

MAS DSE 260: Capstone Project

İlkay ALTINTAŞ, Ph.D.

Lecture 1: Setting up Your Data Process

Today's Topics

1. Reviewing where we are
2. STEP II: Designing the Data Acquisition and Preparation Pipelines
3. Report II Format : DUE 2/1/18

Process Roadmap (260 A)

- ✓ Step 1: Understanding the Challenge
 - ✓ REPORT 1: due 1/18
- Step 2: Designing the Data Acquisition and Preparation Pipelines
 - REPORT 2: due 2/1
- Step 3: Exploring Data
 - PRESENTATION 1: 2/3
 - REPORT 3: due 2/15
- Step 4: Defining Your Hypothesis and Minimum Viable Modeling Product
 - REPORT 4: due 3/1
- Step 5: Creating a Solution Architecture for Modeling and Optimization
 - PRESENTATION 2: 3/3
 - FINAL WINTER REPORT: due 3/16

Asking the Right Question

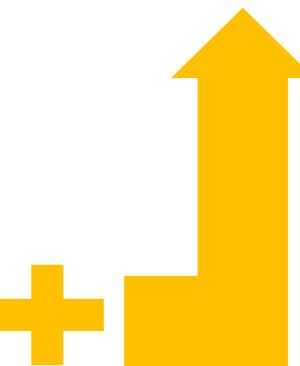


**“A problem well defined
is a problem half
solved.”**

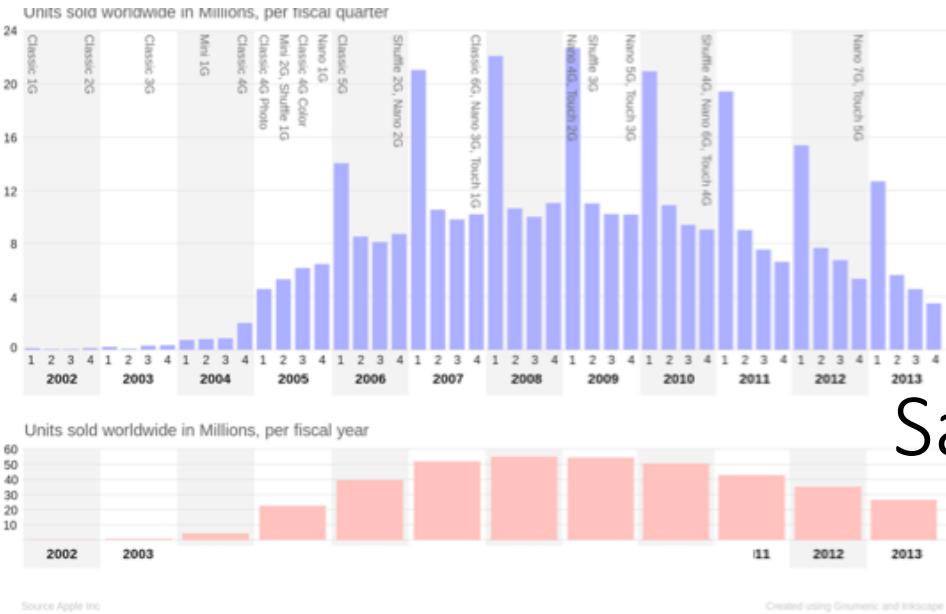
Charles F. Kettering

Define the Problem

Evaluate a new product



Sales figures



Call center logs



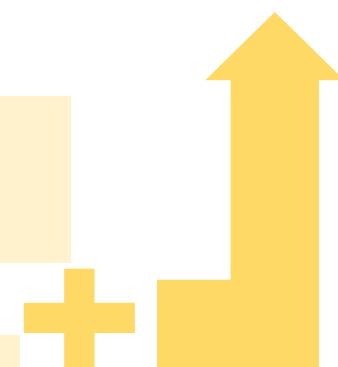


Detect equipment failure

Sensor data

Sensor data

Sensor data



Better targeted
marketing

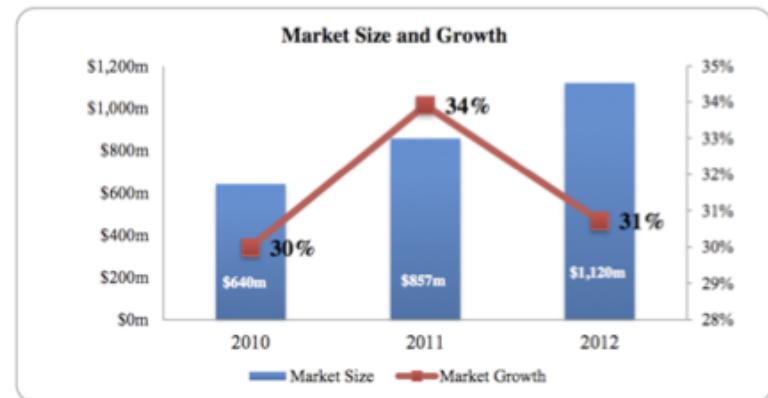


Customer data



Marketing data

1 Market Size and Growth



(ialtintas@ucsd.edu)



Assess the Situation



Risks

Benefits

Contingencies

Regulations

Resources

Requirements

Assess the Situation

Define Goals



Objectives

Criteria

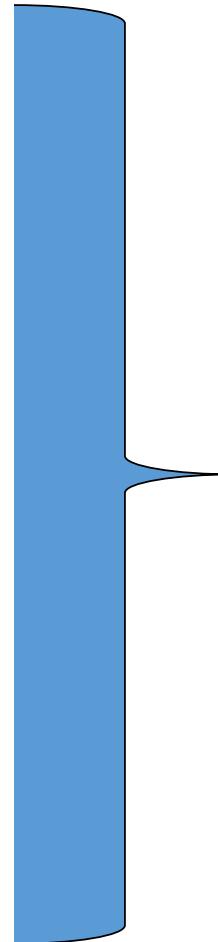
Define the Problem



Assess the Situation



Define Goals



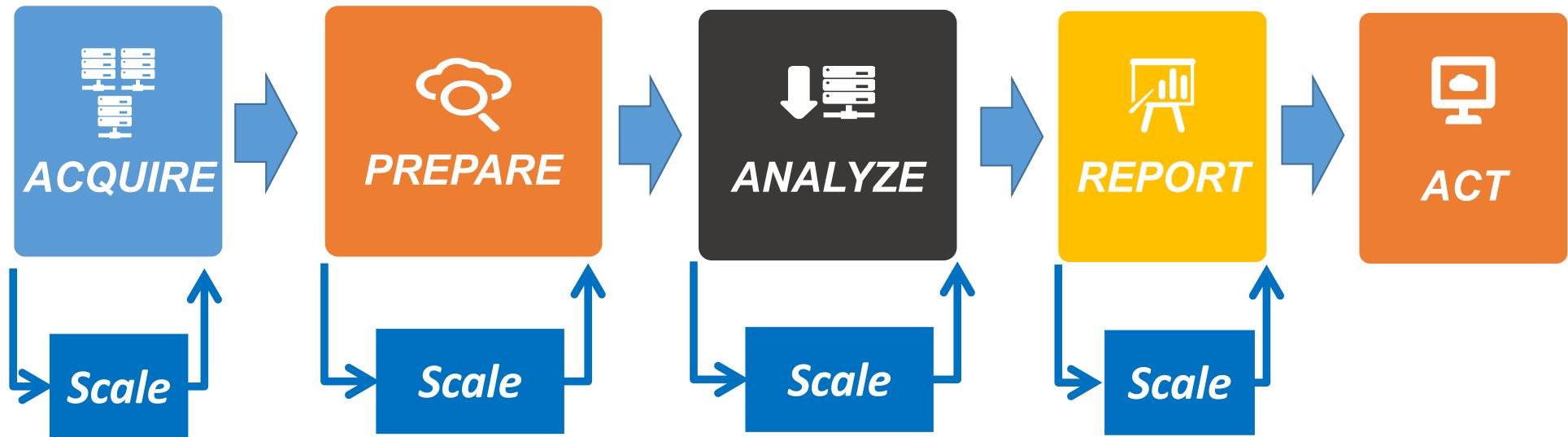
Formulate the Question



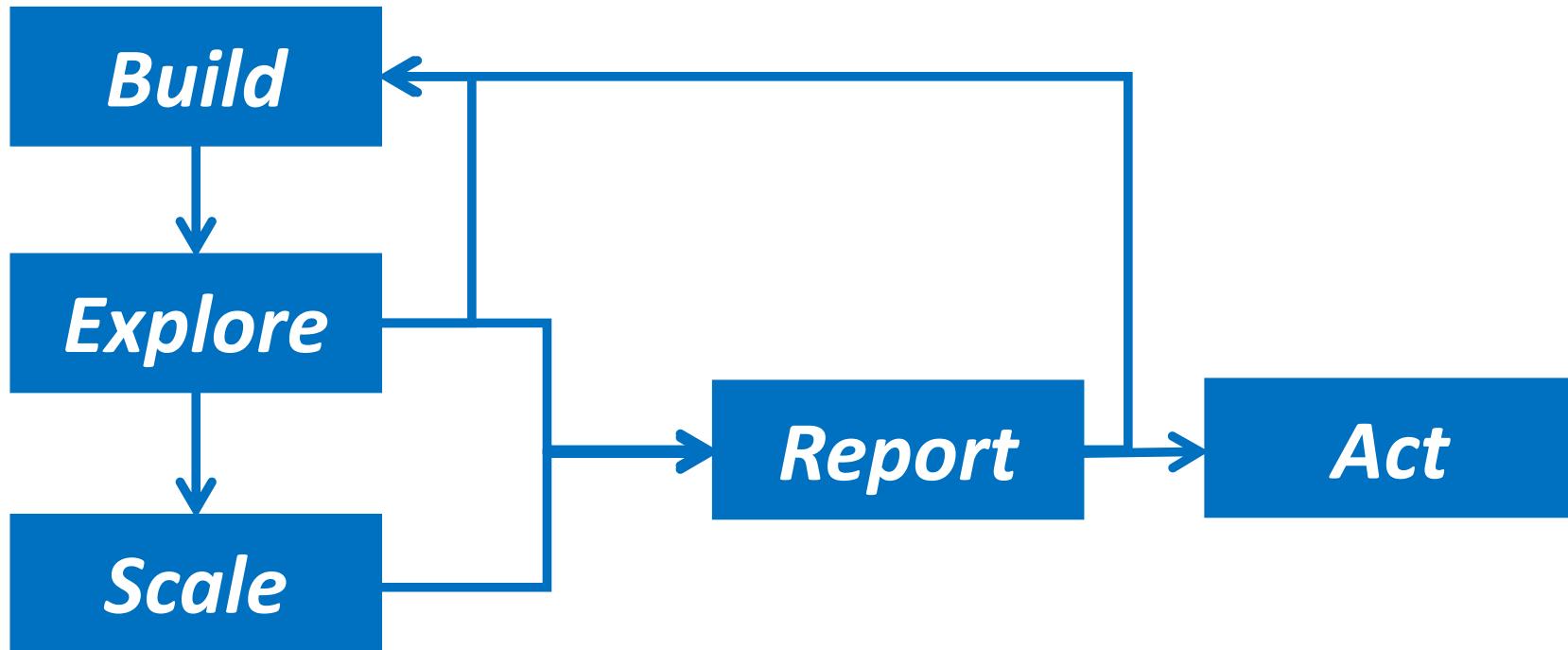
Data Science Process

Data Engineering

Computational Data Science

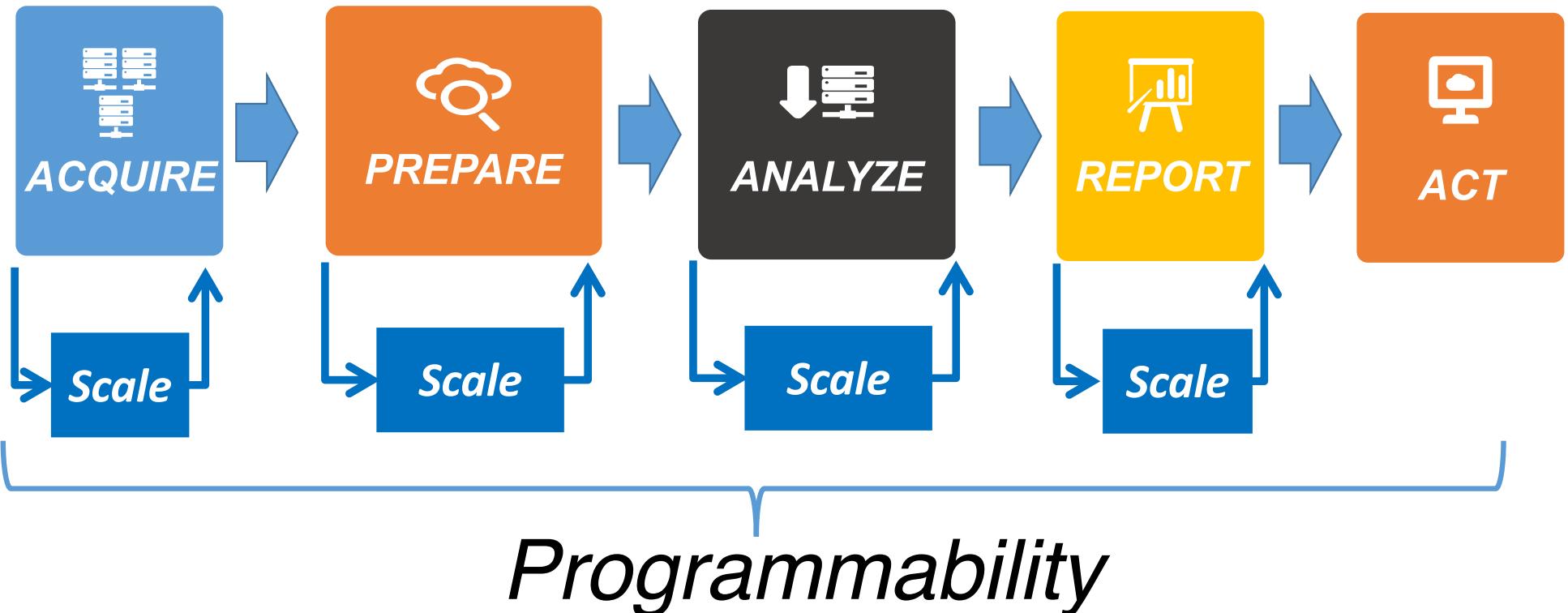


Many iterations and rollbacks between steps.



Data Engineering

Computational Data Science



NEXT:

STEP 2: Designing the Data Acquisition and Preparation Pipelines

Why do we need this step?

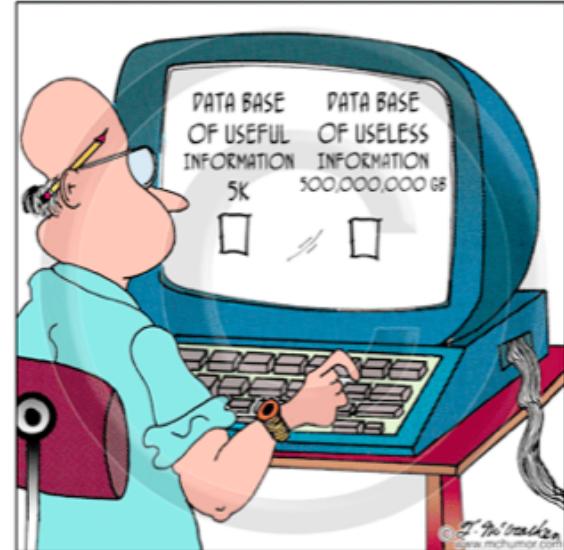
- Data ingest and engineering is the basis for effective modeling
- Many quality issues with data that needs to be cleaned up
 - Not just missing values, but also deeper problems on uncertainty and bias related to data collection
- Defined data pipelines (potentially automated) help with data refresh and updates

If your data pipeline is strong, your modeling will be strong!

Objectives

- Understanding what data is available and useful to you
- Plan for data transport and locate your datasets in the right database, storage or analytics environment
- Model the data
- Identify data integration needs if applicable
- Set up a data pipeline that streamline (and potentially automates) data ingest and refresh
- Design programmatic access to data

McHUMOR.com by T. McCracken



©T. McCracken mchumor.com

Planning for a Good Data Pipeline

- Focus on
 - where your datasets come from (data sources)
 - adding context and model data as you collect it
 - privacy and policy issues with the data
 - data quality assurance
 - where each dataset gets stored
 - how can pieces of data be combined
 - how each datasets get used
 - programmatic access to data



Using the PPODS Approach

- Each step in your data pipelines is a separate pod
- Define success metrics for calling a step done

Metrics for accountability should be built into the process.

Cost Timeline

Planning of deliverables

Purpose Expectations

Pod → sub-process

Defined by:

- Purpose and goal
- Stakeholders
- Expectations
 - Key questions to be answered, production/consumption relationships, needs, dependencies, limits, ...
- Contracts
 - Performance, economic, accuracy, policy, privacy, reproducibility, political, ...
- Knowns
- Known unknowns

Step II Report Guidelines

- Title, team members and advisor(s)
 - Sections:
 - Raw Data Sources
 - Summaries of each dataset description from report 1
 - Table for Dataset name, source location, destination in your data pipeline, data movement and processing scripts and notebooks, and data size
 - Data Exploration, Cleaning, Wrangling and Engineering
 - Data Exploration Summary
 - Data Preprocessing Approach
 - Approach for storing processed and/or integrated data
 - Processed dataset description for each processed dataset including why you want to process it that way
 - Table for processed data sets including processed data set name, input datasets, link to the processing scripts and notebooks, and provisional data size
 - Approach for Feature Engineering and Data Modeling
 - Summary of feature sets
 - Table for feature set including links to input datasets, feature engineering scripts and notebooks, and provisional data size
 - Approach for Data Access
 - Initial design for data querying interfaces
 - Justification for manual vs. programmatic access
 - Data Pipeline
 - Description of the needs, approach, and data access and refresh frequency
 - Logical diagram showing major data pipeline components for data sources and sinks
 - Set up for your data environment
 - Cloud vs. local, database vs. flat files, etc.
 - Bullets for each team member's individual contributions in Step 2
 - Any major updates to Step 1 as a result of data pipeline step
- Keep it to 5-8 pages
 - Due date: 2/1/2018 midnight



You will have more time to iterate. We are doing the first draft. Next, are some reminders on what to start focusing on.

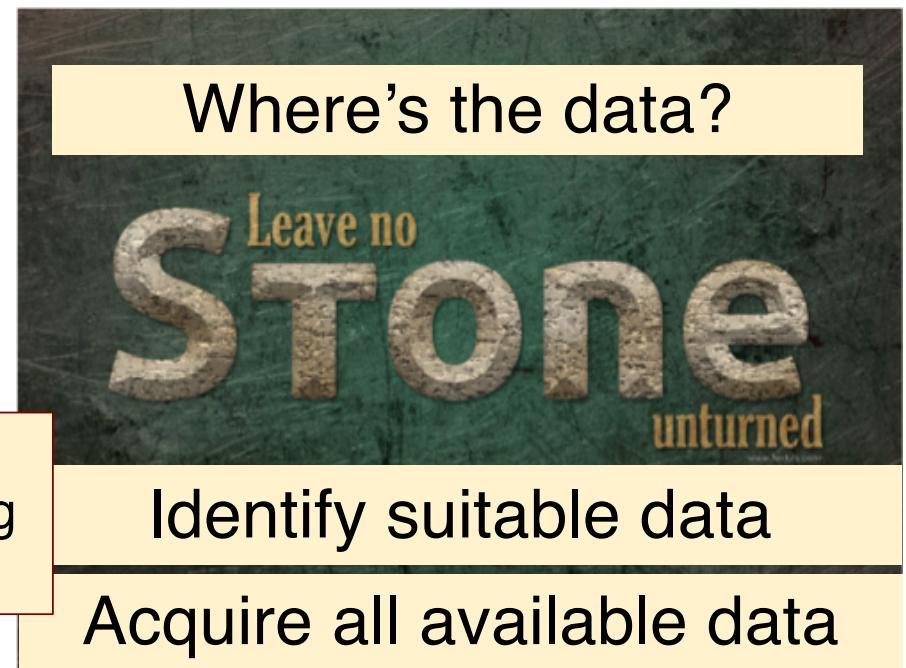
Focus 1: Acquiring Data

After you are done with this part, you should be able to..

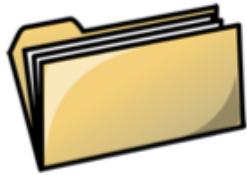
- List techniques and technologies to access and retrieve the data you need
- Describe an example scenario related to your problem accessing your datasets



Data can come from a variety of sources using different technologies



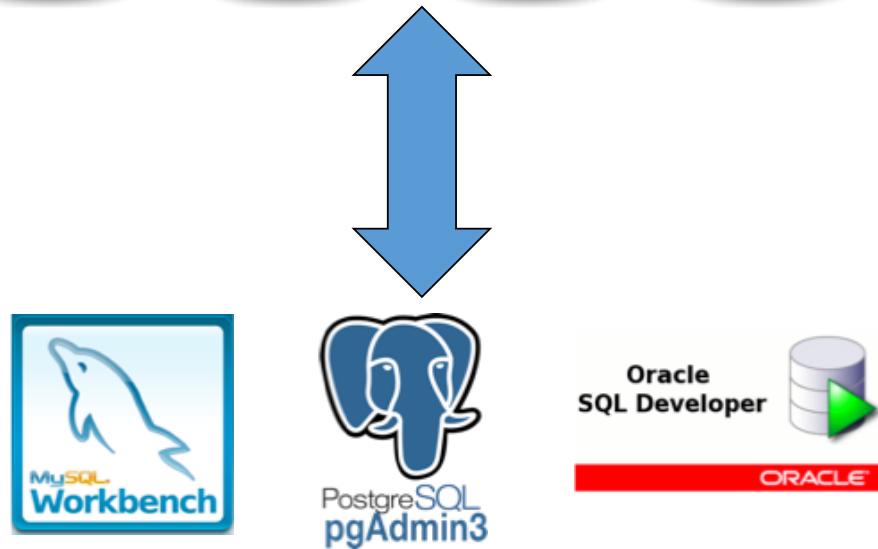
Data comes from many places...



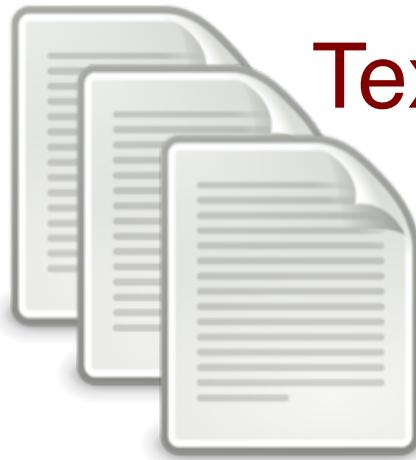
...with many ways to access it.



Traditional databases



SQL and query browsers



Text files



JavaScript



python™



Ruby



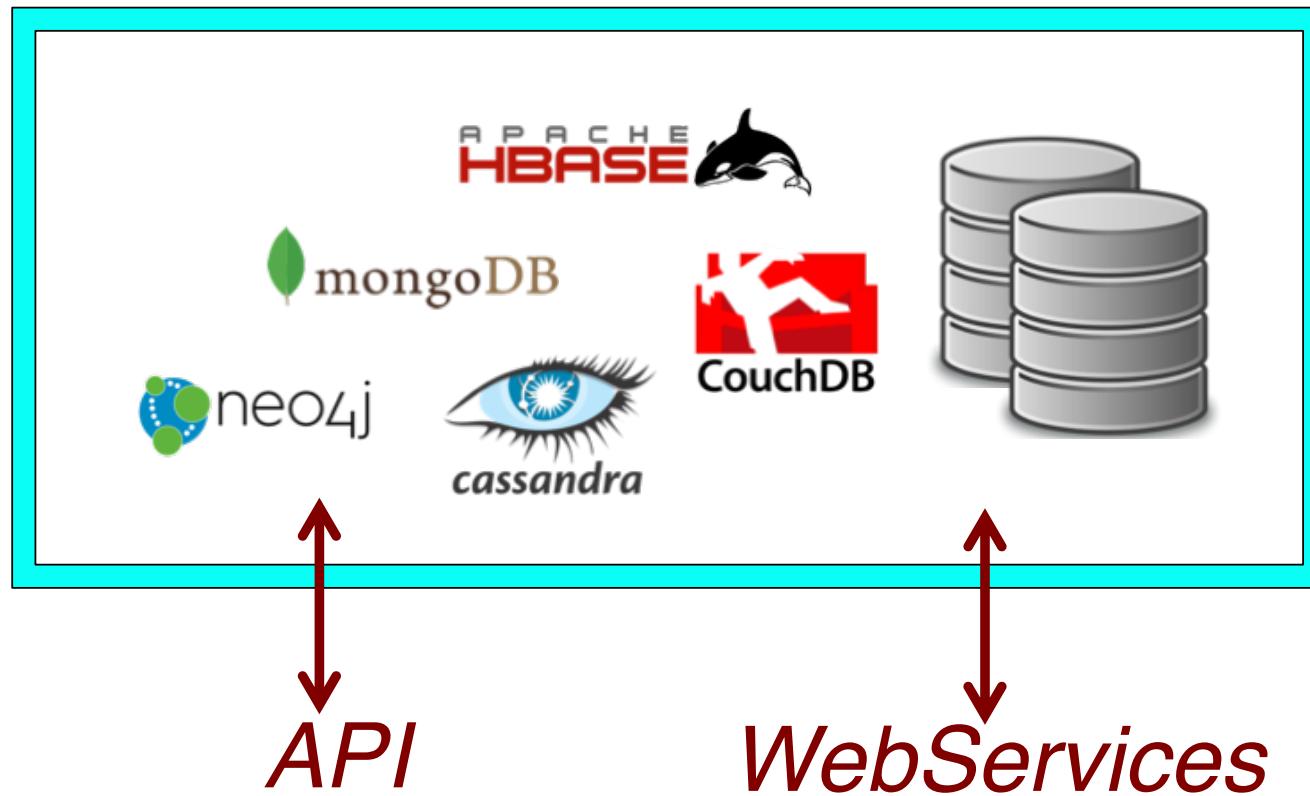
Perl

Scripting languages

Remote data



NoSQL storage



Acquiring Data From WIFIRE

Historical weather

Current weather

Real-time tweets
near fires



WebSocket



REST



Traditional databases

SQL and query browsers



Remote data

Web Services



Text files



NoSQL storage

Scripting languages

Web Services

Programming Interfaces

Focus 2: Exploring Data



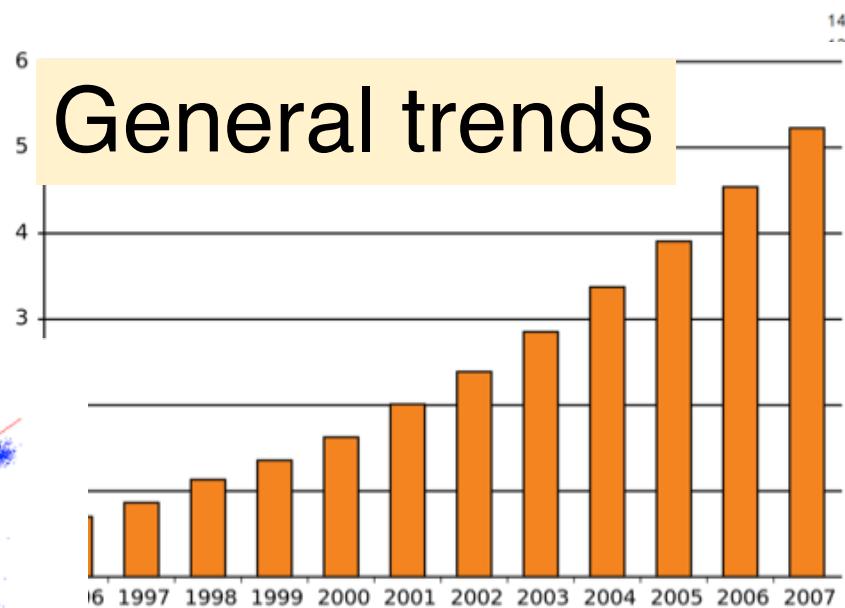
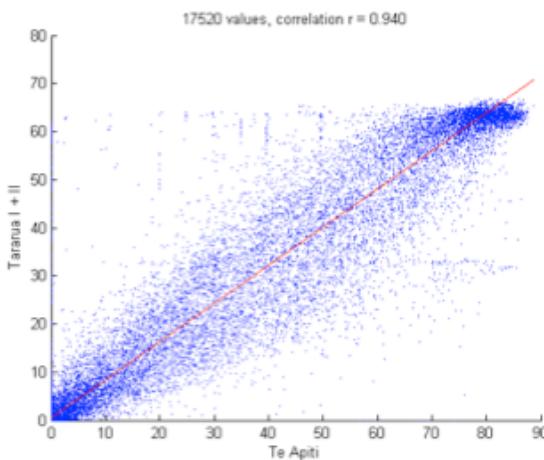
Why Explore?

Goal: Understand your data

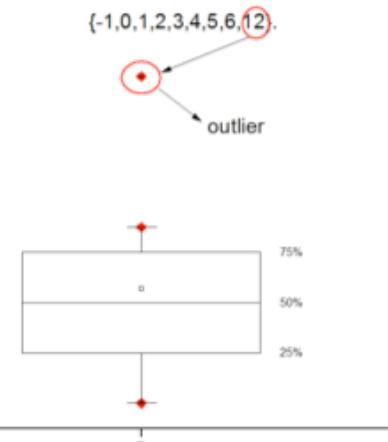


Why Explore?

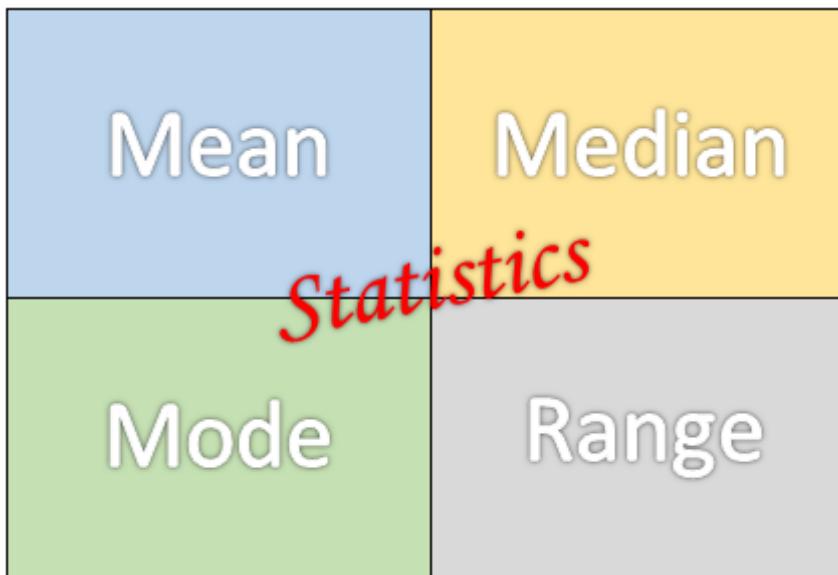
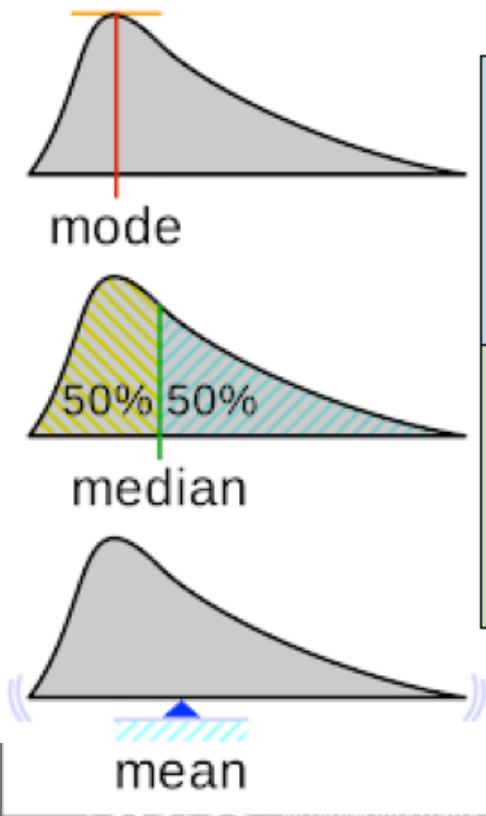
Correlations



Outliers

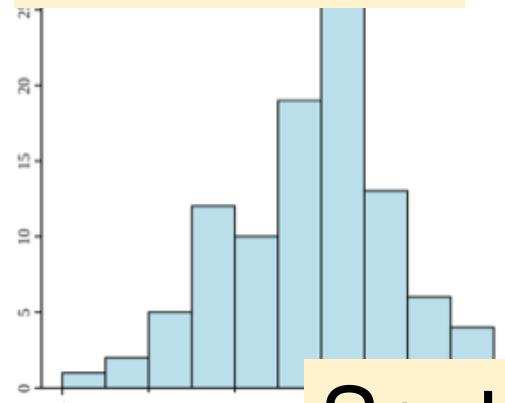


Describe Your Data

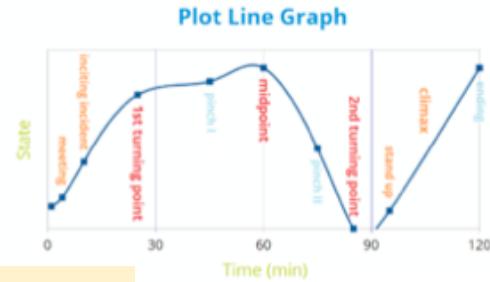


Visualize Your Data

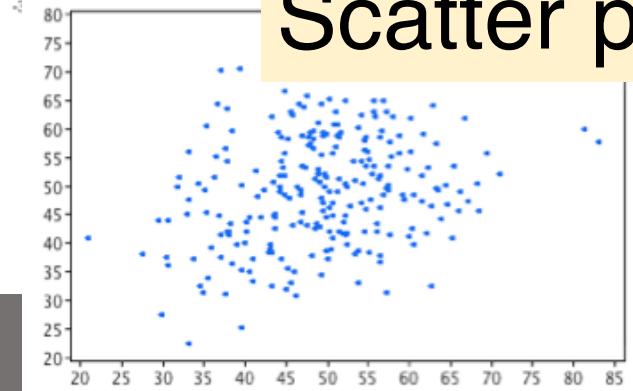
Histogram



Line graphs

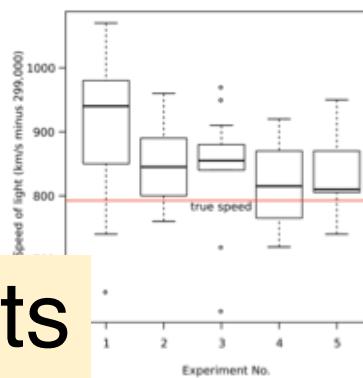


Scatter plots

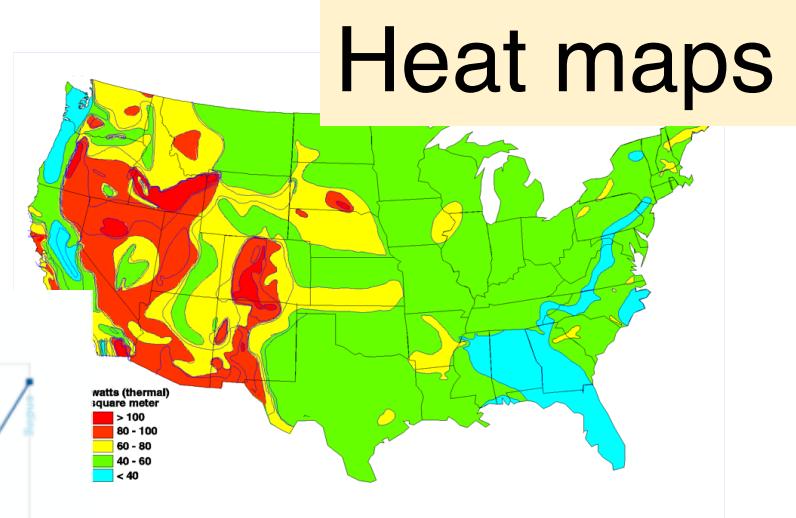


Boxplots

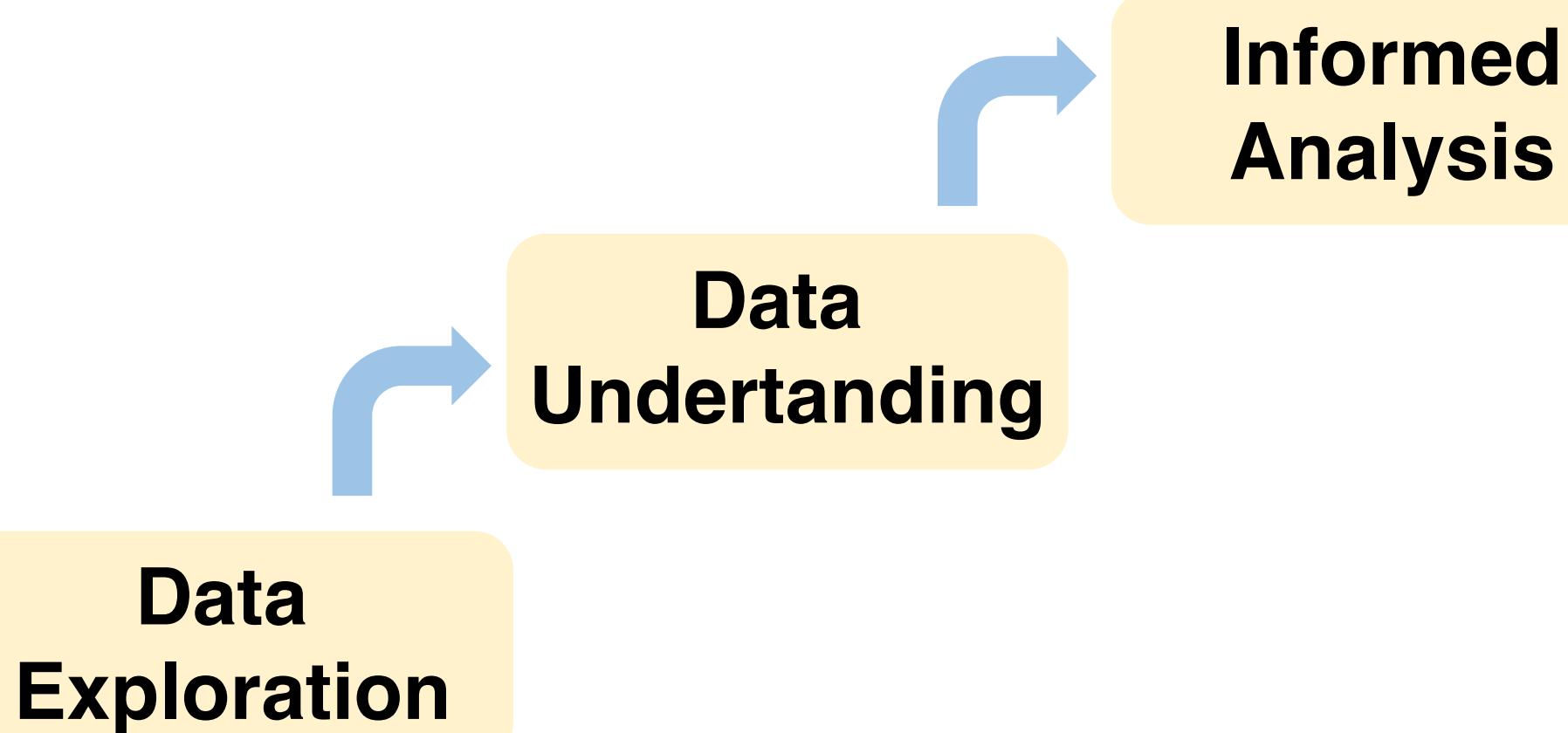
DSL 200 May 2011



Heat maps

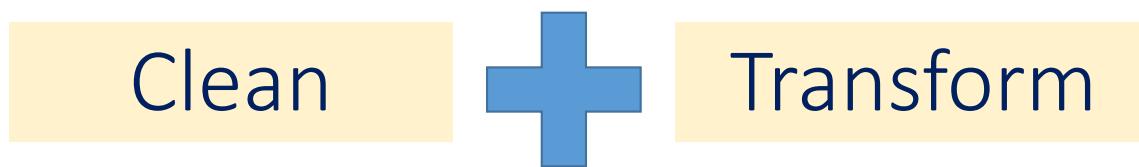


UC San Diego



Focus 3: Pre-processing Data

Describe what is needed to transform raw data to data that can be used for analysis.





Data Quality Issues

Real-world data is messy!

Inconsistent values

Duplicate records

Missing values

Invalid data

Outliers

Addressing Data Quality Issues

Remove data with
missing values

Generate best estimate for
invalid values

Merge duplicate records

Remove outliers

Addressing Data Quality Issues

Remove data with
missing values

Generate best estimate for
invalid values

Merge duplicate records

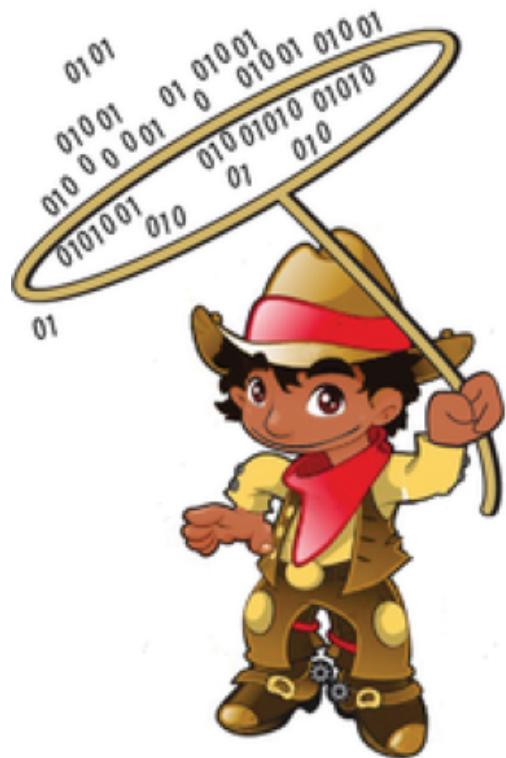
Remove outliers

*Domain
Knowledge*

Getting Data in Shape

Data Munging

Data Preprocessing



Data Wrangling

Data Munging

*Dimensionality
Reduction*

*Data
Manipulation*

Transformation

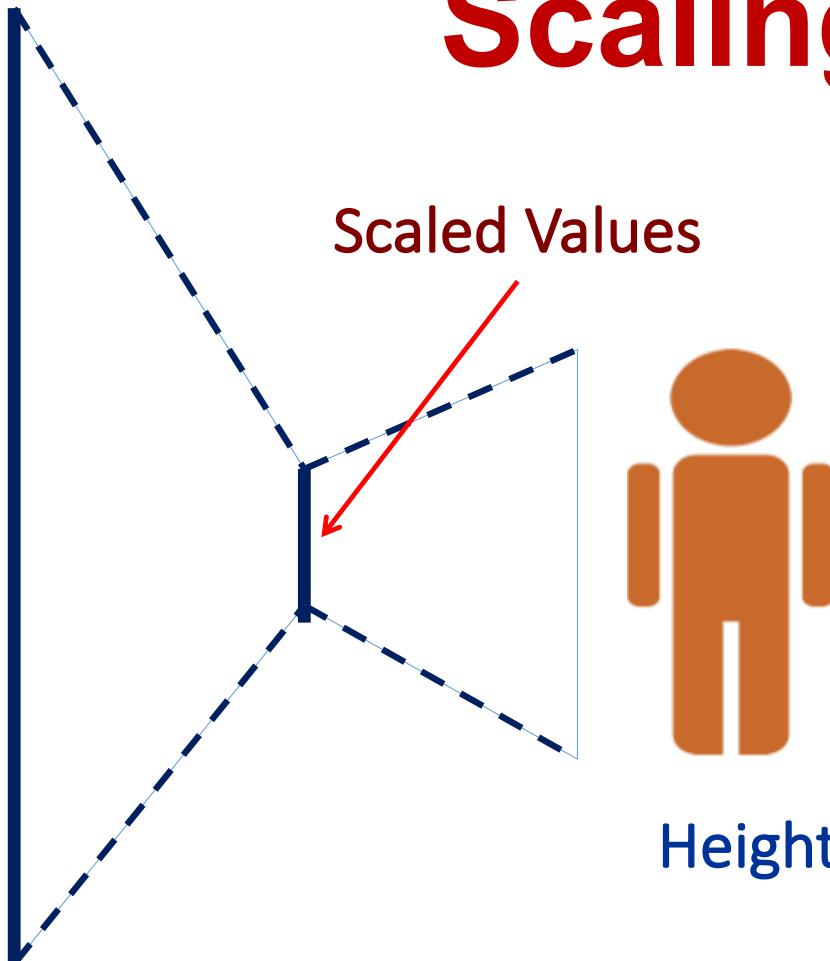
*Feature
Selection*

Scaling

Scaling



Weight

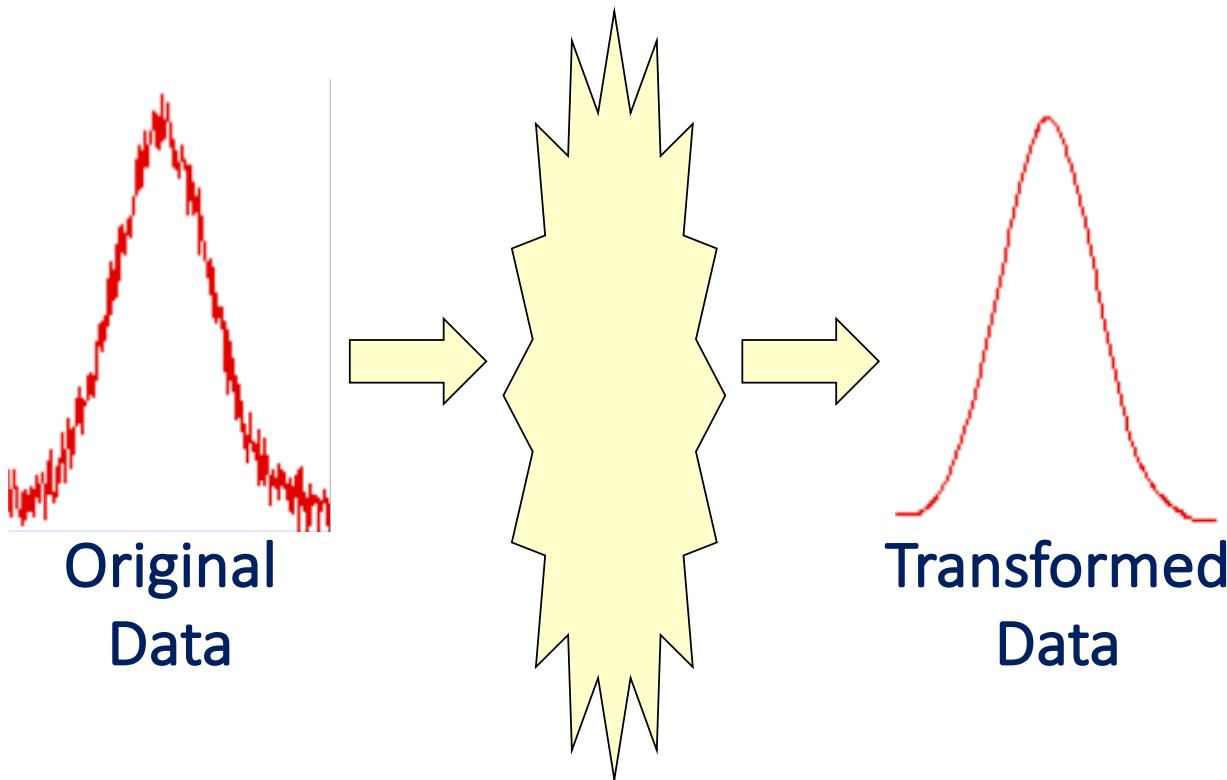


Scaled Values

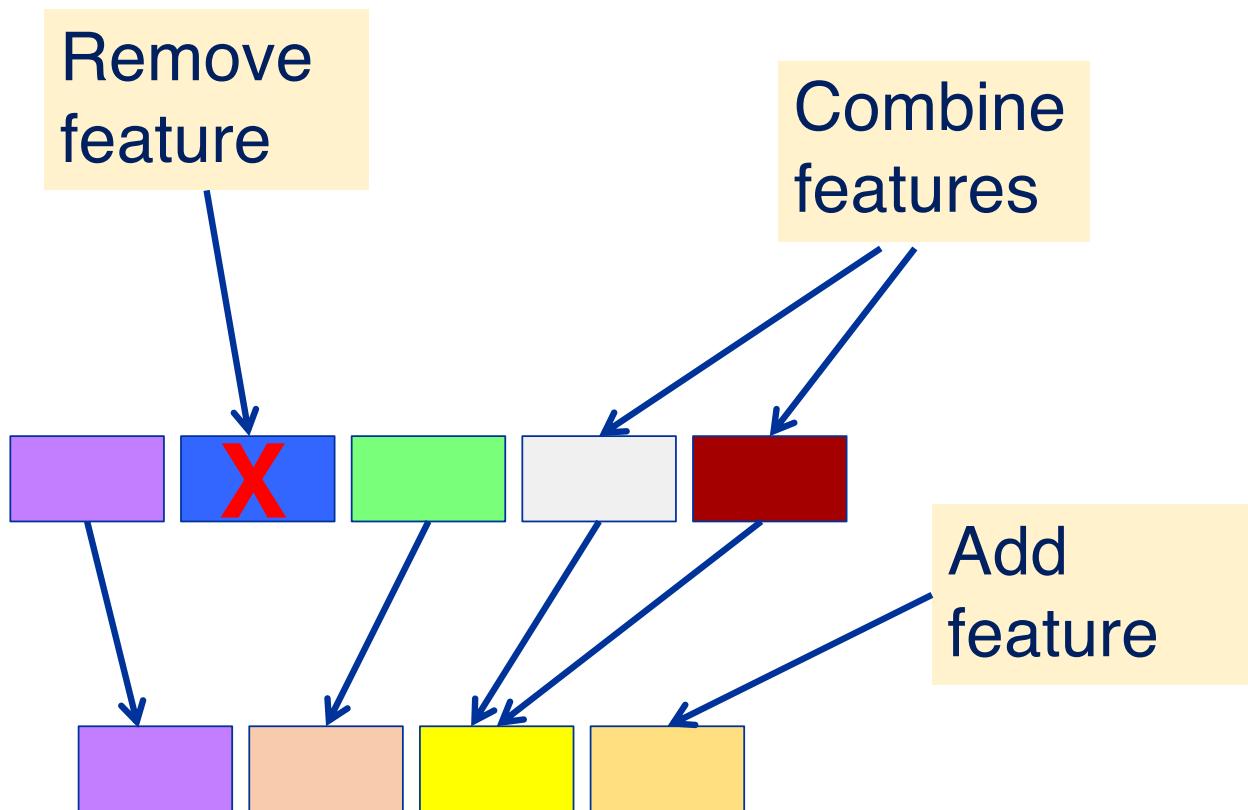


Height

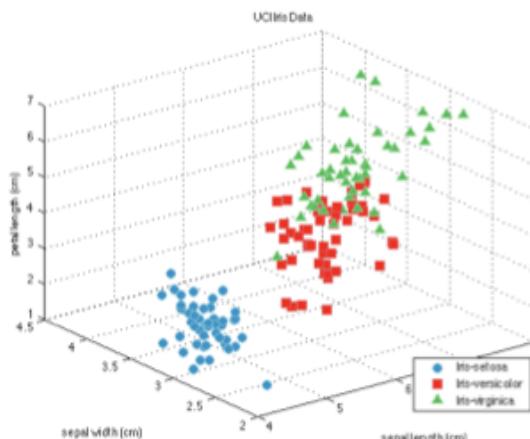
Transformation



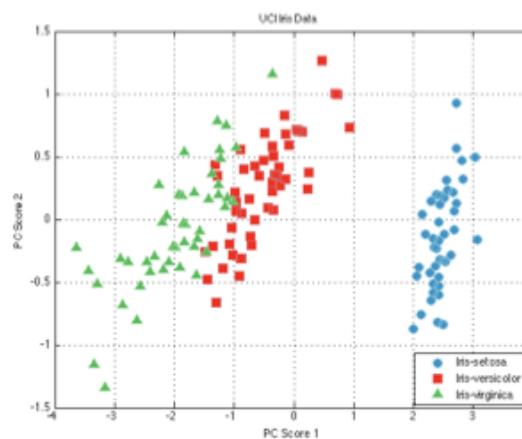
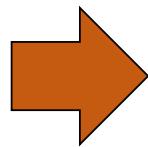
Feature Selection



Dimensionality Reduction

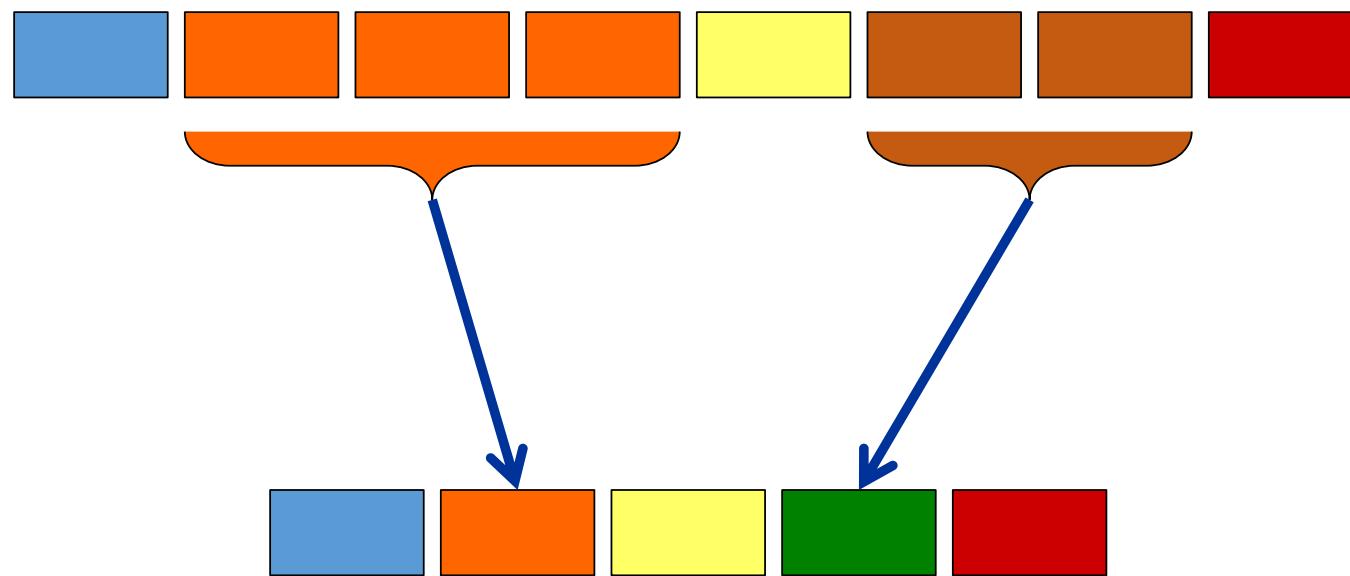


3D



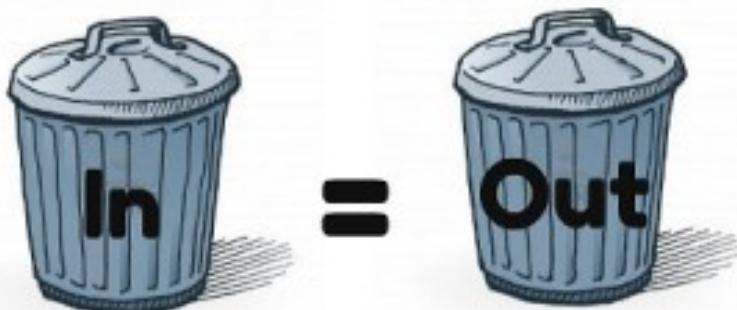
2D

Data Manipulation



Always Remember!

Garbage in = Garbage out

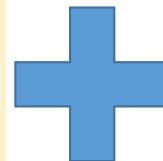


Data preparation is
very important for
meaningful analysis!

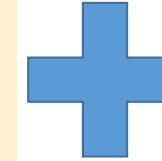
Focus 4: Designing Your Data Pipelines

Draft solution architecture to turn data operations into an ongoing and accountable process.

Solution Architecture
Diagram



Step by step
flow diagram



Environment

Consider...

- Batch vs. streaming or both
- Local vs. cloud or both
- Database vs. object storage or both
- Manual vs. automated
- Refresh frequency
- Data provenance

A UNIX PIPE

cat 0 → sort

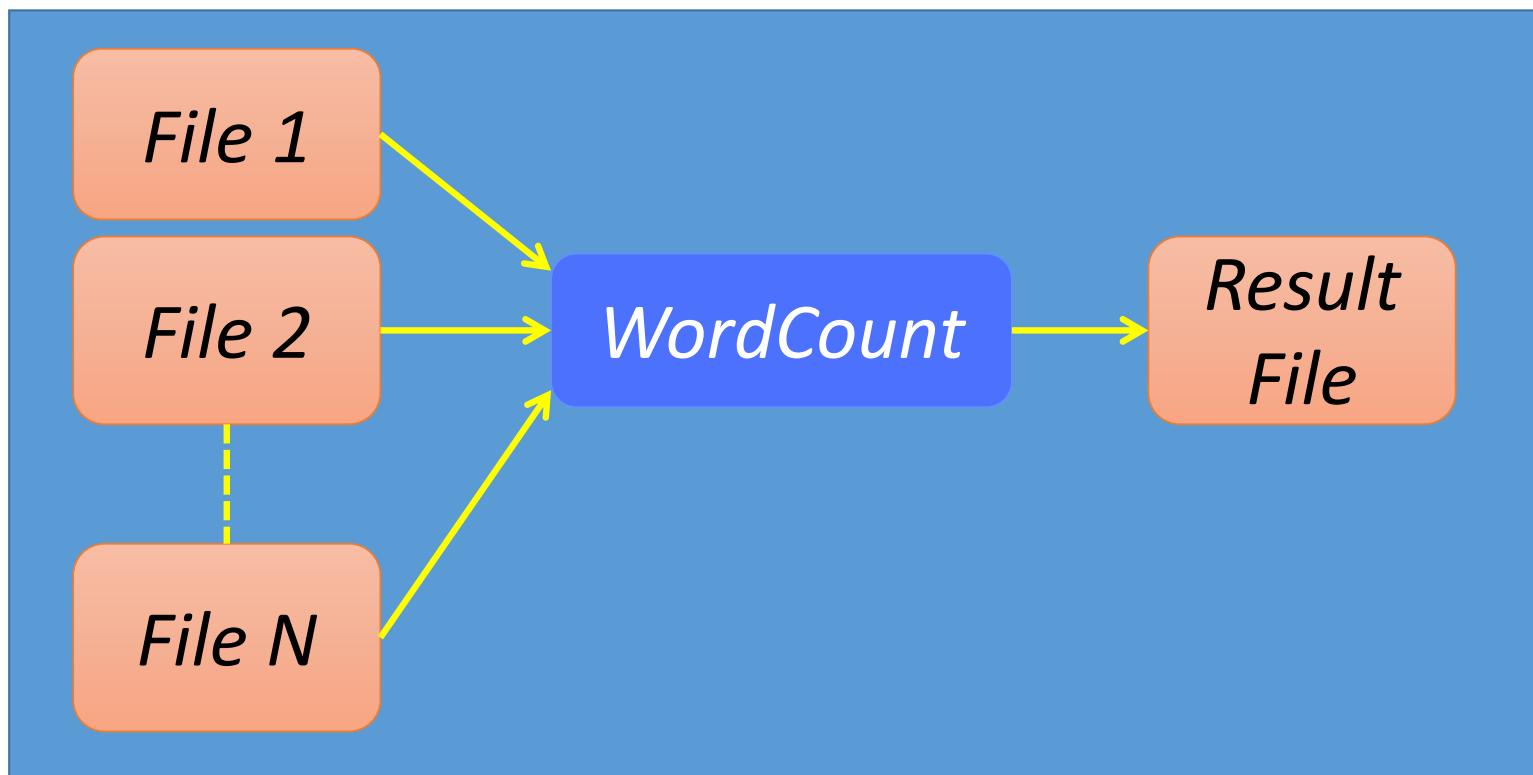
A UNIX pipe provides one-way communication
between two processes on the same computer

Batch Example Data Pipeline



Represents a large number of applications.

Example MapReduce Application: WordCount

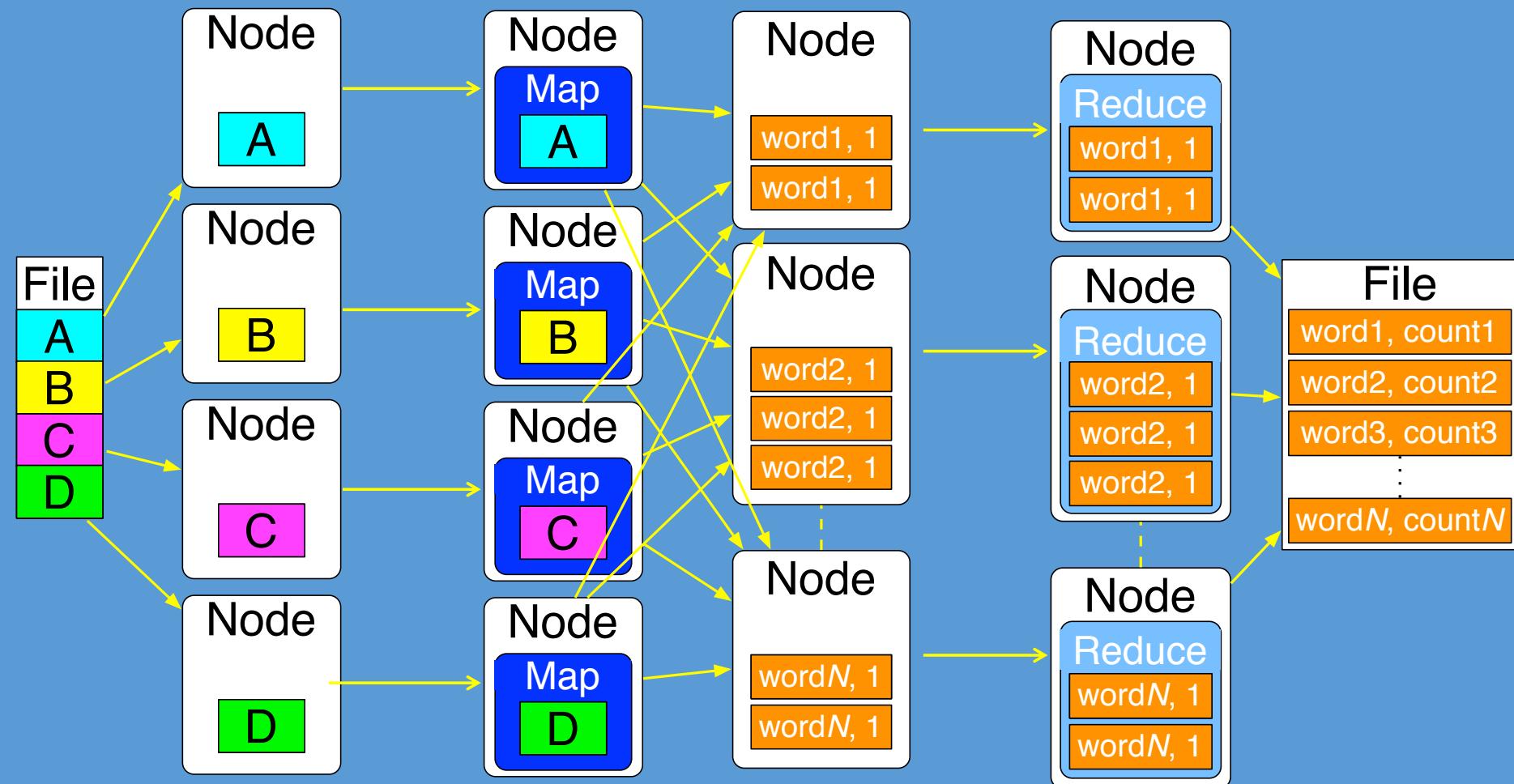


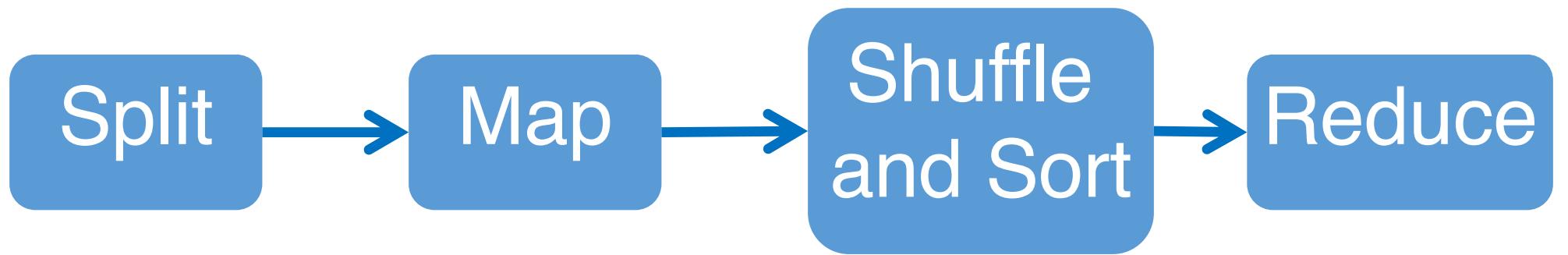
Step 1: Split

Step 2: Map

Step 3: Shuffle and Sort

Step 4: Reduce



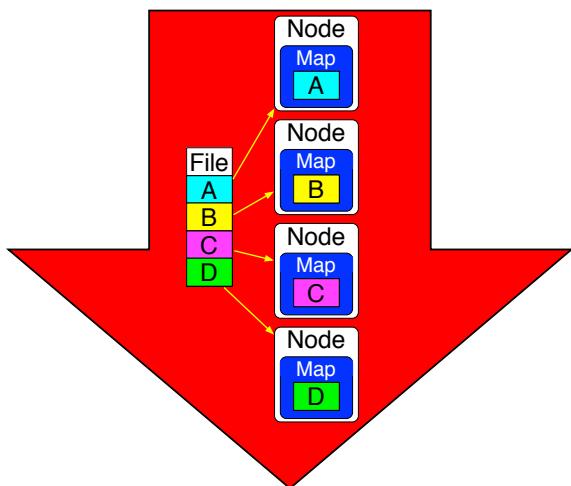




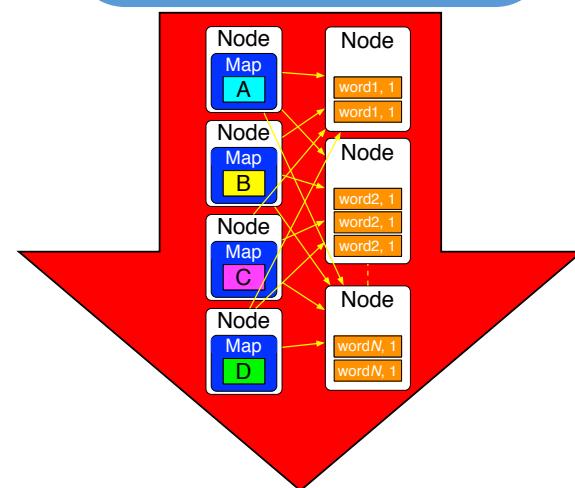
Map

Shuffle and Sort

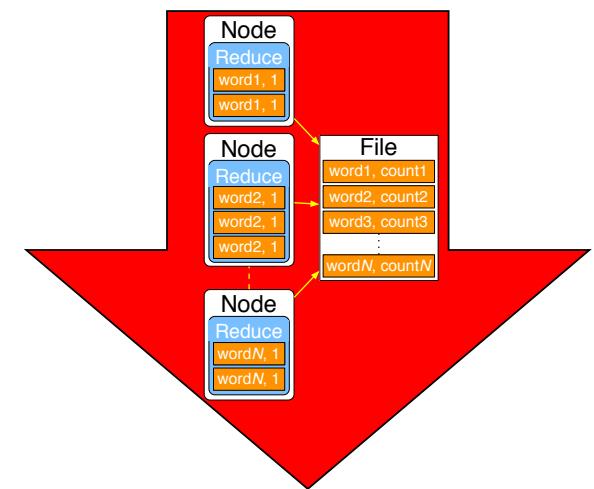
Reduce



Parallelization
over the input



Parallelization over
intermediate data



Parallelization
over data groups

Streaming Data

What is a Data Stream?

Data Stream

A possibly
unbounded
sequence of data
records

Timestamped

Geo-tagged

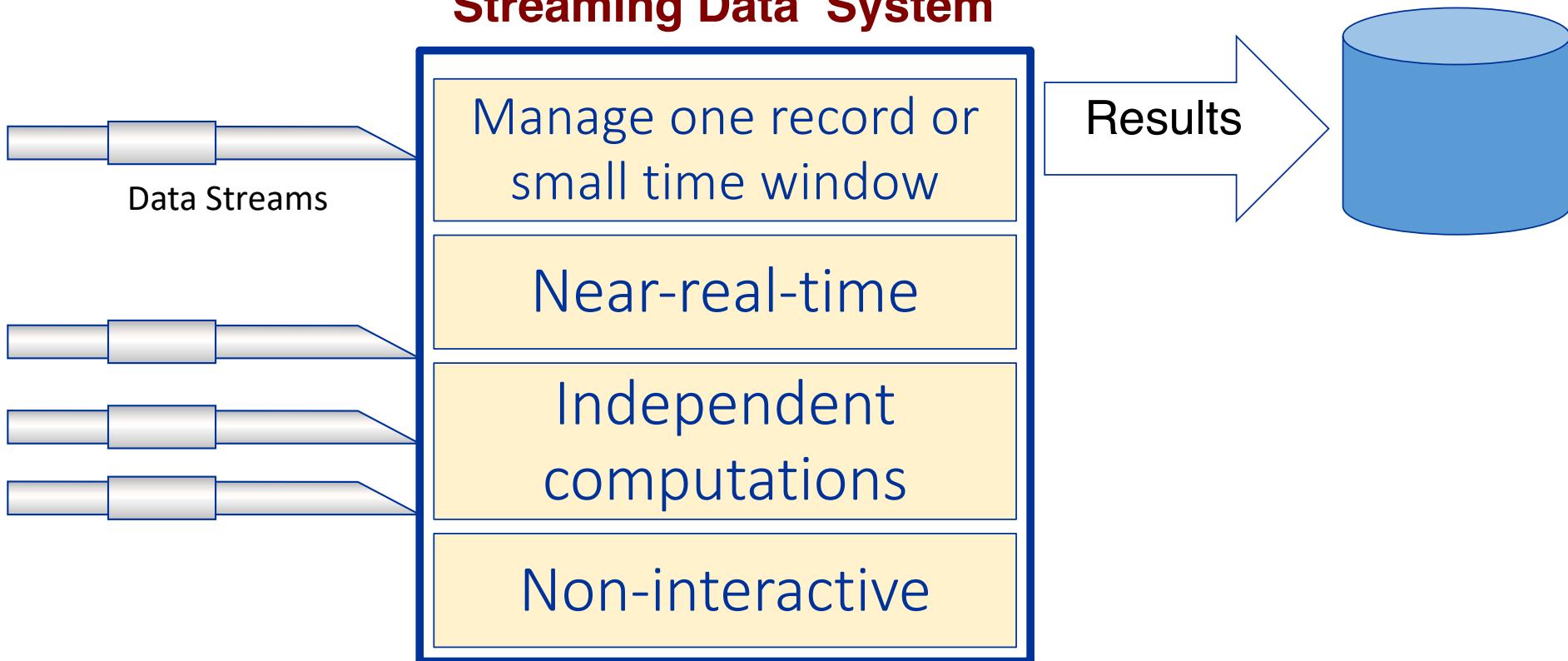


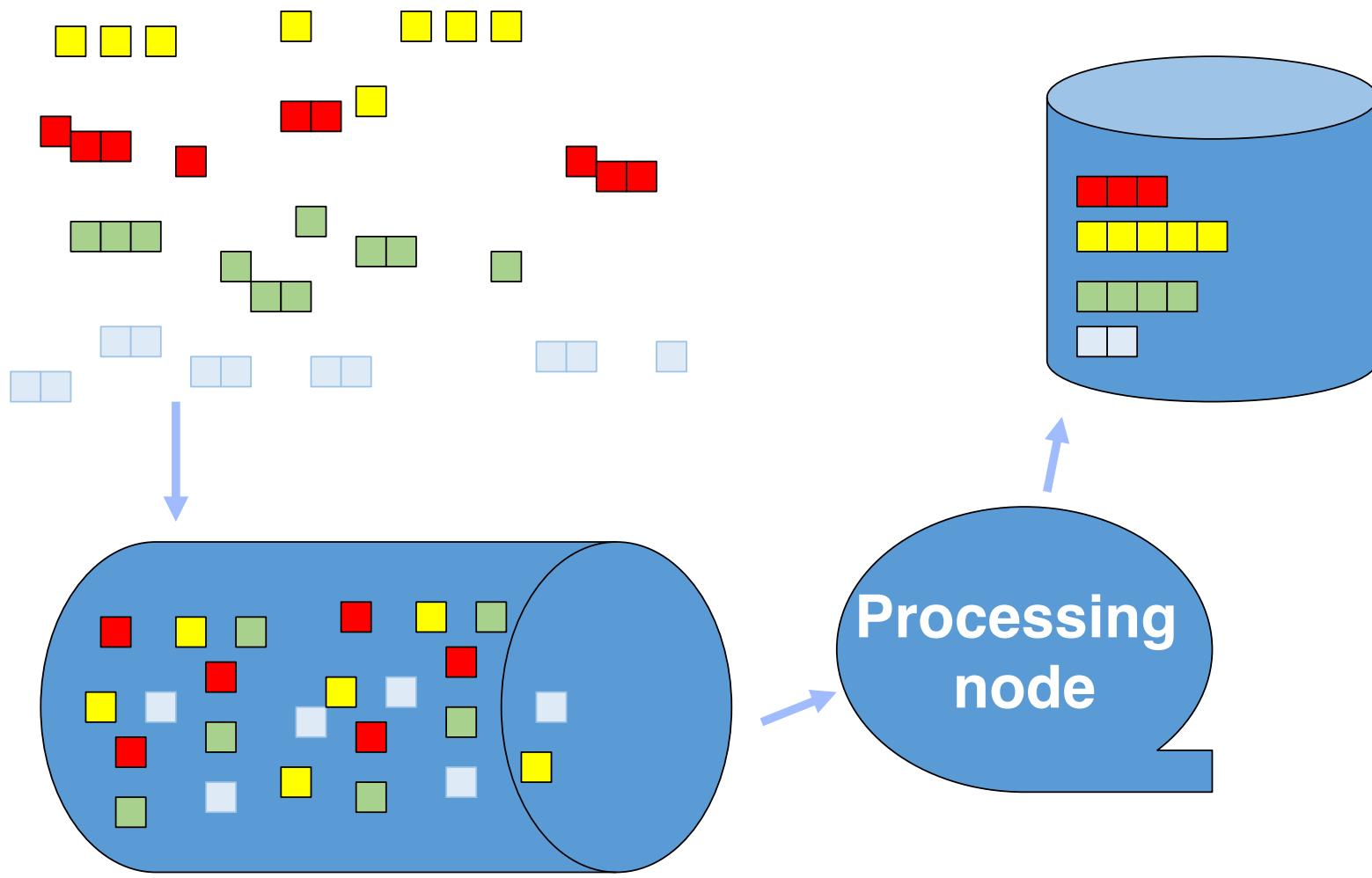
Data Stream

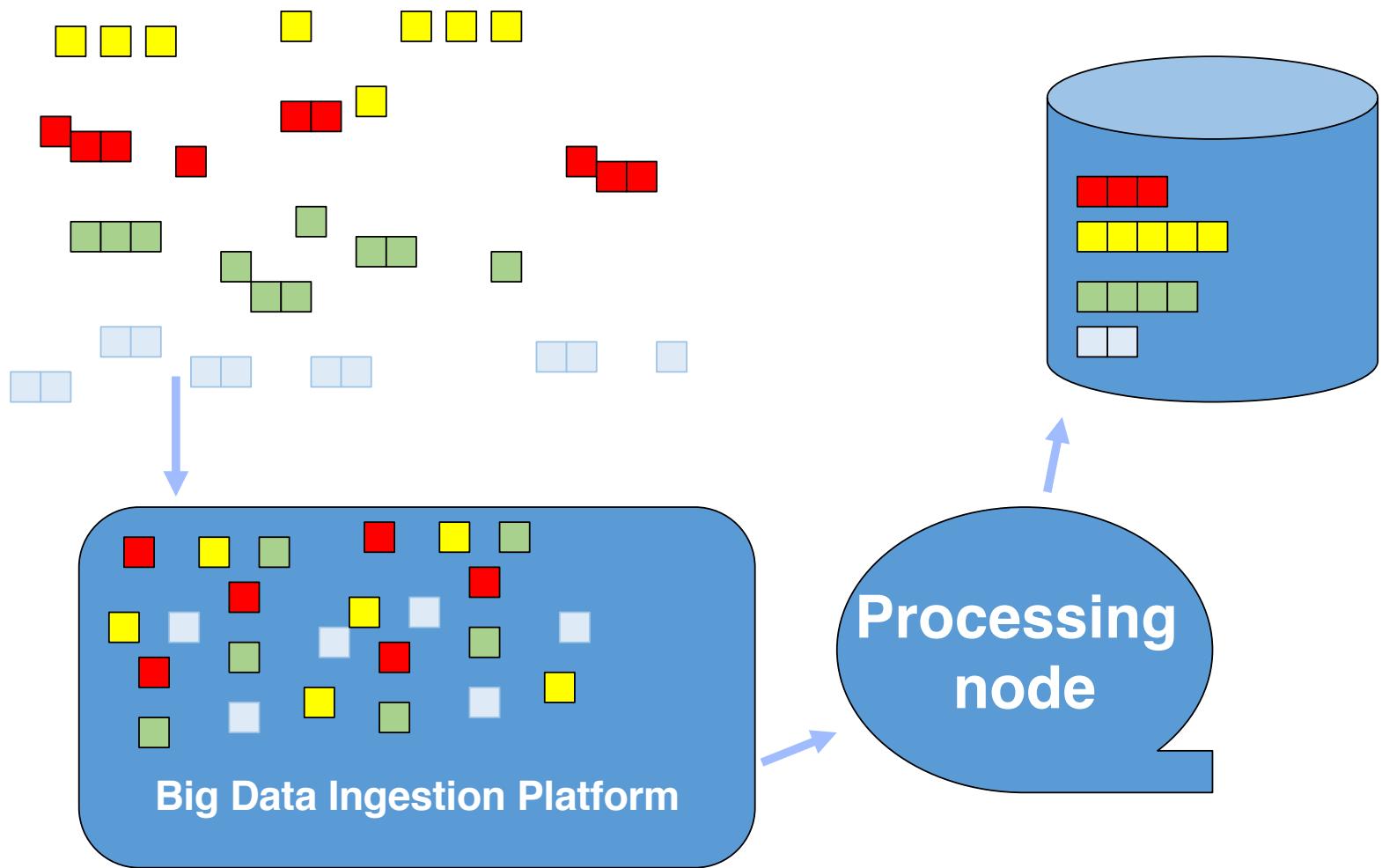
Synchronized
sequence of
events

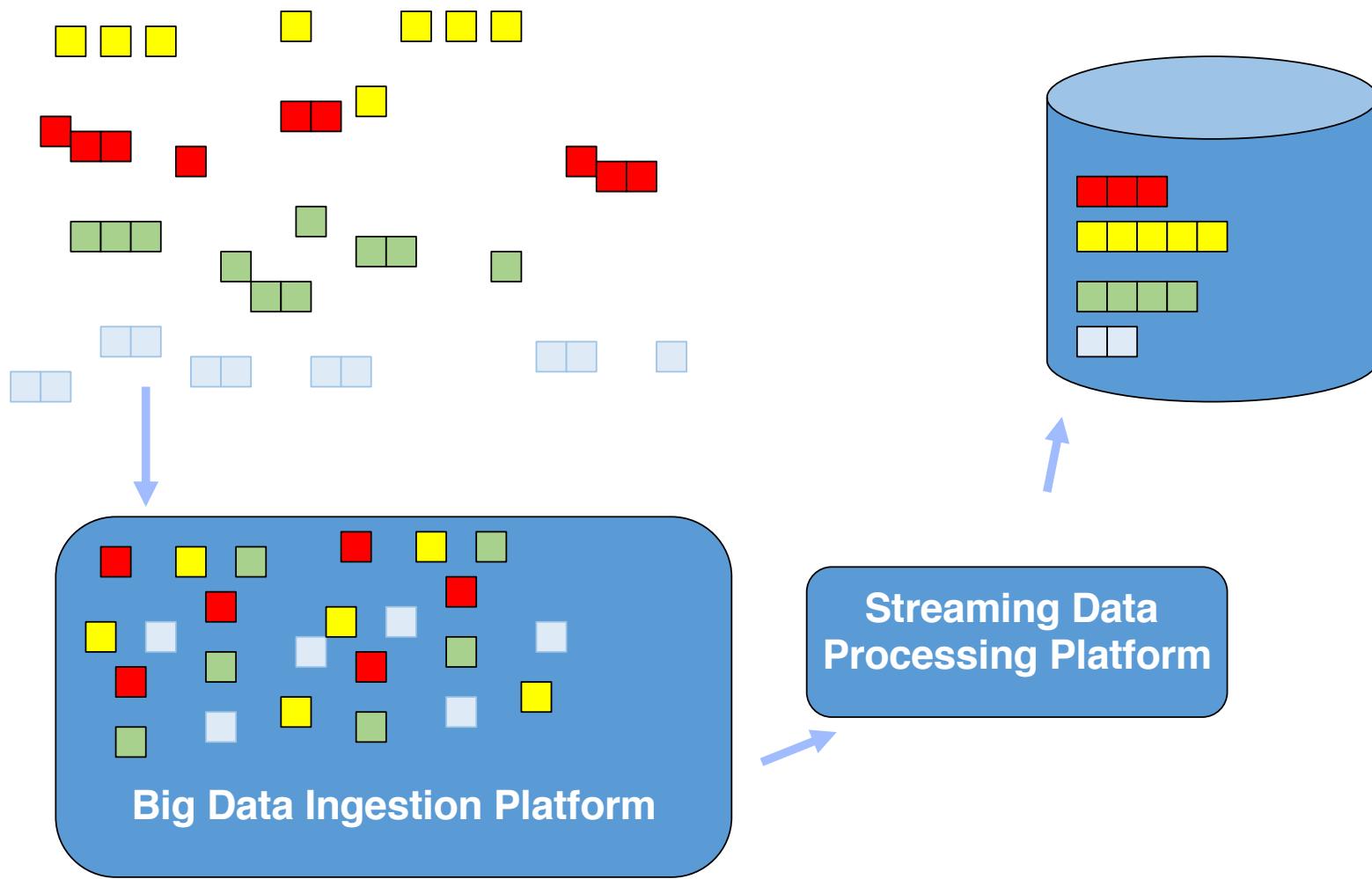
Streaming Data Systems

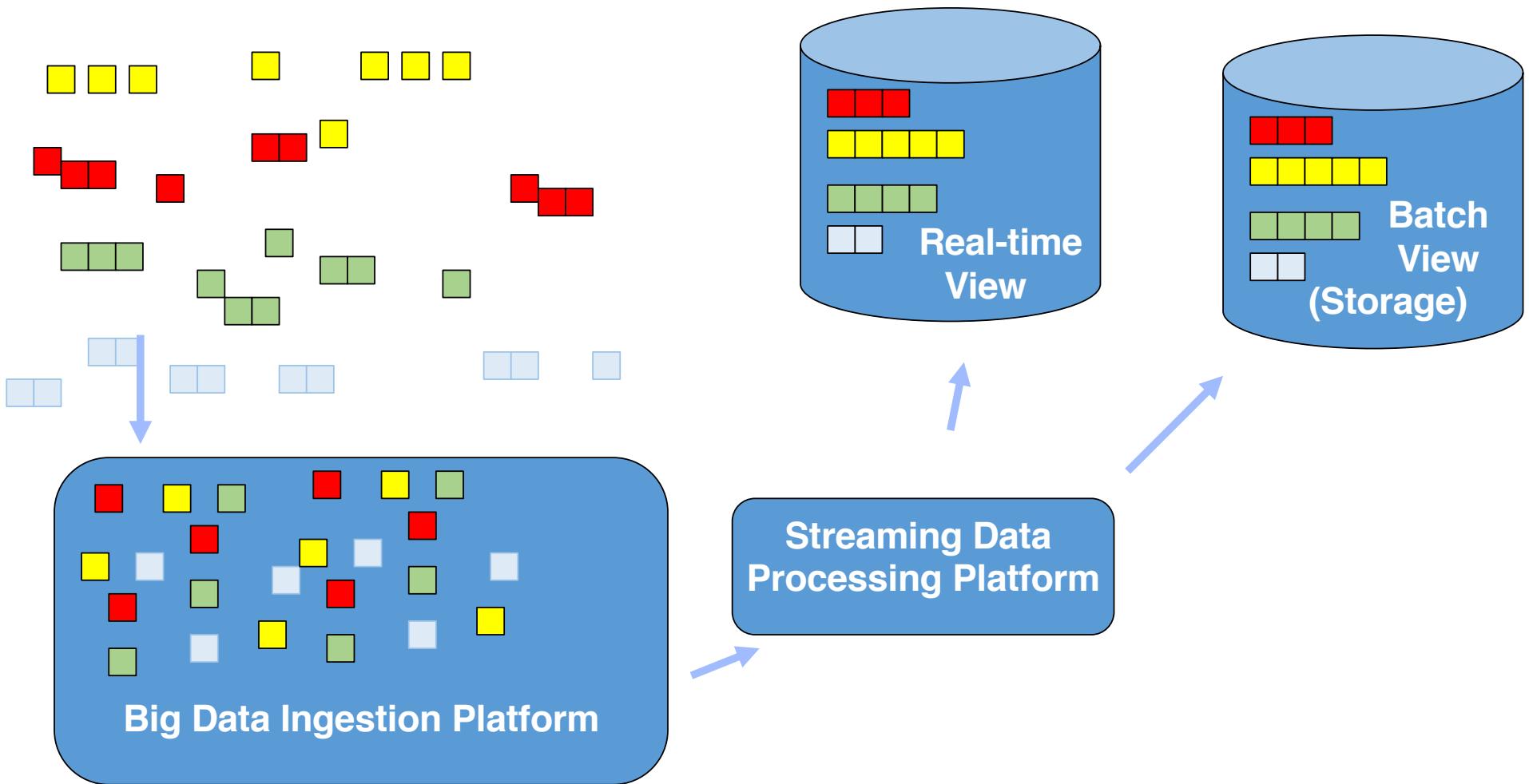
Streaming Data System











Databases vs. Data Lakes

Short definition of a data lake

Click Streams

Social Media

Sensor Data

Sales Transactions

Geo-locations



“ If you think of a datamart as a store of bottled water – cleansed and packaged and structured for easy consumption – the data lake is a large body of water in a more natural state. The contents of the data lake stream in from a source to fill the lake, and various users of the lake can come to examine, dive in, or take samples. ”

*James Dixon
CTO, Pentaho Corporation*

Load data from source



Store raw data



Add data model on read

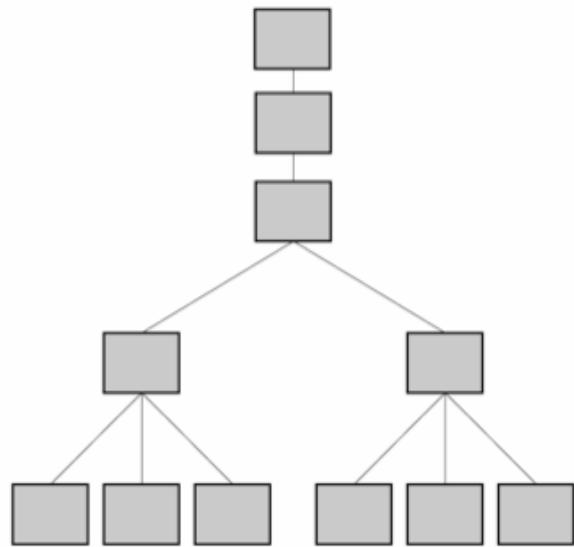
Data Lake

- Schema-on-read

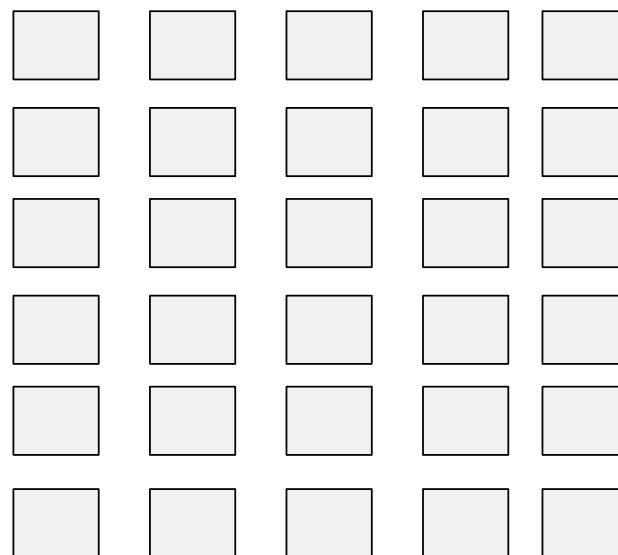
Data Warehouse

- Schema-on-write
- Transform and structure before load

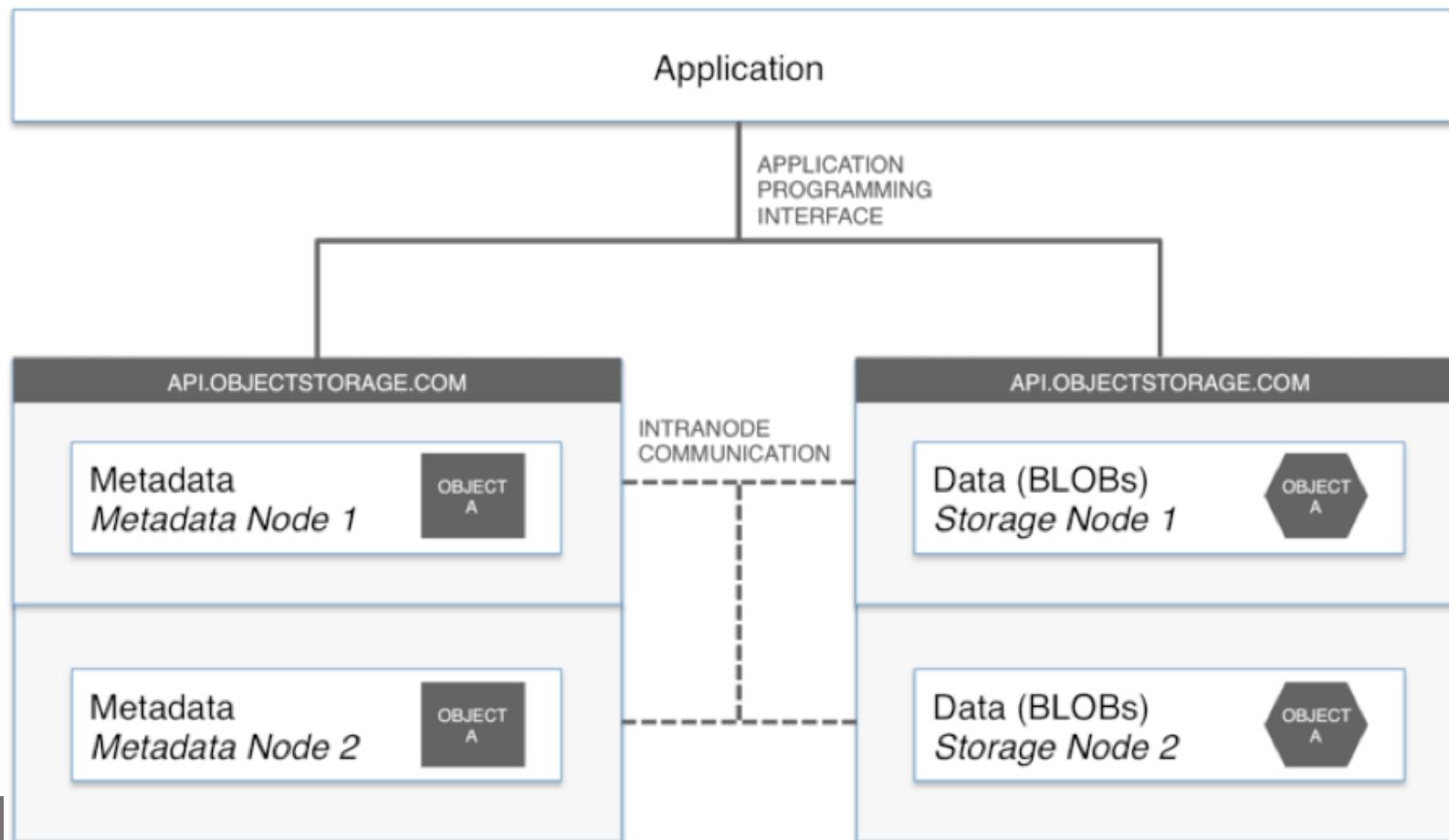
Hierarchical File System



Object storage



Object Storage



Questions?

Ilkay Altintas, Ph.D.
Email: ialtintas@ucsd.edu