

MIE1512 Data Analytics

# Welcome!

## Course Overview

# Overview

- This course is a research seminar that focuses on recent developments in the area of Data Analytics.
- Science, businesses, society and government have been revolutionized by data-driven methods. The increased access to large quantities of digital information has provided new opportunities for innovation.
- A new area of Data Analytics, known as Big Data, is made possible thanks to novel affordable techniques for processing huge amounts of data
  - Focus on Scalable Data Management

# Big Data Landscape 2016

## Infrastructure

**Hadoop On-Premise**  
cloudera, Hortonworks, Pivotal, IBM InfoSphere, splicemachine, bluedata, jethro

**Hadoop in the Cloud**  
amazon, Microsoft Azure, Google Cloud Platform, IBM InfoSphere, CAZENA, altiscale, Duoble, xplenty

**Spark**  
databricks, GridGain, TACHYON NEXUS

**Cluster Services**  
amazon, kubernetes, HPCC SYSTEMS, MESOSPHERE, CoreOS, pepperdata, StackIQ

## Analytics

**Analyst Platforms**  
Palantir, AYASDI, Quid, enigma, Digital Reasoning, ORBITAL INSIGHT

**Analytics Platforms**  
Microsoft, guavus, Datameer, inter|ana

**Data Science Platforms**  
context relevant, DataRobot, Alpine, MODE, ADATA, dataiku, DOMINO, yhat, ALGORITHMIA

**Visualization**  
tableau, Google Cloud Platform, Roambi, QOMDATA, Qlik, CHARTIO

## Applications

**Sales & Marketing**  
RADIUS, Gainsight, bloomreach, Zeta, livefyre, blueyonder, kahuna, Lattice, SAILTHRU, persado, infer, sense, AVISO, ACTIONIQ, QUANTIFIND, ENGAGIO

**Customer Service**  
MEDALLIA, ATTENTIVITY, CLARABRIDGE, STELLA Service, NGDATA, DigitalGenius, Preact, Wiseio, appurri, fuse/machines

**Human Capital**  
gild, Connectifier, textio, entelo, hiQ

**Legal**  
RAVEL, LUDICATA, Everlaw, Brevia, PREMONITION

**NoSQL Databases**  
amazon, DynamoDB, Google Cloud Platform, ORACLE, Microsoft Azure, MarkLogic, mongoDB, DATASTAX, Couchbase, SequoiaDB, redislabs, influxdata

**NewsSQL Databases**  
SAP, Clustrix, Pivotal, paradigm4, memsql, nuODB, MariaDB, VOLTDB, citusdata, deepdb, Trafojet, Cockroach LABS

**BI Platforms**  
Power BI, amazon, DOMO, Wave Analytics, GoodData, birst, platforma, looker, atscale, ACADIA, BISSENSE

**Statistical Computing**  
sas, SPSS, MATLAB

**Log Analytics**  
splunk, sumologic, kibana, CLOUD PHYSICS, loggly

**Social Analytics**  
NETBASE, DATASIFT, trackr, bitly, synthetio, simplereach

**Ad Optimization**  
MediaMath, Integral, OpenX, Adgorithms, LiveIntent, dStillery, DataXu, Appier, TAPAD

**Security**  
CYCLANE, CounterTack, cyberreason, AREA 1 SECURITY, SentinelOne, Recorded Future, Guardian Analytics, FORTSCALE, sift science, Kaybase, feedzai, SICNIFYD

**Vertical AI Applications**  
facebook, Clara, KASIST, lumiata

**Graph Databases**  
neo4j, ORIENTDB, InfiniteGraph

**MPP Databases**  
TERADATA, VERTICA, Netezza, Kognitio, dremio

**Cloud EDW**  
amazon, Google Cloud Platform, Microsoft Azure, Pivotal, snowflake, WATERLILY DATA, Infoworks

**Data Transformation**  
alteryx, TRIFACTA, tamr, StreamSets, Alation

**Data Integration**  
informatica, MuleSoft, snapLogic, BedrockData

**Real-Time**  
amazon, METAMARKETS, confluent, DATATOURNMENT, dataArtisans

**Machine Learning**  
Azure Machine Learning, H2O, Gridspace, SKYYTREE, rapidminer, DATAGRAM, deepsense, VISENZE, PredictionIO, glowfish

**Speech & NLP**  
NarrativeScience, api.ai, NUANCE, Gridspace, semanticmachines, corticalio, MindMeld, IDIBON, yseop

**Horizontal AI**  
IBM Watson, Cortana, sentient, viv, nervana, nara, SI, clarifai

**Publisher Tools**  
outbrain, mixpanel, Chartbeat, yieldbot, Yieldmo

**Govt/ Regulation**  
Socrata, OPENGOV, EN FiscalNote, enigma, PREDPOL, mark43, OpenDataSoft

**Finance**  
affirm, LendingClub, OnDeck, Kreditech, Kabbage, tidemark, ZUORA, Dataminr, Lenddo, KENSHC, AIDYA, iSENTIUM, Quantopian, sentient

**Management / Monitoring**  
New Relic, APPDYNAMICS, amazon, actifio, Numerify, splunk, DATADOG, Yricono, Anodot

**Security**  
TANIUM, Illumio, CODE42, DataGravity, CipherCloud, VECTRA, sqrrl, BlueTalon

**Storage**  
amazon, Google Cloud Platform, Microsoft Azure, panasas, nimblestorage, Qumulo

**App Dev**  
apigee, CRASK, Typesafe, CONCURRENT

**Crowd-sourcing**  
amazon, mechanicalturk, CrowdPower, WorkFusion

**Search**  
hp, Amazon, ORACLE, ENDECA, EXALEAD, Lucidworks, elastic, ThoughtSpot, MAANA, swifttype, Algolia, SINEQUA

**Data Services**  
UO, OPERA, Mu Sigma, DATA SCIENCE, DATA SCIENCE, kaggle, DataKind

**For Business Analysts**  
OrigamiLogic, ClearStory, CIRRO, import io

**SMB / Commerce**  
Google Analytics, AMPLITUDE, RJMetrics, BLUECORE, sumall, granify, Airtable, retention, custora

**Education/ Learning**  
KNEWTON, Clever, Declara, PANORAMA, knowre

**Life Sciences**  
23andMe, Counsyl, RECOMBINE, KYRUS, FLATIRON, zymogen, HealthTap, METABIOTA, ZEPHYR, ovia, Gingerio, transcriptic, Glow, enlitic, AiCure, Atomwise

**Industries**  
OPOWER, eHarmony, RetailNext, STITCH FIX, WorkFusion, BLUE2RIVER, TACHYUS, SwiftKey, Seeg, FarmLogs, HowGood, select, BOXEVER

## Cross-Infrastructure/Analytics

amazon, Google, Microsoft, IBM, SAP, SAS, hp, Autonomy, vmware, talend, TIBCO, TERADATA, ORACLE, NetApp

## Open Source

**Framework**  
hadoop, YARN, Spark, MESOS, TEZ, Flink, CDAP

**Query / Data Flow**  
SLAMDATA, Apache Drill, Google Cloud Dataflow

**Data Access**  
cassandra, HBASE, mongoDB, CouchDB, riak, OPENSCB, kafka, nifi

**Coordination**  
talend, Apache Zookeeper, Apache Ambari

**Real-Time**  
STORM, Spark, APEX, Flink, TACHYON, druid

**Stat Tools**  
R, Scala, NumPy, SciPy

**Machine Learning**  
mlilb, Aerosolve, Apache SINGA, MADlib, CNTK, TensorFlow, jupyter, DL4J

**Search**  
elasticsearch, Solr, Lucene

**Security**  
Apache Ranger, Visualizaton, Corpepin

## Data Sources & APIs

**Health**  
JAWBONE, GARMIN, practicefusion, fitbit, Withings, VALIDIC, netatmo

**IOT**  
UPTAKE, ThingWorx, helium, samsara

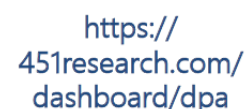
**Financial & Economic Data**  
Bloomberg, DOW JONES, YODLEE, PREMISE, S&P CAPITAL IQ, quandl, xignite, CB INSIGHTS, mattermark, estimate, PLAID

**Air / Space / Sea**  
PLANET LABS, spire, WINDWARD, CRUISE, SKYCATCH

**Location/People/Entities**  
GARMIN, foursquare, InsideView, esri, STREETLINE, CARTOOB, factual, PlaceIQ, Crism Hexagon, placemeter, BASIS, Sense

**Other**  
qualtrics, panjiva, DATA.GOV

**Incubators & Schools**  
DataCamp, INSIGHT, DataElite, METIS, The Data Incubator



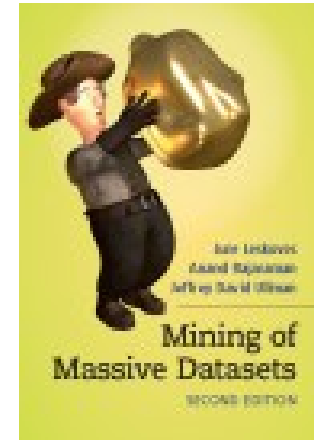




# Overview

- This seminar provides an overview of data analytics concepts, approaches, and techniques, including distributed computations on massive datasets and frameworks for enabling large-scale parallel data processing on clusters of commodity servers. Emphasis is given to algorithmic techniques for analyzing Web Data and Open Data.
- The course evaluation is based on and invigilated lab test, course participation and presentations, and a project.
- The project goal is to prepare reproducible (and potentially publishable) contributions in the area of data analytics, with an emphasis in scalable data preparation and exploration.

# General Textbook Reference



## Mining of Massive Datasets

Jure Leskovec, Anand Rajaraman, Jeff Ullman

<http://infolab.stanford.edu/~ullman/mmds/book.pdf>

# Course Grading

Invigilated Lab Test (notebooks)	20%
Class participation and presentation	15%
Course project	65%

- The presentation is based on the bibliography selected for the project
- All the project deliverables contribute to the project grade (individual grading within a group project)
- The analytical techniques must be selected from the literature, and then applied to an open dataset (the focus is on data preparation and exploration)
- Originality constraint: the project cannot use a dataset AND analytical techniques if already described elsewhere



# Project Schedule

	<i>Week</i>
Form Groups	4-5
Bibliography + PlanV1	5-7
Validation + PlanV2	8-9
Progress + PlanV3	10-11
Final Report	12-13
Presentations	6-13

# Presentation

- Select techniques described in the project literature (see Bibliography deliverable)
- Motivate the problem
- Present the approach (using examples)
- Describe related work, contributions, and relevance

# Form Groups

- Groups must have 5 to 7 members
- The project is graded on a **individual** basis
- Project plans (V1 to V3) represent the evolving group deliverables
- Individual grade is based on the contributions of each member to the group deliverables
  - Planned contribution vs. actual contribution
- Warning: group chemistry and the quality of the group deliverables will affect the individual grades

# Bibliography

- Each group member must review 5 papers in the literature relevant to the project
  - Emphasis on quality of venue, authors, publication
- Each group member must focus on covering a selected paper in depth (directly relevant to the project)
- Suggestions
  - Bibliographic portals (ACM, IEEE, dblp)
  - Look at [aminer.org](http://aminer.org) and the tutorial

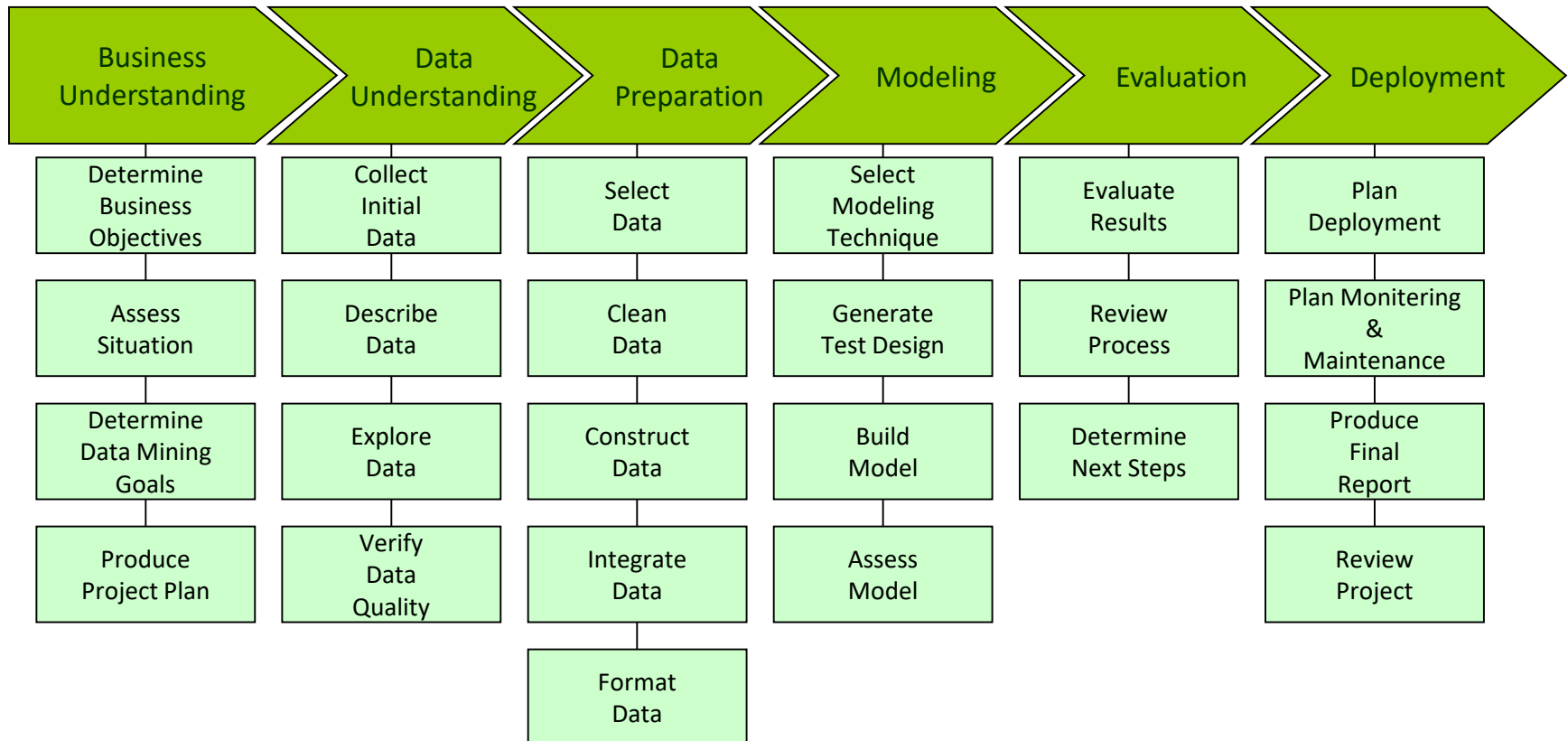
[www.wsdm-conference.org/2016/invited-speakers.html#pe-tang](http://www.wsdm-conference.org/2016/invited-speakers.html#pe-tang)

# Project Plan V1

- Project summary describing techniques and datasets, with references to the bibliography
- Initial (V1) four week long plan of activities, describing the tasks (with estimate of hours) for each project member
  - Select relevant subset of CRISP-DM
- Planning should focus on de-risking dataset selection and data preparation
  - Accuracy of task/estimates for first week is critical
- Must include notebooks with individual work

# CRISP-DM

## Phases and Tasks





# Suggestions

- Understand the Metadata of the Open Data available in a given Domain
- Browse around
  - [datatau.com](http://datatau.com)
  - XML Standard Data ([xml.coverpages.org](http://xml.coverpages.org))
  - Social Data ([datakind.org](http://datakind.org)), Urban Data
  - Developers Data ([github.com](http://github.com))
  - [commoncrawl.org](http://commoncrawl.org)
  - [ckan.org](http://ckan.org)
  - [schema.org](http://schema.org)
  - [dataverse.org](http://dataverse.org)
- Explore Data Challenges
  - [kaggle.com](http://kaggle.com)

# Project Plans V2,V3

- Updates to project summary
  - Opportunity for incorporating instructor feedback
- Updates to project plan
  - Planned vs. Actual
  - Revisions to earlier plans
- Must include notebooks with individual work
- For V2 notebooks should include a full validation of the data preparation
- For V3 you should be almost done

# Final Project Report

- The results of the project should be reported in a 5-10 page long manuscript following ACM or IEEE conference/journal style guidelines
  - Single space, single column, 8.5 X 11 paper, 11pt font
- All the project data manipulations **must be reproducible** and (mostly) completed using notebooks (including documentation and datasets)

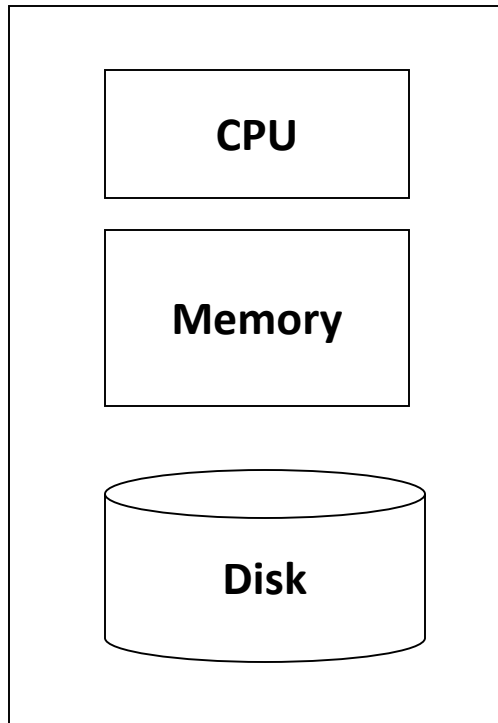
# Course Prerequisite

- MIE253 or equivalent data management course
- Suggested MOOC Option
  - <https://cs.stanford.edu/people/widom/DB-mooc.html>
  - <https://lagunita.stanford.edu/courses/DB/2014/SelfPaced/about>

# MIE1512 Data Analytics

# MapReduce

# Single-node architecture



**Machine Learning, Statistics**

**“Classical” Data Mining**



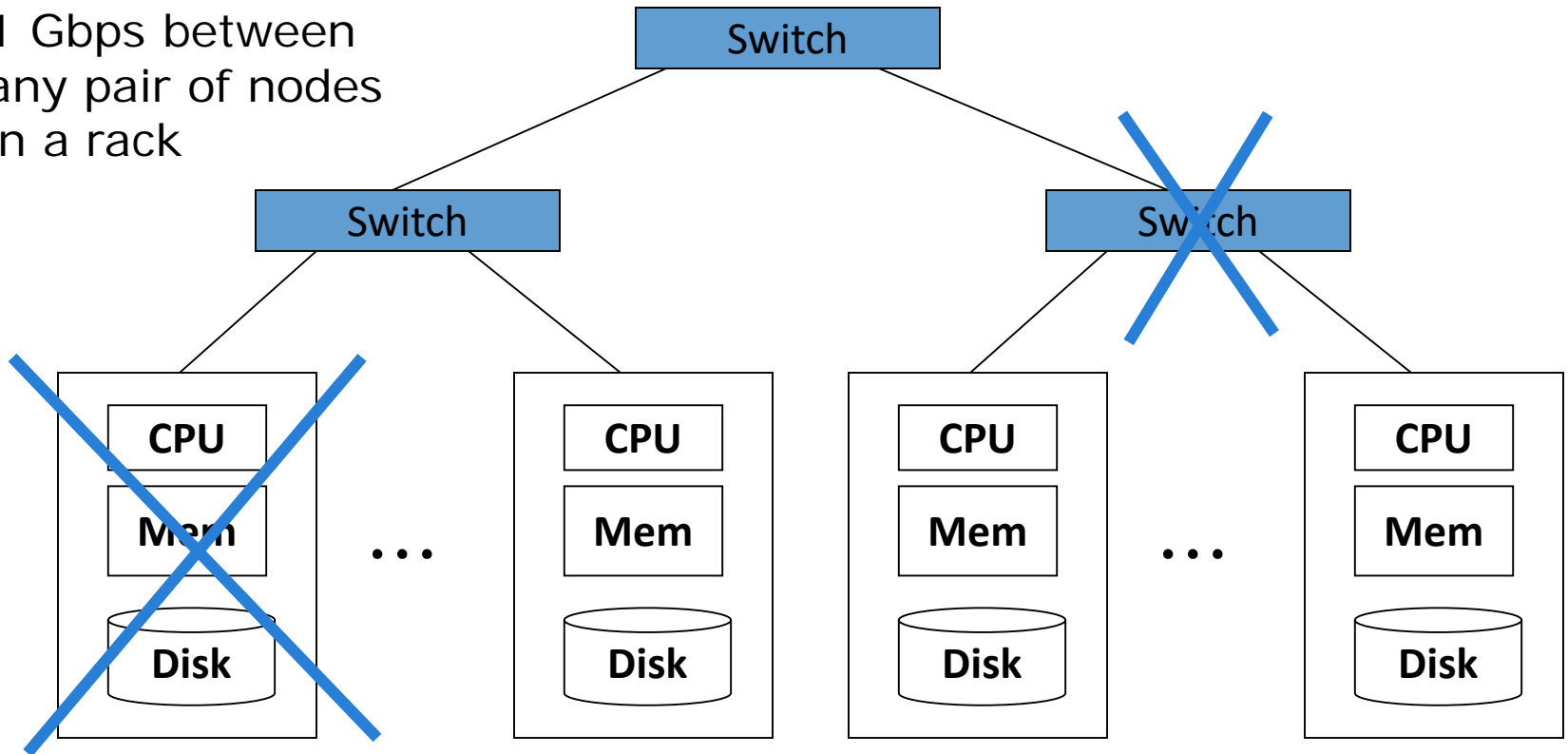
# Commodity Clusters

- Web data sets can be very large
  - Tens to hundreds of terabytes
- Cannot mine on a single server (why?)
- Standard architecture emerging:
  - Cluster of commodity Linux nodes
  - Gigabit ethernet interconnect
- How to organize computations on this architecture?
  - Mask issues such as hardware failure

# Cluster Architecture

2-10 Gbps backbone between racks

1 Gbps between  
any pair of nodes  
in a rack



Each rack contains 16-64 nodes

# Stable storage

- First order problem: if nodes can fail, how can we store data persistently?
- Answer: Distributed File System
  - Provides global file namespace
  - Google GFS; Hadoop HDFS; Kosmix KFS
- Typical usage pattern
  - Huge files (100s of GB to TB)
  - Data is rarely updated in place
  - Reads and appends are common

# Distributed File System

- Chunk Servers
  - File is split into contiguous chunks
  - Typically each chunk is 16-64MB
  - Each chunk replicated (usually 2x or 3x)
  - Try to keep replicas in different racks
- Master node
  - a.k.a. Name Nodes in HDFS
  - Stores metadata
  - Might be replicated
- Client library for file access
  - Talks to master to find chunk servers
  - Connects directly to chunkservers to access data

# Reading

- Jeffrey Dean and Sanjay Ghemawat,  
**MapReduce: Simplified Data Processing on Large Clusters,**  
**CACM 2008 (OSDI 2004)**  
<https://dl.acm.org/citation.cfm?doid=1327452.1327492>
- Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung, **The Google File System, SOSP 2003**  
<https://dl.acm.org/citation.cfm?doid=1165389.945450>

# Warm up: Word Count

- We have a large file of words, one word to a line
- Count the number of times each distinct word appears in the file
- Sample application: analyze web server logs to find popular URLs



## Word Count (2)

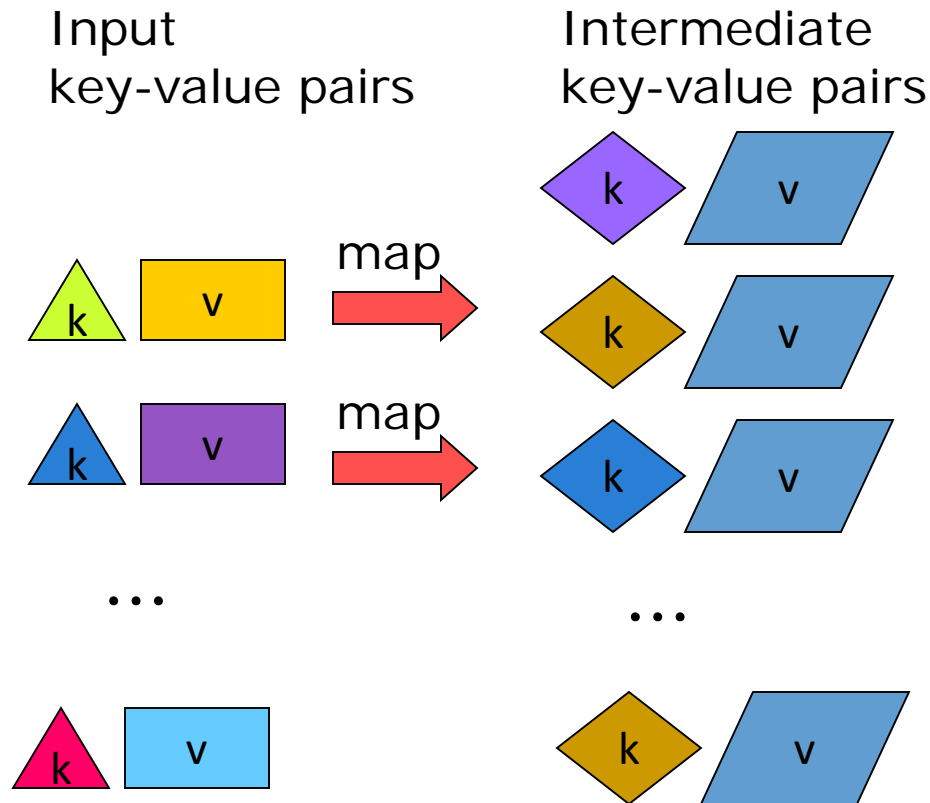
- Case 1: Entire file fits in memory
- Case 2: File too large for mem, but all <word, count> pairs fit in mem
- Case 3: File on disk, too many distinct words to fit in memory

```
-sort datafile | uniq -c
```

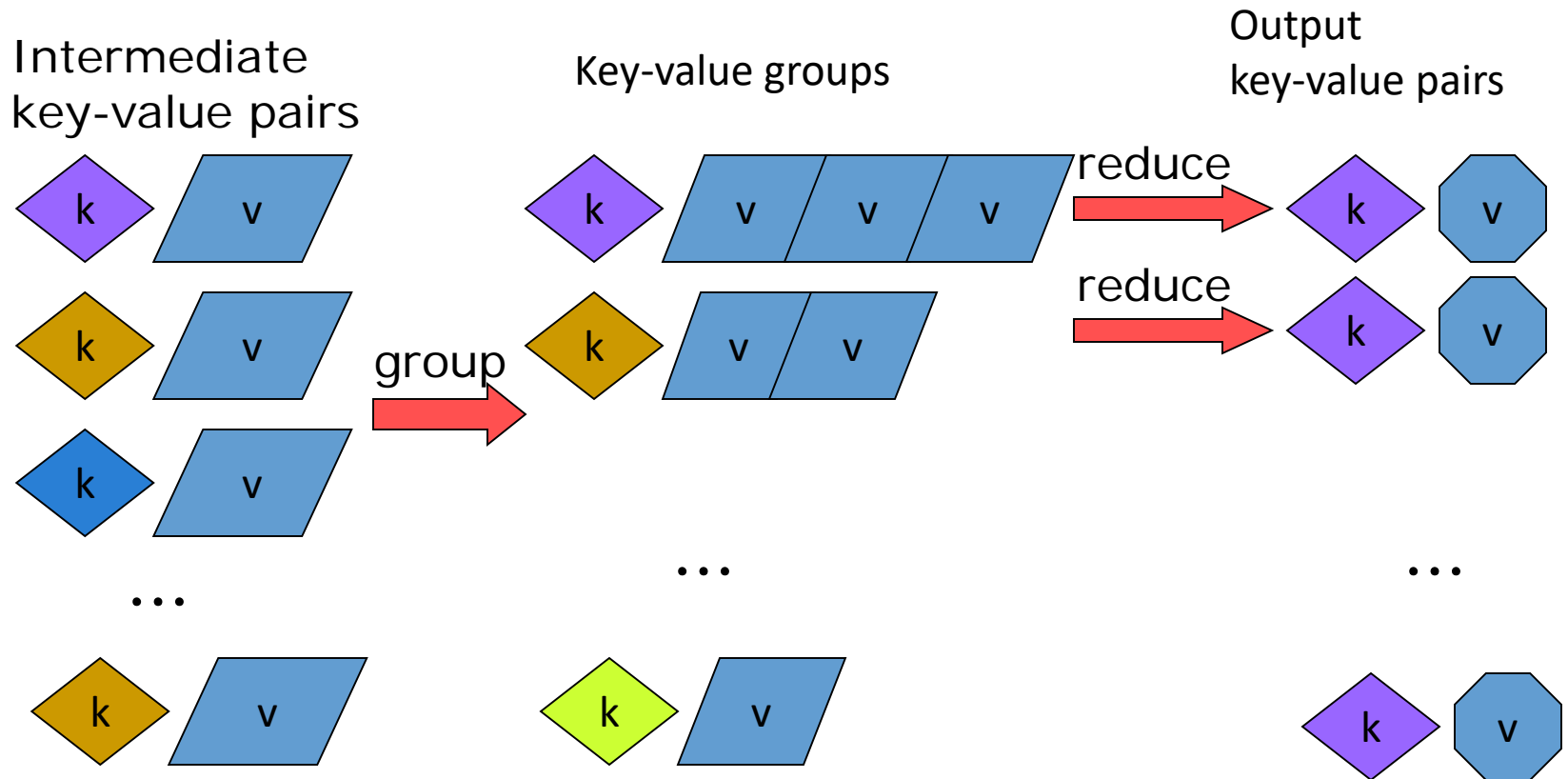
# Word Count (3)

- To make it slightly harder, suppose we have a large corpus of documents
- Count the number of times each distinct word occurs in the corpus
  - `words(docs/*) | sort | uniq -c`
  - where `words` takes a file and outputs the words in it, one to a line
- The above captures the essence of MapReduce
  - Great thing is it is naturally parallelizable

# MapReduce: The Map Step



# MapReduce: The Reduce Step



# MapReduce

- Input: a set of key/value pairs
- User supplies two functions:
  - $\text{map}(k,v) \rightarrow \text{list}(k1,v1)$
  - $\text{reduce}(k1, \text{list}(v1)) \rightarrow v2$
- $(k1,v1)$  is an intermediate key/value pair
- Output is the set of  $(k1,v2)$  pairs

# Word Count using MapReduce

map(key, value):

// key: document name; value: text of document

for each word w in value:

emit(w, 1)

reduce(key, values):

// key: a word; value: an iterator over counts

result = 0

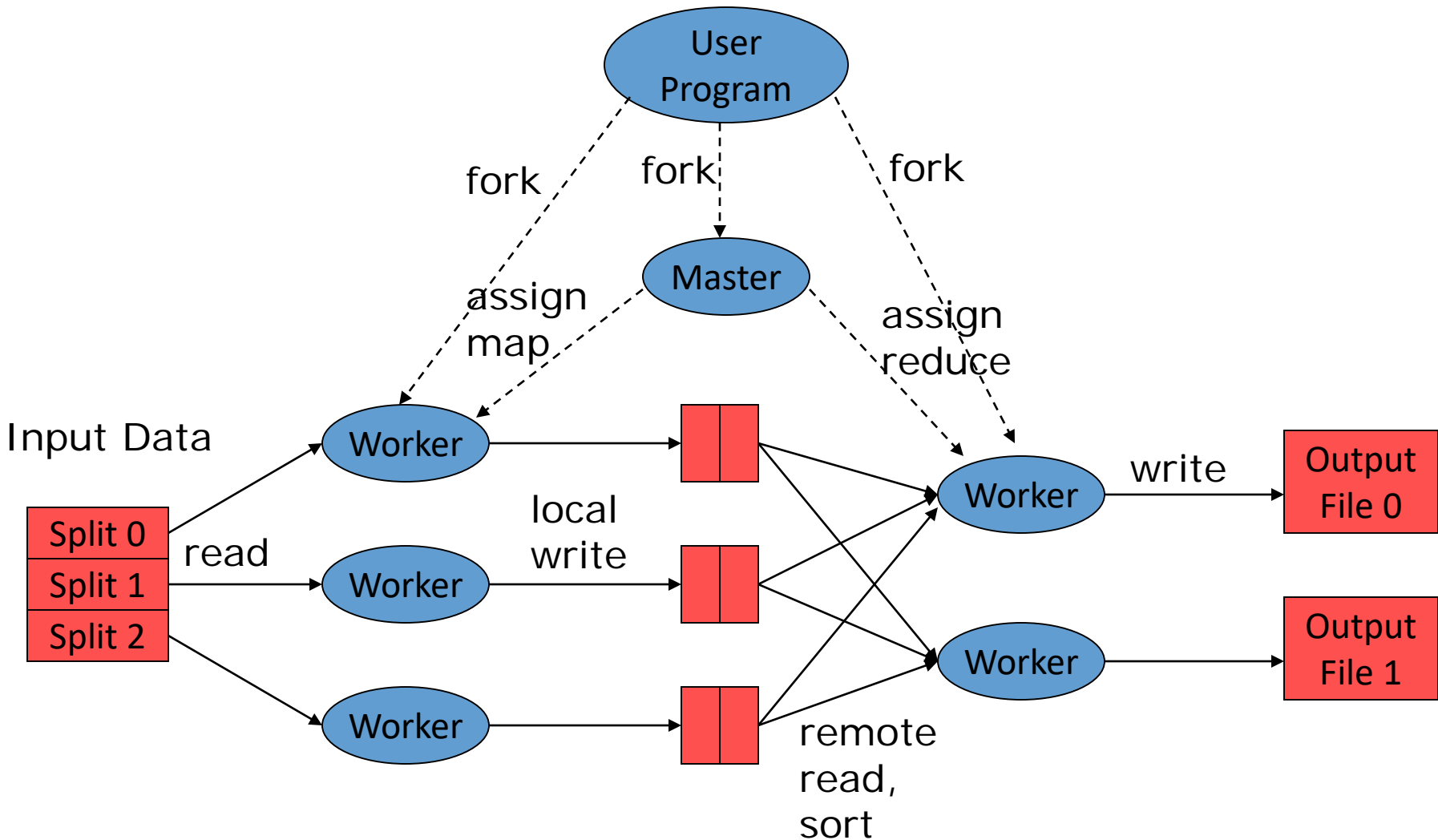
for each count v in values:

result += v

emit(result)



# Distributed Execution Overview



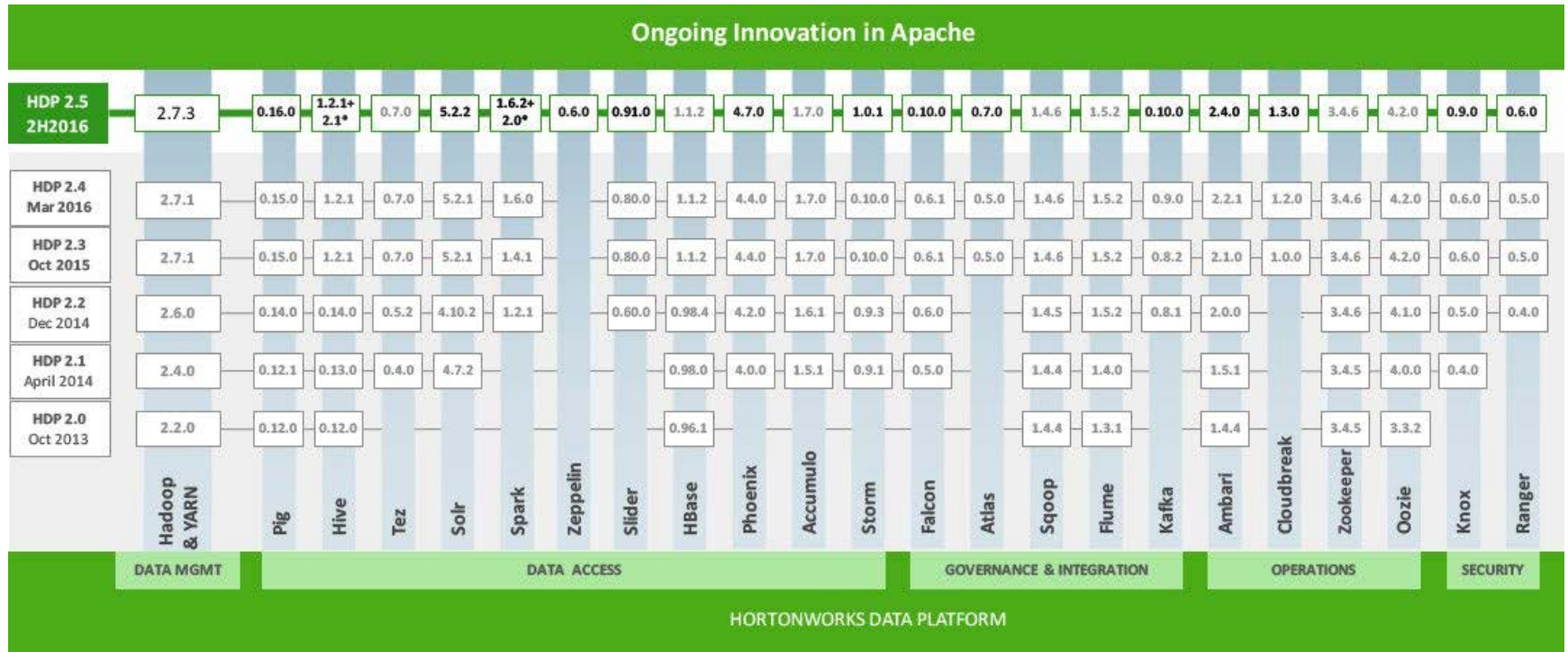
# Partition Function

- Inputs to map tasks are created by contiguous splits of input file
- For reduce, we need to ensure that records with the same intermediate key end up at the same worker
- System uses a default partition function e.g.,  $\text{hash}(\text{key}) \bmod R$
- Sometimes useful to override
  - E.g.,  $\text{hash}(\text{hostname}(\text{URL})) \bmod R$  ensures URLs from a host end up in the same output file

# Implementations

- Created by Google
    - Not available outside Google
  - Hadoop
    - An open-source implementation in Java
    - Uses HDFS for stable storage
    - Available from <http://lucene.apache.org/hadoop/>
  - Several commercial vendors
    - Cloudera
    - Hortonworks
      - Sandbox (VM image) is the easiest way to get started
- <http://hortonworks.com/products/hortonworks-sandbox/>
- Supports Hadoop 2.3, many tutorials

# Hadoop Evolution



\* Spark 1.6.2+ Spark 2.0 – HDP 2.5 support installation of both Spark 1.6.2 and Spark 2.0. Spark 2.0 is Technical Preview within HDP 2.5.  
Hive 1.2.1+ Hive 2.1 – Hive 2.1 is Technical Preview within HDP 2.5.

# Cloud Computing

- Ability to rent cluster computing by the hour
  - Additional services, e.g., persistent storage
  - Eg, Amazon's "Elastic Compute Cloud" (EC2), MS Azure, IBM Bluemix, Google Cloud Platform
- Facilitates scalability and elasticity

# Map-Reduce: Environment

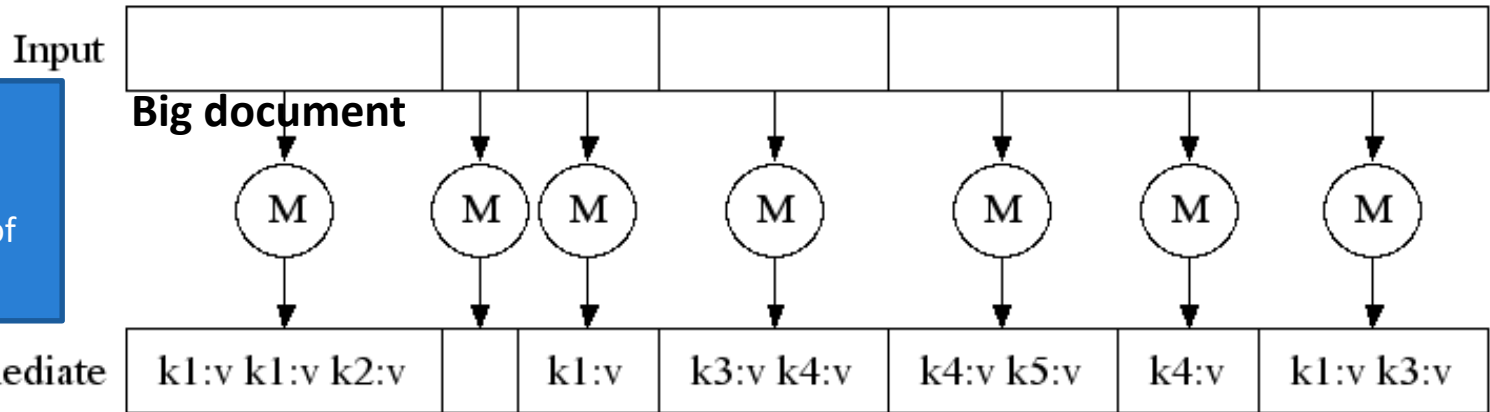
## Map-Reduce environment takes care of:

- Partitioning the input data
- Scheduling the program's execution across a set of machines
- Performing the **group by key** step
- Handling machine failures
- Managing required inter-machine communication

# Map-Reduce: A diagram

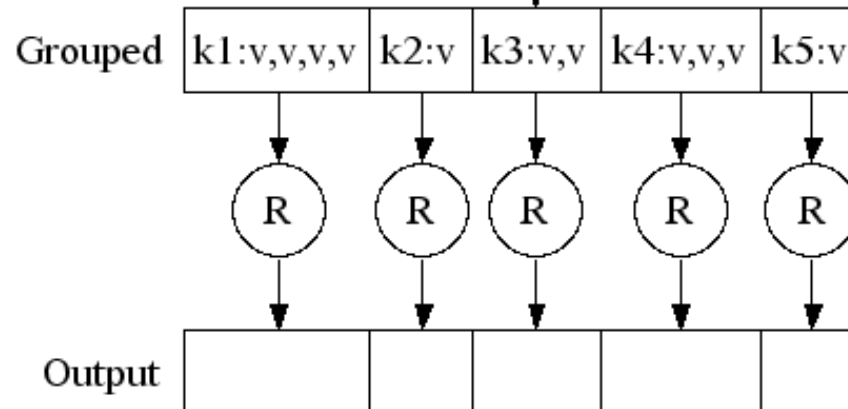
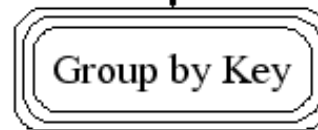
## MAP:

Read input and produces a set of key-value pairs



## Group by key:

Collect all pairs with same key  
(Hash merge, Shuffle, Sort, Partition)



# MR Exercise 1: Host size

- Suppose we have a large web corpus
- Let's look at the metadata file
  - Lines of the form (URL, size, date, ...)
- For each host, find the total number of bytes
  - i.e., the sum of the page sizes for all URLs from that host



# MR Exercise 2: Distributed Grep

- Find all occurrences of the given pattern in a very large set of files

## Exercise 3: Graph reversal

- Given a directed graph as an adjacency list:  
src1: dest11, dest12, ...  
src2: dest21, dest22, ...
- Construct the graph in which all the links are reversed

## Exercise 4: Frequent Pairs

- Given a large set of market baskets, find all frequent pairs
  - Remember definitions from Association Rules lectures