# Behind the Scene: Know whom we should blame?

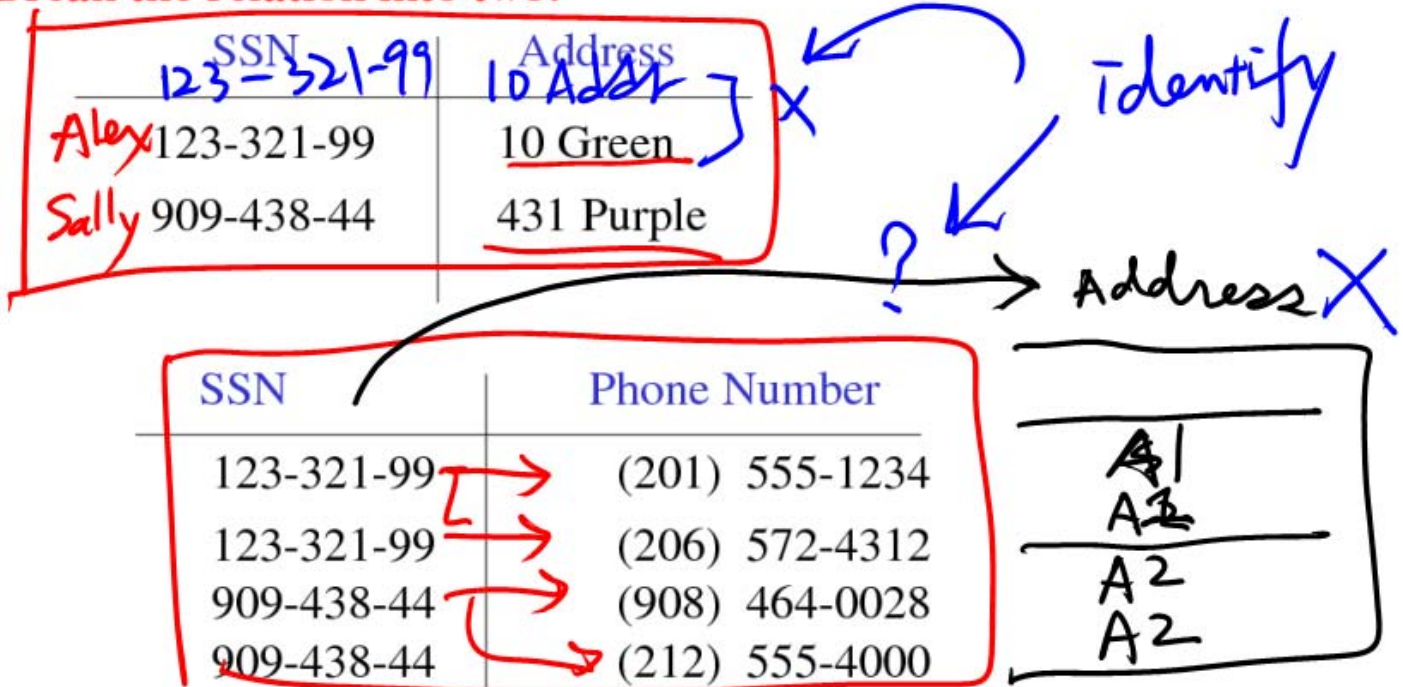| Normal form | Defined by | Brief definition |
|---|---|---|
| First normal form (1NF) | Two versions: E.F. Codd (1970), C.J. Date (2003)[12] | Table faithfully represents a relation and has no "repeating groups" |
| Second normal form (2NF) | E.F. Codd (1971)[13] | No non-prime attribute in the table is functionally dependent on a part (proper subset) of a candidate key |
| Third normal form (3NF) | E.F. Codd (1971)[14]; see also Carlo Zaniolo's equivalent but differently-expressed definition (1982)[15] | Every non-prime attribute is non-transitively dependent on every key of the table |
| Boyce-Codd normal form (BCNF) | Raymond F. Boyce and E.F. Codd (1974)[16] | Every non-trivial functional dependency in the table is a dependency on a superkey |
| Fourth normal form (4NF) | Ronald Fagin (1977)[17] | Every non-trivial multivalued dependency in the table is a dependency on a superkey |
| Fifth normal form (5NF) | Ronald Fagin (1979)[18] | Every non-trivial join dependency in the table is implied by the superkeys of the table |
| Domain/key normal form (DKNF) | Ronald Fagin (1981)[19] | Every constraint on the table is a logical consequence of the table's domain constraints and key constraints |
| Sixth normal form (6NF) | Chris Date, Hugh Darwen, and Nikos Lorentzos (2002)[20] | Table features no non-trivial join dependencies at all (with reference to generalized join operator) |

10

# Our Attack Plan

- Motivation
- Functional dependencies & keys
- Reasoning with FDs and keys
- Desirable properties of schema refinement
- Various normal forms and the trade-offs
  - BCNF, 3rd normal form, 4th normal form, etc.
- Putting all together: how to design DB schema

# Functional Dependencies and Keys

# Better Designs Exist

**Break the relation into two:**

| SSN | Address |
|-----|---------|
| *Alex* 123-321-99 | 10 Green |
| *Sally* 909-438-44 | 431 Purple |

123-321-99  10 Addr  *(handwritten)*

*Identify*

? → Address ✗

| | |
|---|---|
| A1 | |
| A2 | |
| A2 | |
| A2 | |

| SSN | Phone Number |
|-----|--------------|
| 123-321-99 | (201) 555-1234 |
| 123-321-99 | (206) 572-4312 |
| 909-438-44 | (908) 464-0028 |
| 909-438-44 | (212) 555-4000 |

13

# Reminder

Tutorial #2.

→ Today! 4:30-5:30pm

1302 SC.

# P-Fun Topics ?

— Translation ER to Rel. model.

— { RAT

10%

R.A.
↓
SQL
↓
[SQLite]

⇒

R.A.
↓
[SQLite]

— Attr Closure (today) — BCNF (tod.)
— F.D. closuer ( " ) — " " " "

# Functional Dependencies

- A form of constraint (hence, part of the schema)
- Finding them is part of the database design
- Used heavily in schema refinement

Definition:

(SSN) Name  Addr  Phone

10 green
10 gram     Name  Address

If two tuples agree on the attributes

Name

$A_1, A_2, \ldots A_n$

T1  Alex
T2  Alex

then they must also agree on the attributes

Address  $B_1, B_2, \ldots B_m$

Formally:  $A_1, A_2, \ldots A_n \longrightarrow B_1, B_2, \ldots B_m$

14

# Examples

| EmpID | Name | Phone | Position |
|-------|------|-------|----------|
| E0045 | Smith | 1234 | Clerk |
| E1847 | John | 9876 | Salesrep |
| E1111 | Smith | 9876 | Salesrep |
| E9999 | Mary | 1234 | Lawyer |

- EmpID $\longrightarrow$ Name, Phone, Position
- ✓ Position $\longrightarrow$ Phone
- but Phone $\not\longrightarrow$ Position
  1234    clerk, Lawyer

15

# In General

- To check if A ⟶ B violation:

Erase all other columns

*[handwritten: phone ← position]*

| ... | A | ... | B | |
|---|---|---|---|---|
| | X1 | | Y1 | |
| | X2 | | Y2 | |
| | ... | | ... | |

*[handwritten: one ← many / one]*

- check if the remaining relation is many-one
  (called *functional* in mathematics)

16

# Example

| EmpID | Name | Phone | Position |
|-------|------|-------|----------|
| E0045 | Smith | 1234 ← | Clerk ✓ |
| E1847 | John | 9876 ← | Salesrep ✓ |
| E1111 | Smith | 9876 ← | Salesrep ✓ |
| E9999 | Mary | 1234 ← | lawyer ✓ |

*No violation for Pos → Phone*

More examples:

Product:    name → price, manufacturer
Person:     ssn → name, age
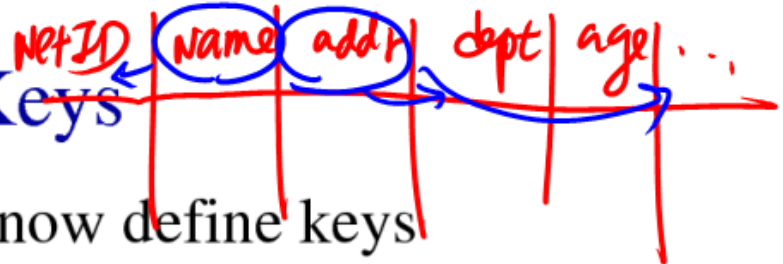Company:  name → stock price, president

# Q: From this, can you conclude phone → SSN?

*a phone is only used by ONE person.*

| SSN | Phone Number |
|---|---|
| 123-321-99 Alex | (201) 555-1234 |
| 123-321-99 Alex | (206) 572-4312 |
| 909-438-44 | (908) 464-0028 |
| 909-438-44 | (212) 555-4000 |
| 123-321-88 Alex Junior | (201) 555-1234 |

F.D, stated at schema design

⟹ assertion

18

# Relation Keys

NetID | Name | addr | dept | age | . . .

- After defining FDs, we can now define keys

- Key of a relation R is a set of attributes that
    - functionally determines all attributes of R   $\{N, A\} \rightarrow \{NetID, dept, age, \}$
    - none of its subsets determines all attributes of R

- Superkey
    - a set of attributes that contains a key

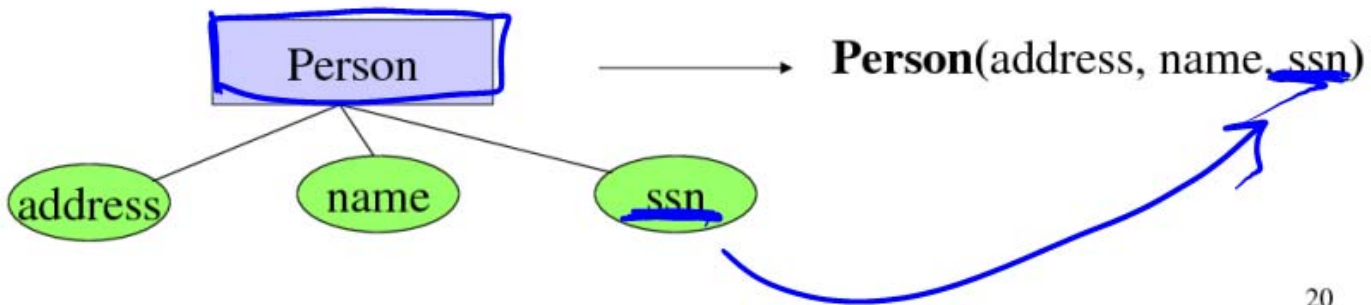- We will need to know the keys of the relations in a DB schema, so that we can refine the schema

|        | (Name, Addr) | (Netid) | (NetID, dept) |
|--------|--------------|---------|---------------|
| key    | ✓            | ✓       | ✗             |
| S. key | ✓            | ✓       | ✓             |

19

# Finding the Keys of a Relation

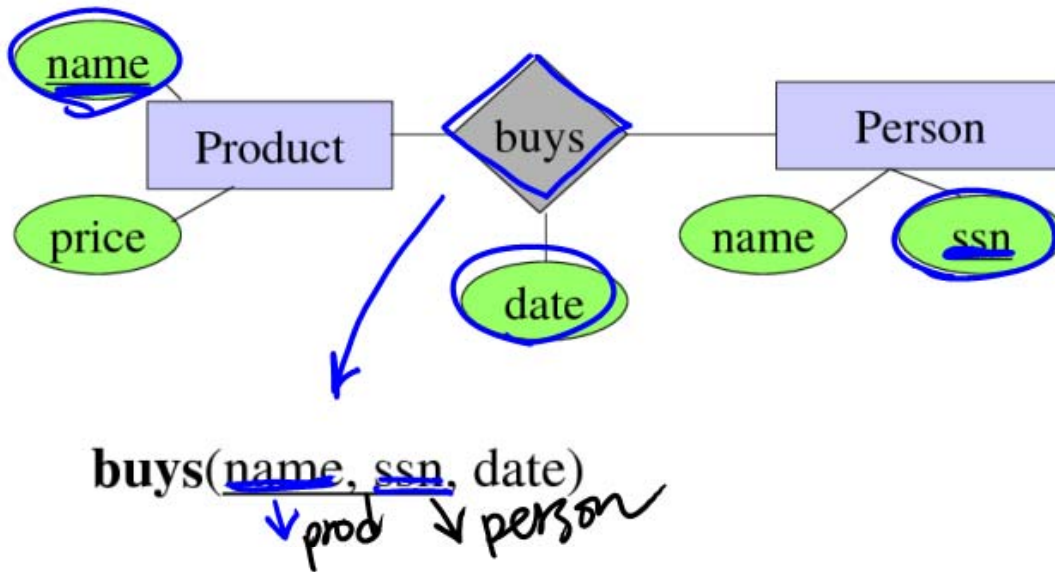Given a relation constructed from an E/R diagram, what is its key?

Rules:
1. If the relation comes from an entity set,
   the key of the relation is the set of attributes which is the
   key of the entity set.



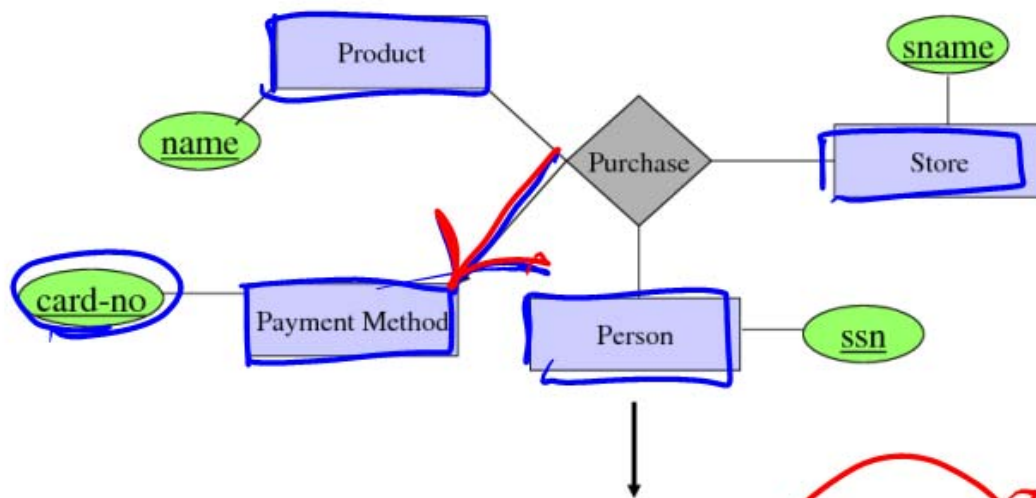**Person**(address, name, ssn)

20

# Finding the Keys

2. If the relation comes from a many-many relationship,
   the key of the relation include the set of all attribute keys in the
   relations corresponding to the entity sets
   (and additional attributes if necessary)



**buys**(name, ssn, date)
    ↓prod  ↓person

21

# Finding the Keys

But: if there is an arrow from the relationship to E, then
we don't need the key of E as part of the relation key.



**Purchase**(name , sname, ssn, card-no)

22

# Finding the Keys

More specific rules:

- Many-one, one-many, one-one relationships
- Multi-way relationships
- Weak entity sets

(Try to find them yourself)

# Reasoning with FDs

1) closure of FD sets
2) closure of attribute sets

# Closure of FD sets

_You_

- Given a relation schema R & a set S of FDs _who gives FDs ?_
  - is the FD f logically implied by S?
- Example
  - R = {A,B,C,G,H,I} 6 attrs.  Name → age,  age → drikable  } Name
  - S = A → B, A → C, CG → H, CG → I, B → H            ↓
  - would A → H be logically implied? _New rule_     drinkable
  - yes (you can prove this, using the definition of FD)
- Closure of S: S+ = all FDs logically implied by S
- How to compute S+?                 _such as_
  - we can use Armstrong's axioms    A → H.

25

# Armstrong's Axioms  Rules

proven by def.

- Reflexivity rule    NetID → dept. addr
  - A1A2...An ➜ a subset of A1A2...An   ⇒) NetID → dept

- Augmentation rule   Name. NetID → dept. addr. Name,
  - A1A2...An ➜ B1B2...Bm, then
    A1A2...An C1C2..Ck ➜ B1B2...Bm C1C2...Ck

- Transitivity rule
  - A1A2...An ➜ B1B2...Bm and    Name → age ⌉
  same B1B2...Bm ➜ C1C2...Ck, then   age → drink ⌋
    A1A2...An ➜ C1C2...Ck    ⇒) Name → drink

26

# Inferring S+ using Armstrong's Axioms

- S+ = S

$$S = \{ \overset{f1}{A \to B}, \overset{f2}{B \to C}, \overset{f3}{AC \to D} \}$$

$S^+$ ?    "form changing"

- Loop
  - foreach f in S, apply reflexivity and augment. rules
  - add the new FDs to S+        so that    $B \to B$
  - foreach pair of FDs in S, apply the transitivity rule    $B \to$
  - add the new FD to S+

- Until S+ does not change any further

① $f_1 . f_2 . Tr \Rightarrow + A \to C.$

② (want to use $AC \to D$)

$f2:$    $A\underline{B} \to A\underline{C}$    $\Rightarrow$  $AA \to D \Rightarrow A \to D$

$f1:$    $AA \to AB$

27

Q1: What do you like best of this cls. that we must keep?

Q2: What .... dislike - - - -

- - - - - - go ?

# Additional Rules

- Union rule *(aug)*
  - X ➔ Y and X ➔ Z, then X ➔ YZ
    (netid)(dept)(netid)(addr)(netid)(dept addr)
  - (X, Y, Z are sets of attributes)

$$XX \rightarrow XY \Rightarrow XX$$
$$XY \rightarrow YZ \rightarrow YZ$$

- Decomposition rule
  - X ➔ YZ, then X ➔ Y and X ➔ Z

- Pseudo-transitivity rule
  - X ➔ Y and YZ ➔ U, then XZ ➔ U

- These rules can be inferred from Armstrong's axioms

28

# Closure of a Set of Attributes

$(name, addr) \longrightarrow ?$

Given a set of attributes $\{A1, ..., An\}$ and a set of dependencies S.
Problem: find all attributes $B$ such that:

(You)

any relation which satisfies S also satisfies:

$A1, ..., An \rightarrow B$

The **closure** of $\{A1, ..., An\}$, denoted $\{A1, ..., An\}^+$,
is the set of all such attributes $B$

?

We will discuss the <u>motivations</u> for attribute closures soon

Is $\{name, addr\}$ a key?

$\{name, addr\}^+ \implies$ all attr.

29

# Algorithm to Compute Closure

Start with X={A1, ..., An}.  $\{name, addr\}$

**Repeat until** X doesn't change  **do**:

if $B_1, B_2, \ldots B_n \longrightarrow C$  is in S, **and**

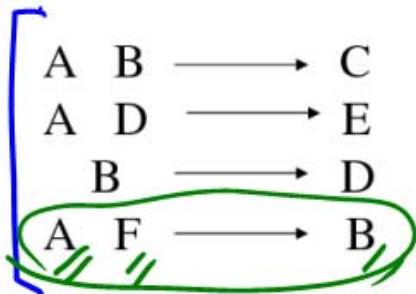$B_1, B_2, \ldots B_n$  are all in X, **and**

C is not in X

**then**

add C to X.  $X = X + \{c\}$

# Example

$R : \langle A, B, C, D, E, F \rangle$    Is $(A, F)$ a key?

A  B ———→ C
A  D ———→ E
   B ———→ D
A  F ———→ B

Name  Addr ———→ age

Closure of {A,B}:    X = {A, B, C, D, E}

Closure of (A, F):    X = {A, F, B, D, C, E}

✓  $(A,F)^+ = \{A, \ldots, F\}$

✗  $(A)^+ = \{ \ldots \}$

✗  $(F)^+ = \{ \ldots \}$

|  | $(A, F)$    $(A)$ ? |
|---|---|
| $AF \rightarrow B$ | A, B, F |
| $B \rightarrow D$ | A, B, D, F |
| $AD \rightarrow E$ | A B D E F |
| $AB \rightarrow C$ | A B C D E F   @Stop |

31

# Usage for Attribute Closure

- Test if X is a superkey
  - compute X+, and check if X+ contains all attrs of R


- Check if X ➜ Y holds
  - by checking if Y is contained in X+

$$Y \subseteq X^{+} \iff X \to Y$$

32

**Desirable Properties of Schema Refinement**

ER

$R \xrightarrow{?} R^*$

original       Normal Form

1) minimize <u>redundancy</u>
2) avoid <u>info loss</u>
3) <u>preserve dependency</u>
4) ensure good query performance

33

# Normal Forms

x set.      string.

x array,    Float.

**First Normal Form** = all attributes are atomic

**Second Normal Form** (2NF) = old and obsolete

SQL.     Ted Codd.

**Boyce Codd Normal Form** (BCNF) ⬅

**Third Normal Form** (3NF)

**Fourth Normal Form** (4NF)

Others...

34

# Boyce-Codd Normal Form

BCNF ? NO

| SSN | addr | phone |
|-----|------|-------|
| Alex | 10 G | 123 |
| Alex | 10 G | 456 |

A simple condition for removing anomalies from relations:

A relation R is in BCNF if and only if:

✓  SSN **Bad** ↛ addr

Whenever there is a nontrivial FD    $A_1, A_2, \ldots A_n \longrightarrow B$
for R, it is the case that $\{ A_1, A_2, \ldots A_n \}$
is a super-key for R.

then
SSN is superkey

In English (though a bit vague):

X ?   SSN → all attr
SSN, addr, phone

Whenever a set of attributes of R is determining another attribute,
it should determine **all** attributes of R.  In Contrast  SSN adar   SSN phone

✓ SSN → addr
✓ SSN a key   ] yes

BCNF ?

| SSN | adar |
|------|------|
| Alex | 10 G |
| Alex | 10 G |

| SSN | phone |
|------|-------|
| Alex | 123 |
| Alex | 456 |

35

# Example

**F.D.** **violate BCNF**

Name ① SSN ② ✗ Phone Number

| Name | SSN | Phone Number |
|------|-----|--------------|
| Fred | 123-321-99 | P1 (201) 555-1234 |
| Fred | 123-321-99 | P2 (206) 572-4312 |
| Joe | 909-438-44 | (908) 464-0028 |
| Joe | 909-438-44 | (212) 555-4000 |

→ { , }

What are the dependencies?

    SSN → Name

What are the keys?

Is it in BCNF?

36

# Decompose it into BCNF

Did u remove the F.D.?

not BCNF

$A \rightarrow B$

T₁ key

BCNF ✓

| SSN | Name |
|-----|------|
| 123-321-99 | Fred |
| 909-438-44 | Joe |

SSN ⟶ Name

(addr)

| SSN | Name | phone |
|-----|------|-------|
| A | B | C |

C A B

A C

T₂

BCNF ?

| SSN | Phone Number |
|-----|--------------|
| 123-321-99 | (201) 555-1234 |
| 123-321-99 | (206) 572-4312 |
| 909-438-44 | (908) 464-0028 |
| 909-438-44 | (212) 555-4000 |

key = (SSN, phone)

37

# What About This?

*BCNF ?*  *Yes*

| Name | Price | Category |
|------|-------|----------|
| Gizmo | $19.99 | gadgets |
| OneClick | $24.99 | camera |

Name ⟶ Price, Category

*key*

38

# BCNF Decomposition

Find a dependency that violates the BCNF condition:

$$A_1, A_2, \ldots A_n \longrightarrow B_1, B_2, \ldots B_m$$

*SSN* *Name* *phone C*

Heuristics: choose $B_1, B_2, \ldots B_m$ "as large as possible"

Decompose:

Others *phone* | A's *SSN* | B's

*SSN* *Name* *SSN*

**R1**     **R2**

Continue until there are no BCNF violations left.

39

# Example Decomposition

$SSN \rightarrow Age \rightarrow D.F.W$  ?

F.D. ✓

Person:

P | Name | SSN | Age | EyeColor | PhoneNumber | DraftWorthy

Functional dependencies:

✓ ① F.D. 1    SSN ⟶ ( Name, Age, Eye Color ) phone.   D.F.W

BNCF: P₁ Person1(SSN, Name, Age, EyeColor)  A+B⁺
      P₂ Person2(SSN, PhoneNumber)   ?   A+C.

A    Ⓑ D.F.W    C

What if we also had an attribute Draft-worthy, and the FD:

P11  (Age, d.w)

② F.D. 2   Age ⟶ Draft-worthy

P12  (SSN, ⋯, age)

40

# BCNF Decomposition: The Algorithm

- Input: relation R, set S of FDs over R $\begin{cases} \text{FD1} & \text{SSN} \rightarrow \dots \\ \text{FD2} & \text{age} \rightarrow \dots \end{cases}$

1) Compute S+

2) Compute keys for R (from ER or from S+)

3) Use S+ and keys to check if R is in BCNF, if not:

     *one*   SSN → age, name, ec,

     a) pick a violation FD f: A → B

                                    SSN → age, …,

     b) expand B as much as possible, by computing A+   D.F.W

     c) create R1 = A union B, R2 = A union (others in R)

     d) compute all FDs over R1, using R and S+,

          then compute keys for R1. Repeat similarly for R2

     e) Repeat Step 3 for R1 and R2

4) Stop when all relations are BCNF or are two-attributes

41

# Q: Is BCNF unique? NO

| SSN | NetID | phone |
|-----|-------|-------|
| ✓ 111 | Alex1 | 123 |
| ○ 222 | Barb1 | 456 |
| ✓ 111 | Alex1 | 321 |

F.D.1   SSN $\xrightarrow{(A)}$ NetID$^{(B)}$   $\overset{(C)}{phone}$

F.D.2   NetID → SSN

① F.D.1   | SSN | NetID |   | SSN | phone |

② F.D.2   | NetID | SSN |   | NI. | phone |

42

# Q: Does BCNF always exist?

All two-attr tables are in BCNF.



$A \rightarrow B$ . A must be key

$B \rightarrow A$ , B must be key

43

# Properties of BCNF

- BCNF removes certain types of redundancy
  - those caused by adding many-many or one-many relations
- For examples of redundancy that it cannot remove, see "multivalued redundancy"
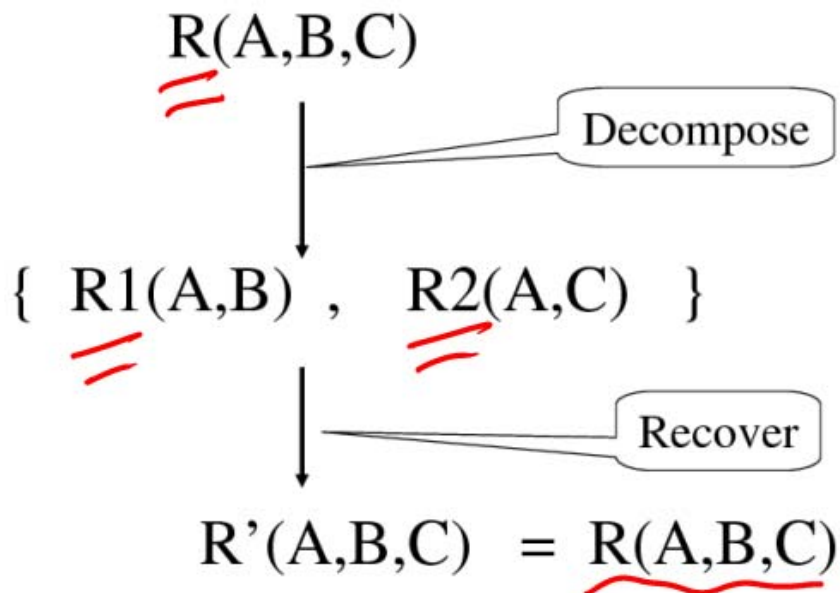- BCNF avoids information loss

Enrollment

| name | dept | course | instr |
|------|------|--------|-------|
| Alex | CS | CS412 | Johnson |
| Ben | EE | CS412 | Johnson |
| Alex | CS | CS423 | Many |

many

SSN addr phone
111 → 123
111 → 321

many

many

44

# Lossless Decompositions

A decomposition is *lossless* if we can recover:

$$R(A,B,C)$$

Decompose

$$\{\ R1(A,B)\ ,\quad R2(A,C)\ \}$$

Recover

$$R'(A,B,C)\ =\ R(A,B,C)$$

R' is in general larger than R.  Must ensure R' = R

# Decomposition Based on BCNF is Necessarily Lossless

*if (Alex, Mb, 'bud light')*

R(A, B, C),     $A \rightarrow C$

*Name  bar  favories*

*Name   favorite*

$C = c'$

$b \neq b'$

BCNF:  R1(A,B),   R2(A,C)

*$R_1$(Name, bar)     $R_2$(Name, Fav)*

Some tuple   (a,b,c)  in R  *(Alex, J, bud)*  (a,b',c') also in R   *(Alex, G, bud)*

decomposes into   (a,b)  in R1        (a,b') also in R1

and     (a,c)  in R2        (a,c') also in R2

Recover tuples in R:  (a,b,c),        (a,b,c'), (a,b',c), (a,b',c') also in R ?

*No extra tuples*

Can  (a,b,c') be a bogus tuple?  What about (a,b',c') ?

46

# However,

- BCNF is not always dependency preserving
- In fact, some times we cannot find a BCNF decomposition that is dependency preserving
- Can handle this situation using 3NF
- See next few slides for example

# Behind the Scene: The Great Debate of '75

*Ted Codd : R. M. 1970*

- The network/COBOL camp:
    - DBTG (Database Task Group, under CODASYL) 1971
    - closely aligned with COBOL
    - DBTG Report would standardize network model
    - Bachman (for network model) got Turing award in 1973
- The relational camp:
    - Codd's paper in 1970
    - resistance even within IBM
    - First implementations, 1973: System R (IBM), INGRES (Berkeley)
    - System R at IBM San Jose Lab

- The "Great Debate" in 1975 SIGMOD conf.

- Codd got Turing award in 1981

48

# Behind the Scene:
## Arguments Against the Other Side?

*Network*

- COBOL/CODASYL → Relational
  - too mathematical (to understand)

*Network*

- Relational → COBOL/CODASYL
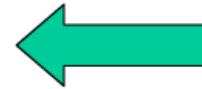  - too complicated (to program)

49

# Normal Forms

**First Normal Form** = all attributes are atomic
**Second Normal Form** (2NF) = old and obsolete

**Boyce Codd Normal Form** (BCNF)
**Third Normal Form** (3NF)
**Fourth Normal Form** (4NF)

Others...

50

# 3NF: A Problem with BCNF

| Unit | Company | Product |
|------|---------|---------|
|      |         |         |

FD's:  Unit $\rightarrow$ Company;    Company, Product $\rightarrow$ Unit

So, there is a BCNF violation, and we decompose.

| Unit | Company |
|------|---------|
|      |         |

Unit $\rightarrow$ Company

| Unit | Product |
|------|---------|
|      |         |

No  FDs

51

# So What's the Problem?

| Unit | Company |
|------|---------|
| Galaga99 | UI |
| Bingo | UI |

| Unit | Product |
|------|---------|
| Galaga99 | databases |
| Bingo | databases |

No problem so far. All *local* FD's are satisfied.

Let's put all the data back into a single table again:

| Unit | Company | Product |
|------|---------|---------|
| Galaga99 | UI | databases |
| Bingo | UI | databases |

**Violates the dependency:   company, product -> unit!**

# Preserving FDs

- What if, when a relation is decomposed, the X of an $X \rightarrow Y$ ends up only in one of the new relations and the Y ends up only in another?

- Such a decomposition is not "dependency-preserving."

- Goal: Always have FD-preserving decompositions

# Solution: 3rd Normal Form (3NF)

A simple condition for removing anomalies from relations:

A relation R is in 3rd normal form if :

Whenever there is a nontrivial dependency $A_1, A_2, ..., A_n \rightarrow B$ for R , then $\{A_1, A_2, ..., A_n\}$ is a super-key for R, or B is part of a key.

# 3NF (General Definition)

- A relation is in Third Normal Form (3NF) if whenever X→A holds, either X is a superkey, or A is a prime attribute.

*Informally: everything depends on the key or is in a key.*

- Despite the thorny technical definitions that lead up to it, 3NF is intuitive and not hard to achieve. *Aim for it in all designs unless you have strong reasons otherwise.*

# 3NF vs. BCNF

- R is in BCNF if whenever X→A holds, then X is a superkey.

- Slightly stronger than 3NF.

- Example: R(A,B,C) with {A,B}→C, C→A
  - 3NF but not BCNF

*Guideline: Aim for BCNF and settle for 3NF*

# Decomposing R into 3NF

- The algorithm is complicated

- 1. Get a "minimal cover" of FDs

- 2. Find a lossless-join decomposition of R (which might miss dependencies)

- 3. Add additional relations to the decomposition to cover any missing FDs of the cover

- Result will be lossless, will be dependency-preserving 3NF; might not be BCNF


- This way equivalent to textbook, but easier to follow.
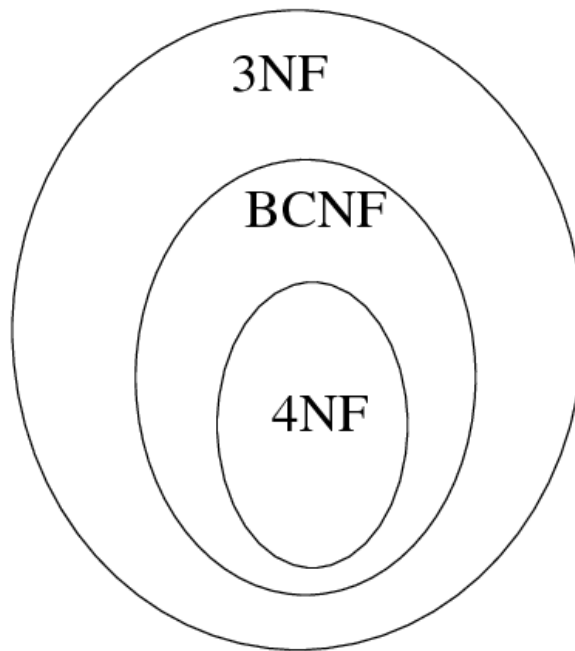
- → *Example 3.27 in textbook.*

# Fact of life...

*Finding a decomposition which is both lossless and dependency-preserving is not always possible.*

# Multi-valued Dependencies and 4NF

we will not cover this.

# Confused by Normal Forms ?



In practice: (1) 3NF is enough, (2) don't overdo it !

# Normalization Summary

- 1NF: usually part of the woodwork
- 2NF: usually skipped
- 3NF: a biggie
  - always aim for this
- BCNF and 4NF: tradeoffs start here
  - in re: d-preserving and losslessness
- 5NF: You can say you've heard of it...

# Caveat

- Normalization is not the be-all and end-all of DB design

- Example: suppose attributes A and B are always used together, but normalization theory says they should be in different tables.

  - decomposition might produce unacceptable performance loss (extra disk reads)

- Plus -- there are constraints other than FDs and MVDs

# Current Trends

*Normalization ≠ Crucial*

- Object DBs and Object-Relational DB's
  - may permit complex attributes
  - 1st normal form unnecessary
- Data Warehouses
  - huge historical databases, seldom or never updated after creation
  - joins expensive or impractical
  - argues against normalization
- Everyday relational DBs
  - aim for BCNF, settle for 3NF

63