# 04 - Spark Development Process

Demi Ben-Ari
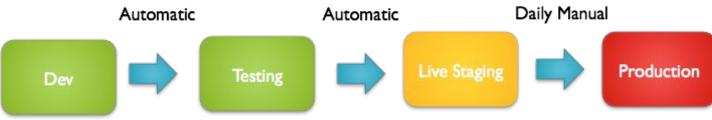
**Panorays**

# Running environments

- **Development – Testing – Production**
  - Don't you need more?
    Be as flexible as you can
- <u>Cluster Utilization</u>
  - **Unified Cluster for all environments** vs. **Cluster per Environment**
  - (Cluster per Data Center)
- Configuration
  - Local Files vs. Distributed

**Panorays**

# DevOps - Keep It Simple Stupid!

- Linux
  - Bash scripts
  - Crontab
- Automation via Jenkins
- Continuous Deployment – with every GIT push



**Panorays**

# Build Automation

- ☐Maven    Sonatype Nexus

  - Sonatype Nexus artifact management

- Jenkins

  - Deploy and Script generation scripts

  - Per Environment Testing

  - Data Validation

  - Scheduled Tasks

**Panorays**

# Workflow Management

- Oozie – Very hard to integrate with Spark
  - XML configuration based and not that convenient
- Azkaban (Haven't tried it)

Chosen:

- Luigi
- Crontab + Jenkins (KISS again)

**Panorays**

# Testing

- **☐Unit**
  - JUnit tests that run on the Spark "Functions"
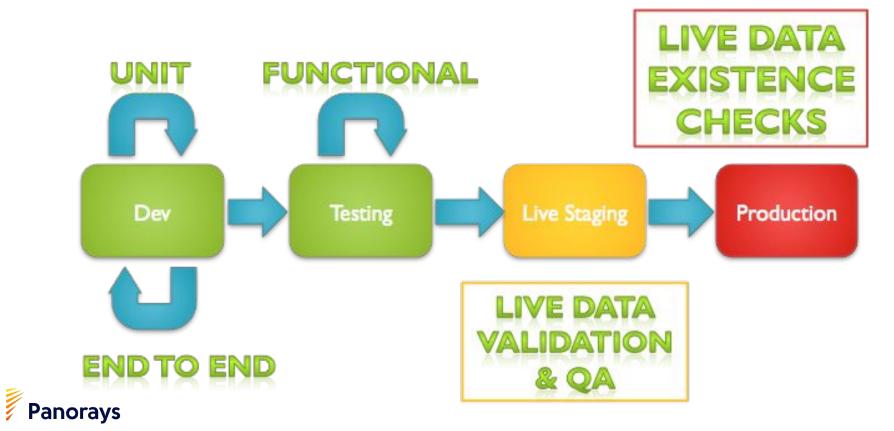- **End to End**
  - Simulate the full execution of an application on a single JVM (local mode) – Real input, Real output
- **Functional**
  - Stand alone application
  - Running on the cluster
  - Minimal coverage – Shows working data flow

**Panorays**

# Testing

# Logging

- ☐Runs by default log4j (slf4j)
- How to log correctly:
  - Separate logs for different applications
  - Driver and Executors log to different locations
  - Yarn logging also exists (Might find problems there too)
- ELK Stack (Logstash - ElasticSearch – Kibana)
  - By Logstash Shippers (Intrusive) or UDP Socket Appender (Log4j2)
  - DO NOT use the regular TCP Log4J appender

**Panorays**

# History Server

- ☐Can be run on all Spark deployments:
    - Stand Alone, YARN, Mesos
- Integrates both with YARN and Mesos
    - ☐In Yarn / Mesos, run history server as a daemon.
- [Documentation](#)

**Panorays**

# Monitoring

- [Graphite](#)
  - Online application metrics
- [Grafana](#)
  - Good Graphite visualization
- Jenkins - Monitoring
  - Scheduled tests
  - Validate result set of the applications
  - Hung or stuck applications
  - Failed application

**Panorays**

# Reporting and Monitoring
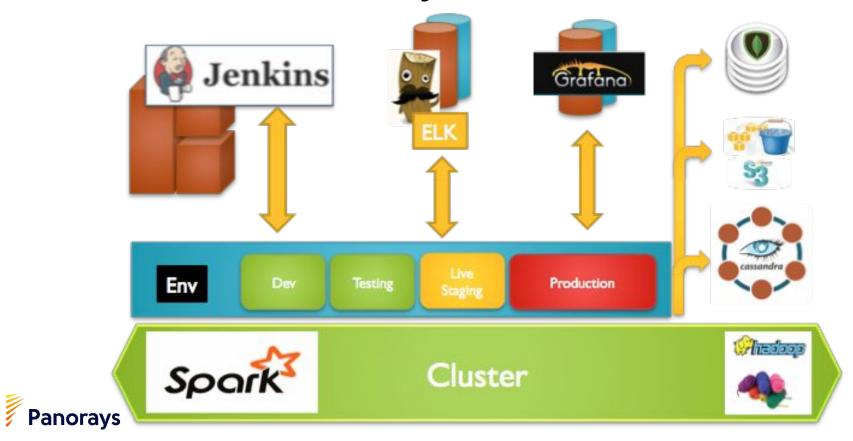
- Grafana + Graphite (InfluxDB) - Example

# More Spark Monitoring Abilities

- Master REST API
- Metrics Library
  - Can be plugged to different data sinks
- OS
  - dstat
  - iostat
  - iotop
- JVM
  - jstack

**Panorays**

# Environment Summary

# Questions?