# Bug bites Elephant?

Test-driven Quality Assurance

in Big Data Application Development

*Dr. Dominik Benz, Inovex GmbH*
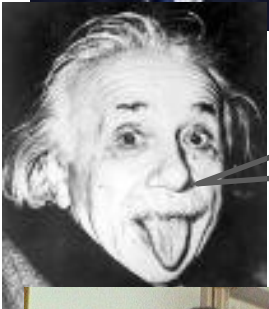
*2013/06/03, Berlin Buzzwords*

the **FitNesse** approach

our Big Data QA **problem**

test **data** definition / selection

result **inspection**

job & **workflow** control

BI **reporting**, web **analytics**, …

**DWH**

~ **1 billion** log events / day,
~ **1 TB** (thrift) logfiles

**Hadoop Cluster**

chains of MR jobs, running on **20 nodes** / 8 cores / 96 GB RAM (CDH)
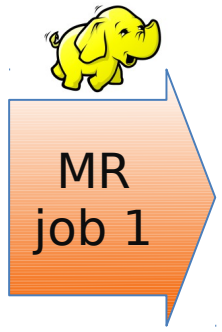
**?** create (sample) input data

**?** inspect (binary) formats

**?** control workflows

Log Files (thrift)

MR job 1

Inter-mediate result (avro)

MR job 2

...

DWH (RDBMS)

| method | tests what? | issues for our usecase |
|--------|-------------|------------------------|
| **JUnit** | isolated functions | no integration, Java syntax |
| **MRUnit** | 1 mapper + 1 reducer | „little" integration, Java syntax |
| **iTest** | hadoop jobs/workflows | Java / Groovy syntax |
| **Scripts/ CLI** | (manual) | „script chaos", |

⊟ **FitNesse** as suitable addition / solution!

*Big Data QA is different!*

the **FitNesse** approach

test **data** definition / selection

job & **workflow** control

result **inspection**

inovex



*„fully integrated standalone wiki and acceptance testing framework"*

„executable" **Wiki**-Pages (returning test results)

(almost) **natural language** test specification

connection to SUT via (Java-)"**Fixtures**"

**Brows er**

```
script |
check  |
num results |
3 |
```

*FitNesse Server*

**Fixtur es**

```
public int
   numResults
{ ... }
```

System under Test

→ „calling java methods from wiki", compare return values

→ Integrates with REST,
Jenkins

# FitNesse — An Exemplary Test



10

```
!path /home/inovex/lib/*.jar


| script | Hadoop |

| upload | viewLog.csv | to hdfs | /testdata/ |

| hadoop job from jar  | viewLog.jar | [...] |

| show    | job output |

| check   | number of output files | 3 |
```

```java
public class Hadoop {

  public boolean uploadToHdfs(String localFile,

                              String remoteFile) {...}


  public boolean hadoopJobFromJar(String jar,

              String input, String output) {...}


  public String jobOutput() {...}


  public String numberOfOutputFiles() {...}
}
```

*Big Data QA is different!*

*Fitnesse Wiki test execution!*

test **data** definition / selection

job & **workflow** control

result **inspection**

inovex

| Table:Log File | | | | |
| --- | --- | --- | --- | --- |
| /home/inovex/custom_logs/input/viewLog.csv | | | | |
| date | user | product | browser | os |
| 2013-03-12 | john | 1 | ff | win |
| 2013-03-12 | john | 2 | ff | win |
| 2013-03-13 | john | 2 | ff | win |
| 2013-03-14 | lisa | 1 | ie | win |
| 2013-03-14 | peter | 1 | ff | lin |
| 2013-03-15 | lisa | 2 | ie | lin |
| 2013-03-15 | peter | 2 | ff | mac |
| 2013-03-16 | lisa | 1 | ff | mac |

▸ Big Data: **Efficient** data transfer among **heterogeneous sources**

▸ Def
lang

| Table:Thrift Log File | | | | |
| --- | --- | --- | --- | --- |
| /home/inovex/viewLog.thrift | de.inovex.thrift.ViewLog | | | |
| date | user | product | browser | os |
| 2013-03-12 | john | 1 | ff | win |
| 2013-03-12 | john | 2 | ff | win |
| 2013-03-13 | john | 2 | ff | win |
| 2013-03-14 | lisa | 1 | ie | win |

- Dev/Test Hadoop Cluster: **Identical Hardware** like Prod, but fewer nodes

- (random/biased) **sampling** e.g. on daily basis

- **Feedback loop:**

  - identify „**special cases**" from real data

  - include them in (manual) data definition

  - Gradually **increase test coverage** / artefact quality

inovex



*FitNesse Wiki test execution!*

*Big Data QA is different!*

*Define CSV / thrift / real-world test data!*

result **inspection**

job & **workflow** control

| script | Shell | |
|--------|-------|--|
| show | exec | hadoop fs -put /home/inovex/custom_logs/input /viewLog.csv /user/inovex/input/ |
| show | exec | hadoop jar /home/inovex/showcase/ViewLogCounter.jar /user/inovex/input/*.csv /user/inovex/output |

‣ Execute arbitrary (shell) commands

‣ Mainly a **wrapper** around
**apache.commons.exec.CommandLine**

inovex

| script | Hadoop | | |
|--------|--------|--------|----------|
| upload | viewLog.csv | to hdfs | /testdata |
| hadoop job from jar; | ViewLogCounter.jar | /testdata /*.csv | output |
| show | job output | | |
| check | number of output files | 3 | |

► **Hide complexity** from test authors

► „define" appropriate **test language** via (Java) method names

► **re-use** other fixtures (Shell, …) internally

- FitNesse allows to group tests into **suites**

- Can be used to simulate MR **processing chains**

- **SetupSuite** / TearDownSuite for creating / destroying test conditions

- Tests can still be executed **individually**

Big Data QA is different!

FitNesse Wiki test execution!

Define CSV / thrift / real-world data!

Use suites & fixtures for jobs/workflows!

result **inspection**

inovex

| Query:Db Select | SELECT cust, prod, count FROM viewcount WHERE cust='john' | |
|---|---|---|
| cust | prod | count |
| john | 2 | 2 |
| john | 1 | 1 |

‣ Validate **RDBMS contents** (via JDBC)

‣ E.g. for checking the **final** result

‣ Or use **Hive** + Hive-Server to query raw data

inovex

| script | Pig Console | MAPREDUCE | fixtures.jar | avro.output.codec=snappy |
|---|---|---|---|---|
| load avro file | /data/viewlog.avro | using alias | viewlog | |
| execute | flights = foreach viewlog generate key.campaign as flight:int | | | |
| execute | filtered = filter flights by flight = 5283 | | | |
| check | number of records from alias | 35 | | |

▸ Execute arbitrary **pig commands** from Wiki page

▸ Inspect e.g. **binary intermediate results** (avro, …)

```java
public class PigConsole extends PigServer {


    public void loadAvroFileUsingAlias(String
                          filename, String alias) {
        this.registerQuery(
            alias + "= LOAD" + filename + "USING" +
            AVRO_STORAGE_LOADER + ";");
    }

}
```

inovex



Fitnesse Master

TestEnvironments

ProjA   ProjB

TestConfigurations

ProjA

de   qs   live

ProjB

de   qs   live

*Import / edit tests remotely*

*Import / edit config remotely*

de   qs   live

ProjA

Dev ProjA Slave

QS ProjA Slave

Live ProjA Slave

inovex

FitNesse Wiki test execution!

Big Data QA is different!

Define CSV / thrift / real-world data!

inovex

Inspect results via Pig/Hive

Use suites & fixtures for jobs/workflows!

# Want more?   Inovex trains you!

- **Android Developer Training** (3 days, Karlsruhe/München)

- **Certified Scrum Developer Training** (5 days, Köln)

- **Hadoop Developer Training** (3 days, Karlsruhe/Köln)

- **Liferay Portal-Developer Training** (4 days, Karlsruhe)

- **Liferay Portal-Admin Training** (3 days, Karlsruhe)

- **Pentaho Data Integration Training** (4 days, München/Köln)

information and registration at

**www.inovex.de/offene-trainings**

Stefan

Kathri
n

Bernha
rd

Jörg

Andre
w

Christi
an

Christia
n

inovex

Fitnesse Master

TestEnvironments

ProjA    ProjB

TestConfigurations

ProjA

de    qs    live

ProjB

de    qs    live

***Import / edit tests remotely***

*Import / edit config remotely*

de

qs

live

Dev ProjA Slave

QS ProjA Slave

Live ProjA Slave
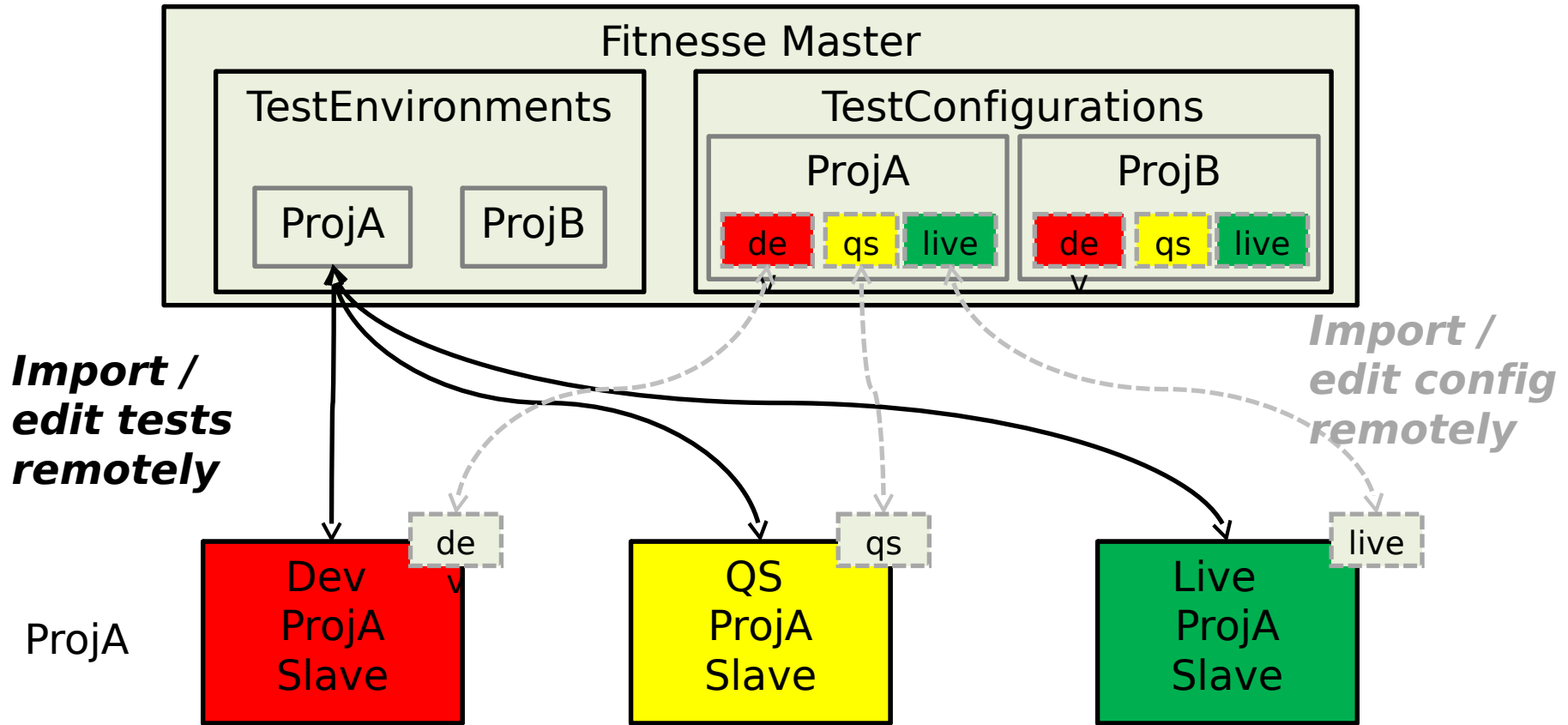
ProjA

- ▸ Download & install FitNesse server

- ▸ Create csv log file

- ▸ Run hadoop job which counts viewed items

- ▸ Inspect Results with Hive

inovex

```
!path /home/inovex/lib/*.jar


| Table:Log File                |

| /home/inovex/viewLog.csv | |

| date          | user   | product | browser | os    |

| 2013-03-12 | john   | 1          | ff          | win |


| script | Hadoop |

| upload | viewLog.csv | to hdfs | /testdata/ |

| hadoop job from jar  | viewLog.jar | [...] |

| show    | job output |

| check   | number of output files | 3 |
```

localhost:8080/FrontPage.HadoopLogFileProcessingTest?test ☆ ▾ ⟳ | Google ▾ | 🔍

FrontPage

# HadoopLogFileProcessingTest

| Tests Executed OK | Test | Edit | Add | Tools |

☑ **Assertions:** 15 right, 0 wrong, 0 ignored, 0 exceptions (0.077 seconds)

▶ *Precompiled Libraries* | Expand All | Collapse All

▶ *setup* | Expand All | Collapse All

| Table:Log File | | | | |
|---|---|---|---|---|
| /home/inovex/viewLog.csv | | | | |
| date | user | product | browser | os |
| 2013-03-12 | john | 1 | ff | win |

| script | Hadoop | | | |
|---|---|---|---|---|
| upload | viewLog.csv | to hdfs | /testdata | |
| hadoop job from jar; | ViewLogCounter.jar | /testdata/*.csv | output | |
| show | job output | | | Starting MR job [...] |
| check | number of output files | 3 | | |