# Intro to Big Data

Darryl Dutton

Email: darryl.dutton@t4g.com

# Agenda and Goal

Continue the conversation to the next level

- **Understand the characteristics of Big Data**
  - Why is Big Data different
  - Common Scenarios
  - Key Ingredients
- **How to address Big Data problems**
  - Storage
  - Processing
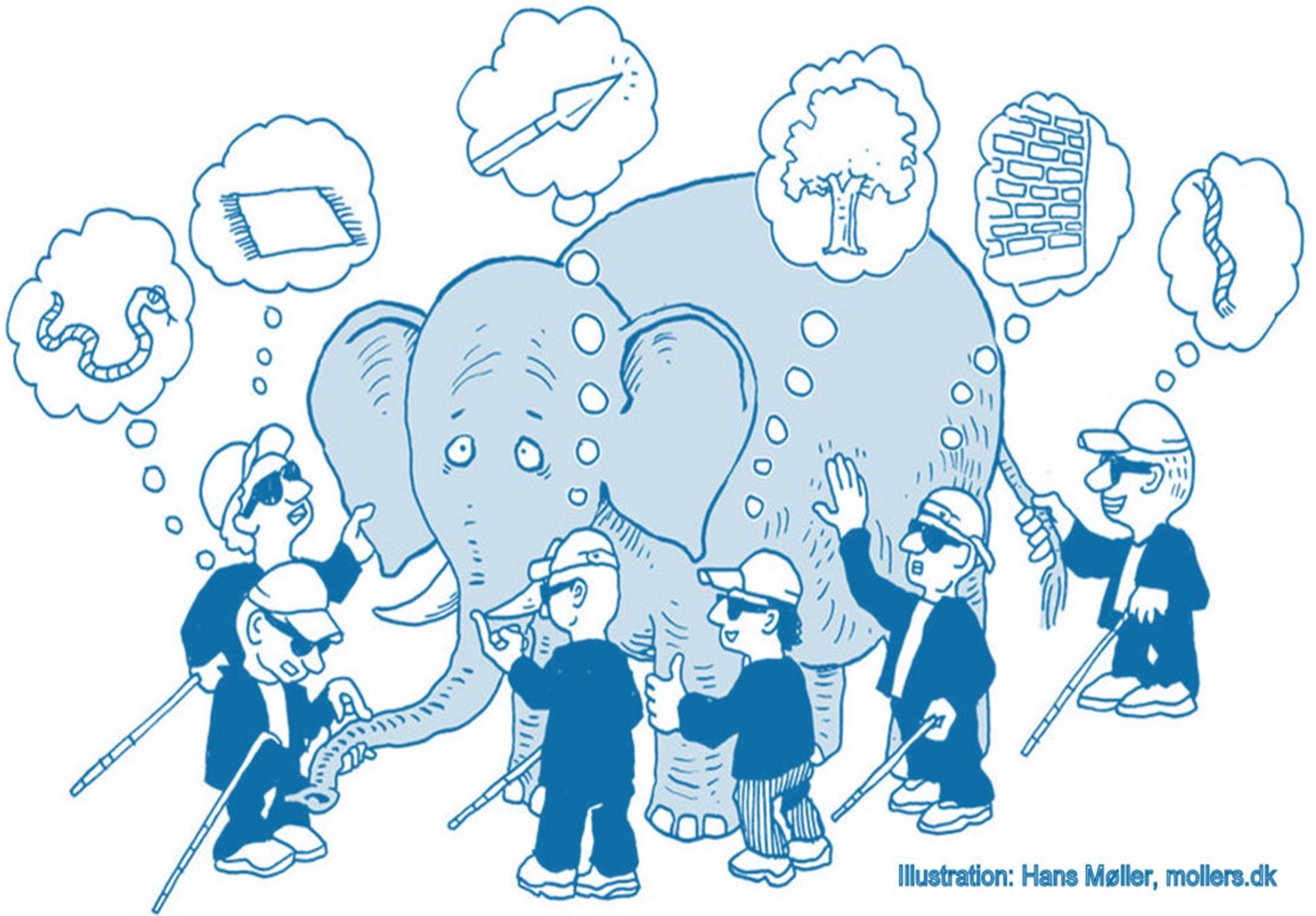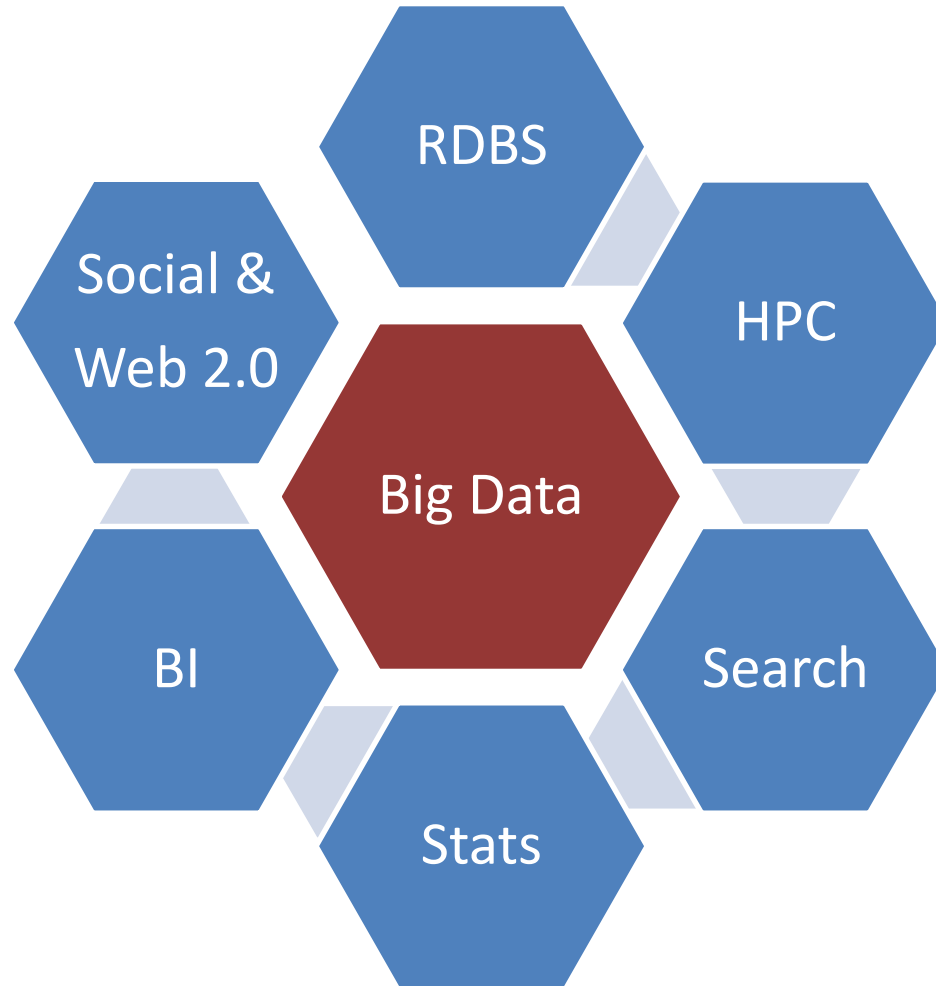  - Platform
  - Coding Pattern

Illustration: Hans Møller, mollers.dk

# Six Perspectives

# Why is Big Data different?

## How big is big….

- Users will generate over 300Gb per year
- Predicate to have about 8 ZETTABYTES data by 2015
  - 8,000,000,000,000,000,000,000 - $10^{21}$
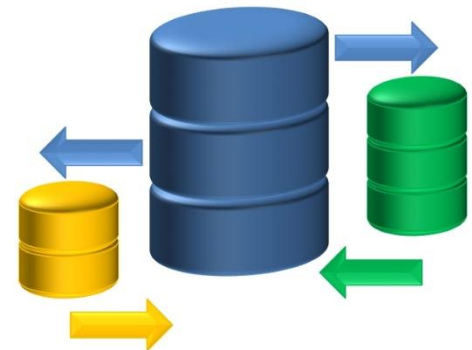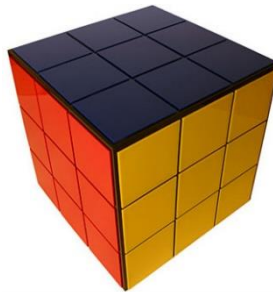
# Why is Big Data different?

The allure of finding useful information...

# Why is Big Data different?

I have my favorite hummer...

- The tools and approaches we know and love become awkward in this scale of data

- It is difficult to capture store, search, share and analyze at a certain point

- Most of this data is unstructured and it comes in oceans and streams

- More exploratory in nature instead of structured

# Why is Big Data different?

## The data sets and collection process is different

- Look outward to leverage external data
- "Collecting data first and ask questions later"
- Applying schema and context later in the process

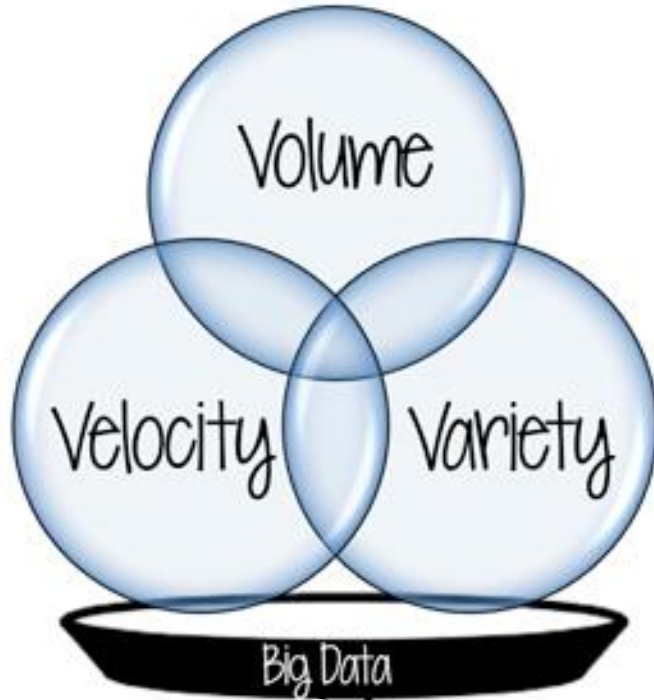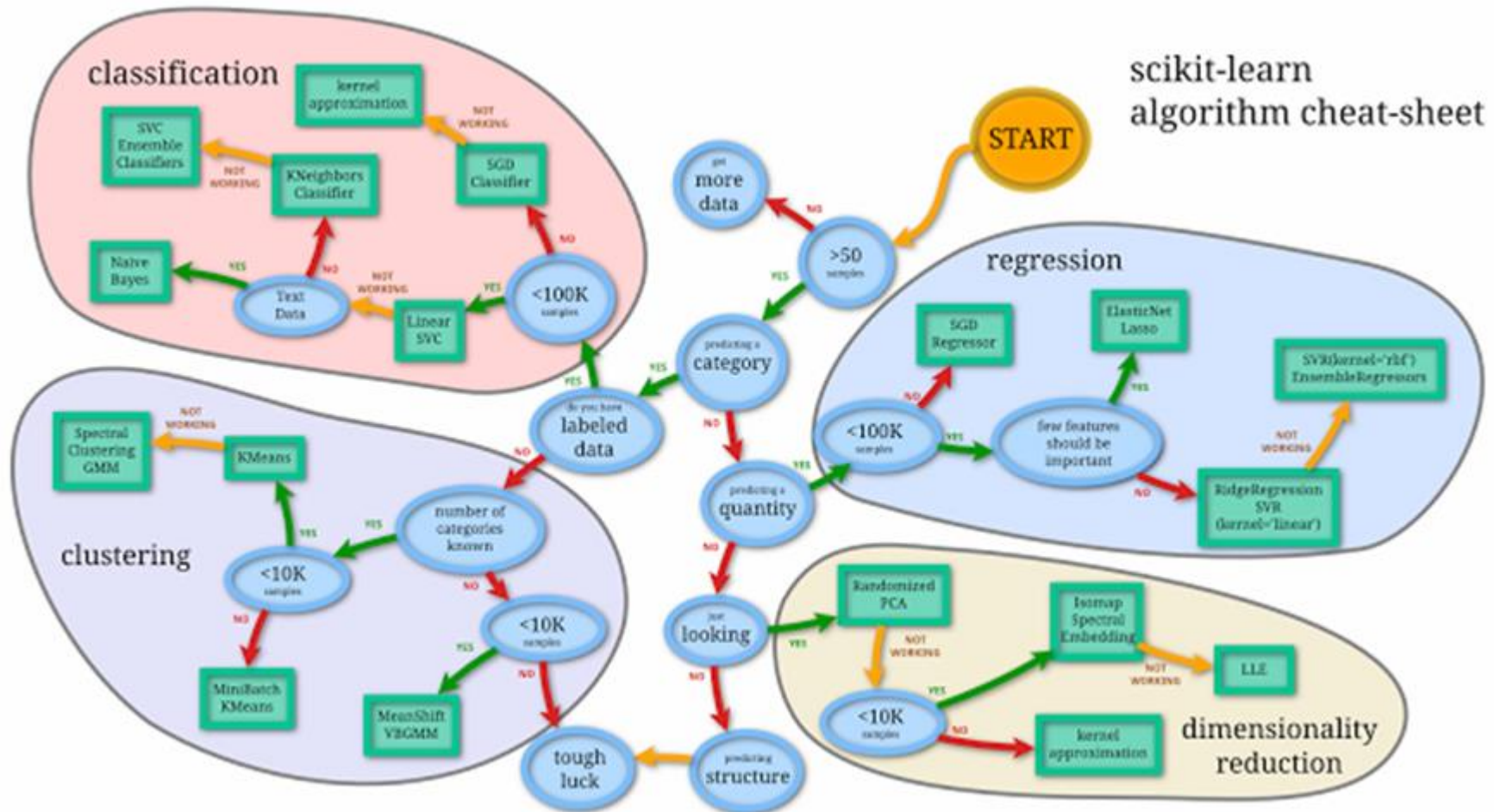# The Characteristics of Big Data

# Some of this is teaching old dog new tricks….



apologies to Disney

**T4G** The Intelligent Application of Technology™

# Big Data is exploratory by nature….



scikit-learn algorithm cheat-sheet

Source: Nishant Chandra http://n-chandra.blogspot.ca/

# Common Scenarios

- Behavioral Analysis
  - customer churn
  - trending
- Sentiment Analysis
  - social media analysis
  - multi channel analysis
- Recommendation Engines
  - cross-sell
  - up-sell
- Fraud Detection
  - clickstream analysis
  - mining

- Risk Mitigation
  - asset portfolio analysis
  - transactional log analysis
- Root Cause
  - network log
  - sensor data analysis
- Marketing Effectiveness
  - campaign ROI
  - Ad targeting
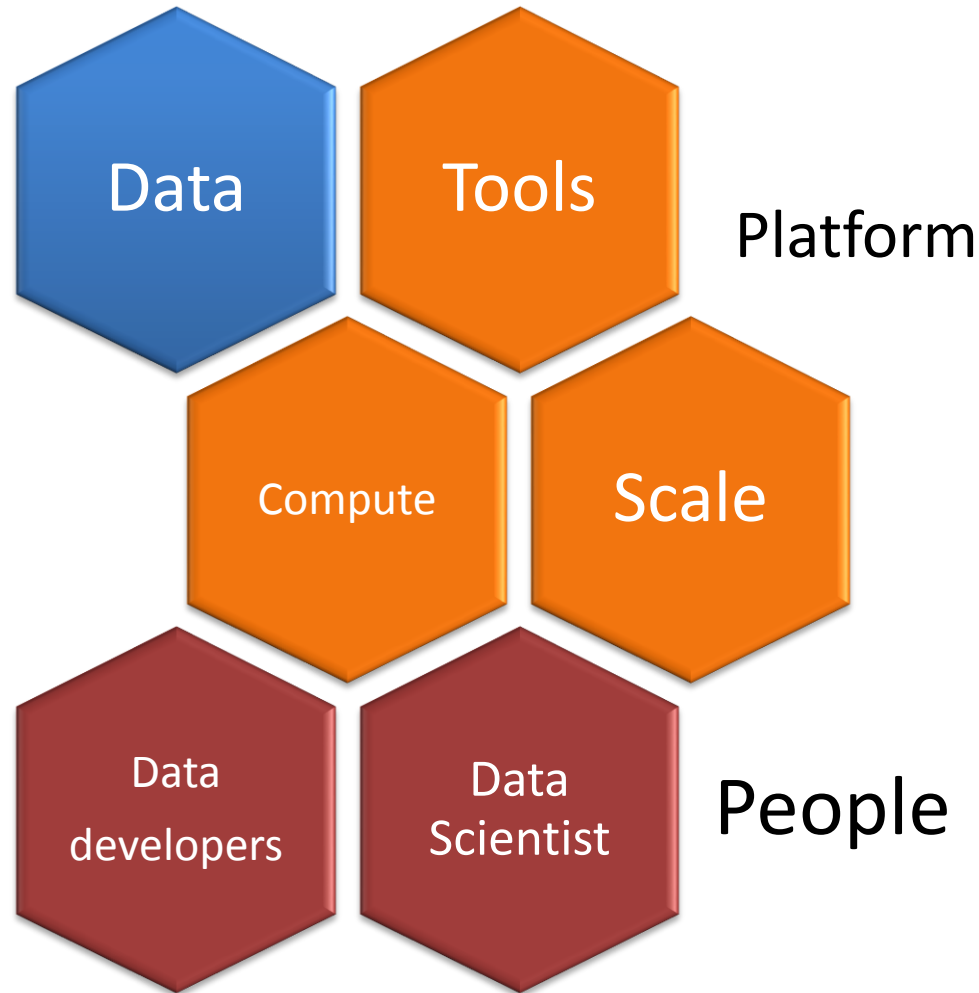
# Common Scenarios



Canada Flu Activity

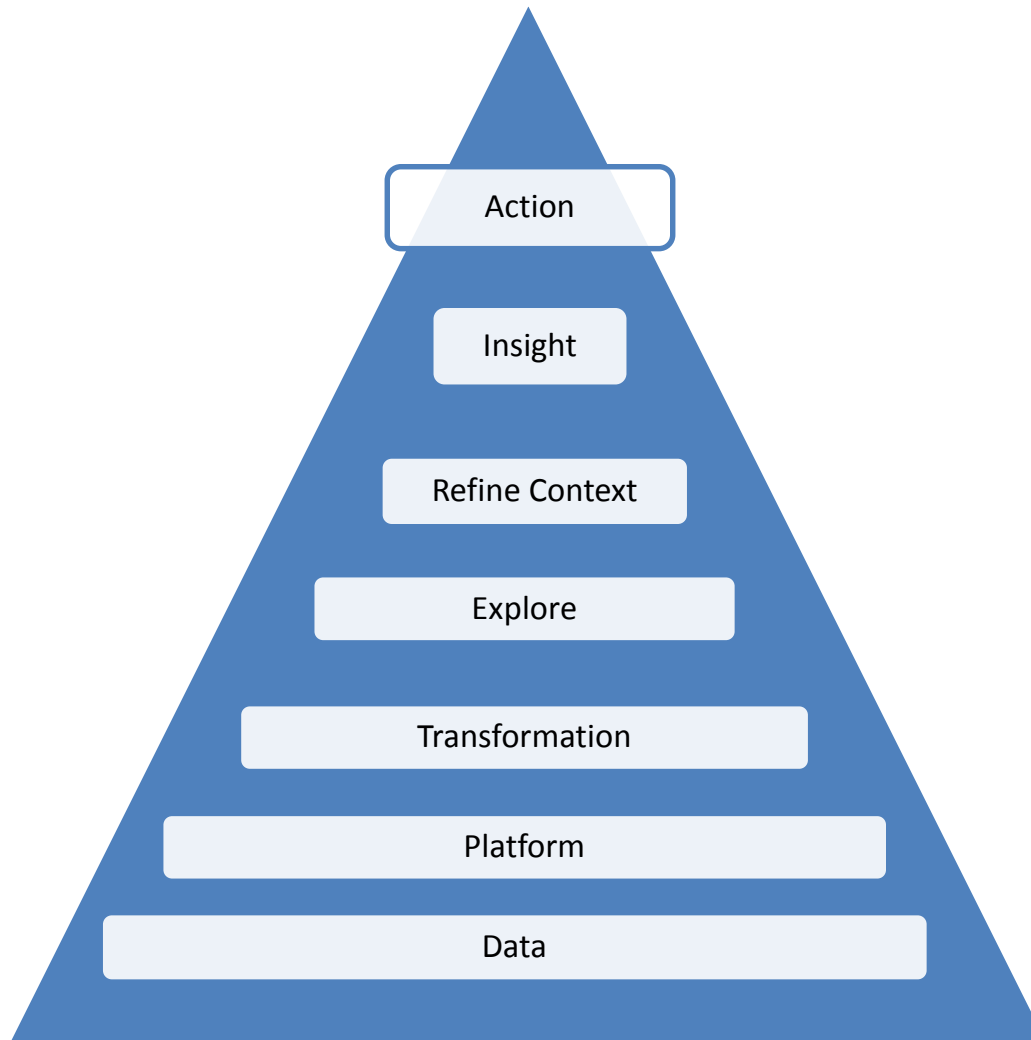Influenza estimate

● Google Flu Trends estimate  ● Canada data

6,644

4,983

3,322

1,661

2004   2005   2006   2007   2008   2009

Canada: Influenza-like illness (ILI) data provided publicly by the Public Health Agency of Canada.

# Key ingredients for a Big Data project



Data

Tools

Platform

Compute

Scale

Data developers

Data Scientist
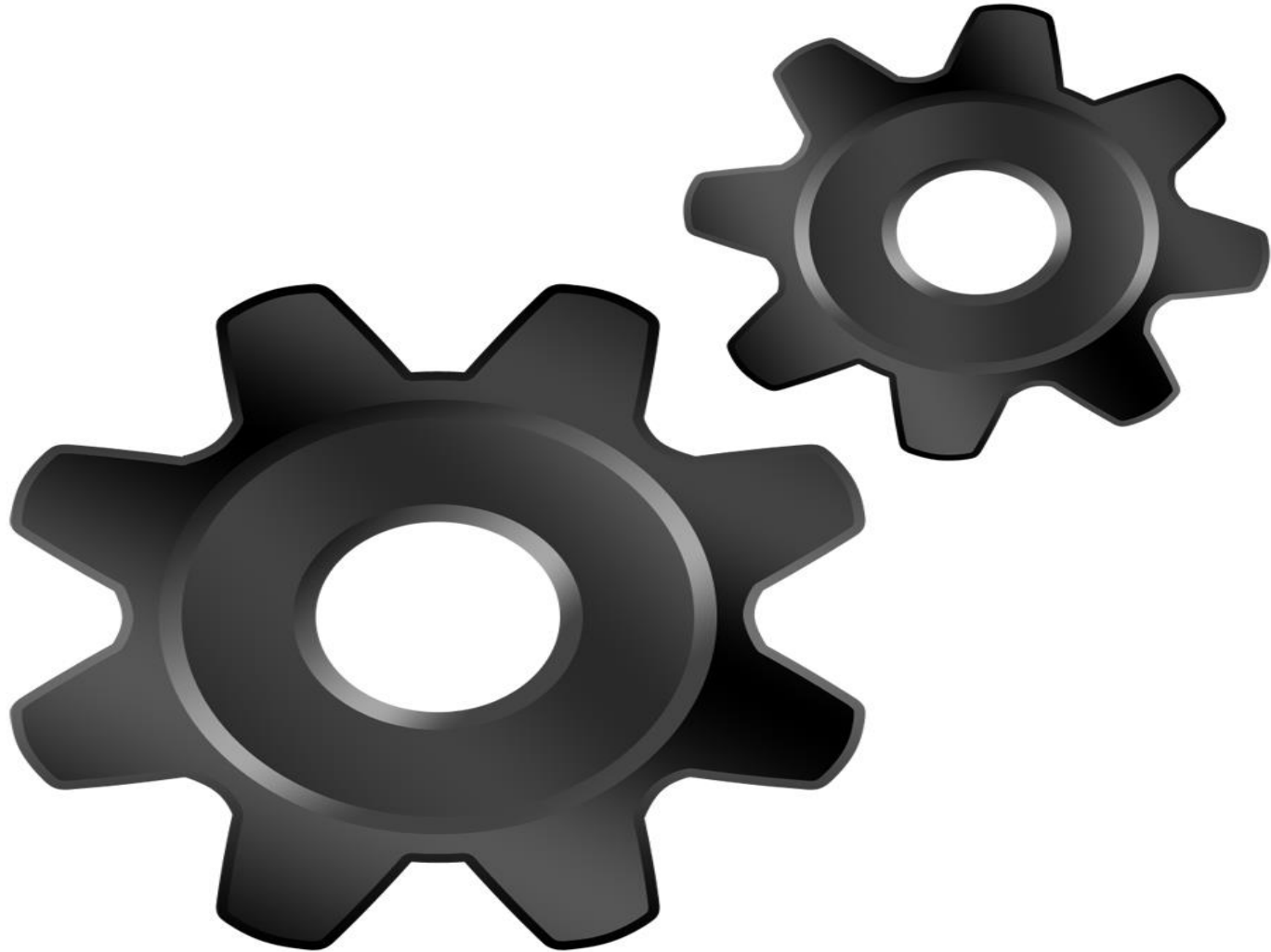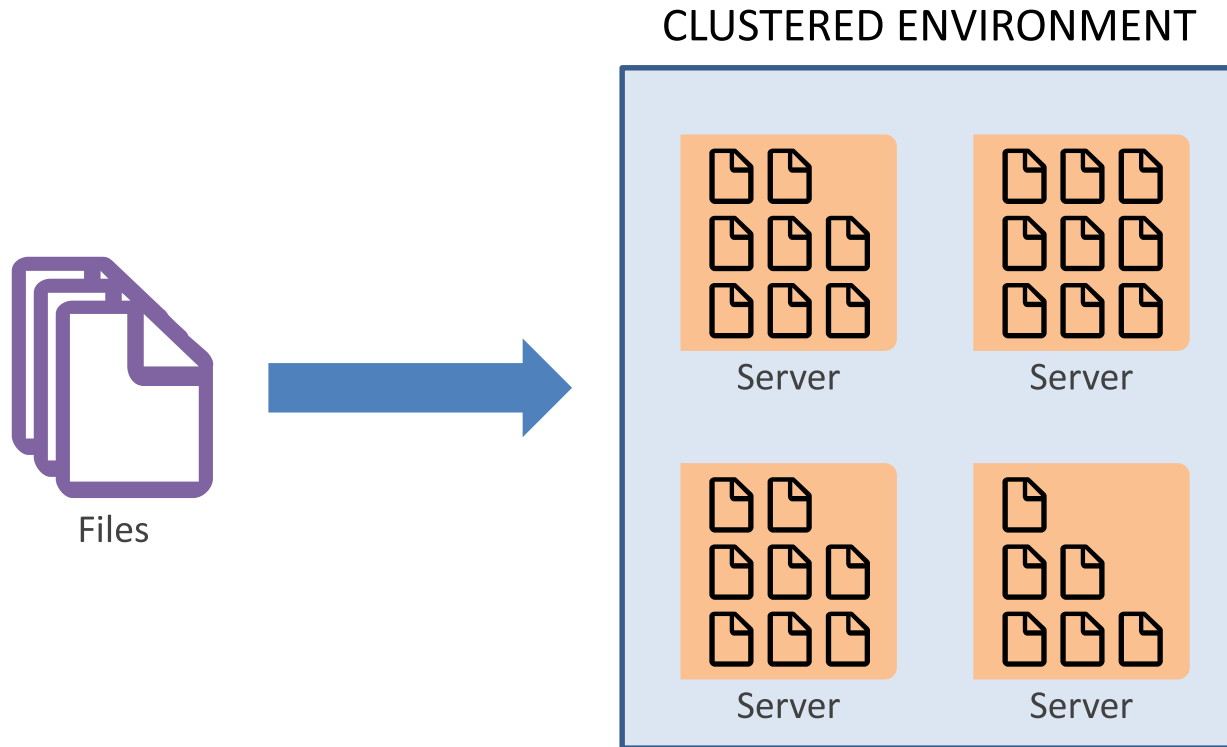
People

# Big Picture, it is about getting to action…

# Let's switch gears and focus on the 'how'

# How do we store big data?



CLUSTERED ENVIRONMENT

Files

Server

Server

Server

Server
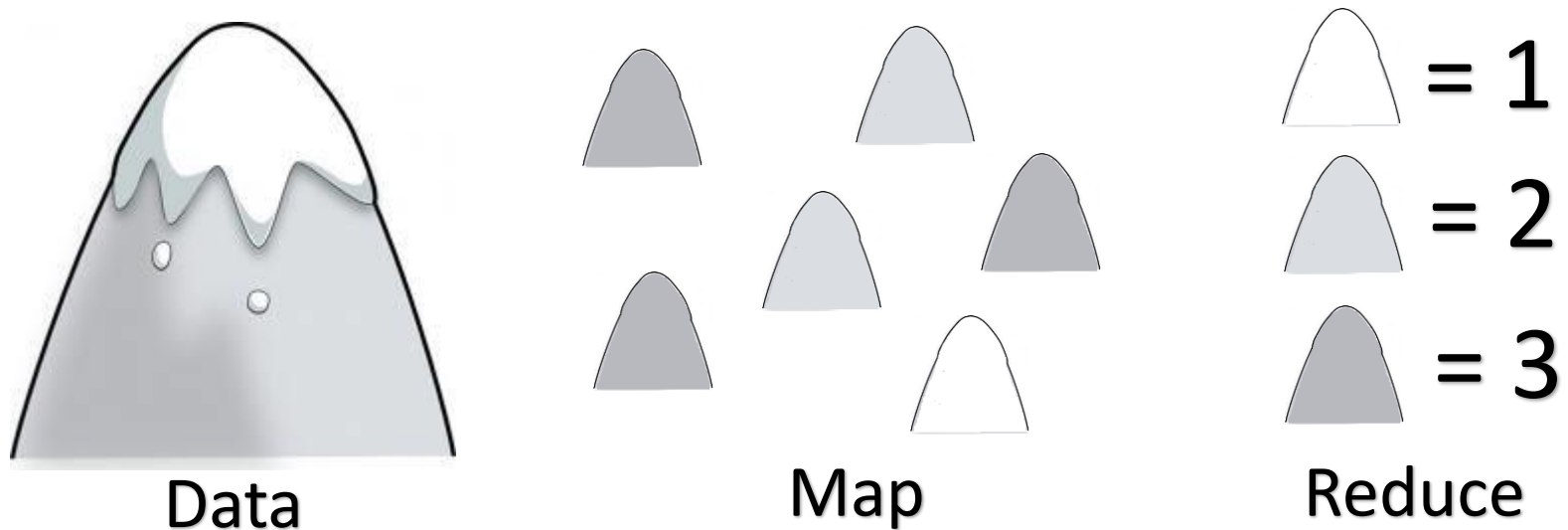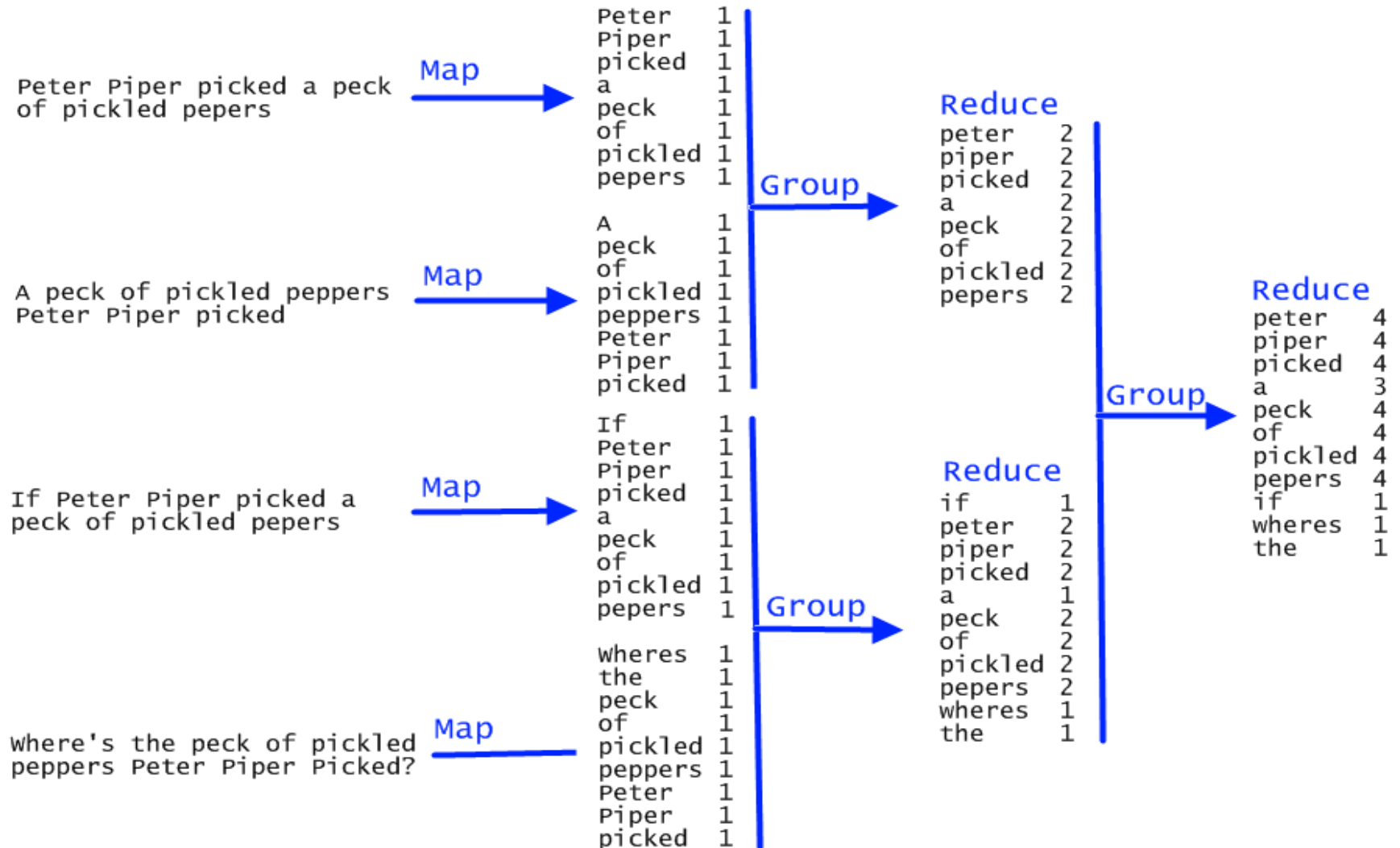
# How do we process big data?

## Map and Reduce - Simple 'divide and conquer'

- Help simplify the complexities of analyzing data

- Data structure is based on <key , value> pairs

- Map applies a function to each element in the list

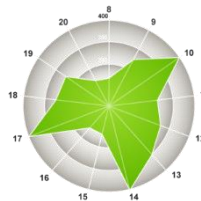- Reduce applies a function to elements with the same key to aggregate back to a smaller list
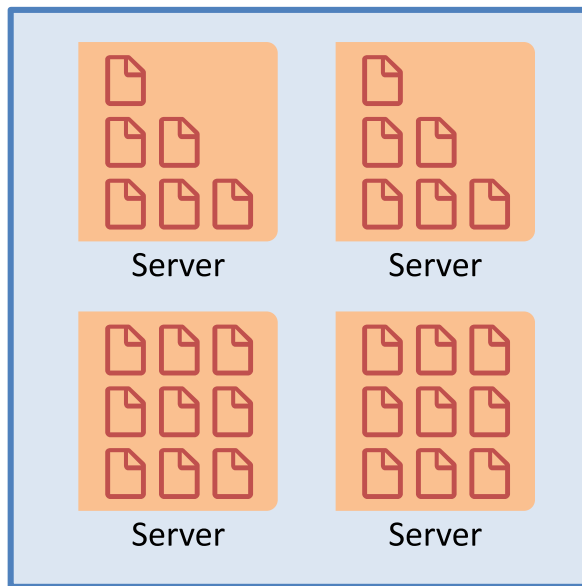


Data          Map          Reduce

= 1
= 2
= 3

# Simple MapReduce example – How many words?

# How do we process big data?

## Running our code….
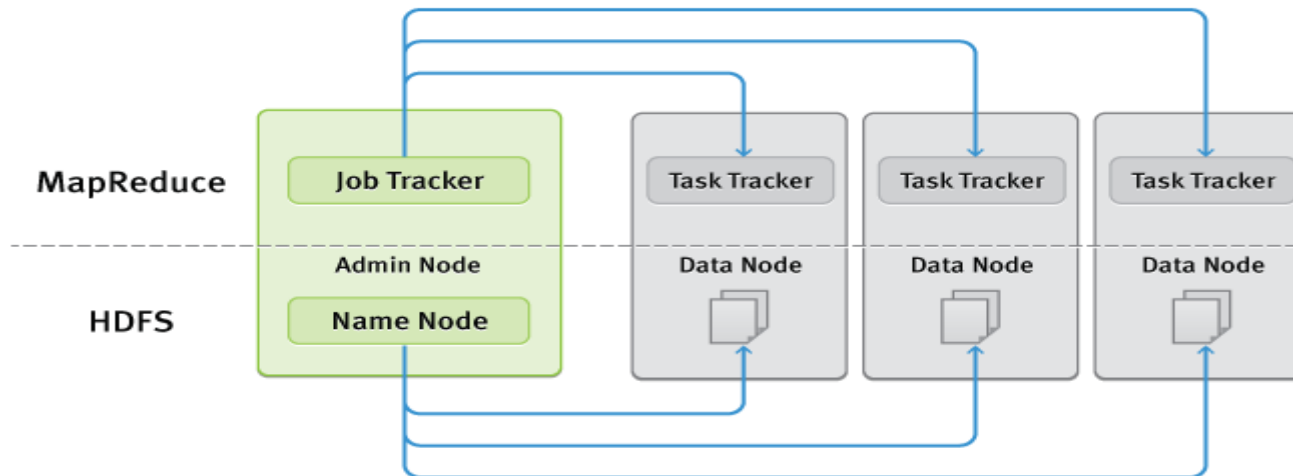
CLUSTERED ENVIRONMENT



Code

```
// Map Reduce function in JavaScript

var map = function (key, value, context) {
var words = value.split(/[^a-zA-Z]/);
for (var i = 0; i < words.length; i++) {
        if (words[i] !== "")
context.write(words[i].toLowerCase(),
1);}
}};

var reduce = function (key, values,
context) {
var sum = 0;
while (values.hasNext()) {
sum += parseInt(values.next());
    }
context.write(key, sum);
};
```
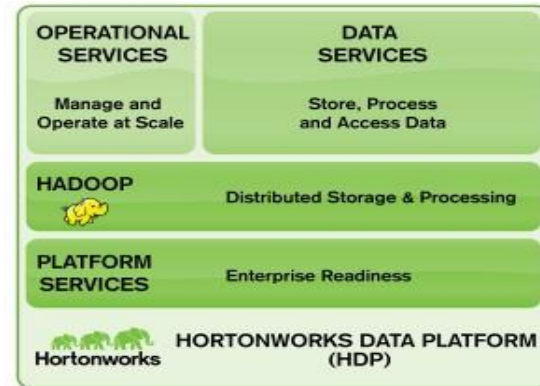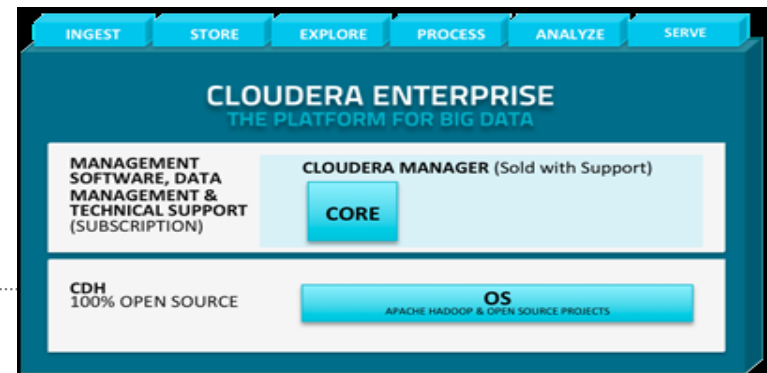
# What is one of frameworks we will use?



- Open source software framework that facilitates big data management and analysis

- Parallelizes data processing across many nodes (computers) in a compute cluster using 'batch' style jobs

- Parallelizing Map and Reduce functions (MapReduce)



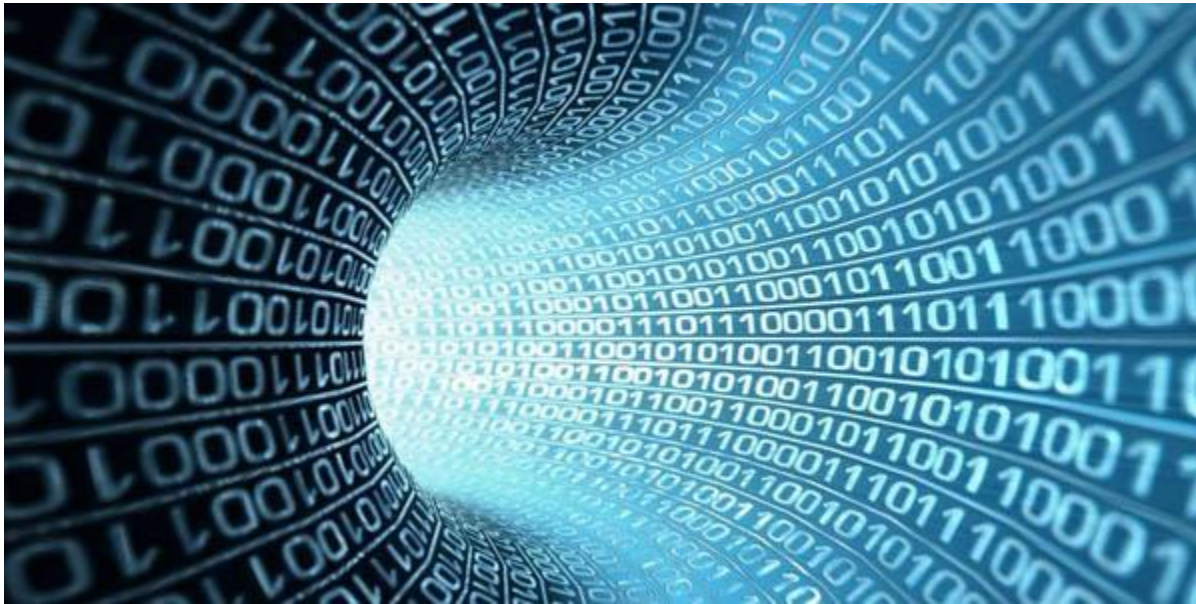| | | | | |
|---|---|---|---|---|
| MapReduce | Job Tracker | Task Tracker | Task Tracker | Task Tracker |
| | Admin Node | Data Node | Data Node | Data Node |
| HDFS | Name Node | | | |

# Our Big Data Partners in Hadoop…

# In Summary

The Big Data conversation will continue..

- Looked at **'why'** Big Data is different
- Briefly looked at **'how'**  Big Data problems are addressed
- The drive to insight and action!

# Thank You…

Darryl.Dutton@T4G.com