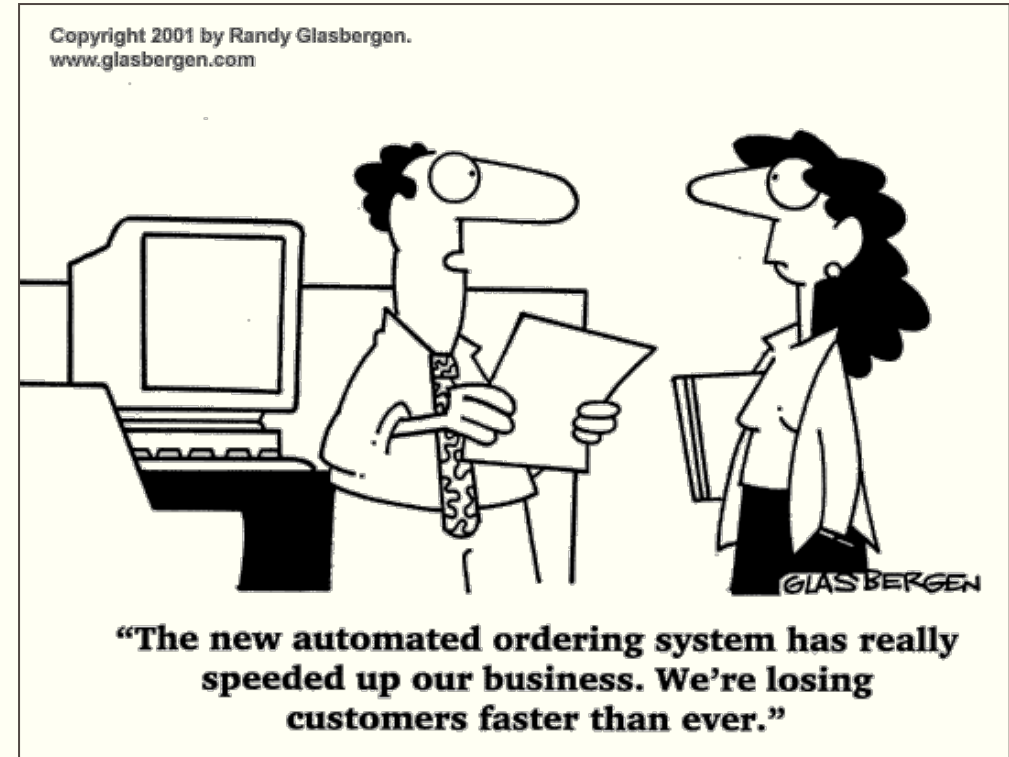


Srini Badri

# Design of relational data models

---

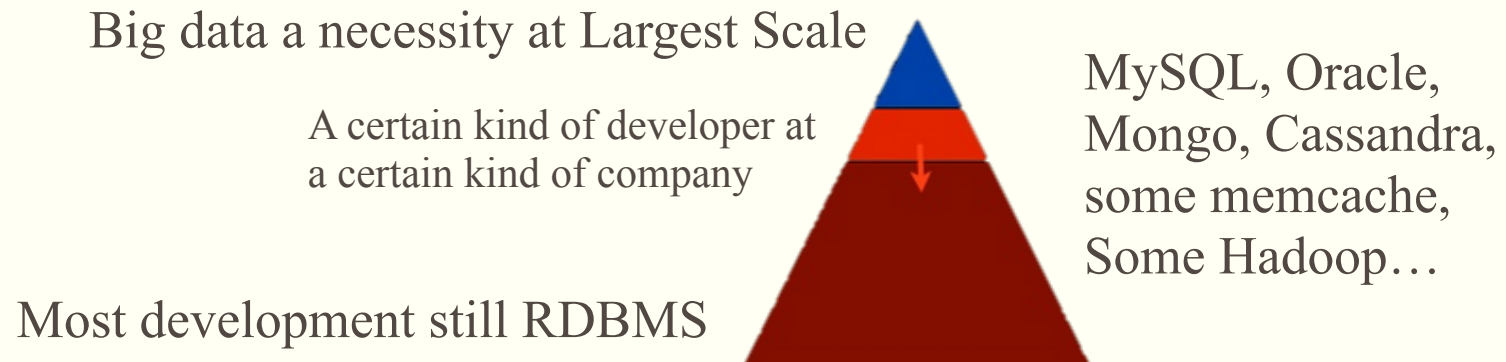
- Functional dependencies
- The normal forms 1NF, 2NF, 3NF, BCNF



# A Big Data Fallacy

---

- Data model design in the era of Big Data is less important?
  - New high-volume data streams
  - Specialized hardware/software
  - Storage issues coped by hardware appliance
- Fact
  - Most data is physically located in DBMS and new special-purpose appliance
  - Data loads, extract, transform preprocessing operations continue as is
  - Database design for quality assurance



# Relational Databases Design

---

- Relational Database design:
  - The grouping of attributes to form “good” relation schemas
- We have assumed schema  $R$  is given
  - Universal Relation:  $R$  could have been a single relation containing all attributes that are of interest
  - Normalization breaks  $R$  into smaller relations
  - $R$  could have been the results of some ad-hoc design of relations, which we then test/convert to normal form

---

- Charges its clients by billing hours spent on each contract
- Hourly billing rate is dependent on employee's position

TABLE 5.1 A SAMPLE REPORT LAYOUT

\* Indicates project leader

# A Table in the Report Format

Table name: RPT_FORMAT				Database name: Ch05_ConstructCo			
	PROJ_NUM	PROJ_NAME	EMP_NUM	EMP_NAME	JOB_CLASS	CHG_HOUR	HOURS
▶	15	Evergreen	103	June E. Arbough	Elect. Engineer	\$84.50	23.8
			101	John G. News	Database Designer	\$105.00	19.4
			105	Alice K. Johnson *	Database Designer	\$105.00	35.7
			106	William Smithfield	Programmer	\$35.75	12.6
			102	David H. Senior	Systems Analyst	\$96.75	23.8
	18	Amber Wave	114	Annelise Jones	Applications Designer	\$48.10	24.6
			118	James J. Frommer	General Support	\$18.36	45.3
			104	Anne K. Ramoras *	Systems Analyst	\$96.75	32.4
			112	Darlene M. Smithson	DSS Analyst	\$45.95	44.0
	22	Rolling Tide	105	Alice K. Johnson	Database Designer	\$105.00	64.7
			104	Anne K. Ramoras	Systems Analyst	\$96.75	48.4
			113	Delbert K. Joenbrood *	Applications Designer	\$48.10	23.6
			111	Geoff B. Wabash	Clerical Support	\$26.87	22.0
			106	William Smithfield	Programmer	\$35.75	12.8
	25	Starflight	107	Maria D. A. P. 30	Programmer	\$35.75	24.6
			115	Travis B. Baywangi	Systems Analyst	\$96.75	45.8
			101	John G. News *	Database Designer	\$105.00	56.3
			114	Annelise Jones	Applications Designer	\$48.10	33.1
			108	Ralph B. Washington	Systems Analyst	\$96.75	23.6
			118	James J. Frommer	General Support	\$18.36	30.5
			112	Darlene M. Smithson	DSS Analyst	\$45.95	41.4

# The Need for Normalization

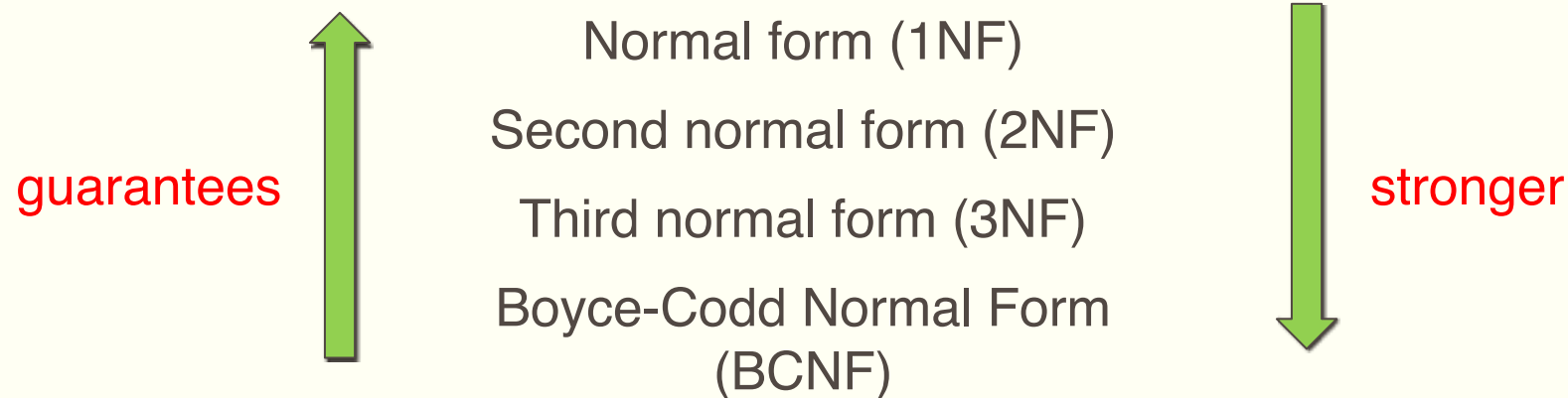
---

- Mixing attributes of multiple entities may cause problems
  - Information is stored redundantly wasting storage
  - Problems with update anomalies
    - Insertion anomalies
    - Deletion anomalies
    - Modification anomalies
- The report may yield different results depending on data anomaly
  - Primary keys
  - Redundancy
  - Possible data inconsistencies
    - E.g. JOB\_CLASS: [Elect Engineer](#), [Elect. Eng.](#), [E. E.](#), [El. Eng.](#)

# Normalization

---

- Process of decomposing unsatisfactory “bad” relations by breaking up their attributes into smaller relations
- Normal Form:
  - Condition using keys and functional dependencies (FDs) of a relation to certify whether a relation schema is in a particular normal form
  - 2NF, 3NF, BCNF based on keys and FDs of a relation schema
  - 4NF based on keys, multi-valued dependencies





# Functional Dependencies (FDs)

---

- FDs are used to specify formal measures of the “goodness” of relational designs
- A set of attributes  $X$  functionally determines a set of attributes  $Y$  if the value of  $X$  determines a unique value  $Y$
- $X \rightarrow Y$  holds if whenever two tuples have the same value for  $X$ , they must have the same value for  $Y$  on all relation instances

$$\forall t_1, t_2 \in r(R), \quad t_1[X] = t_2[X] \implies t_1[Y] = t_2[Y]$$

- FDs are derived from the real-world constraints on the attributes
  - $SSN \rightarrow ENAME$
  - $PNUMBER \rightarrow \{PNAME, PLOCATION\}$
  - $\{SSN, PNUMBER\} \rightarrow HOURS$
  - What if LHS is a key?

## Function Key (Cont'd)

---

- $A \rightarrow C$
- $\langle A, B \rangle \rightarrow C$
- $\langle A, C \rangle \rightarrow D$
- $A \rightarrow D$
- $C \rightarrow D$

A	B	C	D
1	1	3	2
1	2	2	3
1	3	3	2

- Full/Partial FD
  - If removal of any attribute from X means the FD does not hold any more, it is a full dependency; otherwise it's a Partial dependency
- Transitive FD:  $X \rightarrow Y$ 
  - If there is a set of attributes  $Z$  that are neither a primary or candidate key and both  $X \rightarrow Z$  and  $Z \rightarrow Y$  holds.

# Inference Rules for FDs

---

- Given a set of FDs  $F$ , we can infer additional FDs that hold whenever the FDs in  $F$  holds.
- Armstrong's inference rules
  - Reflexive  
if  $Y \subseteq X$ , then  $X \rightarrow Y$
  - Argumentation  
if  $X \rightarrow Y$ , then  $X \cup Z \rightarrow Y \cup Z$  (or sometime written as  $XZ \rightarrow YZ$ )
  - Transitive  
if  $X \rightarrow Y$  and  $Y \rightarrow Z$ , then  $X \rightarrow Z$

# First Normal Form

---

- Tabular format in which
  - All key attributes are defined
  - There are no repeating groups in the table
  - All attributes are dependent on primary key
- Disallow composite attributes, multi-valued attributes and nested relations
- 1NF deals with the “shape” of the tables

# 1NF Normalization

Table name: DATA\_ORG\_1NF

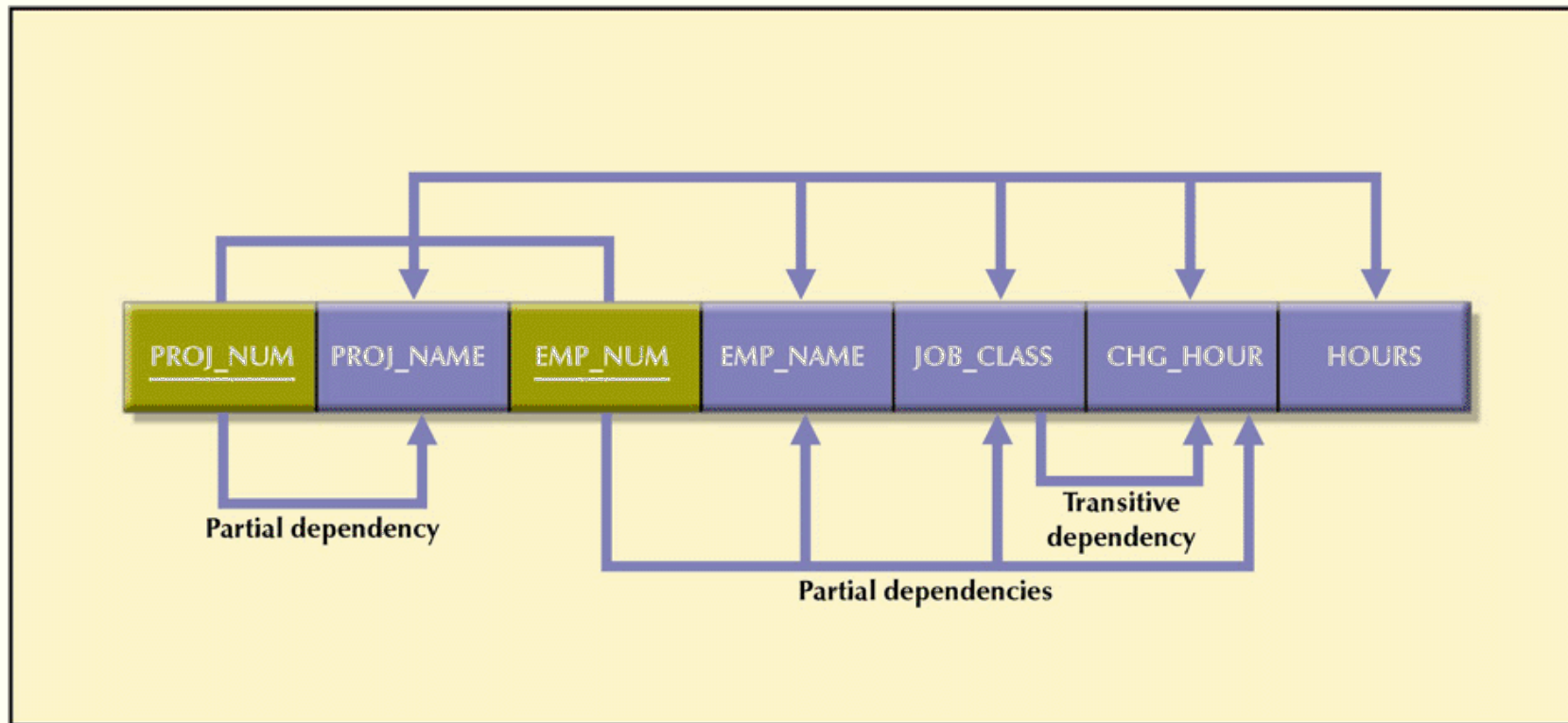
Database name: Ch05\_ConstructCo

	PROJ_NUM	PROJ_NAME	EMP_NUM	EMP_NAME	JOB_CLASS	CHG_HOUR	HOURS
►	15	Evergreen	103	June E. Arbough	Elect. Engineer	\$54.50	23.8
	15	Evergreen	101	John G. News	Database Designer	\$105.00	19.4
	15	Evergreen	105	Alice K. Johnson *	Database Designer	\$105.00	35.7
	15	Evergreen	106	William Smithfield	Programmer	\$35.75	12.6
	15	Evergreen	102	David H. Senior	Systems Analyst	\$96.75	23.8
	18	Amber Wave	114	Annelise Jones	Applications Designer	\$48.10	24.6
	18	Amber Wave	113	James J. Frommer	General Support	\$18.36	45.3
	18	Amber Wave	104	Anne K. Ramoras *	Systems Analyst	\$96.75	32.4
	18	Amber Wave	112	Darlene M. Smithson	DSS Analyst	\$45.95	44.0
	22	Rolling Tide	105	Alice K. Johnson	Database Designer	\$105.00	54.7
	22	Rolling Tide	104	Anne K. Ramoras	Systems Analyst	\$96.75	48.4
	22	Rolling Tide	113	Delbert K. Joenbrood *	Applications Designer	\$48.10	23.6
	22	Rolling Tide	111	Geoff B. Wabash	Clerical Support	\$26.87	22.0
	22	Rolling Tide	106	William Smithfield	Programmer	\$35.75	12.8
	25	Starflight	107	Maria D. Alonzo	Programmer	\$35.75	24.6
	25	Starflight	115	Travis B. Bawangi	Systems Analyst	\$96.75	45.8
	25	Starflight	101	John G. News *	Database Designer	\$105.00	55.3
	25	Starflight	114	Annelise Jones	Applications Designer	\$48.10	33.1
	25	Starflight	106	Ralph B. Washington	Systems Analyst	\$96.75	23.6
	25	Starflight	113	James J. Frommer	General Support	\$18.36	30.5
	25	Starflight	112	Darlene M. Smithson	DSS Analyst	\$45.95	41.4

# FD Diagram

---

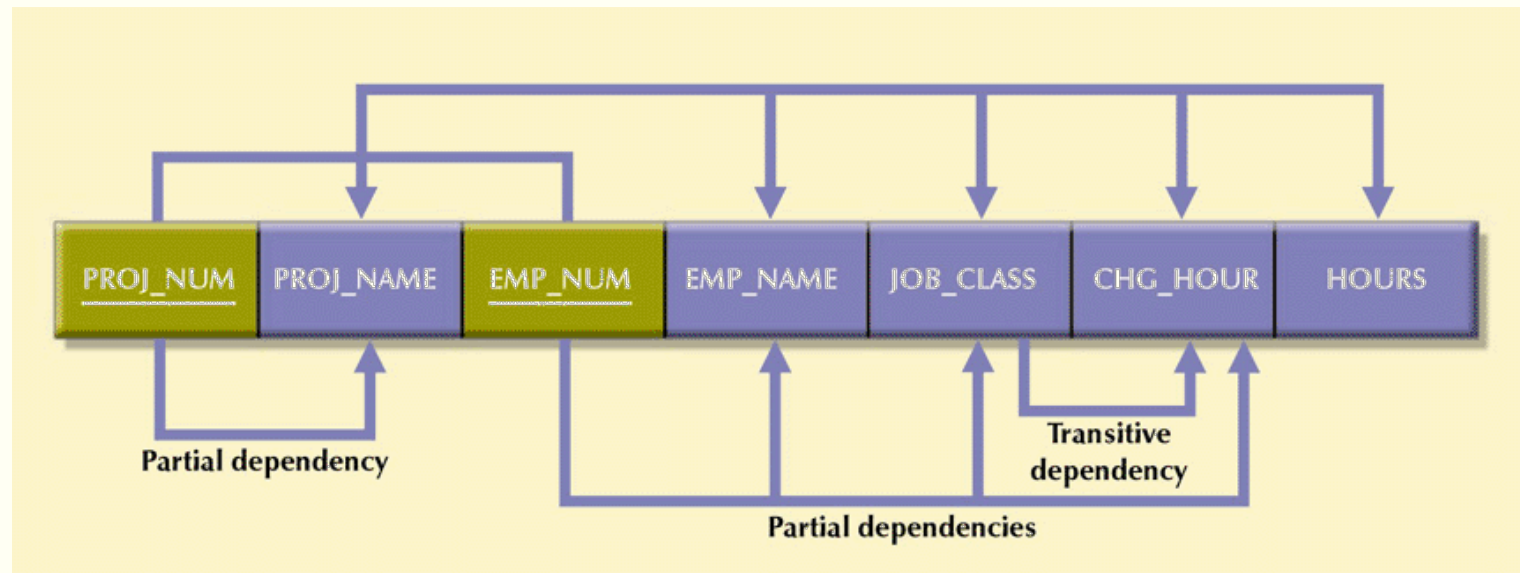
**FIGURE 5.3** A DEPENDENCY DIAGRAM: FIRST NORMAL FORM (1NF)



# Second Norm Form (2NF)

---

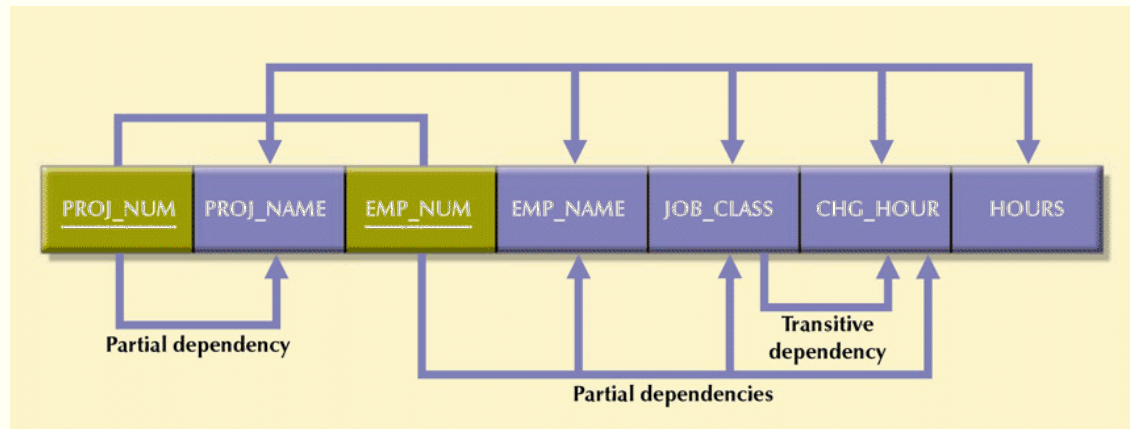
- Table is in 2NF if
  - It is in 1NF
  - It includes no partial dependencies
    - No attribute is dependent on only a portion of the primary key
    - Every attribute A not in PK is fully functionally dependent on PK
- 2NF deals with the relationship between non-key and key attributes



# Conversion 1NF to 2NF

---

- Step 1: Write each key attribute on separate line and then write the original (composite) key on the last line; Each component will become the key in a new table.
- Step 2: Determine which attributes are dependent on which other attributes (remove anomalies)

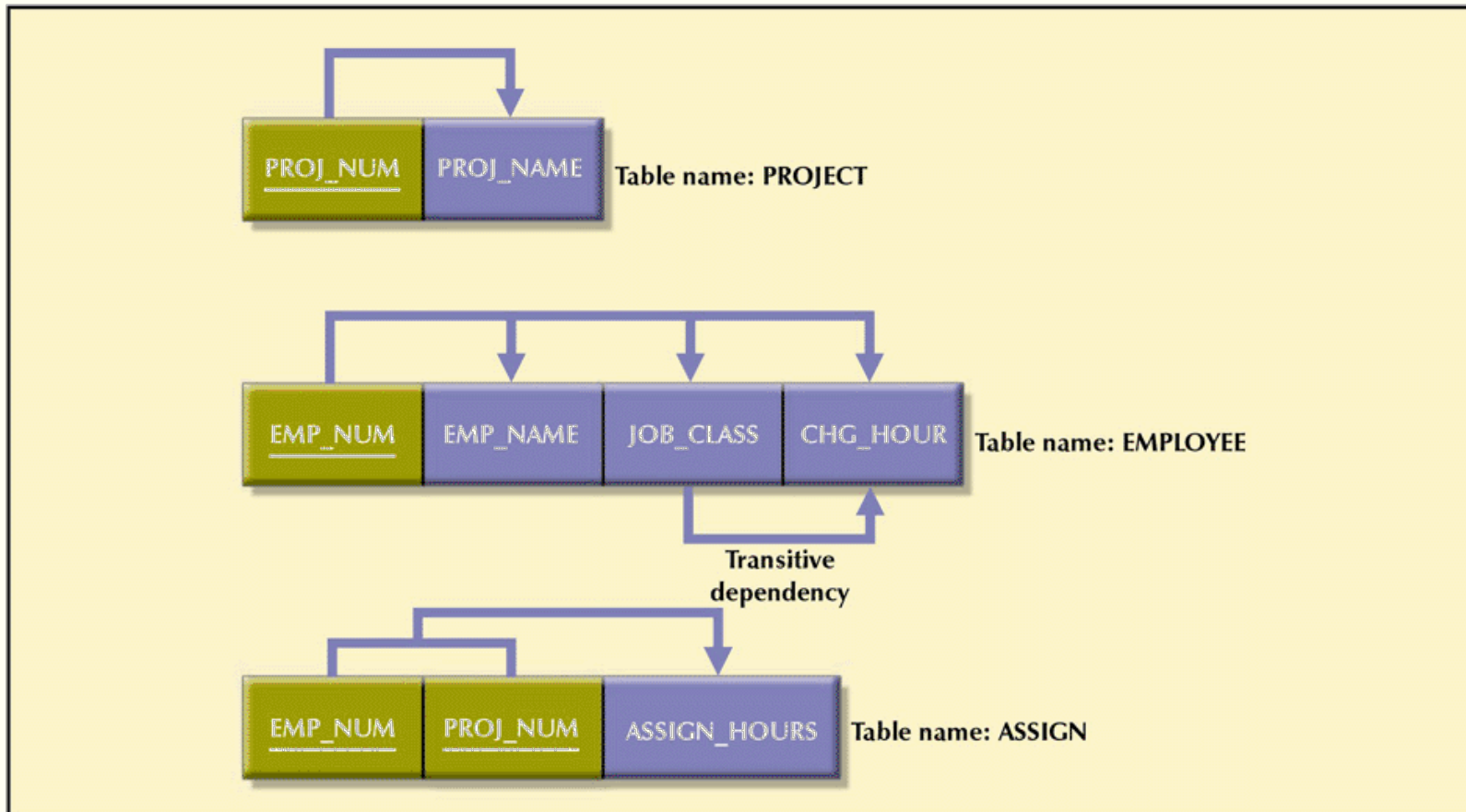




## 2NF Results

---

**FIGURE 5.4** SECOND NORMAL FORM (2NF) CONVERSION RESULTS



# Third Normal Form (3NF)

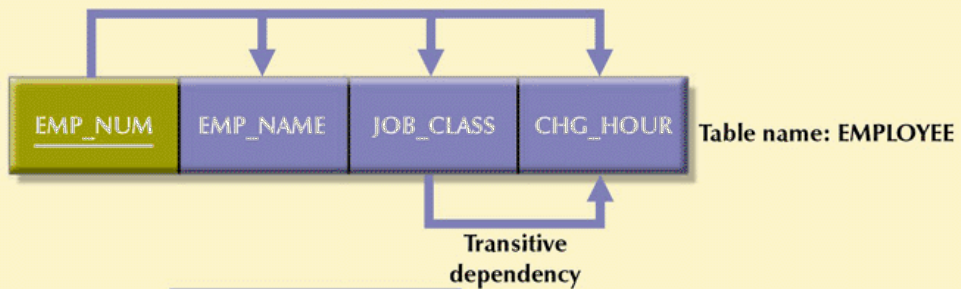
---

- A table is in third normal form (3NF) if
  - It is in 2NF
  - It contains no transitive dependencies
- 3NF removes transitive dependencies

# Conversion to 3NF

---

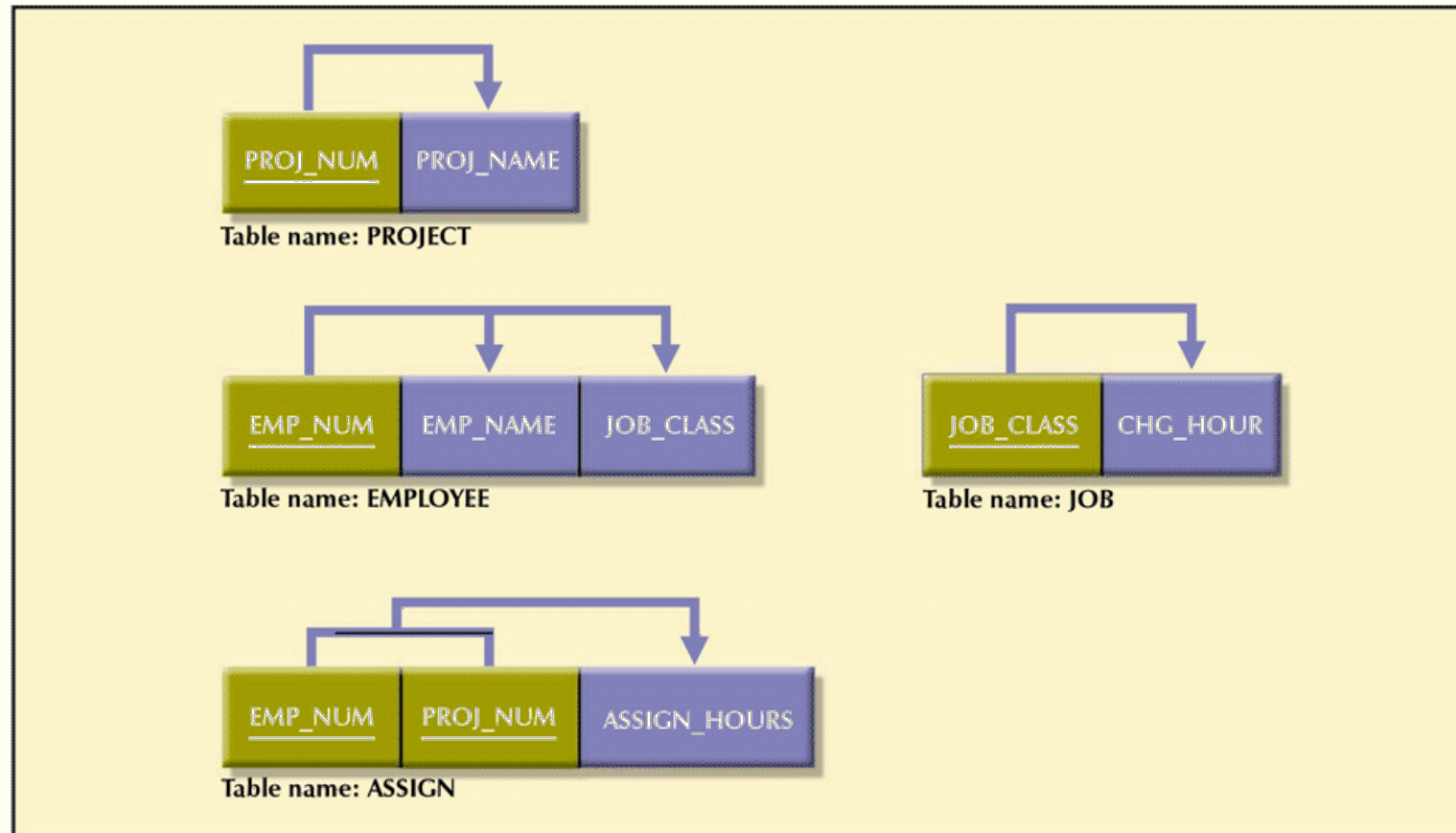
- Step 1: Find new Fact
  - For every transitive dependency  $X \rightarrow Y$ , write fact  $Z$  as a PK for a new table where  $X \rightarrow Z$  and  $Z \rightarrow Y$
- Step 2: Identify the dependent attributes
  - Identify the attributes dependent on each  $Z$  identified in previous step and find the dependency
  - Name the table to reflect its contents and function
- Step 3: Remove  $X \rightarrow Y$  from original table



# 3NF Results

---

**FIGURE 5.5 THIRD NORMAL FORM (3NF) CONVERSION RESULTS**

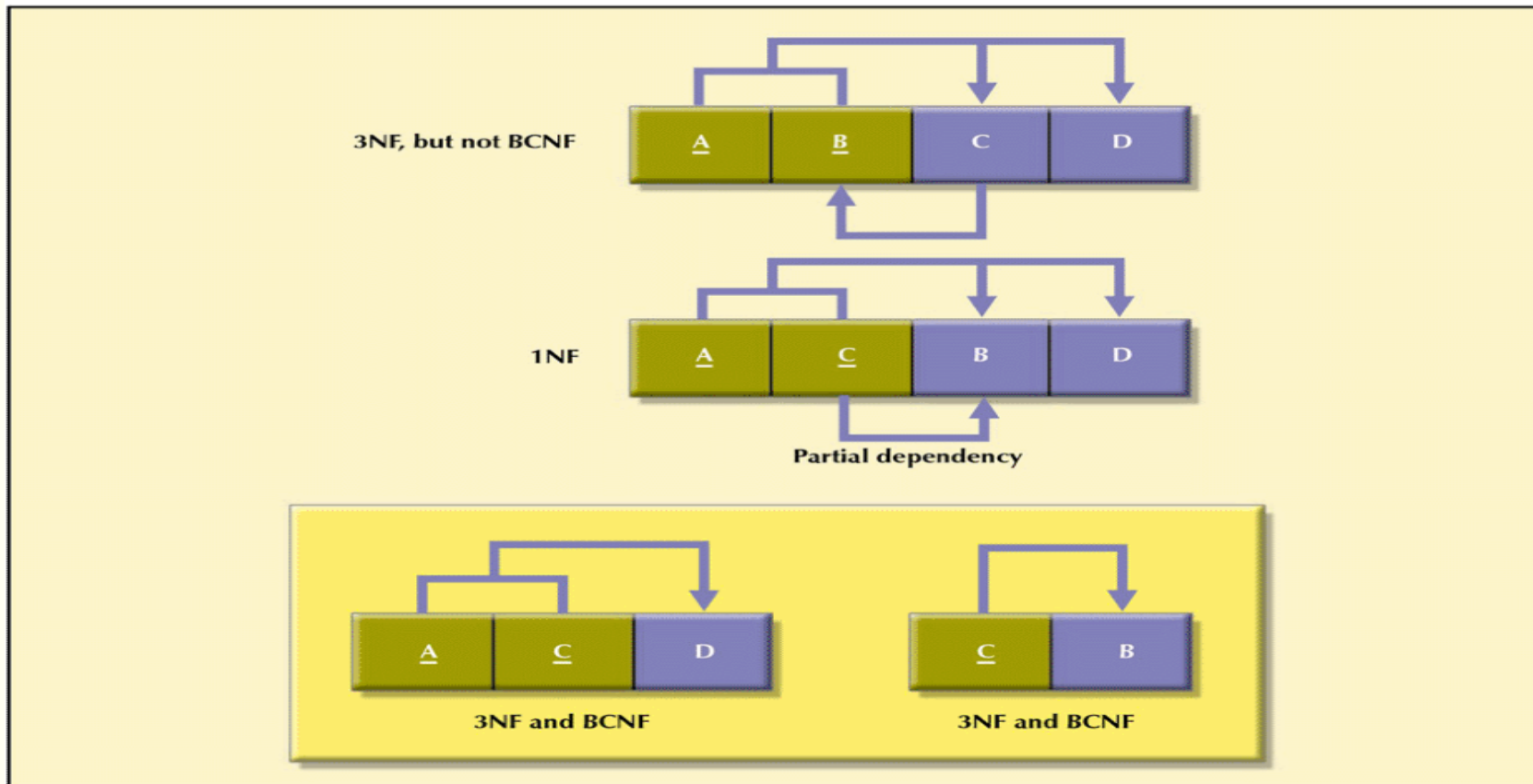


# Boyce-Codd Normal Form (BCNF)

---

- A relation schema  $R$  is in BCNF, a.k.a. 3.5NF, if whenever an FD  $X \rightarrow A$  holds in  $R$ , then  $X$  is a super-key of  $R$

FIGURE 5.8 DECOMPOSITION TO BCNF



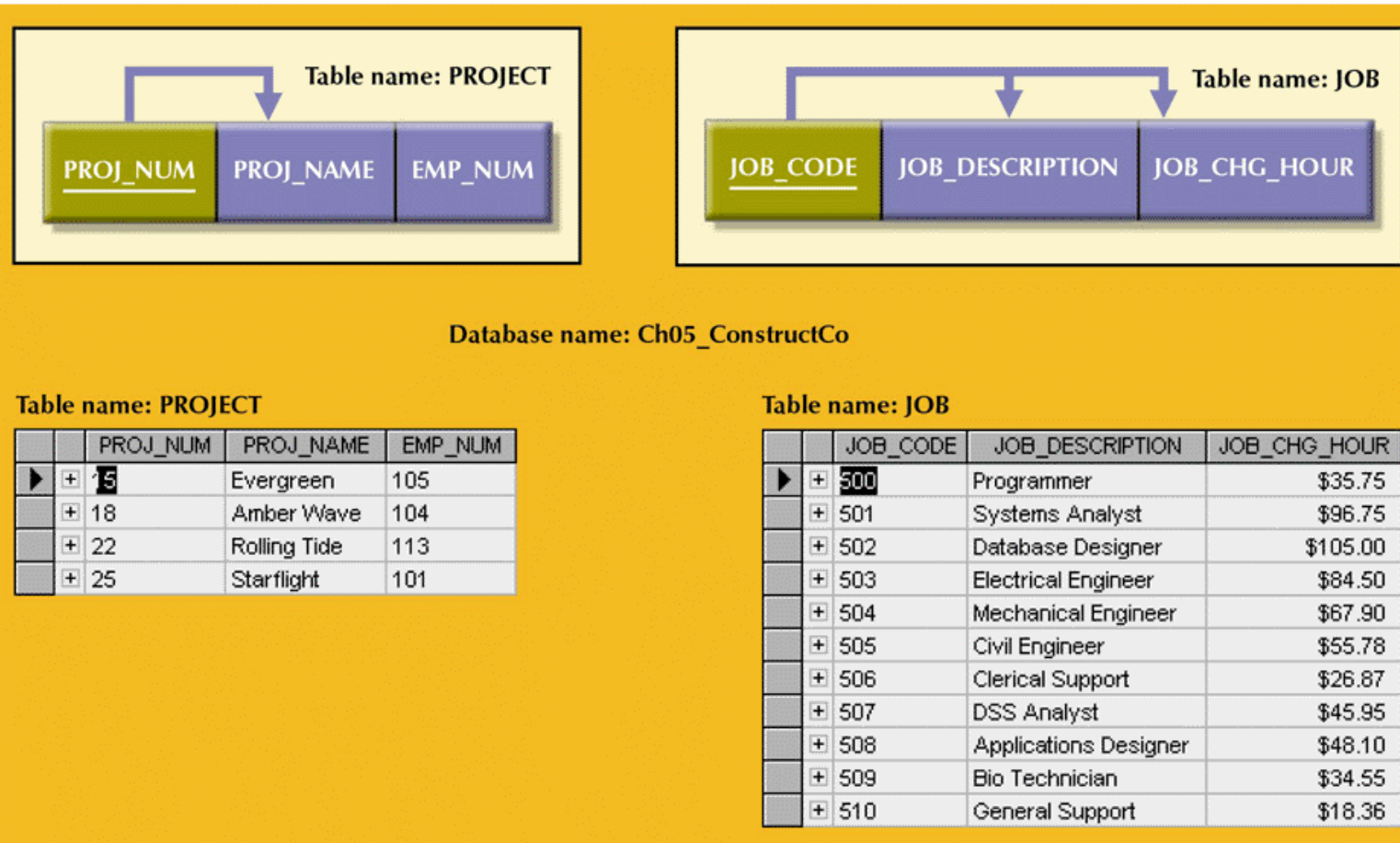
# Normalization

---

- Each normal form is strictly stronger than the previous one
  - Every 2NF relation is in 1NF
  - Every 3NF relation is in 2NF
  - Every BCNF relation is in 3NF
- There exist relations that in 3NF but not in BCNF
- The goal is to have each relation in BCNF (or 3NF)

# The Completed Database

FIGURE 5.6 THE COMPLETED DATABASE





# The Completed Database

FIGURE 5.6 THE COMPLETED DATABASE (CONTINUED)

