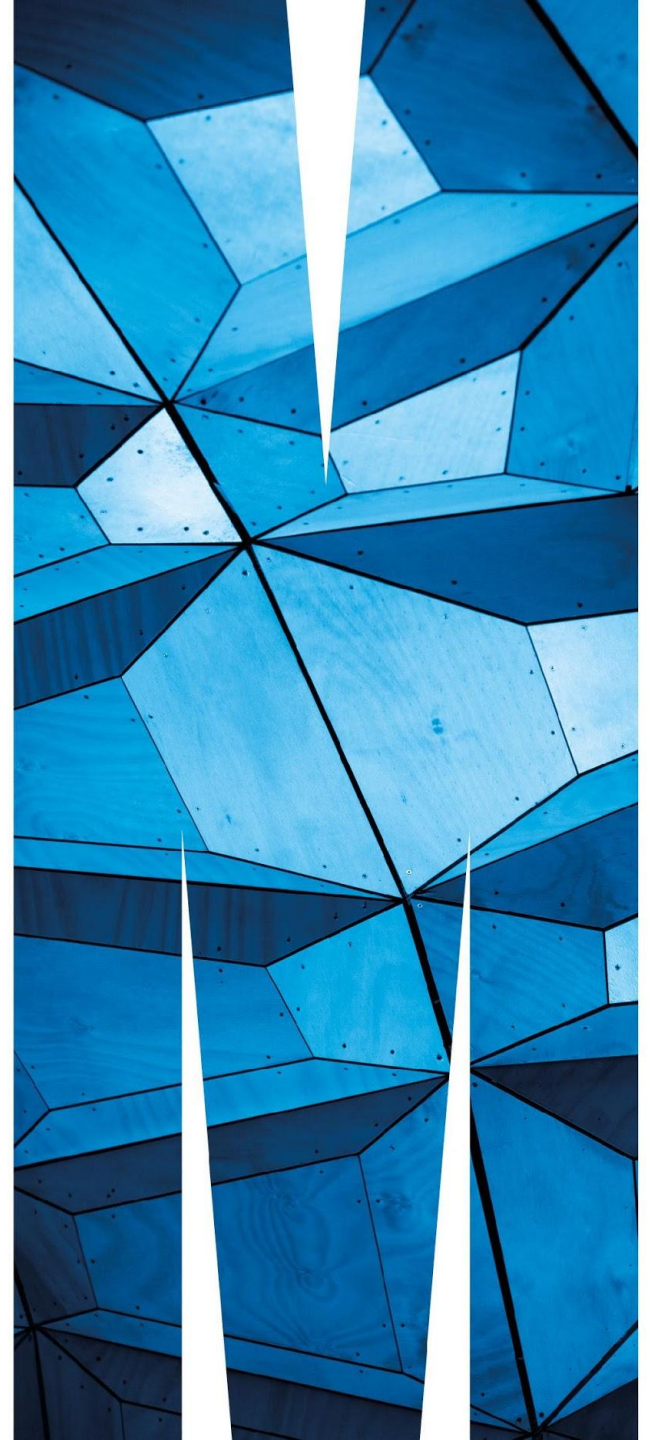




MONASH  
University

MONASH  
INFORMATION  
TECHNOLOGY

# Normalisation



# Data Normalisation

- Relations should be normalised in order to avoid anomalies which may occur when inserting, updating and deleting data.
- Normalisation is a **systematic series of steps** for progressively refining the data model.
- A formal approach to analysing relations based on their primary key (or candidate keys) and functional dependencies.
- Used:
  - as a design technique "bottom up design", and
  - as a way of validating structures produced via "top down design" (ER model converted to a logical model - see next week)

# Sample Data

PROJ_NUM	PROJ_NAME	EMP_NUM	EMP_NAME	JOB_CLASS	CHG_HOUR	HOURS
15	Evergreen	103	June E. Arbough	Elect. Engineer	84.50	23.80
15	Evergreen	101	John G. News	Database Designer	105.00	19.40
15	Evergreen	105	Alice K. Johnson *	Database Designer	105.00	35.70
15	Evergreen	106	William Smithfield	Programmer	35.75	12.60
15	Evergreen	102	David H. Senior	Systems Analyst	96.75	23.80
18	Amber Wave	114	Annelise Jones	Applications Designer	48.10	24.60
18	Amber Wave	118	James J. Frommer	General Support	18.36	45.30
18	Amber Wave	104	Anne K. Ramoras *	Systems Analyst	96.75	32.40
18	Amber Wave	112	Darlene M. Smithson	DSS Analyst	45.95	44.00
22	Rolling Tide	105	Alice K. Johnson	Database Designer	105.00	64.70
22	Rolling Tide	104	Anne K. Ramoras	Systems Analyst	96.75	48.40
22	Rolling Tide	113	Delbert K. Joenbrood *	Applications Designer	48.10	23.60
22	Rolling Tide	111	Geoff B. Wabash	Clerical Support	26.87	22.00
22	Rolling Tide	106	William Smithfield	Programmer	35.75	12.80
25	Starflight	107	Maria D. Alonzo	Programmer	35.75	24.60
25	Starflight	115	Travis B. Bawangi	Systems Analyst	96.75	45.80
25	Starflight	101	John G. News *	Database Designer	105.00	56.30
25	Starflight	114	Annelise Jones	Applications Designer	48.10	33.10
25	Starflight	108	Ralph B. Washington	Systems Analyst	96.75	23.60
25	Starflight	118	James J. Frommer	General Support	18.36	30.50
25	Starflight	112	Darlene M. Smithson	DSS Analyst	45.95	41.40

\* against EMP\_NAME indicates the project leader

# Problems with sample data

- JOB\_CLASS invites **entry errors** eg. Elec. Eng. vs Elect. Engineer vs E.E.
- Table has **redundant data**
  - Details of a charge per hour are repeated for every occurrence of job class
  - Every time an employee is assigned to a project emp name repeated
- Relations that contain redundant information may potentially suffer from several update anomalies
  - Types of update anomalies include:
    - **Insert Anomaly**
      - Insert a new employee only if they are assigned to a project
    - **Delete Anomaly**
      - Delete the only employee assigned to a project?
      - Delete the only employee of a particular job class?
    - **Modification (or update) Anomaly**
      - Update a job class hourly rate - need to update multiple rows

# The Normalisation Process Goals

- Creating valid relations, i.e. each relation meets the properties of the relational model. In particular:
  - Entity integrity
  - Referential integrity
  - No many-to-many relationship
  - Each cell contains a single value (is atomic).
- In practical terms when implemented in an RDBMS:
  - Each table represents a single subject
  - No data item will be unnecessarily stored in more than one table.
  - The relationship between tables can be established (pair of PK and FK is identified).
  - Each table is void of insert, update and delete anomalies.

# Representing a form as a relation

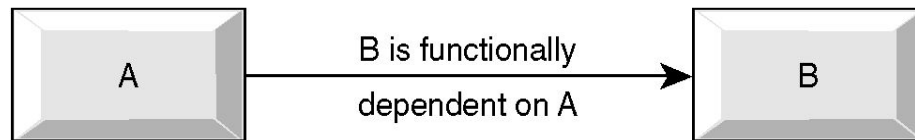
- This process follows a **standard** approach:
  - arrive at a name for the form which indicates what it represents (its subject)
  - determine if any attribute is multivalued (repeating) **for a given entity instance of the forms subject**
    - if an attribute (or set of attributes) appears multiple times then the group of related attributes need to be shown enclosed in brackets to indicate there are multiple sets of these values for each instance
- Looking at our SAMPLE DATA
  - Name: EMPLOYEE\_PROJECT\_ASSIGNMENT
    - simplify name to ASSIGNMENT for lecture
  - ASSIGNMENT (proj\_num, emp\_num, emp\_name, job\_class, chg\_hour, assign\_hours)
  - i.e. the form consists of repeating rows (instances) of assignment data



# Representing a form as a relation

CUSTOMER ORDER			
<b>Order Number:</b> 61384		<b>Order Date:</b> 12/3/2020	
<b>Customer Number:</b> 1273			
<b>Customer Name:</b> Computer Training Centre			
<b>Customer Address:</b> 123 Excellent St Monash, Vic, 3000			
PRODUCT NUMBER	DESCRIPTION	QTY ORDERED	LINE PRICE
M128	Bookcase	4	800
B381	TV Cabinet	2	600
R210	Round Table	3	1500

# Functional Dependency Revisited



- An attribute B is FUNCTIONALLY DEPENDENT on another attribute A, if a value of A determines a single value of B at any one time.
  - $A \rightarrow B$
  - $\text{PRODNO} \rightarrow \text{PRODESC}$
  - $\text{CUSTNUMB} \rightarrow \text{CUSTNAME}$
  - $\text{ORDERNO} \rightarrow \text{ORDERDATE}$ 
    - ORDERNO - independent variable, also known as the DETERMINANT
    - ORDERDATE - dependent variable
- TOTAL DEPENDENCY
  - attribute A determines B AND attribute B determines A
    - $\text{EMPLOYEE-NUMBER} \rightarrow \text{TAX-FILE-NUMBER}$
    - $\text{TAX-FILE-NUMBER} \rightarrow \text{EMPLOYEE-NUMBER}$



# Functional Dependency

- For a **composite** PRIMARY KEY, it is possible to have FULL or PARTIAL dependency.
- FULL DEPENDENCY
  - occurs when an attribute is always dependent on all attributes in the composite PK
  - ORDERNO, PRODNO → QTYORDERED
- Lack of full dependency for multiple attribute key = **PARTIAL DEPENDENCY**
  - ORDERNO, PRODNO  
→ PRODDDESC, QTYORDERED
  - here although qtyordered is **fully dependent** on orderno and prodno, *only* prodno is required to determine proddesc
  - proddesc is said to be **partially dependent** on orderno and prodno

# Functional Dependency

- **TRANSITIVE DEPENDENCY**

- occurs when Y depends on X, and Z depends on Y - thus Z also depends on X ie.  $X \rightarrow Y \rightarrow Z$
- **and** Y is not a candidate key (or part of a candidate key)
- $ORDERNO \rightarrow CUSTNUMB \rightarrow CUSTNAME$

- Dependencies are depicted with the help of a **Dependency Diagram**.
- Normalisation converts a relation into relations of progressively smaller number of attributes and tuples until an optimum level of decomposition is reached - little or no data redundancy exists.
- The output from normalisation is a set of relations that meet all conditions set in the relational model principles.

# Unnormalised Form (UNF)

- The UNF representation of a relation is the representation which you have mapped from your inspection of the form
  - it is a **single** named representation (name is not pluralised)
  - no PK etc have as yet been identified
- **ASSIGNMENT** (proj\_num, emp\_num, emp\_name, job\_class, chg\_hour, assign\_hours)
- **ORDER** (orderno, orderdate, custnumb, custname, custaddress (prodno, proddesc, qtyordered, lineprice))

Can ASSIGNMENT and/or ORDER be called a relation?  
If not, why not?

# First Normal Form

- FIRST NORMAL FORM (part of formal definition of a relation)
  - A RELATION IS IN FIRST NORMAL FORM (1NF) IF:
    - a unique primary key has been identified for each tuple/row.
    - it is a valid relation
      - Entity integrity (no part of PK is null)
      - Single value for each cell ie. no repeating group (multivalued attribute).
    - all attributes are functionally dependent on all or part of the primary key

## UNF to 1NF

- Move from UNF to 1NF by:
  1. identify a unique identifier for the repeating group.
  2. *remove any repeating group* along with the PK of the main relation.
  3. The PK of the new relation resulting from the removal of repeating group will *normally* have a composite PK made up of the PK of the main relation and the unique identifier chosen in 1. above, but this ***must be checked***.

# 1NF to 2NF

- A RELATION IS IN 2NF IF -
  - all non key attributes are functionally dependent on the primary key (simple definition)
    - used by the textbook in examples
  - all non key attributes are functionally dependent on **any candidate key** (general definition)
    - see textbook section 6-3, same as *simple* if only one candidate key
    - **Requirement for our unit**

## 2NF to 3NF

- A RELATION IS IN 3NF IF -
  - all transitive dependencies have been removed
    - check for ***non key attribute dependent on another non key attribute***
- Move from 2NF to 3NF by removing transitive dependencies



# Customer Order Form

CUSTOMER ORDER			
<b>Order Number:</b> 61384		<b>Order Date:</b> 12/3/2020	
<b>Customer Number:</b> 1273			
<b>Customer Name:</b> Computer Training Centre			
<b>Customer Address:</b> 123 Excellent St Monash, Vic, 3000			
PRODUCT NUMBER	DESCRIPTION	QTY ORDERED	LINE PRICE
M128	Bookcase	4	800
B381	TV Cabinet	2	600
R210	Round Table	3	1500

**ORDER** ( orderno, orderdate, custnumb, custname,  
custaddress (prodno, proddesc, qtyordered,  
lineprice))

## EMPLOYEE ON-BOARDING FORM

<b>Employee Number</b>	1123 (office use only)			
<b>First Name</b>	Ada	<b>Last Name</b>	Lovelace	
<b>DOB</b>	1-Jan-1990			
<b>Address</b>	<i>Street No</i>	<i>Street</i>	<i>Suburb</i>	<i>Postcode</i>
	900	Dandenong Rd	Caulfield East	3145
<b>Phone</b>	04113344556 (M), 99031000 (OFFICE)			
<b>Qualifications</b>				
	<b>Degree Name</b>	<b>Institution</b>	<b>Year</b>	
	Bachelor of Computer Science	MIT	2011	
	Master of Information Technology	Monash	2013	
<b>Family Members</b>				
	<b>No</b>	<b>Name</b>	<b>DOB</b>	
	1	Albert Einstein	02-Jan-1992	
	2	Grace Hopper	12-May-1994	
<b>SKILL (tick selected)</b>				
	<b>Skill name</b>			
	Java			
	SQL			
	SPARK			
	Python			

*Assume a phone number may be shared between employees*

# Summary

- Things to remember
  - Represent form as presented, no interpretation, to yield starting point (UNF)
  - Functional dependency
  - Process of removing attributes in relations based on the concept of 1NF, 2NF and 3NF.
    - UNF to 1NF define PK & remove repeating group.
    - 1NF to 2NF remove partial dependency.
    - 2NF to 3NF remove transitive dependency.