# Open Data Science
# Big Data and Infographics

http://training.theodi.org/BigInfographics

## Dr. Dave Tarrant
@davetaz – davetaz@theodi.org

# Big Data

Explain the current trends in big data

Identify a number of "big" datasets

Perform a number of short investigations with "big data"

# Exercise

What is Big Data to you?

# Big Data

Dataset that are too large and complex to manipulate with standard methods or tools.

# Excel

Workbook <span style="color:red">WAS</span> limited to 65,536 rows ($2^{16}$ aka 16-Bit)

64-Bit operating system addressing limit is $2^{64}$

18,446,744,073,709,551,615
q  q   t    b    m    t    h

# What is big data?

**V**olume
**V**elocity
**V**ariety
**V**eracity

# What is big data?

## **V**olume

Velocity

Variety

Veracity

We create around 4 zettabytes of data day.

That's 1 sextillion bytes per day (128-Bit OS required)

# Exercise

6,000,000 rows of data.

Visualise it in 10 minutes…

# Significance

Data stays on the web

We "download" the computer.

The computer is a cluster…

# Amazon public datasets

Explore the connection here to commodity computing and the volume problem.

Come back to the visualisation later

# What is big data?

**V**olume

**V**elocity

**V**ariety

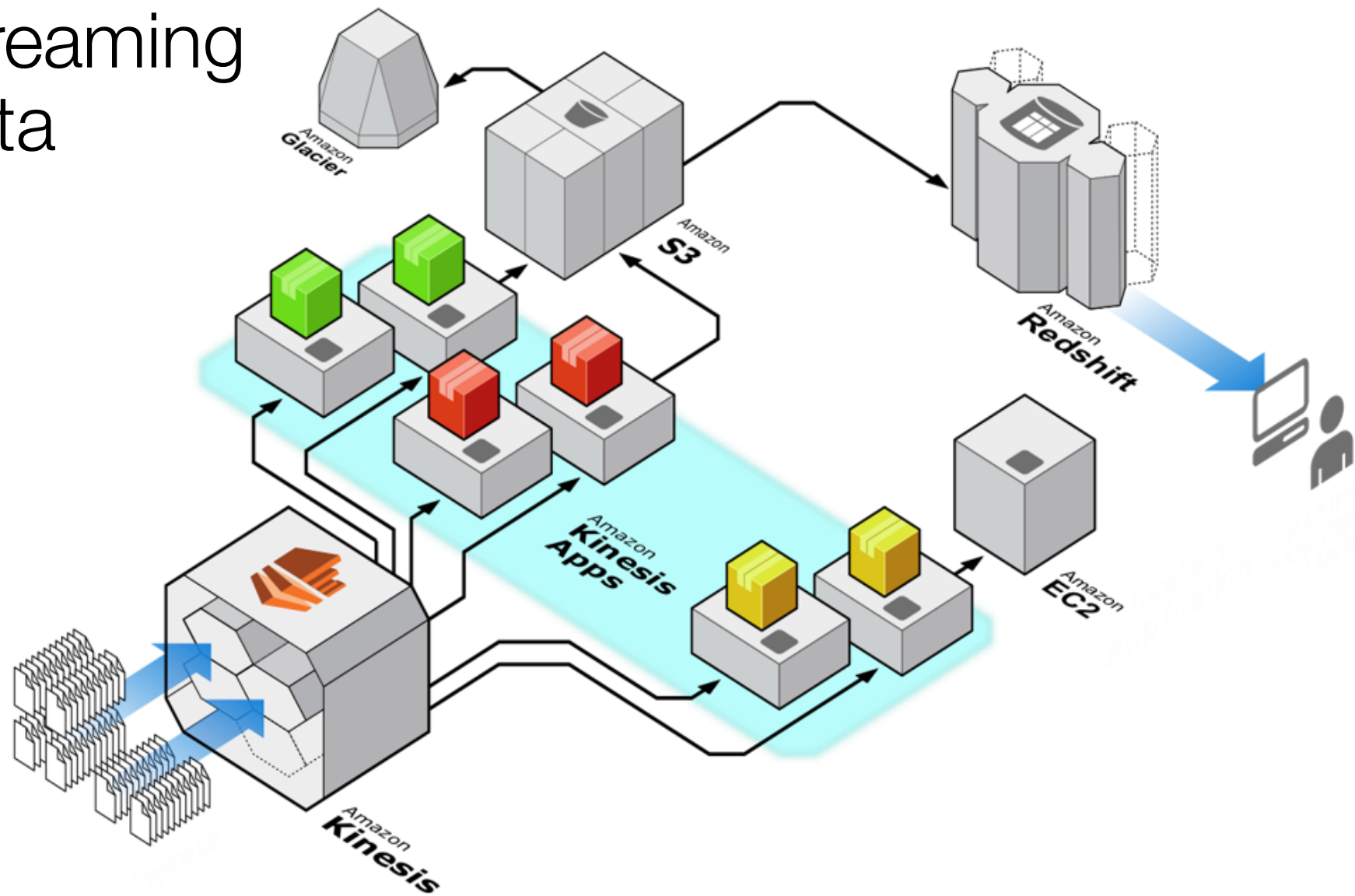**V**eracity

The data is created quicker than we can process it.

# Create a stream processor

You are running a chain of Sushi restairants. You want to get a live view on what people are eating, so you fix an RFID tag to the bottom of each bowl on the conveyor in every shop.

How do you get a live summary view of everythng being consumed? Where is the processing power required, and where isn't it?

# Streaming data

# What is big data?

**V**olume

**V**elocity

**V**ariety

**V**eracity

The data is continuously changing in structure, format and detail.

# Variety in simple data?

## Spend over £25,000 in the Foreign and Commonwealth Office

✉ API

Published by Foreign and Commonwealth Office.  Licensed under **OGL** Open Government Licence.
Openness rating: ★ ★ ★ ☆ ☆

A monthly-updated list of all financial transactions spending over £25,000 made by the Foreign and Commonwealth Office, as part of the Government's commitment to transparency in expenditure.

**Government Spending**

### DATA PACKAGE

Download a copy of all of the cached resources for this dataset

| ZIP | Download all data (523.3 kB) | ⬇ |
|-----|------------------------------|---|

### DATA RESOURCES (65 IN A TIME SERIES)

2015  View Less ⌃

# What is big data?

Volume

Velocity

Variety

The data quality is highly variable and affected by changing perception of truth and fact.

**V**eracity

# Big Data

Taken collectively. All digital data is big data. Looking at a facet might reveal that you are looking at a dataset that only conforms to one or two of the **V**s.

# A few more V's

## **V**alue and **V**iability

More data does not mean better results.

In fact often entirely the opposite is true.

Sample selection is critical to all good statistic studies.

Not being able to control selection may lead to an incorrect conclusion.

# Conclusion

## The majority of datasets are large.

Lots of rows with lots of joins that can be processed. If you know how to exploit computing power available.

# Scaling

# Computing clusters

Injestors
Translators
Indexors
Caches
Validators
Data stores
Visualisers
General purpose units

# British Library (in 2008)

80 terabytes of digitised newspapers,

60 terabytes of web-harvested information, e-journals and books,

25 million pages of digitised C19 literature, broadcast television, digital video and digital maps.


**GOAL:** All 80Tb of newspaper images migrated from TIFF to JP2000

http://www.planets-project.eu/docs/newsletters/Planetarium7_July09.pdf

# Exercise

Design a system to help the BL with their goal?

How long will it take?

**GOAL**: All 80Tb of newspaper images migrated from TIFF to JP2000

# Cloud computing

```
                    ┌──────────┐
                    │  Client  │
                    └────┬─────┘
                         │
                         ▼
                   ┌────────────┐
                   │ Job Tracker│
                   └──┬──────┬──┘
              ┌───────┘      └───────┐
              ▼                      ▼
    ┌──────────────┐        ┌──────────────┐
    │ Task Tracker │        │ Task Tracker │
    └──┬────┬──────┘        └────┬──────┬──┘
       │    │     ┌──────────┐   │      │
       │    └────►│Name Node │◄──┘      │
       │          └──────────┘          │
       ▼          ▼          ▼          ▼
  ┌─────────┐┌─────────┐┌─────────┐┌─────────┐
  │Data Node││Data Node││Data Node││Data Node│
  └─────────┘└─────────┘└─────────┘└─────────┘
```

# Cloud computing

# Open corporates

Separate data store from website
API vs Human readable

```
                    ┌─────────────────────┐
                    │  Human Web Server   │
                    └─────────────────────┘
          ┌──────────────────────────────────────────┐
          │              Load Balancer               │
          └──────────────────────────────────────────┘
     ┌──────────────┐  ┌──────────────┐  ┌──────────────┐
     │ Machine API  │  │ Machine API  │  │ Machine API  │
     │   Server     │  │   Server     │  │   Server     │
     └──────────────┘  └──────────────┘  └──────────────┘
     ┌────────────────────────────────────────────────────┐
     │                  Task/Job tracker                  │
     └────────────────────────────────────────────────────┘
   ┌───────────┐  ┌───────────┐  ┌───────────┐  ┌───────────┐
   │ Data node │  │ Data node │  │ Data node │  │ Data node │
   └───────────┘  └───────────┘  └───────────┘  └───────────┘
```

# Socrata model

| Org #1 web site | Org #2 web site | Org #3 web site |
|---|---|---|

| Load Balancer |
|---|

| Machine API Server | Machine API Server | Machine API Server |
|---|---|---|

| Task/Job tracker |
|---|

| Data node | Data node | Data node | Data node |
|---|---|---|---|

# Databases



Data node

Figure 1: Entity Relationship Diagram (ERD) for the aer_sight database

# Flat = fast

- Something about noSQL and differences

- Using big query…

- But first the dataset..
.

# UK Trade data

# UK Trade Data

**Exports**
Non-EU
150,000 to 200,000
per month

**Imports**
Non-EU
190,000 to 220,000
per month

**Dispatches**
EU
210,000 to 250,000
per month (+estimates)

**Arrivals**
EU
125,000 to 135,000
per month (+estimates)

# Distilled information



| | Imports to uk | Exports from uk |
|---|---|---|
| UK Imports & Exports | | |
| **2009** | | |
| All commodities | | |

United States Of America
Germany
Netherlands
France
People's Republic Of China
Republic Of Ireland
Belgium
Italy
Spain
Norway
Switzerland
Hong Kong
Canada
Japan
South Africa
Sweden
Poland
India
Russia

**Netherlands**
Imports: £21,499,133,940
Exports: £17,554,538,157
Combined: £39,053,672,097
Net: £-3,944,595,783

# Exercise

Q: How have imports and exports on Jet Engines changed over the years?

Design a processing pipeline that can answer this question from the data.

bit.ly/uk_trade

# Stage 1: What the format????

# Stage 2: RTFM

**Table of Contents:**

# Stage 3: Decode 010110100

# Stage 3b: API? 010110100

`https://www.gov.uk/trade-tariff/headings/0101?country=&day=1&month=6&year=2009`



The codes for the same things have changed. Meaning that we have to compare the text! Ahhh.

# Stage 4: API for data?

`https://www.gov.uk/trade-tariff/headings/0101`**`.json`**`?country=&day=1&month=6&year=2009`

```
{
    "goods_nomenclature_item_id": "0101000000",
    "description": "Live horses, asses, mules and hinnies",
    "bti_url": "http://ec.europa.eu/taxation_customs/dds2/ebti/ebti_consultation.jsp?Lang=en&nomenc=0101000000&Expand=true",
    "formatted_description": "Live horses, asses, mules and hinnies",
    "_response_info": {
        "links": [
            {
                "rel": "self",
                "href": "/trade-tariff/headings/0101.json"
            },
            {
                "rel": "chapter",
                "href": "/trade-tariff/chapters/01"
            },
            {
                "rel": "section",
                "href": "/trade-tariff/sections/1"
            }
        ]
    },
    "chapter": {
        "goods_nomenclature_item_id": "0100000000",
```

# Stage 5: Predict scale

(12 * 4) files per year

12 Comcode tables

12 Portcode tables


To answer one query you may have to join 48 tables to 24 others to answer it.

This is not how map reduce and big data work.

# Large databases
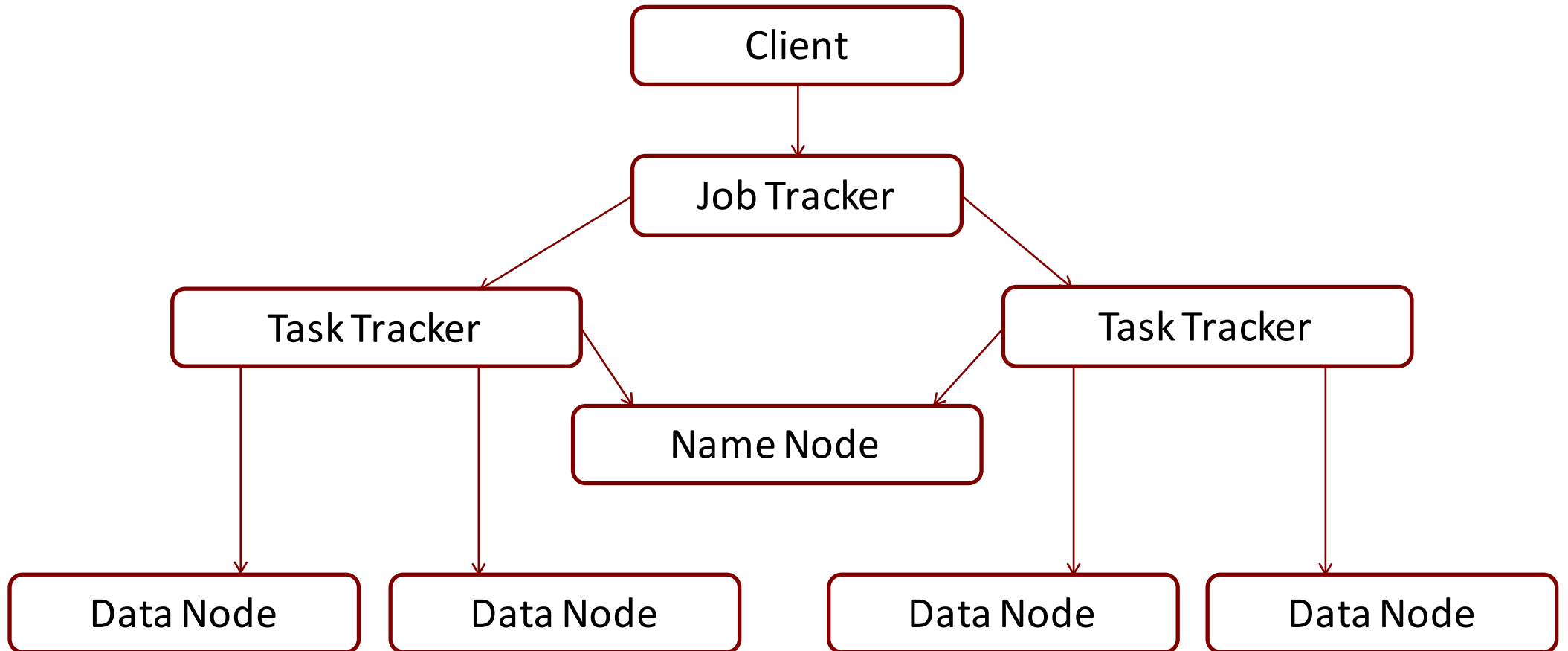
1) Extract data
2) Denormalise
3) Transform

MAP

4) Upload
5) Query

REDUCE

# Cloud computing

# Process pipeline

1) Translate to CSV (exports_makecsv)
2) Filter out supressed data (exports_process_supression)
3) Get ComCode data for that month (get_comcodes)
4) De-Normalise CSV with ComCodes and translate dates to timestamps (expand_csv)
5) Import into Big Query

# DEMO & EXERCISE

Data in Socrata: bit.ly/uk_trade_socrata

# Quesitons

Is the UKTrade data big data?

What are the biggest problems with the data?

How would you change your data to use cloud compute platforms?

# Where are we now?