**Ulster University**

# COM 745
## Databases within big data

**Dr. Joe Rafferty**

ulster.ac.uk

1

---

# An overview of databases
## Introduction

- A key pillar of data science is storage and manipulation of Big Data within databases
- Big data is defined by the three V's; a high Volume of Variable data stored at a high Velocity
- Big data storage relies on a variety of database classes that are optimised for specific roles and for specific data types
- These systems balance strengths and weaknesses which must be considered when producing systems that facilitate big data

**Ulster University**

2

# An overview of databases
## Introduction

- This module aims to introduce an overview of database systems
- A variety of database systems will be covered
  - The underlying concepts behind these databases will be summarised
  - The optimisations each of these systems will be covered
  - Use cases of these systems will be explored
  - Practical skills related to use of these systems will be explained and taught

Ulster
University

3

# An overview of databases
## Introduction

- This module will not explore implementation of these databases systems

**Axiom of reflexivity** [ edit ]

If $Y \subseteq X$ then $X \to Y$

**Axiom of augmentation** [ edit ]

If $X \to Y$, then $XZ \to YZ$ for any $Z$

**Axiom of transitivity** [ edit ]

If $X \to Y$ and $Y \to Z$, then $X \to Z$

Ulster
University

4

## An overview of databases
Introduction

- This module will explore practical skills related to use of these database systems

  SQL SELECT Syntax
  SELECT *column_name*, *column_name*
  FROM *table_name*;

Ulster
University

5

## An overview of databases
Introduction

- A database system is defined as:


**"a comprehensive collection of related data organized for convenient access, generally in a computer."**

Ulster
University

6

# An overview of databases
## Introduction

- A database system is an organised collection of data
- These collections typically facilitate:
  - Insertion modification and deletion of data
  - Retrieval of stored data
  - Administration of the database, such as providing security

- These collections are stored within abstractions and logical arrangements, such as tables

Ulster
University

7

# An overview of databases
## Introduction

- The collections typically employ some indexing strategy to efficiently retrieve the data within
- The collections typically employ some storage rules to satisfy their design goals and indexing strategy
- The collections may have some design aspects to cater for security, data integrity and concurrent access

Ulster
University

8

# An overview of databases
## Introduction

- Databases can reside in the physical world and within computing
- Physical world databases include:
  - An address book
  - A phone book
  - Index Cards
  - The Dewey Decimal System*



Ulster University

9

# An overview of databases
## Introduction

- Generally computing based databases can either be locally hosted or server based
  - Local databases include:
  - A spreadsheet*
  - A CSV file*
  - A locally hosted database system, such as SQLite, Libré Base or MS Access



Ulster University

10

# An overview of databases
## Introduction

- Server based databases are generally performance optimised, shared, resources
- These have the greatest variety of design goals
- Server based databases include:
  - MySQL – A relational DB
  - InfluxDB – A time series DB
  - MongoDB – A document oriented DB
  - Neo4J – a graph DB

Ulster
University

11

# An overview of databases
## Database Variety

- Databases vary to optimise for specific use cases, within a set restrictions or goals
- However, it is possible to use these databases in ways beyond their intended use
- A number of physical world databases will be presented in the upcoming slides. Following this their misuse will be explored.

Ulster
University

12

# An overview of databases
## Database Variety

# An overview of databases
## Database Variety

- These databases have optimisations

- Phonebook
  - Use: retrieving a phone number given an name and contracted address in a combined field.
  - Indexed alphabetically by surname
  - Read only
  - All fields non optional
  - Phone number is the unique key

# An overview of databases
## Database Variety



**Ulster University**

15

# An overview of databases
## Database Variety

- These databases have optimisations

- Address book
  - Use: retrieving and storing a in-depth contact information about a person
  - Indexed alphabetically by forename or surname, as dictated by end user
  - Supports insertion of data but not re-indexing/restructuring
  - All fields are optional, field size varies
  - There is no strictly enforced unique key

**Ulster University**

16

17

## An overview of databases
### Database Variety

- These databases have optimisations
- Index card
  - Use: retrieving and storing user defined information
  - Index is as dictated by end-user, may be weakly enforced by storage container structure
  - Supports insertion of data and re-indexing/restructuring
  - Undefined quantity of fields, variable proportioned size*
  - Fields may vary per entry
  - Supports image data
  - May support multiple indexing strategies
  - No strictly enforced unique key

18

# An overview of databases
## Database Variety

- It is possible to store the same types of information within each of these types of database
- Some types of information are more suited to specific types of database
- Selection of correct database for data type, retrieval strategy and purpose is essential

Ulster
University

19

# An overview of databases
## Real world example

- Storage and processing of sensor data
- Data stored:

  - Metadata about sensors – inserted once per device
    - SensorID
    - Name
    - Location
    - Type
    - Manufacturer

Ulster
University

20

# An overview of databases
## Real world example

- Storage and processing of sensor data
- Data stored:

    - Sensor Data
        - Sensor records – inserted at 6Hz, per device
        - Time
        - SensorID
        - State
        - JSON Data

Ulster
University

21

# An overview of databases
## Real world example

-   Initially stored all data in a "Big Data ready" relational database
    - Inserting, retrieving and modifying metadata individually took approximately 4 milliseconds
    - Inserting, retrieving and modifying metadata in bulk took approximately 30 milliseconds
    - Inserting, retrieving and modifying individual sensor records took approximately 4 milliseconds
    - Retrieving bulk sensor records for a 5 minute window from a single sensor took approximately **4 minutes**

Ulster
University

22

# An overview of databases
### Real world example

- Transitioned sensor records to a time series DB
  - Retrieving bulk sensor records for a 5 minute window from a single sensor took approximately 0.3 seconds

- Selection of the correct database systems, given data types and volume is essential

Ulster
University

23

# An overview of databases

- In the remainder of this module we will cover the following types of database systems
  - Relational databases
  - Document-oriented databases
  - Time series Databases
  - Graph databases
  - Semantic stores

Ulster
University

24

**An overview of databases**

- Next topic:

    Relational databases

Ulster
University

25

**Any Questions?**

Ulster
University

26