

# 5b. Normalization

CSCI 2541 Database Systems & Team Projects

Wood & Chaufournier

# Announcements?

Exam! -1.5 -2.5 weeks away

Relational Algebra HW ✓

SQL and Shopping Cart HW ]  
- ER diagrams  
- SQL

# Last time...

SQL DDL &  
DML

ER  
Entity  
Relationship  
Model

## Normalization

- Bad Schemas -
- Normal Forms -
- Functional Dependencies -

# this time...

# Good Schemas

The ER model can help us design a logical DB structure that matches our business goals

The conceptual schema must be translated into a logical (SQL) schema

How do we judge if a SQL schema is well designed?

# Bad Schemas

Let's track professors and their department

- We will put all the info together in one table so we don't have to worry about joining stuff!

<i>ID</i>	<i>name</i>	<i>salary</i>	<i>dept_name</i>	<i>building</i>	<i>Dept</i> <i>budget</i>
22222	Einstein	95000	<u>Physics</u>	<u>Watson</u>	<u>70000</u>
12121	Wu	90000	Finance	Painter	120000
32343	El Said	60000	History	Painter	50000
45565	Katz	75000	Comp. Sci.	Taylor	100000
98345	Kim	80000	Elec. Eng.	Taylor	85000
76766	Crick	72000	Biology	Watson	90000
10101	Srinivasan	65000	Comp. Sci.	Taylor	100000
58583	Califieri	62000	History	Painter	50000
83821	Brandt	92000	Comp. Sci.	Taylor	100000
15151	Mozart	40000	Music	Packard	80000
33456	Gold	87000	<u>Physics</u>	<u>Watson</u>	<u>70000</u>
76543	Singh	80000	Finance	Painter	120000

<i>Dept</i>	<i>Name</i>	<i>Bl</i>	<i>\$</i>
Phy			
CS			

Why is this a bad idea?

# Bad Schemas

Let's track professors and their department

- We will put all the info together in one table so we don't have to worry about joining stuff!

<i>ID</i>	<i>name</i>	<i>salary</i>	<i>dept_name</i>	<i>building</i>	<i>budget</i>
22222	Einstein	95000	Physics	Watson	70000
12121	Wu	90000	Finance	Painter	120000
32343	El Said	60000	History	Painter	50000
45565	Katz	75000	Comp. Sci.	Taylor	100000
98345	Kim	80000	Elec. Eng.	Taylor	85000
76766	Crick	72000	Biology	Watson	90000
10101	Srinivasan	65000	Comp. Sci.	Taylor	100000
58583	Califieri	62000	History	Painter	50000

**Update Anomalies:** need to modify all repetitive rows

**Insertion Anomalies:** Need to use NULL if we add a department with no instructors

**Deletion Anomalies:** Removing all instructors loses information about the department

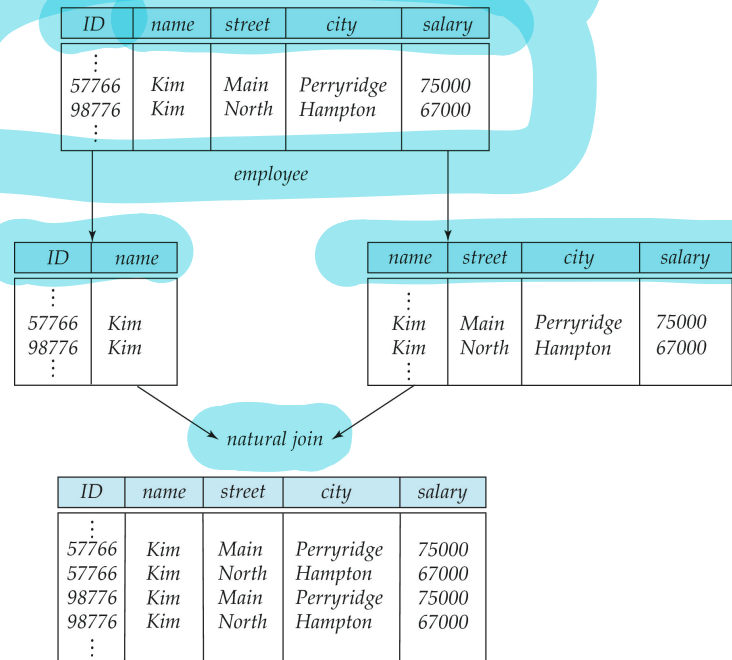
# Splitting Tables

**Decomposing** into separate tables helps resolve this... but there are multiple ways to split tables

- Not all decompositions are good!

## A Lossy Decomposition

results in us losing data if we try to merge back using a join



# What is Normalization?

1. Tests to see how “good” a schema is
2. Normalization algorithms to decompose relations into smaller relations that contain less redundancy
  - This decomposition requires that **no information is lost** and **reconstruction** of the original relations from the smaller relations must be possible.

Normalization should be done when you design your schema and anytime you update it



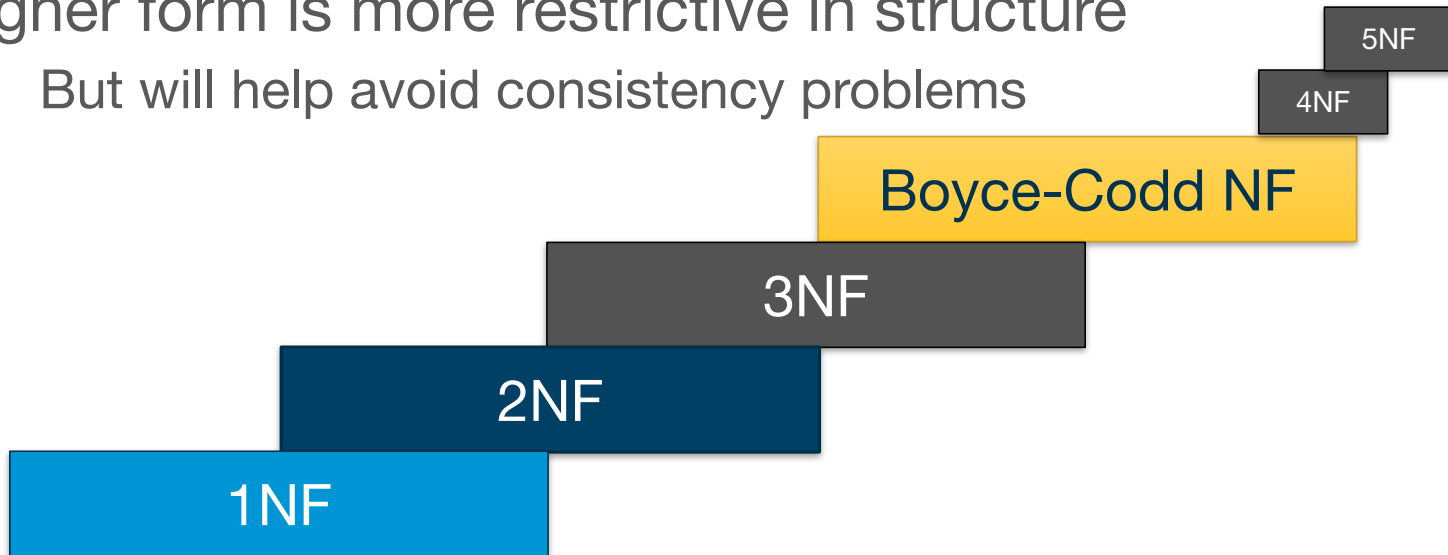
# Normal Forms

Normal forms give us a hierarchy of rules

- No normalization - unconstrained, messy data
- First Normal Form - removes some redundancy
- Second Normal Form - removes more redundancy... etc

Higher form is more restrictive in structure

- But will help avoid consistency problems



# First Normal Form (1NF)

Attributes should be atomic and tables should have no repeating groups

Each field only has one value

No columns repeat the same type of information

No duplicate rows in the table; order doesn't matter

# 1NF Examples

Attributes should be atomic and tables should have no repeating groups

Do these violate 1NF?

*Not atomic!*

①

Customer ID	First Name	Surname	Telephone Number
123	Pooja	Singh	555-861-2025, 192-122-1111
456	San	Zhang	(555) 403-1659 Ext. 53, 182-929-2929
789	John	Doe	555-808-9633

②

Customer ID	First Name	Surname	<u>TNumber1</u>	<u>TNumber2</u>
123	Pooja	Singh	555-861-2025	192-122-1111
456	San	Zhang	(555) 403-1659 Ext. 53	182-929-2929
789	John	Doe	555-808-9633	

TN 3

# 1NF Examples

Attributes should be atomic and tables should have no repeating groups

Do these violate 1NF?

Customer ID	First Name	Surname	Telephone Number
123	Pooja	Singh	555-861-2025, 192-122-1111
456	San	Zhang	(555) 403-1659 Ext. 53; 182-929-2929
789	John	Doe	555-808-9633

Both are bad!

Customer ID	First Name	Surname	TNumber1	TNumber2
123	Pooja	Singh	555-861-2025	192-122-1111
456	San	Zhang	(555) 403-1659 Ext. 53	182-929-2929
789	John	Doe	555-808-9633	

# 1NF Split or Flatten

Attributes should be atomic and tables should have no repeating ~~groups~~ <sup>columns</sup>

## Possible solutions

<u>Customer ID</u>	<u>First Name</u>	<u>Surname</u>	<u>Telephone Number</u>
123	Pooja	Singh	555-861-2025
123	Pooja	Singh	192-122-1111
456	San	Zhang	182-929-2929
456	San	Zhang	(555) 403-1659 Ext. 53
789	John	Doe	555-808-9633

meets  
1NF

OR

<u>Customer ID</u>	<u>First Name</u>	<u>Surname</u>
123	Pooja	Singh
456	San	Zhang
789	John	Doe

<u>Customer ID</u>	<u>Telephone Number</u>
123	555-861-2025
123	192-122-1111
456	(555) 403-1659 Ext. 53
456	182-929-2929
789	555-808-9633

✓

# 1NF Violations

Generally easy to detect:

1. Check for Column names with a number  
(telephone1, telephone2, course1, course2, etc)
2. Make sure that order of rows doesn't matter ✓
3. Have a primary key to enforce uniqueness across rows

# Second Normal Form (2NF)

No value in a table should be dependent on only **part** of a key that uniquely identifies a row

It must be in 1NF and...

We should **not** be able to derive the value of a column based on only **a part of a Candidate Keys**

- Must hold for all Candidate Keys if there are multiple

# Reminder: Key types

## Superkey of R:

- A (**possibly larger than necessary**) set of attributes that is sufficient to uniquely identify each tuple in  $r(R)$

## **Candidate Key** of R: A “minimal” superkey

- A **minimal set** of attributes to denote uniqueness!
- A Candidate Key is a Superkey but opposite may not be true

**Primary Key:** A specific Candidate Key chosen to represent a relation/table



# 2NF Examples

No value in a table should be dependent on only part of a key that uniquely identifies a row

Does this violate 2NF?

<u>Customer ID</u>	First Name	Surname	<u>Telephone Number</u>
123	Pooja	Singh	555-861-2025
123	Pooja	Singh	192-122-1111
456	San	Zhang	182-929-2929
456	San	Zhang	(555) 403-1659 Ext. 53
789	John	Zhang	555-808-9633

# 2NF Examples

No value in a table should be dependent on only part of a key that uniquely identifies a row

Does this violate 2NF?

<u>Customer ID</u>	First Name	Surname	<u>Telephone Number</u>
123	Pooja	Singh	555-861-2025
123	Pooja	Singh	192-122-1111
456	San	Zhang	182-929-2929
456	San	Zhang	(555) 403-1659 Ext. 53
789	John	Zhang	555-808-9633

Yes!

- Our Key is (Customer ID, Telephone Number), but from Customer ID alone we could uniquely identify the name
- We could make  $\text{func}(\text{CustomerID}) \rightarrow (\text{First Name}, \text{Surname})$

In general, better to use the splitting method for 1NF

# 2NF vs 1NF

Why do we care??

1NF

<u>Customer ID</u>	First Name	Surname	<u>Telephone Number</u>
123	Pooja	Singh	555-861-2025
123	Pooja	Singh	192-122-1111
456	San	Zhang	182-929-2929
456	San	Zhang	(555) 403-1659 Ext. 53
789	John	Zhang	555-808-9633

VS

2NF

<u>Customer ID</u>	First Name	Surname
123	Pooja	Singh
456	San	Zhang
789	John	Zhang

<u>Customer ID</u>	<u>Telephone Number</u>
123	555-861-2025
123	192-122-1111
456	(555) 403-1659 Ext. 53
456	182-929-2929
789	555-808-9633

# 2NF vs 1NF

Redundant data can lead to inconsistencies if it is only partially updated!

1NF

<u>Customer ID</u>	First Name	Surname	<u>Telephone Number</u>
123	Pooja	Singh	555-861-2025
123	Pooja	Sing	192-122-1111
456	San	Zhang	182-929-2929
456	San	Zhang	(555) 403-1659 Ext. 53
789	John	Zhang	555-808-9633

VS

2NF

<u>Customer ID</u>	First Name	Surname
123	Pooja	Singh
456	San	Zhang
789	John	Zhang

<u>Customer ID</u>	<u>Telephone Number</u>
123	555-861-2025
123	192-122-1111
456	(555) 403-1659 Ext. 53
456	182-929-2929
789	555-808-9633

# More 2NF Examples

<u>Manufacturer</u>	<u>Model</u>	Price	<u>Manufacturer country</u>
Forte	X-Prime	50	Italy
Forte	Ultraclean	50	Italy
Dent-o-Fresh	EZbrush	65	USA
Brushmaster	SuperBrush	34	USA
Kobayashi	ST-60	22	Japan
Hoch	Toothmaster	18	Germany
Hoch	X-Prime	50	Germany

# More 2NF Examples

## This avoids **Update Anomalies**

- Previously we would have had to scan all tuples if a manufacturer moved to a different country to ensure consistency

<u>Manufacturer</u>	<u>Model</u>	Price
Forte	X-Prime	45
Forte	Ultraclean	50
Dent-o-Fresh	EZbrush	65
Brushmaster	SuperBrush	34
Kobayashi	ST-60	22
Hoch	Toothmaster	18
Hoch	X-Prime	22

<u>Manufacturer</u>	Country
Forte	Italy
Dent-o-Fresh	USA
Brushmaster	USA
Kobayashi	Japan
Hoch	Germany

# Third Normal Form (3NF)

No value should be able to be derived based on another non-key field

It must be in 2NF and...

all **non-prime attributes** depend only on the **candidate keys** and do not have a **transitive dependency** on another key

# 3NF Intuition

No value should be able to be derived based on another non-key field

What is the redundant information in this table?

<u>Customer ID</u>	First Name	Surname	Birthday	Age	Fav Color
123	Pooja	Singh	1/4/1984	37	Blue
456	San	Zhang	3/15/2001	19	Blue
789	John	Zhang	11/12/2006	14	Buff



# 3NF Intuition

No value should be able to be derived based on another non-key field

What is the redundant information in this table?

<u>Customer ID</u>	First Name	Surname	Birthday	Age	Fav Color
123	Pooja	Singh	1/4/1984	37	Blue
456	San	Zhang	3/15/2001	19	Blue
789	John	Zhang	11/12/2006	14	Buff

If we know Birthday, we can calculate Age -> there is an obvious dependency between them! Can remove Age.

# 3NF Intuition

No value should be able to be derived based on another non-key field

What is the redundant information in this table?

<u>Tournament</u>	<u>Year</u>	<u>Winner</u>	<u>Winner's Birthplace</u>
Indiana Invitational	1998	<u>Al Fredrickson</u>	Ohio
Cleveland Open	1999	Bob Albertson	New York
Des Moines Masters	1999	<u>Al Fredrickson</u>	Ohio
Indiana Invitational	1999	Chip Masterson	Kentucky

*Handwritten notes:* Red circles around "Al Fredrickson" and "Ohio" in the first and third rows. Red arrows point from the "Winner" column to the "Winner's Birthplace" column. A red box highlights the first two columns of the first row. A red box highlights the "Winner" and "Winner's Birthplace" columns, with "A. F." and "Ohio" written below them, and vertical ellipses indicating other rows.

# 3NF Intuition

No value should be able to be derived based on another non-key field

What is the redundant information in this table?

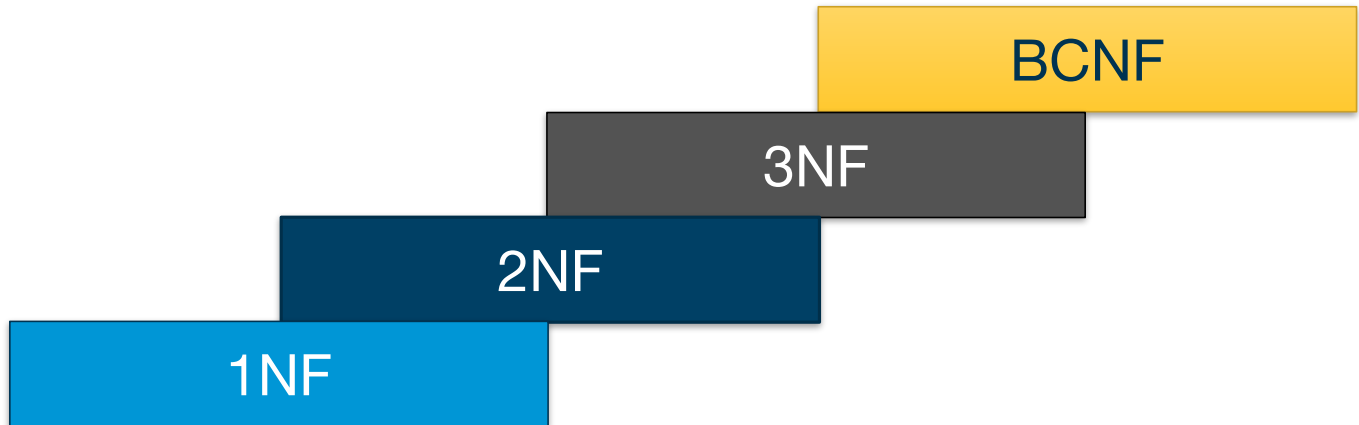
<u>Tournament</u>	<u>Year</u>	<u>Winner</u>	<u>Winner's Birthplace</u>
Indiana Invitational	1998	Al Fredrickson	Ohio
Cleveland Open	1999	Bob Albertson	New York
Des Moines Masters	1999	Al Fredrickson	New Jersey
Indiana Invitational	1999	Chip Masterson	Kentucky

The {Winner's Birthplace} attribute can be determined based on Winner, which is not a Candidate Key for the table. Need to split!

# Normal Form Redundancy

1NF and 2NF - eliminate redundancy **across** rows

3NF, BCNF - also eliminate redundancy **within** rows



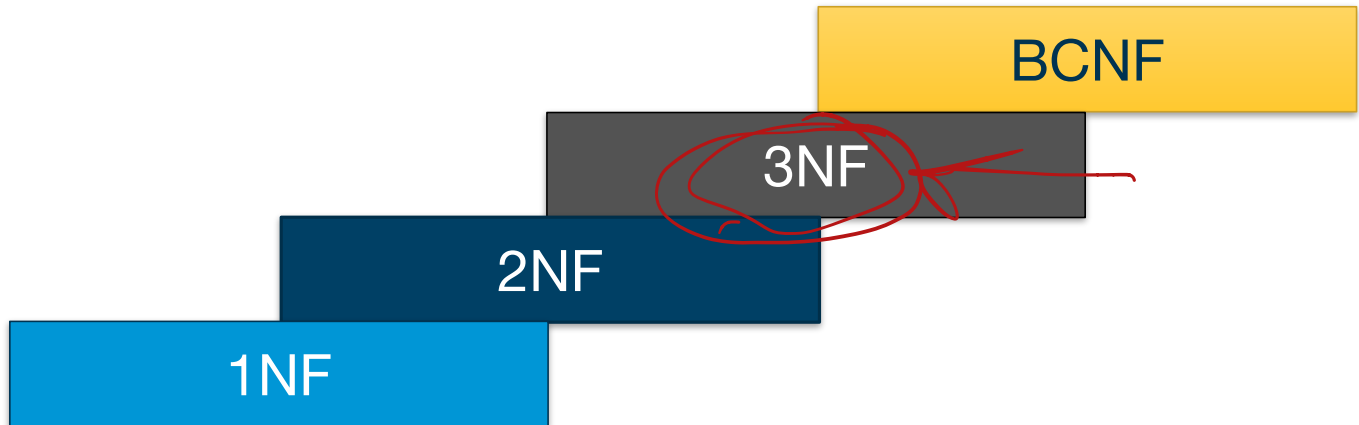
Good and Bad  
Schemas

**Functional  
Dependencies**

**Normal Forms  
based on  
Functional  
Dependencies**

# Dependencies

How can we formally represent dependencies between Attributes in a Relation?



# Functional Dependencies

Use functional dependencies! (abbreviated FD)

We say a set of attributes **X** functionally determines an attribute **Y** if *given the values of X we always know the only possible value of Y.*

- Notation:  $X \rightarrow Y$
- X functionally determines Y
- Y is functionally dependent on X

Example:

-  $GWID \rightarrow Name$

-  $\{GWID, CourseID, Semester, Year\} \rightarrow Grade, ProfName$

# Sets of Functional Dependencies

Some more functional dependencies

- 1 - {GWID} → {NAME, ADDRESS, MAJOR}
- 2 - <sup>CSCI BS.</sup> {MAJOR} → {DEPT\_NAME, DEPT\_CHAIR}

From above dependencies, we can infer

- 3 - {GWID} → {DEPT\_NAME, DEPT\_CHAIR} ✓

We can do math on functional dependencies!

A functional dependency “holds” if it must be true for all legal relations



# Functional Dependency Ops

Armstrong's Axioms: where A, B, C are sets of attributes

- Reflexive rule: if  $B \subseteq A$ , then  $A \rightarrow B$  (if B is subset of A)  $(\text{GWID}, \text{Name}) \rightarrow \text{Name}$
- Augmentation rule: if  $A \rightarrow B$ , then  $CA \rightarrow CB$
- Transitivity rule: if  $A \rightarrow B$ , and  $B \rightarrow C$ , then  $A \rightarrow C$

These rules are

- Sound and complete — generate all functional dependencies that hold.

$\{\text{GWID}\} \rightarrow \{\text{Name, Address, Major}\}$

$\{\text{Major}\} \rightarrow \{\text{Dept\_Name, Dept\_Chair}\}$

$\{\text{GWID, CourseID, Semester, Year}\} \rightarrow \text{Grade}$

$\text{GWID, CourseID} \rightarrow \text{Name, Address, Major, CourseID}$

# Functional Dependency Ops

**Armstrong's Axioms:** where  $A, B, C$  are sets of attributes

- **Reflexive rule:** if  $B \subseteq A$ , then  $A \rightarrow B$  (if  $B$  is subset of  $A$ )
- **Augmentation rule:** if  $A \rightarrow B$ , then  $CA \rightarrow CB$
- **Transitivity rule:** if  $A \rightarrow B$ , and  $B \rightarrow C$ , then  $A \rightarrow C$

These rules are

- Sound and complete — generate all functional dependencies that hold.

Bonus rules to make life easier:

- **Union rule:** If  $\alpha \rightarrow \beta$  holds and  $\alpha \rightarrow \gamma$  holds, then  $\alpha \rightarrow \beta \gamma$  holds.
- **Decomposition rule:** If  $\alpha \rightarrow \beta \gamma$  holds, then  $\alpha \rightarrow \beta$  holds and  $\alpha \rightarrow \gamma$  holds.
- **Pseudotransitivity rule:** If  $\alpha \rightarrow \beta$  holds and  $\gamma \beta \rightarrow \delta$  holds, then  $\alpha \gamma \rightarrow \delta$  holds.

# Definition: Closure of a Set of FD's

Defn. Let F be a set of FD's.

Its closure, F+, is the set of all FD's:

$\{X \rightarrow Y \mid X \rightarrow Y \text{ is derivable from } F \text{ by}$   
Armstrong's Axioms}

Two sets of dependencies F and G are equivalent  
if F+=G+

- i.e., their closures are equal
- i.e., the same sets of FDs can be inferred from each

# Example Closure

What FDs can we infer?

$R = (A, B, C, G, H, I)$

$F = \{A \rightarrow B$

$A \rightarrow C$

$CG \rightarrow H$

$CG \rightarrow I$

$B \rightarrow H\}$

Reflexive rule: if  $\beta \subseteq \alpha$ , then  $\alpha \rightarrow \beta$

Augmentation rule: if  $\alpha \rightarrow \beta$ , then  $\gamma \alpha \rightarrow \gamma \beta$

Transitivity rule: if  $\alpha \rightarrow \beta$ , and  $\beta \rightarrow \gamma$ , then  $\alpha \rightarrow \gamma$

F+

$A \rightarrow B$   
 $B \rightarrow H \Rightarrow A \rightarrow H$

$A \rightarrow C$   
 $AG \rightarrow CG$   
 $CG \rightarrow I \Rightarrow$

# Example Closure

$R = (A, B, C, G, H, I)$

$F = \{$   
   $A \rightarrow B$   
   $A \rightarrow C$   
   $CG \rightarrow H$   
   $CG \rightarrow I$   
   $B \rightarrow H\}$

**Reflexive rule:** if  $\beta \subseteq \alpha$ , then  $\alpha \rightarrow \beta$

**Augmentation rule:** if  $\alpha \rightarrow \beta$ , then  $\gamma \alpha \rightarrow \gamma \beta$

**Transitivity rule:** if  $\alpha \rightarrow \beta$ , and  $\beta \rightarrow \gamma$ , then  $\alpha \rightarrow \gamma$

A few members of  $F^+$  include:

- $A \rightarrow H$ 
  - by transitivity from  $A \rightarrow B$  and  $B \rightarrow H$
- $AG \rightarrow I$ 
  - by augmenting  $A \rightarrow C$  with  $G$ , to get  $AG \rightarrow CG$   
and then transitivity with  $CG \rightarrow I$
- $CG \rightarrow HI$ 
  - by augmenting  $CG \rightarrow I$  to infer  $CG \rightarrow CGI$ ,  
and augmenting of  $CG \rightarrow H$  to infer  $CGI \rightarrow HI$ ,  
and then transitivity

Good and Bad  
Schemas

**Functional  
Dependencies**

**Normal Forms  
based on  
Functional  
Dependencies**

# Redefining 2NF

Using Functional Dependencies and Closures lets us more precisely define our Normal Forms

**Second Normal Form:** For every  $X \rightarrow A$  that holds over relationship schema  $R$ , where  $A$  is a non-prime attribute (i.e.,  $A$  is not an attribute in any candidate key)

1. either  $A \in X$  (it is trivial), or  $AB \rightarrow A$
2.  $X$  is a superkey for  $R$ , or
3.  $X$  is transitively dependent on a super key  $R$

*Easier to think of the opposite:* There cannot be  $X \rightarrow A$  where  $X$  is a partial candidate key for  $R$

- Says nothing about non-prime to non-prime dependencies!

# 2NF Violations

~~X~~ → A

<u>ID</u>	<u>First Name</u>	<u>Surname</u>	<u>Telephone Number</u>
123	Pooja	Singh	555-861-2025
123	Pooja	Singh	192-122-1111
456	San	Zhang	182-929-2929
456	San	Zhang	(555) 403-1659 Ext. 53
789	John	Zhang	555-808-9633

ID ⇒ FN, SN



<u>Tournament</u>	<u>Year</u>	<u>Winner</u>	<u>Winner's Birthplace</u>
Indiana Invitational	1998	Al Fredrickson	Ohio
Cleveland Open	1999	Bob Albertson	New York
Des Moines Masters	1999	Al Fredrickson	Ohio
Indiana Invitational	1999	Chip Masterson	Kentucky

Winner → Birthplace





# 2NF Violations

<u>ID</u>	<u>First Name</u>	<u>Surname</u>	<u>Telephone Number</u>
123	Pooja	Singh	555-861-2025
123	Pooja	Singh	192-122-1111
456	San	Zhang	182-929-2929
456	San	Zhang	(555) 403-1659 Ext. 53
789	John	Zhang	555-808-9633

**ID** -> {**First Name, LastName**}  
Violates 2NF  
since **ID** is a  
partial  
Candidate Key

<u>Tournament</u>	<u>Year</u>	<u>Winner</u>	<u>Winner's Birthplace</u>
Indiana Invitational	1998	Al Fredrickson	Ohio
Cleveland Open	1999	Bob Albertson	New York
Des Moines Masters	1999	Al Fredrickson	Ohio
Indiana Invitational	1999	Chip Masterson	Kentucky

No 2NF  
violation

# Redefining 3NF

Third Normal Form (3NF): For every  $X \rightarrow A$  that holds over relationship schema  $R$ ,

1. either  $A \in X$  (it is trivial), or ✓
2.  $X$  is a ~~superkey~~ for  $R$ , or ✓
3.  $A$  is a member of some <sup>Candidate</sup> key for  $R$  ✓

*Easier to think of the opposite:* There cannot be  $X \rightarrow A$  where  $X$  is not a full candidate key for  $R$

“Every non-key attribute must provide a fact about the Key, the whole Key, and nothing but the Key... so help me Codd”

# 3NF Violations

Birthday → Age

<u>Customer ID</u>	<u>First Name</u>	<u>Surname</u>	<u>Birthday</u>	<u>Age</u>	<u>Fav Color</u>
123	Pooja	Singh	1/4/1984	37	Blue
456	San	Zhang	3/15/2001	19	Blue
789	John	Zhang	11/12/2006	14	Buff



<u>Tournament</u>	<u>Year</u>	<u>Winner</u>	<u>Winner's Birthplace</u>
Indiana Invitational	1998	Al Fredrickson	Ohio
Cleveland Open	1999	Bob Albertson	New York
Des Moines Masters	1999	Al Fredrickson	Ohio
Indiana Invitational	1999	Chip Masterson	Kentucky

Winner → Birthplace



# 3NF Violations

<u>Customer ID</u>	<u>First Name</u>	<u>Surname</u>	<u>Birthday</u>	<u>Age</u>	<u>Fav Color</u>
123	Pooja	Singh	1/4/1984	37	Blue
456	San	Zhang	3/15/2001	19	Blue
789	John	Zhang	11/12/2006	14	Buff

**Birthday** → **Age**  
holds, but  
**Birthday** is not  
a superkey

<u>Tournament</u>	<u>Year</u>	<u>Winner</u>	<u>Winner's Birthplace</u>
Indiana Invitational	1998	Al Fredrickson	Ohio
Cleveland Open	1999	Bob Albertson	New York
Des Moines Masters	1999	Al Fredrickson	Ohio
Indiana Invitational	1999	Chip Masterson	Kentucky

**Winner** → **Birthplace**  
holds, but  
**Winner** is not a  
superkey

# Normal Forms 1-3

**1NF:** Attributes should be atomic and tables should have no repeating groups

- *Prevents messiness within a cell and repetition of rows*

**2NF:** There cannot be  $X \rightarrow A$  where  $X$  is a partial candidate key for  $R$

- Doesn't forbid non-prime to non-prime dependencies
- Prevents repetition of cells across rows

**3NF:** There cannot be  $X \rightarrow A$  where  $X$  is not a full candidate key for  $R$

- Only allows dependencies on Keys
- Prevents repetition of data within a row