr 2023 - Spark 3.4**
PVC-oriented executor pod allocation

**Nov 2018 - Spark 2.4**
Client Mode, Volume Mounts, PySpark and R support

**March 2021 - Spark 3.1**
Spark on Kubernetes generally available Graceful node shutdown, NFS mounts, Dynamic Persistent Volume Claims

**June 2020- Spark 3.3**
Executor Rolling in Kubernetes environment, Support Customized Kubernetes Schedulers

**Spark 3.5**
Upgrade kubernetes-client

# Spark on YARN: architecture & pain points



**Global Spark version and shared libraries**
- You'll have a Spark 2.4 cluster, a Spark 3.0 cluster, a Spark 3.1 cluster.
- Transient clusters are recommended for stability, but increase costs.

**Limited Docker image support***
- Environment is built from AMIs and bash scripts, flaky runtime library installation
- Debugging is painful - there's no way to run Spark locally, environment is subtle

**Resource Overhead**
- Slow startup time
- YARN master node, YARN Node Mgr are JVM processes using a lot of resources.

# Spark on Kubernetes: architecture & benefits



**Native Dockerization**
- Simpler dependency management
- Reliable executions across environments (locally during development, staging, production)
- Faster startup time

**A single long-running cluster**
- Quick to scale up (and down) based on load
- Mix different Spark versions
- Mix Spark and non-Spark apps
- Mix use cases (notebooks, batch/streaming jobs)

**A standard, agnostic infrastructure layer**
- Reduce lock in
- Simplify your operations
- Leverage the open-source tools from the cloud-native ecosystem

# Two ways to run Spark apps on k8s

## Spark-submit

- "Vanilla" way from Spark main open source repo

- Requires Spark distribution on client

- Configs spread between Spark config (mostly) and k8s manifests

- Less pod customization support (improving)

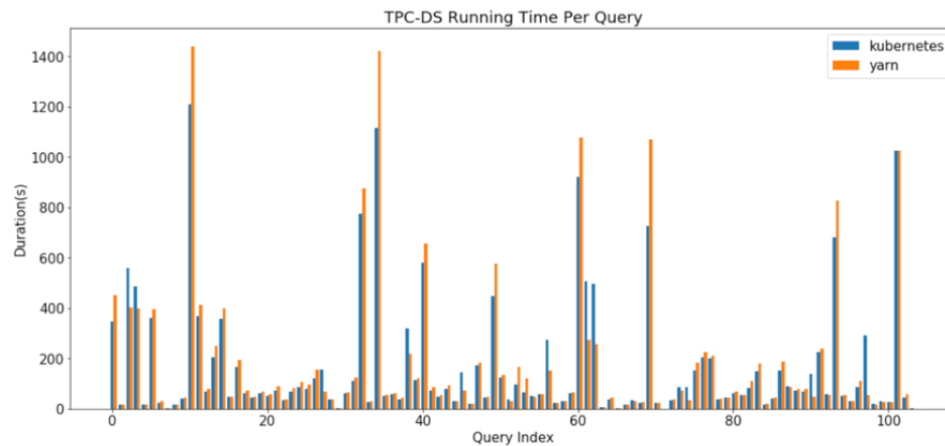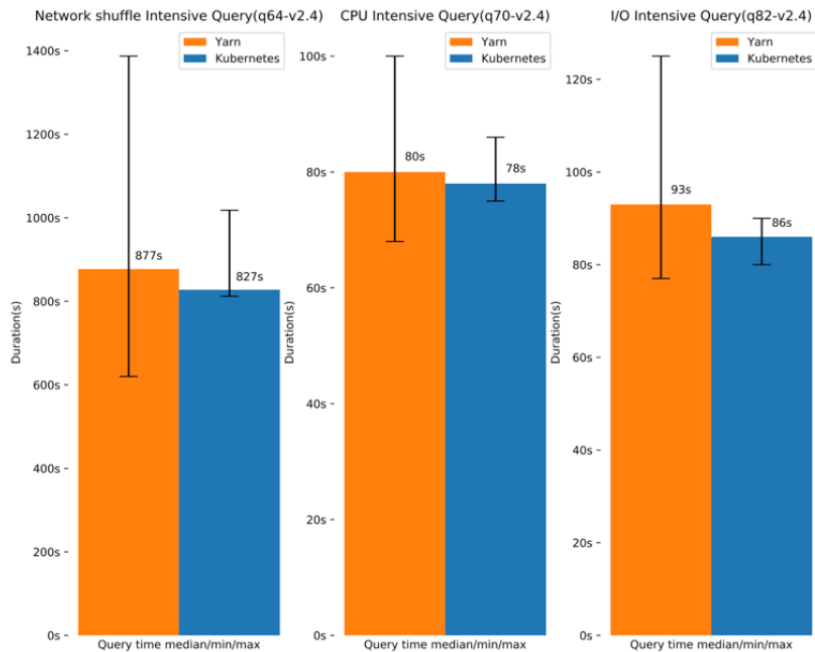- App management is more manual

## spark-on-k8s operator

**Overview**

The Kubernetes Operator for Apache Spark aims to make specifying and running Spark applications as easy and idiomatic as running other workloads on Kubernetes. It uses Kubernetes custom resources for specifying, running, and surfacing status of Spark applications. For a complete reference of the custom resource definitions, please refer to the API Definition. For details on its design, please refer to the design doc. It requires Spark 2.3 and above that supports Kubernetes as a native scheduler backend.

The Kubernetes Operator for Apache Spark currently supports the following list of features:

- Supports Spark 2.3 and up.
- Enables declarative application specification and management of applications through custom resources.
- Automatically runs `spark-submit` on behalf of users for each `SparkApplication` eligible for submission.
- Provides native cron support for running scheduled applications.
- Supports customization of Spark pods beyond what Spark natively is able to do through the mutating admission webhook, e.g., mounting ConfigMaps and volumes, and setting pod affinity/anti-affinity.
- Supports automatic application re-submission for updated `SparkApplication` objects with updated specification.
- Supports automatic application restart with a configurable restart policy.
- Supports automatic retries of failed submissions with optional linear back-off.
- Supports mounting local Hadoop configuration as a Kubernetes ConfigMap automatically via `sparkctl`.
- Supports automatically staging local application dependencies to Google Cloud Storage (GCS) via `sparkctl`.
- Supports collecting and exporting application-level metrics and driver/executor metrics to Prometheus.
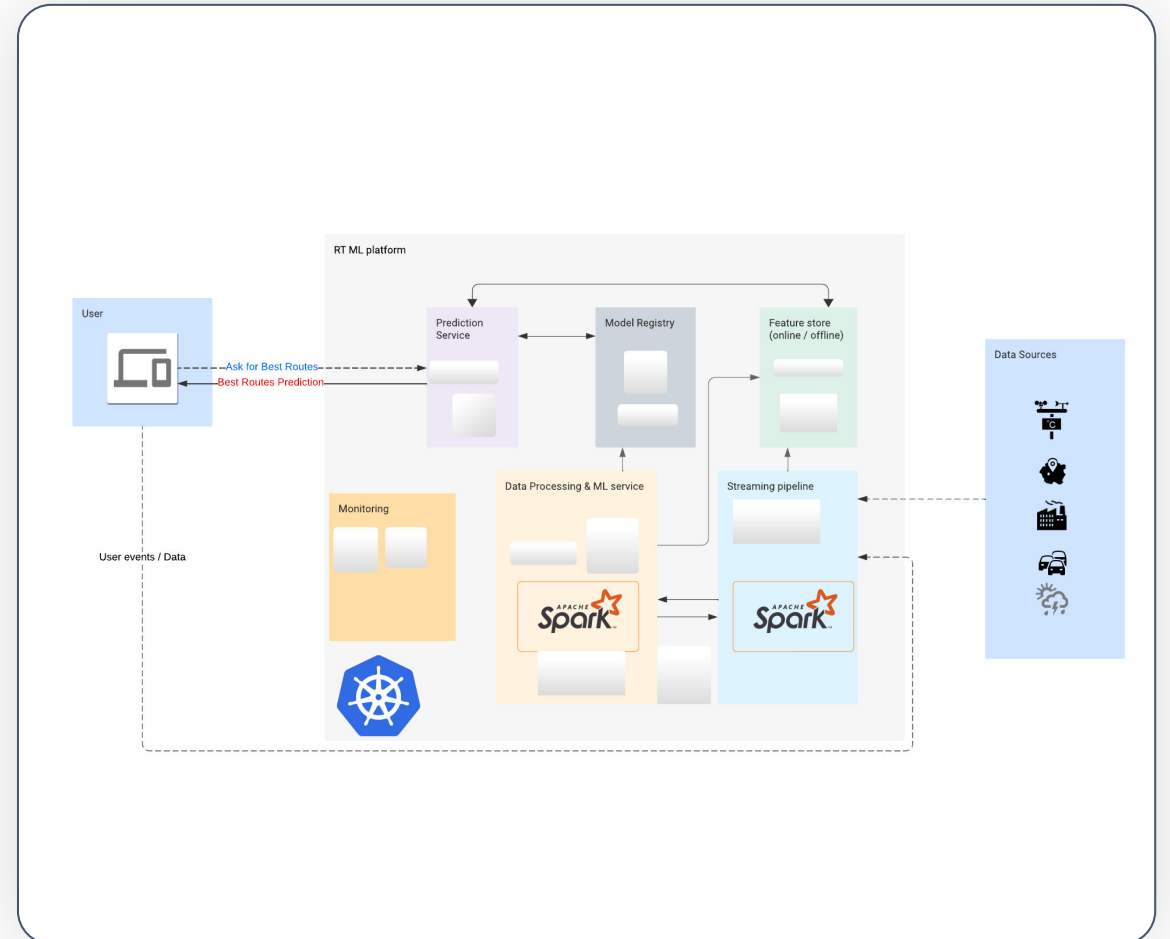
TPC-DS 1T Benchmark

https://aws.amazon.com/fr/blogs/containers/optimizing-spark-performance-on-kubernetes/

# Some challenges of running Spark on kubernetes and solutions

# Challenges in the context of R-L M-L

- Monitoring

- Scalability

- Latency

- Models Training

# Monitoring: logs, logs and more logs

Key information is buried under a lot of noisy one.

- Spark event/driver/executor logs.
- kubernetes logs
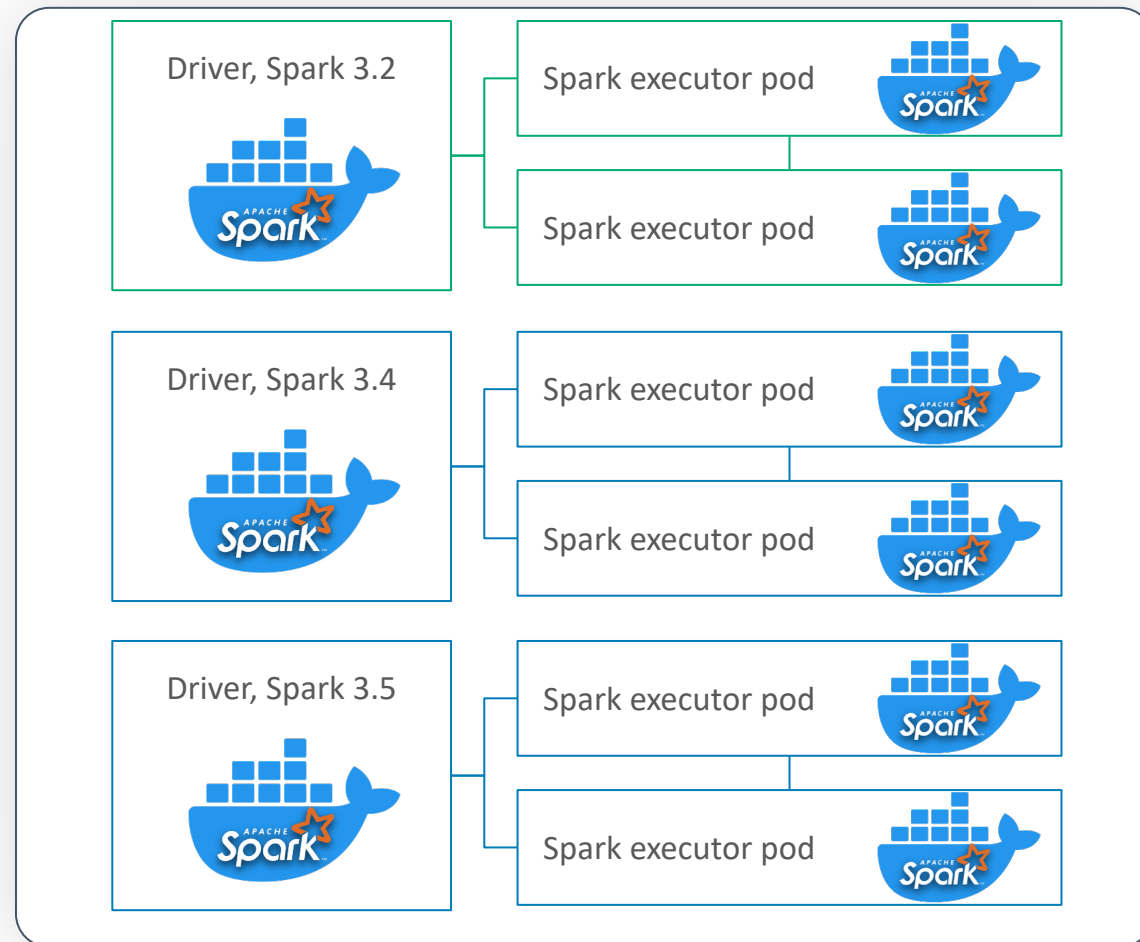- Hard to reconcile with Spark jobs/stages/tasks

Solution

- Logs shipping tools : fluentbit & logstach
- Prometheus: Spark has a built-in Prometheus
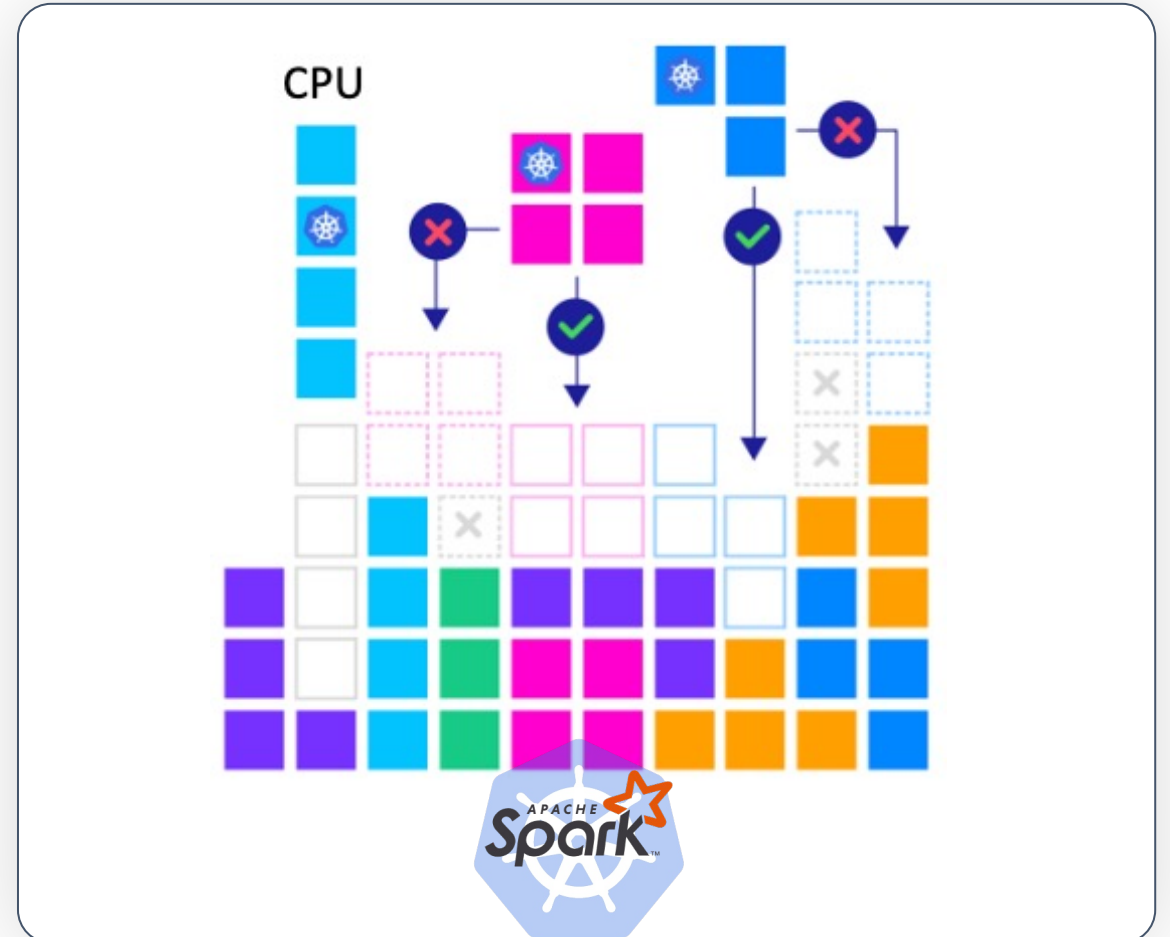
# Scalability

Key factors to consider

- Cluster sizing, infra choice/specs.

- Dynamic Allocation

- Shuffle data (NO external shuffle service YET )

# Scalability: the right sizing

For the sizing:

- Continuous and repeated exercise: know your data sources
- When selecting the cluster focus on enhancing parallelism in relation to the source.
- Streaming is CPU-bounded, State matters too (avoid spills)
- Deep Learning models with relatively long training and inference time: mix CPU with GPU (when required).

# Scalability: Dynamic Allocation ( A two-sided problem)

## Dynamic Allocation in Spark Structured Streaming

- Designed for batch jobs, it is compatible with batch and Spark structured streaming. Works poorly for certain applications !

## Dynamic Allocation (within k8s )

- This feature **may** cause issues with Spark Scalability on k8s !

# Scalability: Shuffle data

External shuffle service for Spark on kubernetes is not supported yet. There are 4 options (+1):

- Cloud Shuffle Storage Plugin for Apache Spark - AWS Glue

- IBM/spark-s3-shuffle: Shuffle plugin for Apache Spark and S3 compatible service

- GitHub - oap-project/remote-shuffle: Spark shuffle plugin for support shuffling data through a remote Hadoop-compatible file system (Intel)

- Apache Spark on Kubernetes - Local Storage (main project)

- AWS S3 CSI driver and High-Performance Storage – S3 Express One Zone, AWS

---

**Future Work**

There are several Spark on Kubernetes features that are currently being worked on or planned to be worked on. Those features are expected to eventually make it into future versions of the spark–kubernetes integration.

Some of these include:

- External Shuffle Service
- Job Queues and Resource Management

---

IBM Spark S3 Shuffle plugin is our choice:

- Supports Spark versions 3.2 to 3.4. Successfully tested with Spark 3.5.

- To support different cloud vendors, the corresponding Hadoop connector needs to be added to the classpath.
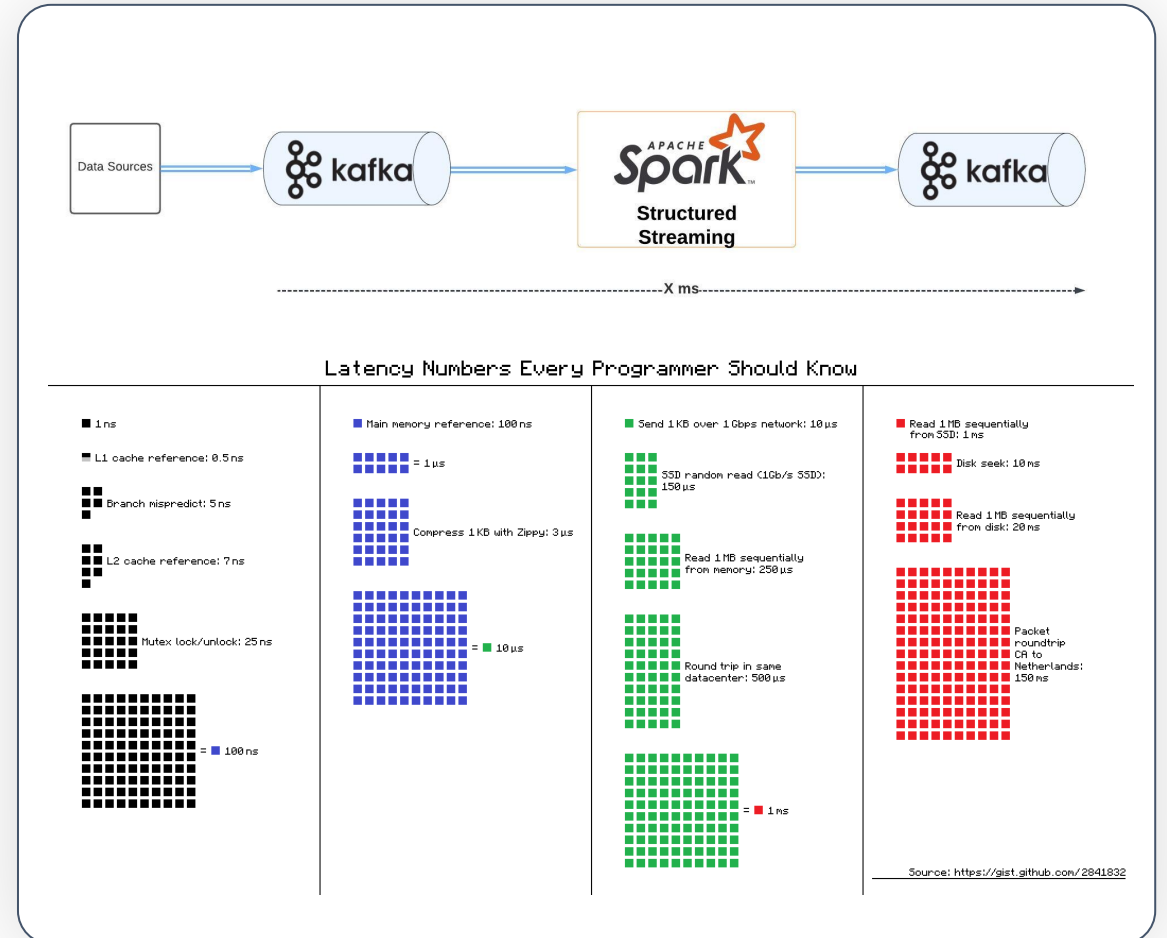
# Latency

## Key factors to consider

- Spark configurations.
- Sub-second latency expectations are a challenge
- Stateful vs Stateless Pipelines

## Consider the following:

- Use only simple computations involving data transformation or enrichment !
- Always use a message bus (e.g., Apache Kafka or Apache Pulsar) and fast key-value stores (e.g., Apache Cassandra or Redis)
- RocksDB state store provider

# M-L Training

Key factors to consider

- Most ML frameworks were designed for single-node environments

- Spark MLlib is lagging behind !

# M-L Training

Consider the following:

- Use TensorFlow, Keras, and PyTorch

- Accelerator, Distributed ML & GPU :
  - Horovod
  - NIVIDIA RAPIDS Accelerator for Spark



Source: RAPIDS Accelerator for Apache Spark (NVIDIA)

# M-L Training (Batch)

With the default scheduler:  workloads experience
higher rates of resource starvation, leading to
performance degradation or failure !

Solutions for the default scheduler:

- Use custom k8s Scheduler support. Enabling
  YARN-like capabilities such as queue, gang
  scheduling, etc
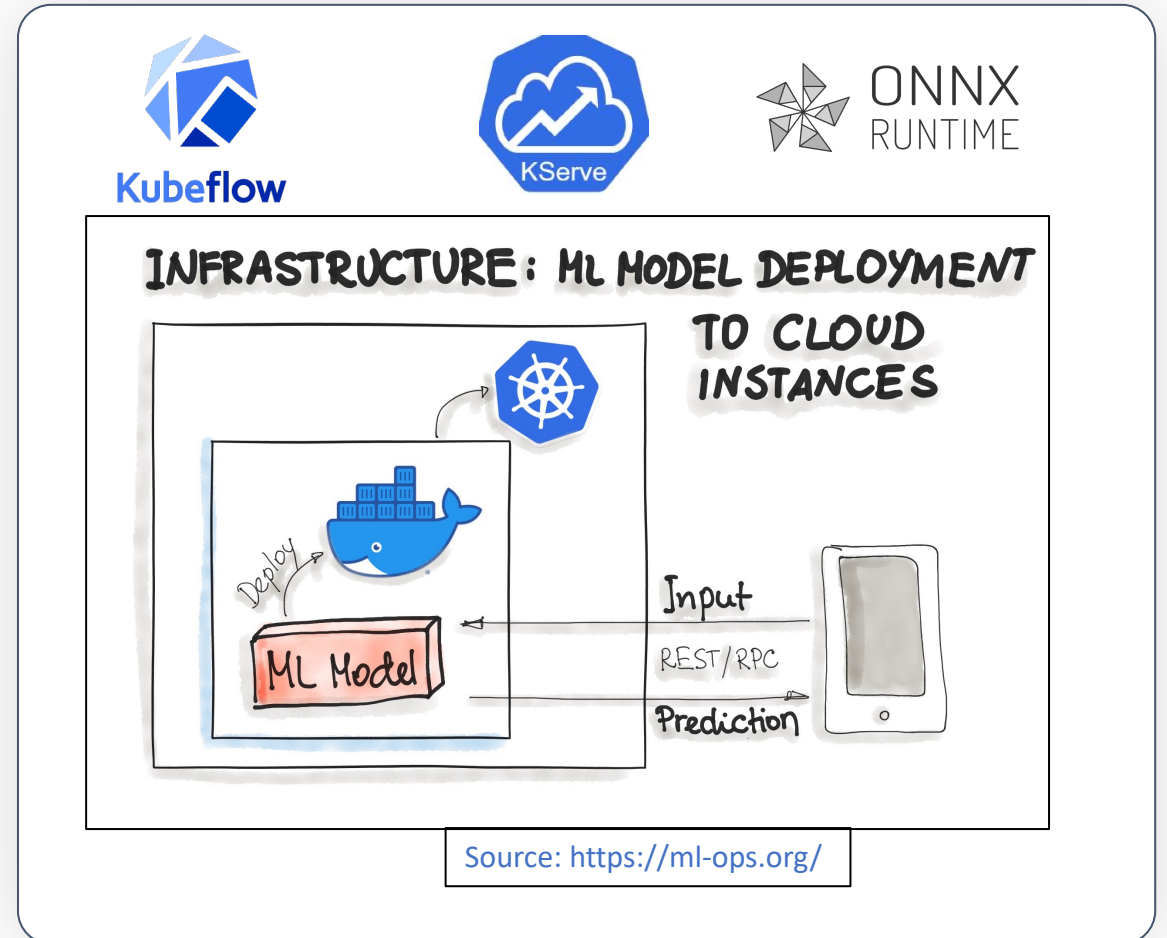
# M-L Deployment & Serving (WIP)

## Key factors to consider

- Package the whole ML tech stack (dependencies) and the code for ML model prediction into a Docker container.

- Model optimization and Model compression



Source: https://ml-ops.org/

# Conclusion and takeaways

# Why Spark on k8s integration is Important for R-L M-L?

- Native Integration

- k8s best practices apply to Spark on k8s for free!

- Scalability, Latency, Fault-tolerance

- Models Training and Serving

- Integration with a rich ecosystems

# Key takeaways

- Use the k8s Spark operator

- Design and build your logging and monitoring stack

- Keep adhering to Spark best practices compatible with k8s

- Use the rich k8s (Monitoring, Mlops, etc) ecosystem

- Contribute to the OSS (share your experiences, code, ideas, challenges )

- Keep Cycling …

Thank you