



London
Business
School

AM05 Data Management
05. Database Connections and Big Data
Dr. David Tilson

Agenda

Database connections and semi-structured data

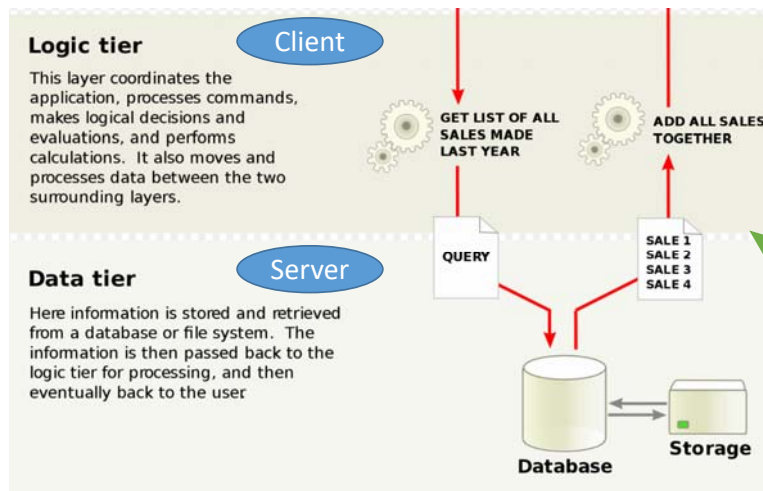
Business analytics and the advent of big data

Technical challenges of big data

Intro to Hadoop and other **Big Data** technologies

Wrapping up the course

Database connection allows client software to talk to database server software (on same machine or over a network)



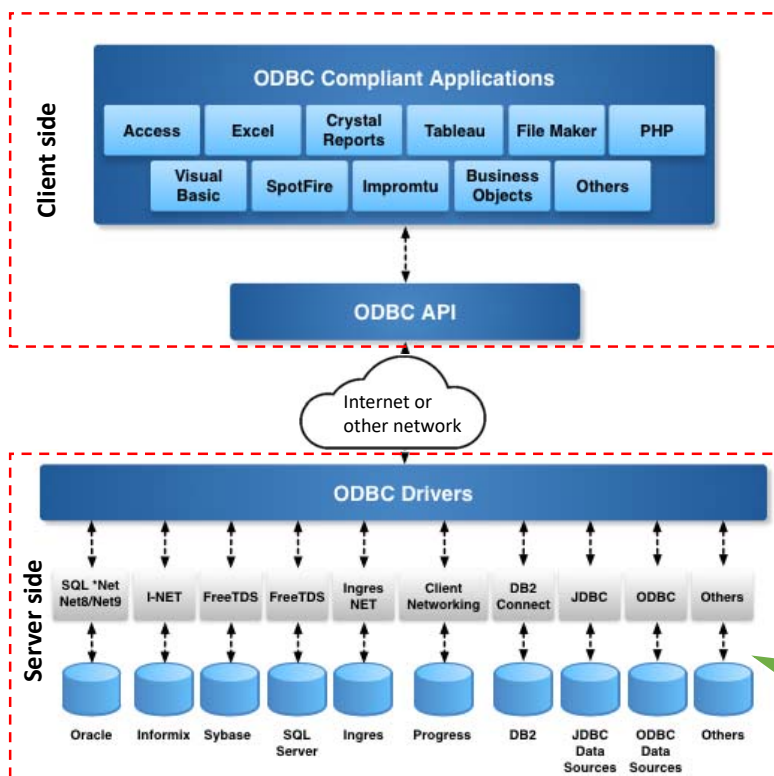
A connection is required to send commands and receive answers, usually in the form of a result set.

Commands can only be performed against a database with an "open and available" connection to it

There are many way of doing this (see chapter 5 of DSC). Some more efficient / secure than others

261

Open Database Connectivity (ODBC) is a standard application programming interface (API) for accessing DBMSs



- **ODBC** aims to be independent of database and operating systems
- ODBC-based applications can be ported to other platforms with few changes to the data access code

Servers could be

- On same machine as client
- On different machine in data center
- In a different corporate data center
- On a cloud provider data service

262

In practice the ODBC API is available as a library on the client side (as are other types of database connectors)

pyodbc

pyodbc is an open source Python module that makes accessing ODBC databases simple. It implements the [DB API 2.0](#) specification but is packed with even more Pythonic convenience.

The easiest way to install is to use pip:

```
pip install pyodbc
```

<http://mkleehammer.github.io/pyodbc/>

Installers

conda install 

<https://anaconda.org/anaconda/pyodbc>

To install this package with conda run:

```
conda install -c anaconda pyodbc
```

Using an ODBC driver

<https://db.rstudio.com/odbc/>

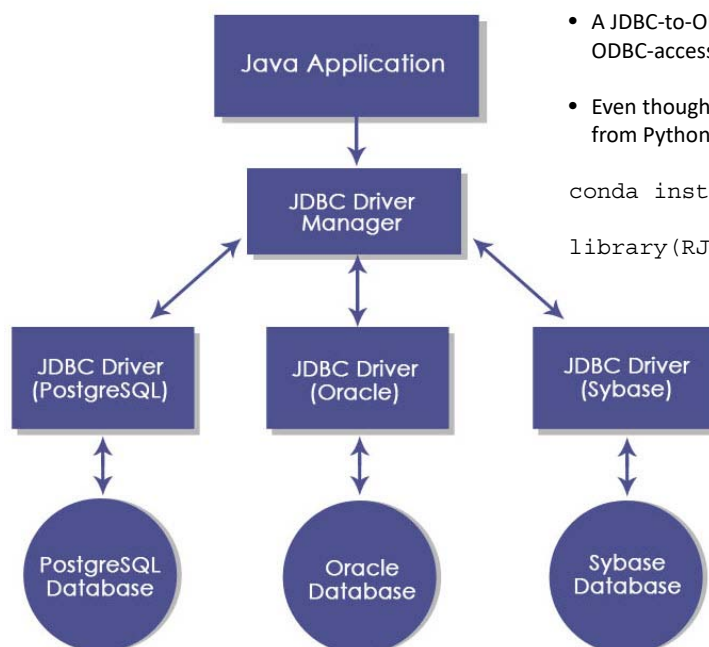
The `odbc` package provides a DBI-compliant interface to [Open Database Connectivity](#) (ODBC) drivers. It allows for an efficient, easy way to setup connection to any database using an ODBC driver, including [SQL Server](#), [Oracle](#), [MySQL](#), [PostgreSQL](#), [SQLite](#) and others. The implementation builds on the [nanodbc](#) C++ library.



ODBC drivers can typically be downloaded from your database vendor, or they can be [downloaded from RStudio](#) when used with RStudio professional products. The `odbc` package works with the DBI

263

Java Database Connectivity (JDBC) is another common API which defines how clients can access RDBMs



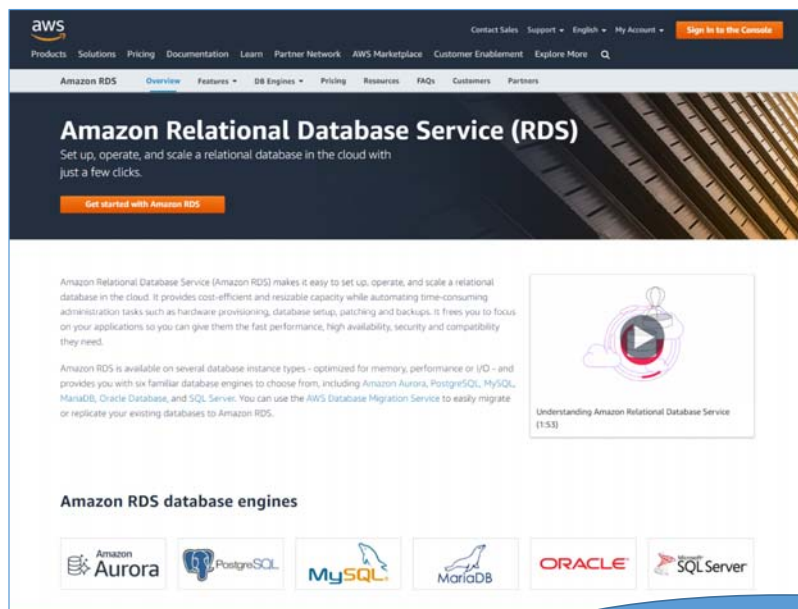
- A JDBC-to-ODBC bridge also enables connections to any ODBC-accessible data source
- Even though JDBC is a Java technology it can be used from Python e.g. to install from anaconda

```
conda install -c conda-forge jaydebeapi
```

library(RJDBC) in R (see rforge.net/RJDBC)

264

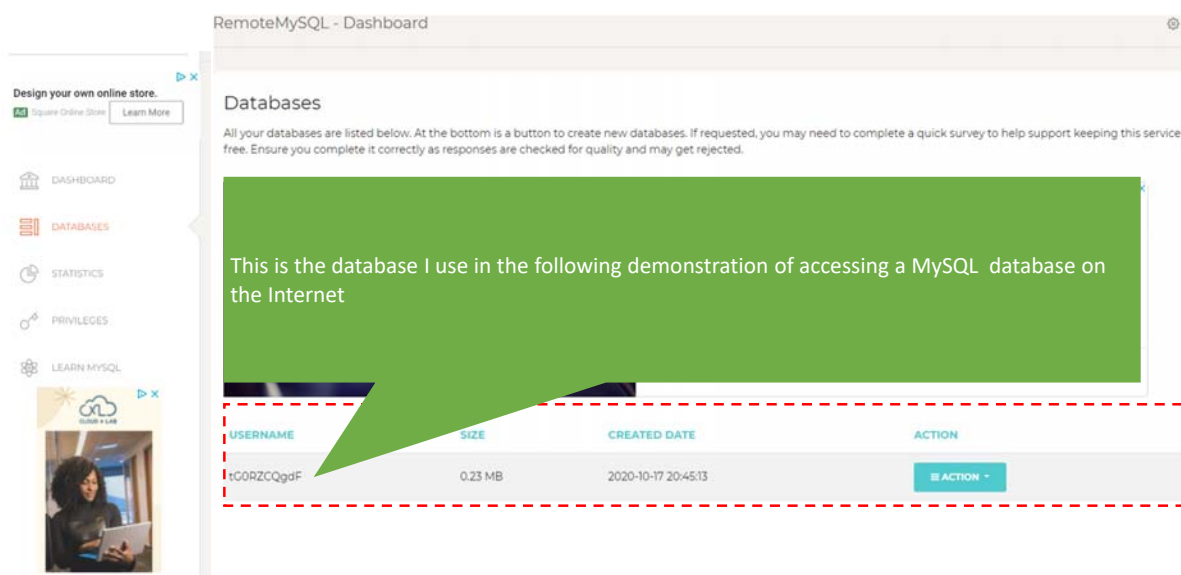
The major cloud infrastructure providers allow you to set up RDMSs in minutes (e.g. on Amazon Web Services)



Google Cloud, Microsoft Azure and many others have similar offerings

265

You can even set up free MySQL databases on-line to play with or test software (but you will get lots of adverts)



<https://remotemysql.com>

Obviously, don't use it for anything important

266

Let's look at example of connecting to an MySQL database over the Internet using python

Programmatic Access to On-Line MySQL Database.ipynb

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

RAM 1.0 GB Disk 1.0 GB Editing

Demonstration of programmatic database access

This simple notebook shows how to establish a connection with a MySQL database created located somewhere on the Internet.

- It is partly based on the code from this tutorial <https://dev.mysql.com/doc/connector-python/en/connector-python-examples.html>
- It also uses the example ap database we used in class
- I am running it from Google Colaboratory <https://colab.research.google.com/notebooks/welcome.ipynb>. Some changes might be required if you are running it in your own environment.

```
[1] # If you are running this for the first time you may have to install the mysql-connector library
# If you are trying this on your own machine you will probably have to install it from the command line
!pip install mysql-connector-python

Collecting mysql-connector-python
  Downloading https://files.pythonhosted.org/packages/6c/1d/e666f7d43496a2315d3963a2fb7f8df84e7793b4ddbf05e46d6db4a8892/mysql_connector_python-8.0.22-cp36-cp36m-macosx_10_10_x86_64.whl (10.0MB)
Requirement already satisfied: protobuf>=3.0.0 in /usr/local/lib/python3.6/dist-packages (from mysql-connector-python) (3.12.4)
Requirement already satisfied: six>=1.9 in /usr/local/lib/python3.6/dist-packages (from protobuf>=3.0.0->mysql-connector-python) (1.15.0)
Requirement already satisfied: setuptools in /usr/local/lib/python3.6/dist-packages (from protobuf>=3.0.0->mysql-connector-python) (50.3.0)
Installing collected packages: mysql-connector-python
Successfully installed mysql-connector-python-8.0.22
```

```
[2] # Import some libraries
import mysql.connector
from mysql.connector import errorcode
import datetime
```

Jupyter notebook is on Canvas

267

JSON and XML are common ways of communicating semi-structured data across networks

JavaScript Object Notation (JSON) Example

```
{"products": [
  {"number": 1, "name": "Zoom X", "Price": 10.00},
  {"number": 2, "name": "Wheel Z", "Price": 7.50},
  {"number": 3, "name": "Spring 10", "Price": 12.75}
]}
```

eXtensible Markup Language (XML) Example

```
<products>
  <product>
    <number>1</number> <name>Zoom X</name> <price>10.00</price>
  </product>
  <product>
    <number>2</number> <name>Wheel Z</name> <price>7.50</price>
  </product>
  <product>
    <number>3</number> <name>Spring 10</name> <price>12.75</price>
  </product>
</products>
```

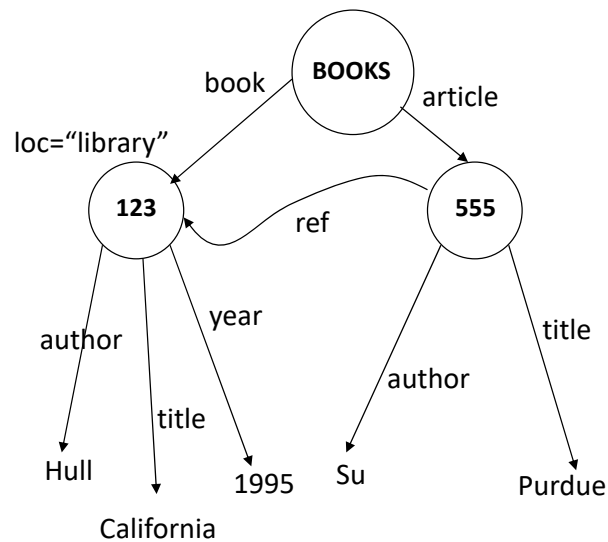
Both JSON and XML provide a **semi-structured** means of representing data hierarchically

Semi-structured data does not obey the formal structure of data models associated with RDBMSs but contains tags to separate semantic elements and enforce hierarchies of records and fields within (aka a self-describing structure)

XML documents define sets of rules for encoding documents in human and machine-readable format

Example of XML document elements

```
<BOOKS>
<book id="123" loc="library">
  <author>Hull</author>
  <title>California</title>
  <year> 1995 </year>
</book>
<article id="555" ref="123">
  <author>Su</author>
  <title> Purdue</title>
</article>
</BOOKS>
```



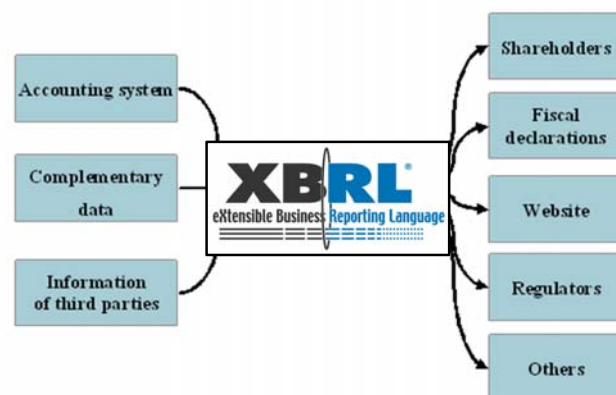
Hundreds of XML-based document formats have been defined e.g. RSS, SOAP for web services, and SOA. Even modern Microsoft Office file formats are XML-based

269

XBRL is an XML-based standard developed for the exchange of financial reports over networks

XBRL (**eXtensible Business Reporting Language**) is a framework for exchanging business information

- XML-based: uses XML syntax/technologies
- Supports business reporting like the definition and exchange of financial information (e.g. financial statements)



Individual fields or whole XML documents can be stored in RDBMSs or other data stores

270

JSON is used extensively on the Web (e.g. Twitter feeds)

Example tweet in JSON format

```
[{"created_at": "Thu Jun 22 21:00:00 +0000 2017",
  "id": 877994604561387500,
  "id_str": "877994604561387520",
  "text": "Creating a Grocery List Manager Using Angular, Part 1: Add & Display",
  "truncated": false,
  "entities": {
    "hashtags": [
      {
        "text": "Angular",
        "indices": [103, 111]
      }
    ],
    "symbols": [],
    "user_mentions": [],
    "urls": [
      {
        "url": "https://t.co/xFox78juL1",
        "expanded_url": "http://buff.ly/2sr60pf",
        "display_url": "buff.ly/2sr60pf",
        "indices": [79, 102]
      }
    ]
  },
  "source": "<a href='\"http://bufferapp.com\"' rel='\"nofollow\"'>@Buffer</a>",
  "user": {
    "id": 772682964,
    "id_str": "772682964",
    "name": "SitePoint JavaScript",
    "screen_name": "SitePointJS",
    "location": "Melbourne, Australia",
    "description": "Keep up with JavaScript tutorials, tips, tricks and articles",
    "url": "http://t.co/cCH13gqeUK",
    "entities": {
      "url": "http://t.co/cCH13gqeUK",
      "entities": {
        "url": {
          "url": "http://t.co/cCH13gqeUK",
          "expanded_url": "http://sitepoint.com/javascript",
          "display_url": "sitepoint.com/javascript",
          "indices": [0, 22]
        }
      },
      "description": {
        "urls": []
      }
    },
    "protected": false,
    "followers_count": 2145,
    "friends_count": 18,
    "listed_count": 328,
    "created_at": "Wed Aug 22 02:06:33 +0000 2012",
    "favourites_count": 57,
    "utc_offset": 43200,
    "time_zone": "Wellington"
  }
}]
```

Individual tweet may have different sets of tags

MySQL database has JSON datatype
Popular MongoDB NoSQL database natively stores JSON documents

See <https://developer.twitter.com/en/docs/tweets/data-dictionary/overview/intro-to-tweet-json>

271

Agenda

Database connections and semi-structured data

Business analytics and the advent of big data

Technical challenges of big data

Intro to Hadoop and other **Big Data** technologies

Wrapping up the course

272

Definitions of big data are a bit squishy

Some definitions

1. “Big data exceeds the reach of commonly used hardware environments and software tools to capture, manage, and process in within a tolerable elapsed time for its user population.”
Gartner’s Merv Adrian in Teradata Magazine
2. “Big data refers to data sets whose size is beyond the ability of typical database software tools to capture, store, manage and analyze.”
McKinsey Global Institute, Big Data: The Next Frontier for Innovation, Competition and Productivity



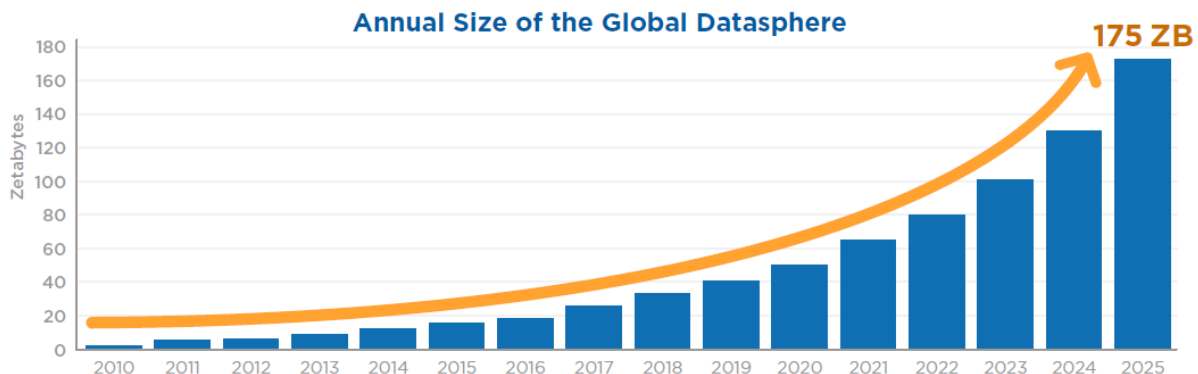
Implies

- What qualifies as “big data” changes as technology advances
- Tomorrow’s “big” will be bigger than today’s

273

Big. . . and getting bigger

Figure 1 – Annual Size of the Global Datasphere



It's more than just lots of data . . . The 3Vs of “Big Data”

Three drivers of increased complexity

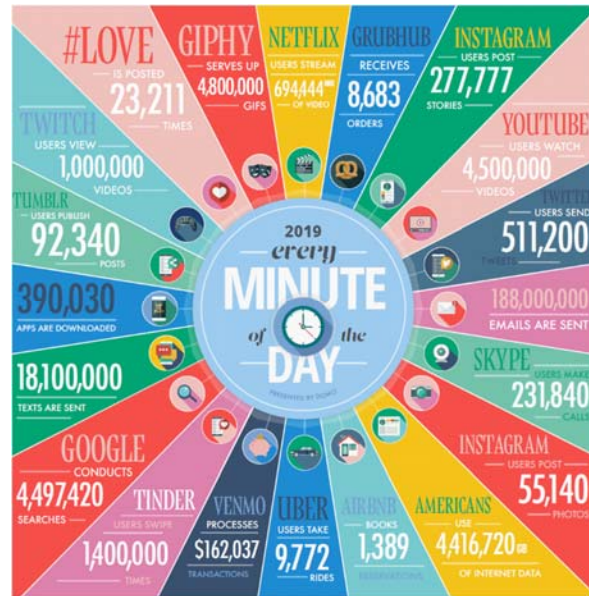
¶ **Volume** – a lot of data to work with

¶ **Velocity** – its coming at you fast,

- Not just periodic reports on structured data
- Continually arriving data needs to be processed for insight
- Moving to real-time analysis

¶ **Variety** – many different sources

- Structured: databases, sensor data
- Semi-structured: Weblogs, Social
- Unstructured: flat-file, images, video, audio



Data source: <https://www.newgenapps.com/blog/big-data-statistics-predictions-on-the-future-of-big-data>

Gartner identified the 3Vs – others added Value, Veracity, Vulnerability (security), Variability, Volatility, Visualization

275

Volume is a moving target . . . but petabytes are certainly big data

Quantities of bytes				
SI prefixes			Binary prefixes (IEC 60027-2)	
Name (Symbol)	Value in Popular Usage	Value in Standard SI	Name (Symbol)	Value
kilobyte (kB)	2 ¹⁰	10 ³	kibibyte (KiB)	2 ¹⁰
megabyte (MB)	2 ²⁰	10 ⁶	mebibyte (MiB)	2 ²⁰
gigabyte (GB)	2 ³⁰	10 ⁹	gibibyte (GiB)	2 ³⁰
terabyte (TB)	2 ⁴⁰	10 ¹²	tebibyte (TiB)	2 ⁴⁰
petabyte (PB)	2 ⁵⁰	10 ¹⁵	pebibyte (PiB)	2 ⁵⁰
exabyte (EB)	2 ⁶⁰	10 ¹⁸	exbibyte (EiB)	2 ⁶⁰
zettabyte (ZB)	2 ⁷⁰	10 ²¹	zebibyte (ZiB)	2 ⁷⁰
yottabyte (YB)	2 ⁸⁰	10 ²⁴	yobibyte (YiB)	2 ⁸⁰

Big data today

¶ Traditional databases can handle tens of terabytes

¶ PB definitely starting to be big data

¶ Velocity – rate at which it is growing means you

- Might soon need big data technologies for large volume
- Need big data technologies to ingest the daily volume

¶ Variety makes things more complex

- Beyond rows and columns
- Messages (text), images, video

276

There are some features that tend to distinguish “big data” sources from traditional ones

¶ Often generated automatically

- Traditional data sources have people taking action: call detail records, bank/retail transactions, shipments, payments
- Many “big data” sources don’t involve people: smart meters, engine sensors

¶ Novel data sources

- Traditional data sources are often transactions, even if done on-line
- New sources include detailed browsing behaviors
- “More of the same” can be new: Smart meter readings every 15 minutes

¶ Not designed to be friendly (i.e. easy to process and analyze)

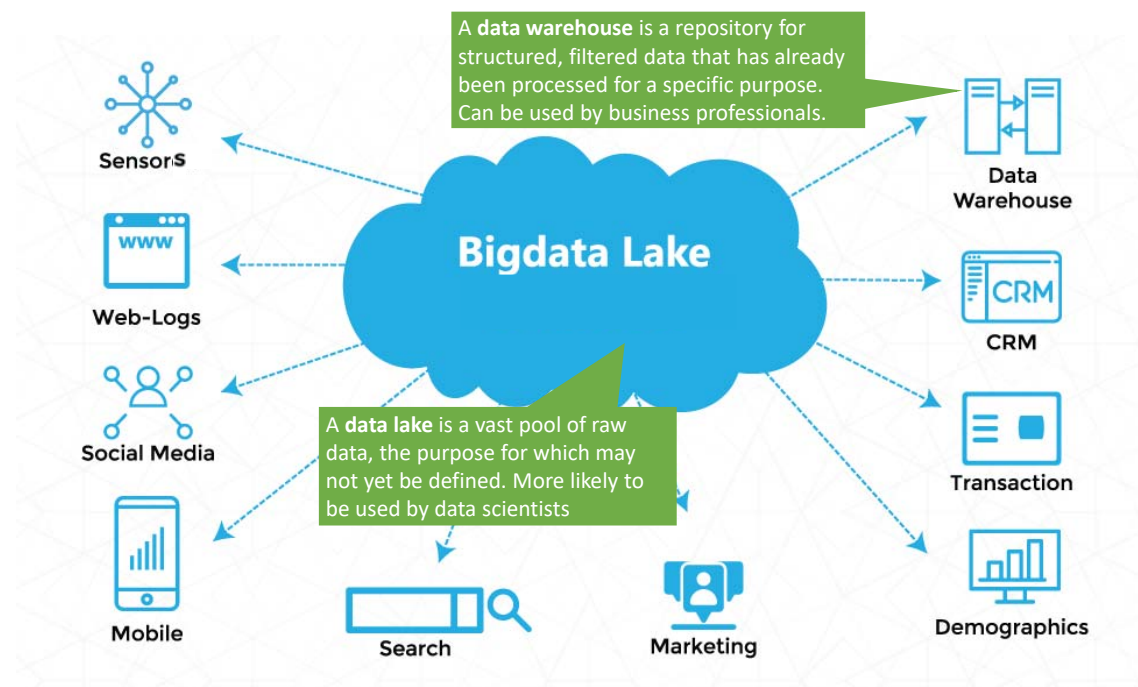
- Traditional sources designed to be friendly (database transaction records)
- Some not designed at all: text streams from social media
 - No standards of grammar, sentence ordering, or vocabulary
 - Get what you get
- May have to wade through mess, junk filled data during analysis

¶ Much of it may be worthless

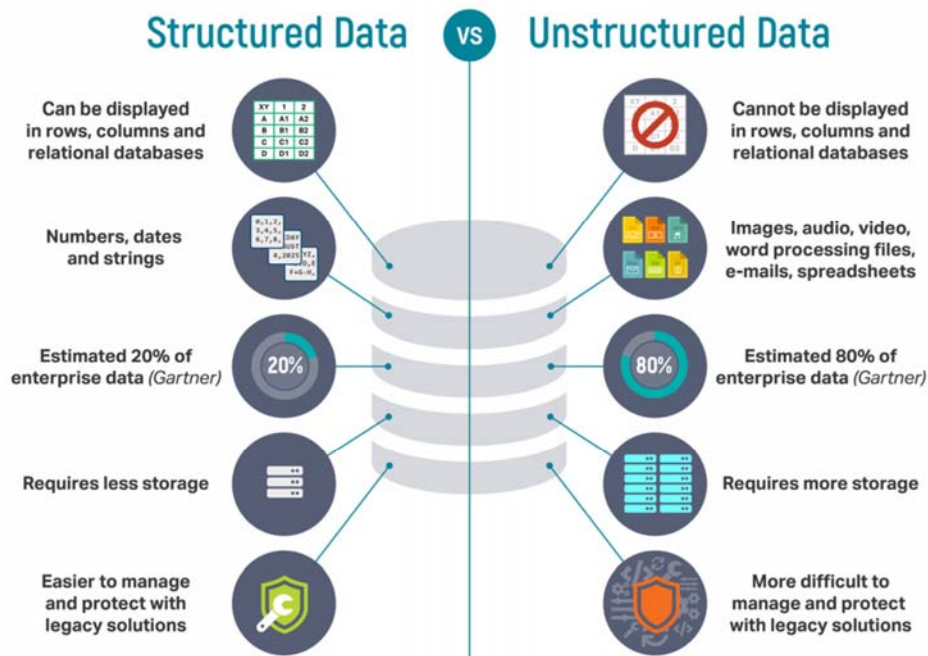
- Traditional sources designed to be 100% relevant and make best use of limited storage and processing power
- Big data tends to capture everything possible and figure out what matters later

277

In a business context Big Data often associated with predictive analytics and user/customer behavior analysis



Much of the increase in volume is tied to increases in variety of data... structured, unstructured (text, logs, images, video, sensors)



<https://www.igneous.io/blog/structured-data-vs-unstructured-data>

279

Web log data is big(gest) source of “Big Data”

Web logs as rich new data source

¶ Detailed web behavior insight

- Research behaviors
- Decisions making
- Purchase paths

¶ Factual info on customer

- Preferences
- Future intentions
- Motivations

¶ Info across browsing sessions

¶ New data for segmentation (e.g. dreamers)

- Less than 2% of website visits result in purchases
- So, only tracking completed transactions means gaining no insight on >98% of visitors



But web logs are not “ready to analyze”

```
4 #Fields: date time s-sitename s-computername s-ip cs-method cs-uri-stem cs-uri-query s-port cs
5 2010-03-24 07:00:01 ZZZZC941948879 RUFFLES 222.222.222.222 GET / - 80 - 220.181.7.113 HTTP/1.1
6 2010-03-24 07:00:23 ZZZZC941948879 RUFFLES 222.222.222.222 GET /2009/12/im_not_mean_im_just
7 2010-03-24 07:00:32 ZZZZC941948879 RUFFLES 222.222.222.222 GET /terminal-blank.gif - 80 -
8 2010-03-24 07:00:32 ZZZZC941948879 RUFFLES 222.222.222.222 GET /grep-options.gif - 80 -
9 2010-03-24 07:00:32 ZZZZC941948879 RUFFLES 222.222.222.222 GET /terminal-cat.gif - 80 - 217.
10 2010-03-24 07:00:32 ZZZZC941948879 RUFFLES 222.222.222.222 GET /terminal-pwd-cd.gif - 80 - 217
11 2010-03-24 07:00:39 ZZZZC941948879 RUFFLES 222.222.222.222 GET /robots.txt - 80 - 95.55.207.95
12 2010-03-24 07:00:39 ZZZZC941948879 RUFFLES 222.222.222.222 GET /rss-chart.xml - 80 - 173.45.21
```

280

Unstructured text is massive source of potential insight



Natural Language
Text Parsing



Meaning

¶ Lots of sources

- Email
- Text messages
- Tweets
- Social media posts
- Instant messaging / chats
- Audio transcriptions (e.g. support calls)

¶ Potential insights

- “Buzz” around product
- Customer sentiment about company / product / service
- Sources of complaints
- Fraud detection
- Legal discovery
- Targeted adverts

281

There are other novel sources of “Big Data”



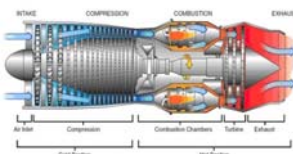
Location information



RFIDs



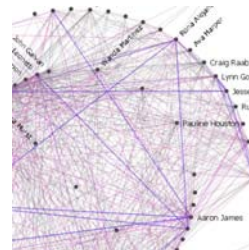
Smart grid / meters



Sensor data



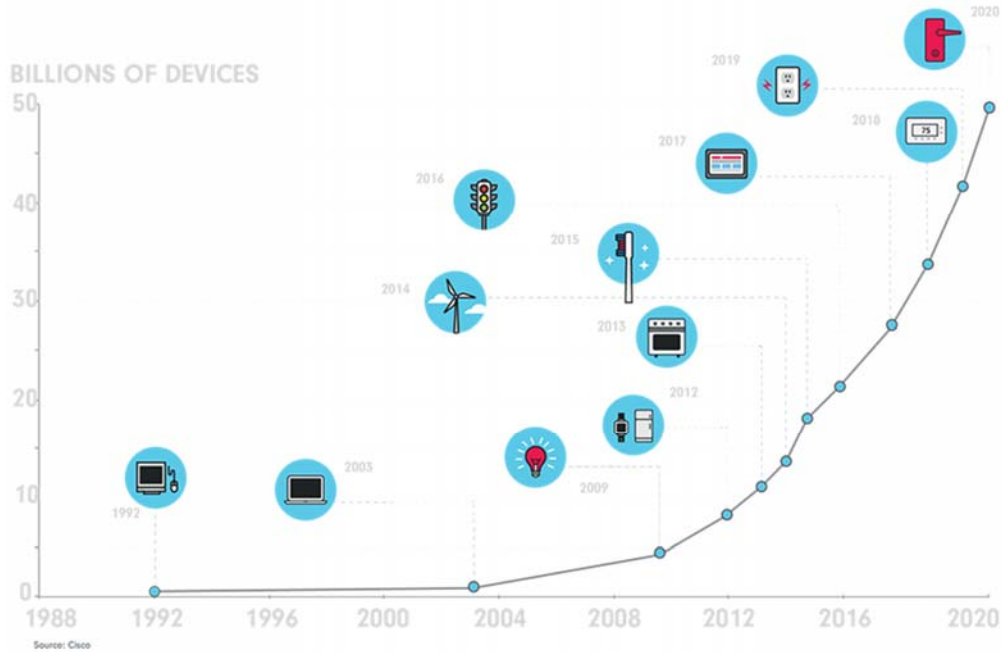
Sensor data



(Social) Network graphs

282

Data from Internet of Things (IoT) growing exponentially



283

Agenda

Database connections and semi-structured data

Business analytics and the advent of big data

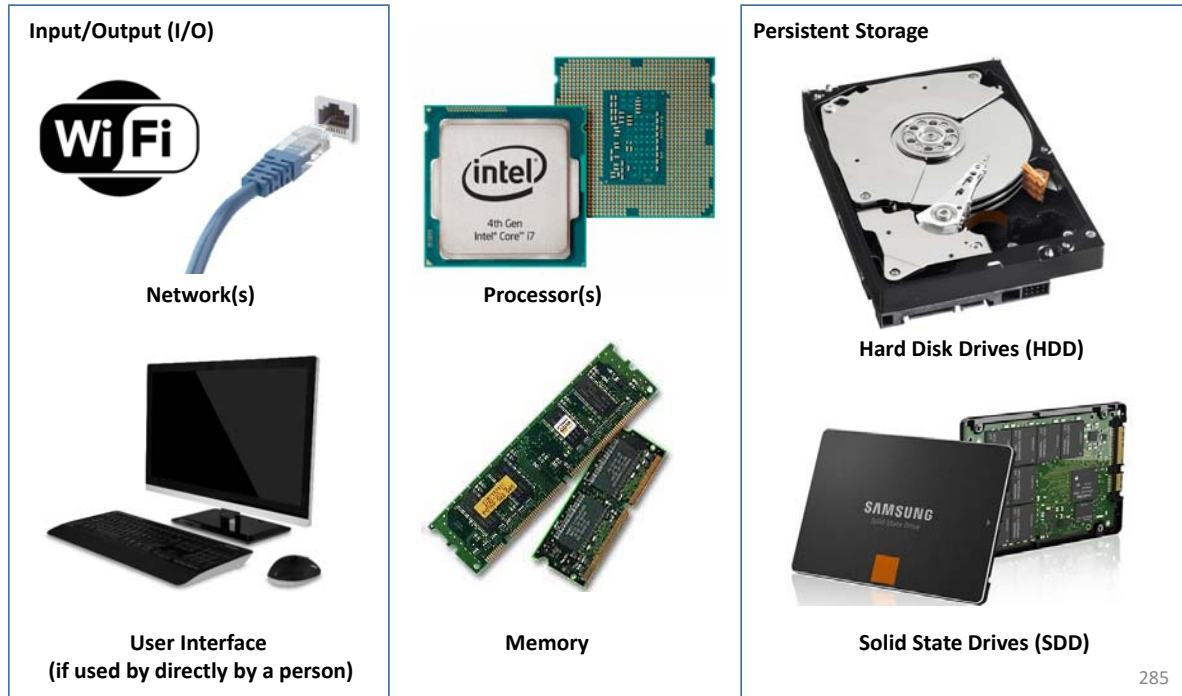
Technical challenges of big data

Intro to Hadoop and other **Big Data** technologies

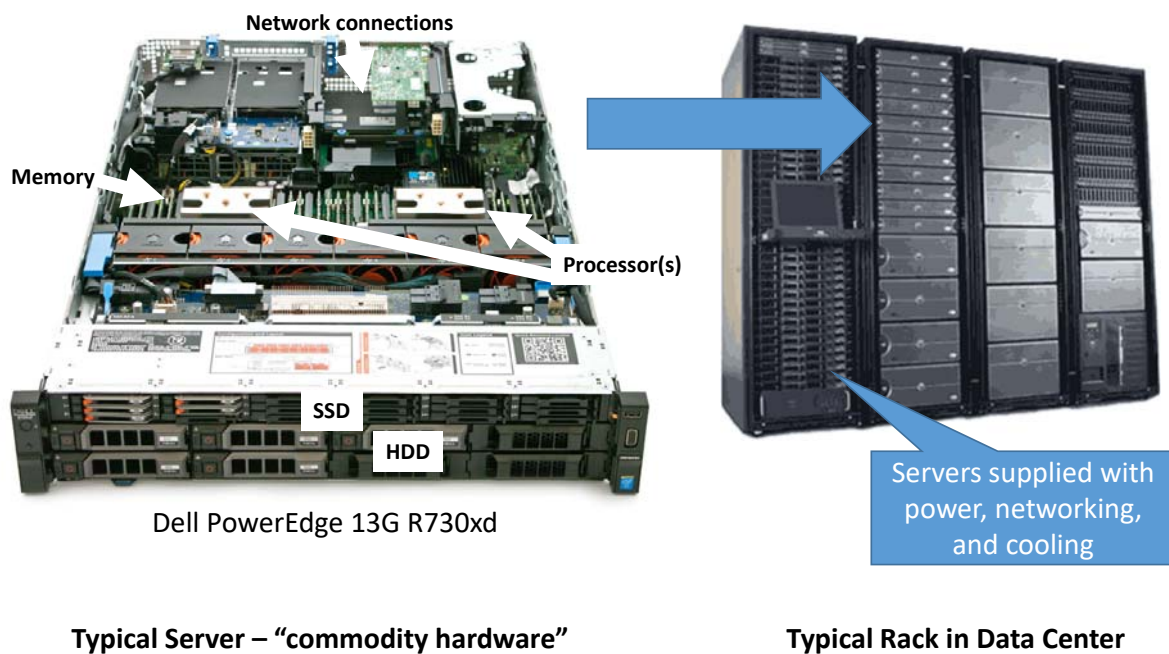
Wrapping up the course

284

Digital computers have several key components coordinated by software

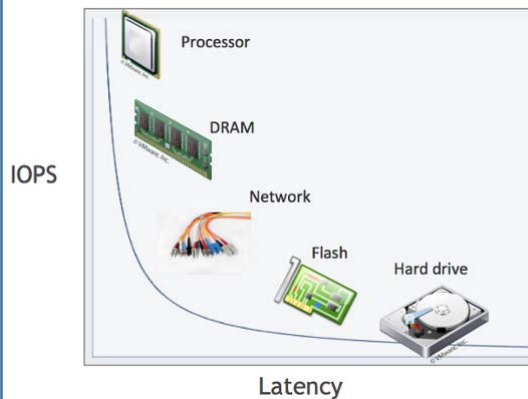


Enterprise IT and Cloud services supported by servers installed in racks



When big data is spread across servers it becomes impractical to bring it all to one machine for processing

Components of computers operate at very different speeds



If Memory = **Minute**
Network = **Weeks***
Flash = **Months**
Disk = **Decades**

* Still want to avoid moving data over networks as it still has to transfer on/off disks

Getting data to processors becomes the bottleneck

Impossible for data to go to the code

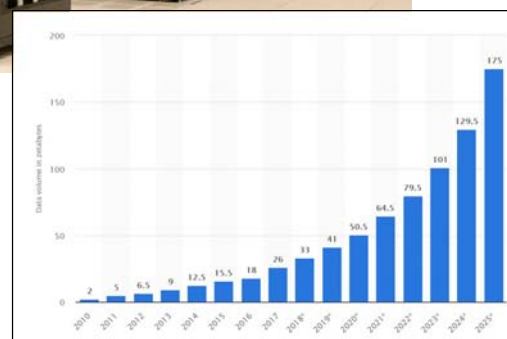


So the code must go to the data

287

Big data does not fit on one computer

Yahoo's Hadoop Cluster



288

Distributed Systems have certain problems

- ¶ Can use multiple machines for single task
- ¶ Programming distributed systems much more complex
 - Synchronizing data exchanges
 - Managing a finite bandwidth
 - Controlling computation timing
- ¶ Vulnerable to failure of one machine



289

Agenda

Database connections and semi-structured data

Business analytics and the advent of big data

Technical challenges of big data

Intro to Hadoop and other **Big Data** technologies

Wrapping up the course

290

What is Apache Hadoop?



¶ Open source software framework designed for storage and processing of large-scale data on clusters of “commodity hardware”

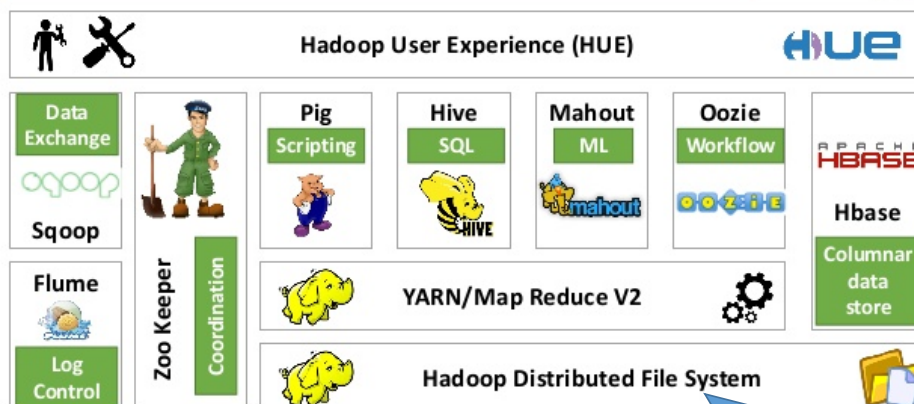
¶ Used for

- Data-intensive text processing
- Assembly of large genomes
- Graph mining
- Machine learning and data mining
- Large scale social network analysis

291

Hadoop designed to alleviate the problems with distributed computing

The Apache Hadoop Stack



¶ Must support partial failure

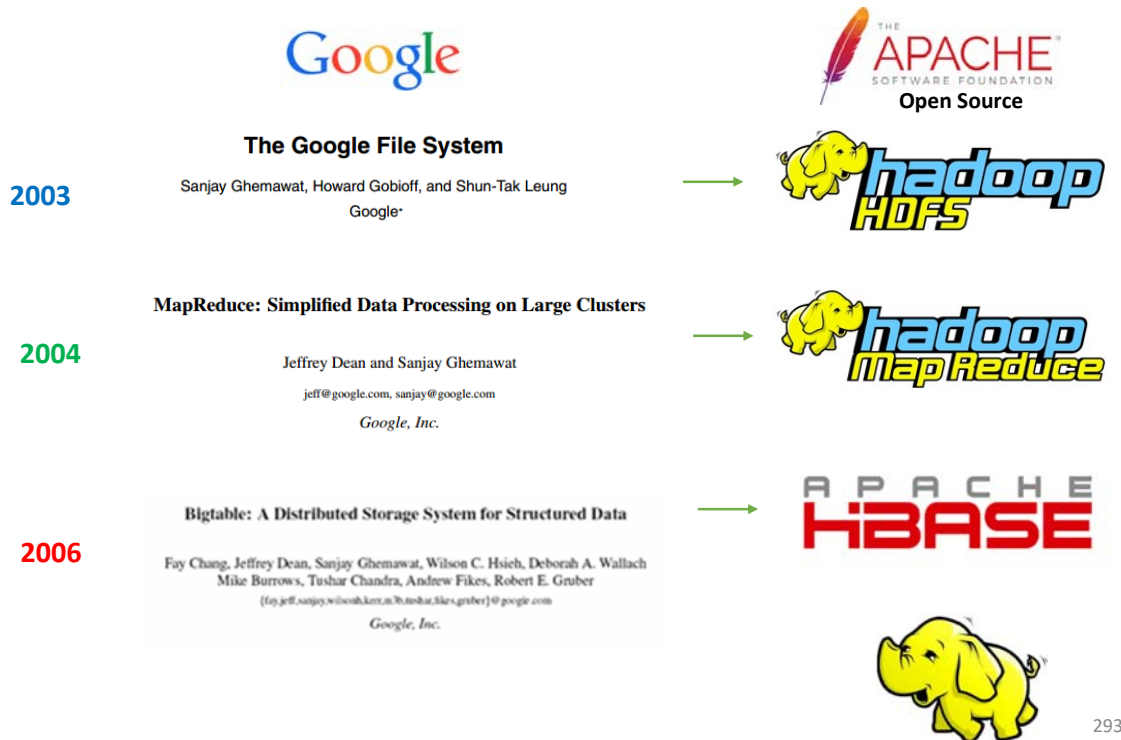
¶ Must be scalable

- Increasing resources should increase load capacity
- Increasing load on the system should result in graceful decline in performance for all jobs . . . rather than system failure

HDFS file system hides the fact that data is stored across a cluster of machines from applications

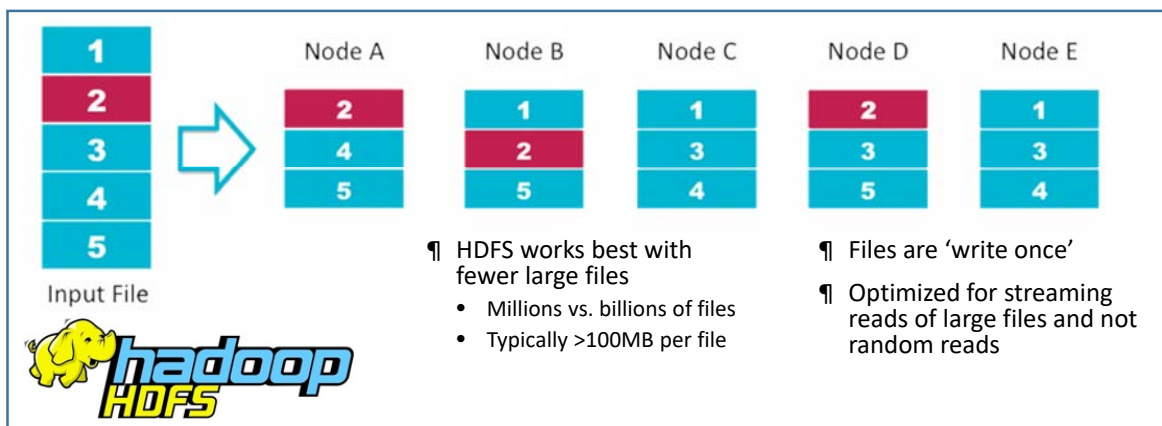
292

Ideas behind Apache Hadoop emerged from Google... one of the first companies to grapple with big data

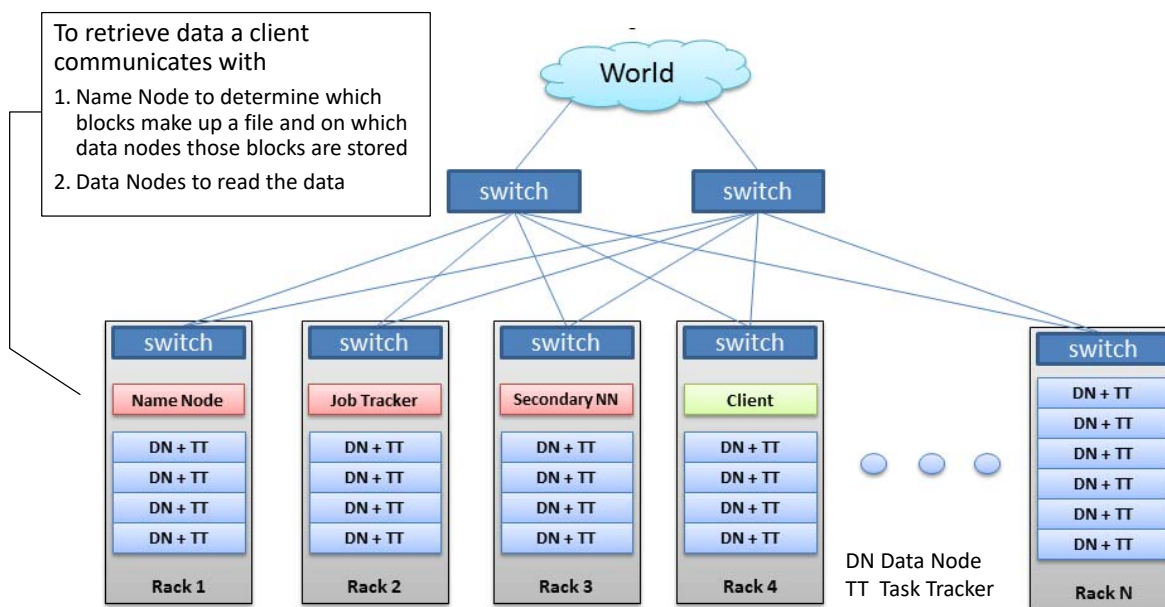


Hadoop Distributed File System (HDFS) provides redundant storage for massive amounts of data

- ¶ Files split into blocks (64 or 128 MB)
 - Split across many machines at load time
 - Replicated across multiple machines (for fault tolerance)
- ¶ **NameNode** keeps track of which blocks make up a file and where they are stored



Hadoop clusters can scale to thousands of nodes

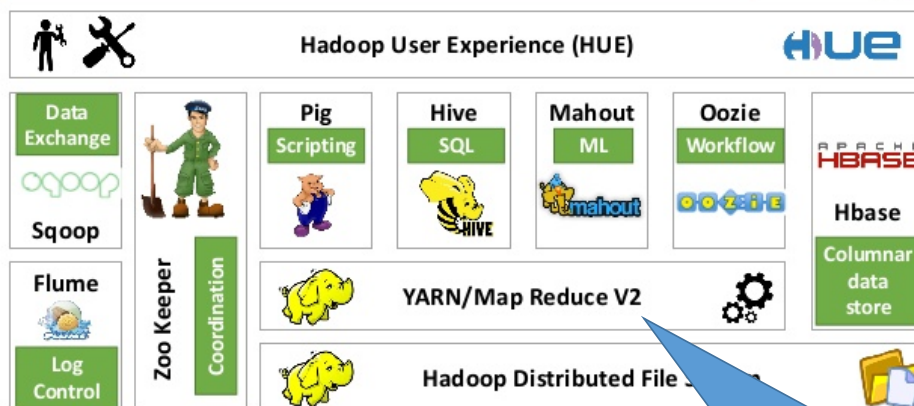


This approach to increasing processing power, memory, and storage is referred to as **scaling-out** ... as opposed to scaling-up by buying or building a single more capable server

295

MapReduce is Hadoop's distributed computing platform

The Apache Hadoop Stack



MapReduce Features

- ¶ Automatic parallelization and distribution of computing across nodes in the cluster
- ¶ Fault-Tolerance
- ¶ Clean abstraction for programmers

Developers \ Analysts can write Map Reduce jobs (e.g. in Java or Python) without worrying about details of the cluster

296

Core concepts in MapReduce distributed computing model

¶ Code goes to the data (parallel computing)

- Each node can perform computations on data it stores
- No need to move data – at least for initial (map) processing
- Tries to minimize slow data transfers between nodes and racks

¶ Applications written in high-level programming language

- No network programming needed
- Handles timing issues
- Fault tolerance built in
 - Failures detected and tasks reassigned to a different nodes
 - Restarting a task does not affect nodes working on other tasks
 - If failed node restarts, it is added back to cluster and assigned new tasks

¶ Need to be a (Java) programmer to use MapReduce directly... not really for data analysts or data scientists

297

MapReduce is a relatively simple but flexible model



¶ Tasks divided into two phases

- Map tasks done on portions of data where it is stored
- Reduce tasks combine data from map tasks to produce final output

¶ Job Tracker allocates work to individual nodes

¶ Key functionality of SQL and other tasks can be implemented in MapReduce

- Select \ filter
- GroupBy
- Joins
- Linear Algebra (matrix math)

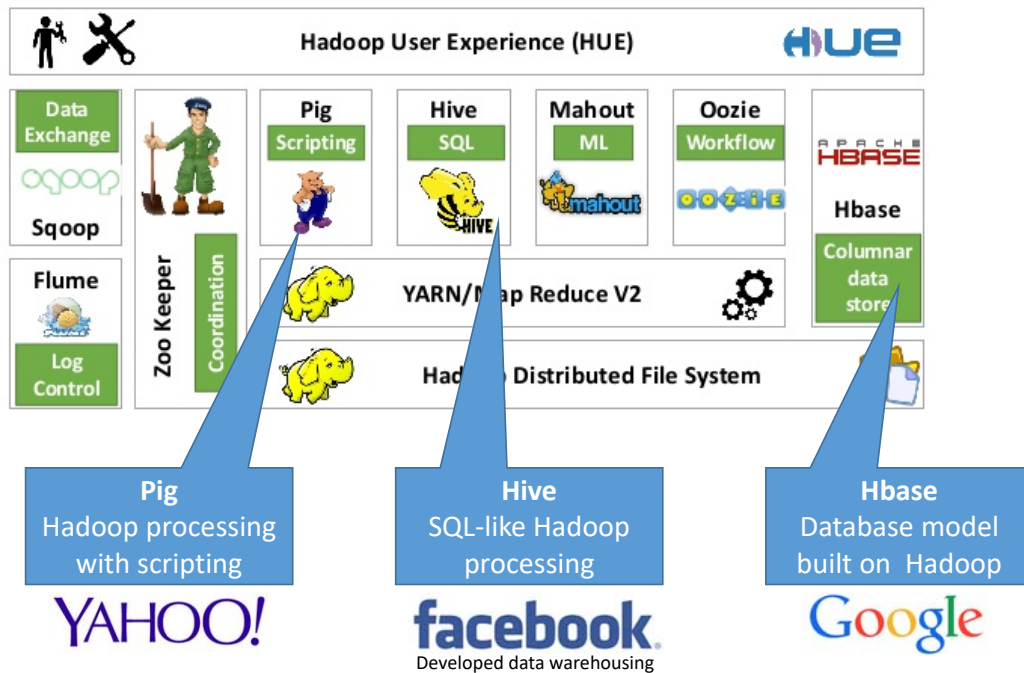
Map Reduce speeds up Big Data processing in two ways

1. Minimizing movement of data across networks
2. Using the power of many processors across the cluster (parallel computing)

298

Other tools in Hadoop ecosystem allow analysts to be more productive while harnessing MapReduce and HDFS

The Apache Hadoop Stack



299

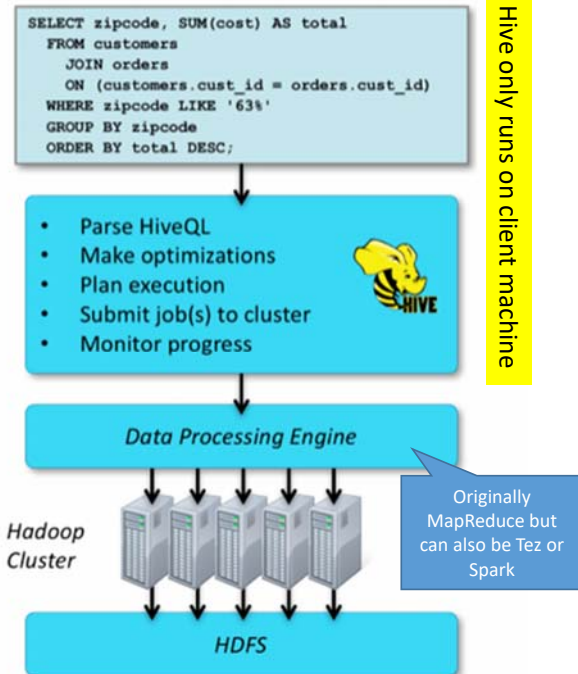
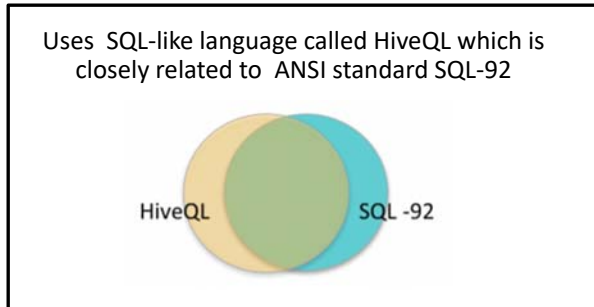
Many large companies rely on Hadoop



Apache Hive allows business analysts with SQL expertise gain insight from big data on Hadoop clusters

Hive runs on client machine

- ¶ Turns HiveQL queries into MapReduce jobs
- ¶ Submits those jobs to the cluster
 - As execution plan (not Java code)
- ¶ Shares many architectural similarities with Pig



301

Hive brings 'big data' to a broader range of potential users

- ¶ More productive than writing MapReduce directly
 - Five lines of HiveQL might be equivalent to 100 lines or more of Java
- ¶ Brings large-scale data analysis to a broader audience
 - No software development (Java) experience required
 - Leverages existing knowledge of SQL
- ¶ Offers interoperability with other systems
 - Extensible through UDFs, JDBC/ODBC, and external scripts (could access via python programs for example)
 - Many business intelligence (BI) tools support Hive

Hive shares many similarities with an RDBMS but there are important differences too

Feature	RDBMS	Hive
Query language	SQL (full)	SQL (subset)
Update individual records	Yes	No*
Delete individual records	Yes	No*
Transactions	Yes	No*
Index support	Extensive	Limited
Latency	Very low	High
Data size	Terabytes	Petabytes
Storage cost	Very high	Very low

Key similarities include concept of tables and SQL

Key differences include

- ¶ In RDBMS, you create tables with rigid structure that must be specified before data is added (“**schema on write**”)
- ¶ With Hive you can store data in HDFS without knowing its format. Only specify data format when you read it (“**schema on read**”)

Pro: Provides far more flexibility

Con: conflict between expected and actual data formats won’t be detected as records added to table. Only discovered when queries are performed

Drives use of RDBMS

Drives use of Hadoop

* Hive now has limited support for UPDATE, DELETE, and transactions. Not industrial strength as yet. Hive has no triggers either.

303

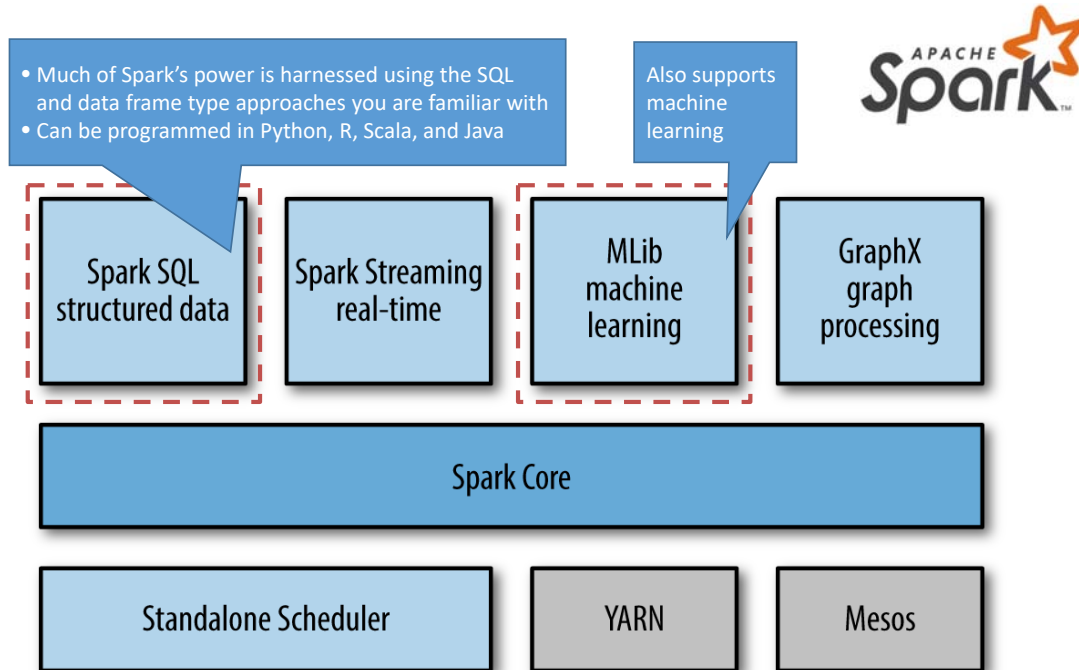
What is Spark?



- ¶ Apache Spark is a cluster computing platform on top of storage layer
 - a fast, general-purpose engine for large-scale data processing and analysis developed by AMPLab at UC Berkeley in 2009
 - Contributed to Apache Software Foundation in 2013
- ¶ Extends MapReduce with support for more components
 - Multi-pass analytics (e.g. machine learning, graph analytics)
 - Real-time streaming processing
 - Interactive ad-hoc analysis
- ¶ Runs in memory
 - Spark offers the ability to run computations in memory (can be 100x faster than MapReduce on some jobs)

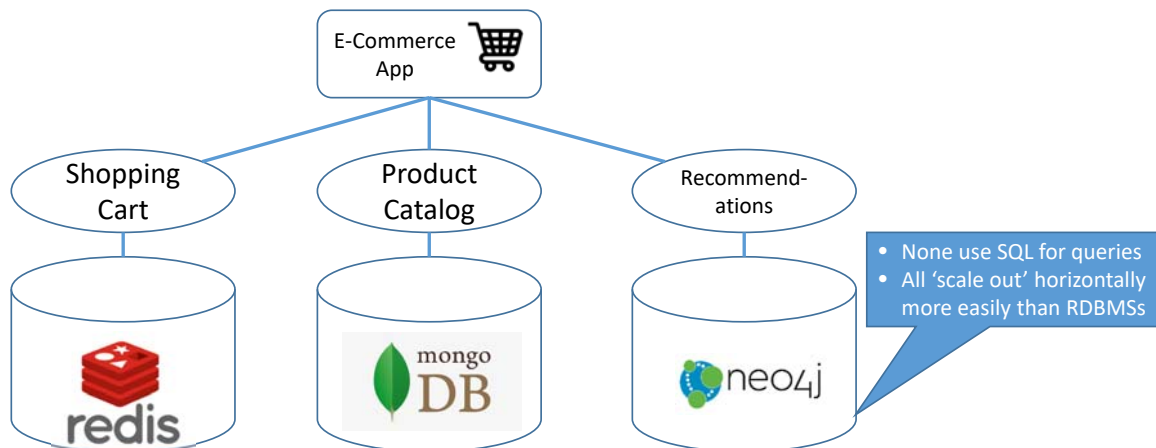
304

Key Spark component builds on the popularity of SQL



305

NoSQL (Not Only SQL) are non-relational databases that play an increasing role in data storage



Redis (Remote Dictionary Server) is an in-memory key-value database

- Stores python-like dictionary data structures
- Very fast in data caching applications

MongoDB is a leading document database – type of NoSQL database with a document-based data model

- Data stored in JSON like format
- Allows querying of semi-structured data
- Ideal in role to support a diverse product portfolio with complex querying and filtering across many product attributes – can populate web page with single query

Neo4j is a leading graph database – type of NoSQL database storing graph structures

- Keeping track of customer purchases is a common use case
- Helps in generating personalized product recommendations
- Can also be used to represent social network graphs

306

Examples of big data applications

- ¶ Recommender systems – Netflix, Amazon, Social Networks
- ¶ Netflix picked director and actors for “House of Cards” based on correlations in its data
- ¶ Wal-Mart uses text analysis and ML to produce better search results -> leading to billions in extra sales
- ¶ Morton Steak House used it to find a great PR opportunity
- ¶ Crime prediction
- ¶ Amex uses it to predict customer churn
- ¶ Insurance company used it to mine adjusters’ reports to recover \$12 million of fraudulent claims
- ¶ Finding relative with DNA sequencing
- ¶ Lots of other examples in health, manufacturing, government, water, energy



307

Agenda

Database connections and semi-structured data

Business analytics and the advent of big data

Technical challenges of big data

Intro to Hadoop and other **Big Data** technologies

Wrapping up the course

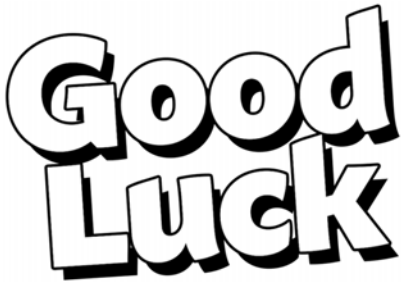
308

Final exam on Thursday

¶ Final exam 2hr 30min window from 9am 22/10 – 9am 23/10

- Open book and notes
- Covers all content of course (but less from Session 5)
- More emphasis on “hands on” topics than last year
- It will cover SQL and data modeling

¶ Any final questions??



Good
Luck

¶ Exam

¶ MAM Programme

¶ Career