Schema Refinement: Dependencies and Normal Forms

Ashraf Aboulnaga

David R. Cheriton School of Computer Science University of Waterloo

CS 348
Introduction to Database Management
Winter 2013

	CS 348	Schema Refinement	Winter 2013	1 / 43
Notes				

Outline

1 Introduction

Design Principles
Problems due to Poor Designs

2 Functional Dependencies

Logical Implication of FDs Attribute Closure

3 Schema Decomposition

Lossless-Join Decompositions
Dependency Preservation

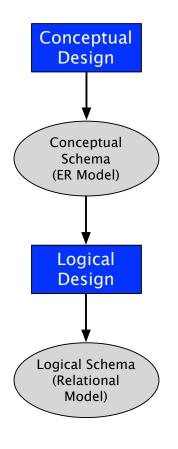
4 Normal Forms based on FDs

Boyce-Codd Normal Form Third Normal Form

	CS 348	Schema Refinement	Winter 2013	2 / 43
T - 4	05 040	Schema Remement	Williter 2013	4 / 43
lotes				

Design Process - Where are we?

CS 348



Step 1 - ER-to-relational mapping

Step 2 - Normalization: "Improving" the design

	CS 348	Schema Refinement	Winter 2013	3 / 43
Votes				

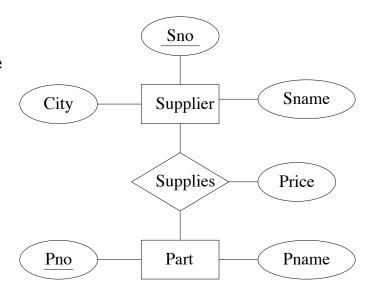
Relational Design Principles

- Relations should have semantic unity
- Information repetition should be avoided
 - Anomalies: insertion, deletion, modification
- Avoid null values as much as possible
 - Certainly avoid excessive null values
- Avoid spurious joins

Notes		CS 348	Schema Refinement	Winter 2013	4 / 43
	Notes				

A Parts/Suppliers Database Example

- Description of a parts/suppliers database:
 - Each type of part has a name and an identifying number, and may be supplied by zero or more suppliers. Each supplier may offer the part at a different price.
 - Each supplier has an identifying number, a name, and a contact location for ordering parts.



	CS 348	Schema Rennement	Winter 2013	5 / 43
Notes				

Parts/Suppliers Example (cont.)

Suppliers

<u>Sno</u>	Sname	City
S1	Magna	Ajax
S2	Budd	Hull

Parts

<u>Pno</u>	Pname
P1	Bolt
P2	Nut
P3	Screw

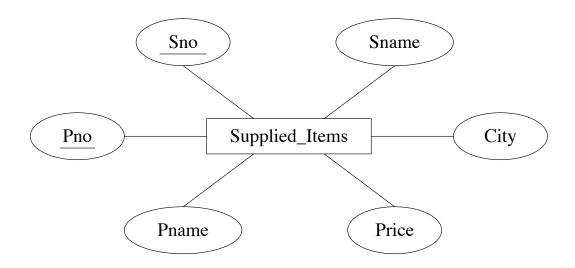
Supplies

<u>Sno</u>	<u>Pno</u>	Price
S1	P1	0.50
S1	P2	0.25
S1	P3	0.30
S2	P3	0.40

An instance of the parts/suppliers database.

	CS 348	Schema Refinement	Winter 2013	6 / 43
Notes				

Alternative Parts/Suppliers Database



An alternative E-R model for the parts/suppliers database.

	CS 348	Schema Refinement	Winter 2013	7 / 43
Notes				

Alternative Example (cont.)

Supplied Items

<u>Sno</u>	Sname	City	<u>Pno</u>	Pname	Price
S1	Magna	Ajax	P1	Bolt	0.50
S1	Magna	Ajax	P2	Nut	0.25
S1	Magna	Ajax	P3	Screw	0.30
S2	Budd	Hull	P3	Screw	0.40

A database instance corresponding to the alternative E-R model.

CS 348	Schema Refinement	Winter 2013	8 / 43
	CS 348	CS 348 Schema Refinement	CS 348 Schema Refinement Winter 2013

Change Anomalies

Consider

- Is one schema better than the other?
- What does it mean for a schema to be good?
- The single-table schema suffers from several kinds of problems:
 - Update problems (e.g. changing name of supplier)
 - Insert problems (e.g. add a new item)
 - Delete problems (e.g. Budd no longer supplies screws)
 - Likely increase in space requirements
- The multi-table schema does not have these problems.

	CS 348	Schema Refinement	Winter 2013	9 / 43
Notes	00010	Sonoma roomonono	Willvoi 2010	0 / 10
110162				

Another Alternative Parts/Supplier Database

Is more tables always better?

Snos Snames Cities City $\underline{\text{Sno}}$ <u>Sname</u> Magna Ajax S1 S2 Budd Hull Prices Inums Inames **Price** <u>Inum</u> <u>Iname</u> 0.50 Bolt I1 0.25 **I**2 Nut 0.30 I3 Screw 0.40

Information about relationships is lost!

	CS 348	Schema Refinement	Winter 2013	10 / 43
Notes				

Designing Good Databases

Goals

- A methodology for evaluating schemas (detecting anomalies).
- A methodology for transforming bad schemas into good schemas (repairing anomalies).
- How do we know an anomaly exists?
 - Certain types of *integrity constraints* reveal regularities in database instances that lead to anomalies.
- What should we do if an anomaly exists?
 - Certain schema decompositions can avoid anomalies while retaining all information in the instances

	CS 348	Schema Refinement	Winter 2013	11 / 43
Notes				

Functional Dependencies (FDs)

Idea: Express the fact that in a relation schema (values of) a set of attributes uniquely determine (values of) another set of attributes.

Definition (Functional Dependency)

Let R be a relation schema, and $X, Y \subseteq R$ sets of attributes. The functional dependency

$$X \rightarrow Y$$

holds on R if whenever an instance of R contains two tuples t and u such that t[X] = u[X] then it is also true that t[Y] = u[Y].

We say that X functionally determines Y (in R).

Notation: $t[A_1, \ldots, A_k]$ means projection of tuple t onto the attributes A_1, \ldots, A_k . In other words, $(t.A_1, \ldots, t.A_k)$.

	CS 348	Schema Refinement	Winter 2013	12 / 43
Notes				

Examples of Functional Dependencies

Consider the following relation schema:

EmpProj
SIN | PNum | Hours | EName | PName | PLoc | Allowance |

• SIN determines employee name

 $SIN \rightarrow EName$

- project number determines project name and location $PNum \rightarrow PName$, PLoc
- allowances are always the same for the same number of hours at the same location

PLoc, Hours \rightarrow Allowance

	CS 348	Schema Refinement	Winter 2013	13 / 43
No	es			

Functional Dependencies and Keys

- Keys (as defined previously):
 - A superkey is a set of attributes such that no two tuples (in an instance) agree on their values for those attributes.
 - A candidate key is a minimal superkey.
 - A primary key is a candidate key chosen by the DBA
- Relating keys and FDs:
 - If $K \subseteq R$ is a superkey for relation schema R, then dependency $K \to R$ holds on R.
 - If dependency $K \to R$ holds on R and we assume that R does not contain duplicate tuples (i.e. relational model) then $K \subseteq R$ is a superkey for relation schema R

	CS 348	Schema Refinement	Winter 2013	14 / 43
Notes				

Closure of FD Sets

How do we know what additional FDs hold in a schema?

- The closure of the set of functional dependencies F (denoted F^+) is the set of all functional dependencies that are satisfied by every relational instance that satisfies F.
- Informally, F^+ includes all of the dependencies in F, plus any dependencies they imply.

	CS 348	Schema Refinement	Winter 2013	15 / 43
Notes				,

Reasoning About FDs

Logical implications can be derived by using inference rules called Armstrong's axioms

- (reflexivity) $Y \subseteq X \Rightarrow X \rightarrow Y$
- (augmentation) $X o Y \Rightarrow XZ o YZ$
- (transitivity) X o Y, $Y o Z \Rightarrow X o Z$

The axioms are

- sound (anything derived from F is in F^+)
- complete (anything in F^+ can be derived)

Additional rules can be derived

- (union) $X \rightarrow Y$, $X \rightarrow Z \Rightarrow X \rightarrow YZ$
- (decomposition) $X o YZ \Rightarrow X o Y$

	CS 348	Schema Refinement	Winter 2013	16 / 43
Notes				,

Reasoning About FDs (example)

```
Example: F = \{ SIN, PNum \rightarrow Hours \}
                        SIN \rightarrow EName
                        PNum → PName, PLoc
                        PLoc, Hours \rightarrow Allowance }
```

A derivation of SIN, PNum \rightarrow Allowance:

- 1 SIN, PNum \rightarrow Hours $(\in F)$
- 2 PNum \rightarrow PName, PLoc ($\in F$)
- 3 PLoc, Hours \rightarrow Allowance $(\in F)$
- 4 SIN, PNum \rightarrow PNum (reflexivity)
- 5 SIN, PNum \rightarrow PName, PLoc (transitivity, 4 and 2)
- 6 SIN, PNum \rightarrow PLoc (decomposition, 5)
- \bigcirc SIN, PNum \rightarrow PLoc, Hours (union, 6, 1)
- 8 SIN, PNum \rightarrow Allowance (transitivity, 7 and 3)

	CS 348	Schema Refinement	Winter 2013	17 / 43
Notes				

Computing Attribute Closures

• There is a more efficient way of using Armstrong's axioms, if we only want to derive the maximal set of attributes functionally determined by some X (called the attribute closure of X).

```
function ComputeX^+(X,F)
begin X^+:=X;
while true do
if there exists (Y \to Z) \in F such that (1) \ Y \subseteq X^+, and (2) \ Z \not\subseteq X^+
then X^+:=X^+ \cup Z
else exit;
return X^+;
```

	CS 348	Schema Refinement	Winter 2013	18 / 43
Not	es			
	, 65			

Computing Attribute Closures (cont'd)

Let R be a relational schema and F a set of functional dependencies on R. Then

Theorem: X is a superkey of R if and only if

$$ComputeX^+(X,F)=R$$

Theorem: $X \to Y \in F^+$ if and only if

$$Y\subseteq ComputeX^+(X,F)$$

	CS 348	Schema Refinement	Winter 2013	19 / 43
Notes				
110000				

Attribute Closure Example

Example:
$$F = \{ SIN \rightarrow EName \\ PNum \rightarrow PName, PLoc \\ PLoc, Hours \rightarrow Allowance \}$$

 $ComputeX^+(\{Pnum,Hours\},F)$:

FD	X^+
initial	Pnum,Hours
Pnum→Pname,Ploc	Pnum,Hours,Pname,Ploc
$PLoc, Hours \rightarrow Allowance$	Pnum, Hours, Pname, Ploc, Allowance

	CS 348	Schema Refinement	Winter 2013	20 / 43
Notes				
				_

Schema Decomposition

Definition (Schema Decomposition)

Let R be a relation schema (= set of attributes). The collection $\{R_1, \ldots, R_n\}$ of relation schemas is a decomposition of R if

$$R = R_1 \cup R_2 \cup \cdots \cup R_n$$

A good decomposition does not

- lose information
- complicate checking of constraints
- contain anomalies (or at least contains fewer anomalies)

	CS 348	Schema Refinement	Winter 2013	21 / 43
Notes				·
110000				

Lossless-Join Decompositions

We should be able to construct the instance of the original table from the instances of the tables in the decomposition

Example: Consider replacing

Marks

Student	Assignment	Group	Mark
Ann	A1	G1	80
Ann	A2	G3	60
Bob	A1	G2	60

by decomposing (i.e. projecting) into two tables

SGM

Student	Group	<u>Mark</u>
Ann	G1	80
Ann	G3	60
Bob	G2	60

AM

Assignment	<u>Mark</u>
A1	80
A2	60
A1	60

Notes

CS 348 Schema Refinement Winter 2013 22 / 43

Notes

Lossless-Join Decompositions (cont.)

But computing the natural join of SGM and AM produces

Student	Assignment	Group	Mark
Ann	A1	G1	80
Ann	A2	G3	60
Ann	A1	G3	60
Bob	A2	G2	60
Bob	A1	G2	60

... and we get extra data (spurious tuples). We would therefore lose information if we were to replace Marks by SGM and AM.

If re-joining SGM and AM would always produce exactly the tuples in Marks, then we call SGM and AM a lossless-join decomposition.

CS 348	Schema Refinement	Winter 2013	23 / 43
Notes			

Lossless-Join Decompositions (cont.)

A decomposition $\{R_1, R_2\}$ of R is lossless if and only if the common attributes of R_1 and R_2 form a superkey for either schema, that is

$$R_1 \cap R_2 o R_1$$
 or $R_1 \cap R_2 o R_2$

Example: In the previous example we had

```
R = \{Student, Assignment, Group, Mark\}, F = \{(Student, Assignment \rightarrow Group, Mark)\}, R_1 = \{Student, Group, Mark\}, R_2 = \{Assignment, Mark\}
```

Decomposition $\{R_1, R_2\}$ is lossy because $R_1 \cap R_2$ (= $\{Mark\}$) is not a superkey of either $\{Student, Group, Mark\}$ or $\{Assignment, Mark\}$

	CS 348	Schema Refinement	Winter 2013	24 / 43
Notes				

Dependency Preservation

How do we test/enforce constraints on the decomposed schema?

Example: A table for a company database could be

and two decompositions

$$D_1 = \{R1[Proj, Dept], R2[Dept, Div]\}$$

 $D_2 = \{R1[Proj, Dept], R3[Proj, Div]\}$

Both are lossless. (Why?)

	CS 348	Schema Refinement	Winter 2013	25 / 43
Νo	tes			

Dependency Preservation (cont.)

Which decomposition is better?

- Decomposition D_1 lets us test FD1 on table R1 and FD2 on table R2; if they are both satisfied, FD3 is automatically satisfied.
- In decomposition D_2 we can test FD1 on table R1 and FD3 on table R3. Dependency FD2 is an interrelational constraint: testing it requires joining tables R1 and R3.

 $\Rightarrow D_1$ is better!

Given a schema R and a set of functional dependencies F, decomposition $D = \{R_1, \ldots, R_n\}$ of R is dependency preserving if there is an equivalent set of functional dependencies F', none of which is interrelational in D.

	CS 348	Schema Refinement	Winter 2013	26 / 43
Notes				
_				

Normal Forms

What is a "good" relational database schema?

Rule of thumb: Independent facts in separate tables:

"Each relation schema should consist of a primary key and a set of mutually independent attributes"

This is achieved by transforming a schema into a normal form.

Goals:

- Intuitive and straightforward transformation
- Anomaly-free/Nonredundant representation of data

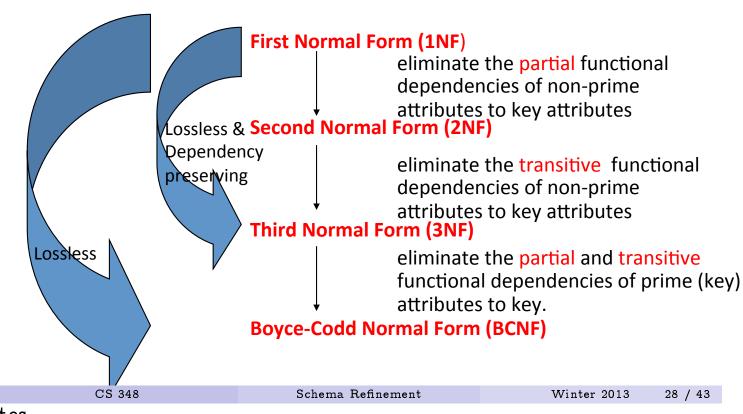
Normal Forms based on Functional Dependencies:

- Boyce-Codd Normal Form (BCNF)
- Third Normal Form (3NF)

	CS 348	Schema Refinement	Winter 2013	27 / 43
Νo	tes			

Normal Forms Based on FDs

1NF eliminates relations within relations or relations as attributes of tuples



	CS 348	Schema Refinement	Winter 2013	28 / 43
Notes				
110000				
-				

Boyce-Codd Normal Form (BCNF) - Informal

- BCNF formalizes the goal that in a good database schema, independent relationships are stored in separate tables.
- Given a database schema and a set of functional dependencies for the attributes in the schema, we can determine whether the schema is in BCNF. A database schema is in BCNF if each of its relation schemas is in BCNF.
- Informally, a relation schema is in BCNF if and only if any group of its attributes that functionally determines *any* others of its attributes functionally determines *all* others, i.e., that group of attributes is a superkey of the relation.

	CS 348	Schema Refinement	Winter 2013	29 / 43
Notes				

Formal Definition of BCNF

Let R be a relation schema and F a set of functional dependencies.

Schema R is in **BCNF** (w.r.t. F) if and only if whenever $(X \to Y) \in F^+$ and $XY \subseteq R$, then either

- (X o Y) is trivial (i.e., $Y \subseteq X$), or
- X is a superkey of R

A database schema $\{R_1, \ldots, R_n\}$ is in BCNF if each relation schema R_i is in BCNF.

	CS 348	Schema Refinement	Winter 2013	30 / 43
Notes				

BCNF and Redundancy

• Why does BCNF avoid redundancy? Consider:

Supplied_Items
Sno | Sname | City | Pno | Pname | Price

• The following functional dependency holds:

Sno \rightarrow Sname, City

- Therefore, supplier name "Magna" and city "Ajax" must be repeated for each item supplied by supplier S1.
- Assume the above FD holds over a schema R that is in BCNF. This implies that:
 - Sno is a superkey for R
 - each Sno value appears on one row only
 - no need to repeat Sname and City values

	CS 348	Schema Refinement	Winter 2013	31 / 43
NTakas	C5 346	Schema Itemnement	Winter 2013	31 / 43
Notes				
				·

Lossless-Join BCNF Decomposition

```
function DecomposeBCNF(R, F)
begin
Result := \{R\};
while some \ R_i \in Result \ and \ (X 	o Y) \in F^+
violate \ the \ BCNF \ condition \ do \ begin
Replace \ R_i \ by \ R_i - (Y - X);
Add \ \{X, Y\} \ to \ Result;
end;
return Result;
```

	CS 348	Schema Refinement	Winter 2013 32 / 43
Not	es		

Lossless-Join BCNF Decomposition

- No efficient procedure to do this exists.
- Results depend on sequence of FDs used to decompose the relations.
- It is possible that no lossless join dependency preserving BCNF decomposition exists
 - Consider $R = \{A, B, C\}$ and $F = \{AB \rightarrow C, C \rightarrow B\}$.

	CS 348	Schema Refinement	Winter 2013	33 / 43
Totes				

BCNF Decomposition - An Example

- $R = \{Sno, Sname, City, Pno, Pname, Price\}$
- functional dependencies:

 ${\tt Sno} \, \rightarrow \, {\tt Sname,City}$

 $\mathtt{Pno} \to \mathtt{Pname}$

 ${\tt Sno,Pno} \to {\tt Price}$

• This schema is not in BCNF because, for example, Sno determines Sname and City, but is not a superkey of R.

	CS 348	Schema Refinement	Winter 2013	34 / 43
Notes				

BCNF Decomposition - An Example (cont.)

Decomposition Diagram:

{Sno, Sname, City, Pno, Pname, Price}

Sno -> Sname, City

{Sno, Pno, Pname, Price}

Pno -> Pname

{Sno, Pno, Price}

{Pno, Pname}

• The complete schema is now

$$R_1 = \{ ext{Sno,Sname,City}\}$$

 $R_2 = \{ ext{Sno,Pno,Price}\}$
 $R_3 = \{ ext{Pno,Pname}\}$

 This schema is a lossless-join, BCNF decomposition of the original schema R.

	CS 348	Schema Refinement	Winter 2013	35 / 43
Notes				

Third Normal Form (3NF)

Schema R is in $\mathbf{3NF}$ (w.r.t. F) if and only if whenever $(X \to Y) \in F^+$ and $XY \subseteq R$, then either

- $(X \to Y)$ is trivial, or
- X is a superkey of R, or
- each attribute in Y X is contained in a candidate key of R

A database schema $\{R_1, \ldots, R_n\}$ is in 3NF if each relation schema R_i is in 3NF.

- 3NF is looser than BCNF
 - allows more redundancy
 - e.g. $R = \{A, B, C\}$ and $F = \{AB \rightarrow C, C \rightarrow B\}$.
- lossless-join, dependency-preserving decomposition into 3NF relation schemas always exists.

	CS 348	Schema Refinement	Winter 2013	36 / 43
Notes				

Minimal Cover

Definition: Two sets of dependencies F and G are equivalent iff $F^+ = G^+$.

There are different sets of functional dependencies that have the same logical implications. Simple sets are desirable.

Definition: A set of dependencies G is minimal if

- \bigcirc every right-hand side of an dependency in F is a single attribute.
- 2 for no $X \to A$ is the set $F \{X \to A\}$ equivalent to F.
- 3 for no $X \to A$ and Z a proper subset of X is the set $F \{X \to A\} \cup \{Z \to A\}$ equivalent to F.

Theorem: For every set of dependencies F there is an equivalent minimal set of dependencies (minimal cover).

	CS 348	Schema Refinement	Winter 2013	37 / 43
Notes				

Finding Minimal Covers

A minimal cover for F can be computed in three steps. Note that each step must be repeated until it no longer succeeds in updating F.

Step 1.

Replace $X \to YZ$ with the pair $X \to Y$ and $X \to Z$.

Step 2.

Remove A from the left-hand-side of $X \to B$ in F if B is in $ComputeX^+(X - \{A\}, F)$.

Step 3.

Remove $X \to A$ from F if $A \in ComputeX^+(X, F - \{X \to A\})$.

	CS 348	Schema Refinement	Winter 2013	38 / 43
Notes				

Dependency-Preserving 3NF Decomposition

Idea: Decompose into 3NF relations and then "repair"

```
function Decompose3NF(R,F)

begin
Result := \{R\};
while some\ R_i \in Result\ and\ (X \to Y) \in F^+
violate\ the\ 3NF\ condition\ do\ begin
Replace\ R_i\ by\ R_i - (Y - X);
Add\ \{X,Y\}\ to\ Result;
end;
N := (a\ minimal\ cover\ for\ F) - (\bigcup_i F_i)^+
for each (X \to Y) \in N do
Add\ \{X,Y\}\ to\ Result;
end;
return\ Result;
```

	CS 348	Schema Refinement	Winter 2013 39 / 43
Note	es		

Dep-Preserving 3NF Decomposition - An Example

- $R = \{Sno, Sname, City, Pno, Pname, Price\}$
- Functional dependencies:

 $Sno \rightarrow Sname, City$

 $\mathtt{Pno} o \mathtt{Pname}$

 $\mathtt{Sno},\mathtt{Pno} o \mathtt{Price}$

Sno, Pname \rightarrow Price

• Following same decomposition tree as BCNF example:

 $R_1 = \{ ext{Sno,Sname,City}\}$ $R_2 = \{ ext{Sno,Pno,Price}\}$ $R_3 = \{ ext{Pno,Pname}\}$

• Minimal cover:

 $\mathtt{Sno} o \mathtt{Sname}$

 $\mathtt{Pno} o \mathtt{Pname}$

 $Sno \rightarrow City$

Sno, Pname ightarrow Price

Add relation to preserve missing dependency

$$R_4 = \{ \text{Sno, Pname, Price} \}$$

	CS 348	Schema Refinement	Winter 2013	40 / 43
Notes				

3NF Synthesis

A lossless-join 3NF decomposition that is dependency preserving can be efficiently computed

```
function Synthesize3NF(R,F)
begin

Result := \emptyset;

F' := a \ minimal \ cover \ for \ F;

for each (X \to Y) \in F' do

Result := Result \cup \{XY\};

if there is no R_i \in Result \ such \ that

R_i \ contains \ a \ candidate \ key \ for \ R then begin

compute \ a \ candidate \ key \ K \ for \ R;

Result := Result \cup \{K\};

end;

return Result;
```

CS 348 Schema Refinement Winter 2013 41 / 43
Notes

3NF Synthesis - An Example

- $R = \{Sno, Sname, City, Pno, Pname, Price\}$
- Functional dependencies:

 $Sno \rightarrow Sname, City$ $Pno \rightarrow Pname$ $Sno, Pno \rightarrow Price$ $Sno, Pname \rightarrow Price$

• Minimal cover:

 $egin{array}{lll} {
m Sno} &
ightarrow {
m Sname} & R_1 = \{ {
m Sno, Sname} \} \ {
m Sno} &
ightarrow {
m City} & R_2 = \{ {
m Sno, City} \} \ {
m Pno} &
ightarrow {
m Pname} & R_3 = \{ {
m Pno, Pname} \} \ {
m Sno, Pname} &
ightarrow {
m R4} = \{ {
m Sno, Pname, Price} \} \end{array}$

- Add relation for candidate key $R_5 = \{$ Sno, Pno $\}$
- Optimization: combine relations R_1 and R_2 (same key)

	CC 040		TTT: 1 0010	10 / 10
_	CS 348	Schema Refinement	Winter 2013	42 / 43
Jotes				

Summary

- Functional dependencies provide clues towards elimination of (some) redundancies in a relational schema.
- Goals: to decompose relational schemas in such a way that the decomposition is
 - (1) lossless-join
 - (2) dependency preserving
 - (3) BCNF (and if we fail here, at least 3NF)

	CS 348	Schema Refinement	Winter 2013	43 / 43
Notes				