

Тестирование в
мире
BIG DATA

by Konstantin Pletenev

**GROW
CONFI
DENTLY**

SQA®
DAYS#23

BIG DATA IS NOT ABOUT THE DATA

THE REVOLUTION IS NOT THAT THERE'S MORE DATA AVAILABLE

THE REVOLUTION IS THAT WE KNOW WHAT TO DO WITH THIS DATA NOW

**—GARY KING
PROFESSOR OF HARVARD UNIVERSITY**





**65 МИЛЛИОНОВ ПОСТОВ В СОЦИАЛЬНЫХ СЕТЯХ
КАЖДЫЙ ДЕНЬ**



**1 ЧАС, 16 МИНУТ КАЖДЫЙ ДЕНЬ ТРАТИТСЯ НА
ПРОСМОТР ВИДЕО С МОБИЛЬНОГО УСТРОЙСТВА**



**100 МИЛЛИОНОВ ПОДПИСЧИКОВ НА
СТРИМИНГОВЫЙ СЕРВИС NETFLIX**

ЧТО ТАКОЕ BIG DATA

ЧТО НЕ BIG DATA

- **НЕ** о традиционных данных, таких как документы и базы данных
- **НЕ** просто еще одно обозначение «Большого количества данных»
- **НЕ** только о данных

ЧТО **ЖЕ** BIG DATA

- **BIG DATA** это большой объем данных обрабатываемый на ежедневной основе
- **BIG DATA** это то, что мы можем делать с этими данными
- **BIG DATA** нужна для принятия лучших и правильных стратегических решений

3 BIG DATA ПРОЕКТА В АССЕНТУРЕ ЗА 2 ГОДА

1. ОБРАБОТКА ДАННЫХ О ПОЛЬЗОВАТЕЛЬСКОЙ АКТИВНОСТИ И ИХ СЕРИАЛИЗАЦИЯ
2. СИСТЕМА ПРЕДОСТАВЛЕНИЯ ВИДЕО КОНТЕНТА И ПОДСЧЕТ КРІ ДЛЯ БИЗНЕС-АНАЛИТИКИ
3. ОБРАБОТКА И ПОДГОТОВКА ДАННЫХ О ПОЛЬЗОВАТЕЛЬСКОЙ АКТИВНОСТИ ДЛЯ DATA SCIENTISTS



BIG DATA V³

VOLUME

VELOCITY

VARIETY

V¹ - VOLUME FACEBOOK STATISTIC

- **100 PB** НА ОДНОМ HDFS КЛАСТЕРЕ
- **300 TB** ДАННЫХ ОБРАБАТЫВАЮТСЯ В HIVE КАЖДЫЙ **ЧАС**
- **500 TB** НОВЫХ ДАННЫХ ДОБАВЛЯЮТСЯ В ХРАНИЛИЩЕ **КАЖДЫЙ** ДЕНЬ

<https://www.facebook.com/notes/facebook-engineering/under-the-hood-hadoop-distributed-file-system-reliability-with-namenode-and-avata/10150888759153920/>

V² - VELOCITY

NETFLIX STATISTIC

- **37% ВСЕЙ** ПРОПУСКНОЙ СПОСОБНОСТИ ИНТЕРНЕТА В СЕВЕРНОЙ АМЕРИКЕ
- АНАЛИЗ ПОСТУПАЮЩИХ ДАННЫХ В **РЕАЛЬНОМ ВРЕМЕНИ**
- **100,000** СЕРВЕРНЫХ ИНСТАНЦИЙ

<https://www.linkedin.com/pulse/netflix-dominates-fixed-internet-video-bandwidth-mobile-jeff-nelson>

V³ - VARIETY

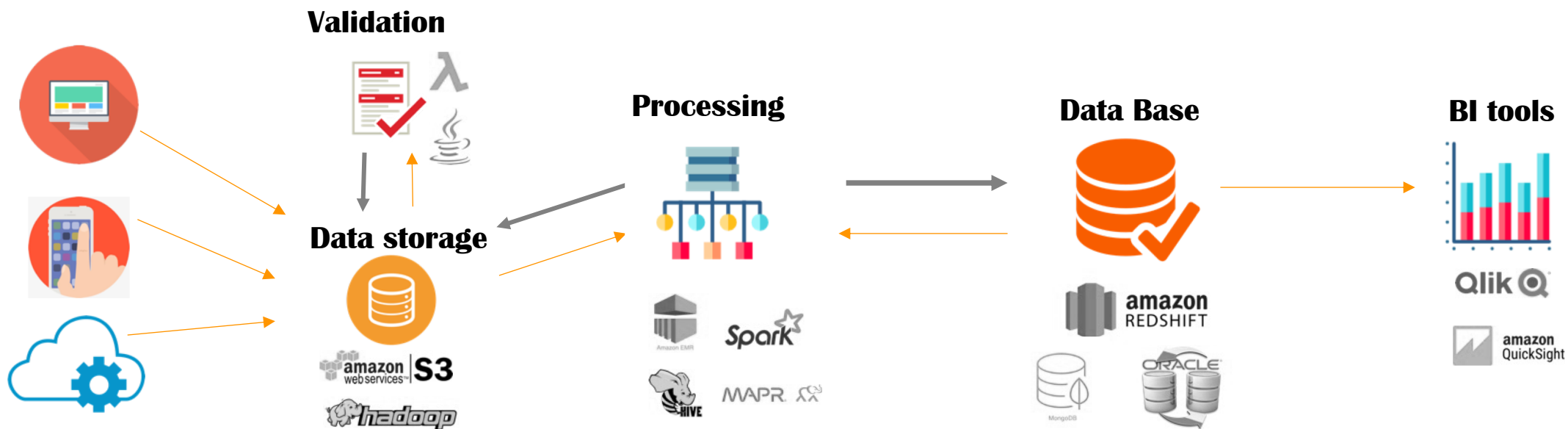
INSTAGRAM STATISTIC

- **300 МВ** ИЗОБРАЖЕНИЙ ЗАГРУЖАЕТСЯ КАЖДУЮ МИНУТУ
- **4.7 МИЛЛИОНА** ЛАЙКОВ КАЖДЫЙ ДЕНЬ
- **1 МИЛЛИОН** ПРОСМОТРОВ СТОРИС В ДЕНЬ

<https://www.forbes.com/sites/bernardmarr/2018/03/16/the-amazing-ways-instagram-uses-big-data-and-artificial-intelligence/#1524dcf25ca6>

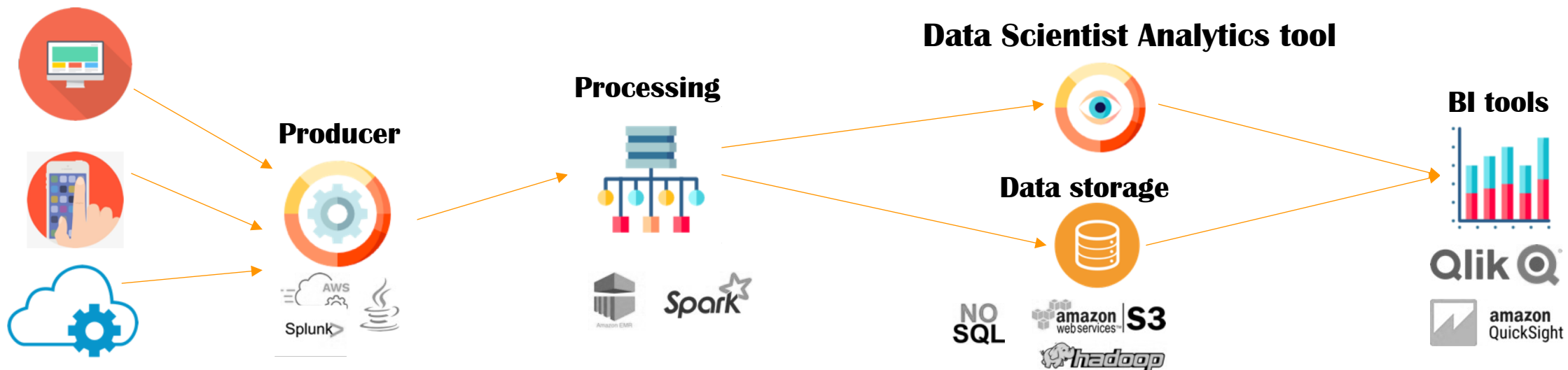
ПРИМЕР BIG DATA PIPELINE

Ежедневная обработка / Batch Pipeline

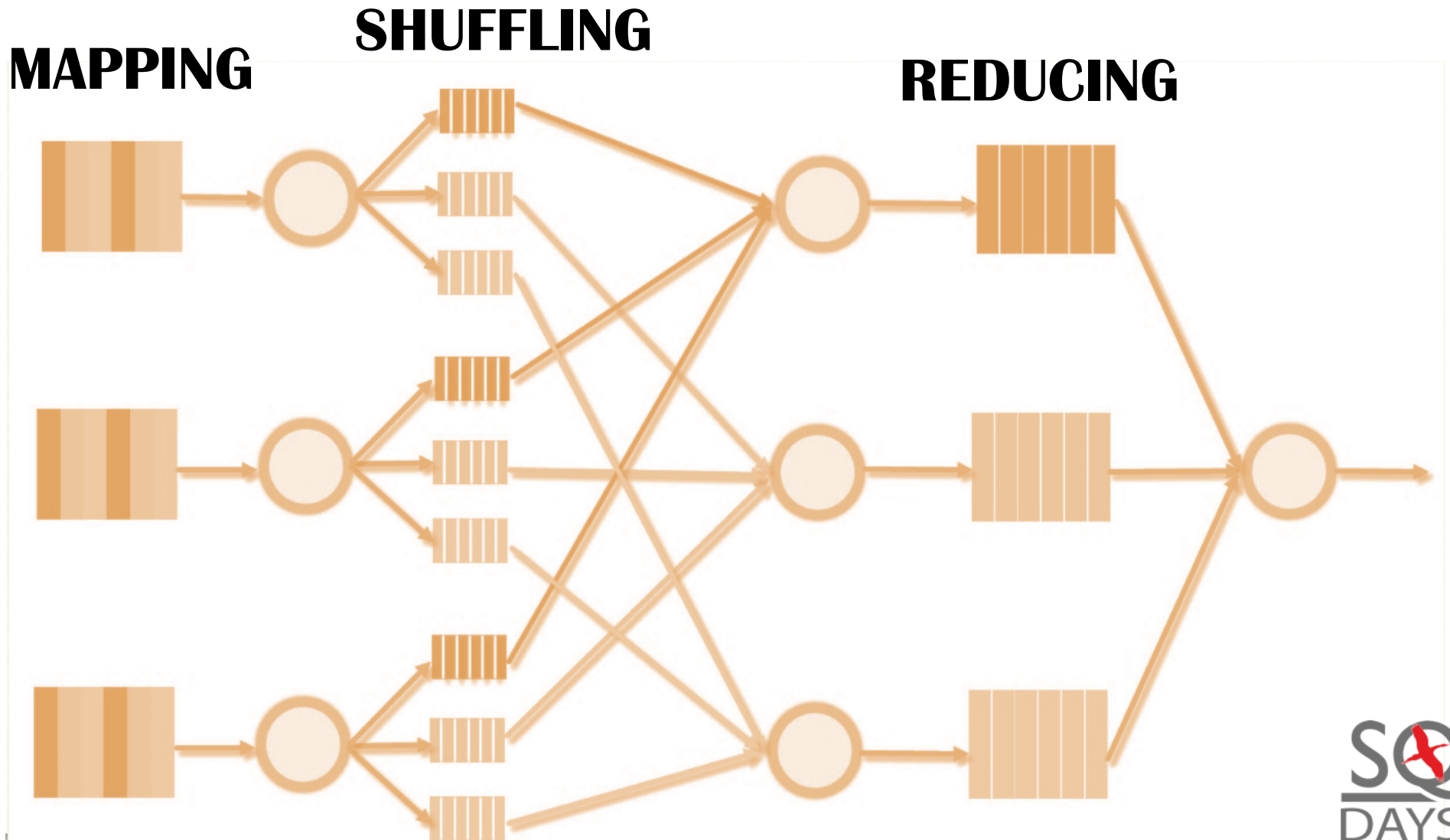


ПРИМЕР BIG DATA PIPELINE

Обработка в реальном времени / Real Time Pipeline



BIG DATA MAP REDUCE



ТЕСТИРОВАНИЕ

BIG DATA

VS

ТРАДИЦИОННОЕ

ТЕСТИРОВАНИЕ

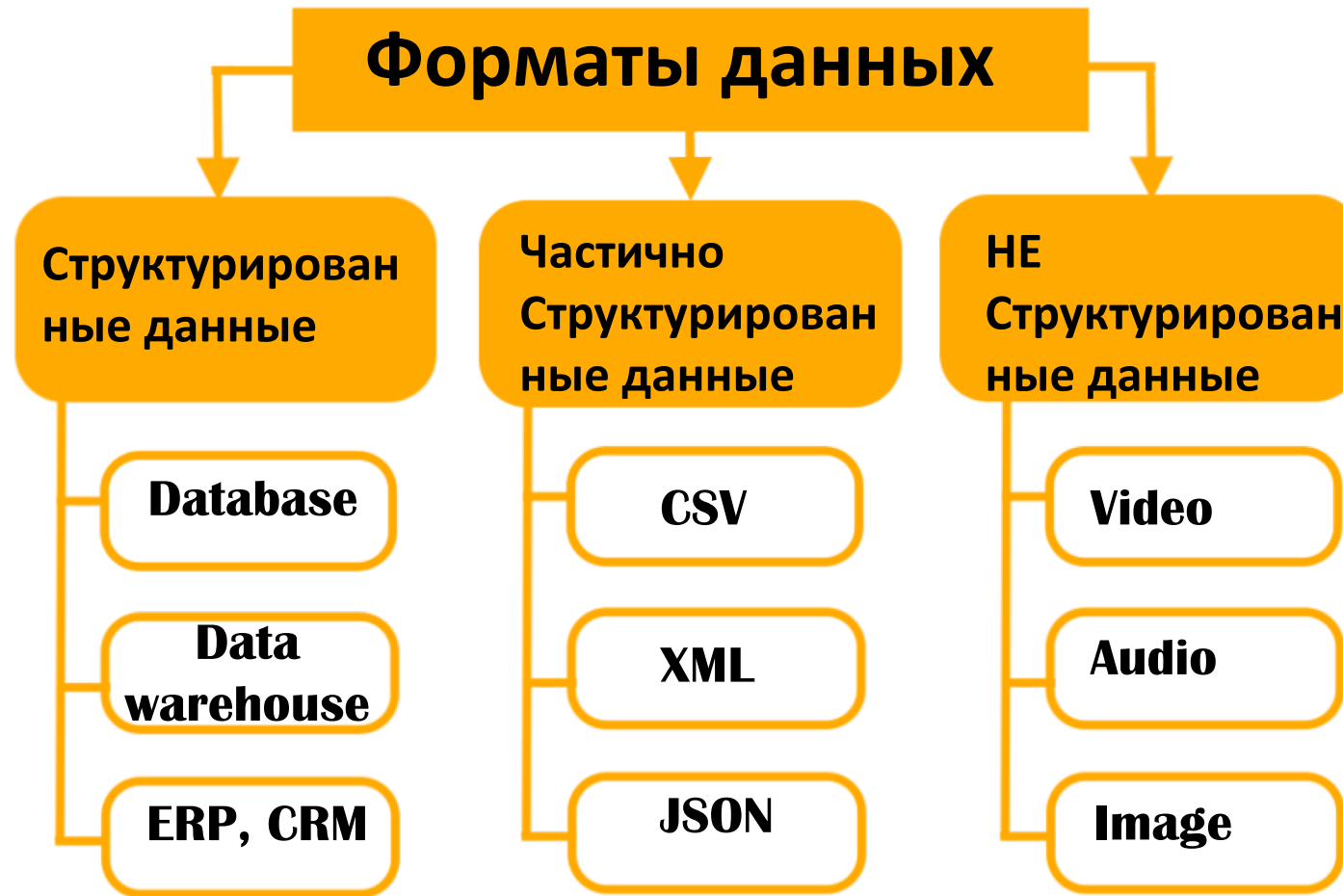
Область	Традиционное тестирование	BIG DATA тестирование
Данные	Тестировщики работают с структурированными данными	Тестировщики работают одновременно и с структурированными данными и с неструктурированными данными
	Подходы к тестированию определены и проверены временем	Подходы к тестированию сфокусированы больше на R&D
	Тестировщики могут применять "Sampling" техники для тестирования	"Sampling" техника в Big data не работает или требует очень больших затрат
	Возможно ручное тестирование	Возможно только автоматизированное тестирование

Область	Традиционное тестирование	BIG DATA тестирование
Инфраструктура	Испытуемая система может быть запущена локально	Требует отдельного окружения для запуска из-за объема данных и количества файлов (HDFS)
Инструменты валидации	Тестировщики имеют множество инструментов, начиная от Excel и заканчивая инструментами UI автоматизации	Нет специальных инструментов именно для тестирования, тестировщики используют тот же набор инструментов, что и разработчики (MapReduce, HIVE)
	Инструменты могут быть использованы без специальных знаний	Требуется знание программирования и прохождения тренингов по работе с данными

ШАГИ К УСПЕХУ

1. ПОНЯТЬ СВОИ ДАННЫЕ
2. ПОНЯТЬ СВОЙ ПРОЦЕСС
(BATCH TIME/REAL TIME)
3. АВТОМАТИЗИРОВАТЬ, АВТОМАТИЗИРОВАТЬ
И ЕЩЕ РАЗ АВТОМАТИЗИРОВАТЬ
4. ПРИГОТОВЬТЕСЬ К БОЛЬШОМУ ОБЪЕМУ
НАГРУЗОЧНОГО ТЕСТИРОВАНИЯ

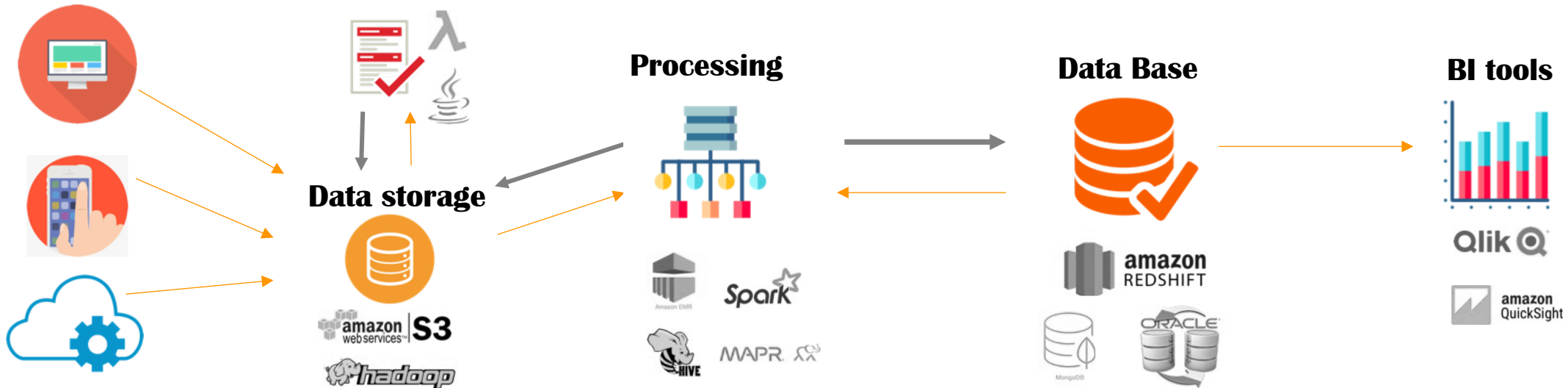
ПОНЯТЬ СВОИ ДАННЫЕ



ПОНЯТЬ СВОИ ПРОЦЕССЫ

Ежедневная обработка / Batch Pipeline

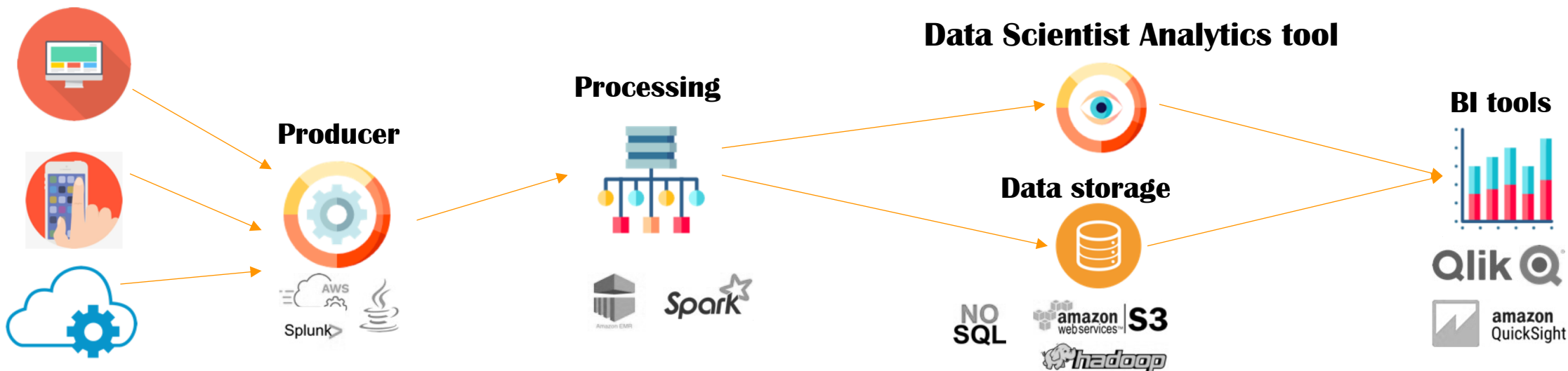
1. ГЕНЕРАЦИЯ ДАННЫХ
2. ПРОВЕРКА ДОСТОВЕРНОСТИ ДАННЫХ
3. ПРОВЕРКА ОБРАБОТКИ ДАННЫХ
4. ПРОВЕРКА ВЫВОДА



ПОНЯТЬ СВОИ ПРОЦЕССЫ

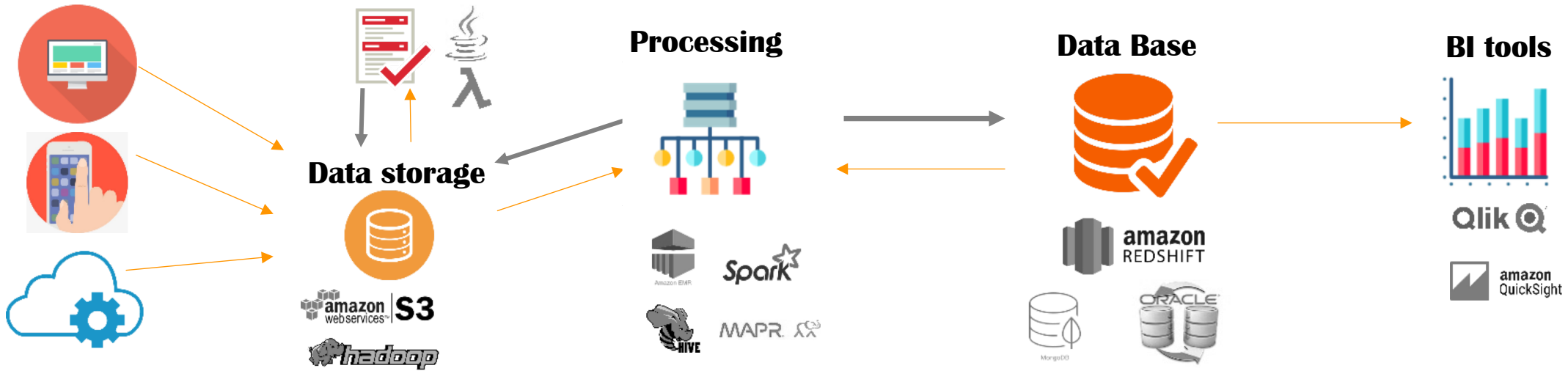
Обработка в реальном времени / Real Time Pipeline

1. СИМУЛЯЦИЯ АКТИВНОСТЕЙ
2. АНАЛИЗ СКОРОСТИ ОБРАБОТКИ ДАННЫХ
3. АНАЛИЗ МЕТРИК
4. ПРОВЕРКА ВЫВОДА



TOOLSET ДЛЯ ТЕСТИРОВАНИЯ

Batch Pipeline



Data generation

Acceptance testing

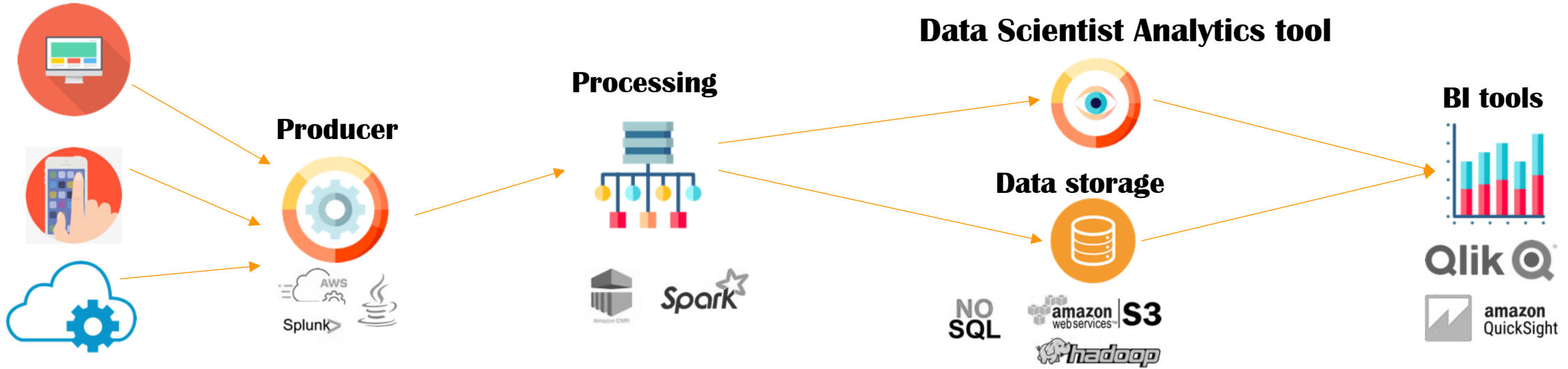
Data validation testing

UI testing



TOOLSET ДЛЯ ТЕСТИРОВАНИЯ

Real Time Pipeline



Activity simulation



Performance testing



Acceptance testing



UI testing



VDD СЦЕНАРИИ ДЛЯ ПРИЁМОЧНОГО ТЕСТИРОВАНИЯ

@Daily @VOD @Guest @Content @1day

Scenario: Total views of the guest users by content type = '%VOD%'

Given following data was ingested and aggregated

partition_date	action_type	user_id	platform_type	count	content_type
@date	accessesviews	@guest_1	@platform_0	1	content_1
@date	accessesviews	@guest_2	@platform_0	1	content_2
@date	accessesviews	@guest_4	@platform_0	1	content_1
@date	accessesviews	@guest_3	@platform_0	1	content_1
@date	accessesviews	@user_id_1	@platform_1	1	content_2

When I execute following <query>

.....

```
SELECT views
FROM agg_views_daily
WHERE partition_date = '@date' and user_type like '@guest' and contenttype like '@content_1'
```

.....

Then I expect to receive following results

views
3

BDD СЦЕНАРИИ ДЛЯ ВАЛИДАЦИИ ДАННЫХ

@Negative @Mandatory

Scenario: Check existing of mandatory fields for Join Flow

Given Generate Join csv with this data

TENANTID	USERNAME	UNIQUECONTRACT	TIMESTAMP
mandatory_fields_for_Join		1000000000000	20160301000001
mandatory_fields_for_Join	test_username		20160301000001
mandatory_fields_for_Join	test_username	1000000000000	

When I execute following query

.....

```
SELECT REASON
FROM JOIN_REJECTED
WHERE BODY like '%mandatory_fields_for_Join%'
ORDER BY REASON
```

.....

Then I expect to receive following results

reason
Mandatory field 'timestamp' is empty
Mandatory field 'uniquecontract' is empty
Mandatory field 'username' is empty

BDD И DDT ДЛЯ ГЕНЕРАЦИИ ДАННЫХ

index	group	id#	partition_date	user_type	sessionid	action_type	count	platform_type	currency	revenue
287	agg_purchases	purchases (7)	date+2	user_id_1	session_0	purchase	1	platform_0	EUR	12
288	agg_purchases	purchases (7)	date+2	user_id_8	session_0	purchase	2	platform_0	EUR	5
289	agg_purchases	purchases (7)	date+2	user_id_9	session_0	purchase	1	platform_0	USD	5
290	agg_purchases	purchases (7)	date+2	user_id_11	session_11	purchase	1	platform_11	EUR	2

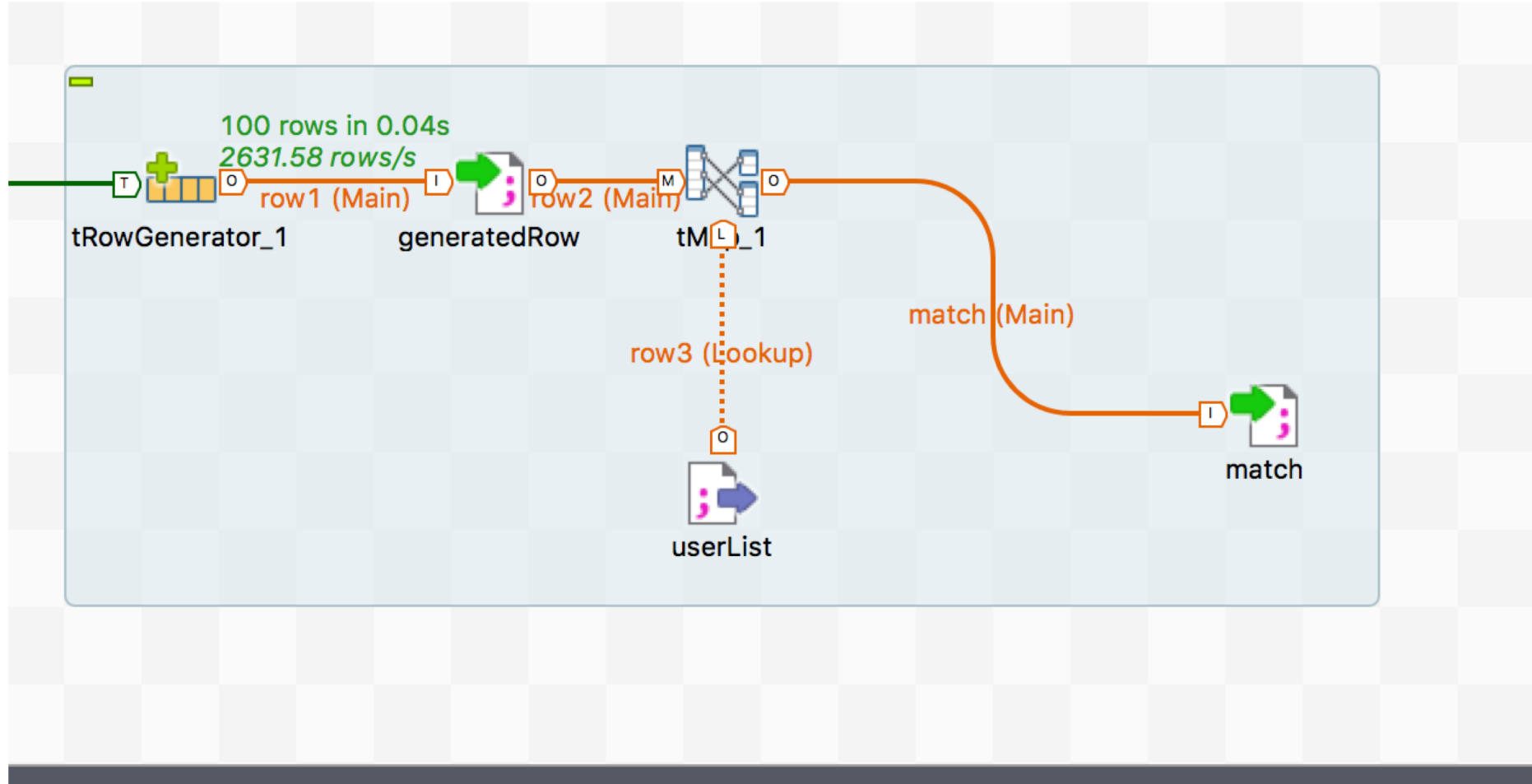
@Revenues @Daily @Platform @Currency

Scenario: A_162, Revenues (USD, GBP, EUR) by currency = (EUR and USD) and platform and state for daily agg

Given following data was ingested and aggregated

partition_date	action_type	user_id	region	platform_type	currency	revenue	count
@date+2	purchase	@user_id_1	@region_1	@platform_0	EUR	12	1
@date+2	purchase	@user_id_8	@region_8	@platform_0	EUR	5	1
@date+2	purchase	@user_id_8	@region_8	@platform_0	EUR	5	1
@date+2	purchase	@user_id_9	@region_9	@platform_0	USD	5	1
@date+2	purchase	@user_id_11	@region_11	@platform_11	EUR	2	1

TALEND OPEN STUDIO



TALEND OPEN STUDIO



Talend Open Studio for Big Data - tRowGenerator - tRowGenerator_1

Schema	Functions	Prev
Column	Type	Functions Environment variables Pre
id	Integer	Numeric.random(int,int) min value=>10000 ; max val...
timestamp	Date	TalendDate.getRandomDate(... min=>"2017-01-01" ; max=...
name	String	TalendDataGenerator.getFirst...
action	String	...
age	Integer	Numeric.random(int,int) min value=>12 ; max value=...
gender	String	... "M" , "F"

Columns ▾ Number of Rows for RowGenerator

Function parameters **Preview**

Number of Rows for Preview

	id	timestamp	name	action	age	gender
1	17116	10-04-2017	Martin	purchase	87	F
2	42918	10-08-2017	Warren	purchase	44	M
3	49261	07-01-2017	Thomas	purchase	66	F
4	43025	05-08-2017	George	view	16	F
5	77766	18-04-2017	Theodore	login	92	M
6	67162	17-06-2017	Bill	login	39	F
7	18906	07-03-2017	Herbert	view	34	F
8	56546	28-09-2017	Theodore	login	68	F
9	14045	12-08-2017	Woodrow	login	55	F
10	57630	08-05-2017	William	view	82	M



TALEND OPEN STUDIO



Talend Open Studio for Big Data - tMap - tMap_1

Find :

Auto map!

row2

Column
id
timestamp
action
token

row3

Expr. key	Column
row2.id	id
	name
	gender
	age

match

row2.action == "login"

Expression	Column
row2.id	id
row3.name	name
row3.gender	gender
row3.age	age

Schema editor | Expression editor

Column	Key	Type	<input checked="" type="checkbox"/> Nullab	Date Pattern (Ctrl+Spac Length)	Precision	Default	Comment
id	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>				
timestamp	<input type="checkbox"/>	Date	<input checked="" type="checkbox"/>	"dd-MM-yyyy"			
action	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>				
token	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>				

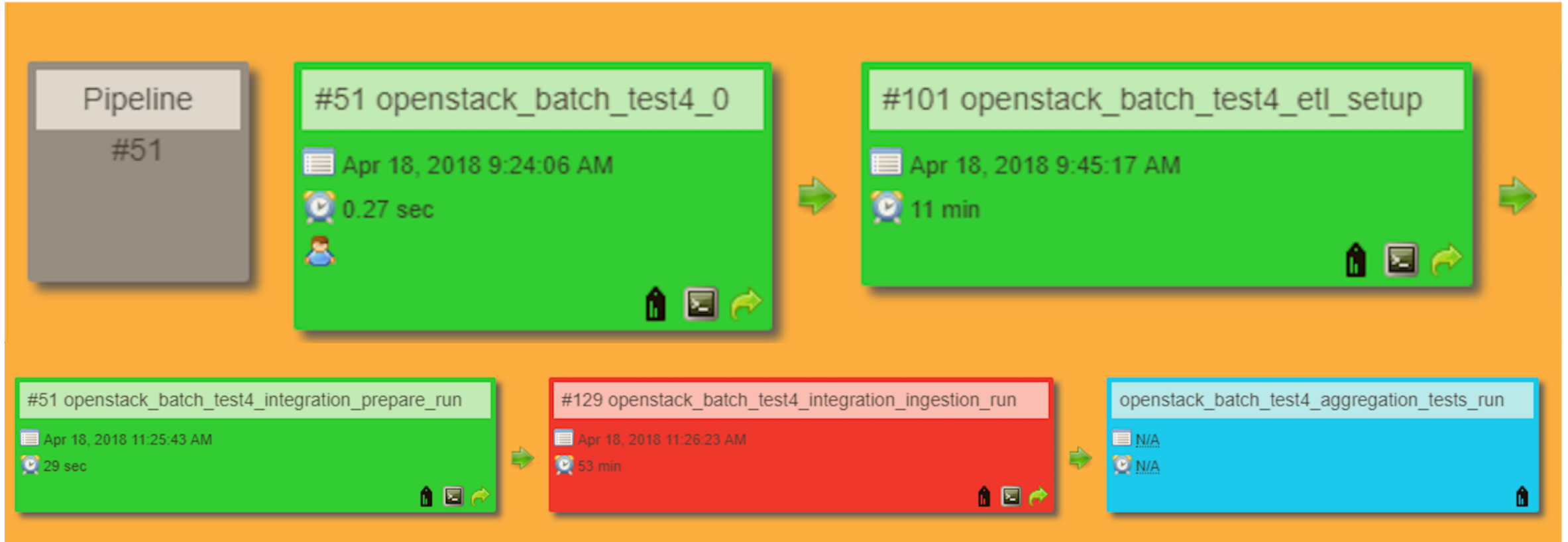
Column	Key	Type	<input checked="" type="checkbox"/> Nullab	Date Pattern (Ctrl+Spac Length)	Precision	Default	Comment
id	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>				
name	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>				
gender	<input type="checkbox"/>	String	<input checked="" type="checkbox"/>				
age	<input type="checkbox"/>	Integer	<input checked="" type="checkbox"/>				

TALEND OPEN STUDIO

The screenshot displays a grid of 13 data integration components in Talend Open Studio. Each component is represented by an icon and a label. The components are arranged in two rows. The first row contains seven components, and the second row contains six components. Each component has a small icon in the top right corner indicating its status: a red exclamation mark for errors and a yellow triangle for warnings.

Component Name	Status
tSSH_1	Error
tSortRow_1	Warning
tMap_1	Warning
tAggregateRow_1	Warning
tUnite_1	Error
tLoop_1	Warning
tREST_1	Warning
tJava_1	Warning
tSplunkEventCollector_1	Warning
tHDFSExist_1	Error
tUniqRow_1	Error
tHiveInput_1	Error
tFuzzyMatch_1	Error

CONTINUOUS INTEGRATION



BDD ОТЧЕТЫ

Features Statistics

The following graphs show passing and failing statistics for features

Scenarios



Feature	Steps						Scenarios			Features	
	Passed	Failed	Skipped	Pending	Undefined	Total	Passed	Failed	Total	Duration	Status
A_20 Accesses	15	0	0	0	0	15	5	0	5	1s 323ms	Passed
A_21 Registered Accesses	15	0	0	0	0	15	5	0	5	328ms	Passed
A_25 Accesses	15	0	0	0	0	15	5	0	5	304ms	Passed
A_26 Accesses with View	12	0	0	0	0	12	4	0	4	22s 783ms	Passed
A_27 Accesses without View	30	0	0	0	0	30	9	0	9	12s 995ms	Passed
A_28 Failed Login	15	0	0	0	0	15	5	0	5	1s 790ms	Passed
A_29 Successful Login	30	0	0	0	0	30	10	0	10	3s 292ms	Passed
A_42 Users Access More Platforms	3	0	0	0	0	3	1	0	1	081ms	Passed
A_316 Lost Users	2	0	0	0	0	2	1	0	1	071ms	Passed
A_317 Reconnected Users	3	0	0	0	0	3	1	0	1	056ms	Passed
A_318 Loyal Users	3	0	0	0	0	3	1	0	1	064ms	Passed



CONTACTS

Visit us @ Accenture.lv



Konstantin Pletenev

Test Automation Engineering Specialist
konstantin.pletenev@accenture.com

Тестирование в мире **BIG DATA**

SQA[®]
—————
DAYS#23