

Course Wrap-up and Review

COMP9313: Big Data Management

Big Data: What is it?

- **Wikipedia:**

“Big data is a field that treats ways to analyze, systematically extract information from, or otherwise deal with data sets that are too large or complex to be dealt with by traditional data-processing application software”

- **Oxford English dictionary:**

“Extremely large data sets that may be analysed computationally to reveal patterns, trends, and associations, especially relating to human behaviour and interactions”

- **Apache Hadoop¹:**

“Datasets which could not be captured, managed, and processed by general computers within an acceptable scope”

¹Chen et al. 2014. Big data: A survey. Mobile networks and applications, 19(2), pp.171-209.

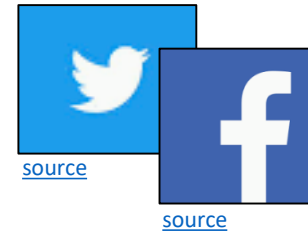
Examples of Big Data



e-commerce



Open Data



Social Media



Healthcare

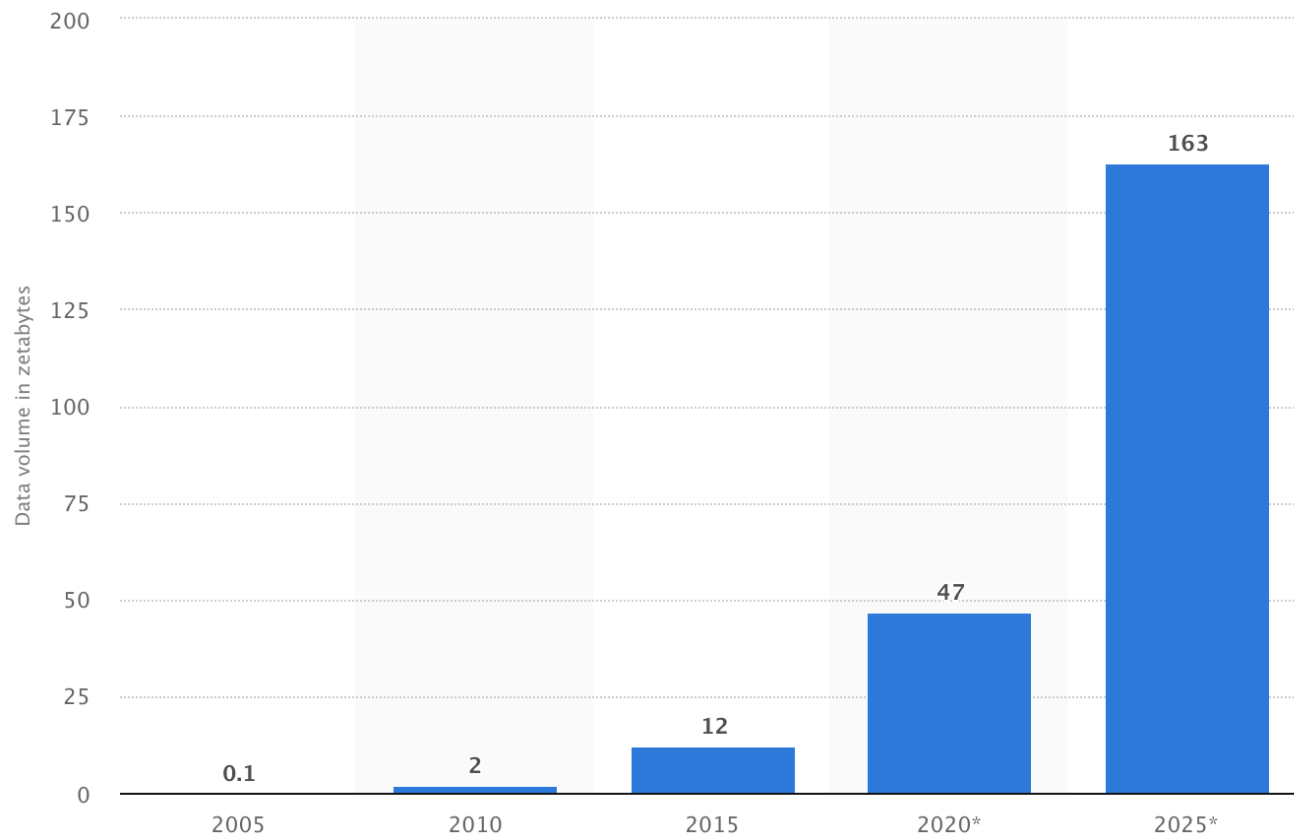


Internet of Things



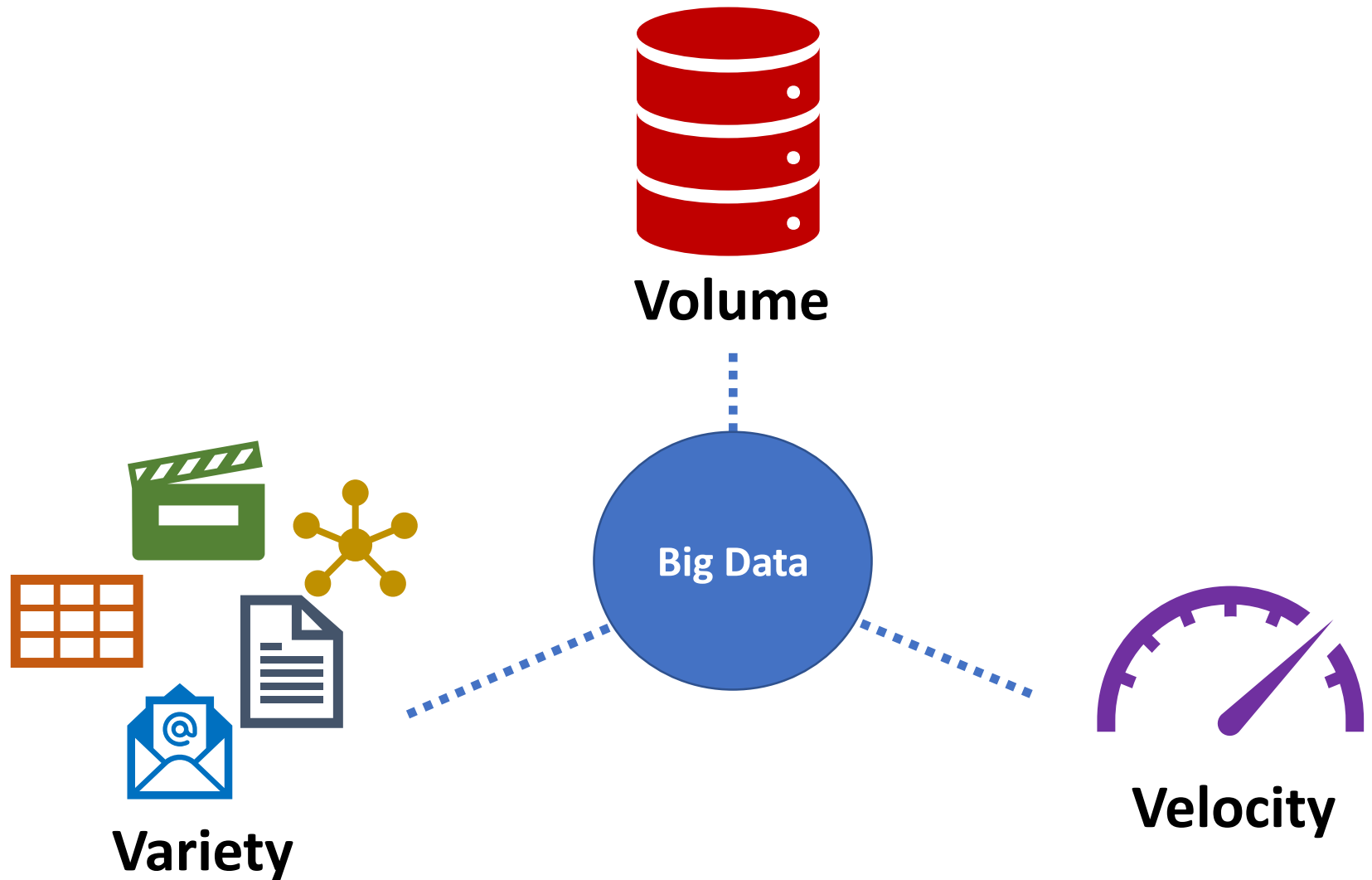
Financial Sector

Volume of data/information created worldwide from 2005 to 2025

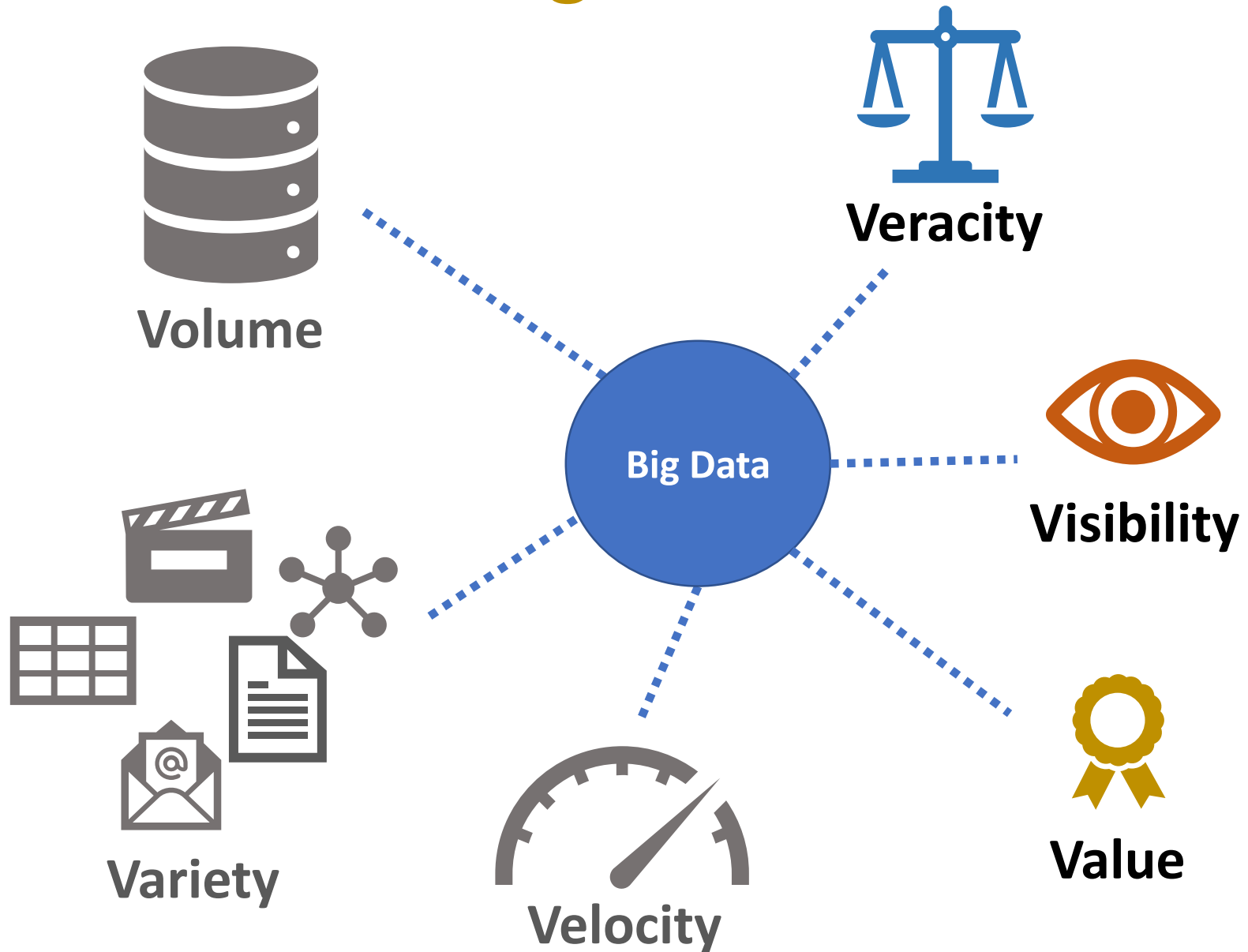


1 zettabyte = 1 million petabytes = 1 billion terabytes

The 3 Vs of Big Data



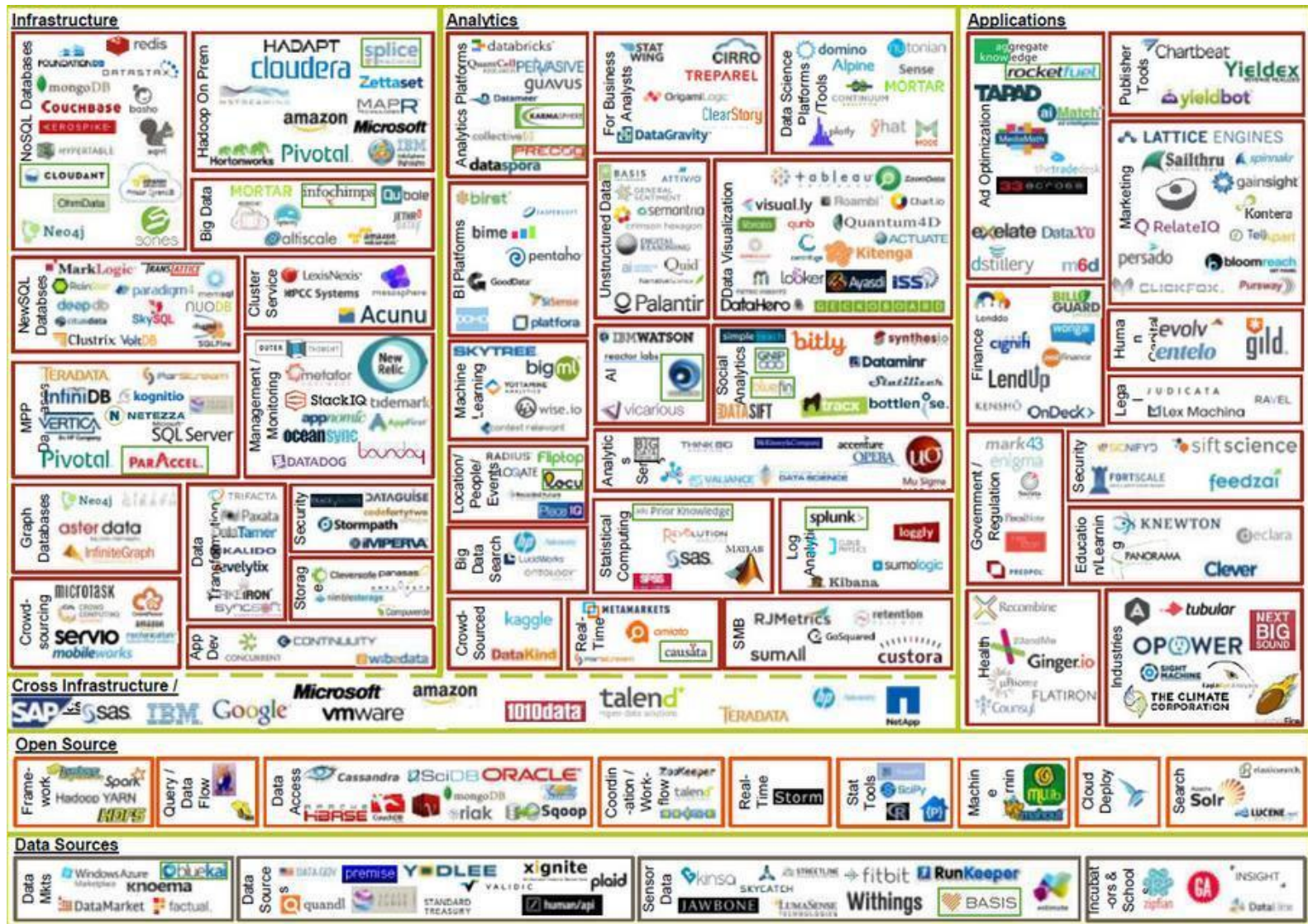
Other Vs for Big Data



Expansion of the original Vs

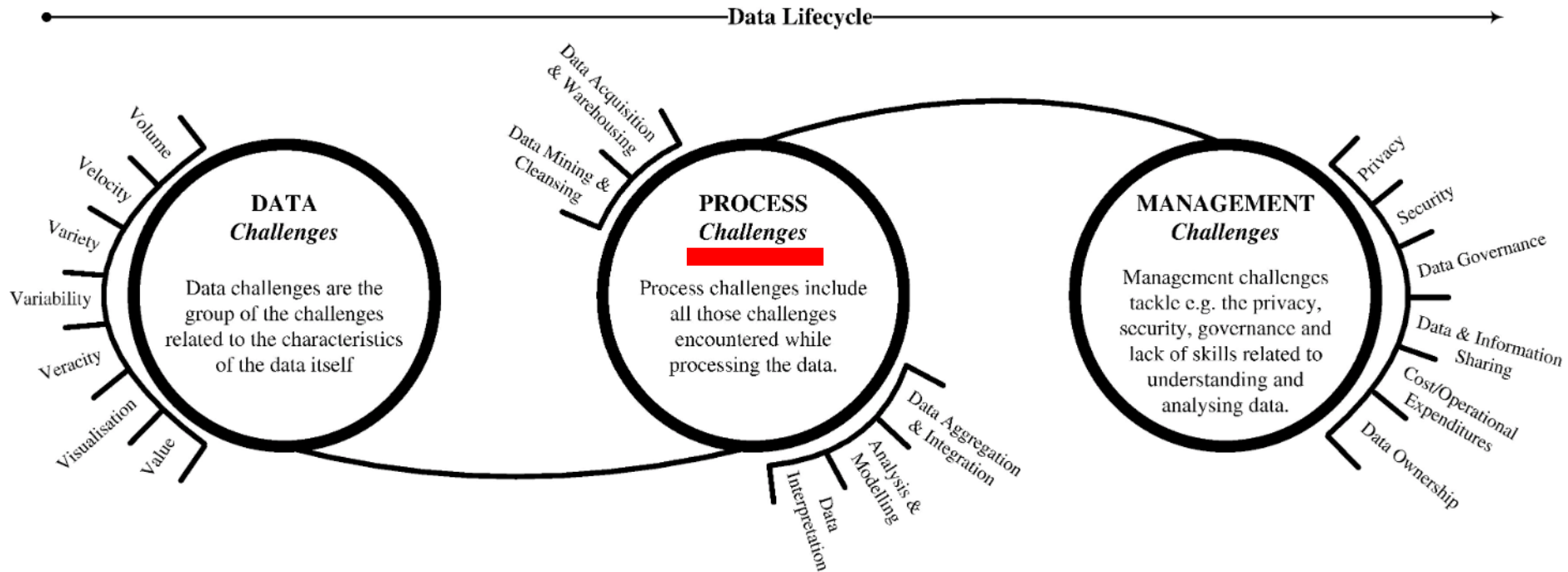
- Vagueness
- Validity
- **Value**
- Variability
- **Variety**
- **Velocity**
- Venue
- **Veracity**
- Viability
- Vincularity
- Viscosity
- **Visibility**
- Visible
- Visualization
- Vitality
- Vocabulary
- Volatility
- **Volume**

Examples of Big Data Technologies

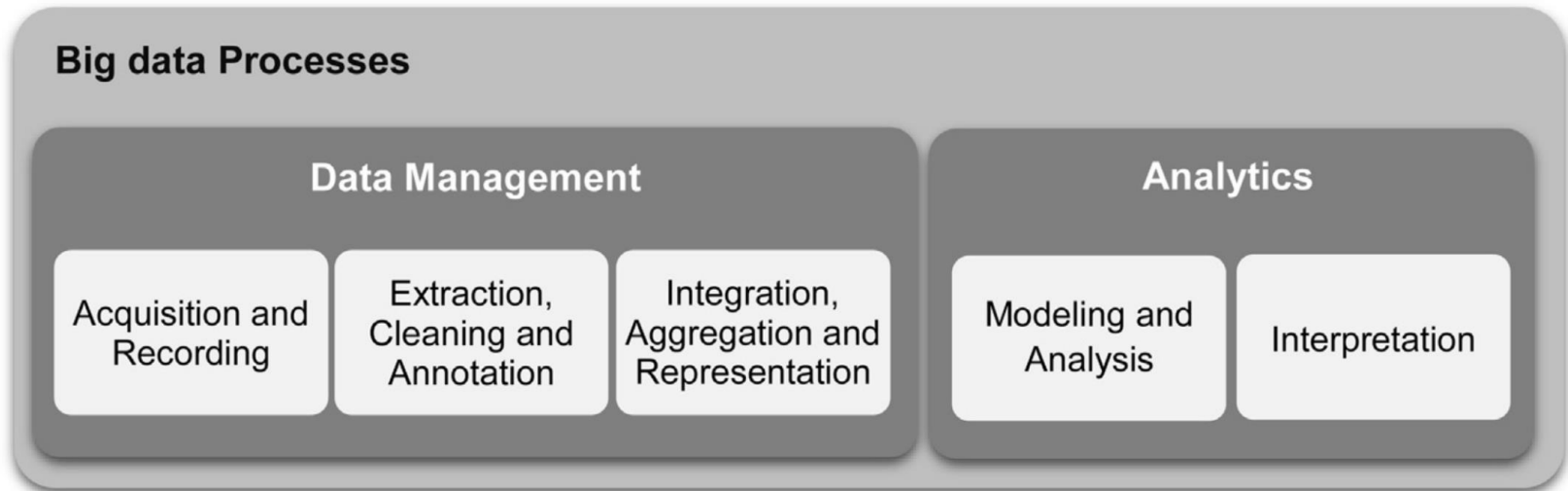


Big Data Processes and Management

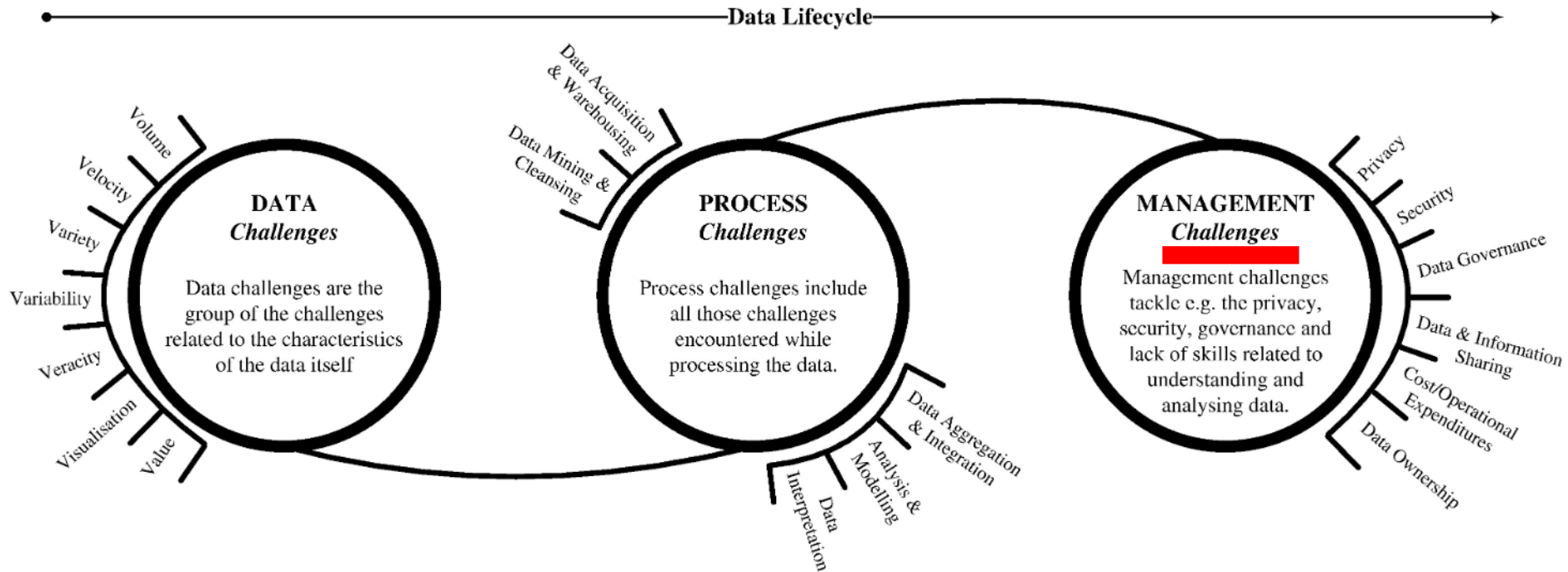
Big Data Lifecycle



Big Data Processes

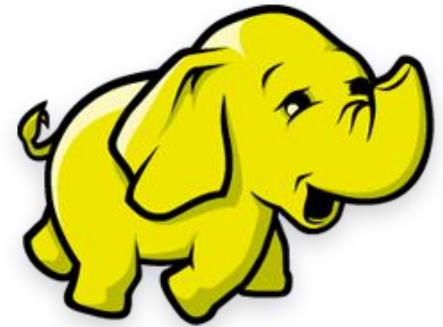


Data Lifecycle



Apache Hadoop and MapReduce

Hadoop



[source](#)

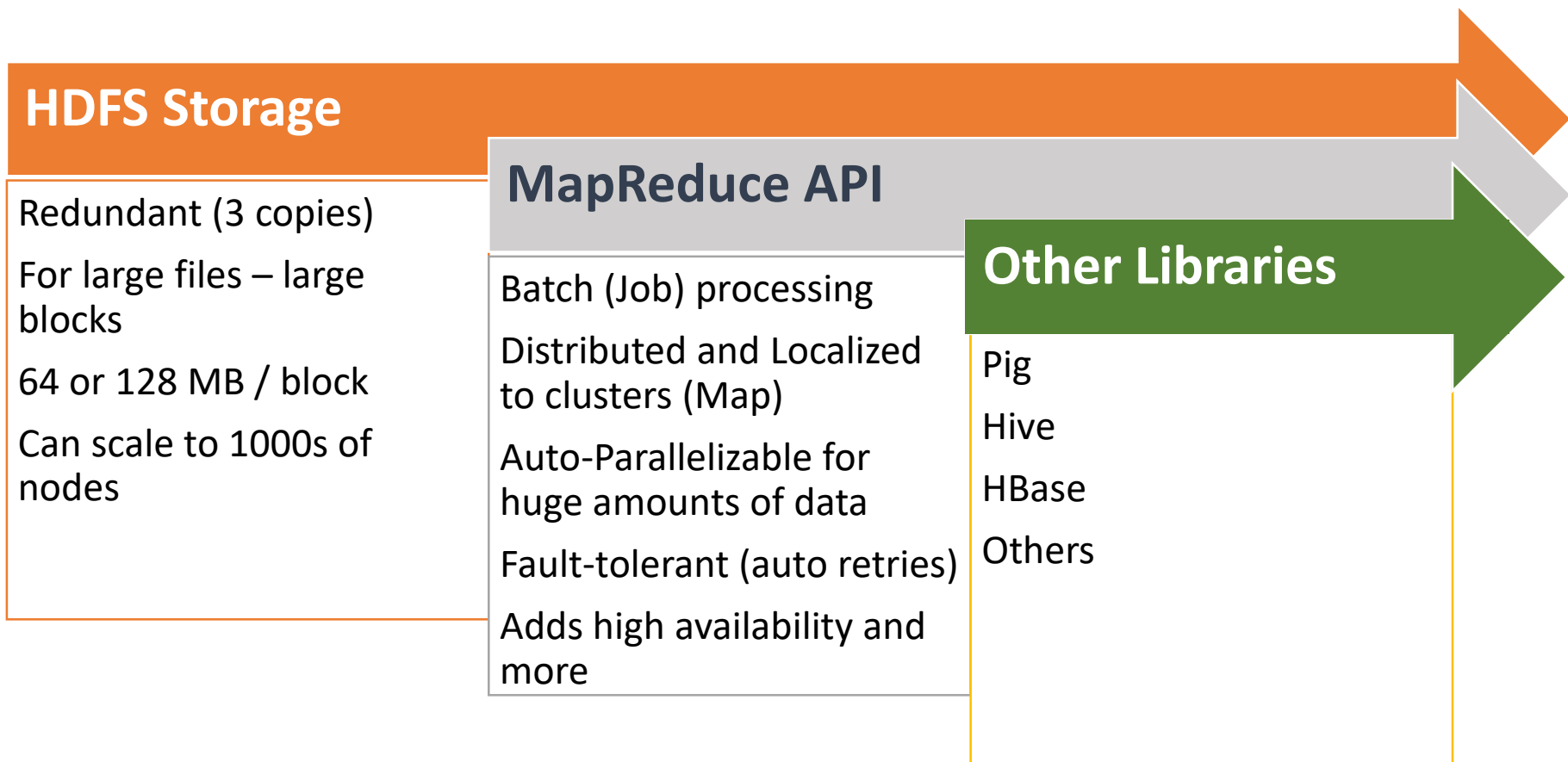
- Open-source data storage and processing platform
- Before the advent of Hadoop, storage and processing of big data was a big challenge
- Massively scalable, automatically parallelizable
 - Based on work from Google
 - Google: GFS + MapReduce + BigTable (Not open)
 - Hadoop: HDFS + Hadoop MapReduce + Hbase (opensource)
- Named by Doug Cutting in 2006 (worked at Yahoo! at that time), after his son's toy elephant

What's offered by Hadoop?

- Redundant, fault-tolerant data storage
- Parallel computation framework
- Job coordination
- Programmers do not need to worry about:
 - Where are files located?
 - How to handle failures and data loss?
 - How to distribute computation?
 - How to program for scaling?



Core Parts of a Hadoop Distribution



What is MapReduce?

- Originated from Google, [OSDI'04]
 - MapReduce: Simplified Data Processing on Large Clusters
 - Jeffrey Dean and Sanjay Ghemawat ([paper](#))
- Programming model for parallel data processing
- Hadoop can run MapReduce programs written in various languages:
e.g. Java, Ruby, Python, C++
- For large-scale data processing
 - Exploits large set of commodity computers
 - Executes process in distributed manner
 - Offers high availability

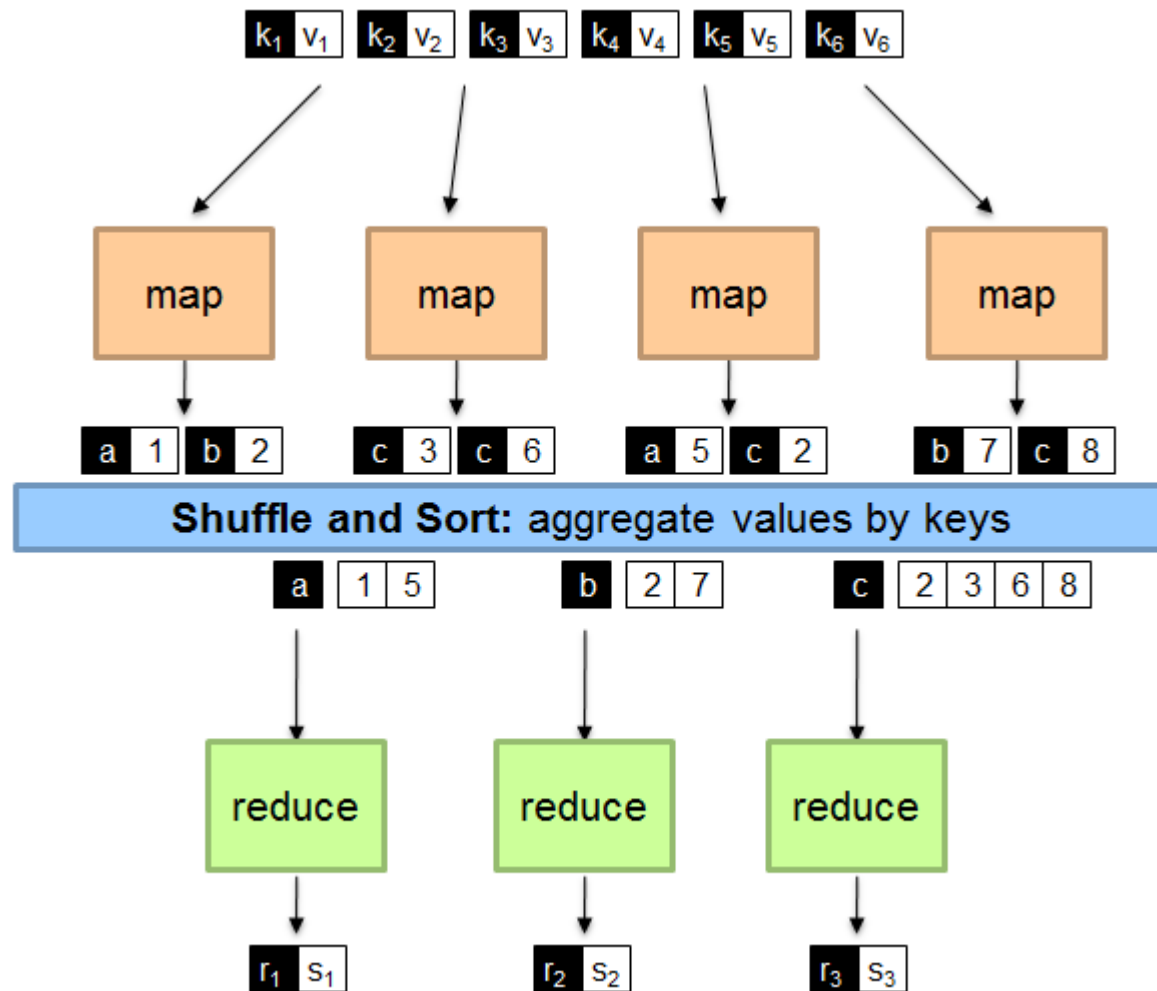
The Idea of MapReduce (1)

- Inspired by the map and reduce functions in functional programming
- We can view map as a transformation over a dataset
 - This transformation is specified by the function f
 - Each functional application happens in **isolation**
 - The application of f to each element of a dataset can be parallelized in a straightforward manner
- We can view reduce as an aggregation operation
 - The aggregation is defined by the function g
 - Data locality: elements in the list must be “brought together”
 - If we can group elements of the list, also the reduce phase can proceed in parallel
- The framework coordinates the map and reduce phases:
 - Grouping intermediate results happens in parallel

The Idea of MapReduce (2)

- Handles scheduling
 - Assigns workers to map and reduce tasks
- Handles “data distribution”
 - Moves processes to data
- Handles synchronization
 - Gathers, sorts, and shuffles intermediate data
- Handles errors and faults
 - Detects worker failures and restarts
- Everything happens on top of a distributed file system (HDFS)
- You don’t know:
 - Where mappers and reducers run
 - When a mapper or reducer begins or finishes
 - Which input a particular mapper is processing
 - Which intermediate key a particular reducer is processing

WordCount - Mapper



Map and Reduce Functions

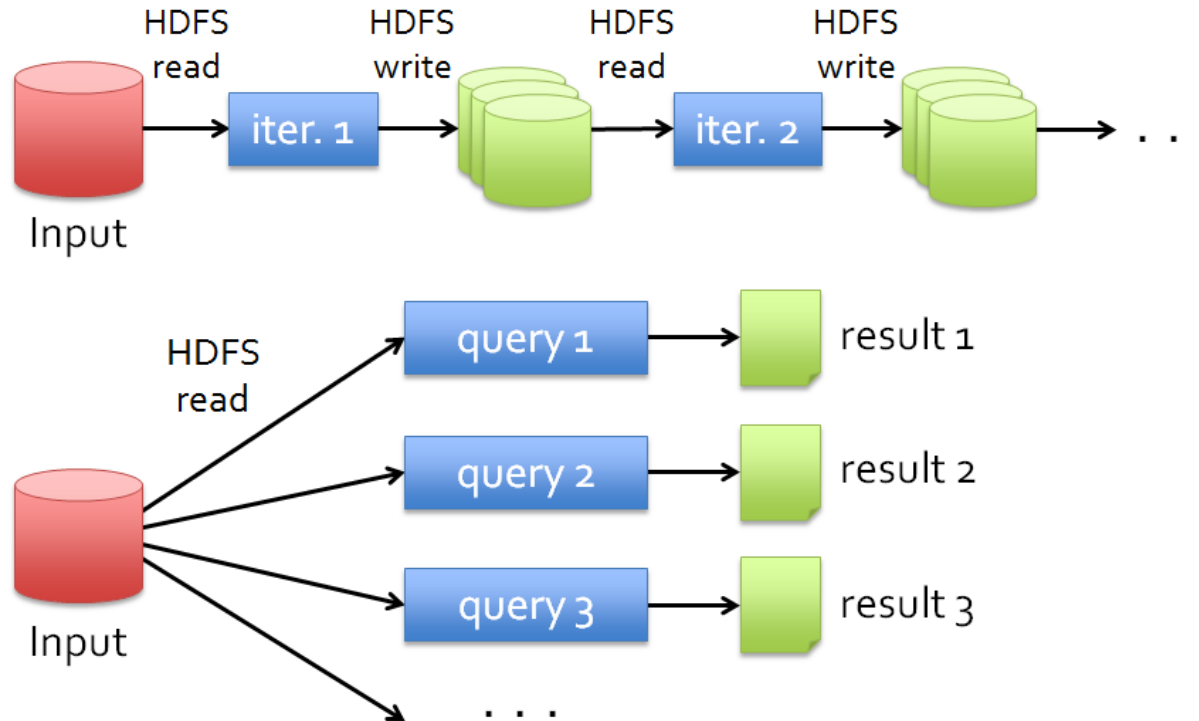
- Programmers specify two functions:
 - **map** (k_1, v_1) \rightarrow list [$\langle k_2, v_2 \rangle$]
 - Map transforms the input into key-value pairs to process
 - **reduce** ($k_2, \text{list } [v_2]$) \rightarrow [$\langle k_3, v_3 \rangle$]
 - Reduce aggregates the list of values for each key
 - All values with the same key are sent to the same reducer
 - list [$\langle k_2, v_2 \rangle$] will be grouped according to key k_2 as ($k_2, \text{list } [v_2]$)
- The MapReduce environment takes in charge of everything else...
- A complex program can be decomposed as a succession of Map and Reduce tasks

Apache Spark

Motivation of Spark

- MapReduce greatly simplified big data analysis on large, unreliable clusters. It is great at one-pass computation
- But as soon as it got popular, users wanted more:
 - More **complex**, multi-pass analytics (e.g. ML, graphs)
 - More **interactive** ad-hoc queries
 - More **real-time** stream processing
- All 3 need faster **data sharing** across parallel jobs
 - One reaction: specialized models for some of these apps, e.g.,
 - Pregel (graph processing)
 - Storm (stream processing)

Data Sharing in MapReduce



Slow due to replication, serialization, and disk IO

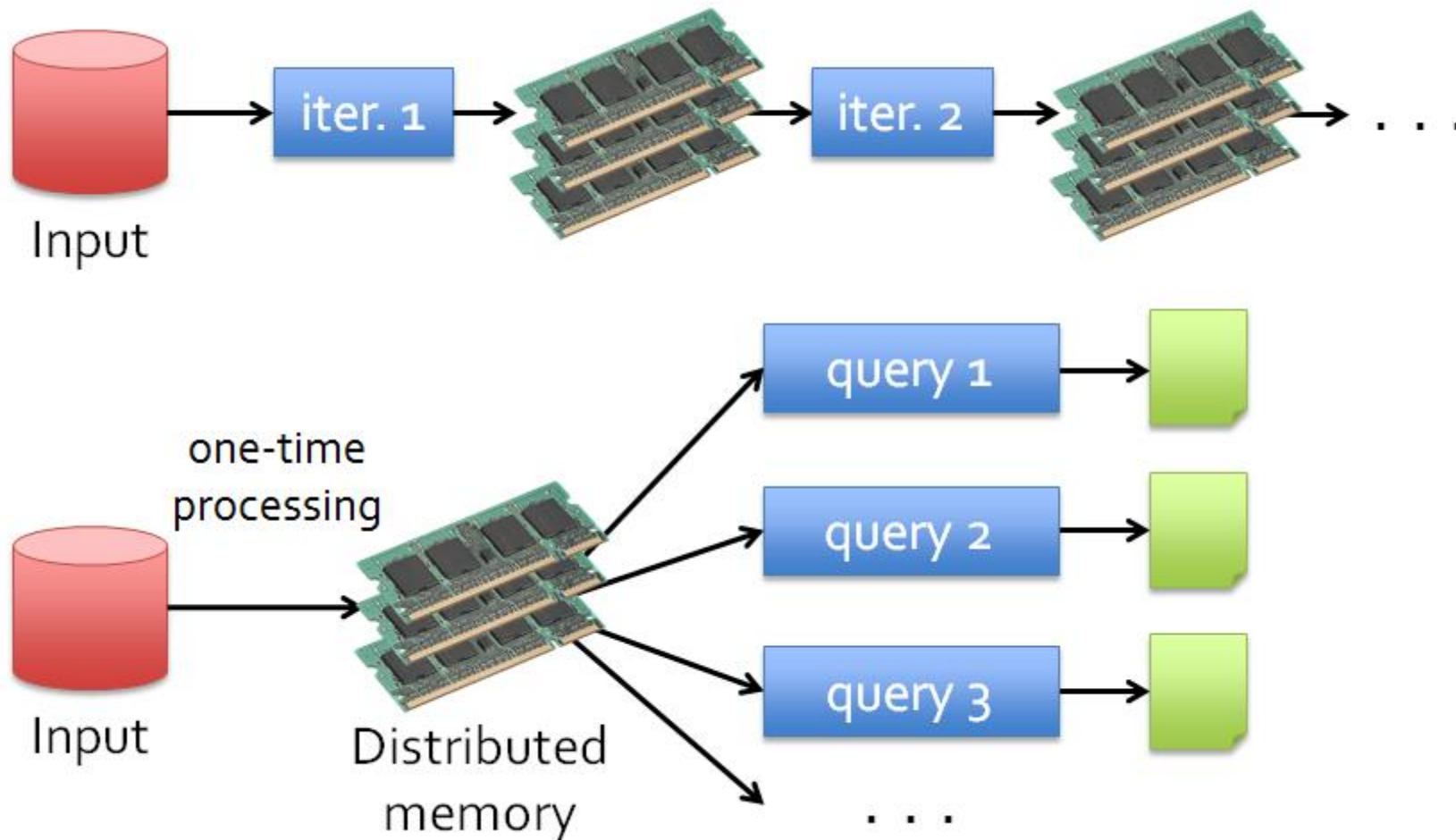
- Complex apps, streaming, and interactive queries all need one thing that MapReduce lacks:

Efficient primitives for **data sharing**

Goals of Spark

- Keep more data in-memory to improve the performance!
- Extend the MapReduce model to better support two common classes of analytics apps:
 - Iterative algorithms (machine learning, graphs)
 - Interactive data mining
- Enhance programmability:
 - Integrate into Scala programming language
 - Allow interactive use from Scala interpreter

Data Sharing in Spark with RDDs (Resilient Distributed Dataset)



10-100 × faster than network and disk

Solution: Resilient Distributed Datasets

- *Resilient Distributed Datasets (RDDs)*
 - Distributed collections of objects that can be cached in memory across cluster
 - Manipulated through parallel operators
 - Automatically recomputed on failure based on lineage
- RDDs can express many parallel algorithms, and capture many current programming models
 - Data flow models: MapReduce, SQL, ...
 - Specialized models for iterative apps: Pregel, ...

RDD Traits

- **In-Memory**, i.e. data inside RDD is stored in memory as much (size) and long (time) as possible
- **Immutable** or **Read-Only**, i.e. it does not change once created and can only be transformed using transformations to new RDDs
- **Lazy evaluated**, i.e. the data inside RDD is not available or transformed until an action is executed that triggers the execution
- **Cacheable**, i.e. you can hold all the data in a persistent "storage" like memory (default and the most preferred) or disk (the least preferred due to access speed)
- **Parallel**, i.e. process data in parallel
- **Typed**, i.e. values in a RDD have types, e.g. `RDD[Long]` or `RDD[(Int, String)]`
- **Partitioned**, i.e. the data inside an RDD is partitioned (split into partitions) and then distributed across nodes in a cluster (one partition per JVM that may or may not correspond to a single node)

RDD Operations



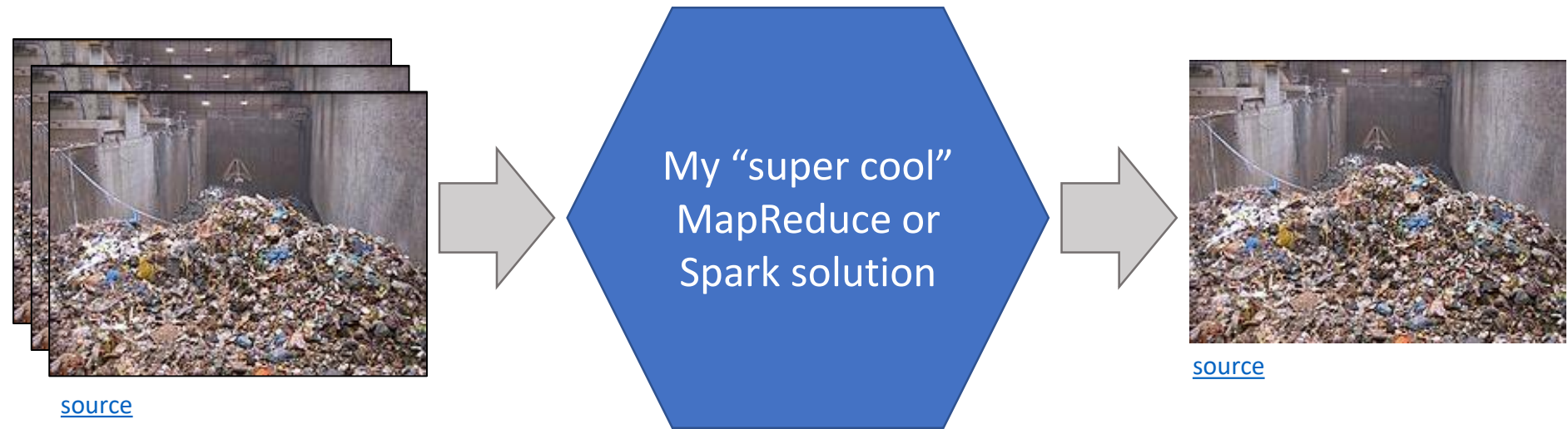
- **Transformation:** returns a new RDD
 - Nothing gets evaluated when you call a Transformation function, it just takes an RDD and return a new RDD
 - Transformation functions include *map*, *filter*, *flatMap*, *groupByKey*, *reduceByKey*, *aggregateByKey*, *filter*, *join*, etc.
- **Action:** evaluates and returns a new value
 - When an Action function is called on a RDD object, all the data processing queries are computed at that time and the result value is returned
 - Action operations include *reduce*, *collect*, *count*, *first*, *take*, *countByKey*, *foreach*, *saveAsTextFile*, etc.

How Spark Works?

- User application create RDDs, transform them, and run actions
- This results in a DAG (Directed Acyclic Graph) of operators
- DAG is compiled into stages
- Each stage is executed as a series of Task (one Task for each Partition)

Data Curation

Garbage in...



**“around 20% of records in any data source
are garbage.” *M. Stonebraker***

**Data Curation is a must in any
Big Data project!**

What is data curation?

“Data curation is the process of identifying which data sources are needed, putting that data in the context of the business so that business users can interact with it, understand it, and use it to create their analysis.”

Data Quality Dimensions

- Accessibility
- Appropriate amount of data
- Believability
- Completeness
- Concise representation
- Ease of manipulation
- Free-of-error
- Interpretability
- Objectivity
- Relevancy
- Reputation
- Security
- Timeliness
- Understandability
- Value-added

Data Curation

- Ingestion
- Validation
- Transformation
- Correction
- Consolidation
- Visualization

NoSQL Technologies

What is NoSQL?

- The name stands for Not Only SQL
- Does not use SQL as querying language
- Class of non-relational data storage systems
- The term NOSQL was introduced by Eric Evans when an event was organized to discuss open source distributed databases
- It's not a replacement for a RDBMS but compliments it
- All NoSQL offerings relax one or more of the ACID properties (will talk about the CAP theorem)

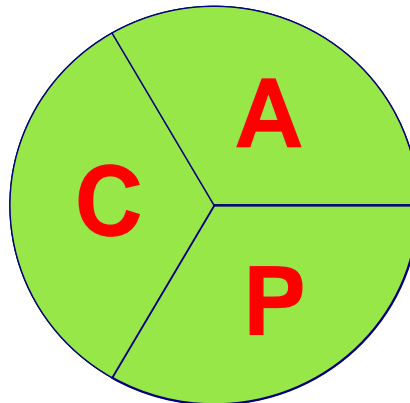


Why NoSQL?

- Web apps have different needs (than the apps that RDBMS were designed for)
 - Low and predictable response time (latency)
 - Scalability & elasticity (at low cost!)
 - High availability
 - Flexible schemas / semi-structured data
 - Geographic distribution (multiple datacenters)
- Web apps can (usually) do without
 - Transactions / strong consistency / integrity
 - Complex queries

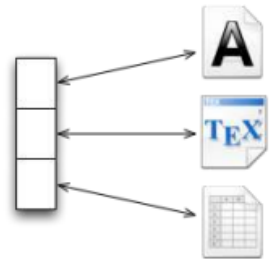
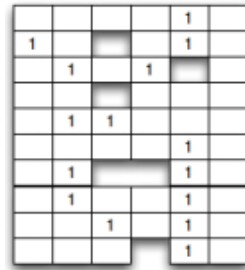
CAP Theorem

- Consider these three properties of a distributed system (sharing data)
 - Consistency:
 - all copies have same value
 - Availability:
 - reads and writes always succeed
 - Partition-tolerance:
 - system properties (consistency and/or availability) hold even when network failures prevent some machines from communicating with others



No SQL Taxonomy

- Key-Value stores
 - Simple K/V lookups (Distributed Hash Table (DHT))
- Column stores
 - Each key is associated with many attributes (columns)
 - NoSQL column stores are actually hybrid row/column stores
 - Different from “pure” relational column stores!
- Document stores
 - Store semi-structured documents (JSON)
- Graph databases
 - Neo4j, etc.
 - Not exactly NoSQL
 - can't satisfy the requirements for High Availability and Scalability/Elasticity very well



Elasticsearch



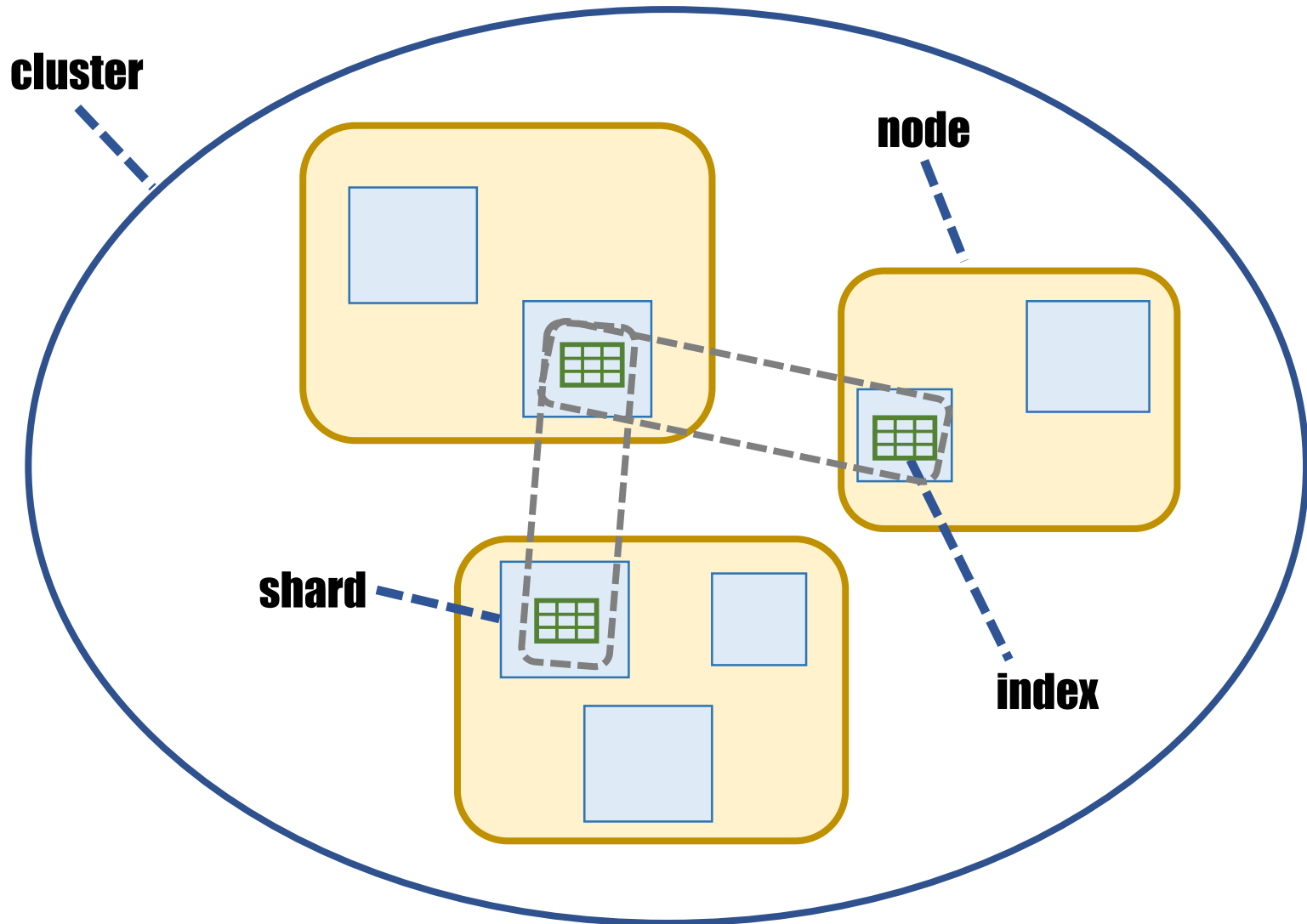
- Open source search engine based on Apache Lucene
- Initial release in 2010
- Provides a distributed, full-text search engine with a REST APIs
 - E.g. GET `http://localhost:9200/person/student/8871`

Elasticsearch

- Document oriented (JSON as serialization format for documents)
- Developed in Java (cross platform)
- Focused on scalability – distributed by design
- Highly efficient search



Elasticsearch Ecosystem



Data mining vs. Process mining

Data mining

- association rules
- graphs
- sequences (of items)
- clusters

Process mining

- process models
- control flows
- decision points
- process execution data

Assessment

- Assessment:
 - 40% formal written exam: individual assessment.
 - 45% on assignment work
 - Assignment1 on MapReduce(individual) 20%
 - Assignment2 on Spark (individual) 25%
 - 15% on 5 online quizzes (WebCMS-based quiz system, 'open' test)
- Final Mark = quizzes + assignments + exam (**No Hurdle**)

The Final Exam

- Duration: Two hours
- Lab Based Exam
- Types of questions:
 - Multiple Choice
 - Short Answers
 - Written Answers

The Final Exam (examples of Questions)

- Multiple Choice Question:

Which of the “Vs” listed below were the three first Vs proposed to characterize big data?

[A] Value, Volume, Velocity.

[B] Volume, Velocity, Viability.

[C] Validity, Variability, Volume.

[D] Velocity, Volume, Variety.

The Final Exam (examples of Questions)

- Short Answer:

In Hadoop Security by default it is setup as **blankA**. When you activate Hadoop in Secure each user and service needs to be authenticated by **blankB** in order to use Hadoop services.

Write down the word(s) that could fill in **blank A** and **blank B**.

The Final Exam (examples of Questions)

- **Written Answer:**

Question: In the context of Apache Spark, in your own words explain what is a Narrow Transformation.

The Final Exam (examples of Questions)

Describe the output of the following Pig code snippet.

```
students = LOAD 'gs://mybucket_pig/students.txt' USING PigStorage(',') as
(zid:chararray, name:chararray);
marks = LOAD 'gs://mybucket_pig/marks.txt' USING PigStorage(',') as
(zid:chararray, m1:int, m2:int);

joined = JOIN students by zid, marks by zid;

finaltable = FOREACH joined GENERATE students::zid, students::name, marks::m1,
marks::m2;

DESCRIBE finaltable;
DUMP finaltable;
```

Supplementary Exam Policy

- Supp Exam is only available to students who:
 - DID NOT attend the final exam
 - Have a good excuse for not attending
 - Have documentation for the excuse
- Submit special consideration within 72 hours (via myUNSW with supporting docs)
- Everybody gets exactly one chance to pass the final exam. For CSE supplementary assessment policy, follow the link in the course outline.

Warning... Warning

Huge Thanks Coming
Your Way...

Thanks for the Teaching Team

- Course Administrator
 - Maisie Badami
- Tutors
 - Mohammad Ali Yaghub Zade Fard
 - Maisie Badami
 - Sara Mumtaz
 - AliReza Tabebordbar
 - Maryam Ghafouri



That's all Folks!