

s

a

n

a

o

m

**Building a data  
analytics platform  
with Hadoop,  
Python and R**

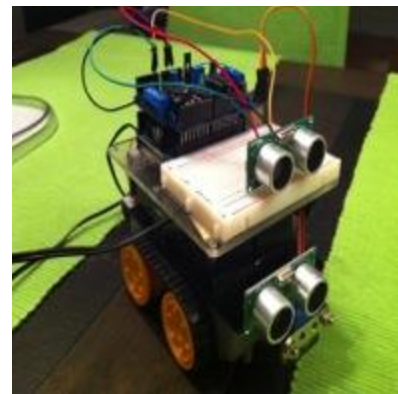
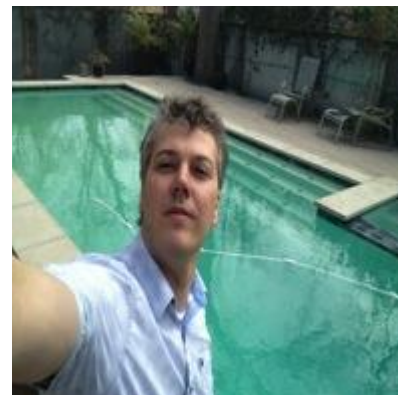
# Agenda

- Me
- Sanoma
- Past
- Present
- Future

# /me

## @skieft

- Software architect for Sanoma
- Managing the data and search team
- Focus on the digital operation
  
- Work:
  - Centralized services
  - Data platform
  - Search
  
- Like:
  - Work
  - Water(sports)
  - Whiskey
  - Tinkering: Arduino, Raspberry Pi, soldering stuff



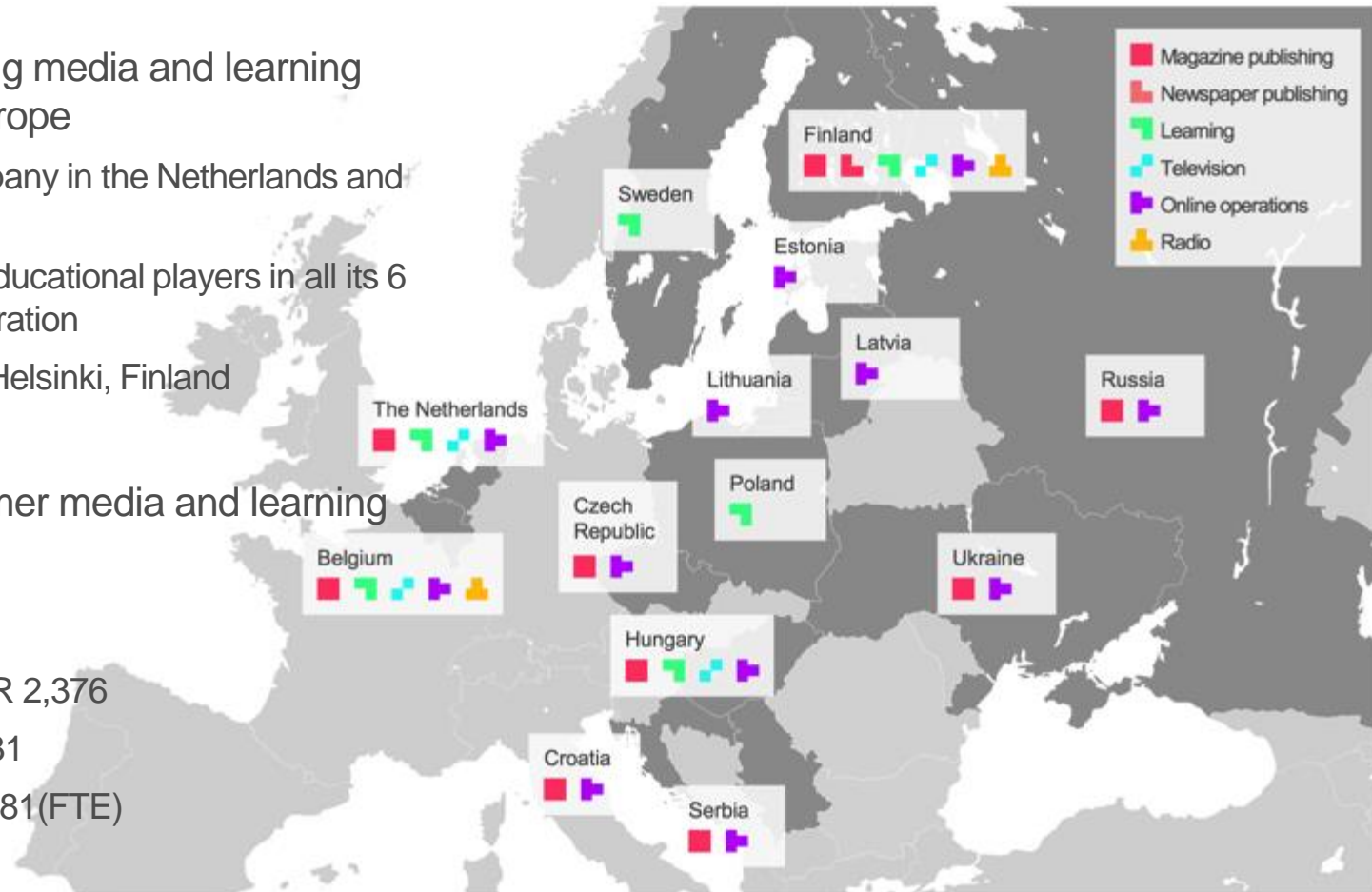
# Sanoma: Market leader in chosen businesses and markets

- One of the leading media and learning companies in Europe
  - #1 media company in the Netherlands and Finland
  - Among top 2 educational players in all its 6 markets of operation
  - Head office in Helsinki, Finland

- Focus on consumer media and learning

- 2012 financials

- Netsales mEUR 2,376
- EBIT mEUR 231
- Personnel 10,381(FTE)



Over 300  
magazines







# Past



# History

< 2008 2009 2010 2011 2012 2013





# Self service

# Self service levels

Personal	Departmental	Corporate
<b>Full self service</b>	<b>Support with publishing dashboards and data loading</b>	<b>Full service and support on dashboard</b>
Information is created by end users with little or no oversight. Users are empowered to integrate different data sources and make their own calculations.	Information has been created by end users and is worth sharing, but has not been validated.	Information that has gone through a rigorous validation process can be disseminated as official data.



# History

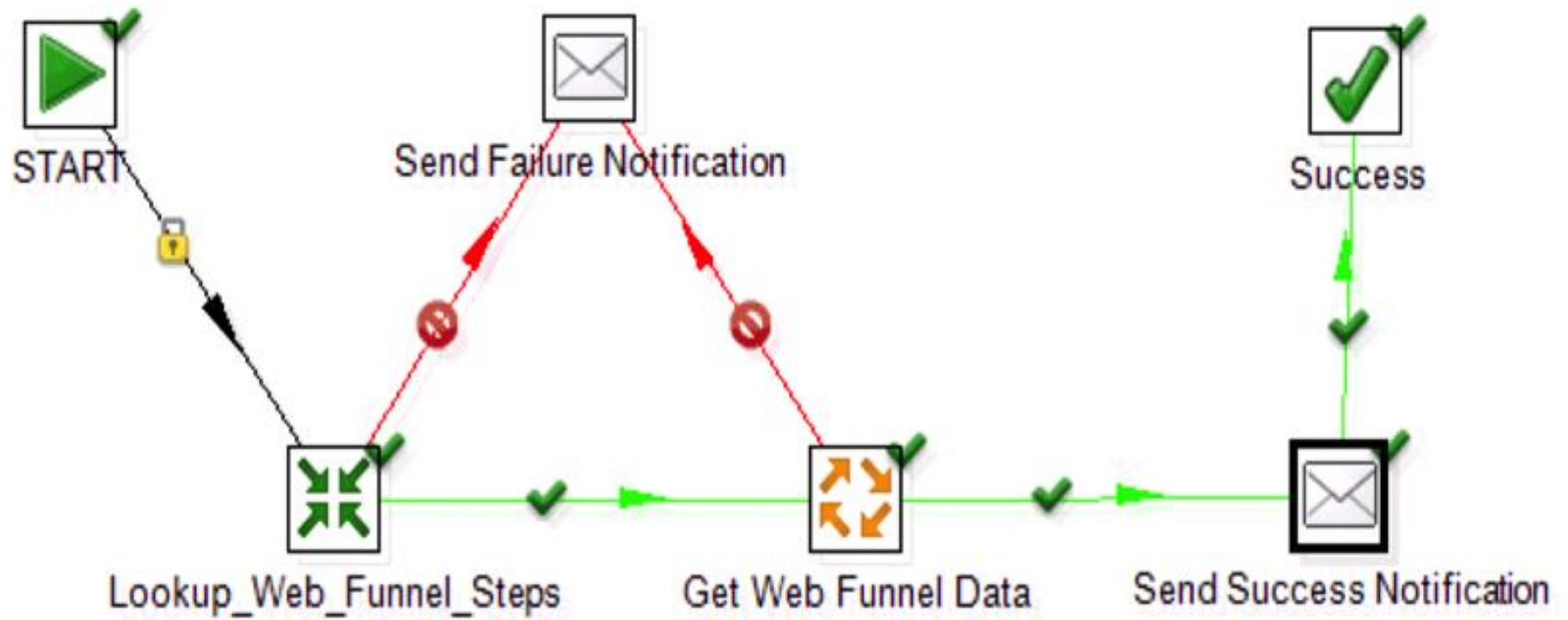
< 2008 2009 2010 2011 2012 2013

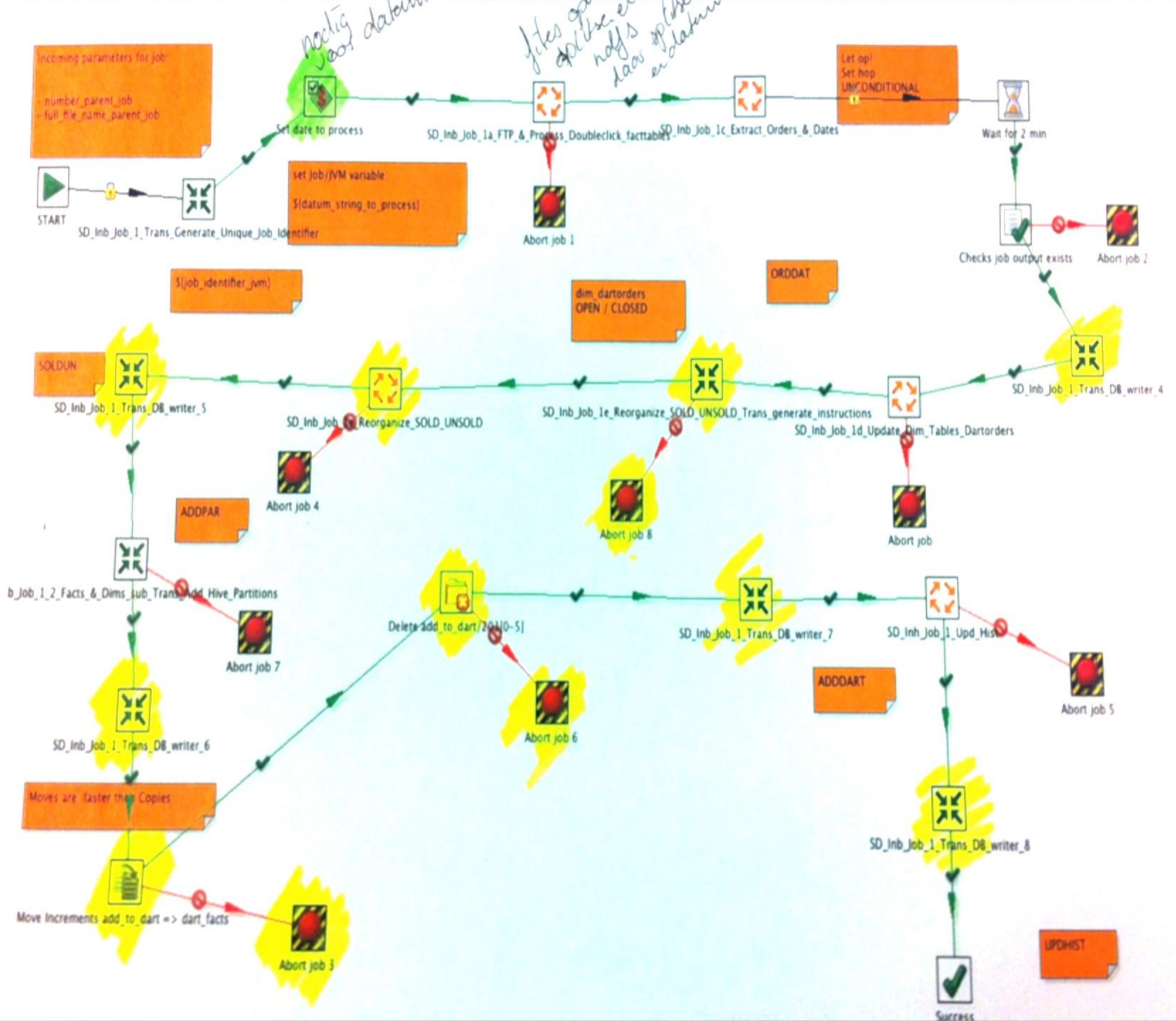


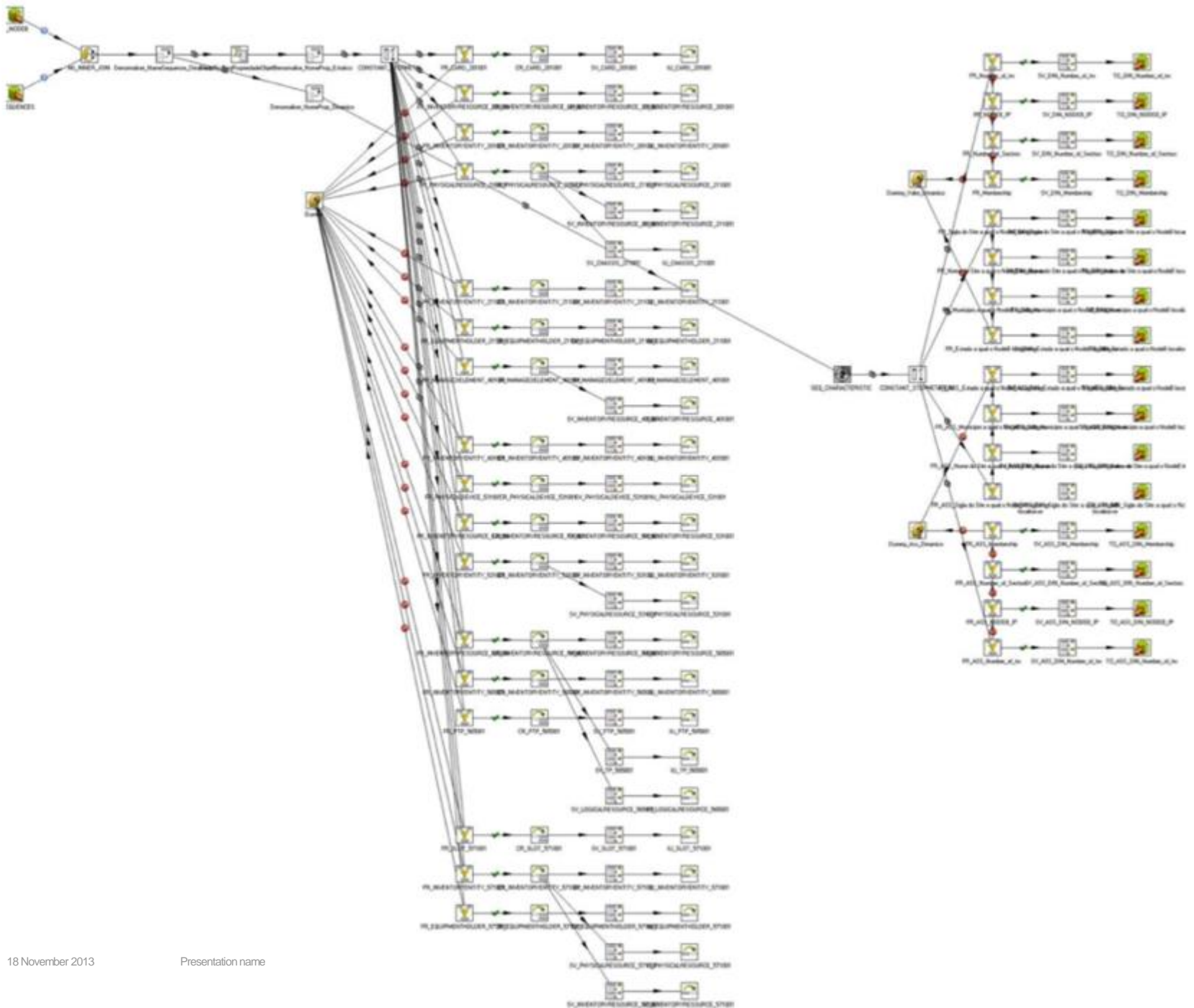
# Glue: ETL

- EXTRACT
- TRANSFORM
- LOAD











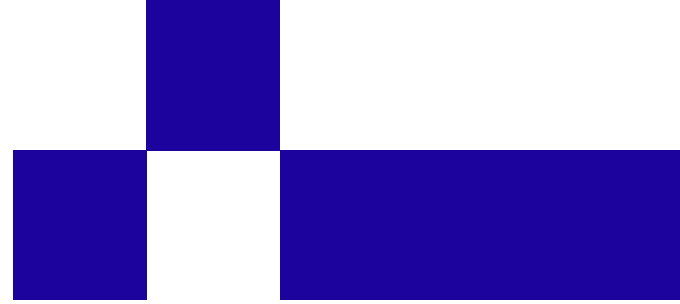


Shell

**#fail**



# Learnings



- Traditional ETL tools don't scale and are not effective for Big Data sources
- Big Data projects are not BI projects
- Doing full end-to-end integrations and dashboard development doesn't scale
- Qlikview was not good enough as the front-end to the cluster
- Hadoop requires developers not BI consultants

# History

< 2008 2009 2010 2011 2012 2013

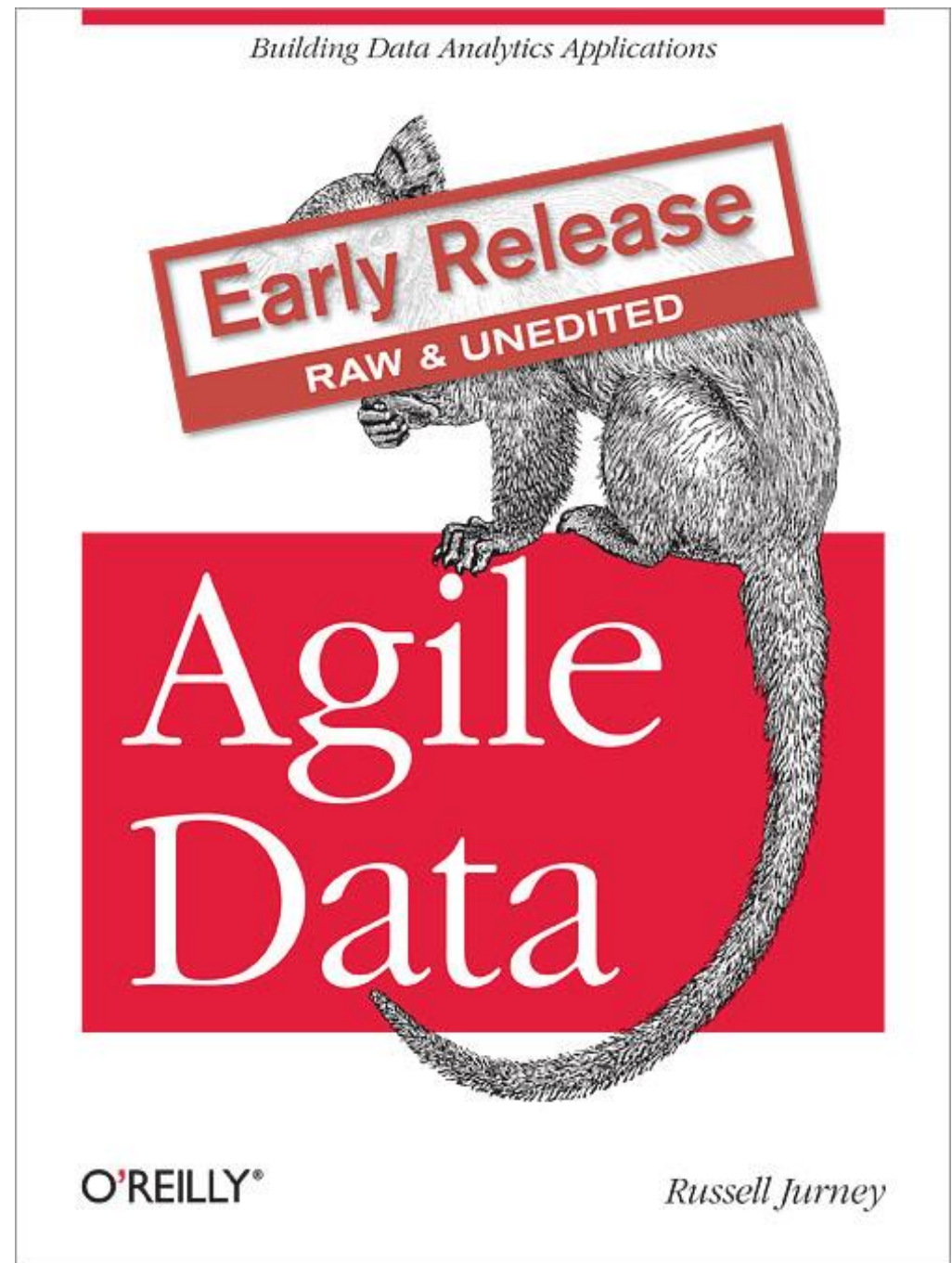


RESET



# Russell Journey

## Agile Data



**New glue**

# ETL Tool features

- Processing
- Scheduling
- Data quality
- Data lineage
- Versioning
- Annotating

	Standard for 2013 (Hourly)	Enterprise
Self-Administration Capable	✓	✓
Reporting	✓	✓
Dashboard	✓	✓
Analysis	✓	✓
HTML 5 Visualization	✓	✓
Data Federation	✓	✓
Data Integration (ETL)		Option
OLAP Server		Option
Multi-Tenancy		Option
Audit Logging		Details
Technical Support	Details	Details
Free Online Support	✓	✓
Full Professional Support	Option	✓
Deployment Support	Option	Option
Payment Terms	Hourly or Discounted Annual Subscription*	Annual Subscription





etl.sh

```
1  #!/bin/bash
2  set -ex
3
4  echo \
5  'open ftp.example.com
6  user username password
7  cd my-files
8  find /' > /tmp/ftp-list-command.txt
9
10 lftp -f /tmp/ftp-list-command.txt | grep "\.gz" > /tmp/full-names-on-ftp.txt
11
12 hadoop fs -ls -R /staging/my-files/ > /tmp/full-names-on-hadoop.txt
13 hadoop fs -ls -R /sources/my-files/ >> /tmp/full-names-on-hadoop.txt
14
15 cat /tmp/full-names-on-hadoop.txt | grep "\.gz" | awk '{print $8}' | awk -F '/' '{print $NF}' > /tmp/basenames-on-hadoop.txt
16
17 cat /tmp/full-names-on-ftp.txt | fgrep -v -f /tmp/basenames-on-hadoop.txt > /tmp/missing-files-on-hadoop.txt
18
19 echo \
20 'open ftp.example.com
21 user username password
22 ' > /tmp/ftp-get-files-command.txt
23 awk '{print "get -c " $0}' /tmp/missing-files-on-hadoop.txt >> /tmp/ftp-get-files-command.txt
24
25
26 cd /local-data/staging/my-files/
27 lftp -vf /tmp/ftp-get-files-command.txt
28 cd -
29
30 while read f
31 do
32     filename=$(echo $f | awk -F '/' '{print $NF}')
33     year=${filename:0:4}
34     month=${filename:4:2}
35     day=${filename:6:2}
36     site=$(echo $f | awk -F '-' '{print $NF}' | sed s/\.gz//g)
37
38     hadoop fs -mkdir "/staging/my-files/$year/$month/$day/$site"
39     hadoop fs -put "/local-data/staging/my-files/$filename" "/staging/my-files/$year/$month/$day/$site/$filename.upload"
40     hadoop fs -mv "/staging/my-files/$year/$month/$day/$site/$filename.upload" "/staging/my-files/$year/$month/$day/$site/$filename"
41 done < /tmp/missing-files-on-hadoop.txt
```

etl.jy

```
1 client = FTPClient(FTP_HOST, FTP_USER, FTP_PASSWORD)
2 remote_files = ftp_file_list(client, FTP_ROOT_DIR)
3
4 fs = hfs.FileSystem.get(hconf.Configuration())
5 hadoop_files = hdfs_file_list(fs, HDFS_ROOT_DIR_STAGING)
6 hadoop_files.extend(hdfs_file_list(fs, HDFS_ROOT_DIR_SOURCES))
7
8 missing_files = filter(lambda f: f not in hadoop_files, remote_files)
9
10 filename_pattern = re.compile(r"(\d{4})(\d{2})(\d{2})\d{4}-\d{12}-(\d{6})\.gz")
11
12 for missing_file in missing_files:
13     temp_file = File(FTP_LOCAL_DOWNLOAD_PATH, missing_file)
14     client.download(missing_file, temp_file)
15
16     year, month, day, site = filename_pattern.match(missing_file[1].name).groups()
17     hdfs_dir = hfs.Path('%s/%s/%s/%s/%s' % (HDFS_UPLOAD_BASE_PATH, year, month, day, site))
18
19     hdfs_dir = hfs.Path('%s/%s/%s/%s/%s' % tuple(filename_pattern.match(missing_file).groups()))
20
21     if (not fs.exists(hdfs_dir)):
22         fs.mkdirs(hdfs_dir)
23
24     fs.copyFromLocalFile(hfs.Path(temp_file.path), hfs.Path(remotedir, remotefilename + '.upload'))
25     fs.rename(hfs.Path(remotedir, remotefilename + '.upload'), hfs.Path(remotedir, remotefilename))
26
27     os.unlink(localfile.path)
28
29
30 fs.close()
31 client.disconnect(False)
32
```

# Processing - Jython

- No JVM startup overhead for Hadoop API usage
- Relatively concise syntax (Python)
- Mix Python standard library with any Java libs



# Scheduling - Jenkins

- Flexible scheduling with dependencies
- Saves output
- E-mails on errors
- Scales to multiple nodes
- REST API
- Status monitor
- Integrates with version control



# Jenkins

# ETL Tool features

- Processing – **Bash & Jython**
- Scheduling – **Jenkins**
- Data quality
- Data lineage
- Versioning – **Mecurial (hg)**
- Annotating – **Commenting the code** 😊

Feature	Option 1	Option 2
Self-Administration Using	✓	✓
Reporting	✓	✓
Dashboards	✓	✓
Analysis	✓	✓
HTML 5 Visualization	✓	✓
Data Federation	✓	✓
Data Integration (ETL)		Option
OLAP Server		Option
Multi-Tenancy		Option
Audit Logging		Details
Technical Support	Details	Details
Free Online Support	✓	✓
Full Professional Support	Option	✓
Deployment Support	Option	Option
Payment Terms	Hourly or Discounted Annual Subscription*	Annual Subscription



# Processes

# Independent jobs

**Source (external)**



HDFS upload + move in place

**Staging (HDFS)**



MapReduce + HDFS move

**Hive-staging (HDFS)**

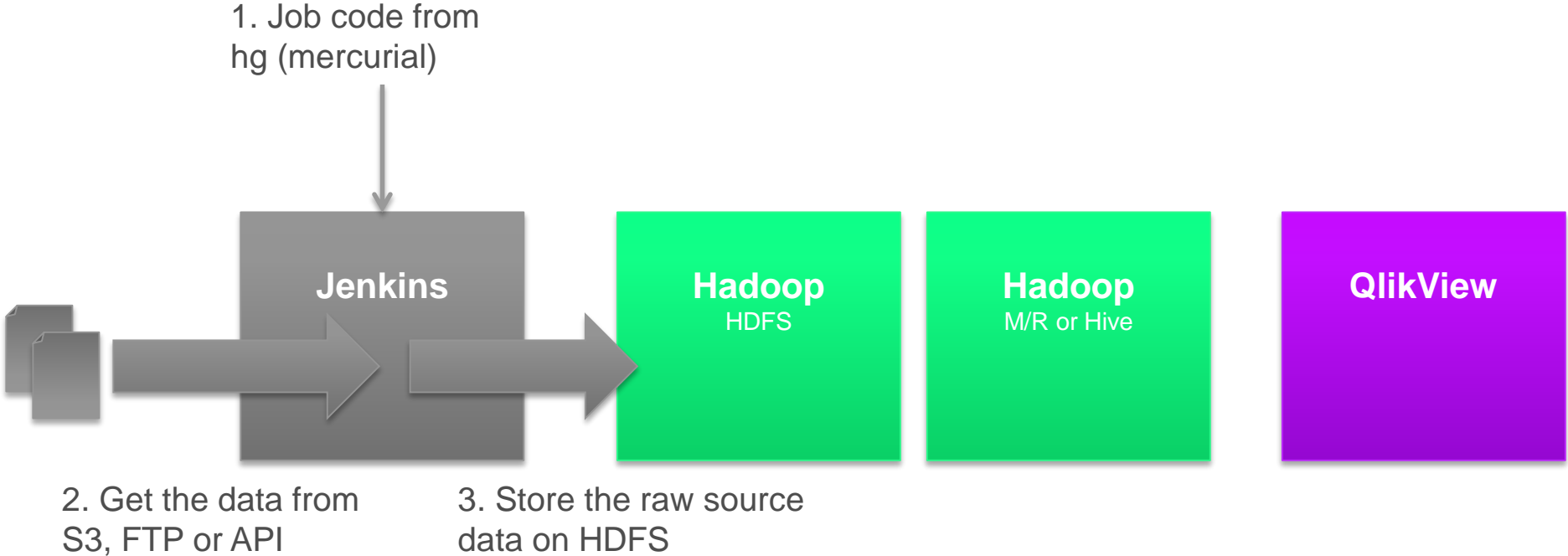


Hive map external table + SELECT INTO

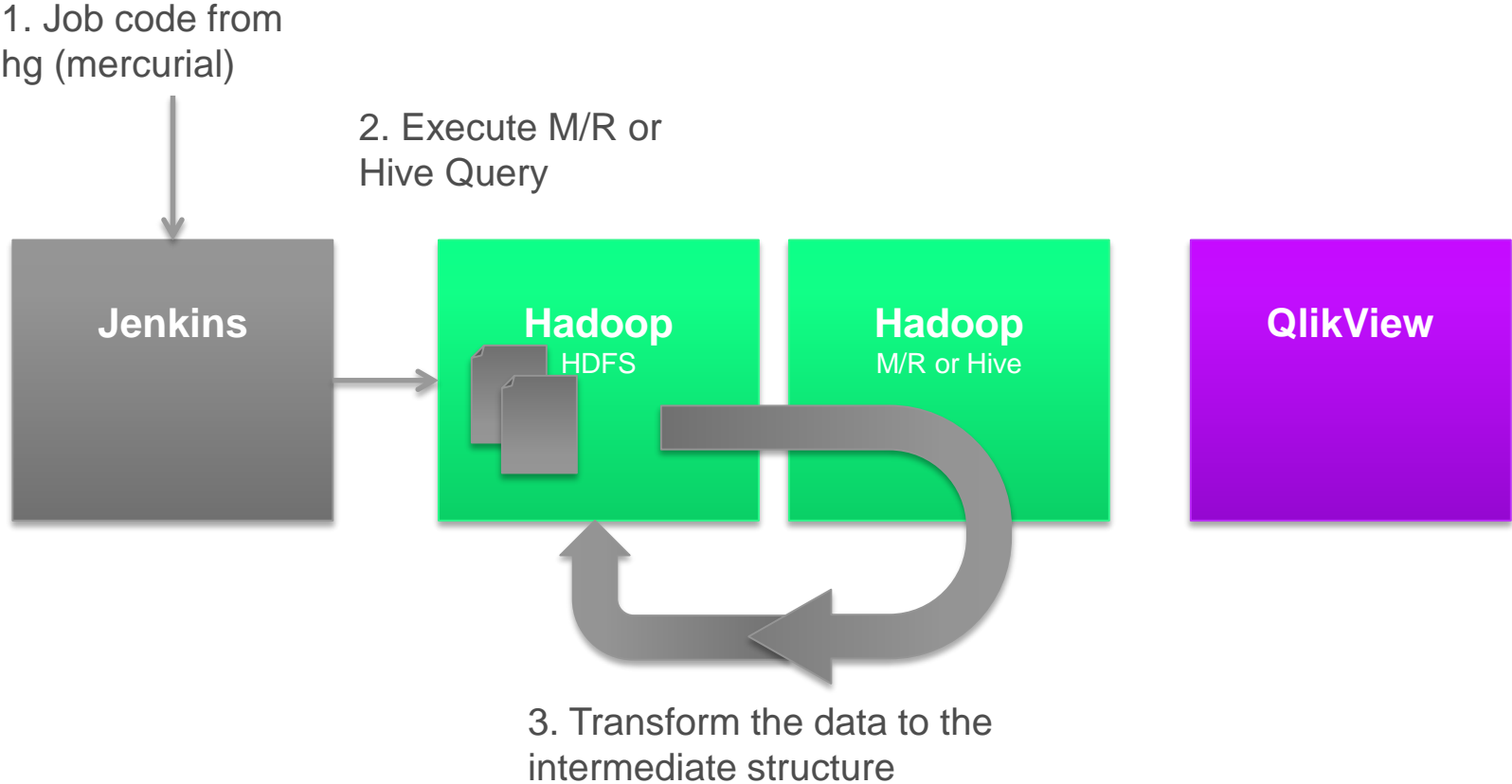
**Hive**



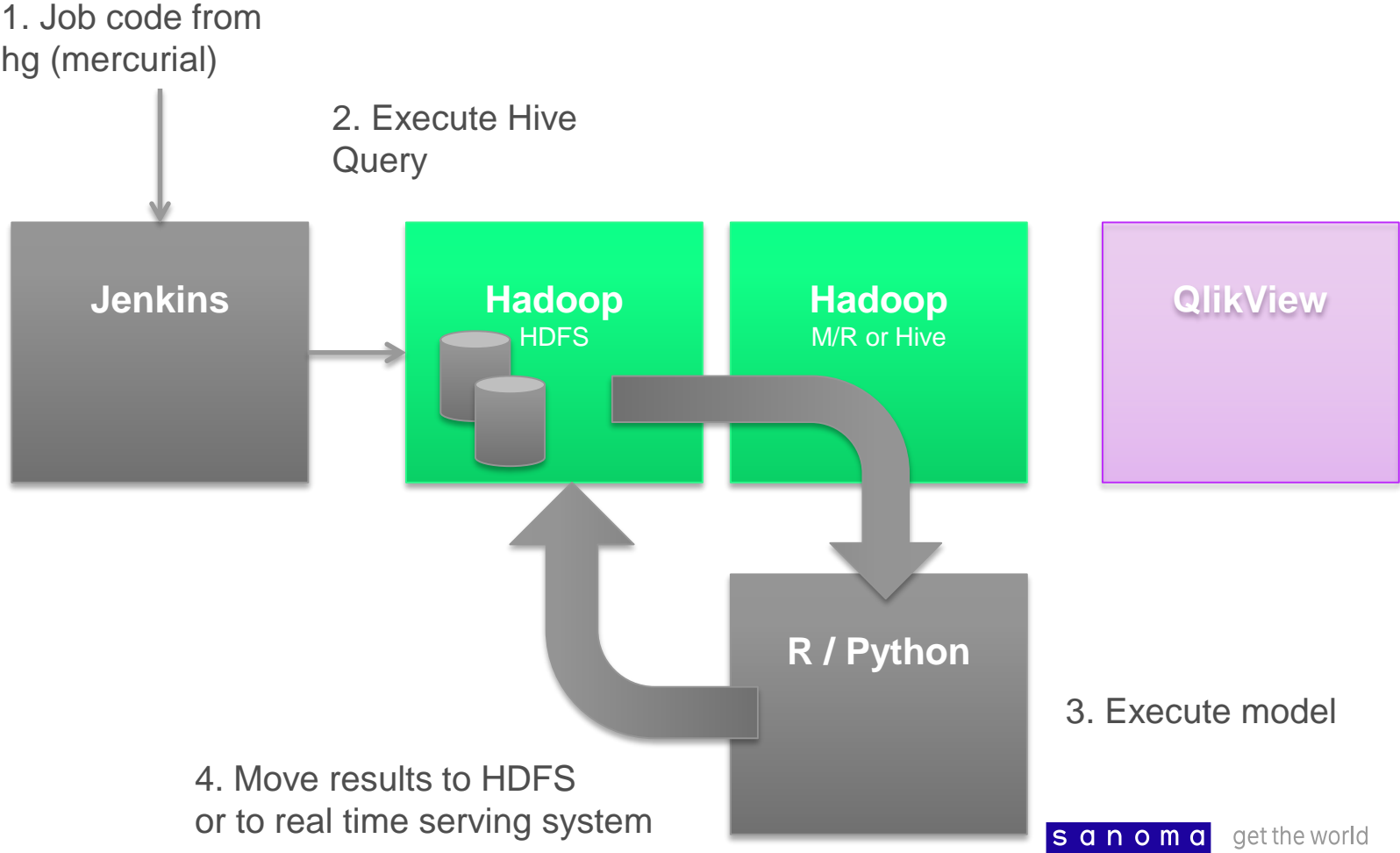
# Typical data flow - Extract



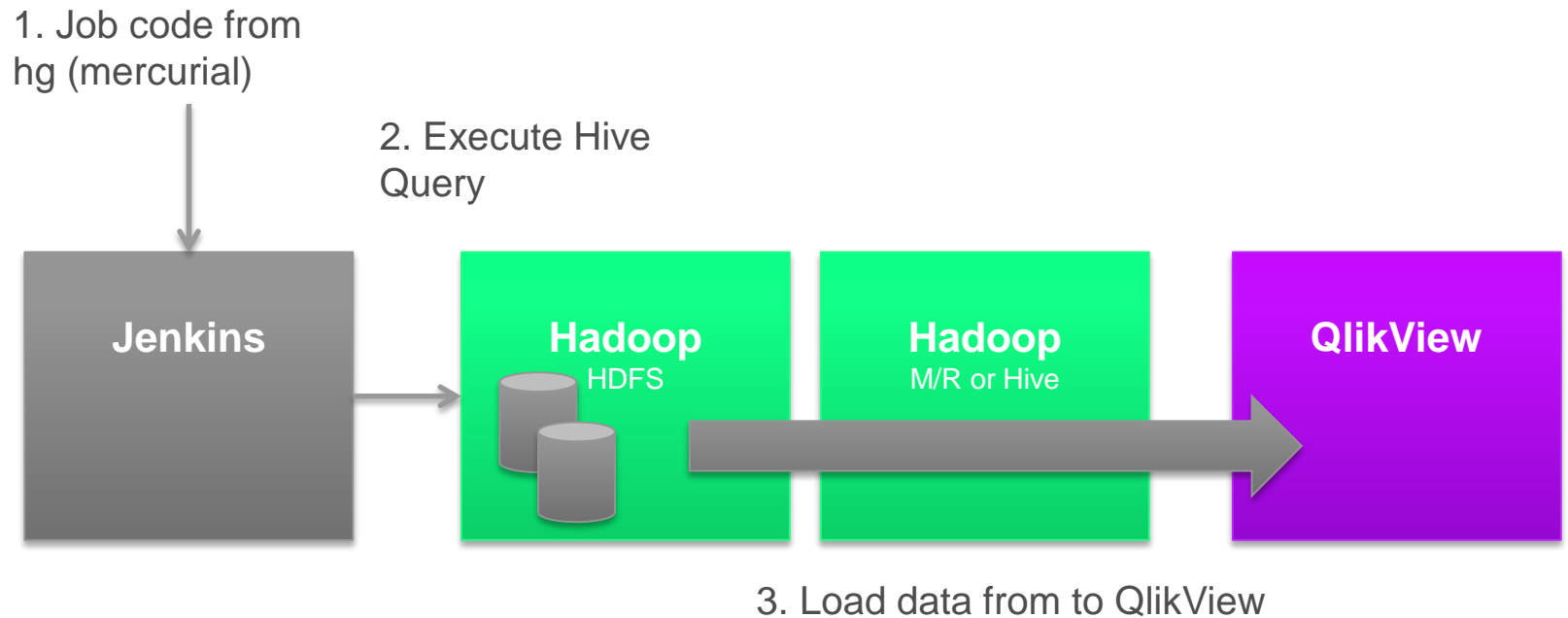
# Typical data flow - Transform



# Typical data flow - Load



# Typical data flow - Load



# Out of order jobs

- At any point, you don't really know what 'made it' to Hive
- Will happen anyway, because some days the data delivery is going to be three hours late
- Or you get half in the morning and the other half later in the day
- It really depends on what you do with the data
- This is where metrics + fixable data store help...

# Fixable data store

- Using Hive partitions
- Jobs that move data from staging create partitions
- When new data / insight about the data arrives, drop the partition and re-insert
- Be careful to reset any metrics in this case
- Basically: instead of trying to make everything transactional, repair afterwards
- Use metrics to determine whether data is fit for purpose

MIXED BUSINESS  
SELF SERVICE  
GENUINE SAVINGS

GR  
FR  
VE  
SA  
L



My Queries

Saved Queries

History

Settings

# Query Editor

```
1 SELECT to_date(creation_tstamp) dt
2     ,COUNT(*) clicks
3 FROM kis_clicks
4 GROUP BY to_date(creation_tstamp)
```

HUE

IZATION

Parameterization

IFICATION

on completion

Execute

Execute Into QV

Save as...

Explain

or create a

New query



# Hadoop job scheduling

- Schedulers to spread the load on your cluster

- CapacityScheduler:

vs

- FairScheduler:

- The default since: CDH 4.1

- Use scheduling pools to separate workloads. ETL vs User based vs Consuming Applications

# Quotas

- Since user can break stuff, easily.. ..and SQL skills get rusty

```
SELECT *  
FROM views  
LEFT JOIN clicks  
WHERE day = '2013-09-26'
```

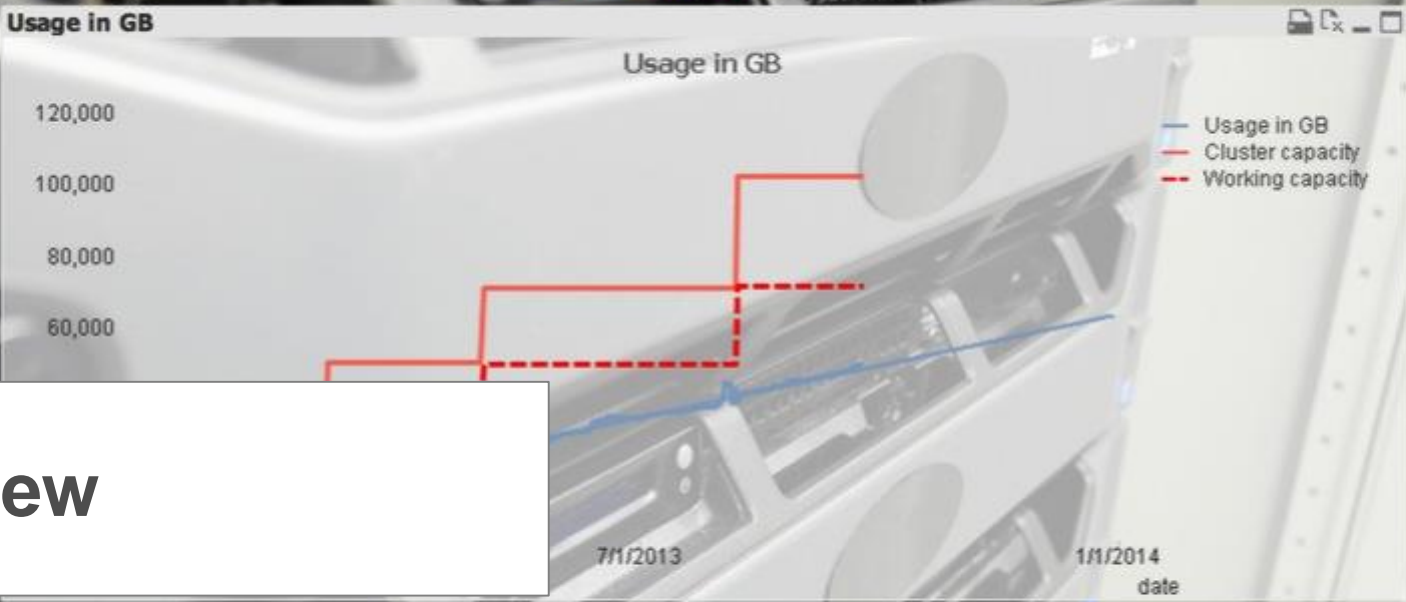
Views table is 10TB and  $36 \times 10^9$  rows

Clicks table is 100GB and  $36 \times 10^6$   
ROWS

**Or use Strict mode = 1 in hive**

Year selection  
2010 2011 2012 2013

ations  
2013



# QlikView

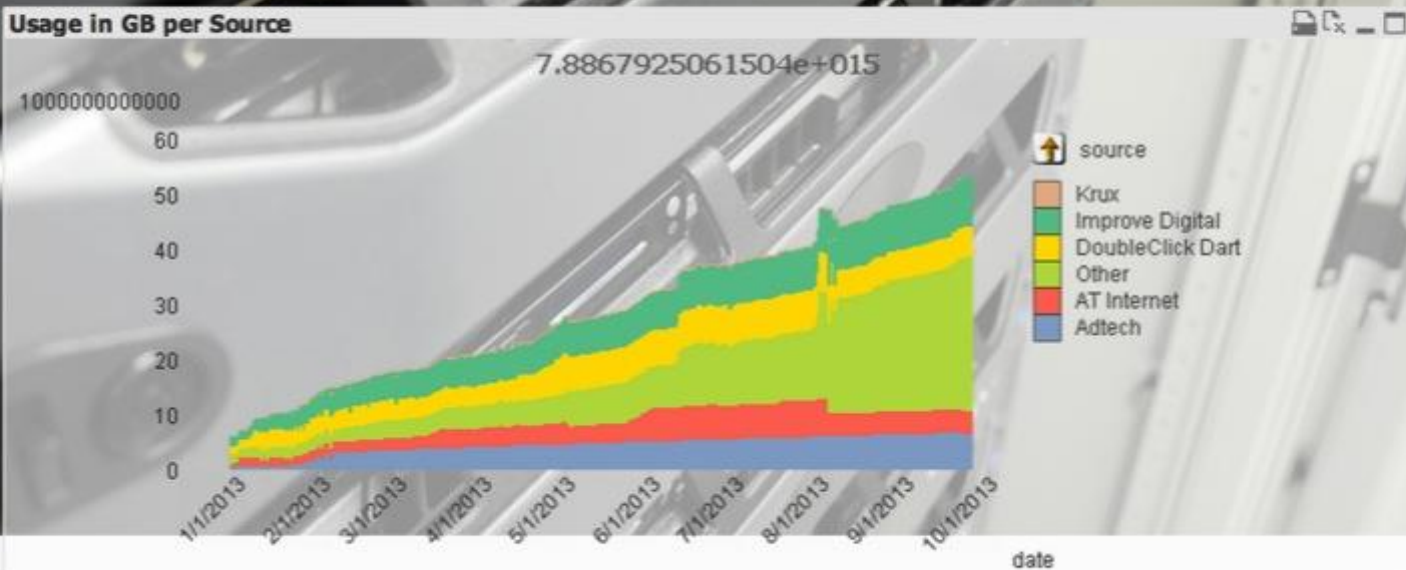
ers  
derzoek

events

bank  
agina  
agina\_dms  
doubler  
racker

mp  
erloop

320.67472614586  
345.45265582427  
29.804923847821  
132.05656672259  
378.84532864576  
36.164873433978



Sources  
source

Adtech  
AT Internet  
Other

```

conn = JDBC("
classPath = paste("/opt/cloudera/parcels/CDH/lib/hive/lib/hive-
"/opt/cloudera/parcels/CDH/lib/hive/lib/hive-
"/opt/cloudera/parcels/CDH/lib/hive/lib/libth
"/opt/cloudera/parcels/CDH/lib/hive/lib/hive-
"/opt/cloudera/parcels/CDH/lib/hive/lib/libft
lib/slf4j-
lib/commo

```

# R Studio Server

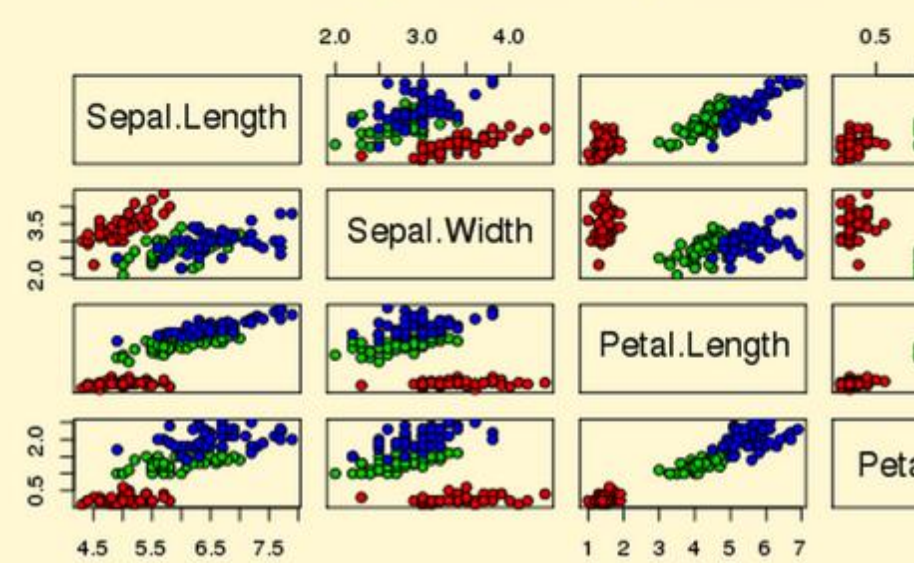
```

range(x))
range(y))
] / xdelta
] / ydelta
scale, yscale)
in[1] / scale - xdelta)
in[2] / scale - ydelta)
), numeric(0),
range(x) + c(-1, 1) * xadd, ylim = range(y) + c(-1, 1) * yadd,
", ann = FALSE)
see next plot: |

```

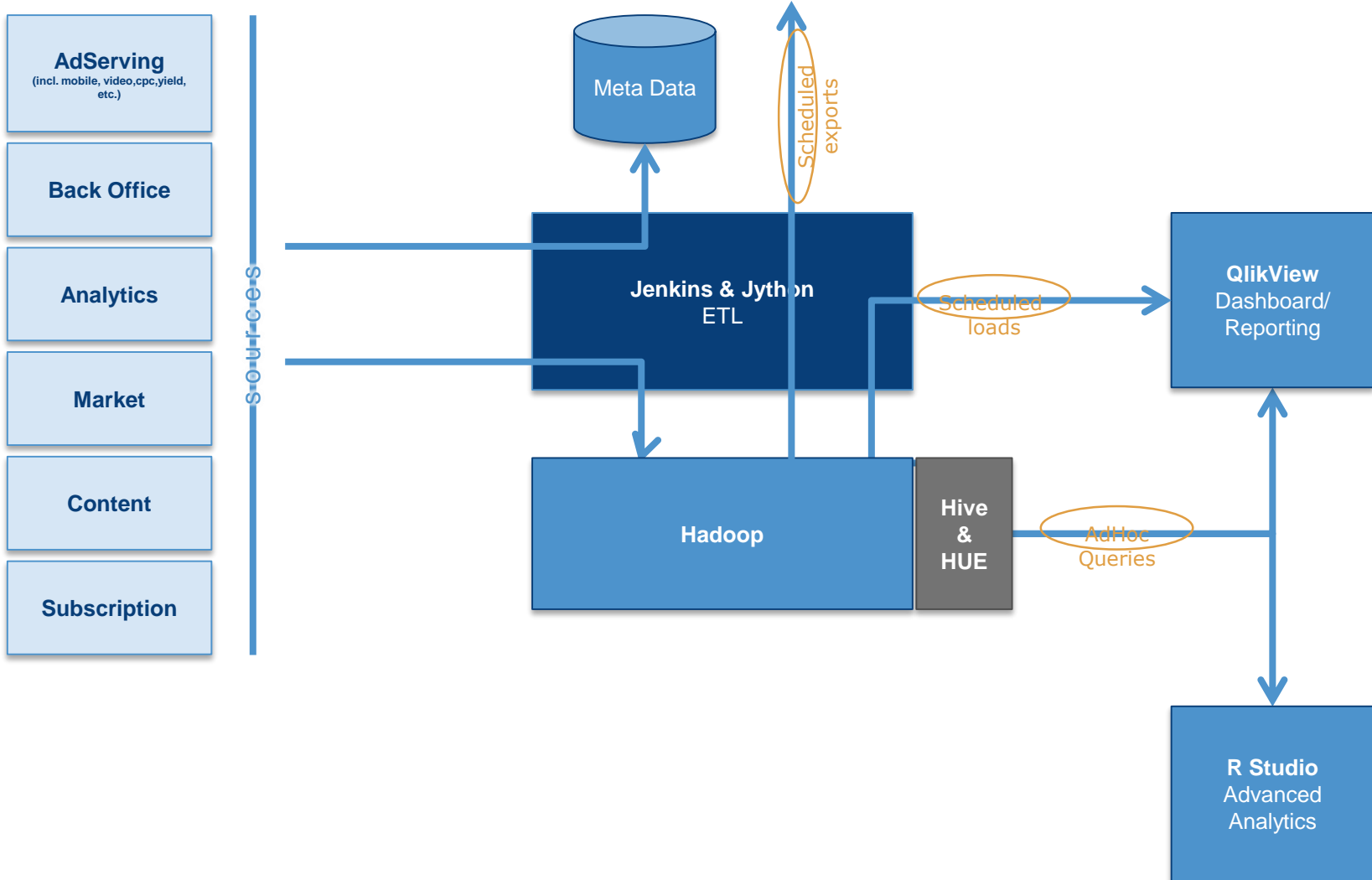
Workspace		History	
Import Dataset			
Values			
conn		JDBCConnection	[1]
drv		JDBCDriver	[1]
g		factor	[1000]
lev		numeric	[12]
n			100
opar		list	[1]
pie.sales		numeric	[6]

## Edgar Anderson's Iris Data



# Architecture

# High Level Architecture



# Sanoma Media The Netherlands Infra

## Colocation

Dev., test and acceptance VMs

Big Data platform

DC1

- NO SLA (yet)
- Limited support BD9-5
- No Full system backups
- Managed by us with systems department help

## Managed Hosting

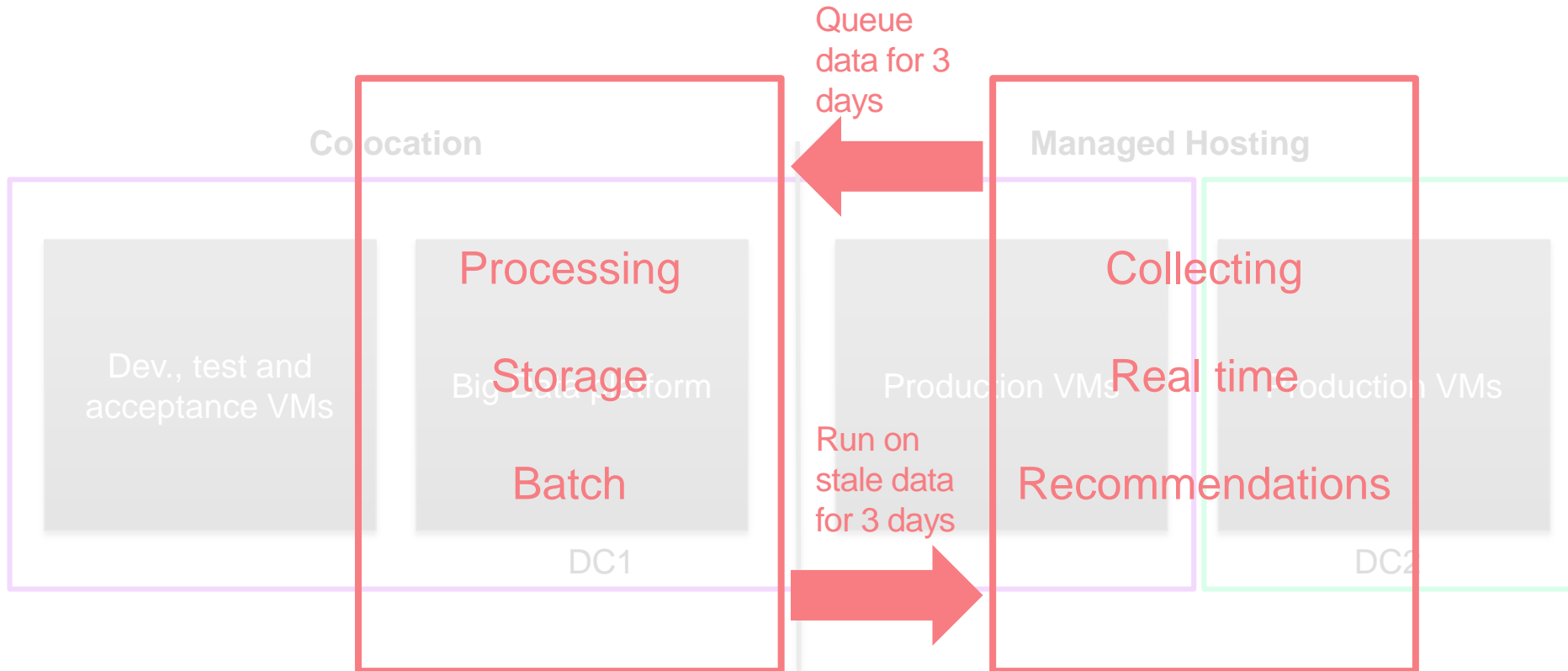
Production VMs

Production VMs

DC2

- 99,98% availability
- 24/7 support
- Multi DC
- Full system backups
- High performance EMC SAN storage
- Managed by dedicated systems department

# Sanoma Media The Netherlands Infra



- NO SLA (yet)
- Limited support BD9-5
- No Full system backups
- Managed by us with systems department help

- 99,98% availability
- 24/7 support
- Multi DC
- Full system backups
- High performance EMC SAN storage
- Managed by dedicated systems department



# Present

# Current state – Use case

- A/B testing + deeper analyses
- Ad auction price optimization
- Recommendations
- Search optimizations

# Current state - Usage

- Main use case for reporting and analytics
- Sanoma standard data platform, used by other Sanoma countries too: Finland, Hungary, ..
- ~ 100 Users: analysts & developers
- 25 daily users
- 43 source systems, with 125 different sources
- 300 tables in hive

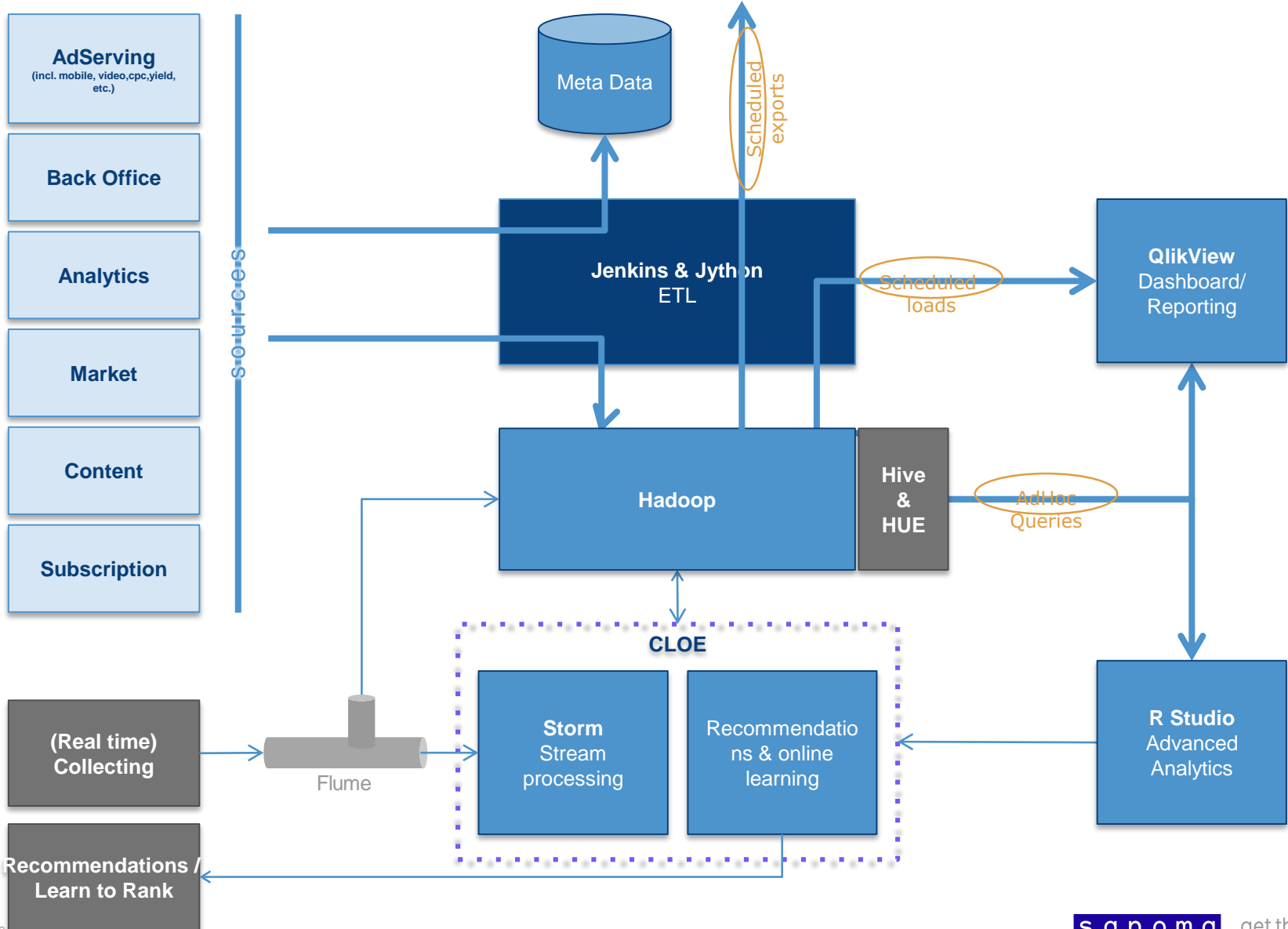
# Current state – Tech and team

- Team:
  - 1 product manager
  - 2 developers
  - 2 data scientists
  - 1 Qlikview application manager
  - ½ architect
- Platform:
  - 30-50 nodes
  - > 300TB storage
  - ~2000 jobs/day
- Typical data node / task tracker:
  - 2 system disks (RAID 1)
  - 4 data disks (2TB, 3TB or 4TB)
  - 24GB RAM

# Current State - Real time

- Extending own collecting infrastructure
- Using this to drive recommendations, user segmentation and targeting in real time
- Moving from Flume to Kafka
- First production project with Storm

# High Level Architecture

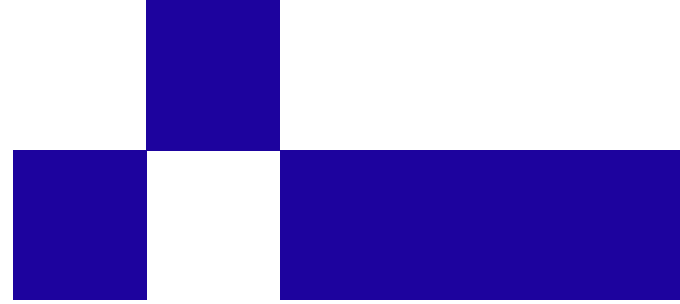




# Future



# What's next



- Cloudera search (Solr Cloud on Hadoop)
  - Index creation
  - External Ranker
  - Easier scaling and maintainance
- Moving some NLP (Natural Language Processing) and Image recognition workload to hadoop
- Optimizing Job scheduling (Fair Scheduling Pools)
- Automated query optimization tips for analysts
- Full roll out R integration, with rmr2
- More support for: cascading, scalding, pig, etc.



s a n

a

o

m

Thank you!  
Questions?

s a n o m a