


Big Data Engineering

Conclusions and Recap

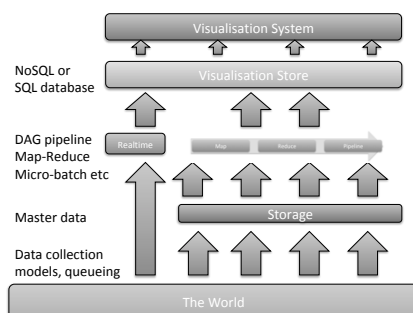
 © Paul Fremont 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <https://creativecommons.org/licenses/by-nc-sa/4.0/>

Contents

- Understanding the bigger picture
- What are the different components
- Message queueing and collection systems
- Map-Reduce and DAG systems
- Realtime Systems
- Fast databases for speed
- Visualisation and Dashboards

 © Paul Fremont 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <https://creativecommons.org/licenses/by-nc-sa/4.0/>

The big picture



 © Paul Fremont 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <https://creativecommons.org/licenses/by-nc-sa/4.0/>

The big picture

- You have *immutable* master data
- You create a set of processes to:
 - Collect that data
 - Store master data
 - Process data
 - Visualise and present
- Some of those processes act on batch and others on real-time data

 © Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <https://creativecommons.org/licenses/by-nc-sa/4.0/>

How to choose the components?

- Two main approaches:
 - Best of breed
 - Choose the best available component in each space
 - Stack
 - Choose a curated stack that a team or organization is providing/selling/supporting

 © Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <https://creativecommons.org/licenses/by-nc-sa/4.0/>

Approach

- Minimise the pain
 - Choose what you need when you need it
 - Don't over engineer

 © Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <https://creativecommons.org/licenses/by-nc-sa/4.0/>

How do I ingest data?

- File transfer
- Live stream
 - Sockets
 - Syslog
 - Messaging system
- From existing databases

 © Paul Fremont 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <https://creativecommons.org/licenses/by-nc-sa/4.0/>

How do I store data?

- HDFS
- NoSQL database only
 - Mongo / HBase / Cassandra
- zFS / GlusterFS / NFS etc
- Apache Parquet, CSV, or speci

 © Paul Fremont 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <https://creativecommons.org/licenses/by-nc-sa/4.0/>


How do I process data?

- Simple Map Reduce
- Hive / Pig
- DAG
- Pipeline
- etc

 © Paul Fremont 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <https://creativecommons.org/licenses/by-nc-sa/4.0/>

How do I visualise data

- From a SQL database?
- From a NoSQL database?
- Generate charts in Python Spark?
- Etc?

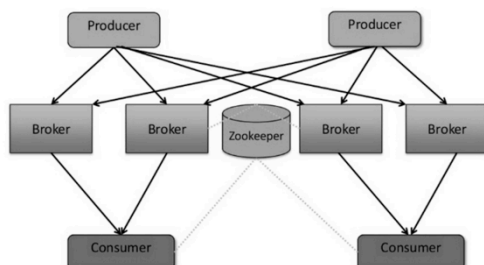
 © Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Collection / Queuing systems

- Two ways of making the choice
 - The protocol
 - The middleware
- Protocols
 - ZeroMQ, MQTT, AMQP, STOMP, Kafka Protocol, Rendezvous, etc
- Middleware
 - Kafka, Apollo, Mosquitto, QPid, WSO2, etc

 © Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

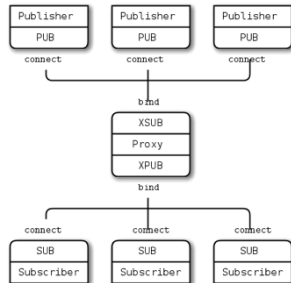
Apache Kafka




 © Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <http://creativecommons.org/licenses/by-nc-sa/4.0/>

Source: <http://www.slideshare.net/charmallo/>

ZeroMQ



 © Paul Fremont 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <https://creativecommons.org/licenses/by-nc-sa/4.0/>

Processing approaches

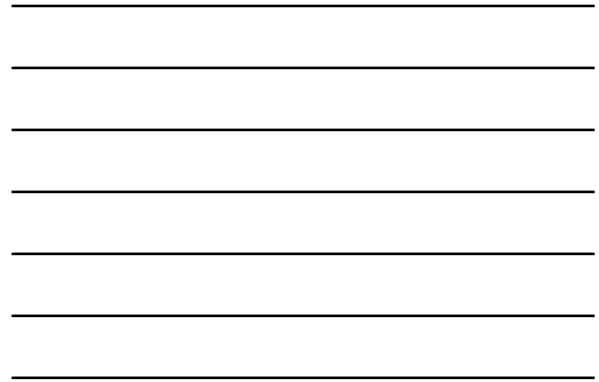
- Covered in detail already
- Hadoop
- Spark
- Tez
- etc

 © Paul Fremont 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <https://creativecommons.org/licenses/by-nc-sa/4.0/>

Cluster Management

- Spark
- YARN
- Mesos
- Kubernetes
- etc

 © Paul Fremont 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <https://creativecommons.org/licenses/by-nc-sa/4.0/>



Visualisation approaches

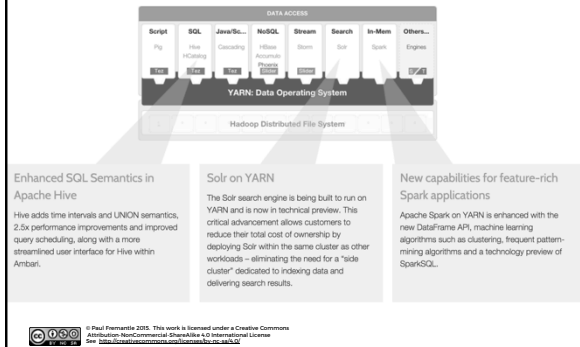


Fortune top 10 big data companies

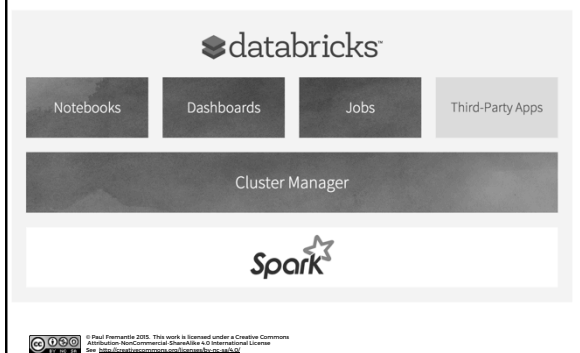
fortune.com/2014/06/13/these-big-data-companies-are-ones-to-watch/



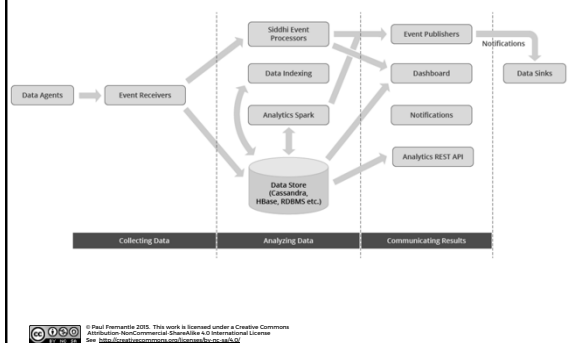
Hortonworks

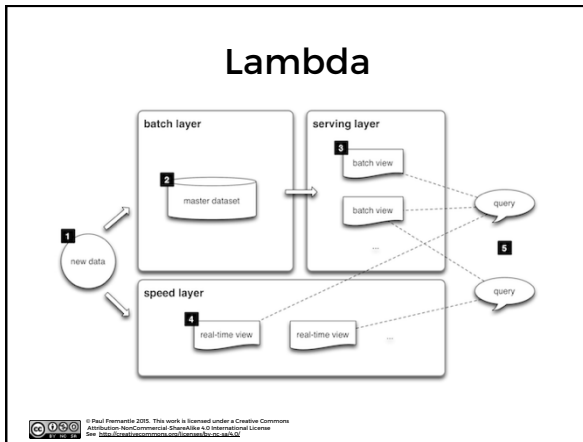


Databricks



WSO2 DAS





The real answer

You are on the bleeding edge
–Expect to have some pain

© Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <https://creativecommons.org/licenses/by-nc-sa/4.0/>

Questions?

© Paul Fremantle 2015. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. See <https://creativecommons.org/licenses/by-nc-sa/4.0/>
