Continuous Integration for Spark Apps



Continuous Integration for Spark Apps

because testing things is pretty important.



Hi, I'm Sean!



Sean McIntyre

Software Architect

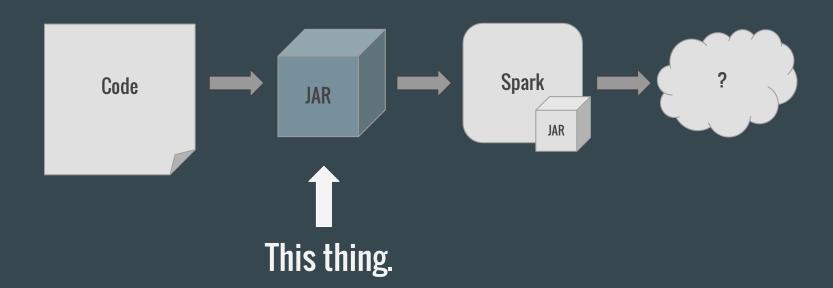
- 2 Berkeley Street Suite 600 Toronto ON M5A 4J5
- **&** 416 203 3003 x 377
- 416 906 7894
- @ smcintyre@uncharted.software
- www.uncharted.software

I'm writing
Spark libraries at
Uncharted Software

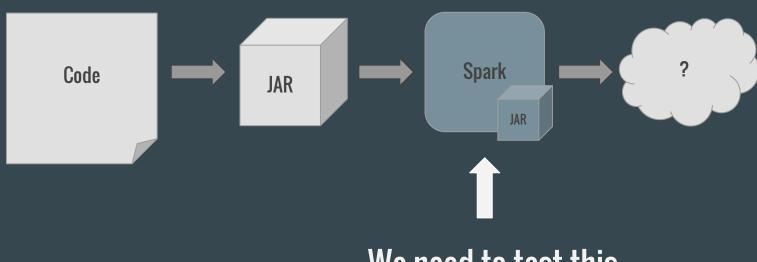


Why is it hard to test Spark apps?

What is a Spark app?

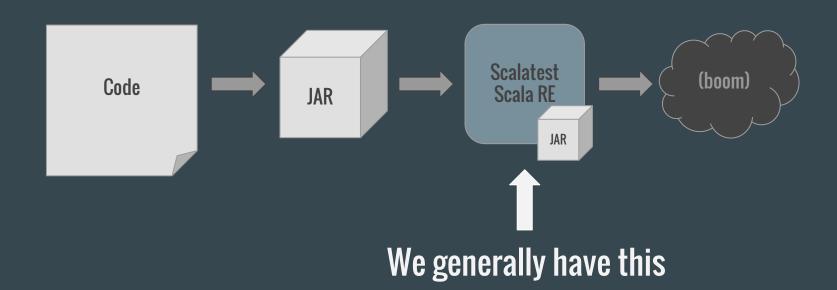


And...

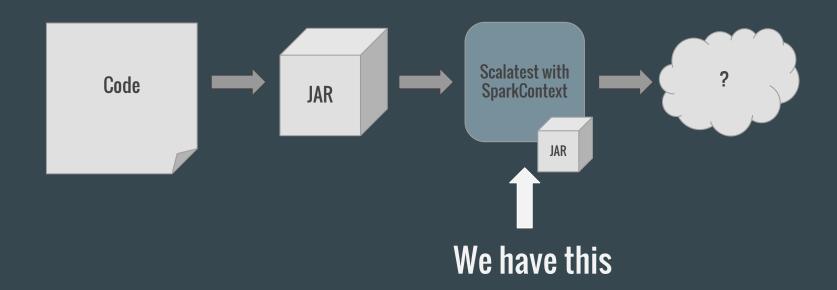


We need to test this

But...

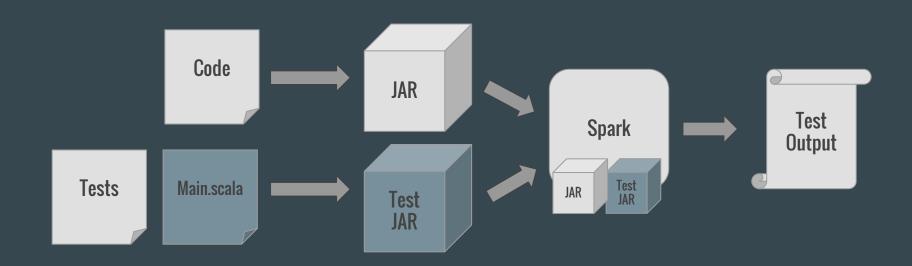


Or at best...



Can we test a Spark app?

Step 1: Squish Scalatest into Spark



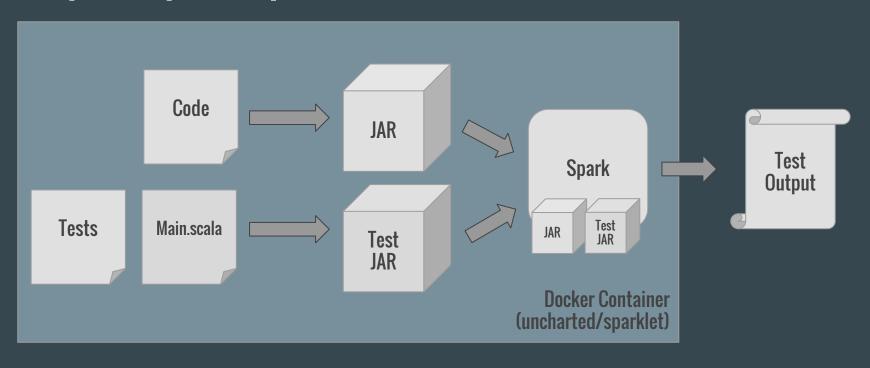
Main.scala

```
import org.scalatest.tools.Runner

object Main {
    def main(args: Array[String]): Unit = {
        val testResult = Runner.run(Array("-o", "-R", "build/classes/test"))
        if (!testResult) {
            System.exit(1) // exit with an error code if a test failed
        }
    }
}
```

Continuous Integration?

Step 2: Squish Spark and Test JAR into Docker



test-environment.sh

```
docker run \
  -e GRADLE_OPTS="-Dorg.gradle.daemon=true" \
  -v $(pwd)/src/test/resources/log4j.properties:/usr/local/spark/conf/log4j.properties \
  -v $(pwd):/opt/mycode \
  -it \
  --workdir="/opt/mycode" \
  uncharted/sparklet:1.5.2 bash
```

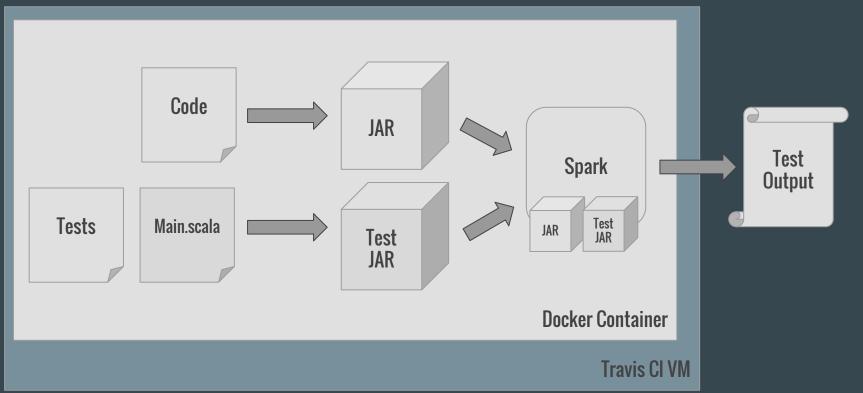
test-environment.sh

```
docker run \
-e GRADLE OPTS="-Dorg.gradle.daemon=true" \
-v $(pwd)/src/test/resources/log4j.properties:/usr/local/spark/conf/log4j.properties \
-v $(pwd):/opt/mycode \
-it \
--workdir="/opt/mycode" \
uncharted/sparklet:1.5.2 bash
```

build.gradle

```
// create an Exec task to test code via spark-submit
task test(overwrite: true, type: Exec, dependsOn: [jar, testJar]) {
  executable = 'spark-submit'
 args = [
    //the --packages flag allows us to place our dependencies on the path
    "--packages", "org.scalatest:scalatest ${scalaBinaryVersion}:2.2.5",
   //the --jars flag allows us to place our code on the path
    "--jars", "/opt/${repoName}/build/libs/${artifactName}-${version}.jar",
    //Main, from before
    "--class", "my.company.project.Main",
    "build/libs/${artifactName}-${version}-tests.jar"
```

Step 3: Squish Docker container into Travis Cl



.travis.yml

sudo: required
language: bash

```
- docker
 - env | grep TRAVIS > travis.env
 - echo "GRADLE OPTS=-Dorg.gradle.daemon=false" >> travis.env
 - echo "CI NAME=travis ci" >> travis.env
 - echo "CI=true" >> travis.env
 - echo "TRAVIS=true" >> travis.env
 - echo "CONTINUOUS INTEGRATION=true" >> travis.env
 - echo "DEBIAN FRONTEND=noninteractive" >> travis.env
 - echo "HAS JOSH K SEAL OF APPROVAL" >> travis.env
 - docker run --env-file travis.env -v $(pwd)/src/test/resources/log4j.properties:/usr/local/spark/conf/log4j.properties -v $(pwd):
/opt/salt --rm --entrypoint=/startup.sh --workdir="/opt/salt" uncharted/sparklet:1.5.1 ./gradlew coverage coveralls
```

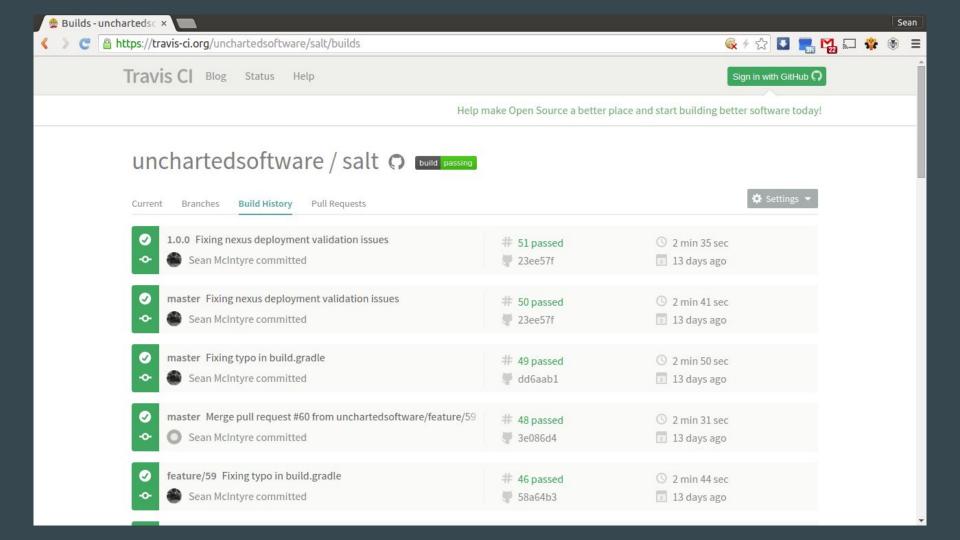
Voila!

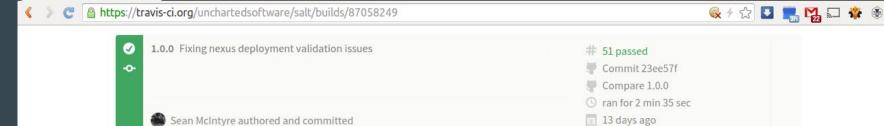
build passing

Examples

Check out Salt (https://github.com/unchartedsoftware/salt) for a working example.

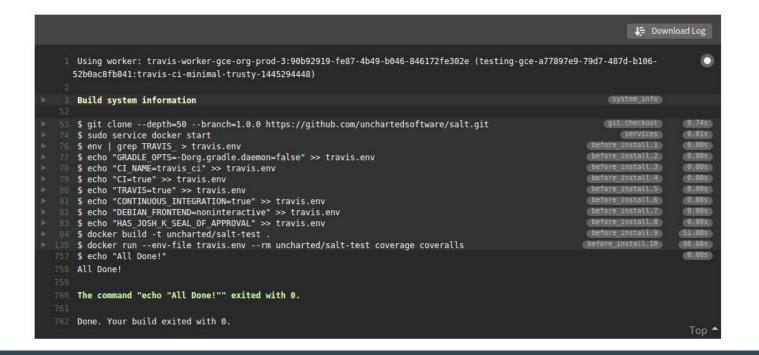
Salt's build is also connected to Coveralls for code coverage reports :)





Sean

Build #51 - uncharte ×



Coming soon...

- Non-linear Spark Pipeline
- Sparkplug

http://blog.uncharted.software

Questions?

