# CBD2204: Week 1

Takis Zourntos

# overview

- data is created constantly, at an ever-increasing rate
- mobile phones, social media, medical imaging technologies, all of these create data, which must be stored and processed
- devices and sensors (e.g., IoT technologies) generate information that needs real-time (deadline-driven) processing
- two major challenges:
  - keeping up with the rate of data generation
  - analyzing the vast amount data, which may not be structured

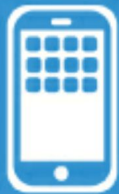leading the way in big data are industries such as:

- credit card companies (monitoring every purchase made by each customer)
- mobile phone companies (analyzing, for example, subscriber calling patterns)
- social media companies (where the data about users has inherent value)

But what makes something "Big Data"?

1. huge volume of data
2. complexity/diversity of data types and structures
3. speed of new data creation and growth ("high velocity data")

**Big Data is sometimes described as having the three Vs: volume, variety and velocity**
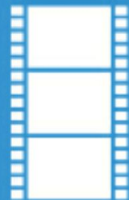
# What's Driving Data Deluge?

**Mobile Sensors**

**Social Media**

**Video Surveillance**

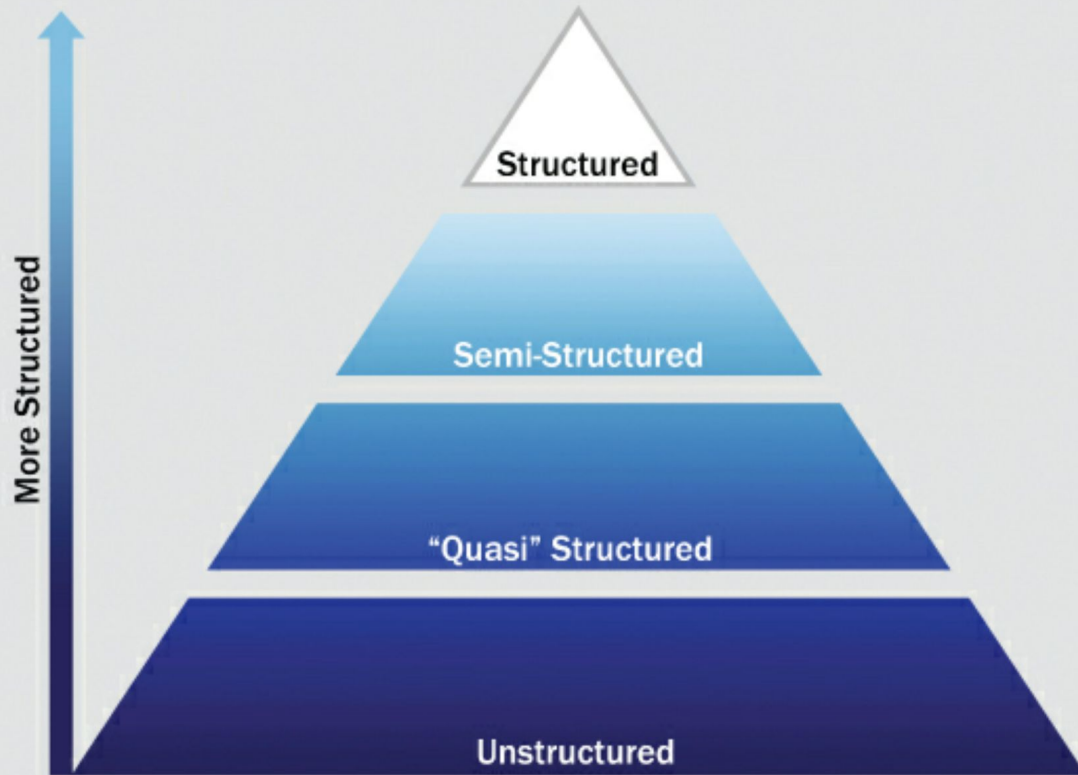**Video Rendering**

**Smart Grids**

**Geophysical Exploration**

**Medical Imaging**

**Gene Sequencing**

Big Data Characteristics: Data Structures
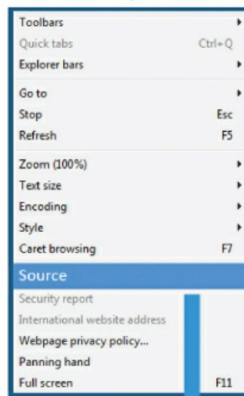Data Growth Is Increasingly Unstructured

# data structures

- **structured data**: data containing defined type and format (think: spreadsheet, CSV files, relational databases)
- **semi-structured data**: textual data files with a discernable pattern, enables parsing (e.g., XML files)
- **quasi-structured data**: textual data with erratic data formats that needs effort to format (e.g., clickstream data containing inconsistencies in data values and formats)
- **unstructured data**: data with no inherent structure such as text documents, images and video

| Fiscal Year | Number of Sites | Peak (July) Participation | Meals Served | Total Federal Expenditures [2] |
|---|---|---|---|---|
| | ------------Thousands------------ | | --Mil.-- | ---Million $--- |
| 1969 | 1.2 | 99 | 2.2 | 0.3 |
| 1970 | 1.9 | 227 | 8.2 | 1.8 |
| 1971 | 3.2 | 569 | 29.0 | 8.2 |
| 1972 | 6.5 | 1,080 | 73.5 | 21.9 |
| 1973 | 11.2 | 1,437 | 65.4 | 26.6 |
| 1974 | 10.6 | 1,403 | 63.6 | 33.6 |
| 1975 | 12.0 | 1,785 | 84.3 | 50.3 |
| 1976 | 16.0 | 2,453 | 104.8 | 73.4 |
| TQ 3] | 22.4 | 3,455 | 198.0 | 88.9 |
| 1977 | 23.7 | 2,791 | 170.4 | 114.4 |
| 1978 | 22.4 | 2,333 | 120.3 | 100.3 |
| 1979 | 23.0 | 2,126 | 121.8 | 108.6 |
| 1980 | 21.6 | 1,922 | 108.2 | 110.1 |
| 1981 | 20.6 | 1,726 | 90.3 | 105.9 |
| 1982 | 14.4 | 1,397 | 68.2 | 87.1 |
| 1983 | 14.9 | 1,401 | 71.3 | 93.4 |
| 1984 | 15.1 | 1,422 | 73.8 | 96.2 |
| 1985 | 16.0 | 1,462 | 77.2 | 111.5 |
| 1986 | 16.1 | 1,509 | 77.1 | 114.7 |
| 1987 | 16.9 | 1,560 | 79.9 | 129.3 |
| 1988 | 17.2 | 1,577 | 80.3 | 133.3 |
| 1989 | 18.5 | 1,652 | 86.0 | 143.8 |
| 1990 | 19.2 | 1,692 | 91.2 | 163.3 |

**SUMMER FOOD SERVICE PROGRAM 1]**
(Data as of August 01, 2011)

example of structured data

example of semi-structured data

example of quasi-structured data

example of unstructured data

types of data repositories

| Data Repository | Characteristics |
|---|---|
| Spreadsheets and data marts ("spreadmarts") | Spreadsheets and low-volume databases for recordkeeping<br>Analyst depends on data extracts. |
| Data Warehouses | Centralized data containers in a purpose-built space<br>Supports BI and reporting, but restricts robust analyses<br>Analyst dependent on IT and DBAs for data access and schema changes<br>Analysts must spend significant time to get aggregated and disaggregated data extracts from multiple sources. |
| Analytic Sandbox (workspaces) | Data assets gathered from multiple sources and technologies for analysis<br>Enables flexible, high-performance analysis in a nonproduction environment; can leverage in-database processing<br>Reduces costs and risks associated with data replication into "shadow" file systems<br>"Analyst owned" rather than "DBA owned" |

**Analytical Approach**

Exploratory

Explanatory

**Time**

Past

Future

Data Science

Business Intelligence

### Predictive Analytics and Data Mining (Data Science)

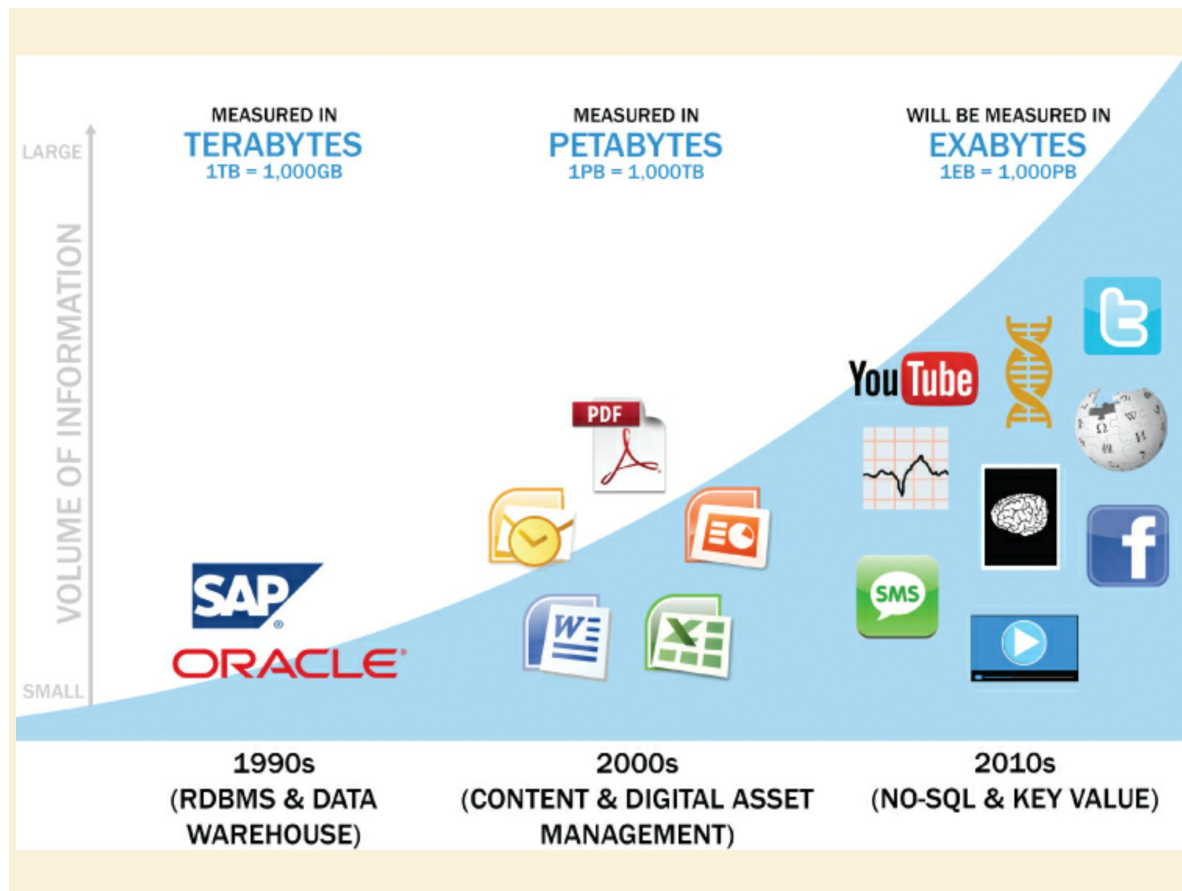| Typical Techniques and Data Types | • Optimization, predictive modeling, forecasting, statistical analysis<br>• Structured/unstructured data, many types of sources, very large datasets |
|---|---|
| Common Questions | • What if...?<br>• What's the optimal scenario for our business?<br>• What will happen next? What if these trends continue? Why is this happening? |

### Business Intelligence

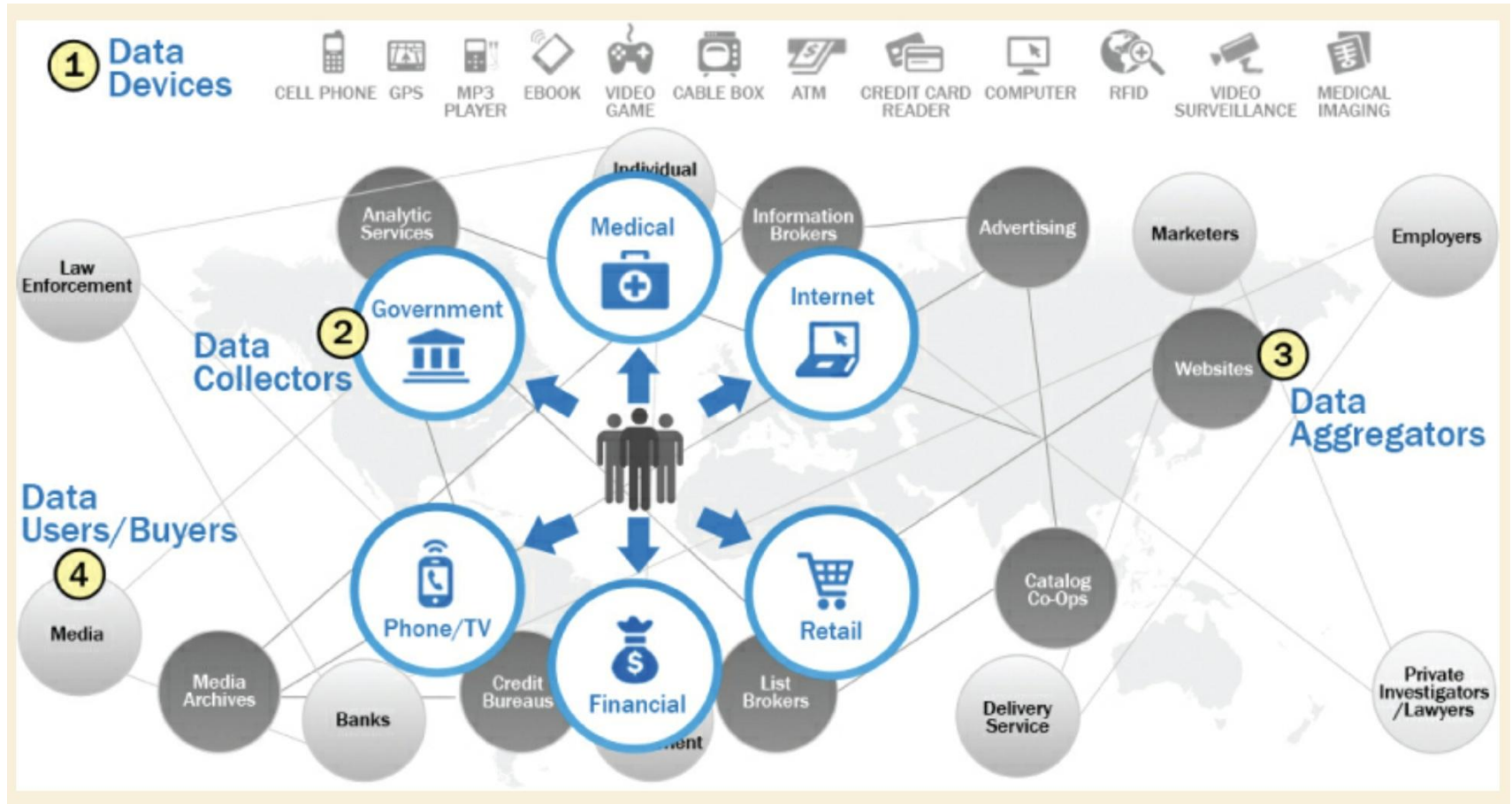| Typical Techniques and Data Types | • Standard and *ad hoc* reporting, dashboards, alerts, queries, details on demand<br>• Structured data, traditional sources, manageable datasets |
|---|---|
| Common Questions | • What happened last quarter?<br>• How many units sold?<br>• Where is the problem? In which situations? |

typical analytical architecture



EDW: Enterprise Data Warehouse

rise of big data sources

emerging Big Data ecosystems

Main Reference:

*Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data*, EMC Education Services

Amazon.ca link:
https://www.amazon.ca/Data-Science-Big-Analytics-Discovering-ebook/dp/B00RXHVQF6/ref=tmm_kin_swatch_0?_encoding=UTF8&qid=&sr=