

# Schema design and dependencies

Remy Wang, 4/22/25

<b>name</b>	<b>location</b>	<b>breed</b>	<b>kind</b>	<b>job</b>
casa	LA	tabby	cat	<b>NULL</b>
kira	seattle	tuxedo	cat	<b>NULL</b>
remy	LA	<b>NULL</b>	<b>NULL</b>	prof
vincent	LA	<b>NULL</b>	<b>NULL</b>	TA

<b>name</b>	<b>location</b>	<b>breed</b>	<b>kind</b>	<b>job</b>
casa	LA	tabby	cat	<b>NULL</b>
kira	seattle	tuxedo	cat	<b>NULL</b>
remy	LA	<b>NULL</b>	<b>NULL</b>	prof
vincent	LA	<b>NULL</b>	<b>NULL</b>	TA

1 table stores 1 kind of data

<b>name</b>	<b>location</b>	<b>salary</b>	<b>graduate</b>	<b>funding</b>
remy	LA	\$30	<b>NULL</b>	\$10
vincent	LA	\$20	2025	<b>NULL</b>

<b>name</b>	<b>location</b>	<b>salary</b>	<b>graduate</b>	<b>funding</b>
remy	LA	\$30	<b>NULL</b>	\$10
vincent	LA	\$20	2025	<b>NULL</b>

<b>name</b>	<b>location</b>	<b>salary</b>	<b>funding</b>
remy	LA	\$30	\$10

<b>name</b>	<b>location</b>	<b>salary</b>	<b>graduate</b>
vincent	LA	\$20	2025

<b>name</b>	<b>location</b>	<b>salary</b>	<b>funding</b>
remy	LA	\$30	\$10

<b>name</b>	<b>location</b>	<b>salary</b>	<b>graduate</b>
vincent	LA	\$20	2025

payroll

<b>name</b>	<b>location</b>	<b>salary</b>
remy	LA	\$30
vincent	LA	\$20

prof

<b>name</b>	<b>funding</b>
remy	\$10

student

<b>name</b>	<b>graduate</b>
vincent	2025

each kind of data has its own table

<b>name</b>	<b>location</b>	<b>salary</b>	<b>course</b>
remy	LA	\$30	143
remy	LA	\$30	240
remy	LA	\$30	249

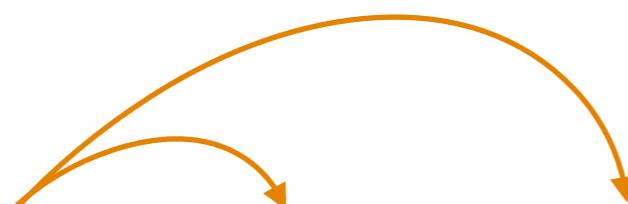
<b>name</b>	<b>location</b>	<b>salary</b>	<b>course</b>
remy	LA	\$30	143
remy	LA	\$30	240
remy	LA	\$30	249

redundancy!

(complicates updates & deletes)

each piece of information stored once

determines



<b>name</b>	<b>location</b>	<b>salary</b>	<b>course</b>
remy	LA	\$30	143
remy	LA	\$30	240
remy	LA	\$30	249

$\text{name} \rightarrow \text{location, salary}$

# functional dependency

$$X \rightarrow Y$$

the values of X *uniquely determines* Y

$$\forall t, t' \in R : \pi_X(t) = \pi_X(t') \implies \pi_Y(t) = \pi_Y(t')$$

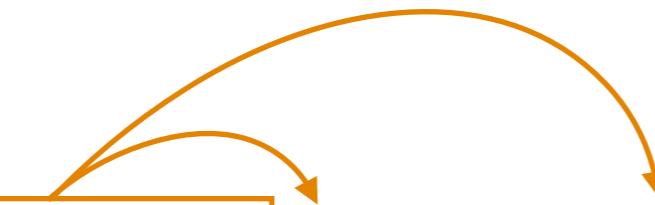
<b>name</b>	<b>location</b>	<b>salary</b>	<b>course</b>
remy	LA	\$30	143
remy	LA	\$30	240
remy	LA	\$30	249
dan	seattle	\$50	344
dan	seattle	\$50	444



$$\forall t, t' \in R : \pi_{\text{name}}(t) = \pi_{\text{name}}(t') \implies \pi_{\text{salary}}(t) = \pi_{\text{salary}}(t')$$

name → salary

determines



<b>first n.</b>	<b>last n.</b>	<b>location</b>	<b>salary</b>	<b>course</b>
remy	w	LA	\$30	143
remy	w	LA	\$30	240
remy	w	LA	\$30	249
dan	s	seattle	\$50	344
dan	s	seattle	\$50	444
dan	o	zurich	\$50	101
dan	o	zurich	\$50	113

determines

<b>first n.</b>	<b>last n.</b>	<b>location</b>	<b>salary</b>	<b>course</b>
remy	w	LA	\$30	143
remy	w	LA	\$30	240
remy	w	LA	\$30	249
dan	s	seattle	\$50	344
dan	s	seattle	\$50	444
dan	o	zurich	\$50	101
dan	o	zurich	\$50	113

$$\{\text{first n., last n.}\} \rightarrow \{\text{location, salary}\}$$

Check  $X \rightarrow Y$  using SQL?

Check  $X \rightarrow Y$  using SQL?

```
SELECT * FROM R
GROUP BY X
HAVING COUNT(Y) > 1
```

# Trivial FDs?

A	B	C	D	E

# Trivial FDs?

A	B	C	D	E

$A \rightarrow A, B \rightarrow B, \dots$

$AB \rightarrow A, AB \rightarrow B, \dots$

# Trivial FDs?

A	B	C	D	E

$$A \rightarrow A, B \rightarrow B, \dots$$

$$AB \rightarrow A, AB \rightarrow B, \dots$$

$$X \supseteq Y \implies X \rightarrow Y$$

<b>name</b>	<b>job</b>	<b>location</b>	<b>salary</b>	<b>tax %</b>
remy	prof	LA	\$30	20
dan	prof	seattle	\$50	15
vincent	TA	LA	\$20	10

$\text{job} \rightarrow \text{salary}$

$\text{name} \rightarrow \text{location}$

$\text{location}, \text{salary} \rightarrow \text{tax \%}$

<b>name</b>	<b>job</b>	<b>location</b>	<b>salary</b>	<b>tax %</b>
remy	prof	LA	\$30	20
dan	prof	seattle	\$50	15
vincent	TA	LA	\$20	10

$\text{job} \rightarrow \text{salary}$

$\text{name} \rightarrow \text{location} \quad \Rightarrow \quad \text{name, job} \rightarrow ?$

$\text{location, salary} \rightarrow \text{tax \%}$

<b>name</b>	<b>job</b>	<b>location</b>	<b>salary</b>	<b>tax %</b>
remy	prof	LA	\$30	20
dan	prof	seattle	\$50	15
vincent	TA	LA	\$20	10

$\text{job} \rightarrow \text{salary}$

$\text{name} \rightarrow \text{location} \quad \Rightarrow \quad \text{name, job} \rightarrow \text{tax \%}$

$\text{location, salary} \rightarrow \text{tax \%}$

<b>name</b>	<b>job</b>	<b>location</b>	<b>salary</b>	<b>tax %</b>
N	J	L	S	T
N	J	?	?	?
vincent	TA	LA	\$20	10

$\text{job} \rightarrow \text{salary}$

$\text{name} \rightarrow \text{location} \quad \Rightarrow \quad \text{name, job} \rightarrow \text{tax \%}$

$\text{location, salary} \rightarrow \text{tax \%}$

# FD inference

Given

$$X_1 \rightarrow Y_1$$

$$X_2 \rightarrow Y_2$$

$$X_3 \rightarrow Y_3$$

...

Does this hold?

$$X \rightarrow Y$$

# Armstrong's axioms

$$Y \subseteq X \implies X \rightarrow Y$$

$$X \rightarrow Y \implies XZ \rightarrow YZ$$

$$X \rightarrow Y \wedge Y \rightarrow Z \implies X \rightarrow Z$$

# Armstrong's axioms

X	Y	Z
x	y	z
x	?	z

$$Y \subseteq X \implies X \rightarrow Y$$

Reflexivity

$$X \rightarrow Y \implies XZ \rightarrow YZ$$

Augmentation

$$X \rightarrow Y \wedge Y \rightarrow Z \implies X \rightarrow Z$$

Transitivity

job → salary

name → location

name, job → tax %

location, salary → tax %

job, name → salary, name

name, salary → location, salary

name, salary → tax %

# FD closure

given

$$X_1 \rightarrow Y_1$$

$$X_2 \rightarrow Y_2$$

$$X_3 \rightarrow Y_3$$

...

compute  $X^+$  :

start w/  $X$

repeat until no change:

add  $Y_i$  to  $X$  if  $X_i \subseteq X$

job → salary

name → location       $\implies$       name, job → tax %

location, salary → tax %

$\{ \text{name, job} \}^+ = \{ \text{name, job} \}$

Why FDs?

# FDs cause anomalies



first n.	last n.	location	salary	course
remy	w	LA	\$30	143
remy	w	LA	\$30	240
remy	w	LA	\$30	249

$X \subseteq \{A_1, \dots, A_5\}$  is a **superkey** if  $\forall i : X \rightarrow A_i$

i.e.  $X^+ = \{A_1, \dots, A_5\}$

A1	A2	A3	A4	A5

PK

<b>name</b>	<b>location</b>	<b>salary</b>
remy	LA	\$30
vincent	LA	\$20

A diagram illustrating a query operation on a table. The table has four columns: **first n.**, **last n.**, **location**, and **salary**. The first two columns are highlighted with an orange border. Two orange curved arrows point from the **last n.** column towards the **location** and **salary** columns, indicating a projection or selection operation.

<b>first n.</b>	<b>last n.</b>	<b>location</b>	<b>salary</b>
remy	w	LA	\$30
dan	s	seattle	\$50
dan	o	zurich	\$50

A diagram illustrating a relationship between the first and last names of individuals and their location. A curved orange arrow originates from the 'first n.' and 'last n.' columns and points to the 'location' column.

<b>first n.</b>	<b>last n.</b>	<b>location</b>	<b>salary</b>	<b>course</b>
remy	w	LA	\$30	143
remy	w	LA	\$30	240
remy	w	LA	\$30	249
dan	s	seattle	\$50	344
dan	s	seattle	\$50	444
dan	o	zurich	\$50	101
dan	o	zurich	\$50	113

<b>name</b>	<b>job</b>	<b>location</b>	<b>salary</b>	<b>tax %</b>
remy	prof	LA	\$30	20
dan	prof	seattle	\$50	15
vincent	TA	LA	\$20	10

$\text{job} \rightarrow \text{salary}$

$\text{name} \rightarrow \text{location} \quad \Rightarrow \quad \text{name, job} \rightarrow \text{tax \%}$

$\text{location, salary} \rightarrow \text{tax \%}$

$X \subseteq \{A_1, \dots, A_5\}$  is a **superkey** if  $\forall i : X \rightarrow A_i$

A **key** is a minimal superkey

(no longer a superkey if removing anything)

revisit examples

**to find a key:** guess  $X$  from small to large, check if  $X^+ = \{A_1, \dots, A_5\}$

$$\text{job}^+ = \{?\} \quad \text{salary}^+ = \{?\} \quad \text{name}^+ = \{?\}$$

$$\text{location}^+ = \{?\} \quad \text{tax}^+ = \{?\}$$

$$\{\text{name}, \text{job}\}^+ = \{?\}$$

$\text{job} \rightarrow \text{salary}$

$\text{name} \rightarrow \text{location} \implies \text{name, job} \rightarrow \text{tax \%}$

$\text{location, salary} \rightarrow \text{tax \%}$

**to find a key:** guess  $X$  from small to large, check if  $X^+ = \{A_1, \dots, A_5\}$

any key must contain **name, job**

job  $\rightarrow$  salary

name  $\rightarrow$  location  $\implies$  name, job  $\rightarrow$  tax %

location, salary  $\rightarrow$  tax %

more exercises

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

$$A \rightarrow BC$$

$$B \rightarrow AC$$

$$AB \rightarrow C$$

Decompose relations, finally!

Boyce-Codd Normal Form:

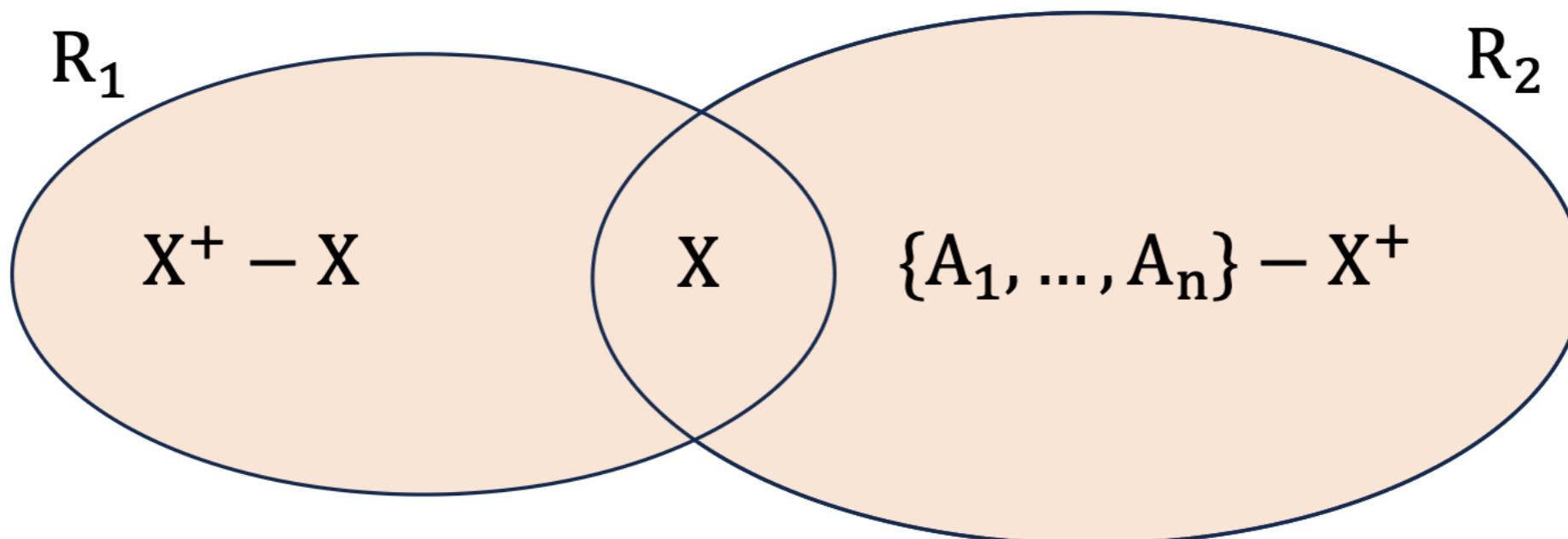
$\forall X \rightarrow Y \models R : Y \subseteq X \vee X \text{ is a superkey}$

i.e.  $\forall X : X^+ = X \vee X^+ = \{A_1, \dots\}$

# Decomposition

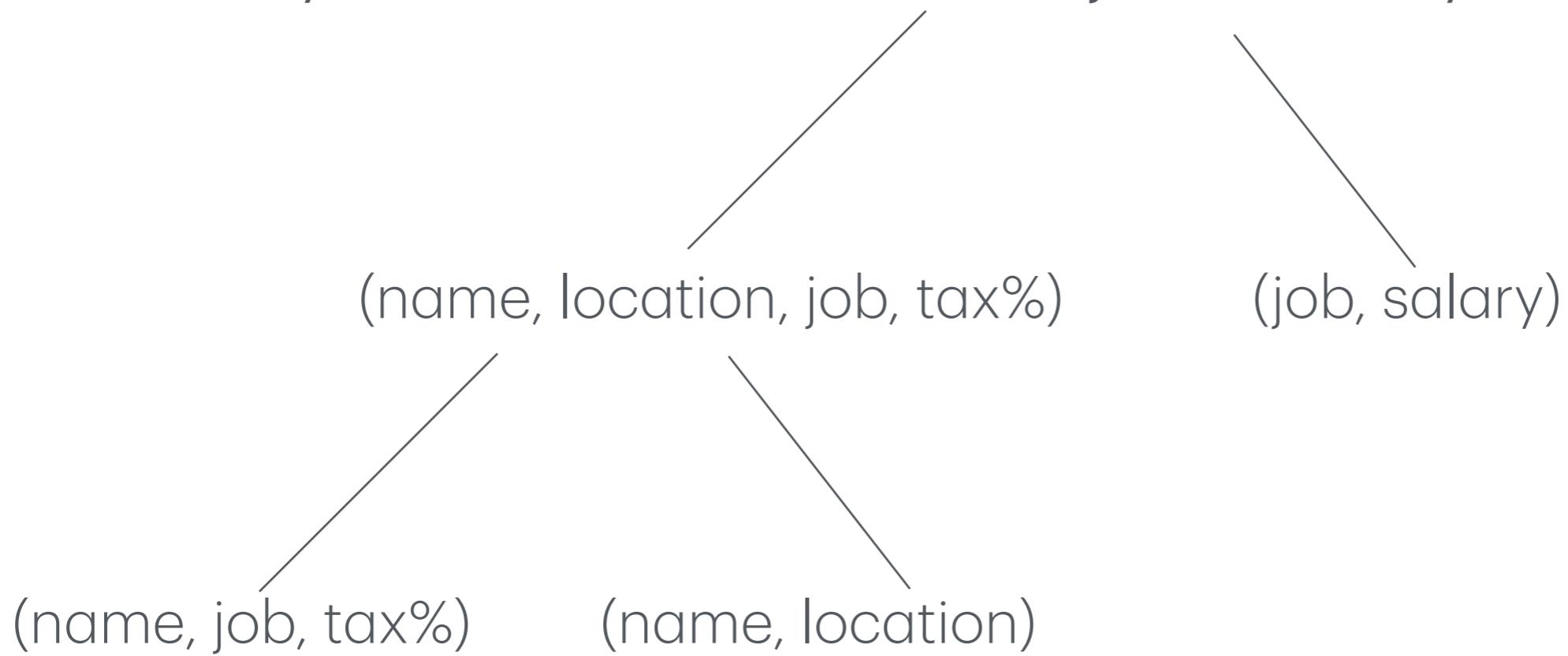
repeat

find  $X \rightarrow Y$  violating BCNF  
"factor out"  $X^+$  from R



$\text{job} \rightarrow \text{salary}$     $\text{name} \rightarrow \text{location}$     $\text{location, salary} \rightarrow \text{tax \%}$

Payroll(name, location, job, salary, tax%)



$\text{job} \rightarrow \text{salary}$     $\text{name} \rightarrow \text{location}$     $\text{location, salary} \rightarrow \text{tax \%}$

Payroll(name, location, job, salary, tax%)



Why care about preserving FDs?

Find tax rate for ("LA", \$50)

(name, job) (job, salary) (name, location) (location, salary, tax%)

(name, job, tax%) (name, location) (job, salary)

# Other considerations

security & privacy

compliance

geolocation

performance