

Overview

Prof. Hyuk-Yoon Kwon

<https://sites.google.com/view/seoultech-bigdata>

Contents

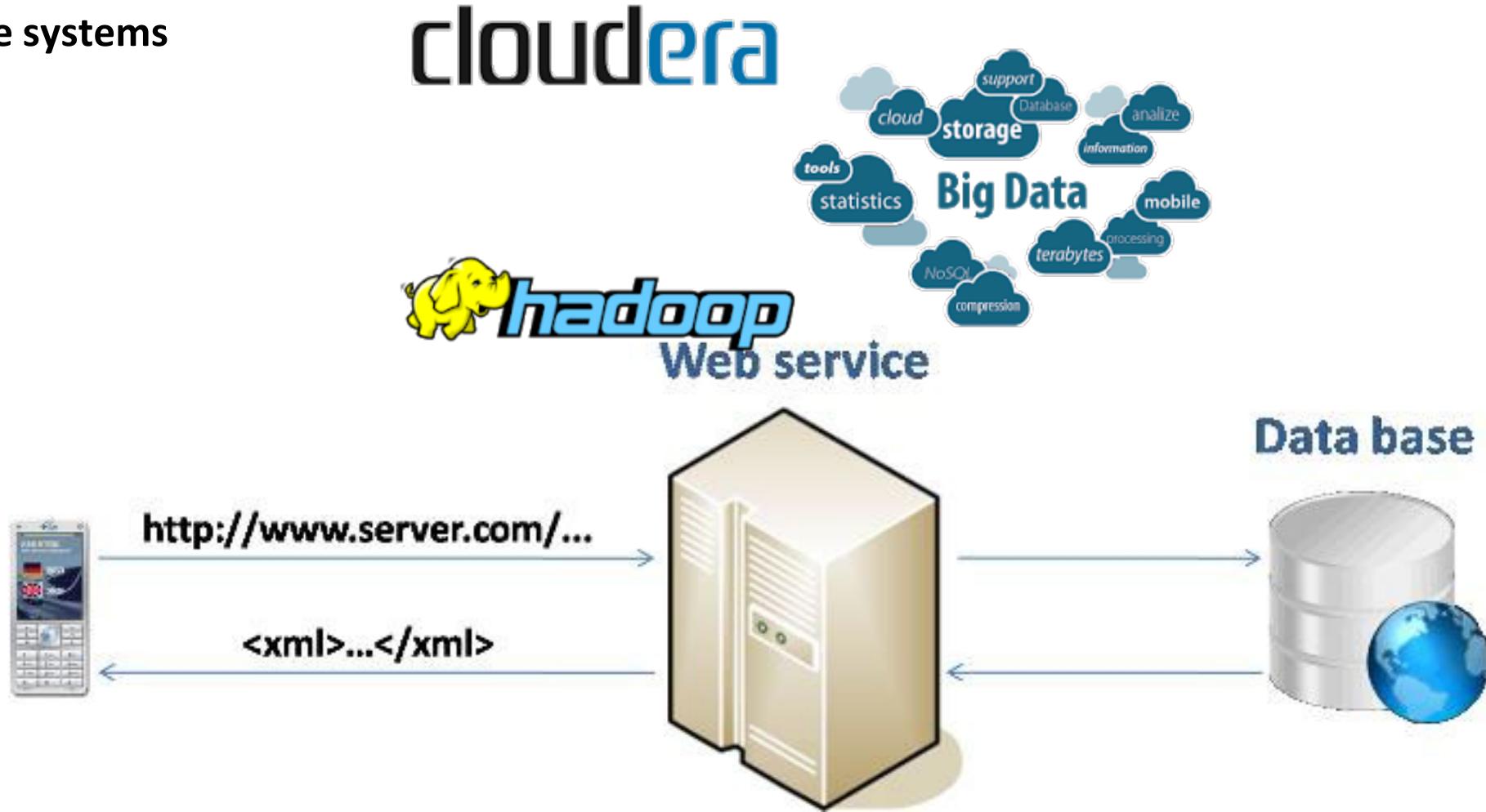
■ Part1: Course Introduction

- Overall course policy
- Grading
- Schedule

■ Part2: Introduction to Big Data Management Systems

Course Theme

- To learn and practice many big software development skills focusing on big data management
- To learn how to deal with the **big data management systems** and how to build real-scale services based on the systems



The Aim of Course

1. To learn big data management systems and practice them using real data sets
2. To make real services based on big data management and analysis



Course Policy: Grading

■ Written Exams (70%)

- Midterm (30%)
- Final (40%)

■ Group projects (30%)

- Idea proposal presentation (10%)
- Final presentation (20%)

Group Projects

■ The requirements are twofold:

1. To use **big data management** systems for storing and managing data efficiently in a big data management environment
2. To conduct meaningful **analysis** based on the managed data. The analysis should move beyond simple summaries to generate valuable insights (e.g., trend detection, predictive modeling, comparative analysis). The final results should be presented in a clear and effective manner (e.g., interactive Web services with visualization).

■ Team formation

- One team will be formed with 3 or 4 students
- The results will be noticed in e-class

■ Presentation: Offline or Online (as you preferred)

- Intermediate presentation (near mid-term exam)
- Final presentation (near final exam)

Course Policy: Lectures

■ Flipped learning

- (Online lecture contents) The contents will be uploaded **by the previous Tuesday**
- (Offline lecture) **Every Monday 2:00 PM ~ (Assigned lecture time: from 2PM to 5PM)**
 - Summary (will be given in English)
 - Q&A (Korean questions are allowed)

■ Lecture place

- Laboratory: Frontier 501

Course Policy: Attendance

■ Attendance is strongly recommended!

- Important points are announced only during the lecture
- Exam problems are inferred from the lecture

■ But, attendance itself is not included in the score

■ Attendance confirmation is done by e-class

- <http://eclasse.seoultech.ac.kr>

Course Policy: Understanding

■ Encourage questions!

- Questions in Korean are allowed
- After the class, use emails and Q&A in e-class aggressively

■ For the questions that break the flow of class or need to discuss

- Please visit me during office hours (or any time if you make reservations)

Course Assumption

■ **Required: Basic programming knowledge**

- Any kind of programming languages – Python, Java, or C/C++

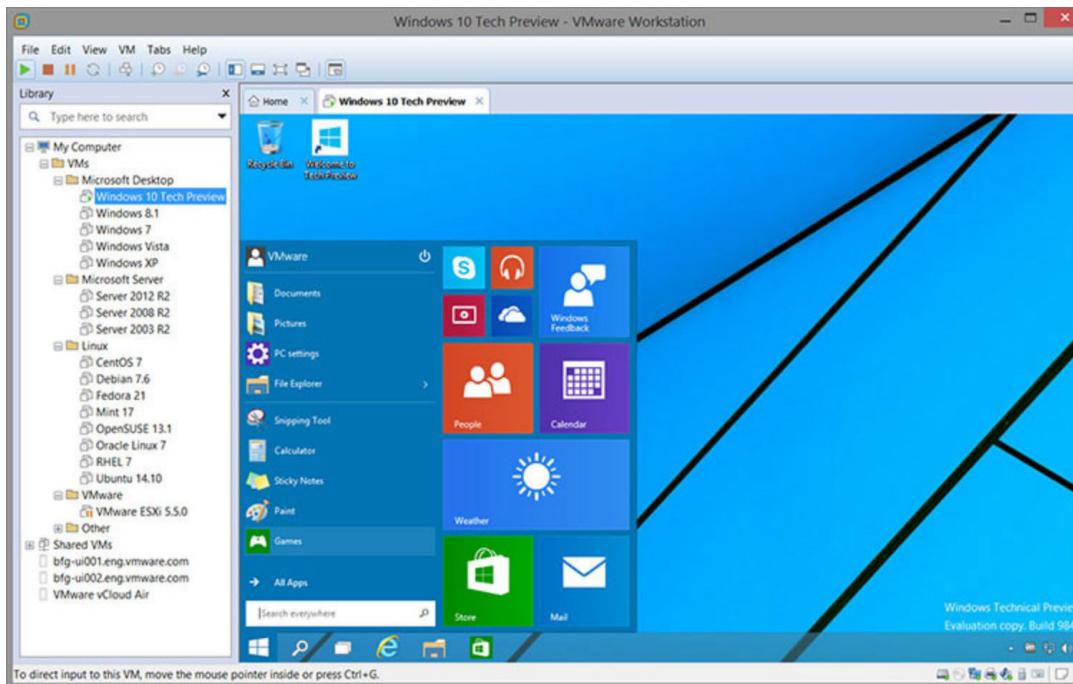
■ **Required: Database management (ITM411)**

■ **This course is designed for ITM students!**

Environments for Practices

■ Cloudera Big Data Management Environments

- The instance for the virtual machine will be provided
- Every software and tools were installed in the instance



■ Other Big Data Management and Analysis Tools

Textbooks

■ The following textbooks are recommended, but are not necessary

[1] Hadoop: The Definitive Guide (4th Edition) — Tom White (O'Reilly)

- Hadoop basics (HDFS, MapReduce, YARN), with introductions to ecosystem projects (Hive, Pig, HBase).
- Considered the "*Bible of Hadoop*"

[2] Hadoop in Practice (2nd Edition) — Alex Holmes

- Covers Hadoop with ecosystem components (YARN, Hive, HBase, Pig, Sqoop, Spark, Storm).
- Good balance between fundamentals and applied recipes.

[3] Big Data: Principles and Paradigms — Rajkumar Buyya et al.

- Academic textbook covering distributed file systems, MapReduce, Spark, Hive, and NoSQL.
- Broader than just Hadoop, often used in university courses.

Getting Help

■ E-class: <http://eclass.seoultech.ac.kr>

- Complete schedule of lectures, exams, and assignments
- Lecture slides, assignments, exams, solutions
- Clarifications to assignments
- Q&A

■ Office hours

- Wednesday, 2:00-4:00pm, Frontier 614
- You can schedule 1:1 appointment in advance even not for the office hours

Course Schedule

▣ 강의계획서 정보

주차	강의내용	강의 진도 계획
1 (영문)	Overview	Flipped learning 강의방법, 과제, 평가
2 (영문)	NoSQL Introduction: Hadoop Architecture	Flipped learning
3 (영문)	Importing data from databases to Hadoop	Flipped learning
4 (영문)	Big data modeling and management	Flipped learning
5 (영문)	Big data partitioning	Flipped learning
6 (영문)	Big data processing: spark basics	Flipped learning
7 (영문)	Intermediate presentation	Presentation
8 (영문)	Mid-term exam	Exam

9 (영문)	Spark Basics	Flipped learning
10 (영문)	Spark working with RDDs	Flipped learning
11 (영문)	Spark Aggregation Data	Flipped learning
12 (영문)	Spark Parallel Processing	Flipped learning
13 (영문)	Spark Algorithms	Flipped learning
14 (영문)	Final exam	Exam
15 (영문)	Final presentation	Presentation

Big Data-Driven AI Laboratory



Big Data-Driven AI Laboratory (BIGBASE)

학부생 연구원 (석사과정 연계), 석/박사과정, 박사후 과정 모집 중입니다.

We are looking for post-doc researchers, MS/Ph.D. students, and undergraduate students who are passionate about 1) Data Engineering (ML Ops and Data Ops), 2) AI-focused Data Analysis, and 3) Large-Scale Data Management (Cloud or Distributed).

We are targeting global top-level research-top conferences (SIGMOD, NeurIPS, AAAI, ICDE, CIKM, ICDM) and top journals (TKDE, TII, TCC) in the field of AI and Big Data.

If you are interested in joining our Lab, please send an email ([hyukyoon.kwon \[at\] seoultech.ac.kr](mailto:hyukyoon.kwon@seoultech.ac.kr)) including 1) CV, 2) cover letter, and 3) transcript. Then, we may have an interview online or offline.

Research Area



- Data-Driven AI
 - Anomaly detection and forecasting in time-series ([NeurIPS2024](#), [IEEE TII2024](#), [AAAI2025](#))
 - Representation learning in graphs ([SIGMOD2025](#), [IEEE TKDE2024](#), [ICDM2023](#))
 - Special-purpose natural language processing ([CIKM2025b](#))
 - Image segmentation and video highlight detection ([Applied Soft Computing 2022](#))
- AI-focused Practical Analysis
 - Fairness of datasets and AI models ([CIKM2025a](#))
 - Continual learning ([IEEE BigData2024a](#), [IEEE BigData2024b](#))
 - Self-training ([IEEE TII2024](#))
 - Federated learning ([ICDE2025](#))
 - Multimodal learning
 - Large language model ([Knowledge-Based Systems 2024](#))
- Scalable Data Computing
 - Data scraping ([Data & Knowledge Engineering 2024](#))
 - Edge-Cloud computing ([IEEE TII2025](#))
 - Distributed data pipeline and ingestion

<https://bigdata.seoultech.ac.kr>

Introduction to Big Data Management Systems

Course contents

1	Introduction	Course Introduction
2	Introduction to Hadoop and the Hadoop Ecosystem	Introduction to Hadoop
3	Hadoop Architecture and HDFS	
4	Importing Relational Data with Apache Sqoop	Importing and Modeling Structured Data
5	Introduction to Impala and Hive	
6	Modeling and Managing Data with Impala and Hive	
7	Data Formats	
8	Data Partitioning	
9	Capturing Data with Apache Flume	Ingesting Streaming Data
10	Spark Basics	
11	Working with RDDs in Spark	Distributed Data Processing with Spark
12	Aggregating Data with Pair RDDs	
13	Writing and Deploying Spark Applications	
14	Parallel Processing in Spark	
15	Spark RDD Persistence	
16	Common Patterns in Spark Data Processing	
17	Spark SQL and DataFrames	
18	Conclusion	Course Conclusion

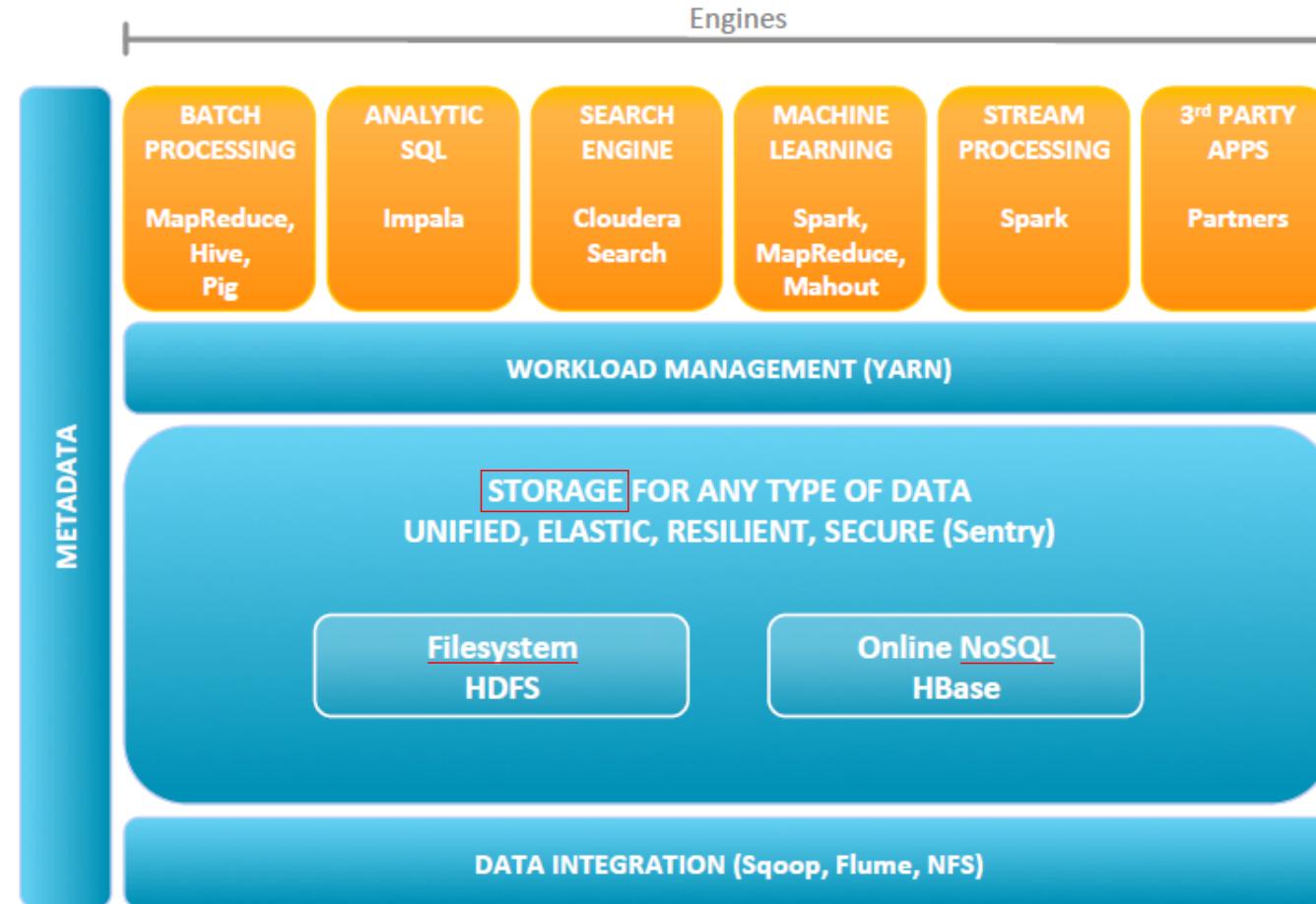
Course objectives

During this course, you will learn

- How the **Hadoop Ecosystem** fits in with the data processing lifecycle
processing, management, storing data
- How **data is distributed**, stored and processed in a Hadoop **cluster**
- How to use **Sqoop** and **Flume** to **ingest data** *collect, store big data*
- How to process distributed data with **Spark**
- Best practices for data storage
SQL-like
- How to model structured data as tables in **Impala** and **Hive**
- How to choose a data storage format for your data usage patterns

CDH (Cloudera's Distribution including Apache Hadoop)

- 100% open source, enterprise-ready distribution of Hadoop and related projects
- The most complete, tested, and widely-deployed distribution of Hadoop
- Integrates all the key Hadoop ecosystem projects
- Available as RPMs and Ubuntu, Debian, or SuSE packages, or as a tarball

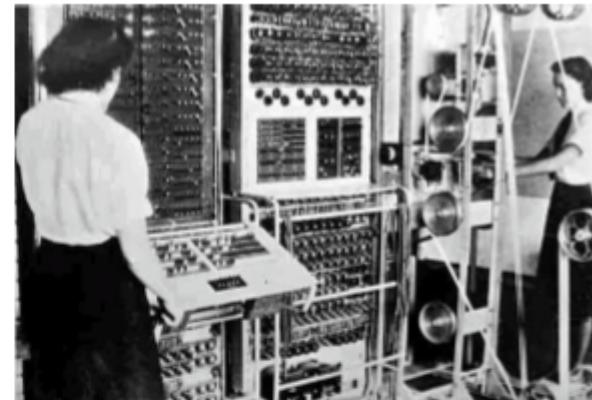


Traditional Large-Scale Computation

why distributed system is needed?

- Traditionally, computation has been processor-bound

- Relatively small amounts of data
- Lots of complex processing



- The early solution: bigger computers *(super computer)* \Rightarrow single computer

- Faster processor, more memory
- But even this couldn't keep up
cannot store all generated data ...

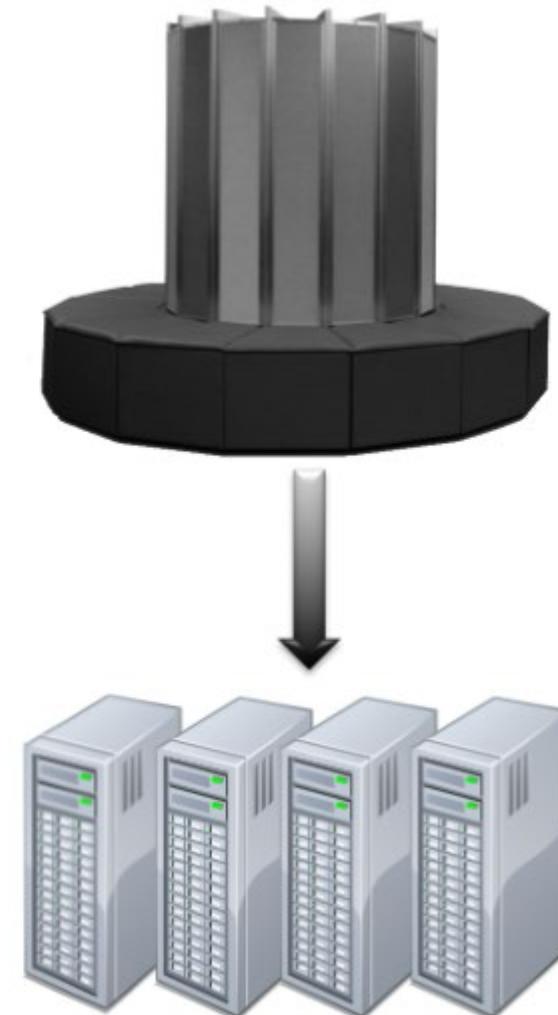


Distributed Systems

- The better solution: more computers
 - Distributed systems – use multiple machines for a single job *multiple paradigm shift*

“In pioneer days they used oxen for heavy pulling, and when one ox couldn’t budge a log, we didn’t try to grow a larger ox. We shouldn’t be trying for bigger computers, but for *more systems* of computers.”

– Grace Hopper



Challenges with Distributed Systems

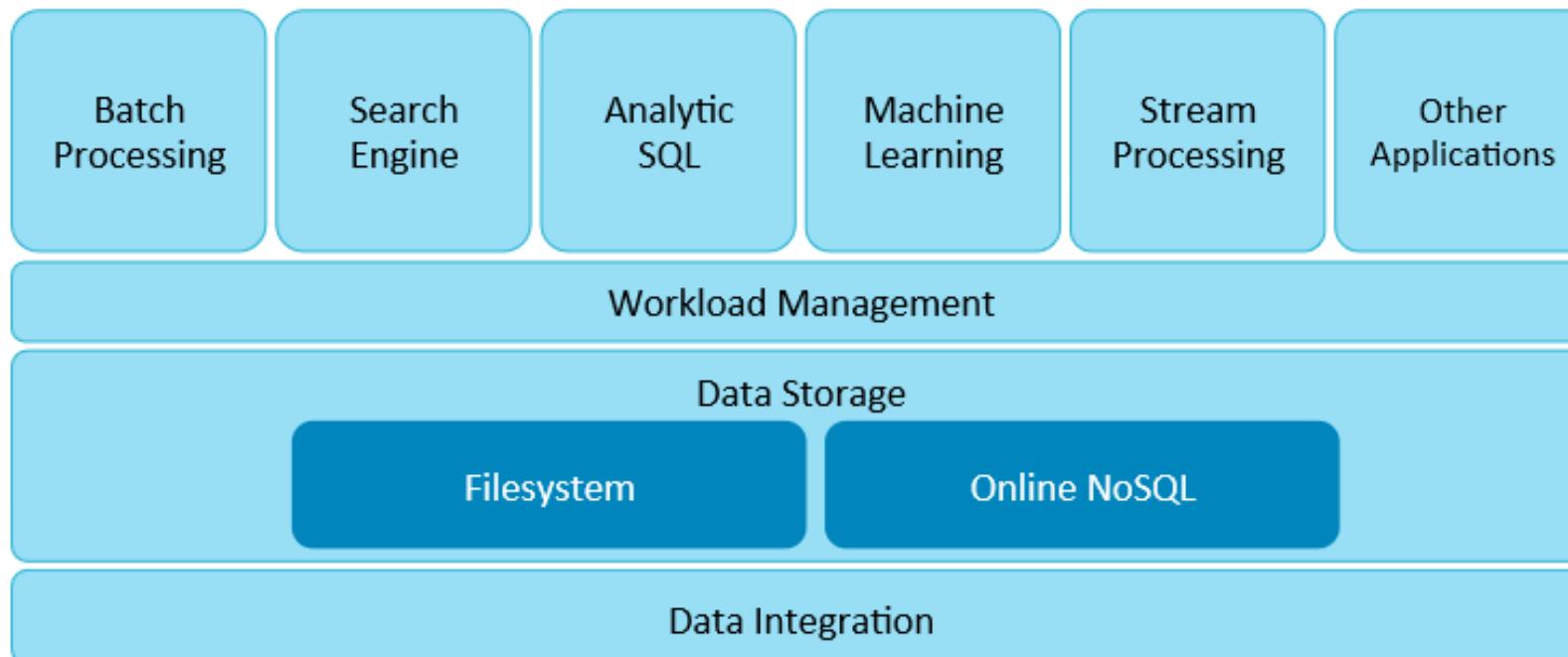
- Challenges with distributed systems
 - Programming complexity
 - Keeping data and processes in sync
 - Finite bandwidth (*network*)
 - Partial failures \Rightarrow replacement is required.

- The solution?
 - Hadoop!

What is Apache Hadoop?

Core < Storage
processing

- Scalable and economical data storage, processing and analysis
 - Distributed and fault-tolerant
 - Harnesses the power of industry standard hardware
- Heavily inspired by technical documents published by Google



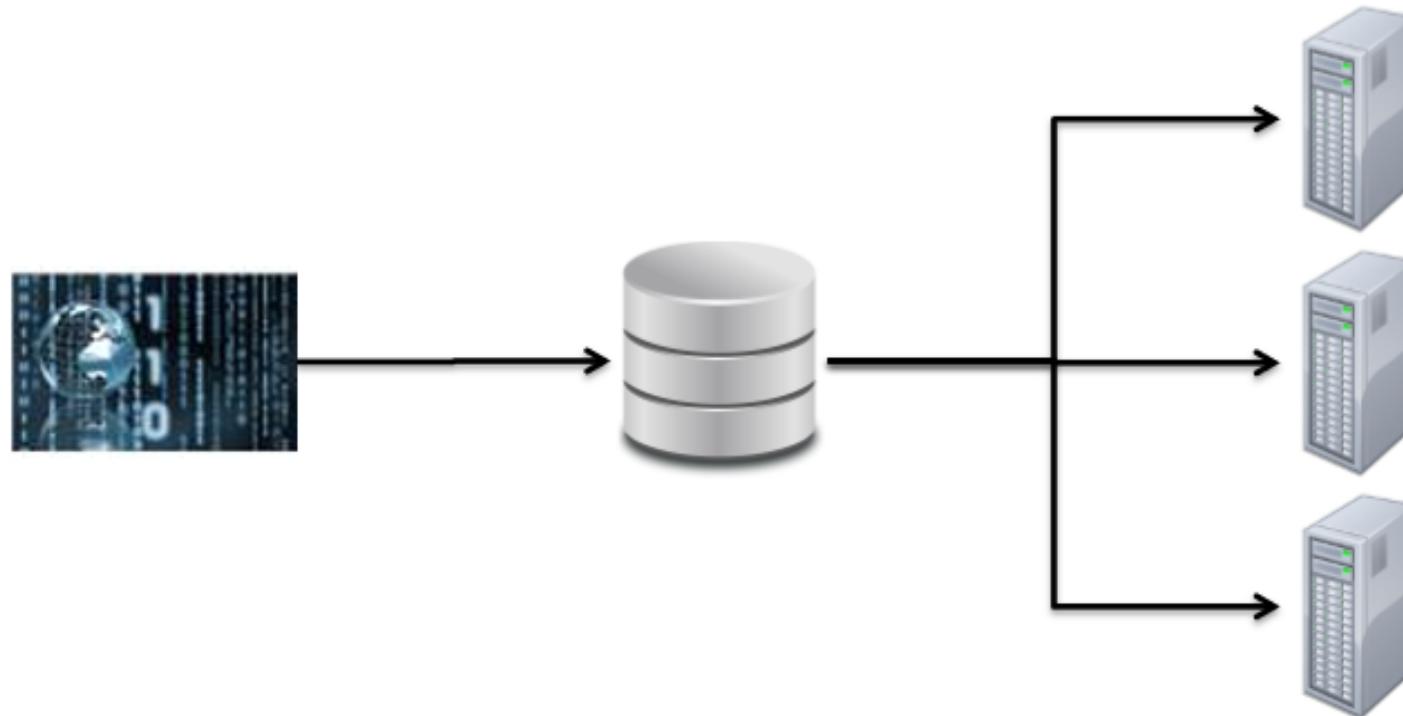
Common Hadoop Use Cases

- Extract/Transform/Load (ETL)
 - Text mining
 - Index building
 - Graph creation and analysis
Shortest path
 - Pattern recognition
 - Collaborative filtering ⇒ recommendation
 - Prediction models
 - Sentiment analysis ⇒ 감성분석 (감정.주관)
 - Risk assessment
-
- What do these workloads have in common? Nature of the data...

big data {
- Volume : lots of data. large scale
- Velocity : be able to process them as fast as it comes.
- Variety :
 { Structured data from RDBMS
 { semi-structured data (log files...)
 { unstructured data (text, HTML, PDF ...)
 { multimedia data } not homogeneous

Distributed Systems: The Data Bottleneck (1)

- Traditionally, data is stored in a central location
- Data is copied to processors at runtime
- Fine for limited amounts of data



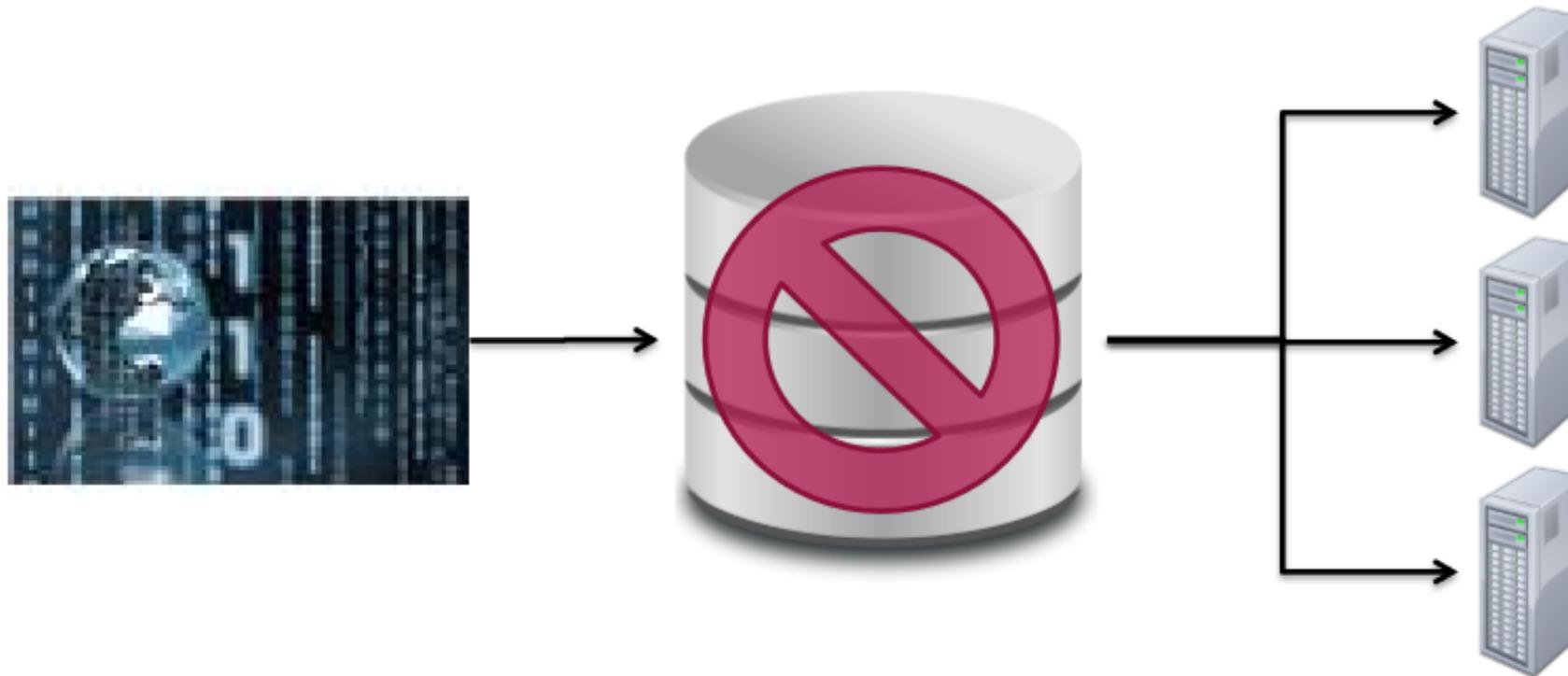
Distributed Systems: The Data Bottleneck (2)

- Modern systems have much more data

- terabytes+ a day
- petabytes+ total

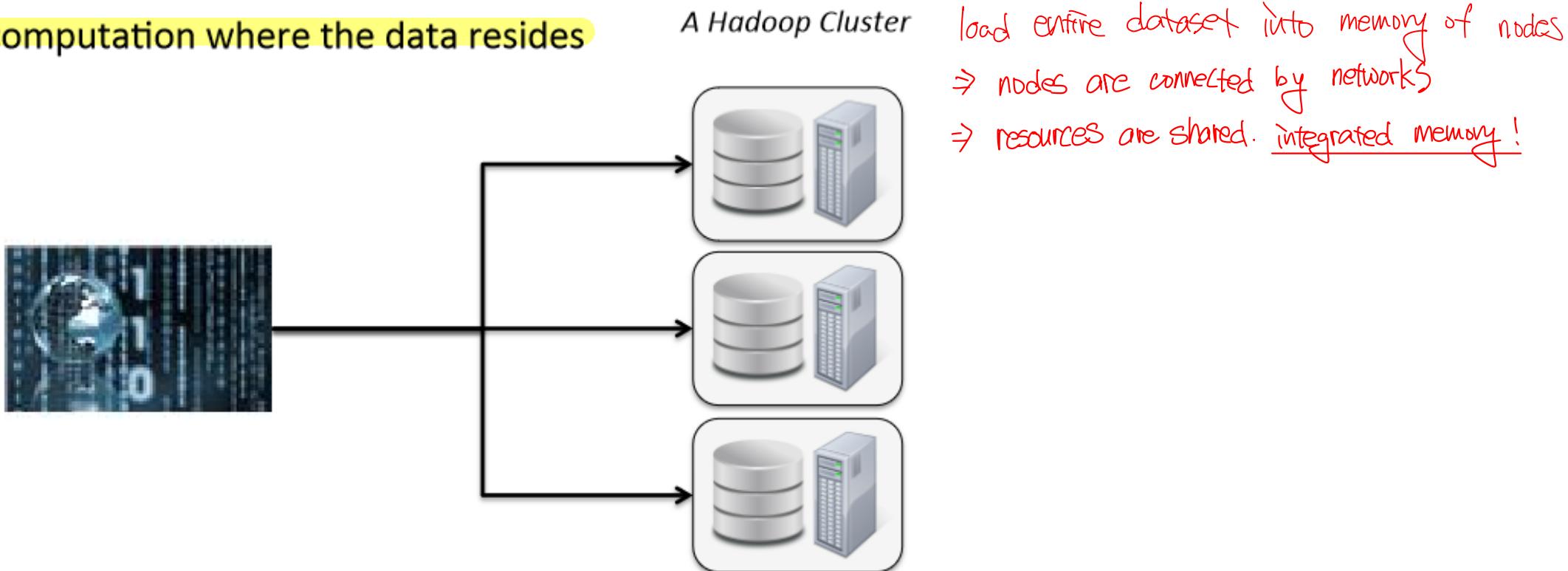
{ data cannot be stored in single storage

- We need a new approach...

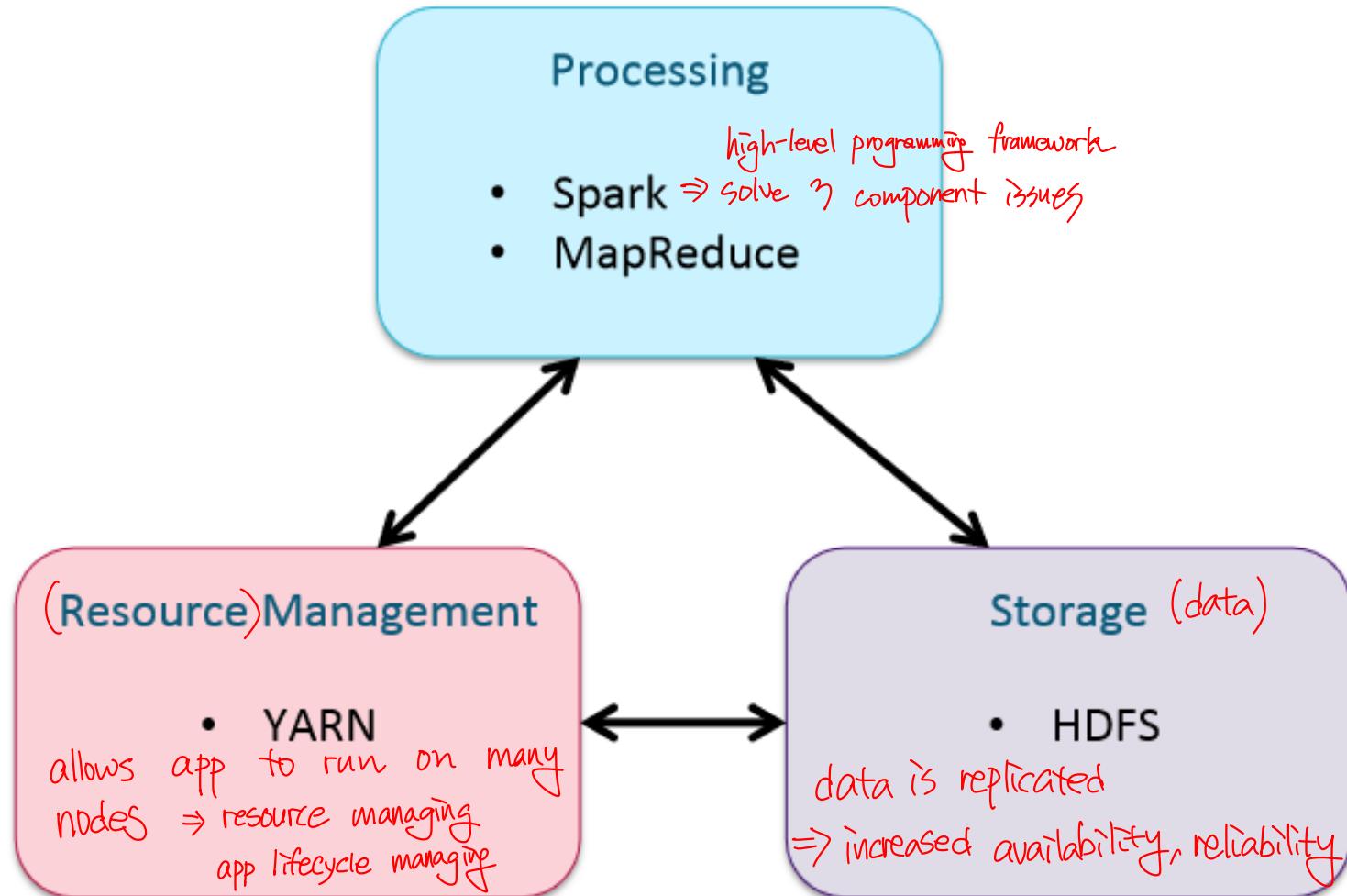


Big Data Processing with Hadoop

- Hadoop introduced a **radical new approach:** ** data → subject*
 - Bring the program to the data rather than the data to the program
- Based on two key concepts
 - Distribute data when the data is stored
 - Run computation where the data resides



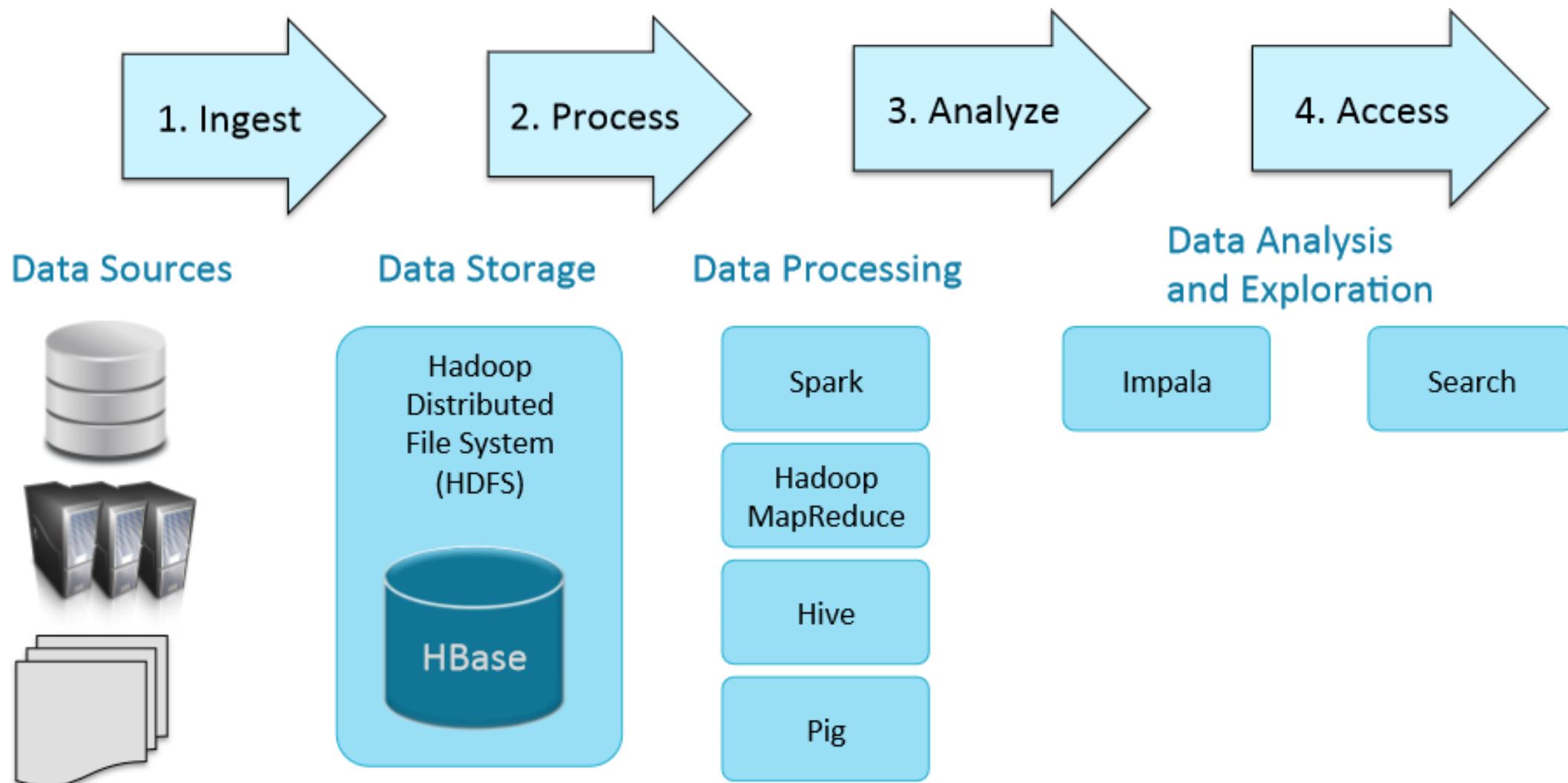
Core Hadoop



A Hadoop Cluster



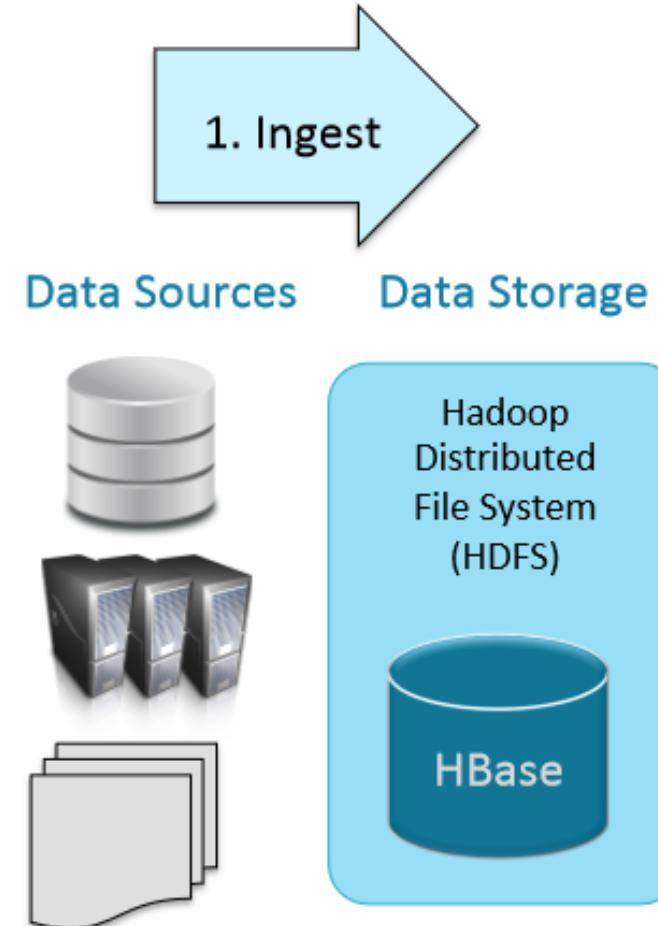
Big Data Processing



Data Ingest and Storage

- Hadoop typically ingests data from many sources and in many formats

- Traditional data management systems, e.g. databases
- Logs and other machine generated data (event data)
- Imported files



Data Storage

- **Hadoop Distributed File System (HDFS)**
 - HDFS is the storage layer for Hadoop *lowest layer*
 - Provides inexpensive reliable storage for massive amounts of data on industry-standard hardware
 - Data is distributed when stored
 - Covered later in this course
- **Apache HBase: The Hadoop Database**
 - A NoSQL distributed database built on HDFS
 - Scales to support very large amounts of data and high throughput
 - A table can have thousands of columns
 - Covered in depth in *Cloudera Training for Apache HBase*



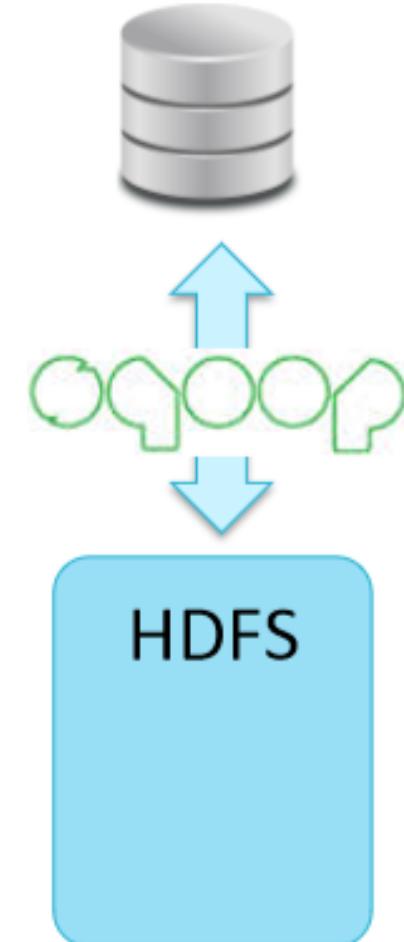
Data Ingest Tools (1)

- **HDFS**

- Direct file transfer

- **Apache Sqoop**

- High speed import to HDFS from Relationship Database (and vice versa)
 - Supports many data storage systems
 - e.g. Netezza, Mongo, MySQL, Teradata, Oracle
 - Covered later in this course



Data Ingest Tools (2)

Apache Flume

- Distributed service for ingesting streaming data
- Ideally suited for event data from multiple systems
 - For example, log files

generated in real-time

IoT devices

✓ 61?

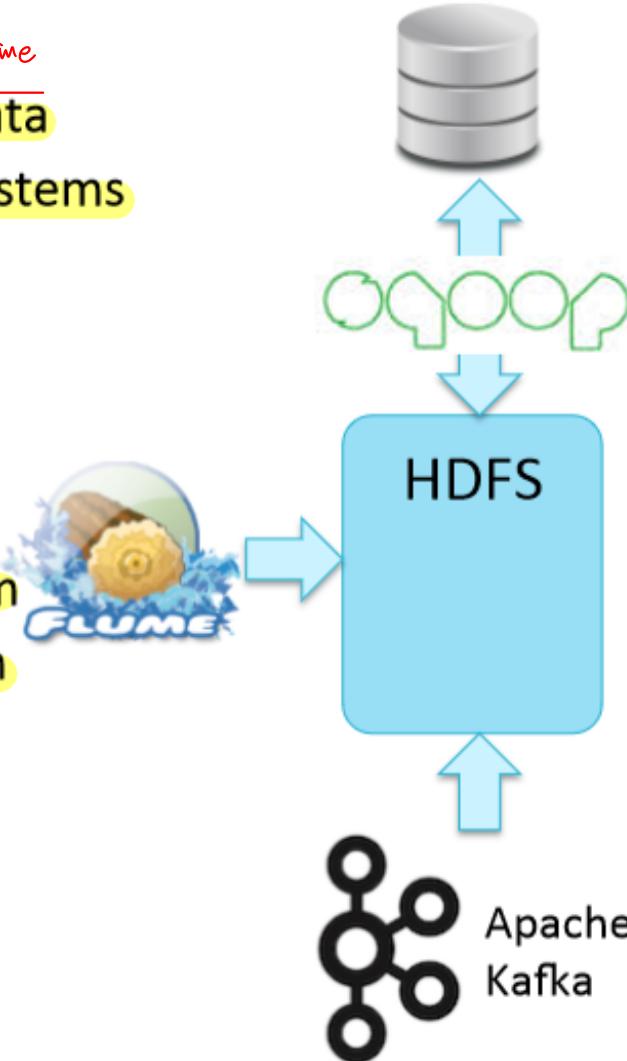
Kafka

- A high throughput, scalable messaging system
- Distributed, reliable publish-subscribe system
- Integrates with Flume and Spark Streaming

log aggregation, large scale messaging

Streaming processing

producer ?



Apache Spark: An Engine For Large-scale Data Processing

- **Spark is large-scale data processing engine**
 - General purpose
 - Runs on Hadoop clusters and data in HDFS
- **Supports a wide range of workloads**
 - Machine learning
 - Business intelligence
 - Streaming
 - Batch Processing
- **This course uses Spark for data processing**



Hadoop MapReduce: The Original Hadoop Processing Engine

- Hadoop MapReduce is the original Hadoop framework
 - Primarily Java based
- Based on the MapReduce programming model
- The core Hadoop processing engine before Spark was introduced
- Still the dominant technology
 - But losing ground to Spark fast
- Many existing tools are still built using MapReduce code
- Has extensive and mature fault tolerance built into the framework



Apache Pig: Scripting for MapReduce

- Apache Pig builds on Hadoop to offer high-level data processing
 - This is an alternative to writing low-level MapReduce code
 - Pig is especially good at joining and transforming data
- The Pig interpreter runs on the client machine
 - Turns Pig Latin scripts into MapReduce or Spark jobs
 - Submits those jobs to a Hadoop cluster
 - Covered in Cloudera *Data Analyst Training*



```
people = LOAD '/user/training/customers' AS (cust_id, name);
orders = LOAD '/user/training/orders' AS (ord_id, cust_id, cost);
groups = GROUP orders BY cust_id;
totals = FOREACH groups GENERATE group, SUM(orders.cost) AS t;
result = JOIN totals BY group, people BY cust_id;
DUMP result;
```

Impala: High Performance SQL

- Impala is a high-performance SQL engine
 - Runs on Hadoop clusters
 - Data stored in HDFS files
 - Inspired by Google's Dremel project
 - Very low latency – measured in milliseconds
 - Ideal for interactive analysis
- Impala supports a dialect of SQL (**Impala SQL**)
 - Data in HDFS modeled as database tables
- Impala was developed by Cloudera
 - 100% open source, released under the Apache software license
- Impala is used for data analysis in this course



Apache Hive: SQL on MapReduce

- **Hive is an abstraction layer on top of Hadoop**
 - Hive uses a SQL-like language called HiveQL
 - Similar to Impala SQL
 - Useful for data processing and ETL
 - Impala is preferred for ad hoc analytics
- **Hive executes queries using MapReduce**
 - Hive on Spark is available for early adopters; not yet recommended for production
- **Hive can optionally be used for data analysis in this course**



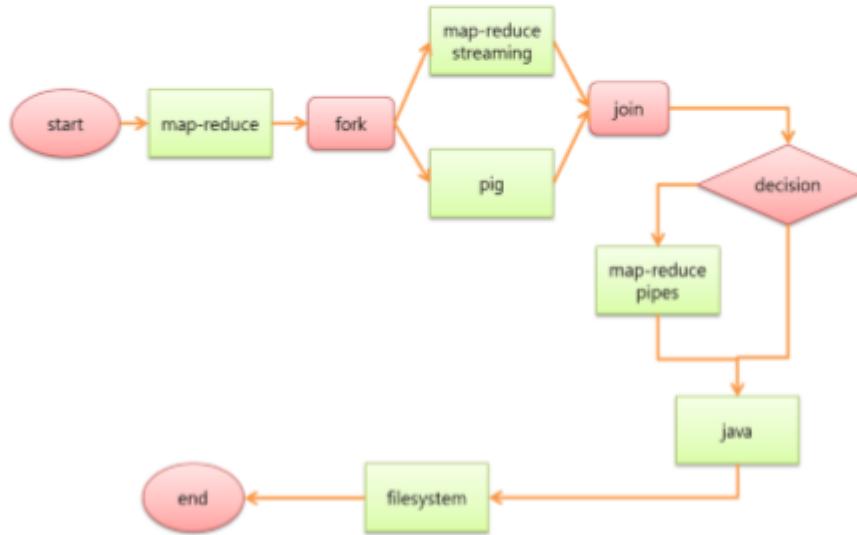
Hue: The UI for Hadoop

- Hue = **Hadoop User Experience**
- Hue provides a **Web front-end to a Hadoop**
 - Upload and browse data
 - Query tables in Impala and Hive
 - Run Spark and Pig jobs and workflows
 - Search
 - And much more
- Makes Hadoop easier to use
- Hue is 100% open-source
- Created by Cloudera
 - Open source, released under Apache license
- Hue is used throughout this course



Apache Oozie: Workflow Management

- **Oozie**
 - Workflow engine for Hadoop jobs
 - Defines dependencies between jobs
- **The Oozie server submits the jobs to the server in the correct sequence**



Apache Sentry: Hadoop Security

- Sentry provides fine-grained access control (authorization) to various Hadoop ecosystem components
 - Impala
 - Hive
 - Cloudera Search
 - HDFS
- In conjunction with Kerberos authentication, Sentry authorization provides a complete cluster security solution
- Created by Cloudera
 - Now an open-source Apache project



Summary

- Hadoop is a framework for distributed storage and processing
- Core Hadoop includes HDFS for storage and YARN for cluster resource management
- The Hadoop ecosystem includes many components for
 - Ingesting data (Flume, Sqoop, Kafka)
 - Storing data (HDFS, HBase)
 - Processing data (Spark, Hadoop MapReduce, Pig)
 - Modeling data as tables for SQL access (Impala, Hive)
 - Exploring data (Hue, Search)
 - Protecting Data (Sentry)
- This course introduces most of the key Hadoop infrastructure

데이터 수집 계층

Flume : 로그 데이터 스트리밍 수집.

Sqoop : RDB \leftrightarrow Hadoop 데이터 이관

Kafka : 분산 메시지 큐. 실시간 데이터 스트리밍 처리 학습.
중간 레이어 data 인정적 전달. Spark Streaming 등이 예시.

저장 계층

HDFS : 대용량 데이터 분산 저장.

HBase : HDFS 기반 NoSQL DB

처리 계층.

Map Reduce : batch 기반 기본 분산 처리 엔진
Map \rightarrow Shuffle \rightarrow Reduce.

Spark : 메모리 기반 분산 처리 엔진. MapReduce보다 빠름.
batch, Streaming, ML, SQL 등 범용 지원
Kafka와 결합. \rightarrow 실시간 분석

Pig : Map Reduce 추상화한 스크립트 언어. (Pig Latin)
스크립트 작성시 배부자로 Map Reduce 실행

분석/쿼리 계층

Hive : HDFS 데이터에 대해 HiveQL 쿼리 가능

Impala : Hive와 유사. 실시간/대화형 쿼리에 특화.

MapReduce 대신 자체 엔진으로 빠른 응답

UI, 검색

Hue : 웹기반 GUI. Hive, Impala, HDFS 검색. Workflow 관리

Search : HDFS/HBase 데이터를 검색 엔진처럼 indexing. 빠르게 조회

보안

Sentry : Hadoop 권한관리 framework



Welcome and Enjoy!