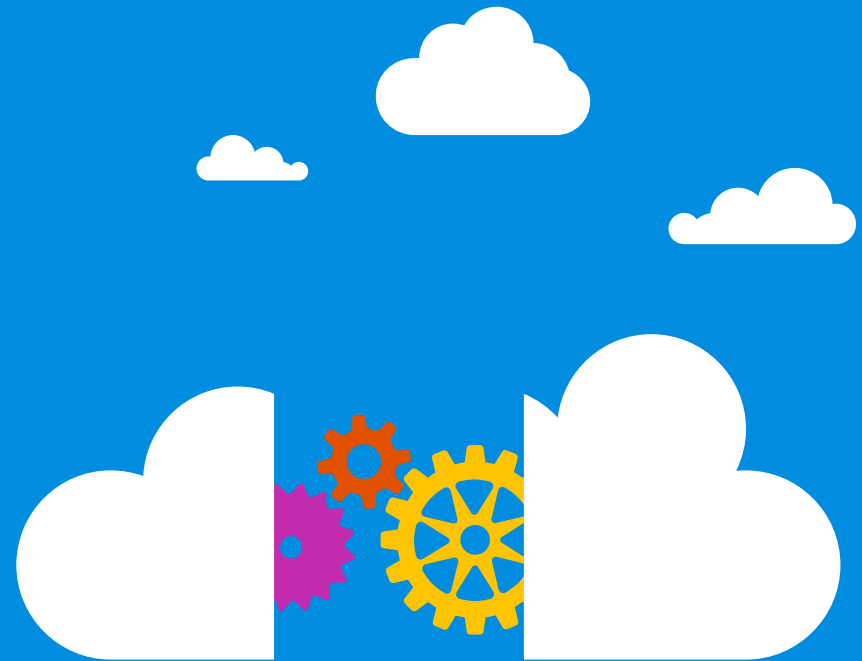


Solving low latency query over big data with Spark SQL

Julien Pierre



Agenda

- Intro
- Project's mission
- Use case
- Learnings
- Implementation details
- The future of Spark in ASG

Who am I?



Julien Pierre
朱力安

- Senior Product Manager
- ASG Shared Data team
- Beijing
- Working on:
 - Spark
 - Driving Microsoft internal Spark community
- Worked on:
 - Office 365 Admin
 - Office 365 Customer Facing Reporting
 - Open sourced the Office 365 Splunk App

What we do

Client Data Fluency

Office

Skype

Bing

Modern Data Capability

Instrumentation & Ingestion

Processing & Storage

Reporting & Analytics

Information Management

Mobile-First Analytics Experience

Experimentation

Project mission

Deliver interactive analytics capabilities for ASG teams to build great products

Increase analyst and engineer productivity

Capitalize on our data

Leverage open technology to run a stable service

Use Case

Where does Spark SQL help us?



Data issue troubleshooting use case



Kuntao Yu, engineer

Customer found an issue with the data

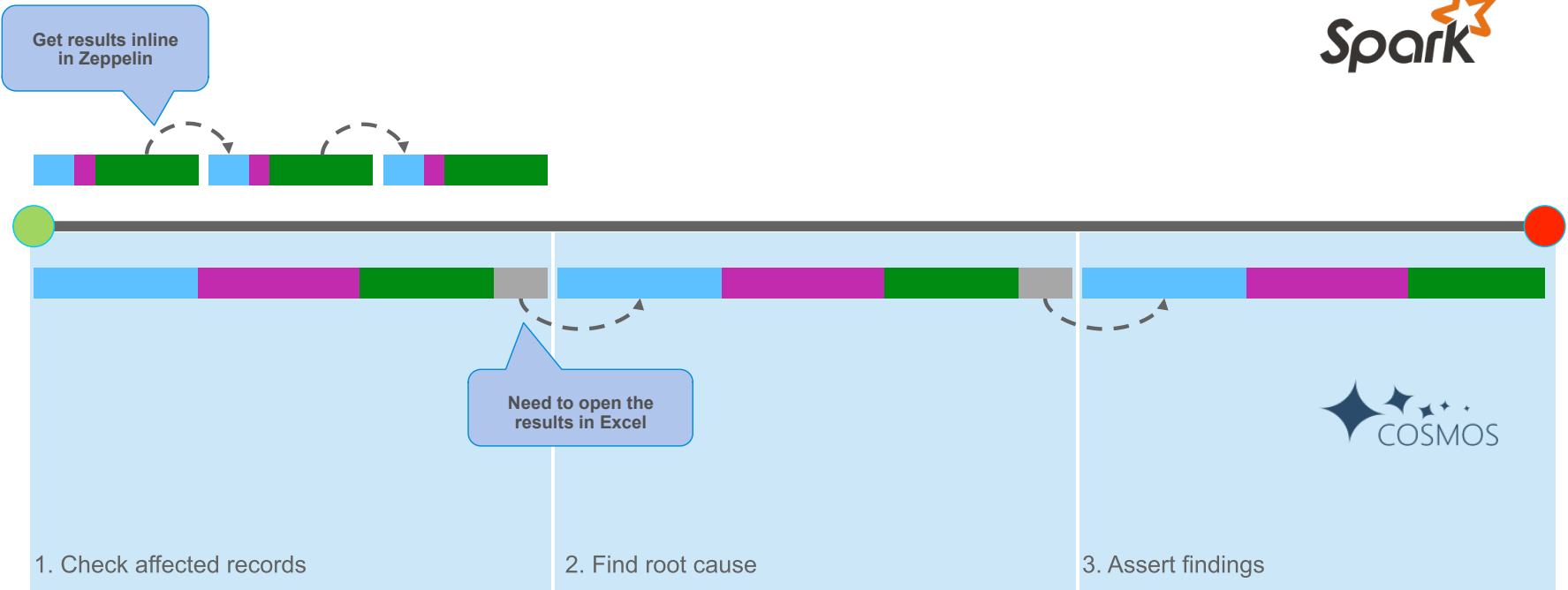
“Why is the Created Date field = 0001/01/01?”

Is it by design or a bug?

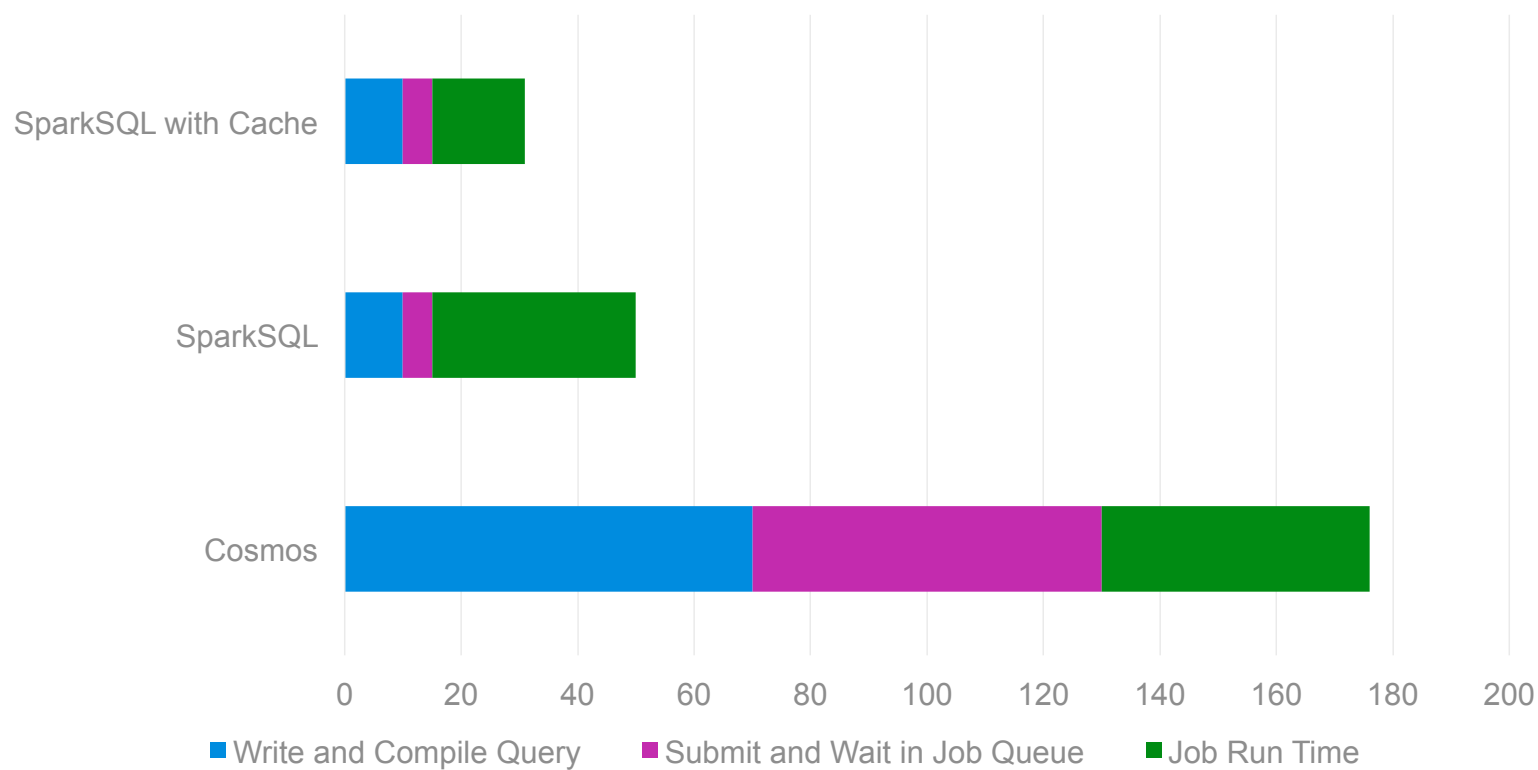
71.2 TB of Office Data in Spark

What did Kuntao have to do?

- Writing & compiling query
- Submitting & waiting in queue
- Job runtime
- Check data in Excel



Elapsed time



What did we learn in the
process?

Learnings

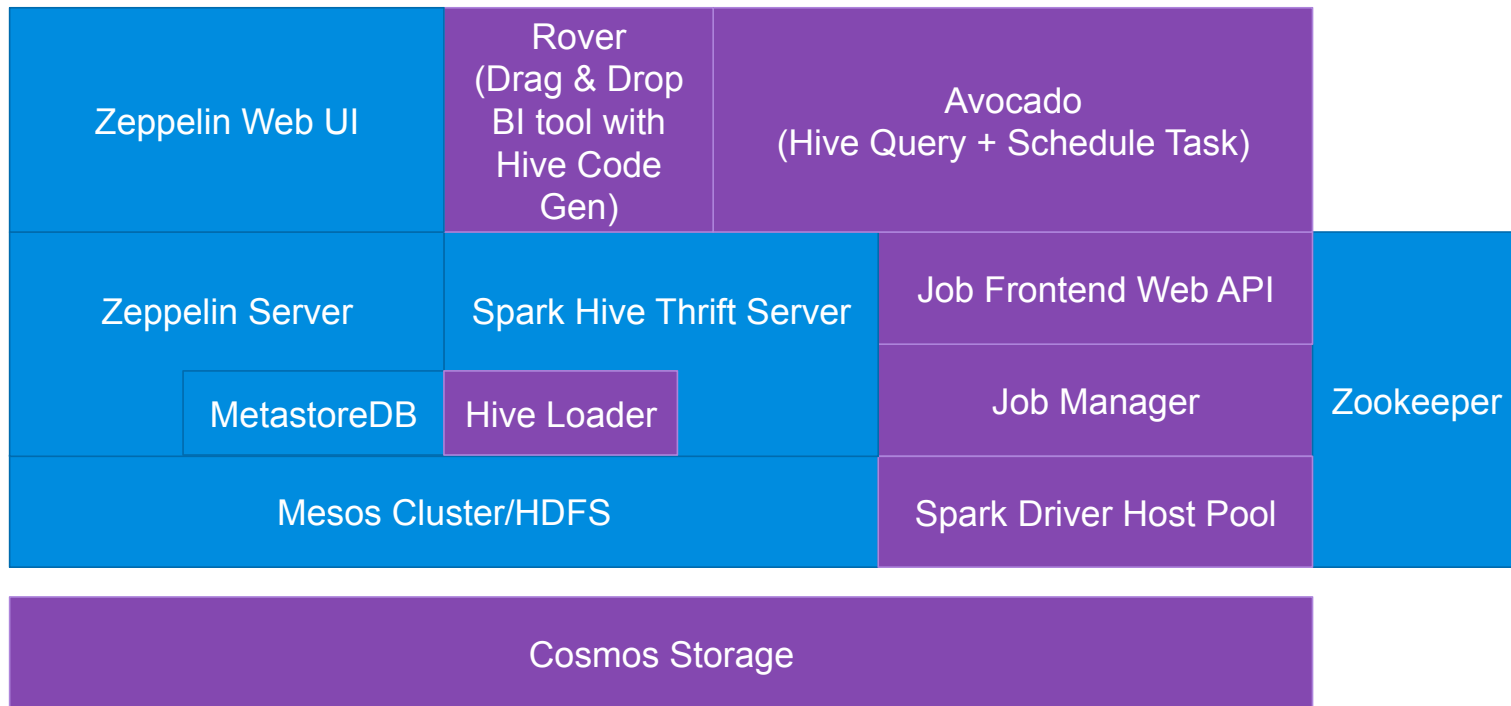
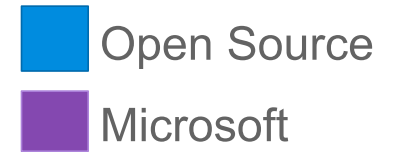
Enable self-service

Minimize time to insight

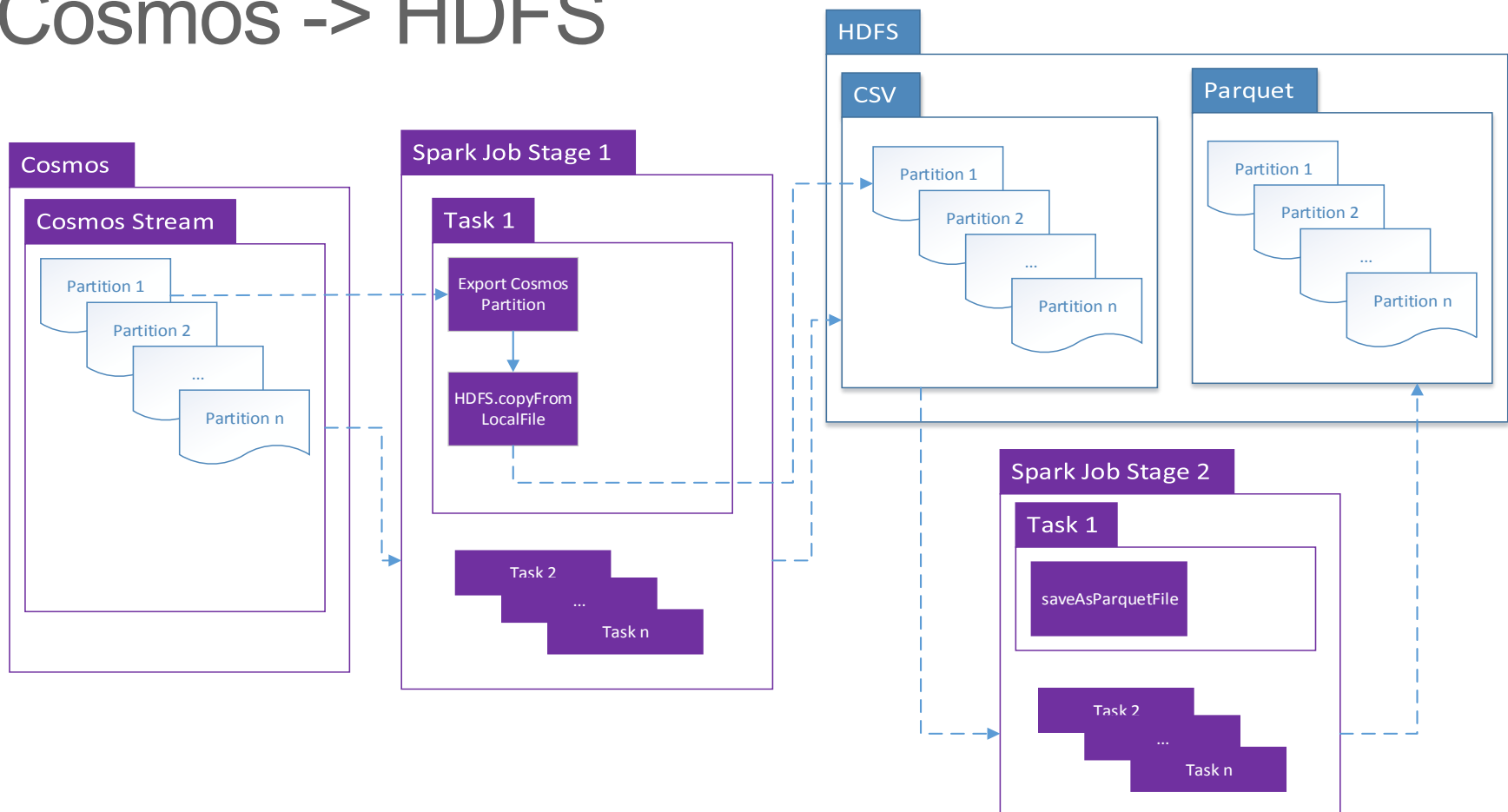
Notebooks at scale

How did we do it?

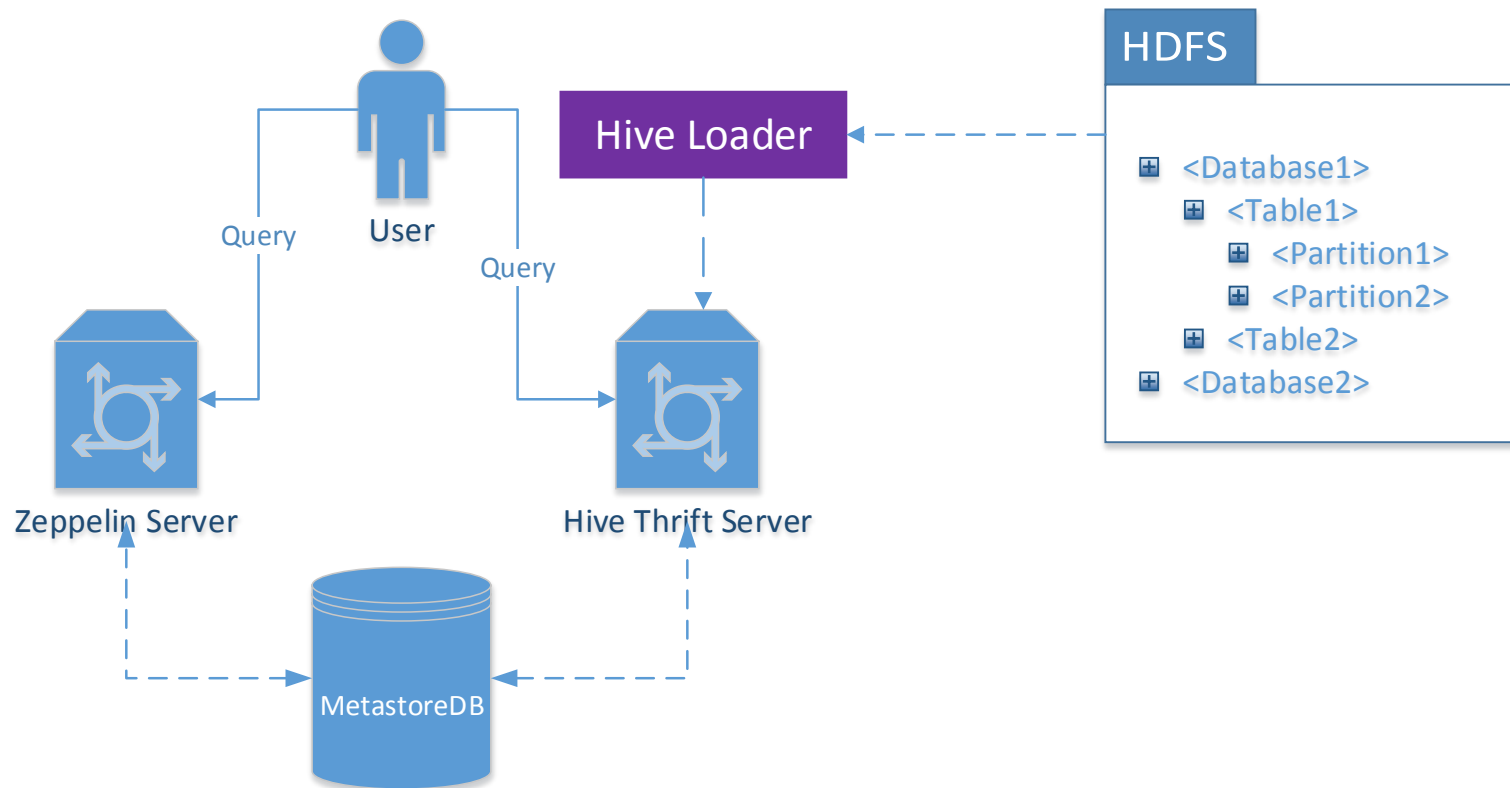
Architecture



Cosmos -> HDFS

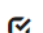



Hive Loader



User Experience in Avocado Task

STEPS spark transformers demo



● Source path SCOPE Source path  

Destination ... SCOPE

ADVANCED

```
1 // Name of the stream/streamset used as input
2 // {{run_date}} can be used in the streamset, plus_seconds/minus_seconds can be appended to adjust it
3 // e.g. /stream_path/%Y/%m/%d/data_%Y_%m_%d_%h.ss?run_date={{run_date | minus_seconds: 86400}}
4 /bing/searchlog/%Y/%m/searchlog_%Y_%m_%d.ss?run_date={{run_date}}
```

STEPS spark transformers demo

Source path SCOPE Destination path  

Destination ... SCOPE

ADVANCED

```
1 //HDFS path to target parquet file
2 //Format: hdfs://<HDFSServer>/<TeamName>/<TableName>/parquet/ssdate=<Date>
3 hdfs://hdfs1/bing/searchLog/parquet/querydate={{run_date}}
```

User Experience in Avocado Task (Cont.)

Schedules [+ Add](#)

Status	Schedule ID	Interval	Parameters	Run Date	Recent Executions	
<input checked="" type="checkbox"/> On <input type="checkbox"/> Off	Stale 33855	1 day		2015-05-07 00:00:00		<input type="button" value="▶ Backfill"/> <input type="button" value="🗑"/>

Backfill Schedules

Status	Schedule ID	Start Date	End Date	Parameters	Run Date	Recent Executions	
<input checked="" type="checkbox"/> On <input type="checkbox"/> Off	Running... 35010	2015-04-27 00:00:00	2015-04-27 00:00:00		2015-04-27 00:00:00		<input type="button" value="🗑"/>
<input checked="" type="checkbox"/> On <input type="checkbox"/> Off	Running... 34859	2015-03-01 00:00:00	2015-03-31 00:00:00		2015-03-21 00:00:00		<input type="button" value="🗑"/>
<input checked="" type="checkbox"/> On <input type="checkbox"/> Off	Running... 34860	2015-02-06 00:00:00	2015-02-28 00:00:00		2015-02-27 00:00:00		<input type="button" value="🗑"/>
<input type="checkbox"/> On <input checked="" type="checkbox"/> Off	OK 34895	2015-04-01 00:00:00	2015-04-01 00:00:00		2015-04-02 00:00:00		<input type="button" value="🗑"/>

User Experience in Avocado Task (Cont.)

More Info

Jobid	20150507-202344-Avocado Cosmos to Spark Service-2a94c979-3234-4f77-be24-590acf0eb9b9
Status	RUNNING
Driverurl	http://sparkjob/proxy/?serverName=CO3828:4040
Stdout	http://sparkjob/logs/stdout-20150507-202344-Avocado Cosmos to Spark Service-2a94c979-3234-4f77-be24-590acf0eb9b9
Stderr	http://sparkjob/logs/stderr-20150507-202344-Avocado Cosmos to Spark Service-2a94c979-3234-4f77-be24-590acf0eb9b9

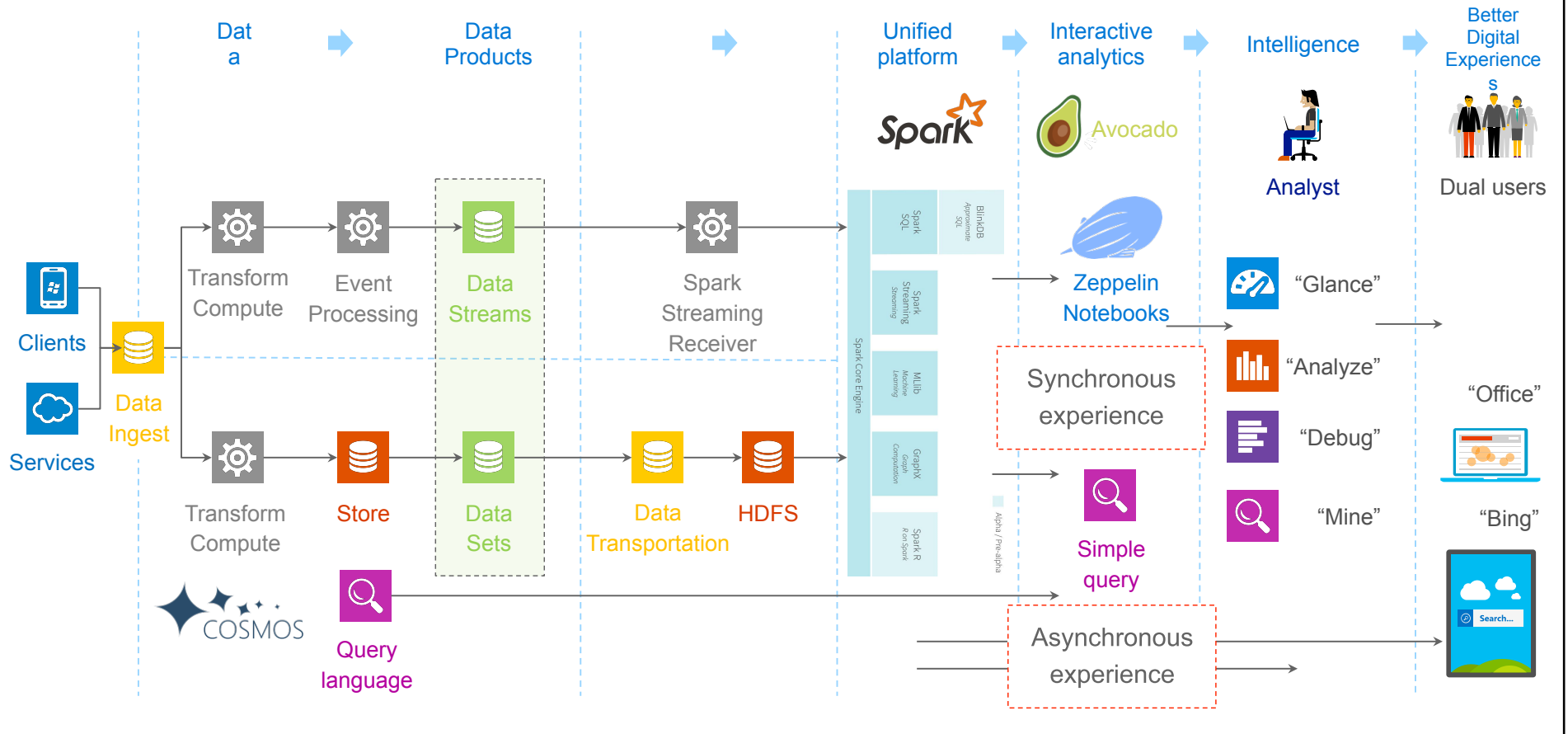
User Experience in Avocado Task (Cont.)

More Info

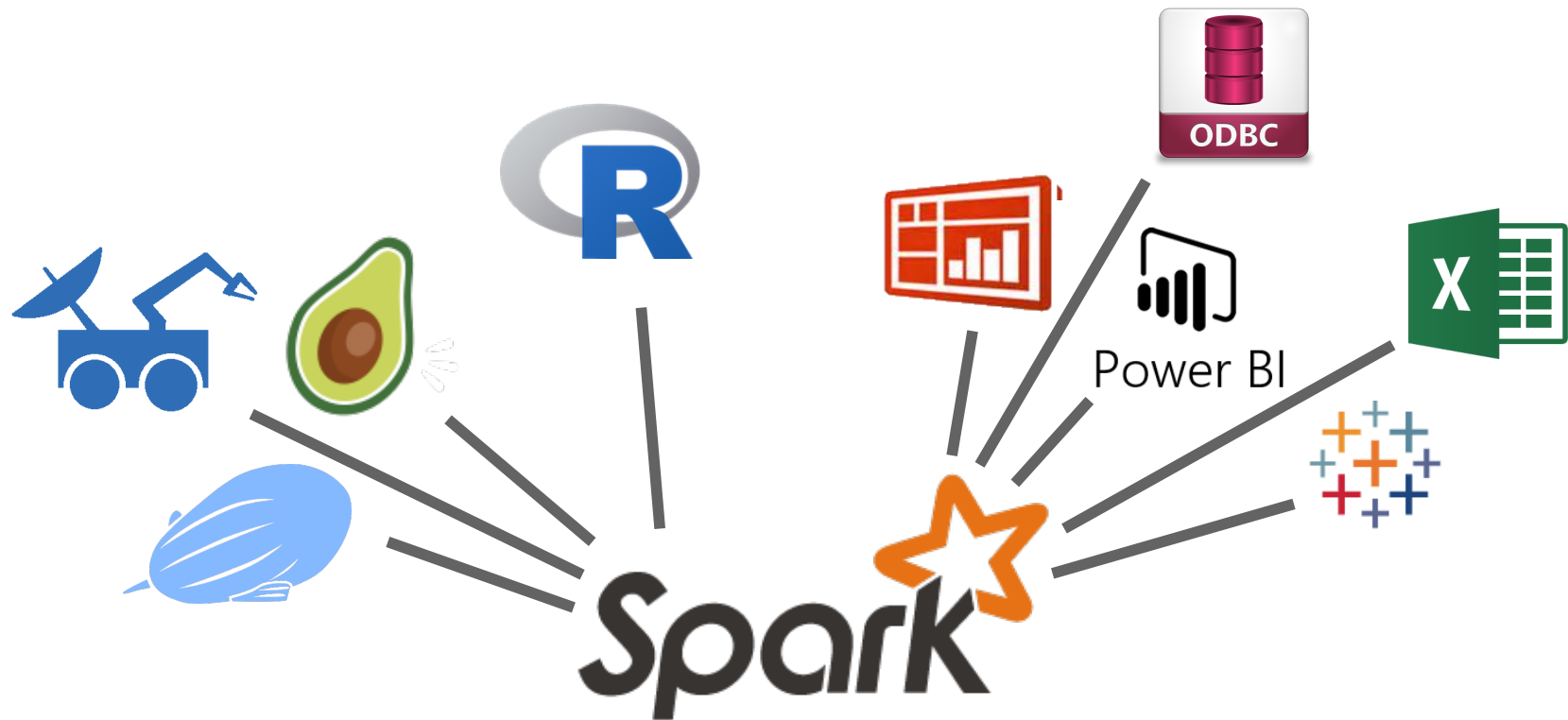
Jobid	20150507-184341-Avocado Cosmos to Spark Service-d170a297-e5e6-4bea-9972-888a71b167d0
Status	COMPLETED
Driverurl	Driver Url is unavailable after job finished, please check Jobhistoryurl
Stdout	http://sparkjob/logs/stdout-20150507-184341-Avocado Cosmos to Spark Service-d170a297-e5e6-4bea-9972-888a71b167d0
Stderr	http://sparkjob/logs/stderr-20150507-184341-Avocado Cosmos to Spark Service-d170a297-e5e6-4bea-9972-888a71b167d0
Jobhistoryurl	http://sparkjob/history/Framework-20150507_114940_1201_70221-52

What next?

Analytics as a Service in ASG



The #DataAnalystHub



Questions?

 @JulienJTPierre

 <https://linkedin.com/in/julienjtpierre>

