

# CI/CD в большом on-premise Datalake-проекте

Никита Благодарный  
telegram @nblagodarnyy

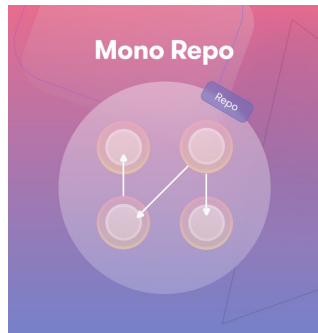
Александра Чекмарева (Китченко)  
telegram @sasha\_kitchen



# Немного вводных

# Про что этот доклад?

## Как структурировать и делать CI/CD

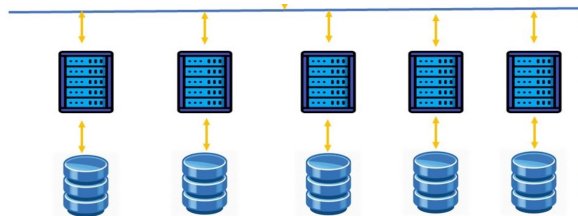
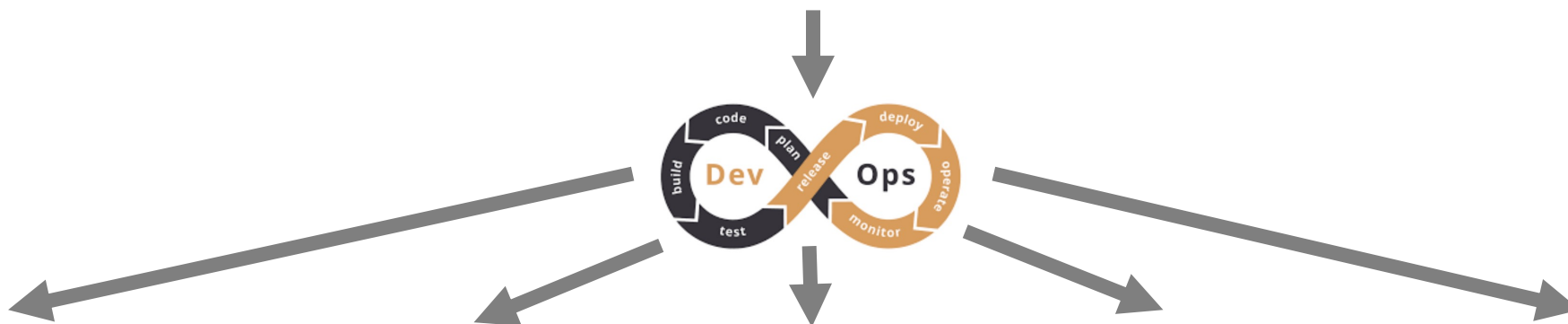


Statistic Statistic

Refresh Refresh on selection Settings

Overview <> properties <> py <> txt <> xml

Extension	Lines CODE	Count
sql (SQL files)	150268	2804x
scala (SCALA files)	98141	1627x
yaml (YAML files)	79769	2225x
py (PY files)	37789	427x
ddl (DDL files)	18012	94x





## Data Engineer, 17 лет в IT, из них 14 в области BI/DE

- 2007-2011 Сбербанк, разработка OLTP-решений на Oracle, ХД/Reporting на стеке Microsoft
- 2012-2016 Сбербанк, разработка ЦХД на Oracle / Teradata / Informatica PC
- 2016-2018 КРОК, проекты на Oracle / Informatica / Hadoop (HDP, Arenadata)/Spark/Hive/NiFi
- 2018-2019 Газпромнефть, построение DataLake (Arenadata, MS SQL, Informatica, Spark/Hive)
- 2019-.... ЦРПТ, DataLake (Vanilla Hadoop/Spark/Hive/HBase/ClickHouse/Postgres/etc)



## DevOps Engineer, 10 лет в IT, последние 6 лет работаю DevOps/MLOps на проектах Big Data

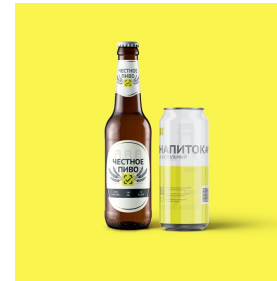
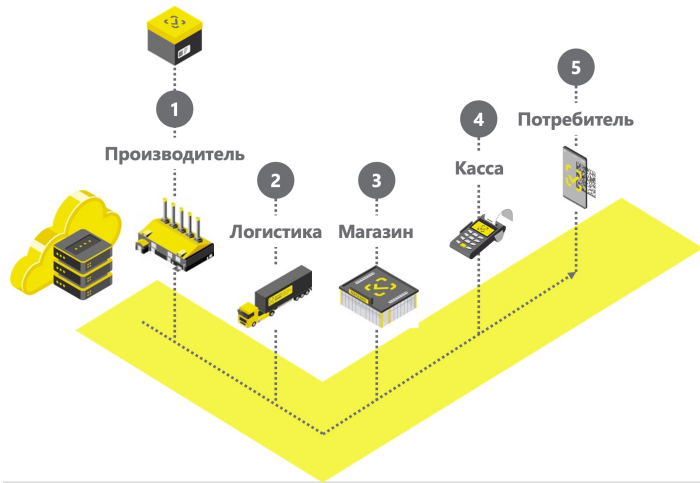
- 2018-2021 TELE2, DevOps, разработка решений машинного обучения (CI/CD, gitlab, Jenkins, airflow, kubernetes, python, sql, pyspark, scala, hive, sqoop, bteq, Hadoop, Teradata)
- 2021-2023 ЦРПТ (Честный Знак), DevOps Аналитической платформы (ansible, gitlab ci, kubernetes, helm, bash, python, maven, Hadoop, ClickHouse, Postgres, HAproxy )
- 2023-... НЛМК, MLOps, разработка DSML-платформы (CI/CD, gitlab, Minio, S3, airflow, kubernetes, helm, python, Hadoop, Computer Vision, Machine learning)



# О компании Честный Знак

## Создаем систему цифровой маркировки и прослеживания товаров "Честный Знак" в России и странах ЕАЭС

- Основана в 2017 году
- > 750 сотрудников
- Офисы разработки в Москве, Питере, Пензе, Ульяновске, Нижнем Новгороде
- Разнообразные товарные группы (Табак, Молоко, Вода, Пиво, Лекарства, Одежда, Обувь, ...)



# О технологическом ландшафте

APACHE PHOENIX

APACHE HBASE



cassandra

4 кластера (180 машин)  
~5 Pb хранения



Кластер 410 машин  
146 Tb RAM  
9,5 Pb хранения



ClickHouse

Кластер 74 машины  
800 Tb хранения



PostgreSQL



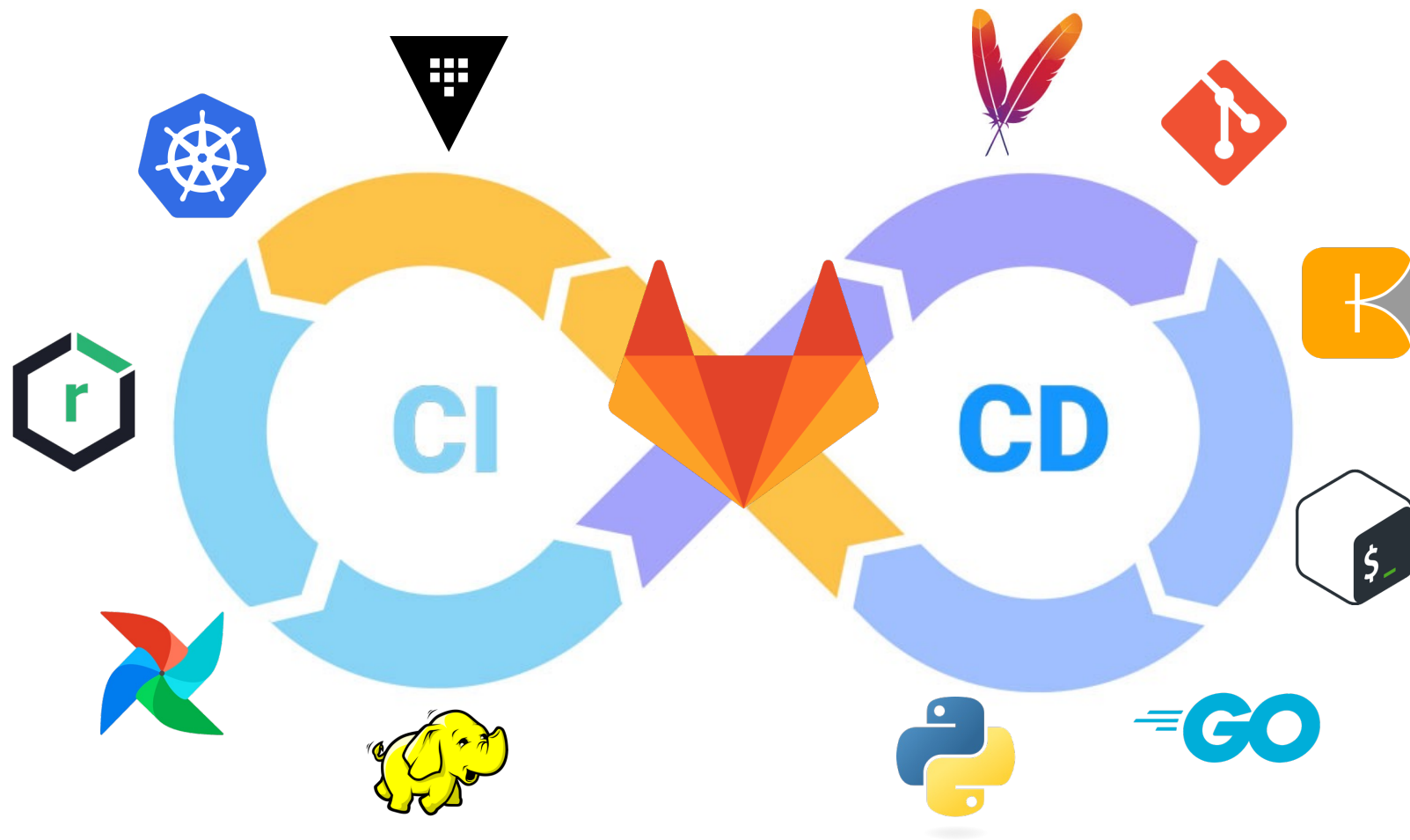
2 кластера Patroni  
~2 Tb хранения



Кластер 18 машин



# О технологическом стеке CI/CD



# Структура проекта в Git

# Структура проекта

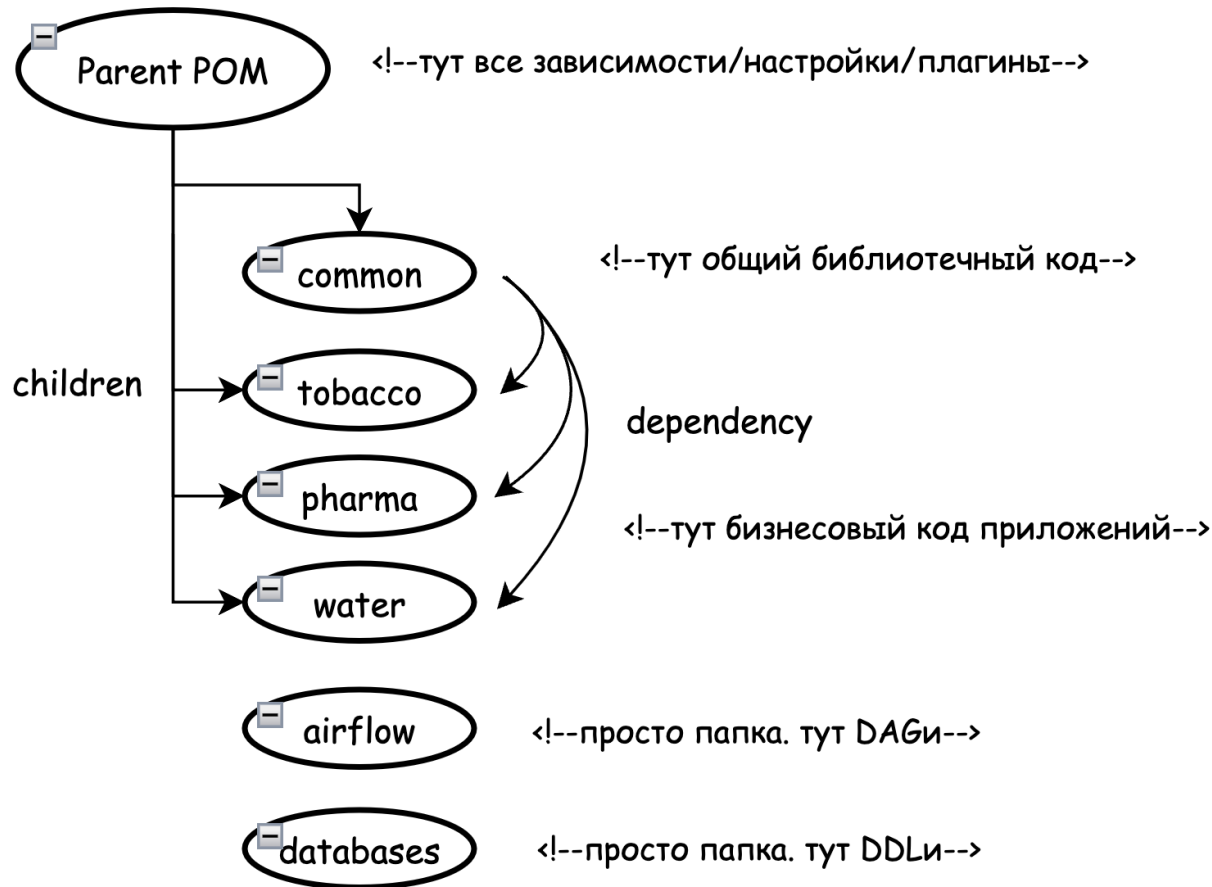




# Структура проекта – иерархия/зависимости



Модель проекта



```
<?xml version="1.0" encoding="UTF-8"?>
<project xmlns="http://maven.apache.org/POM/4.0.0"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:schemaLocation="http://maven.apache.org/POM/4.0.0 http://
  <modelVersion>4.0.0</modelVersion>

  <groupId>ru.crpt.analytics</groupId>
  <artifactId>datalake_etl</artifactId>
  <packaging>pom</packaging>
  <version>${revision}</version>
  <modules>
    <module>common</module>
    <module>clean-data</module>
    <module>tobacco</module>
    <module>pharma</module>
    <module>service_providers</module>
    <module>legprom</module>
    <module>milk</module>
    <module>water</module>
    <module>furs</module>
    <module>master-data</module>
  </modules>
```

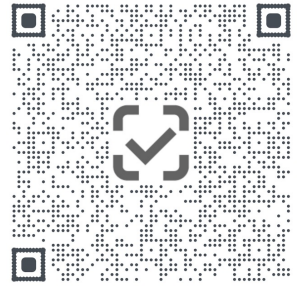
# Структура проекта – parent POM



```
m pom.xml (datalake_etl) ×
2 <project xmlns="http://maven.apache.org/POM/4.0.0"
81 <properties>
95 <!--Scala-->
96 <scala.version.short>2.12</scala.version.short>
97 <scala.version>2.12.15</scala.version>
98 <scala.compile.target.jvm>jvm-1.17</scala.compile.target.jvm>
99 <scalatest.version>3.2.18</scalatest.version>
100 <scalatest.maven.plugin.version>2.2.0</scalatest.maven.plugin.version>
101 <scoverage.plugin.version>2.0.4</scoverage.plugin.version>
102 <scoverage.scalacPluginVersion>2.0.11</scoverage.scalacPluginVersion>
103 <scoverage.aggregate>>true</scoverage.aggregate>
104 <scalastyle.maven.plugin.version>1.0.0</scalastyle.maven.plugin.version>
105
```

```
6 mvn clean compile -Dscala.version="2.12.14"
```

# Структура проекта – parent POM



```
<dependencyManagement>
  <dependencies>
    <dependency>
      <groupId>org.apache.spark</groupId>
      <artifactId>spark-core_${scala.version.short}</artifactId>
      <version>${spark.version}</version>
      <type>test-jar</type>
    </dependency>
    <dependency>
      <groupId>org.apache.spark</groupId>
      <artifactId>spark-catalyst_${scala.version.short}</artifactId>
      <version>${spark.version}</version>
      <type>test-jar</type>
    </dependency>
    <dependency>
      <groupId>org.apache.spark</groupId>
      <artifactId>spark-sql_${scala.version.short}</artifactId>
      <version>${spark.version}</version>
      <type>test-jar</type>
    </dependency>
    <dependency>
      <groupId>com.github.mrpowers</groupId>
      <artifactId>spark-fast-tests_${scala.version.short}</artifactId>
      <version>${com.github.mrpowers.version}</version>
    </dependency>
  </dependencies>
</dependencyManagement>
```

- Централизация версий зависимостей
- Явное указание версии для транзитивных зависимостей

A second, and very important use of the dependency management section is to control the versions of artifacts used in transitive dependencies. As an example consider these projects:

Project A:

1. `<project>`
2. `<modelVersion>4.0.0</modelVersion>`

# Структура проекта – parent POM



```
m pom.xml (datalake_etl) x
2 <project xmlns="http://maven.apache.org/POM/4.0.0"
450 <build>
478 <plugins>
556 <!-- enable scalatest -->
557 <plugin>
558 <groupId>org.scalatest</groupId>
559 <artifactId>scalatest-maven-plugin</artifactId>
560 <version>${scalatest.maven.plugin.version}</version>
561 <configuration>
562 <reportsDirectory>${project.build.directory}/surefire-reports</reportsDirectory>
563 <junitxml>./</junitxml>
564 <filereports>TestSuite.txt</filereports>
565 <parallel>>false
```

```
m pom.xml (datalake_etl) x
2 <project xmlns="http://maven.apache.org/POM/4.0.0"
704 <repositories>
712 <repository>
713 <id>java-dev-team</id>
714 <name>crpt-nexus</name>
715 <url>${nexus.repo.url}/repository/java-dev-team</url>
716 <layout>default</layout>
717 </repository>
718
719 </repositories>
720
721 <distributionManagement>
722
723 <!--Release repo-->
724 <repository>
725 <uniqueVersion>>false</uniqueVersion>
726 <id>crpt-nexus-release</id>
727 <name>Common nexus release repository</name>
728 <!--suppress UnresolvedMavenProperty -->
729 <url>${nexus.repo.url}/repository/analytics-team-release</url>
730 </repository>
731
732 <!--Snapshot repo-->
733 <snapshotRepository>
734 <uniqueVersion>>false</uniqueVersion>
735 <id>crpt-nexus-snapshots</id>
736 <name>Common nexus snapshot repository</name>
737 <!--suppress UnresolvedMavenProperty -->
738 <url>${nexus.repo.url}/repository/java-dev-team-snapshot</url>
739 </snapshotRepository>
740 <downloadUrl>${nexus.repo.url}/repository/java-dev-team-snapshot</downloadUrl>
741 </distributionManagement>
```

# Структура проекта – child POM

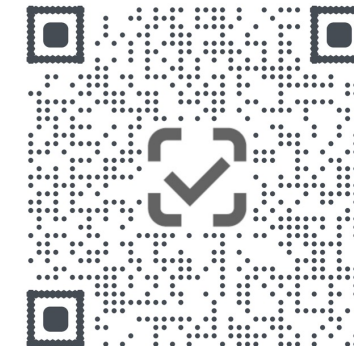


```
m pom.xml (datalake-common-test) x
2 <project xmlns="http://maven.apache.org/POM/4.0.0"
5 <parent>
7 <groupId>ru.crpt.analytics</groupId>
8 <version>${revision}</version>
9 </parent>
10 <modelVersion>4.0.0</modelVersion>
11
12 <properties>
13 <testcontainers.scala.version>0.40.9</testcontainers.scala.version>
14 <kafka.client.version>2.7.0</kafka.client.version>
15 <phoenix.thirdparty.version>1.1.0</phoenix.thirdparty.version>
16 <com.github.mrpowers.version>0.23.0</com.github.mrpowers.version>
17 <assembly.skipAssembly>>true</assembly.skipAssembly>
18 </properties>
```

```
m pom.xml (datalake-gismt) x
2 <project xmlns="http://maven.apache.org/POM/4.0.0" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
69 <build>
70 <plugins>
71 <plugin>
72 <groupId>org.codehaus.mojo</groupId>
73 <artifactId>build-helper-maven-plugin</artifactId>
74 <version>3.5.0</version>
75 <executions>
76 <execution>
77 <id>add-source</id>
78 <phase>generate-sources</phase>
79 <goals>
80 <goal>add-source</goal>
81 </goals>
```



# Структура проекта – CI-friendly



```
m pom.xml (datalake_etl) x
1  <?xml version="1.0" encoding="UTF-8"?>
2  m↓ <project xmlns="http://maven.apache.org/POM/4.0.0"
3      xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
4      xsi:schemaLocation="http://maven.apache.org/POM/4.0.0 http:
5      <modelVersion>4.0.0</modelVersion>
6
7      <groupId>ru.crpt.analytics</groupId>
8      <artifactId>datalake_etl</artifactId>
9      <packaging>pom</packaging>
10     ⚡ <!--Don't change it -->
11     <version>${revision}</version> Blagodarnyi, 18.12.2020, 15:08
12     <modules>
13         <module>common</module>
14         <module>clean-data</module>
15         <module>tobacco</module>
16         <module>pharma</module>
```

Starting with Maven 3.5.0-beta-1 you can use the `${revision}`, `${sha1}` and/or `${changelist}` as placeholders for the version in your pom file.

```
mvn -Drevision=1.0.0-SNAPSHOT clean package
```

# Разные версии Spark

# Структура проекта – разные версии Spark



Overview Active Stale **All**

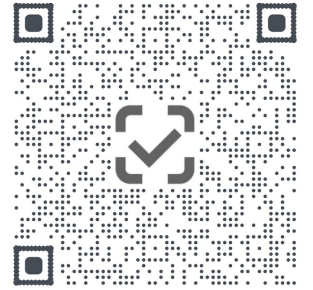
master

master  default  protected

master\_spark\_3.5



# Структура проекта – разные версии Spark



```
m pom.xml (datalake_etl) x
2   <project xmlns="http://maven.apache.org/POM/4.0.0"
44   <profiles>
45     <profile>
46       <id>spark-3.3</id>
47       <activation>
48         <activeByDefault>true</activeByDefault>
49       </activation>
50       <properties...>
67     </profile>
68
69     <profile>
70       <id>spark-3.5</id>
71       <activation>
72         <activeByDefault>>false</activeByDefault>
73       </activation>
74       <properties>
75         <!--Spark-->
76         <spark.version>3.5.1</spark.version>
77         <spark.short.version>3.5</spark.short.version>
78         <!--Spark catalyst for compile different versions-->
79         <org.codehaus.janino.version>3.1.9</org.codehaus.janino.version>
80         <spark.xml.version>0.17.0</spark.xml.version>
81         <!--Enforced versions for Spark-->
82         <fasterxml.jackson.version>2.15.2</fasterxml.jackson.version>
83         <fasterxml.jackson.module.scala.version>2.15.2</fasterxml.jackson.module.scala.version>
84         <fasterxml.jackson.databind.version>2.15.2</fasterxml.jackson.databind.version>
85         <fasterxml.jackson.yaml.version>2.15.2</fasterxml.jackson.yaml.version>
86         <fasterxml.jackson.avro.version>2.15.2</fasterxml.jackson.avro.version>
87         <revision>1.0.0_${spark.version}</revision>
88         <specific.sources>3_5_0</specific.sources>
89         <hadoop.version>3.4.0</hadoop.version>
90         <iceberg.version>1.5.2</iceberg.version>
91       </properties>
92       <dependencies>
93         <dependency>
94           <groupId>org.apache.spark</groupId>
95           <artifactId>spark-protobuf_${scala.version.short}</artifactId>
96           <version>${spark.version}</version>
97         </dependency>
98       </dependencies>
99     </profile>
100  </profiles>
Никита, 27.02.2024, 14:02 • BUILD_SPARK_35
```

```
mvn clean package -P spark-3.3
```

```
mvn clean package -P spark-3.5
```

# Структура проекта – разные версии Spark



```
ProtoParser.scala x
1 package ru.crpt.analytics.dataLake.base
2
3 import org.apache.spark.sql.protobuf.functions.from_protobuf
4
5 object ProtoParser { new *
6
7 } You, Moments ago • Uncommitted changes
8
```

Cannot resolve symbol protobuf

Create case class 'protobuf' ↵ ↵ ↵ More actions... ↵ ↵





# Структура проекта – разные версии Spark



```
└─ hooks
  └─ spark_33 [datalake-spark-hook-3.3]
    └─ src
      └─ main
        ├── resources
        └─ scala
          └─ ru.crpt.analytics.spark_hook
            └─ VersionSpecificTransform
      └─ test
        └─ pom.xml
  └─ spark_35 [datalake-spark-hook-3.5]
    └─ src
      └─ main
        ├── resources
        └─ scala
          └─ ru.crpt.analytics.spark_hook
            └─ VersionSpecificTransform
      └─ test
        └─ pom.xml
```




```
└─ gismt [datalake-gismt]
  └─ src
    └─ main
      ├── resources
      ├── scala
      └─ spark_specific_3_3_0
        └─ ru.crpt.analytics.datalake.gismt
          └─ protoparser
            └─ VersionSpecificTransform
      └─ spark_specific_3_5_0
        └─ ru
          └─ crpt
            └─ analytics
              └─ datalake
                └─ gismt
                  └─ protoparser
                    └─ VersionSpecificTransform
    └─ test
      └─ pom.xml
```

```
<plugin>
  <groupId>org.codehaus.mojo</groupId>
  <artifactId>build-helper-maven-plugin</artifactId>
  <version>3.5.0</version>
  <executions>
    <execution>
      <id>add-source</id>
      <phase>generate-sources</phase>
      <goals>
        <goal>add-source</goal>
      </goals>
      <configuration>
        <sources>
          <source>src/main/spark_specific_${specific.sources}</source>
        </sources>
      </configuration>
    </execution>
    <execution>
      <id>add-test-source</id>
      <phase>generate-test-sources</phase>
      <goals>
        <goal>add-test-source</goal>
      </goals>
      <configuration>
        <sources>
          <source>src/test/spark_specific_${specific.sources}</source>
        </sources>
      </configuration>
    </execution>
  </executions>
</plugin>
```




# Процессы в команде


# Процессы – code review


S sa-review-and-release   

Четверг


 **Григорий**

HotFix  
MR:  
[https://git.██████████-/merge\\_requests/7987](https://git.██████████-/merge_requests/7987)  
Описание: МДЛП. ██████████h\_type.  
Исправлена работа ДАГа 10:42

 1

 **Александр** ██████████

HotFix  
MR:  
[https://git.██████████-/merge\\_requests/7988](https://git.██████████-/merge_requests/7988)  
Задача: <https://jira.██████████-17731>  
Описание: Объединение дубликатов в ██████████ 13:25

 1

# Процессы – release

Passed Сергей [redacted] created pipeline for commit 8d69c945 finished 2 days ago

Related merge request !7837 to merge feature [redacted]

latest merge request 55 Jobs 54 minutes 34 seconds, queued for 9 seconds

Pipeline Needs Jobs 55 **Tests 2682**

## Summary

2682 tests

0 failures

0 errors

# Процессы – release

The screenshot shows a Telegram chat interface for a channel named "sa-review-and-release". At the top, there are icons for video call, information, messages, and a group of 38 members. A green checkmark icon with the number "1" is visible. The main message is from a user named "Александр" (redacted) and contains the following text:

**@room** Планируется релиз ветки master версии  
2024.08.15 в 13:50

gismt:  
[redacted] 17731: Объединение дубликатов [redacted]  
[redacted] silver

The message is timestamped "13:42" and has a rocket icon with the number "1" below it. At the bottom right, there are icons for "+3" and three profile pictures. The date "Вчера" (Yesterday) is centered at the bottom of the chat area.



# Процессы – release

Bulk edit

New merge request



Open 23

**Merged** 141

Closed 10

All 174

Recent searches ▾

Target-Branch =

master ×

Label =

~code-changes ×



Updated date ▾



**-17731**

!7988 · created 2 days ago by

Александр

Merged



0

updated 2 days ago

code-changes

**-17731**

!7964 · created 1 week ago by

Александр

Merged



8✓ Approved

0

updated 3 days ago

code-changes

**-17641 sales , added field**

!7981 · created 3 days ago by

Михаил

Merged



8✓ Approved

0

updated 3 days ago

code-changes

# Процессы – release

🏠 [redacted] / ⚙️ datalake\_etl / Tags

Filter by tag name



Updated date ▾



New tag

Tags give the ability to mark specific points in history as being important

📅 2024.08.15

4335dc46 · Merge branch [redacted] into 'master' · 2 days ago



Create release



gismt:

[redacted]-17731: Объединение дубликатов [redacted]

📅 2024.08.14

f4599d76 · Merge branch [redacted] into 'master' · 3 days ago



Create release



gismt:

[redacted]-17731: Некорректный статус [redacted]

Объединение [redacted]

📅 2024.08.13

9f0495a0 · Merge branch [redacted] into 'master' · 3 days ago



Create release



# Pipelines

- 1. MR: test**
- 2. Tag: test & build & deploy (prod)**
- 3. Run pipeline manual: test & build & deploy (test)**

# Merge Request



latest merge request 28 Jobs 43 minutes 26 seconds, queued for 11 seconds

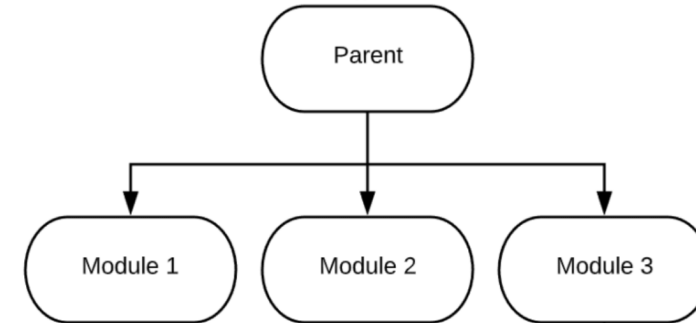
Pipeline Needs Jobs 28 Tests 1118

test

✓ mr:test 28

- ✓ mr:test: [beer]
- ✓ mr:test: [common-apps]
- ✓ mr:test: [common-test]
- ✓ mr:test: [common]
- ✓ mr:test: [data-quality]
- ✓ mr:test: [dispenser\_workers]

# Maven™



Example of parallel:matrix:

```
deploystacks:  
  stage: deploy  
  script:  
    - bin/deploy  
  parallel:  
    matrix:  
      - PROVIDER: aws  
        STACK:  
          - monitoring  
          - app1  
          - app2  
      - PROVIDER: ovh  
        STACK: [monitoring, backup, app]  
      - PROVIDER: [gcp, vultr]  
        STACK: [data, processing]  
  environment: $PROVIDER/$STACK
```

# Merge Request

latest merge request 1 Jobs 16 minutes 3 seconds, queued for 1 seconds

Pipeline Needs Jobs 1 Tests 49

Group jobs by Stage Job dependencies

test

✓ mr:test:milk ↻



Example of `rules:changes:`

```
docker build:  
script: docker build -t my-image:$CI_COMMIT_REF_SLUG .  
rules:  
  - if: $CI_PIPELINE_SOURCE == "merge_request_event"  
    changes:  
      - Dockerfile  
    when: manual  
    allow_failure: true
```

```
docker build alternative:  
variables:  
  DOCKERFILES_DIR: 'path/to/dockerfiles'  
script: docker build -t my-image:$CI_COMMIT_REF_SLUG .  
rules:  
  - if: $CI_PIPELINE_SOURCE == "merge_request_event"  
    changes:  
      - $DOCKERFILES_DIR/**/*
```

# Merge Request

Broken DAG: [/home/airflow/gcs/dags/airflow\_utils.py] 'DEPLOYMENT\_SETUP'

Broken DAG: [/home/airflow/gcs/dags/examples/dbt\_example.py] 'DEPLOYMENT\_SETUP'

Broken DAG: [/home/airflow/gcs/dags/examples/add\_gcp\_connections.py] 'DEPLOYMENT\_SETUP'

Broken DAG: [/home/airflow/gcs/dags/examples/kubernetes\_sample.py] 'DEPLOYMENT\_SETUP'

### DAGs

Search:

	ⓘ	DAG	Schedule	Owner	Recent Tasks ⓘ	Last Run ⓘ
	On	airflow_monitoring	None	airflow	1	2020-12-17 11:52 ⓘ
	On	bigquery_connection_check	@once	airflow	1	2020-01-16 00:00 ⓘ

« < 1 > »

test:airflow

test:airflow

← variables

← connections



## Airflow prod cluster

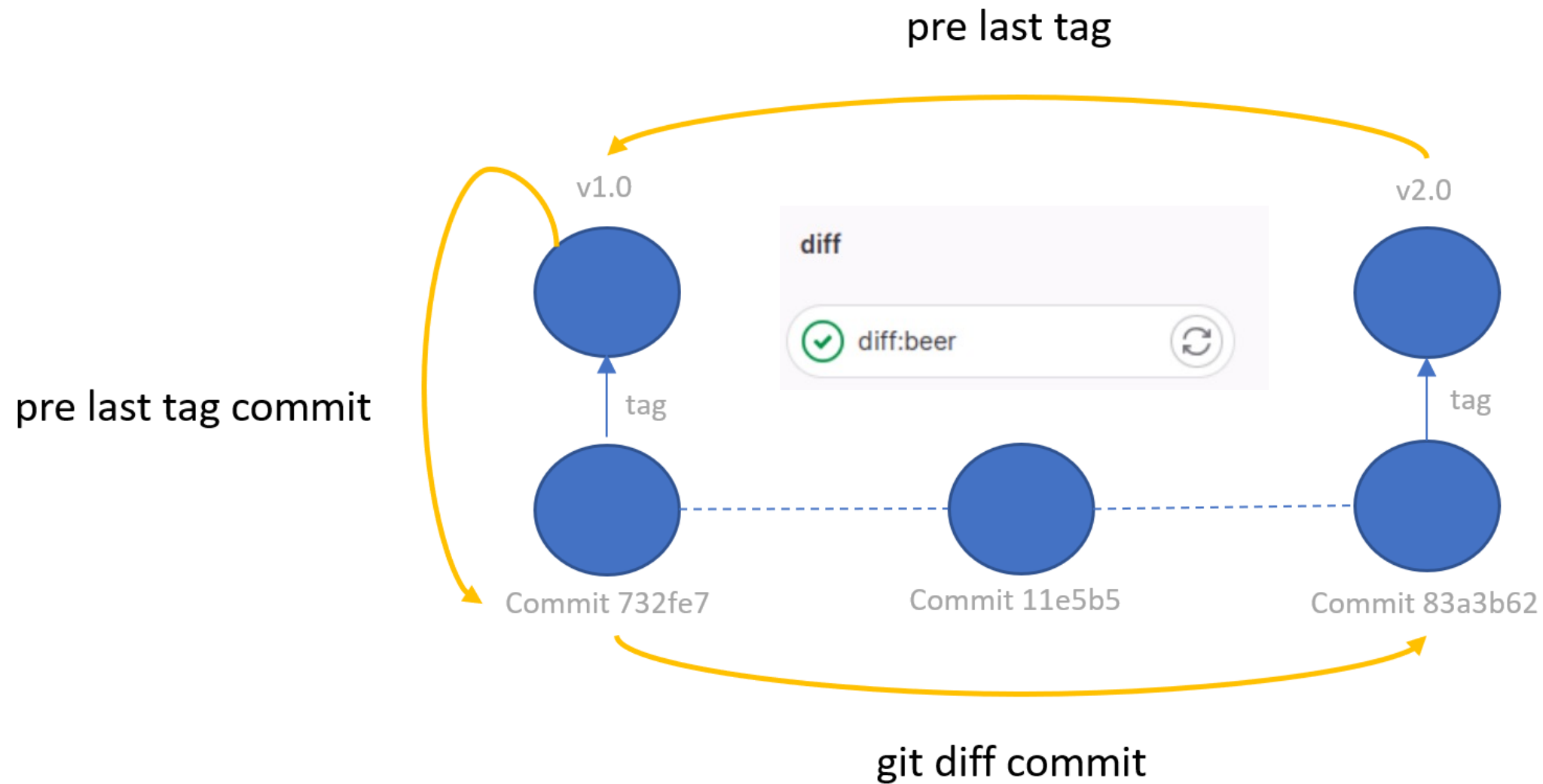
← variables

← connections

## Smoke test Dag's



# Tag



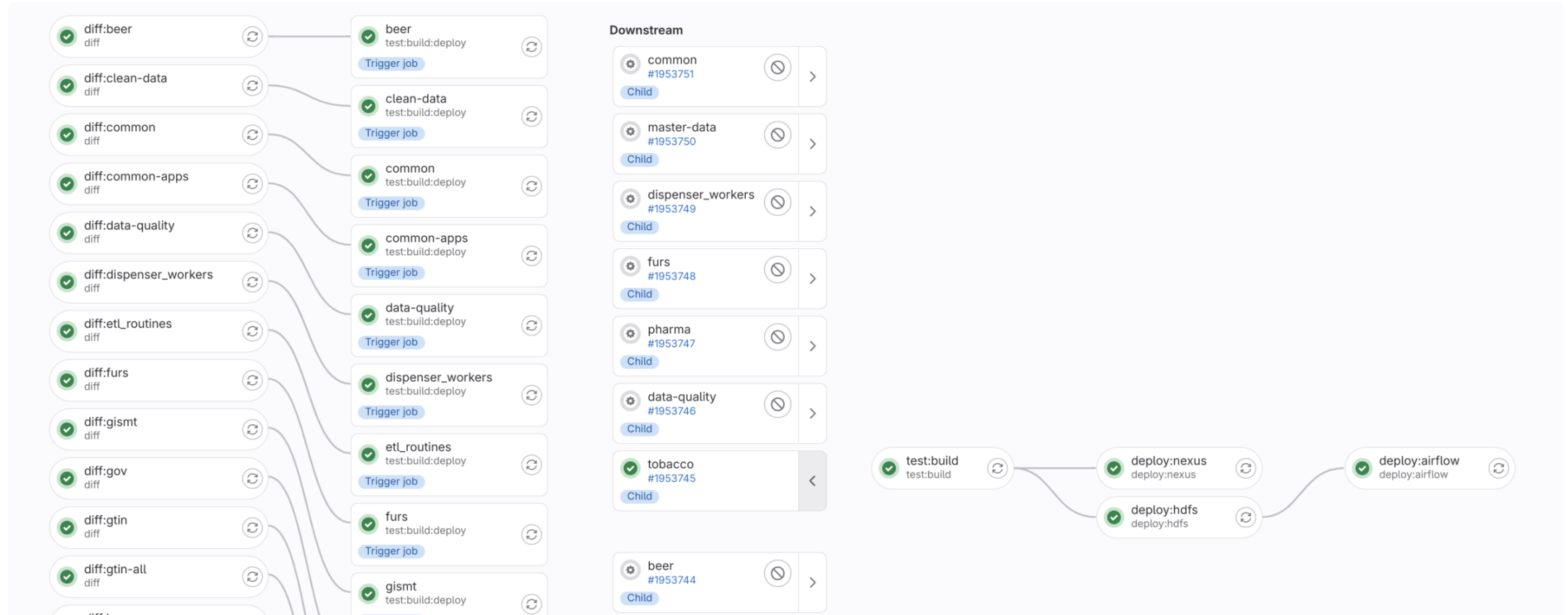
проверка на наличие изменений в модуле

# Tag

Pipeline Needs Jobs 56 Tests 0

Group jobs by Stage Job dependencies Show dependencies

Tip: Hover over a job to see the jobs it depends on to run.



# Manual

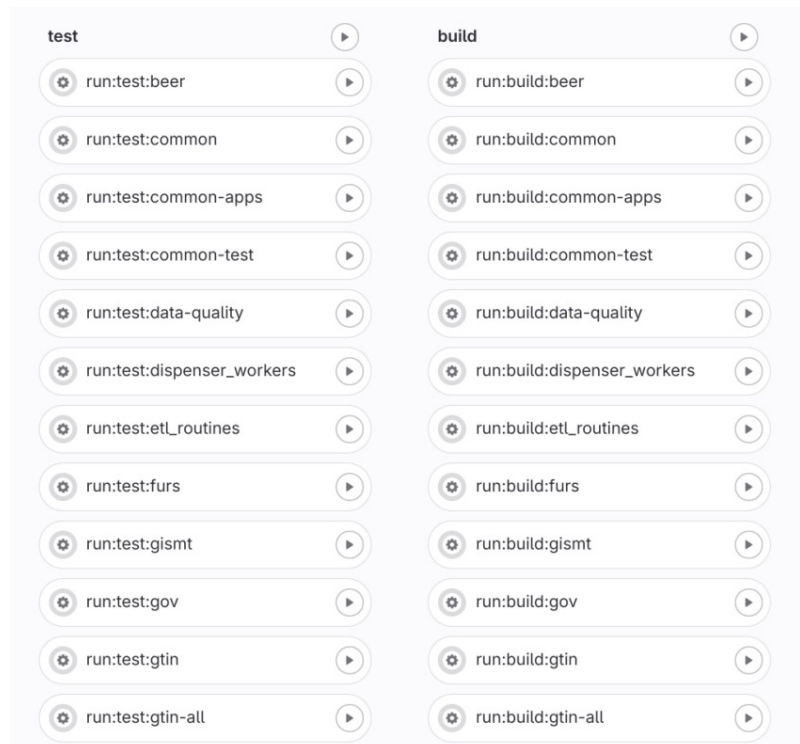
anst > datalake\_etl > Pipelines > #1879122

Pipeline Needs Jobs 115 Tests 0

Group jobs by

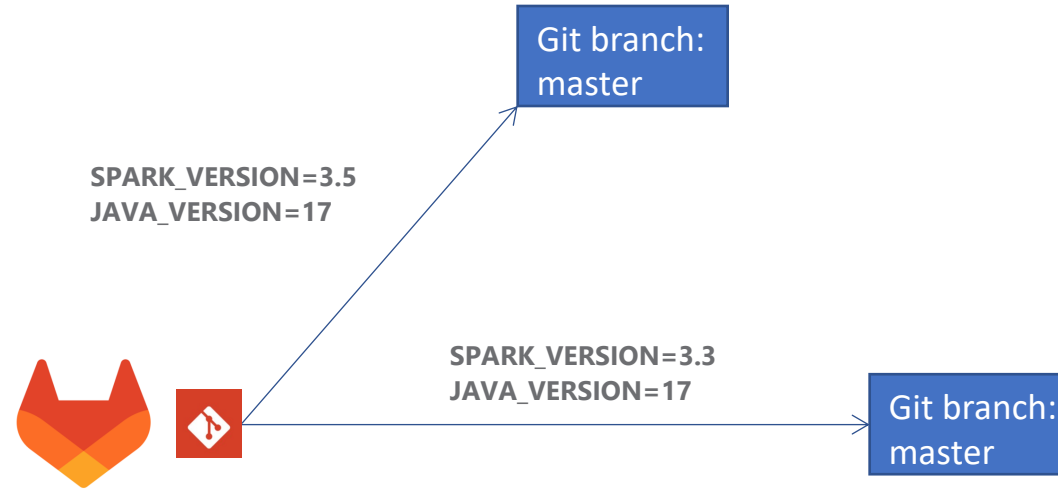
image	test	build	deploy:hdfs-test	deploy:hdfs-offline
<input type="checkbox"/> image:java-tooling	<input type="checkbox"/> run:test:beer	<input type="checkbox"/> run:build:beer	<input type="checkbox"/> run:deploy:hdfs-test:beer	<input type="checkbox"/> run:deploy:hdfs-offline:beer
<input type="checkbox"/> image:liberica-openjdk-debian:8	<input type="checkbox"/> run:test:common	<input type="checkbox"/> run:build:common	<input type="checkbox"/> run:deploy:hdfs-test:common	<input type="checkbox"/> run:deploy:hdfs-offline:common
<input type="checkbox"/> image:liberica-openjdk-debian:17.0.5	<input type="checkbox"/> run:test:common-apps	<input type="checkbox"/> run:build:common-apps	<input type="checkbox"/> run:deploy:hdfs-test:common-apps	<input type="checkbox"/> run:deploy:hdfs-offline:common-apps
	<input type="checkbox"/> run:test:common-test	<input type="checkbox"/> run:build:common-test	<input type="checkbox"/> run:deploy:hdfs-test:common-test	<input type="checkbox"/> run:deploy:hdfs-offline:common-test
	<input type="checkbox"/> run:test:data-quality	<input type="checkbox"/> run:build:data-quality	<input type="checkbox"/> run:deploy:hdfs-test:data-quality	<input type="checkbox"/> run:deploy:hdfs-offline:data-quality
	<input type="checkbox"/> run:test:dispenser_workers	<input type="checkbox"/> run:build:dispenser_workers	<input type="checkbox"/> run:deploy:hdfs-test:dispenser_workers	<input type="checkbox"/> run:deploy:hdfs-offline:dispenser_workers
	<input type="checkbox"/> run:test:etL_routines	<input type="checkbox"/> run:build:etL_routines	<input type="checkbox"/> run:deploy:hdfs-test:etL_routines	<input type="checkbox"/> run:deploy:hdfs-offline:etL_routines
	<input type="checkbox"/> run:test:furs	<input type="checkbox"/> run:build:furs	<input type="checkbox"/> run:deploy:hdfs-test:furs	<input type="checkbox"/> run:deploy:hdfs-offline:furs
	<input type="checkbox"/> run:test:gismt	<input type="checkbox"/> run:build:gismt	<input type="checkbox"/> run:deploy:hdfs-test:gismt	<input type="checkbox"/> run:deploy:hdfs-offline:gismt

## Maven cache & Gitlab CI/CD cache



```
cache:  
  key: master33cache  
  policy: pull  
  paths:  
    - .m2/repository/
```

# Spark



## Downstream

	common #2653437 Child		>
	softdrinks #2653436 Child		<

test:build

test:build

test:build\_3.5

deploy:nexus

deploy:nexus

deploy:nexus\_3.5

deploy:hdfs

deploy:hdfs

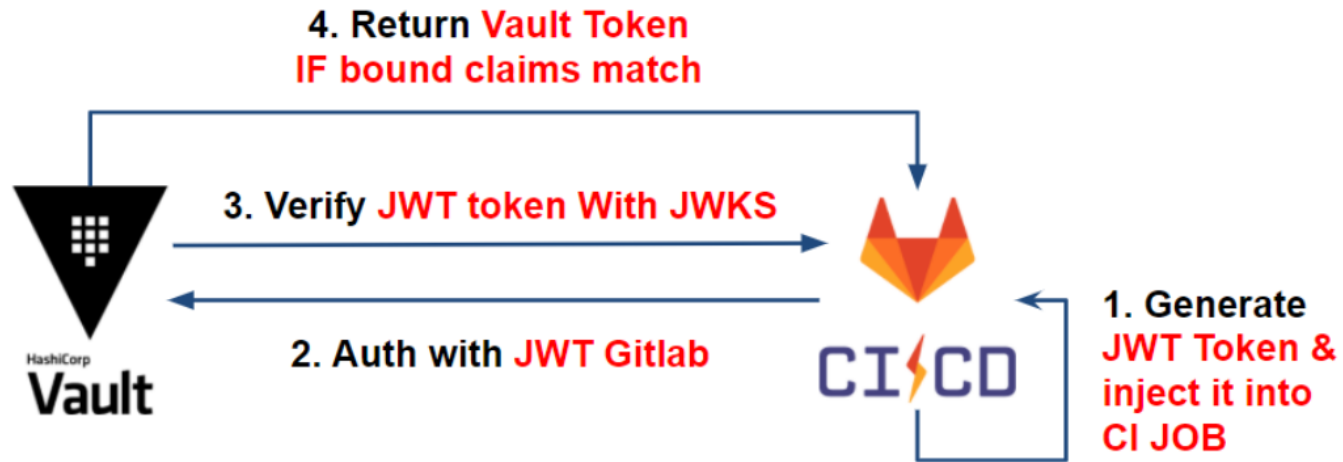
deploy:hdfs\_3.5

deploy:airflow

deploy:airflow

deploy:airflow\_3.5

# Security

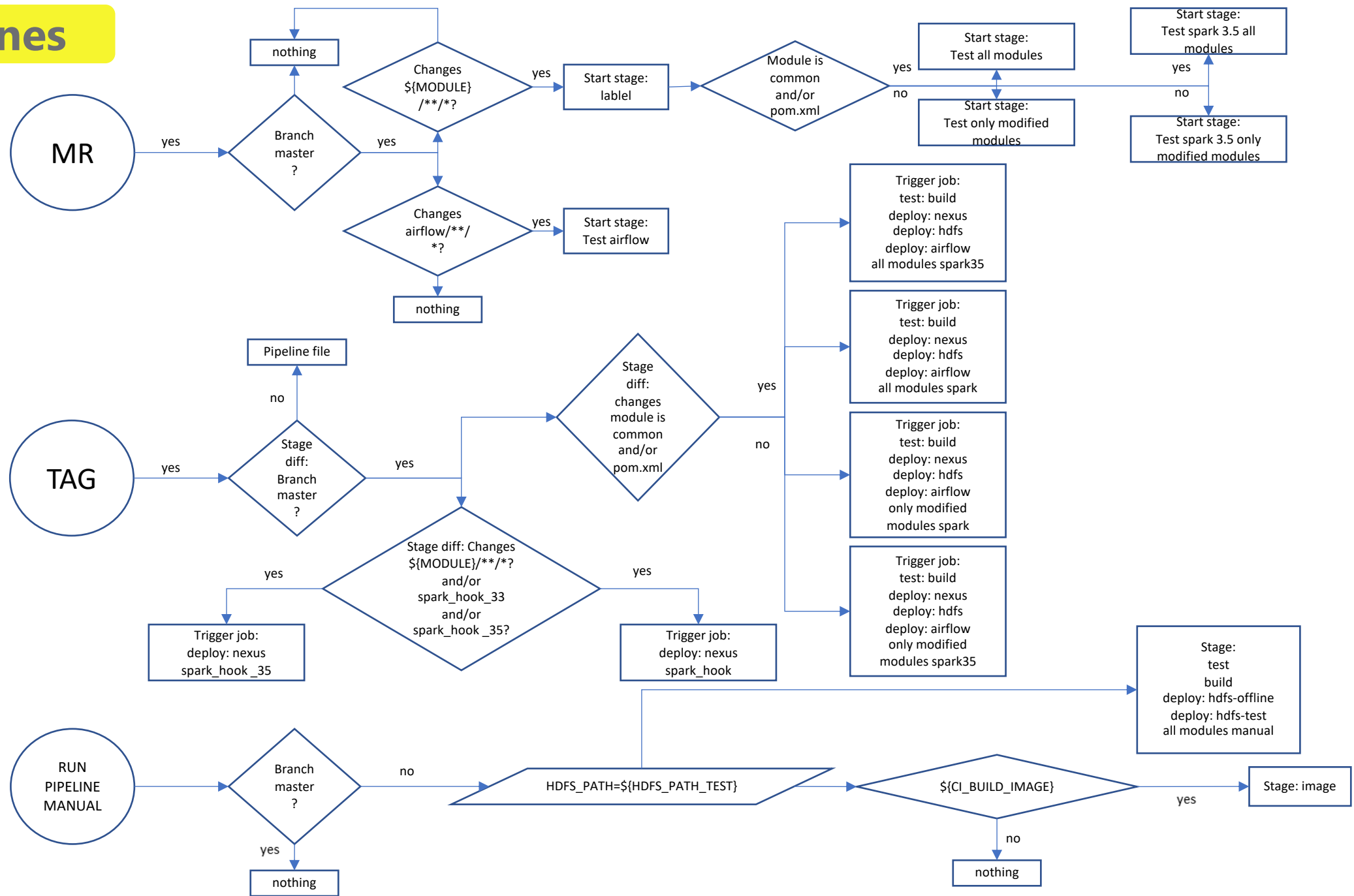


Members 3

Filter members

Account	Source	Access granted	Max role	Expiration
Administrator @root <span>It's you</span>	<a href="#">demo</a>	5 minutes ago	Owner	Expiration date
User 0 @user0	<b>Inherited</b> <a href="#">demo</a>	1 minute ago by Administrator	Reporter	Expiration date
User 1 @user1	Direct member	54 seconds ago by Administrator	Guest	Expiration date

# Pipelines



# Оптимизации I/O



## Оптимизации I/O

```
val res = result.columns.foldLeft(  
  .select(columns.map(col): _*)
```



```
res.toJSON.show(false)    You, 25.0
```

```
val expDf = spark.read.json(getCla  
val expected = expDf.columns.foldL  
  .select(columns.map(col): _*)
```

```
checkAnswer(res, expected)
```

# Оптимизации I/O

```
worker:a doing lots of uninteresting work +0ms
worker:b doing some work +0ms
http listening +23ms
worker:a doing lots of uninteresting work +424ms
worker:a doing lots of uninteresting work +307ms
worker:b doing some work +814ms
worker:b doing some work +58ms
worker:a doing lots of uninteresting work +312ms
worker:a doing lots of uninteresting work +647ms
worker:a doing lots of uninteresting work +469ms
worker:b doing some work +1s
worker:a doing lots of uninteresting work +797ms
worker:a doing lots of uninteresting work +153ms
worker:b doing some work +1s
worker:a doing lots of uninteresting work +491ms
worker:a doing lots of uninteresting work +323ms
worker:b doing some work +602ms
```

**Что дальше?**

# Развитие CI/CD

```
<project xmlns="http://maven.apache.org/POM/4.0.0"
  <build>
    <plugins>
      <!-- enable scalatest -->
      <plugin>
        <groupId>org.scalatest</groupId>
        <artifactId>scalatest-maven-plugin</artifactId>
        <version>${scalatest.maven.plugin.version}</version>
        <configuration>
          <reportsDirectory>${project.build.directory}/surefire-reports</reportsDirectory>
          <junitxml>./</junitxml>
          <filereports>TestSuite.txt</filereports>
          <!--Parallel doesn't work - Only one SparkContext should be running in this JVM (see SPARK-2243)-->
          <parallel>>false</parallel>
```

# Post scriptum

## Post Scriptum

**Мы ищем Data Engineer-ов на наш стек. Приходите ко мне в Telegram @nblagodarnuu**

# Post Scriptum

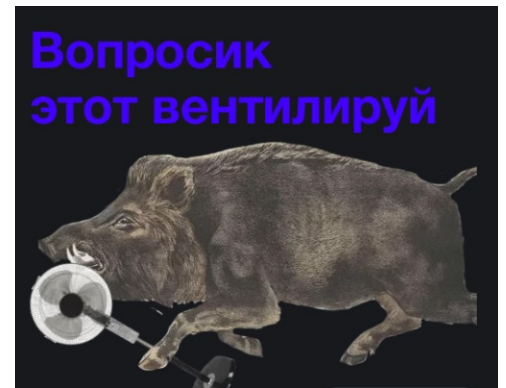
DE-meetup в СПб 11/09



- <https://www.baeldung.com/maven-dependencymanagement-vs-dependencies-tags>
- [https://maven.apache.org/guides/introduction/introduction-to-dependency-mechanism.html#Dependency\\_Management](https://maven.apache.org/guides/introduction/introduction-to-dependency-mechanism.html#Dependency_Management)
- <https://maven.apache.org/maven-ci-friendly.html>



# Вопросы?



**СПАСИБО  
ЗА ВНИМАНИЕ!**

---

Никита Благодарный  
telegram [@nblagodarnyy](#)

Александра Чекмарева (Китченко)  
telegram [@sasha\\_kitchen](#)



**ЧЕСТНЫЙ  
ЗНАК**



центр развития  
перспективных технологий



**SmartData**

2024

# Post Scriptum

DE-meetup в СПб 11/09

