



BIG DATA WEEK
A GLOBAL FESTIVAL OF DATA





Data Lake and the rise of the Microservices



BIG DATA WEEK
A GLOBAL FESTIVAL OF DATA



About Me

- Developer (forcibly) turned Product Manager
- Been designing system architectures as well as products for the hosting market during the last 8 years.
- Passionated about distributed computing, HPC environments and supercomputers

 @alexandrubordei

 alex@bigstep.com

 @bigstepinc

bigstep



BIG DATA WEEK
A GLOBAL FESTIVAL OF DATA



Big Data technologies for mainstream and vice-versa

- Due to the cap on CPU frequency, **the horizontal** is the only dimension left to grow into.
- Client-server architecture outdated.
- All components of an application must be as independent as possible and as scalable as possible.
- Big data technologies increasingly used in general purpose applications
- In-memory technologies are orders of magnitude faster than the others.
- Docker promotes and simplifies large scale application management using low-overhead containers instead of VMs
- Mesos used with Docker and some additional services creates a Distributed OS



BIG DATA WEEK
A GLOBAL FESTIVAL OF DATA



Data as artefacts

- Just like archeological artefacts, old data can yield new insights if correlated in a novel way or analysed with a new technology.
- Throwing away data because it is of no use today might cripple the business tomorrow.



[Source: Tori Randall, Ph.D. prepares a 550-year old Peruvian child mummy for a CT scan](#)

The Data Lake

- Promote data exploration, innovation
- Enable automatic decision making, to Uberize the business
- Perform new, deeper analytics by focusing on correlations between diverse data sources: clickstream, social media, machine data, documents, audio/video, etc.
- Stores data in its original format
- Stores structured and unstructured data together



BIG DATA WEEK
A GLOBAL FESTIVAL OF DATA



A Data Lake \neq A giant universal database

The problem is domain knowledge.

Sample TFL dataset:

08/287, 09:30:05, JUN, 1, 1008, & G1, & F1, 08/287a1=0000, 08/287a2=0000, 08/287b1=0011, 08/287c1=0000, 08/287c2=0000, 00/000 =1111, 00/000 =1111, 00/000 =1111, G1, G1, 1, 2, 888, 104, 88, 1, 1, 1, 1, 1, 0, 0, 0, 1100000000000000

06/010, 09:30:05, JUN, 1, 785, & G2, & F2, 06/010t1=1111, 06/010w1=1111, 06/010n1=1111, 06/010p1=1111, 06/010e1=1111, 00/000 =0000, 06/010s1=0000, 06/010q1=0000, G1, G1, 1, 2, 200, 128, 128, 1, 1, 1, 1, 1, 0, 1, 1, 1100000000000000

09/011, 09:30:05, JUN, 1, 1165, & G3, & F1, 09/011s1=0000, 09/066m1=0000, 09/011p1=0000, 09/011r1=1111, 00/000 =0000, 09/210a1=1111, 00/000 =0000, 00/000 =0000, G1, G1, 1, 2, 565, 88, 104, 1, 1, 1, -2, 0, -2, 0, 0, 1100000000000000



BIG DATA WEEK
A GLOBAL FESTIVAL OF DATA



A Data Lake \neq A giant universal database

- New scientist needs:
 - Domain knowledge
 - Clean data
 - Easy to work with environment
 - Data documentation up to date
 - Domain knowledge
 - ETL job in place to get new data
- Knowledge is not just in the form of data it is also in the form of code
- HDFS = Just the storage medium

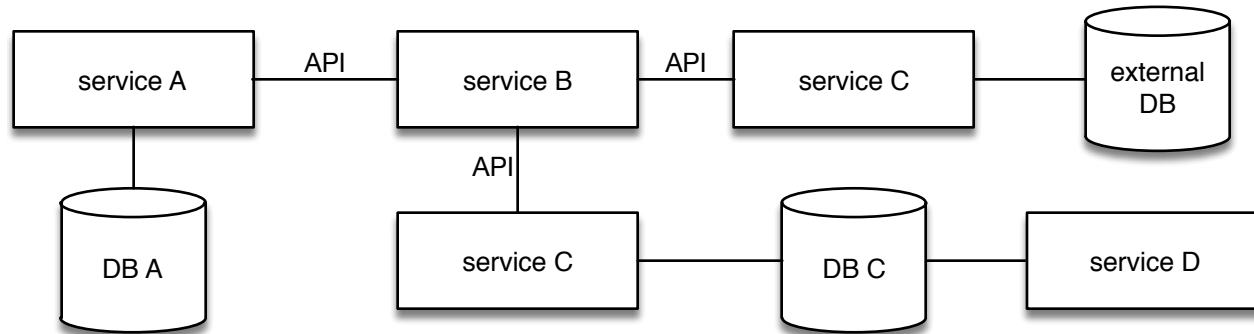


BIG DATA WEEK
A GLOBAL FESTIVAL OF DATA

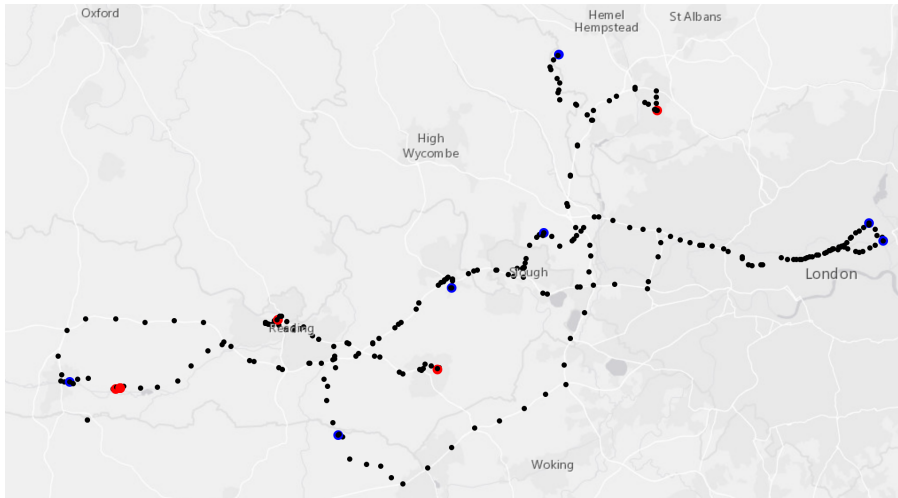


A Data Lake = A collection of Data Services

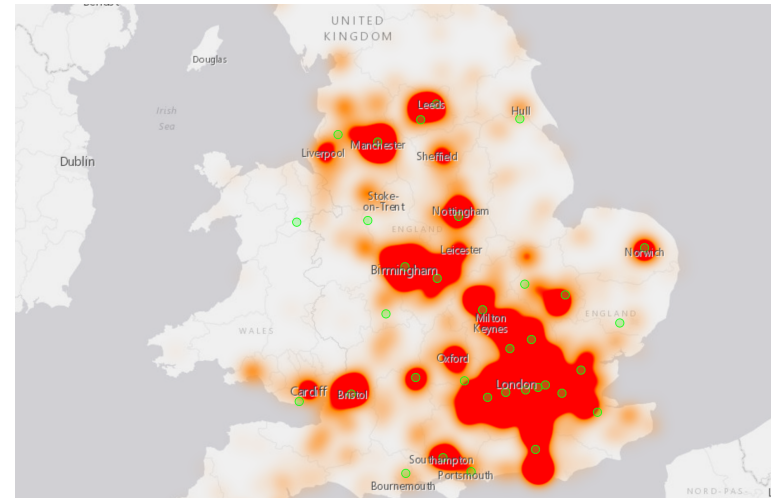
- A service provides an API to the rest of the services
- A service is development by a team that has enough domain knowledge to import and clean the data
- A service could simply store the result of it services periodically back into the shared DB
- Some services offer simply processing services and do not store data at all



An example - A big logistics company

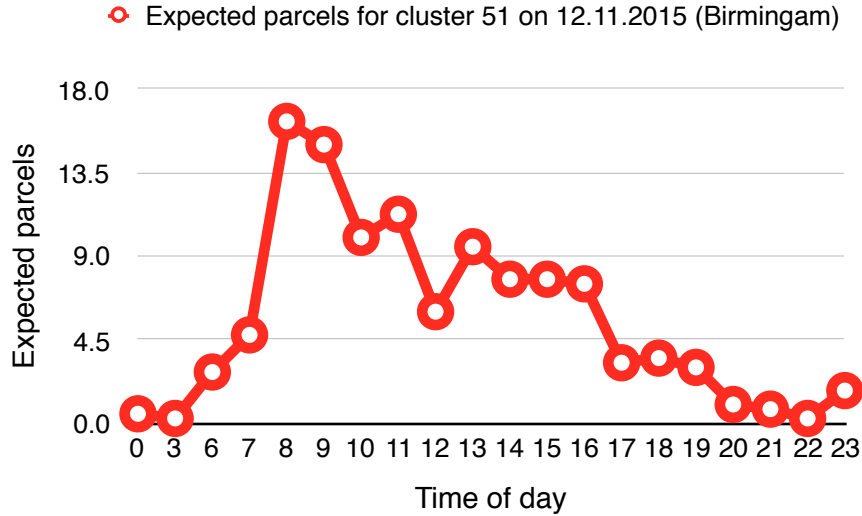


driver service

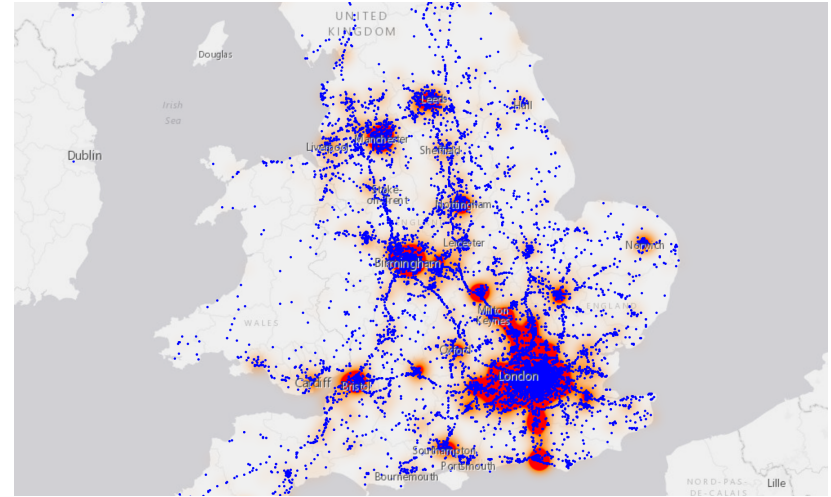


pickup and drop-off service

An example - A big logistics company



cluster service



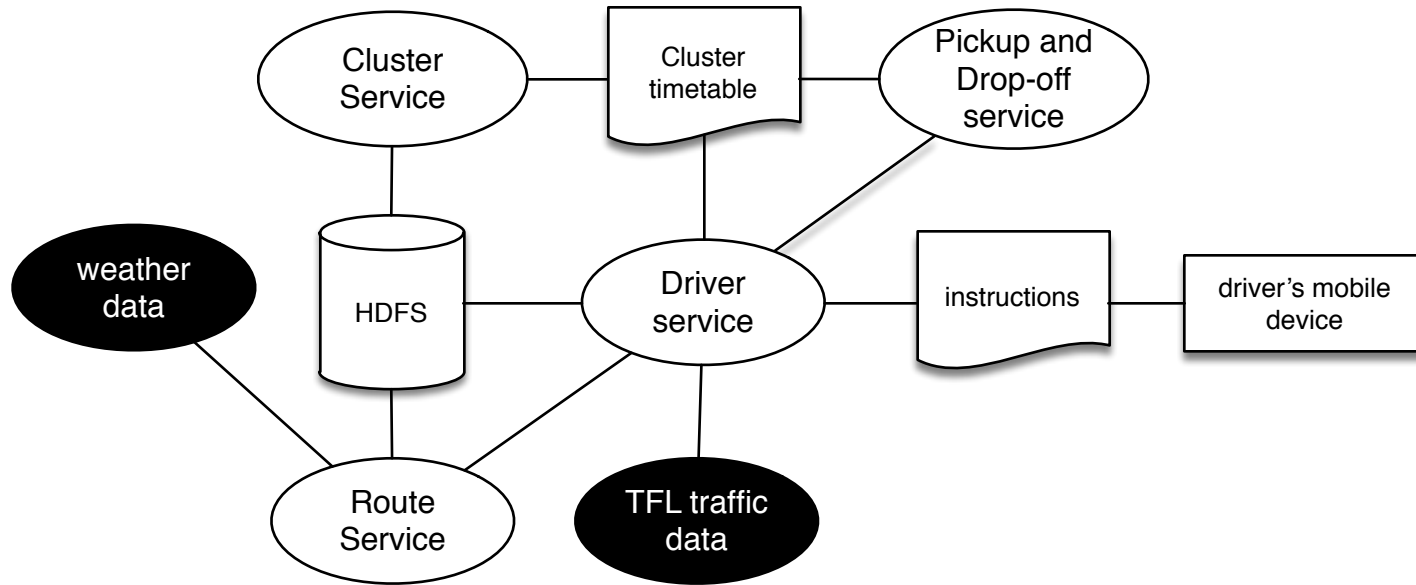
route service



BIG DATA WEEK
A GLOBAL FESTIVAL OF DATA

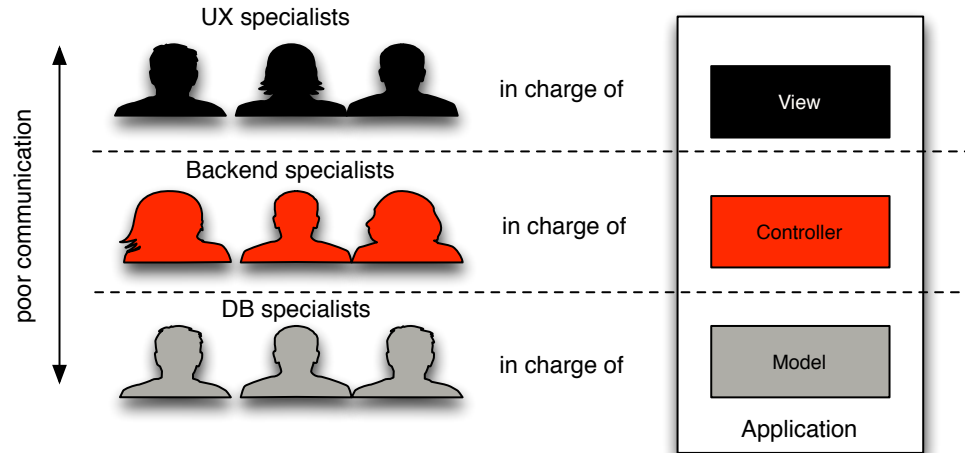


An example - A big logistics company



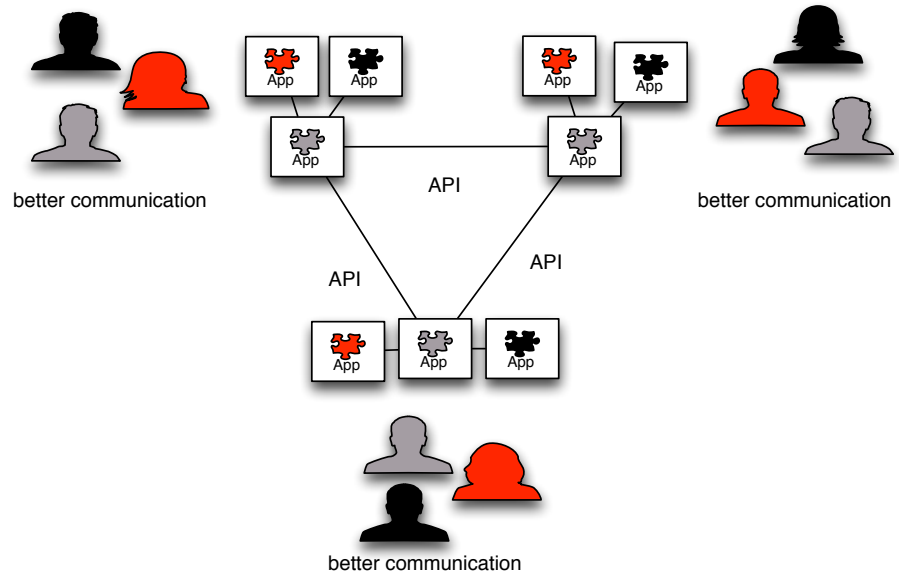
Data Services - It's about the teams and not the technology

- Conway law: “[...] organizations which design systems ... are constrained to produce designs which are copies of the communication structures of these organizations”



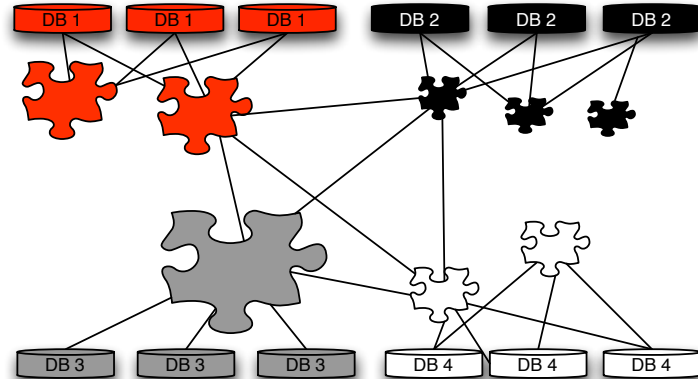
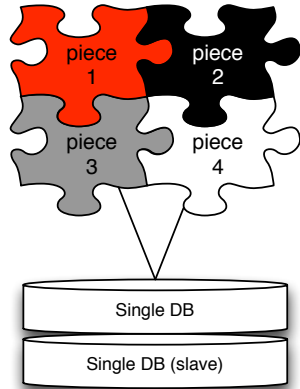
Data teams

- A data service has its own release cycle
- Ultra-specialisation is reduced
- Communication among members of the same team is better
- Each Service has a clear owner



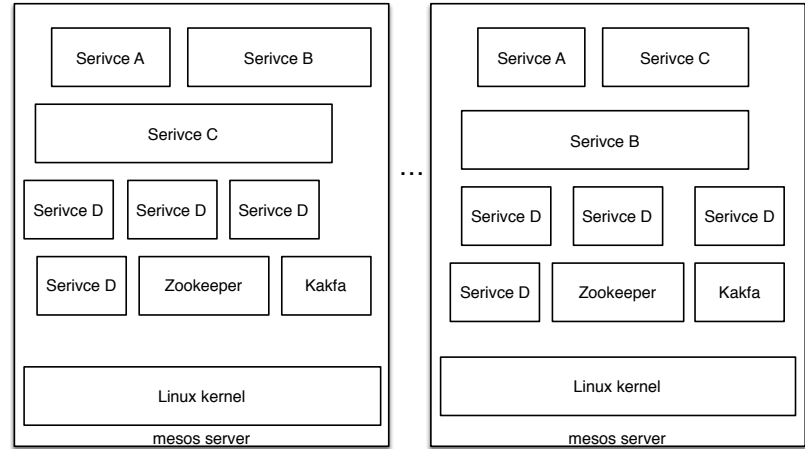
Polyglot persistence

- The data does not have to reside in the same place necessarily (same HDFS cluster)

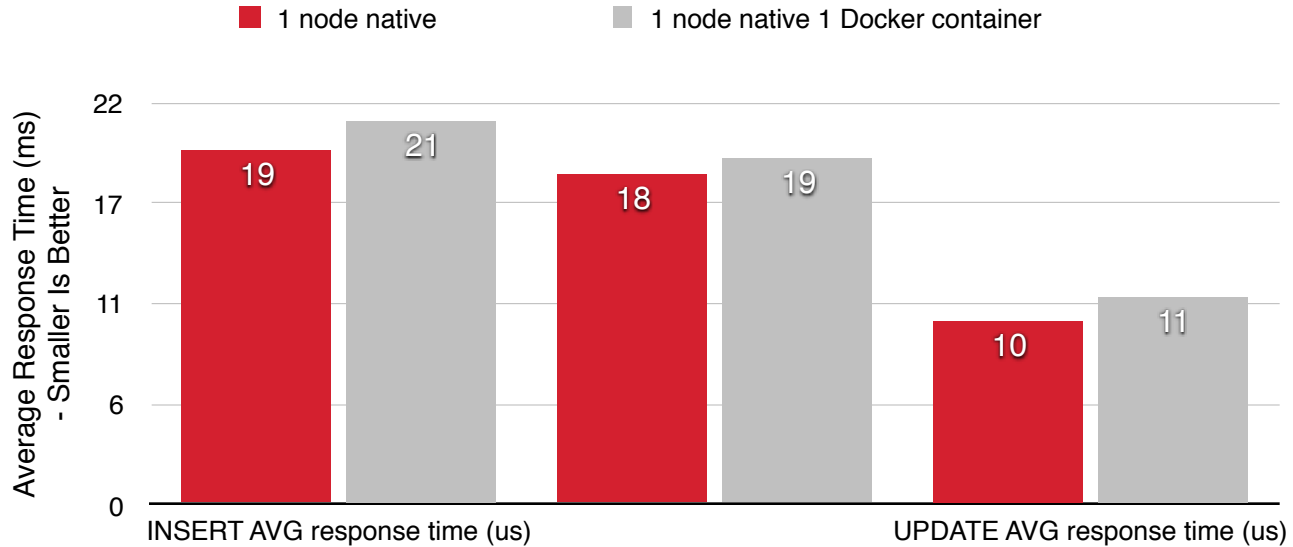


Data Service Isolation - Docker

- A Docker container is neither a VM nor a VPS
- Application level virtualisation
- Offers highest levels of consolidation
- Same kernel, apps are isolated processes
- No performance overhead
- Instant deployment
- Single app per container
- Git-like deployment method with branches and repositories.
- Isolated network stack (network namespaces)
- Isolated zero-copy UnionFS (chroot)
- CPU and Memory isolation (cgroups)



Docker vs Native - Latency



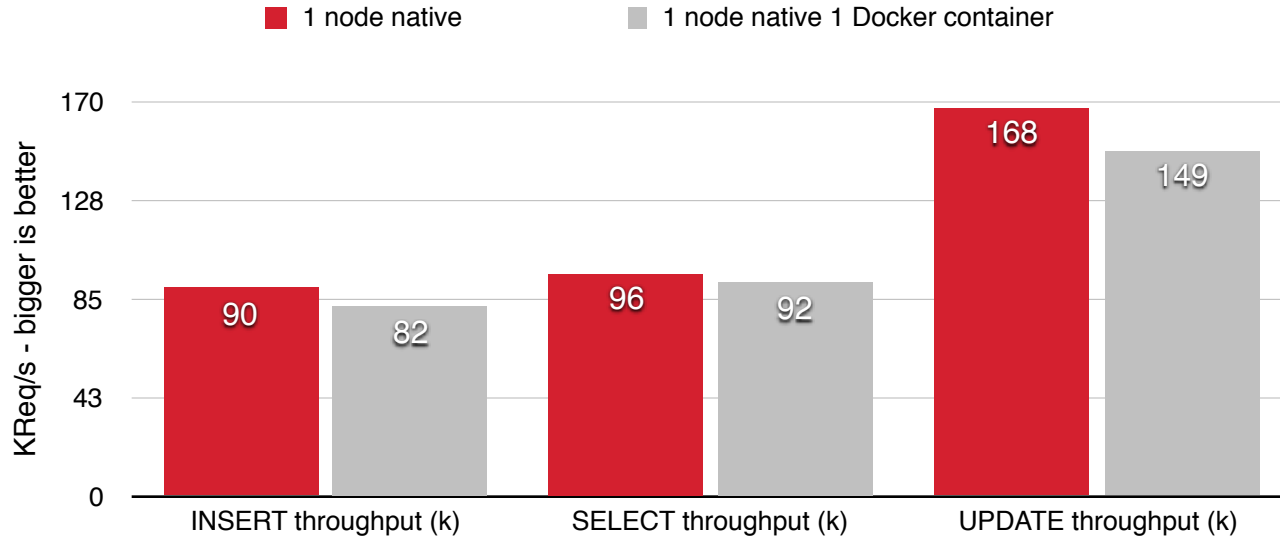
Source: Bigstep's Cassandra Benchmark 2015



BIG DATA WEEK
A GLOBAL FESTIVAL OF DATA



Docker vs Native - Throughput



Source: Bigstep's Cassandra Benchmark 2015

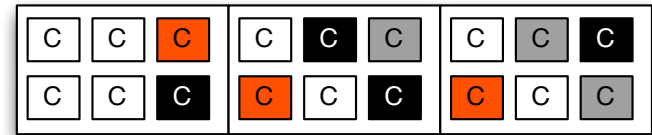


BIG DATA WEEK
A GLOBAL FESTIVAL OF DATA



Mesos & Marathon

- Allows an app's environment to be software defined.
- Docker (currently) knows only about 1 host
- Orchestration layer for Docker containers
- Out of the box load-balancing
- Monitors and restarts containers if failed
- API driven
- Useful for creating high performance, distributed, fault tolerant architectures.

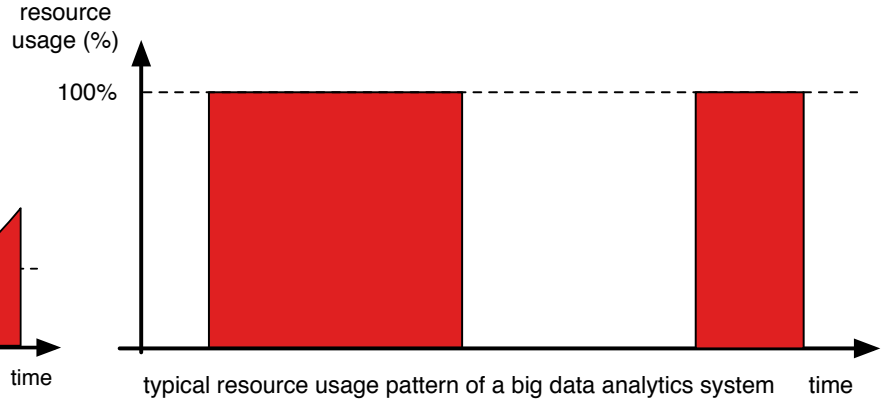
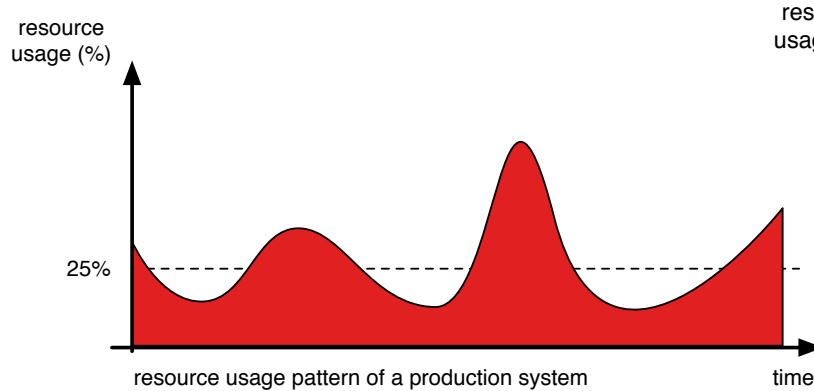


BIG DATA WEEK
A GLOBAL FESTIVAL OF DATA



Embrace streaming

- In business reaction time matters
- Resource usage patterns for streaming resemble those of web-centric systems, and need consolidation for efficiency as well as high availability



Spark with Mesos

- Spark & Spark Streaming are great candidates for building data microservices as they are very fast and easy to use
- Spark can use Mesos as a resource manager
- Spark needs YARN to access Secure HDFS YARN on Mesos: Myriad

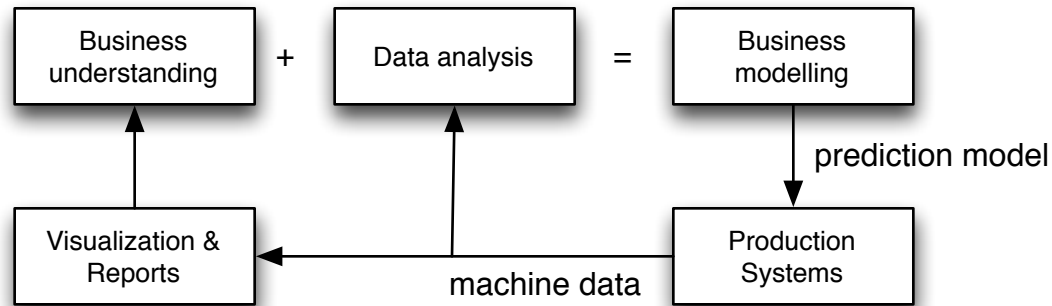


BIG DATA WEEK
A GLOBAL FESTIVAL OF DATA



Conclusions

- **Data (micro-)services** allows building a **data ecosystem** within your organisation. A team is a provider of data to other teams.
- An **agile data environment** enables an **agile business**. New tools must be inserted quickly into the mix. (Eg: found out about Looker today, why not try it on the data).
- There are methods to improve consolidation ratios with 40% while preserving performance of data services



About Bigstep - Data Lake as a service

bigstep

- High performance, bare metal cloud purpose built for big data
- Automatically deployed (managed and unmanaged) big data software stacks
- HDFS as a Service
- Managed Data Services platform based on Docker
- Purely on-demand: bare metal instances get deployed in 2 minutes, can be deleted anytime
- Locally attached drives (12 or 24 2TB drives as well as locally attached NVMe)
- SDN controlled Layer 2 networking (40Gbps per instance, cut through)
- Distributed SSD based storage fabric



elastic

DATASTAX



Couchbase



Datameer

splunk

EXASOL



BIG DATA WEEK
A GLOBAL FESTIVAL OF DATA



I'm all ears!

 @alexandrubordei

 alex@bigstep.com

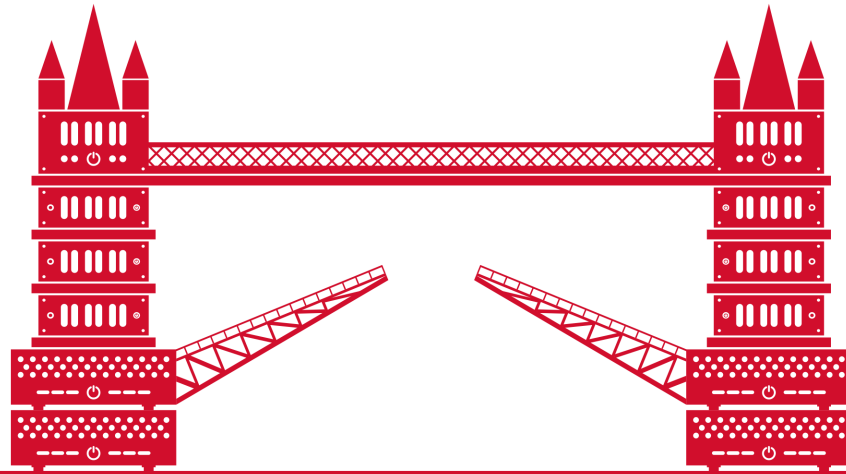
 @bigstepinc



BIG DATA WEEK
A GLOBAL FESTIVAL OF DATA

bigstep





BIG DATA WEEK
A GLOBAL FESTIVAL OF DATA

Bridging data events all over the world.

