

CLICKSTREAM ANALYSIS WITH SPARK – UNDERSTANDING VISITORS IN REAL-TIME

Dr. Josef Adersberger
QAware GmbH, Germany



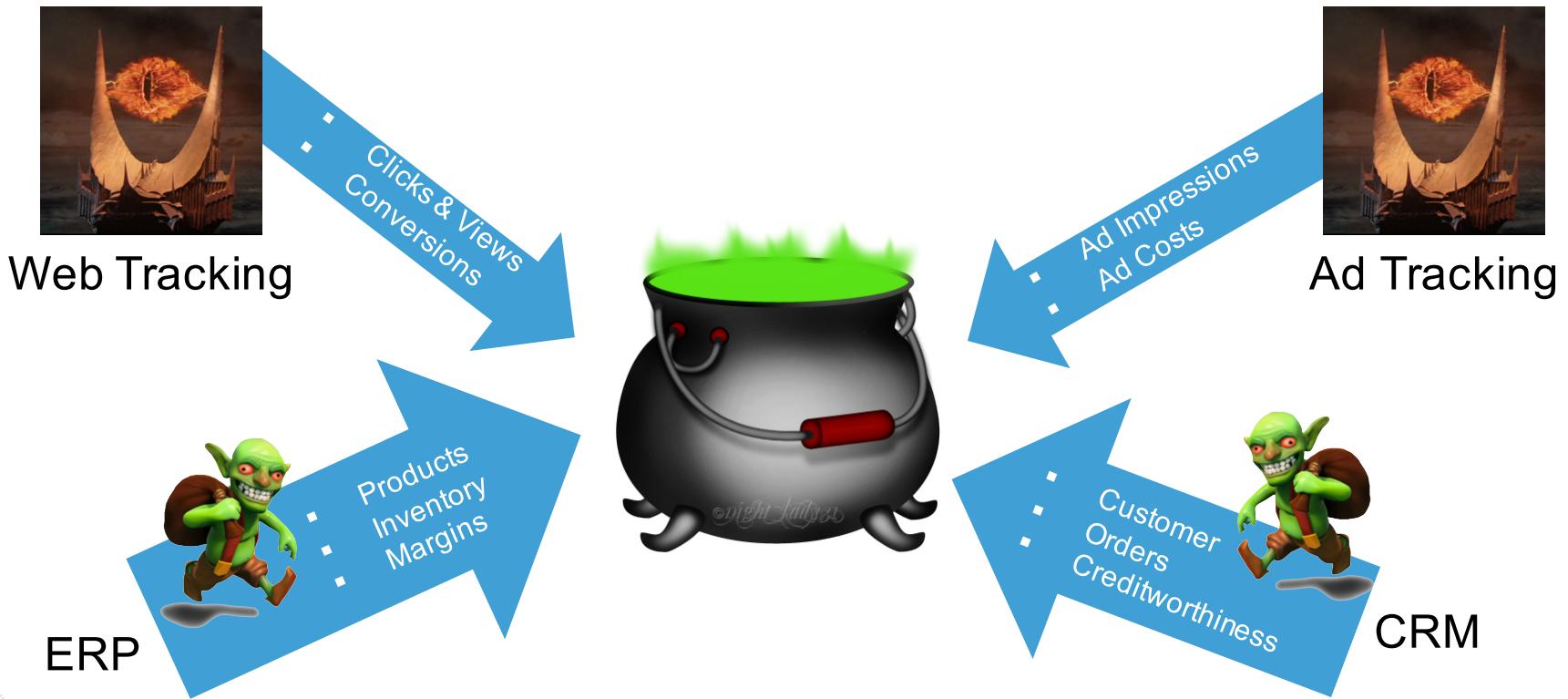
SPARK SUMMIT EAST
DATA SCIENCE AND ENGINEERING AT SCALE
FEBRUARY 16-18, 2016 NEW YORK CITY

THE CHALLENGE

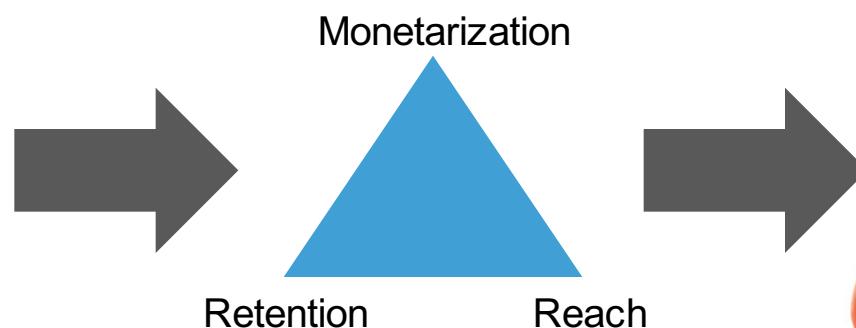


SPARK SUMMIT EAST
2016

One Kettle to Rule 'em All



One Kettle to Rule 'em All



- steer ...
- Campaigns
 - Offers
 - Contents



SPARK SUMMIT EAST
2016

... in real-time



SPARK SUMMIT EAST
2016

THE CONCEPTS

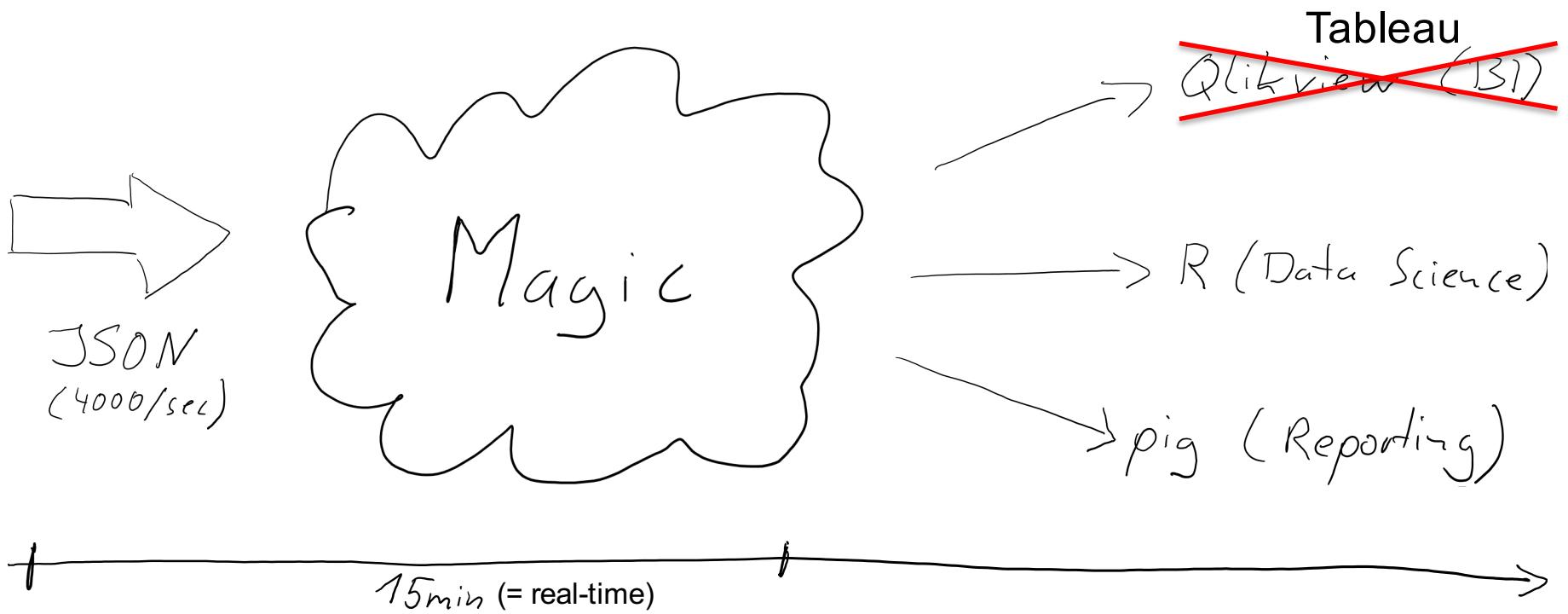


by Randy Paulino



SPARK SUMMIT EAST
2016

The First Sketch





Magic

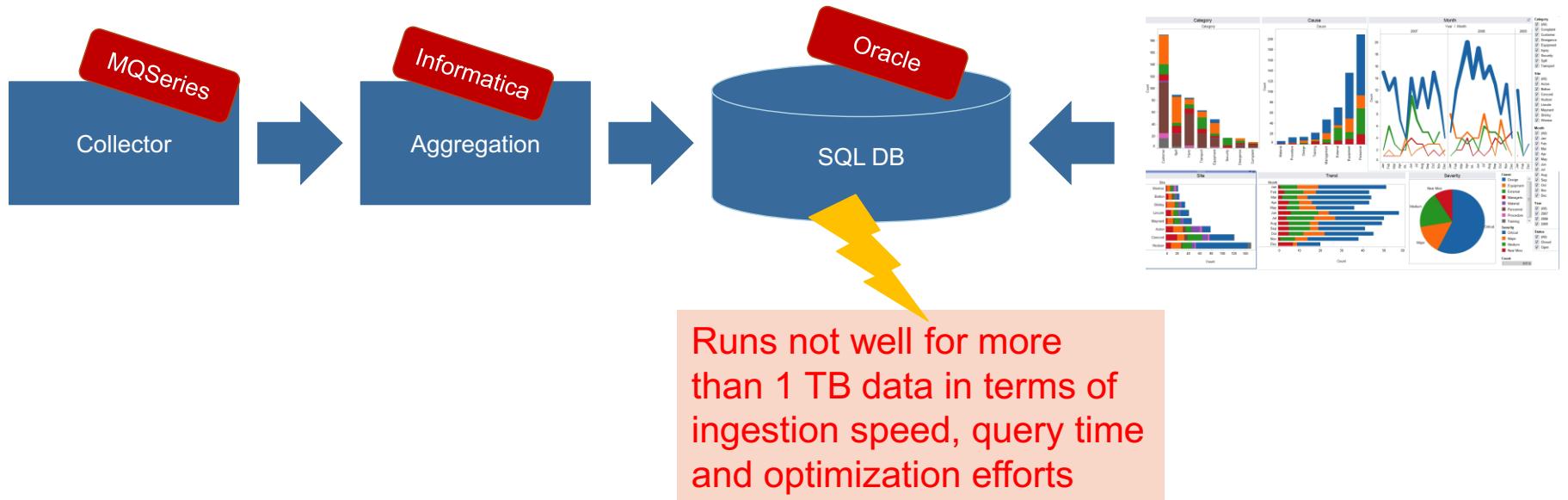


THE ARCHITECTURE



SPARK SUMMIT EAST
2016

„Larry & Friends“ Architecture





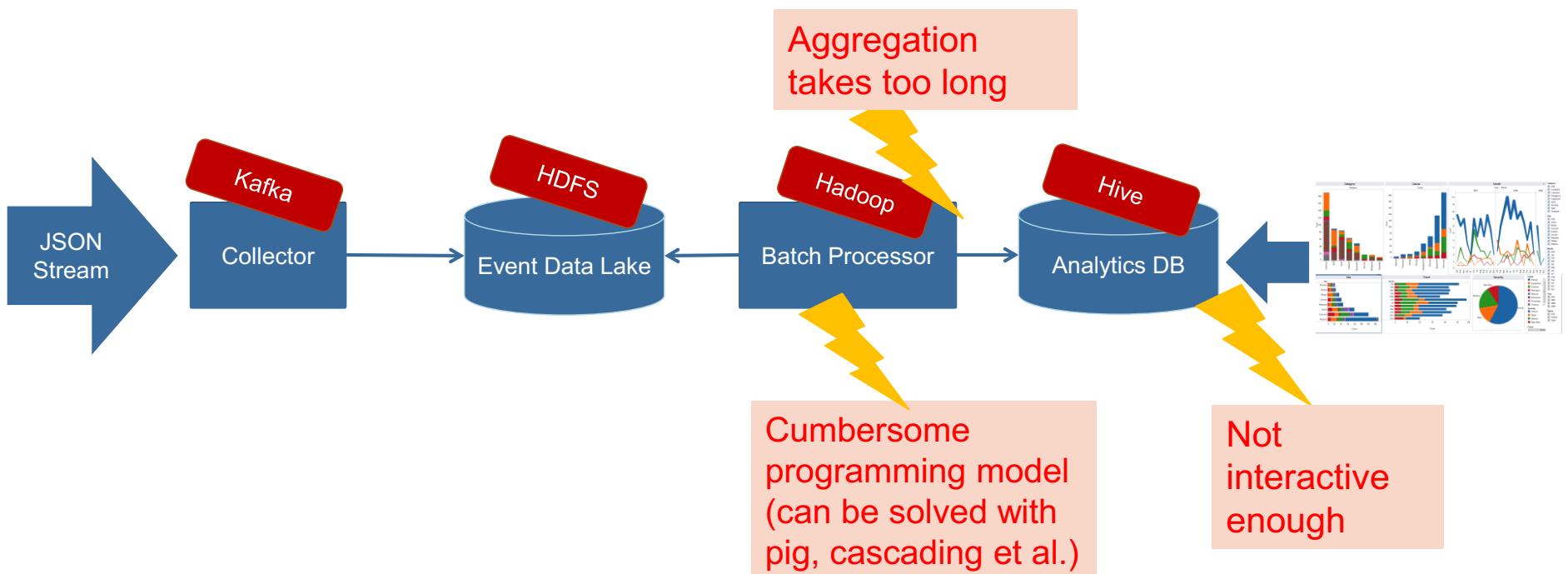
Nope.
Sorry, no Big Data.

by adweek.com



SPARK SUMMIT EAST
2016

„Hadoop & Friends“ Architecture



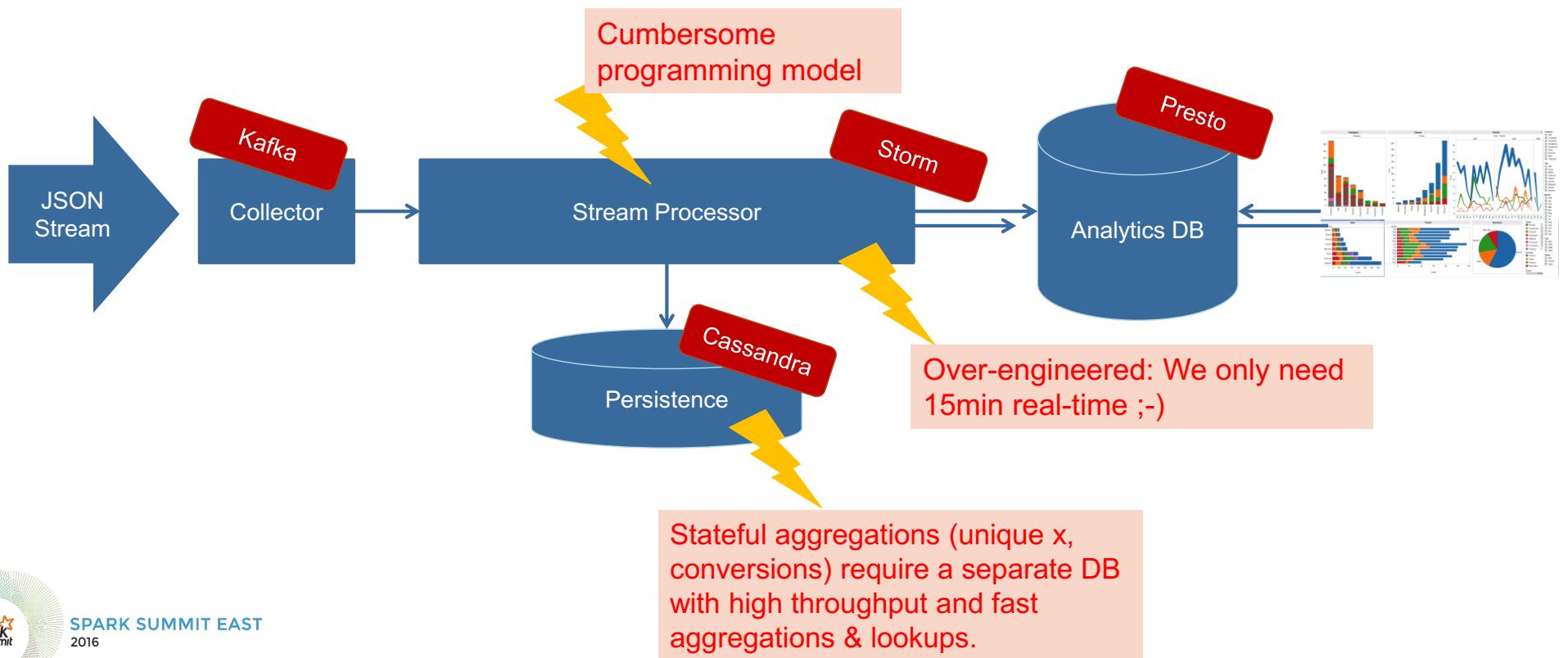


Nope.
Too sluggish.

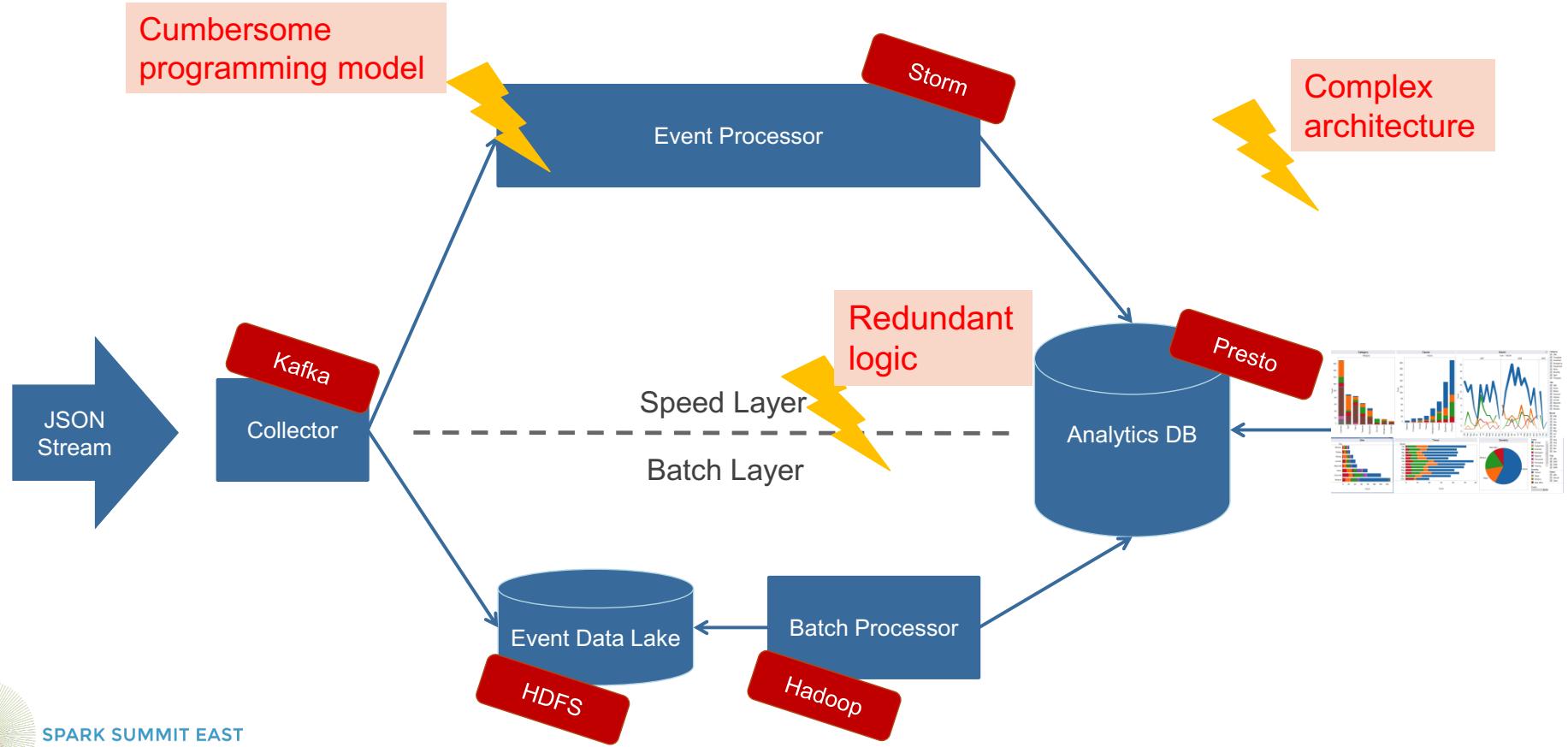


SPARK SUMMIT EAST
2016

κ -Architecture

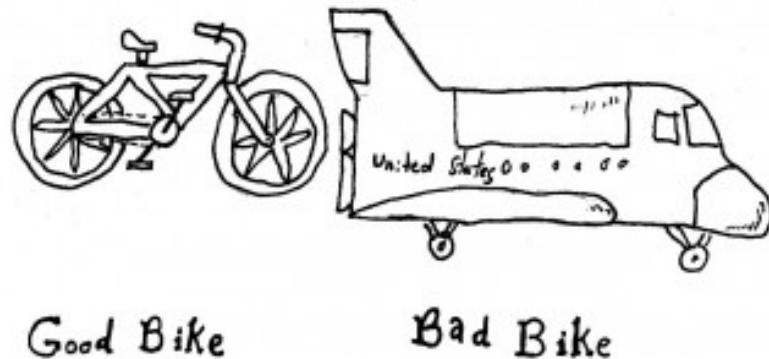


λ -Architecture



SPARK SUMMIT EAST
2016

Feels Over-Engineered...



The background image shows the interior of the Sagrada Família cathedral in Barcelona, featuring its iconic organic, column-like structures and intricate ceiling designs. A blue color filter is applied across the entire image.

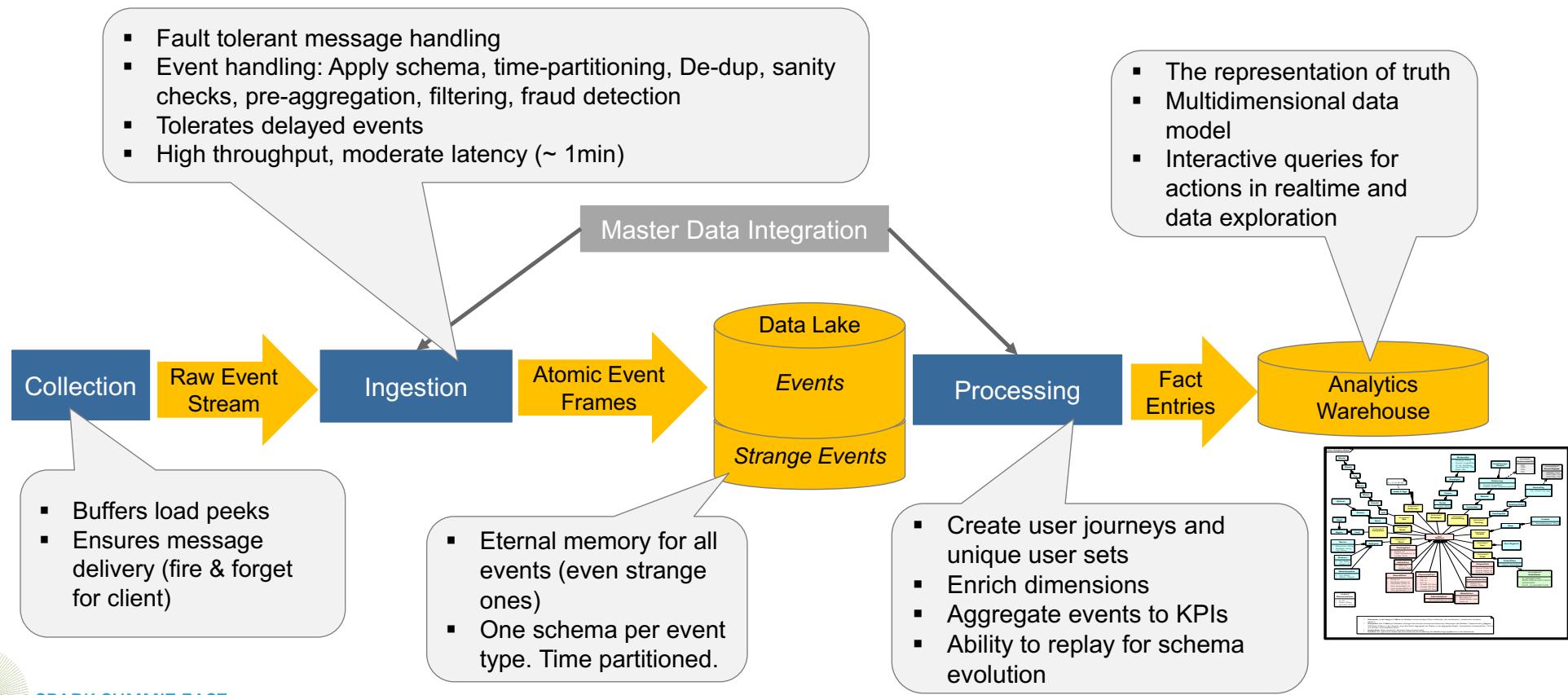
The Final Architecture*

*) Maybe called μ -architecture one day ;-)

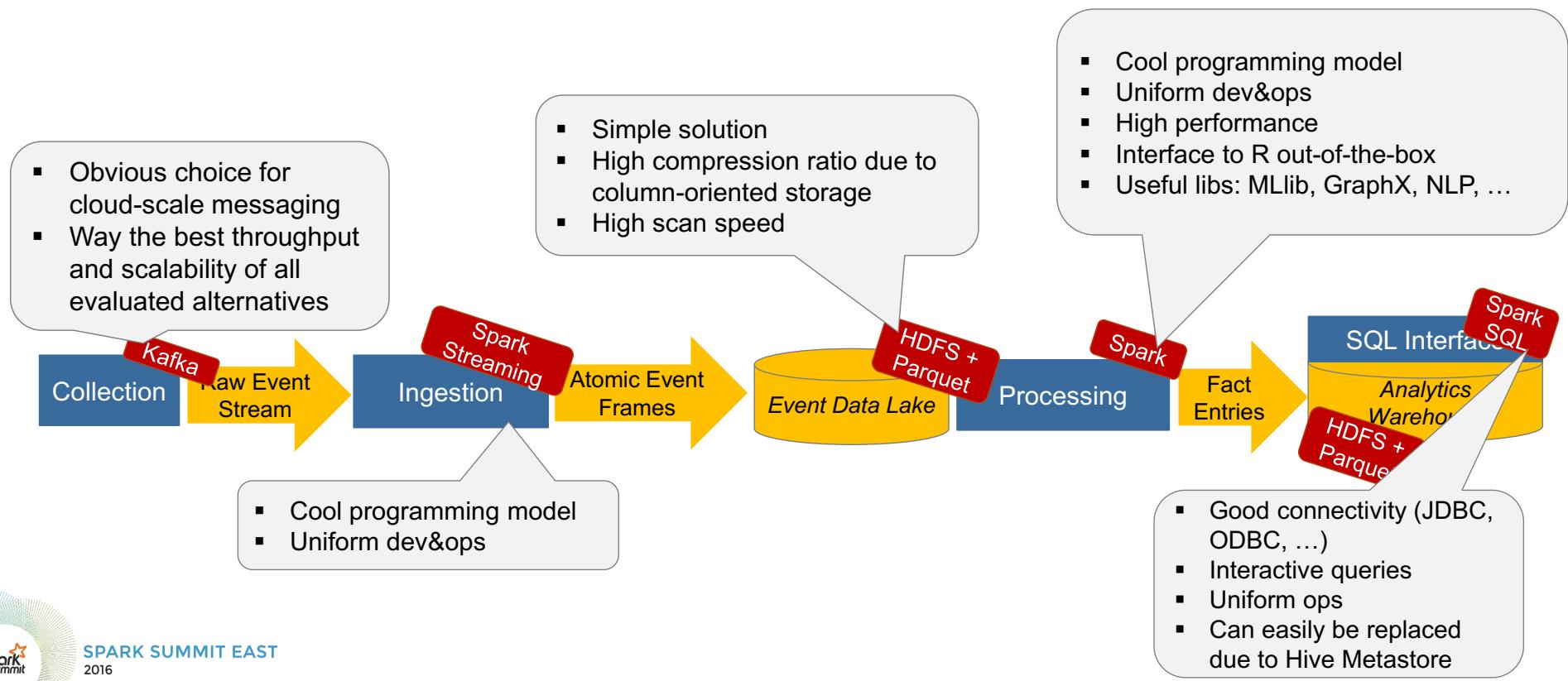


SPARK SUMMIT EAST
2016

Functional Architecture



Series Connection of Streaming and Batching - all based on Spark.



LESSONS LEARNED

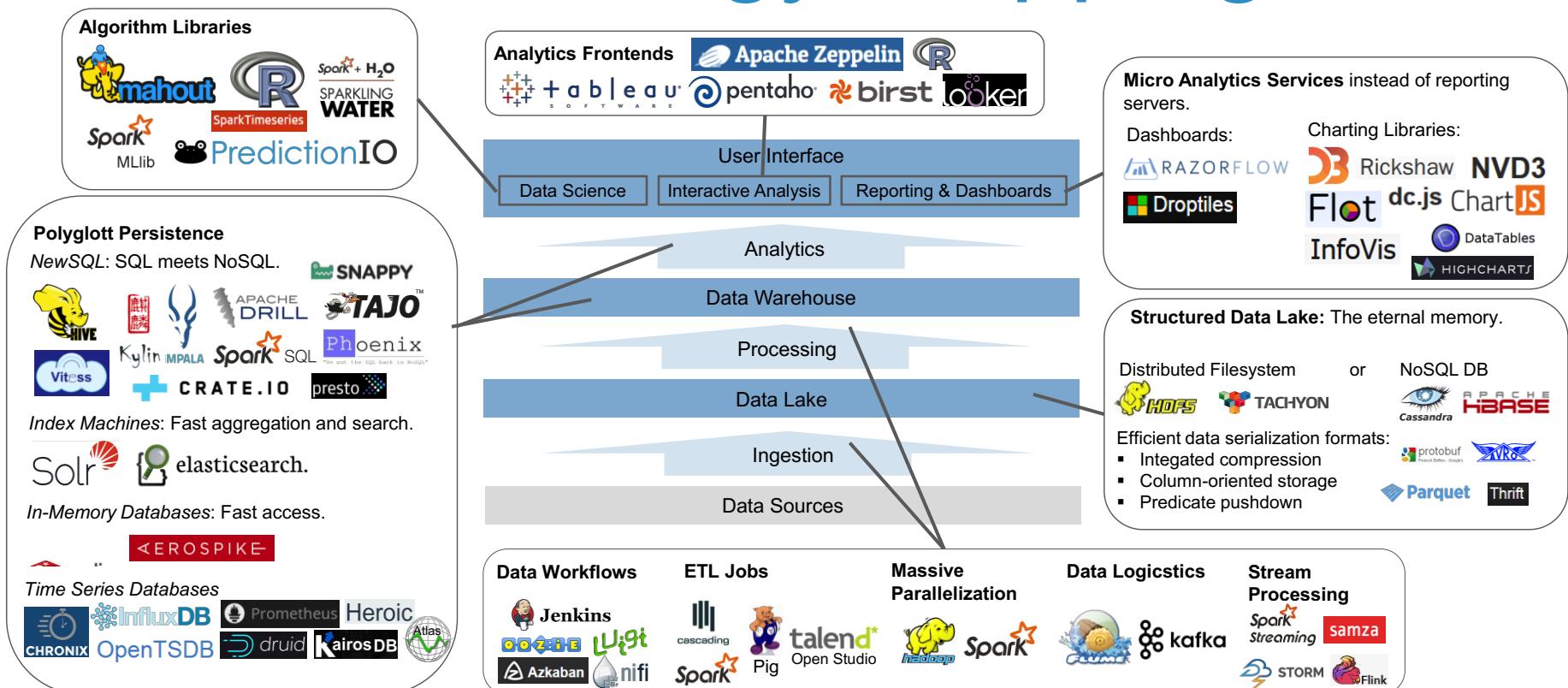


by: <http://hochmeister-alpin.at>



SPARK SUMMIT EAST
2016

Technology Mapping

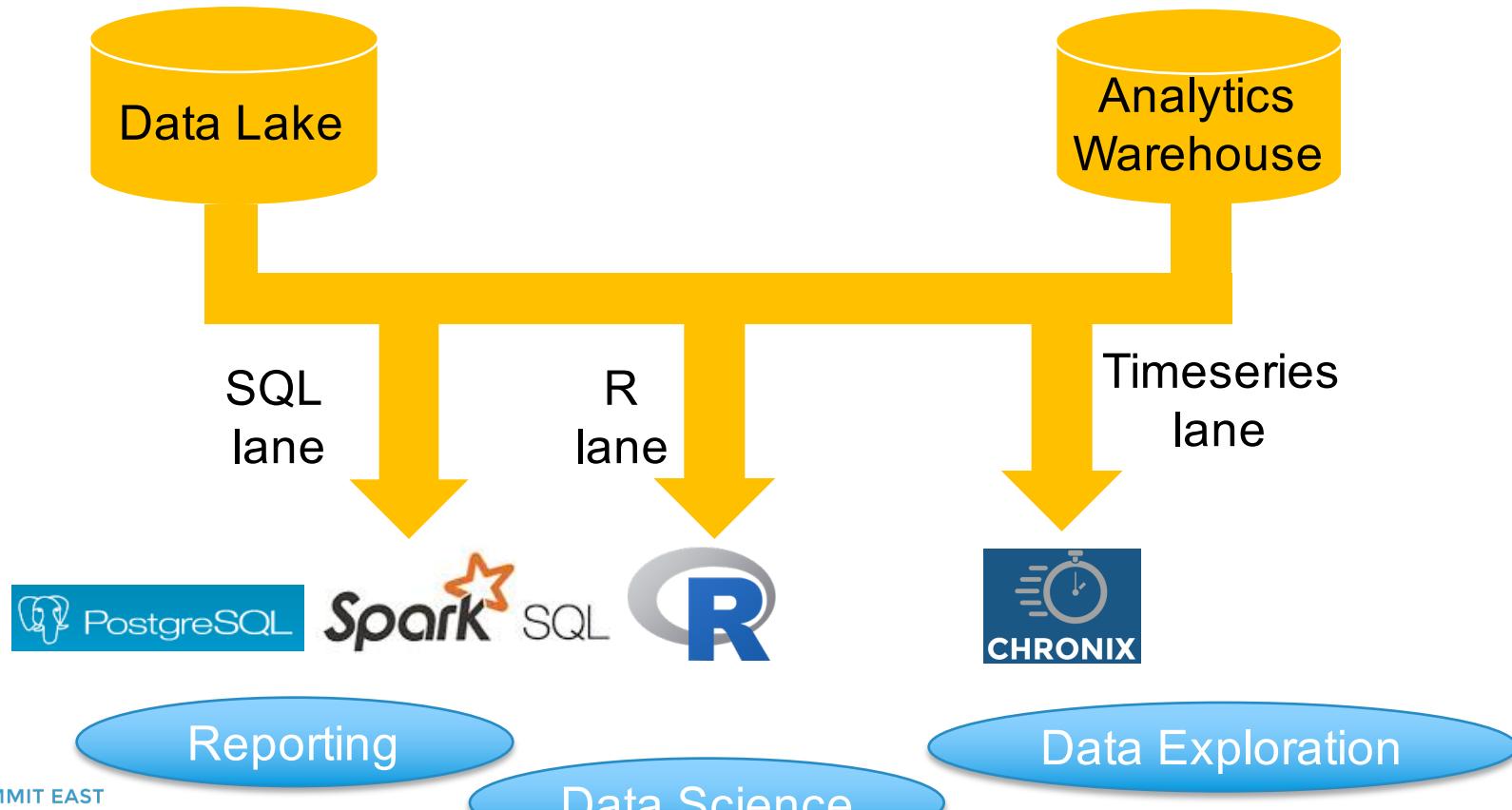


<https://github.com/qaware/big-data-landscape>

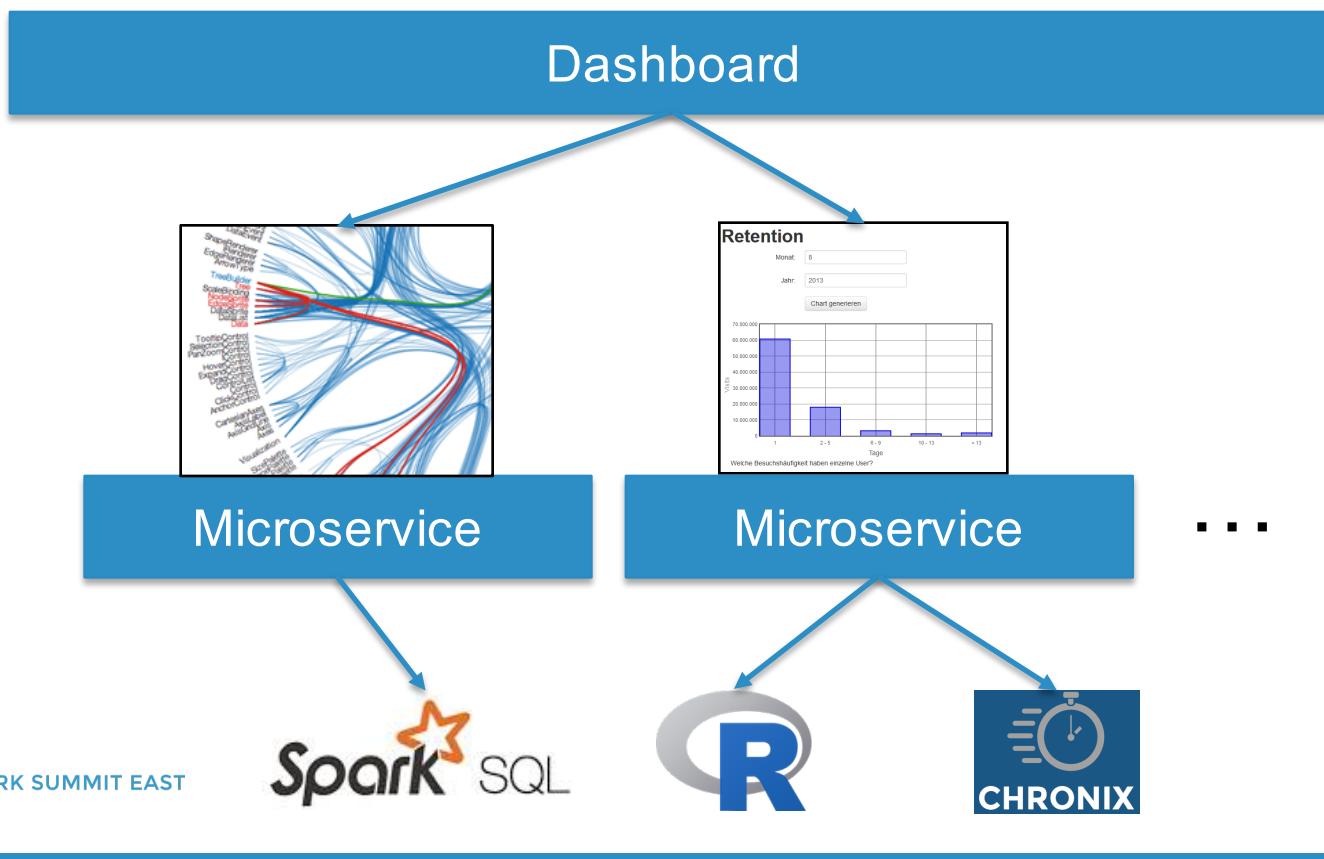


SPARK SUMMIT EAST
2016

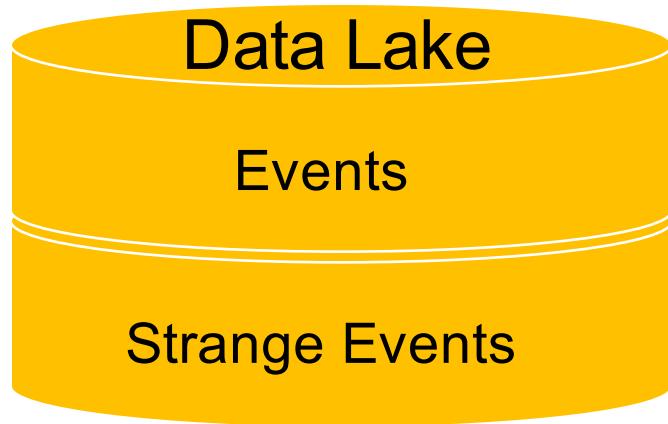
Polyglott Analytics



Micro Analytics Services



No Retention Paranoia



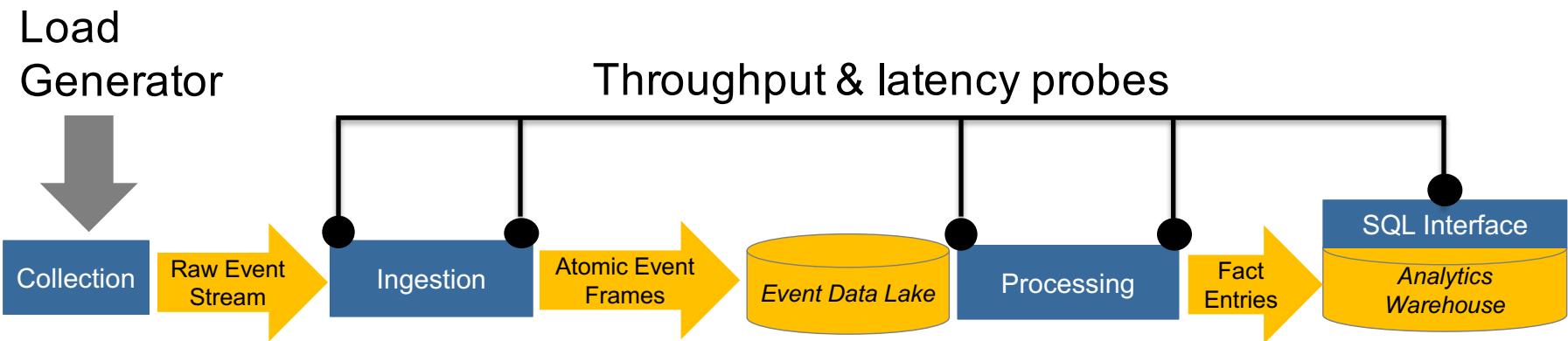
- Eternal memory
- Close to raw events
- Allows replays and refills into warehouse



- Aggressive forgetting with clearly defined retention policy per aggregation level like:
- 15min:30d
 - 1h:4m
 - ...



Continuous Tuning



In Numbers

Overall dev effort until the first release: **250 person days**

Dimensions: **10**

KPIs: **26**

Integrated 3rd party systems: **7**

Inbound data volume per day: **80GB**

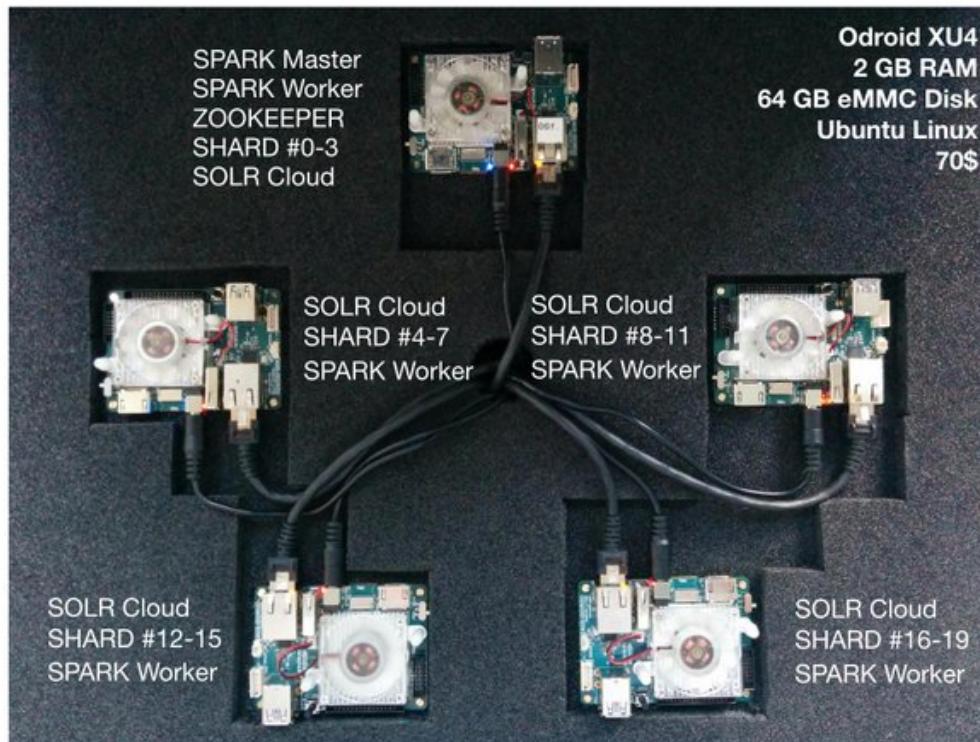
New data in DWH per day: **2GB**

Total price of cheapest cluster which is able to handle production load:



SPARK SUMMIT EAST
2016

Your own Datacenter (less than 500\$)



40 Cores
10 GB RAM
320 GB total disk

ODROID
Hardkernel
Odroid XU4
2 GB RAM
64 GB eMMC Disk
Ubuntu Linux
70\$



THANK YOU.



@adersberger



josef.adersberger@qaware.de

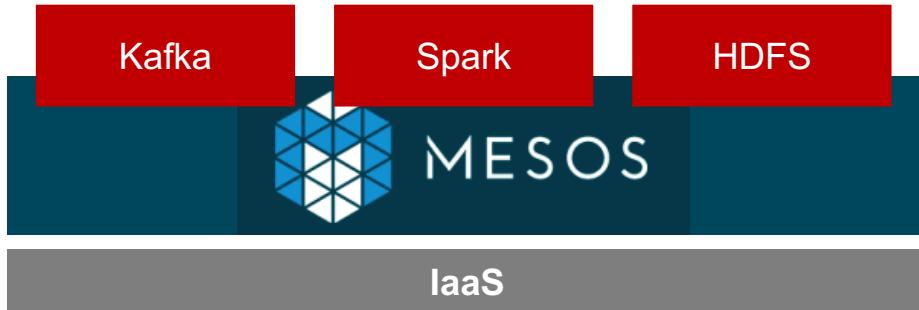


Download slides



SPARK SUMMIT EAST
DATA SCIENCE AND ENGINEERING AT SCALE
FEBRUARY 16-18, 2016 NEW YORK CITY

Bonus Topic: Roadmap



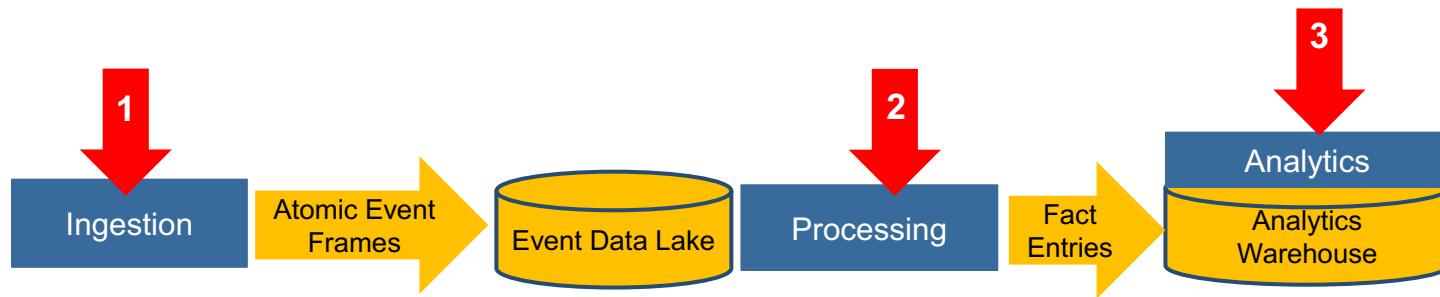
Simplify Ops with Mesos



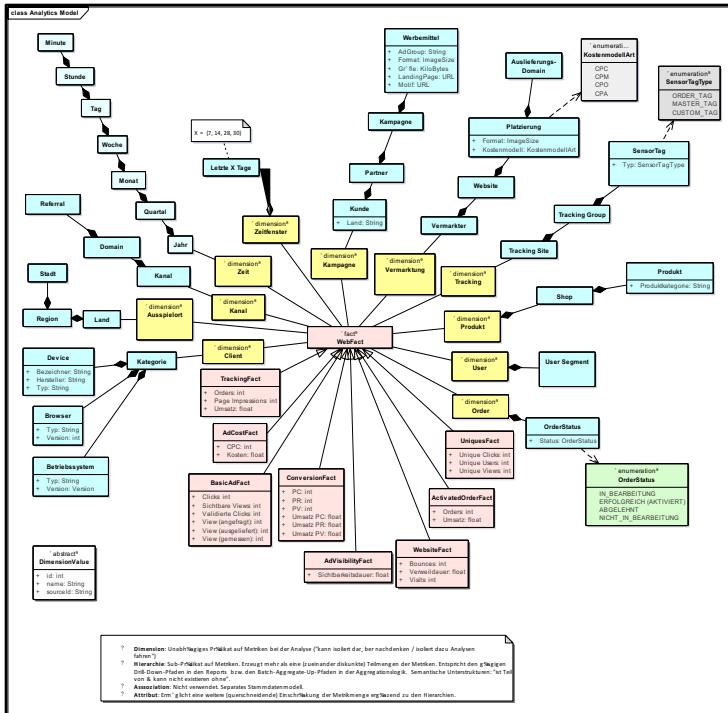
Faster aggregation & easier updates
with Spark-on-Solr

<http://qaware.blogspot.de/2015/06/solr-with-sparks-or-how-to-submit-spark.html>

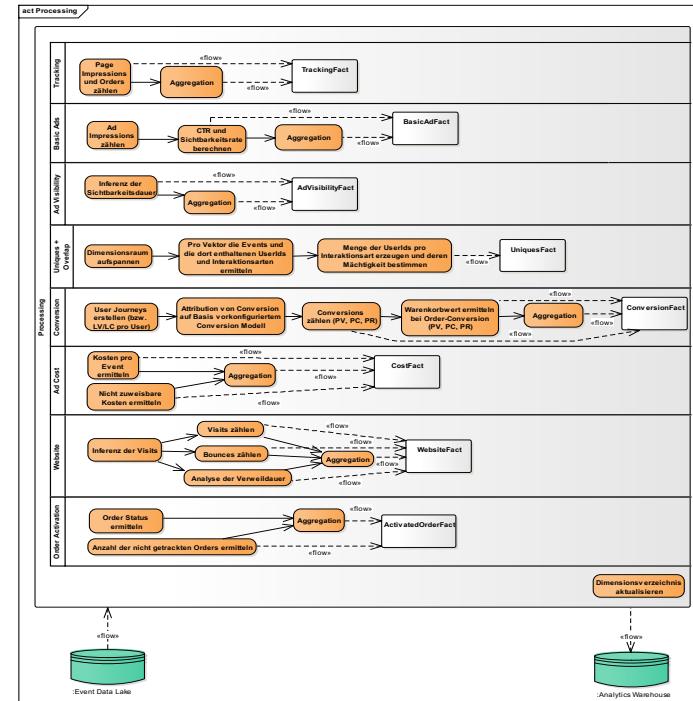
Bonus Topic: Smart Aggregation



Architecture follows requirements



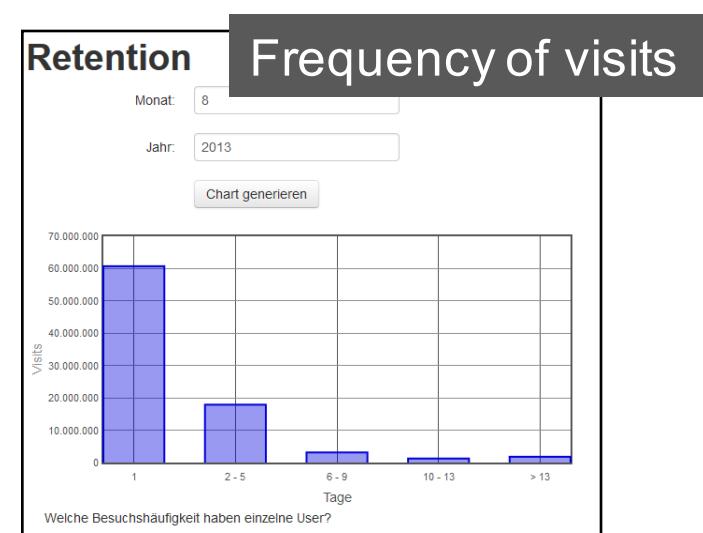
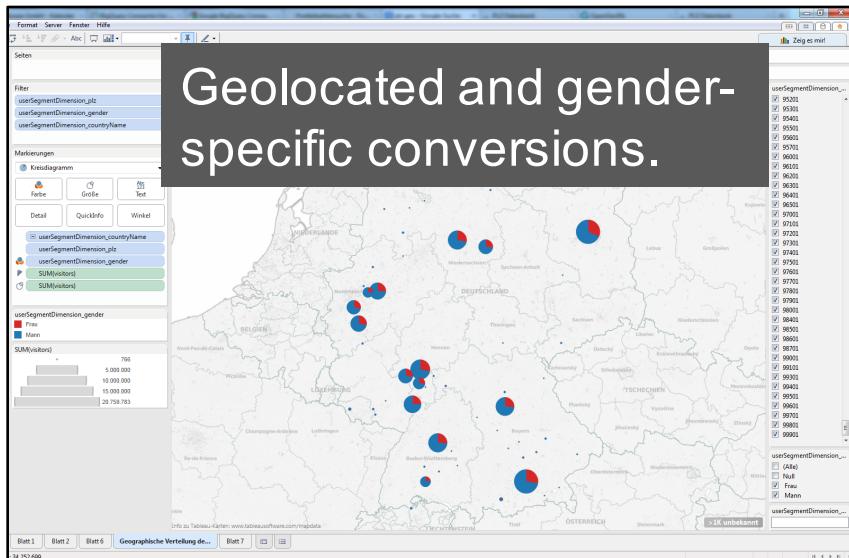
Multidimensional Data Model



Data Processing Workflow



Sample results



SPARK SUMMIT EAST
2016