

---

# Relational Algebra 2

Week 5

---

# Relational Algebra (So far)

- Basic operations:
  - Selection ( $\sigma$ ) Selects a subset of rows from relation.
  - Projection ( $\pi$ ) Deletes unwanted columns from relation.
  - Cross-product ( $\times$ ) Allows us to combine two relations.
  - Set-difference ( $-$ ) Tuples in reln. 1, but not in reln. 2.
  - Union ( $\cup$ ) Tuples in reln. 1 and tuples in reln. 2.

Also,

- Rename ( $\rho$ ) Changes names of the attributes
  - Intersection ( $\cap$ ) Tuples in both reln. 1 and in reln. 2.
- Since each operation returns a relation, *operations can be composed!* (Algebra is “closed”.)
  - Use of temporary relations recommended.

# Additional Operations

We define additional operations that do not add any power to the relational algebra, but that simplify common queries.

- Natural join
- Conditional Join
- Equi-Join
- Division

All joins are really special cases of conditional join

Also, we've already seen “Set intersection”:

$$r \cap s = r - (r - s)$$

# Quick note on notation

*good\_customers*

<i>customer-name</i>	<i>loan-number</i>
Patty	1234
Apu	3421
Selma	2342
Ned	4531

*bad\_customers*

<i>customer-name</i>	<i>loan-number</i>
Seymour	3432
Marge	3467
Selma	7625
Abraham	3597

If we have two or more relations which feature the same attribute names, we could confuse them. To prevent this we can use dot notation.

For example

*good\_customers.loan-number*

# Natural-Join Operation: Motivation

Very often, we have a query and the answer is not contained in a single relation. For example, I might wish to know where Apu banks.

The classic relational algebra way to do such queries is a cross product, followed by a selection which tests for equality on some pair of fields.

$$\sigma_{\text{borrower.l-number} = \text{loan.l-number}}(\text{borrower} \times \text{loan}))$$

While this works...

- it is unintuitive
- it requires a lot of memory
- the notation is cumbersome

<i>borrower</i>		<i>loan</i>	
<i>cust-name</i>	<i>l-number</i>	<i>l-number</i>	<i>branch</i>
Patty	1234	1234	Dublin
Apu	3421	3421	Irvine

<i>cust-name</i>	<i>borrower.l-number</i>	<i>loan.l-number</i>	<i>branch</i>
Patty	1234	1234	Dublin
Patty	1234	3421	Irvine
Apu	3421	1234	Dublin
Apu	3421	3421	Irvine

<i>cust-name</i>	<i>borrower.l-number</i>	<i>loan.l-number</i>	<i>branch</i>
Patty	1234	1234	Dublin
Apu	3421	3421	Irvine

Note that in this example the two relations are the same size (2 by 2), this does not have to be the case.

So, we have a more intuitive way of achieving the same effect, the natural join, denoted by the  $\bowtie$  symbol

# Natural-Join Operation: Intuition

Natural join combines a cross product and a selection into one operation. It performs a selection forcing equality on *those attributes that appear in both relation schemes*. Duplicates are removed as in all relation operations.

So, if the relations have one attribute in common, as in the last slide (“*l-number*”), for example, we have...

$$\text{borrower} \bowtie \text{loan} = \sigma_{\text{borrower.l-number} = \text{loan.l-number}}(\text{borrower} \times \text{loan}))$$

There are two special cases:

- If the two relations have no attributes in common, then their natural join is simply their cross product.
- If the two relations have more than one attribute in common, then the natural join selects only the rows where all pairs of matching attributes match. (let's see an example on the next slide).

**A**

<i>l-name</i>	<i>f-name</i>	<i>age</i>
Bouvier	Selma	40
Bouvier	Patty	40
Smith	Maggie	2

**B**

<i>l-name</i>	<i>f-name</i>	<i>ID</i>
Bouvier	Selma	1232
Smith	Selma	4423

Both the *l-name* and the *f-name* match, so select.

Only the *f-names* match, so don't select.

Only the *l-names* match, so don't select.

<i>l-name</i>	<i>f-name</i>	<i>age</i>	<i>l-name</i>	<i>f-name</i>	<i>ID</i>
<b>Bouvier</b>	<b>Selma</b>	40	<b>Bouvier</b>	<b>Selma</b>	1232
Bouvier	<b>Selma</b>	40	Smith	<b>Selma</b>	4423
<b>Bouvier</b>	Patty	2	<b>Bouvier</b>	Selma	1232
Bouvier	Patty	40	Smith	Selma	4423
Smith	Maggie	2	Bouvier	Selma	1232
<b>Smith</b>	Maggie	2	<b>Smith</b>	Selma	4423

We remove duplicate attributes...

<i>l-name</i>	<i>f-name</i>	<i>age</i>	<i>l-name</i>	<i>f-name</i>	<i>ID</i>
Bouvier	Selma	40	Bouvier	Selma	1232

The natural join of *A* and *B*

Note that this is just a way to visualize the natural join, we don't really have to do the cross product as in this example

$$A \bowtie B =$$

<i>l-name</i>	<i>f-name</i>	<i>age</i>	<i>ID</i>
Bouvier	Selma	40	1232

# Natural-Join Operation

- Notation:  $r \bowtie s$
- Let  $r$  and  $s$  be relation instances on schemas  $R$  and  $S$  respectively. The result is a relation on schema  $R \cup S$  which is obtained by considering each pair of tuples  $t_r$  from  $r$  and  $t_s$  from  $s$ .
- If  $t_r$  and  $t_s$  have the same value on each of the attributes in  $R \cap S$ , a tuple  $t$  is added to the result, where
  - $t$  has the same value as  $t_r$  on  $r$
  - $t$  has the same value as  $t_s$  on  $s$

- Example:

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

$$S = (E, B, D)$$

- Result schema =  $(A, B, C, D, E)$
- $r \bowtie s$  is defined as:

$$\pi_{r.A, r.B, r.C, r.D, s.E} (\sigma_{r.B = s.B \wedge r.D = s.D} (r \times s))$$



# Natural Join Operation – Example

- Relation instances  $r, s$ :

A	B	C	D
$\alpha$	1	$\alpha$	a
$\beta$	2	$\gamma$	a
$\gamma$	4	$\beta$	b
$\alpha$	1	$\gamma$	a
$\delta$	2	$\beta$	b

$r$

B	D	E
1	a	$\alpha$
3	a	$\beta$
1	a	$\gamma$
2	b	$\delta$
3	b	$\epsilon$

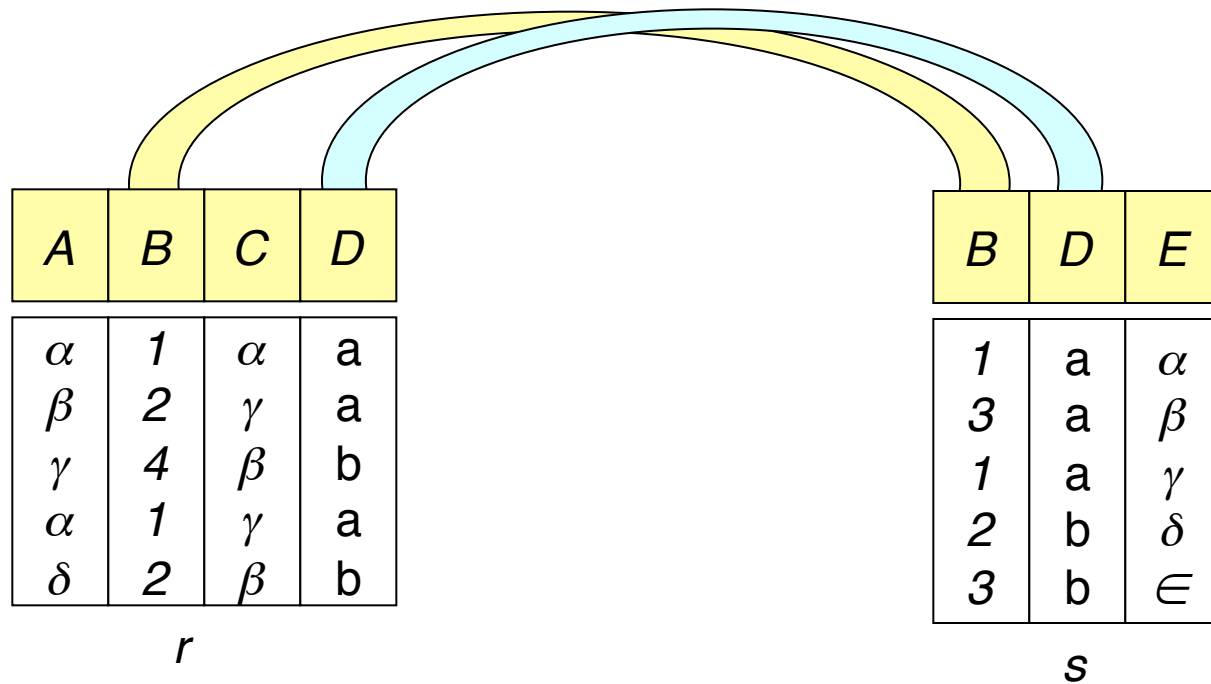
$s$

$r \bowtie s$

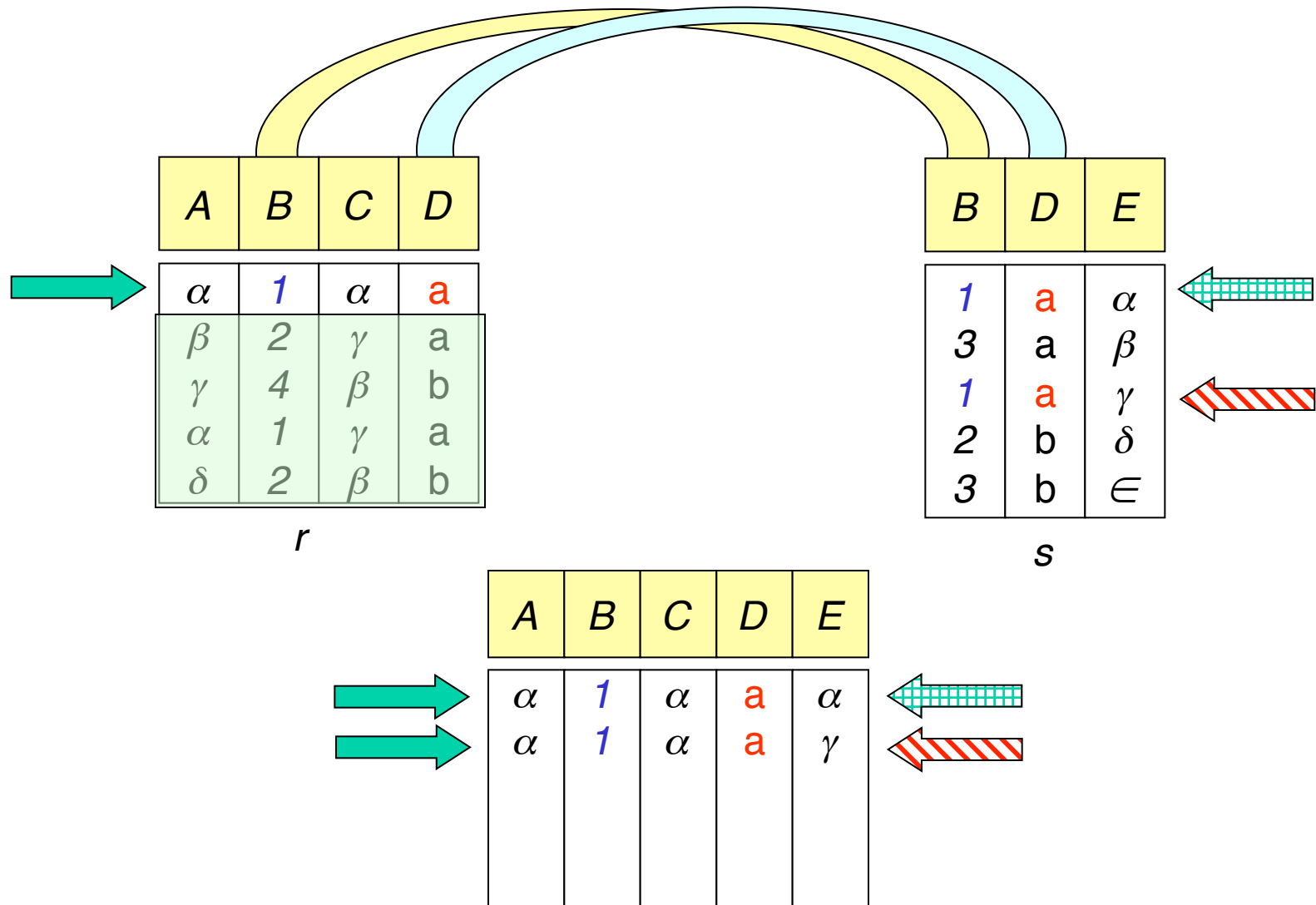
A	B	C	D	E
$\alpha$	1	$\alpha$	a	$\alpha$
$\alpha$	1	$\alpha$	a	$\gamma$
$\alpha$	1	$\gamma$	a	$\alpha$
$\alpha$	1	$\gamma$	a	$\gamma$
$\delta$	2	$\beta$	b	$\delta$

How did we get here?

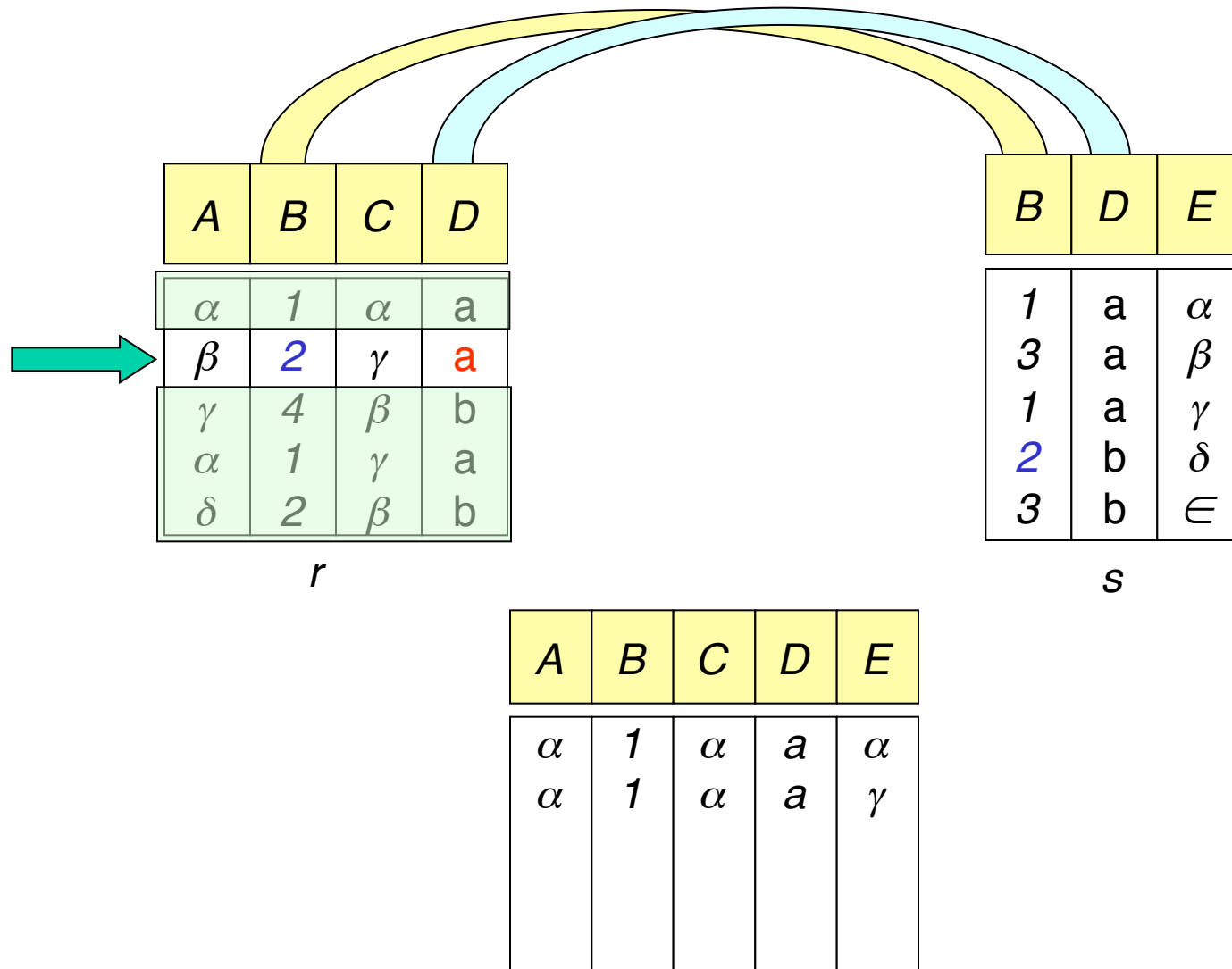
Lets do a trace over the next few slides...



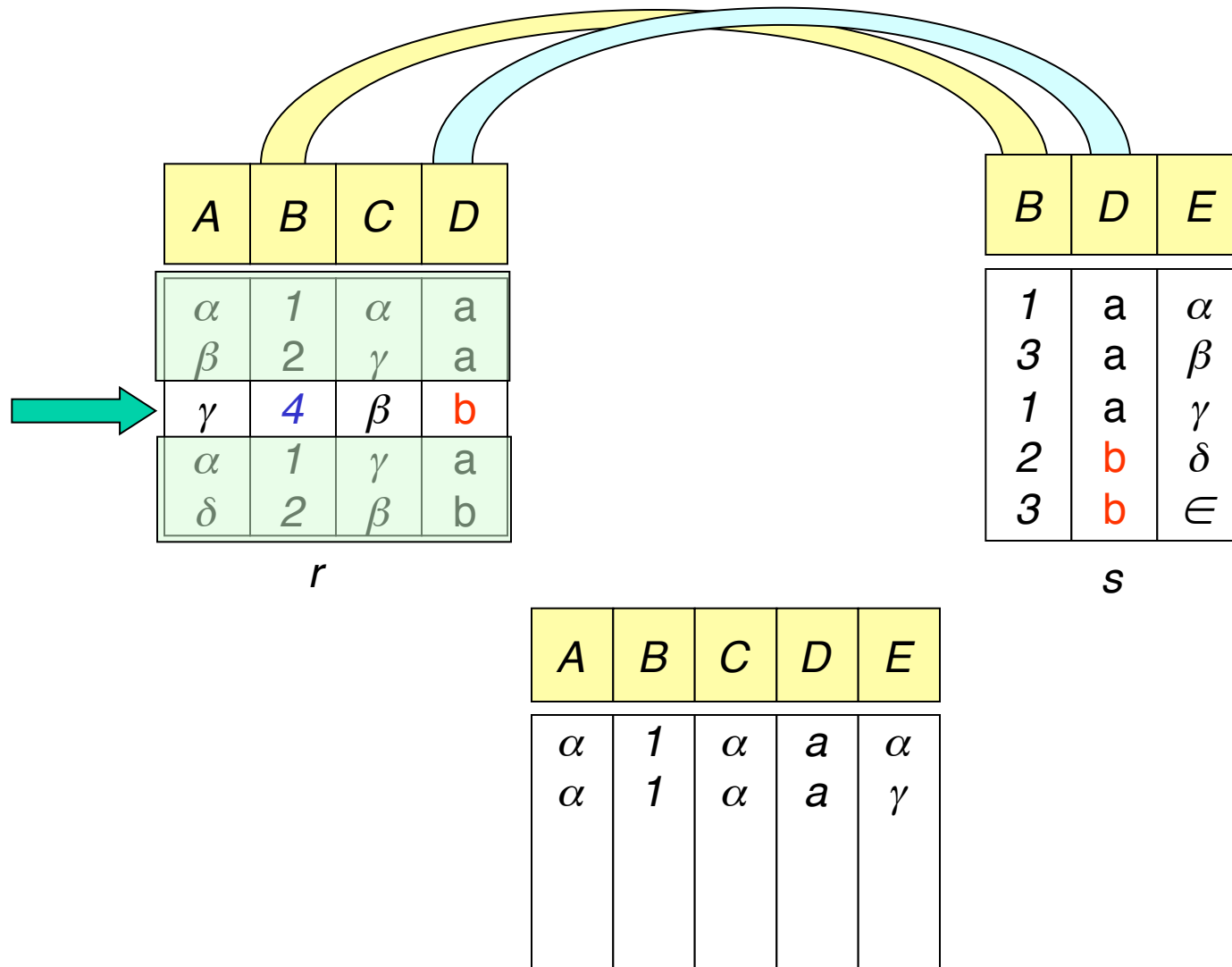
First we note which attributes the two relations have in common<sub>10</sub>..



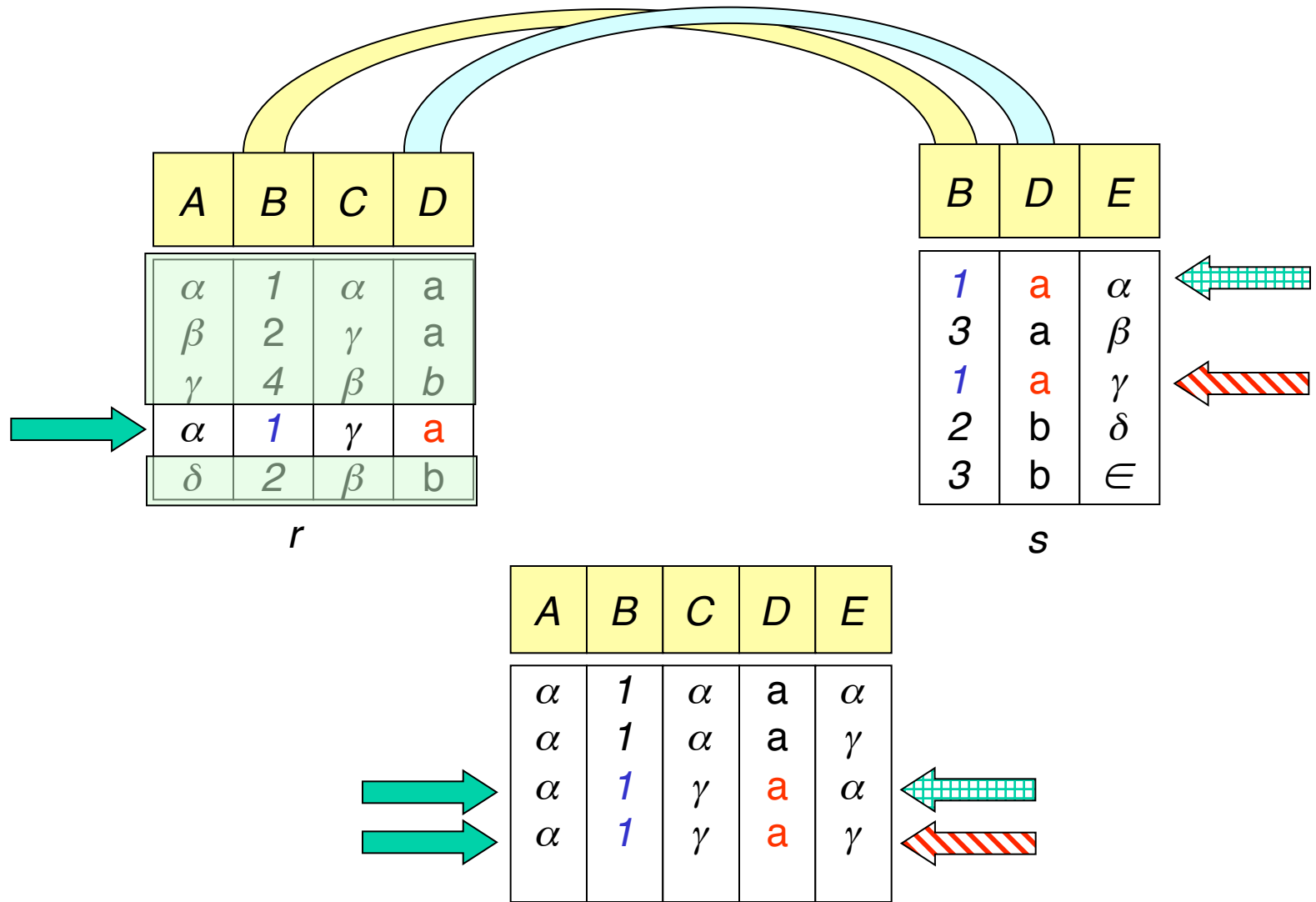
There are two rows in  $S$  that match our first row in  $r$ , (in the relevant attributes) so both are joined to our first row...



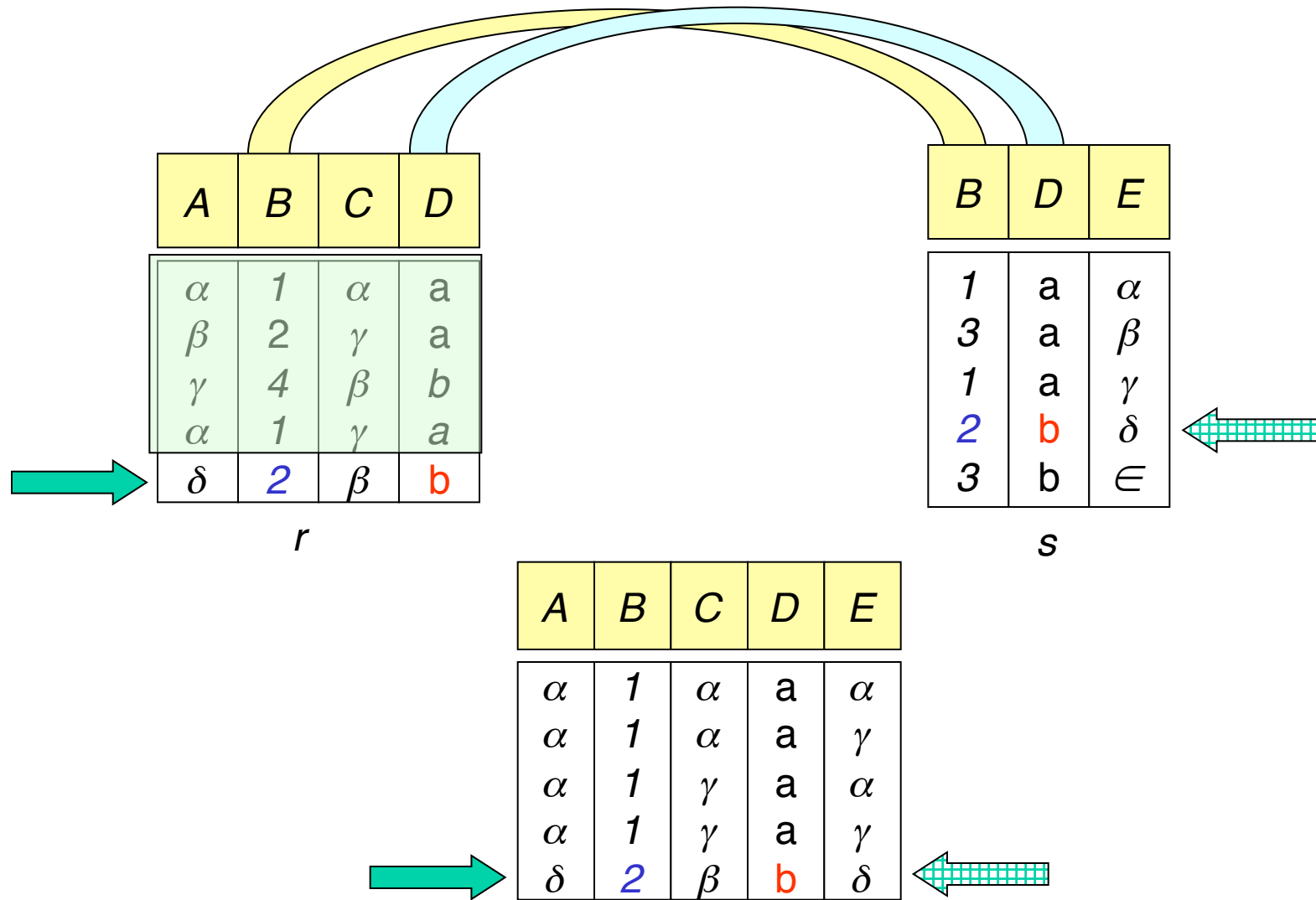
...there are no rows in  $S$  that match our second row in  $r$ , so do nothing...



...there are no rows in  $S$  that match our third row in  $r$ , so do nothing...



There are two rows in  $s$  that match our fourth row in  $r$ , so both are joined to our fourth row...



There is one row that matches our fifth row in  $r$ ,.. so it is joined to our fifth row and we are done!

# Natural Join on Sailors Example

S1

<u>sid</u>	sname	rating	age
22	dustin	7	45.0
31	lubber	8	55.5
58	rusty	10	35.0

R1

<u>sid</u>	<u>bid</u>	<u>day</u>
22	101	10/10/96
58	103	11/12/96

$S1 \bowtie R1 =$

sid	sname	rating	age	bid	day
22	dustin	7	45.0	101	10/10/96
58	rusty	10	35.0	103	11/12/96



## Earlier We Saw...

Query: Find the name of the sailor who reserved boat 101.

$$Temp = \rho(sid \rightarrow sid1, S1) \times \rho(sid \rightarrow sid2, R1)$$

$$Result = \pi_{Sname}(\sigma_{sid1=sid2 \wedge bid=101}(Temp))$$

\* Note my use of “temporary” relation Temp.

# Query revisited using natural join

Query: Find the name of the sailor who reserved boat 101.

$$\text{Result} = \pi_{Sname}(\sigma_{bid=101}(S1 \bowtie R1))$$

*Or*

$$\text{Result} = \pi_{Sname}(S1 \bowtie \sigma_{bid=101}(R1))$$

What's the difference between these two approaches?

# Conditional-Join Operation:

The conditional join is actually the most general type of join. I introduced the natural join first only because it is more intuitive and... natural!

Just like natural join, conditional join combines a cross product and a selection into one operation. However instead of only selecting rows that have equality on those attributes that appear in both relation schemes, we allow selection based on any predicate.

$$r \bowtie_c s = \sigma_c(r \times s)$$

Where  $c$  is any predicate  
the attributes of  $r$  and/or  $s$

Duplicate rows are removed as always, but duplicate columns are not removed!

# Conditional-Join Example:

We want to find all women that are younger than their husbands...

*r*

<i>l-name</i>	<i>f-name</i>	<u><i>marr-Lic</i></u>	<i>age</i>
Simpson	Marge	777	35
Lovejoy	Helen	234	38
Flanders	Maude	555	24
Krabappel	Edna	978	40

*S*

<i>l-name</i>	<i>f-name</i>	<u><i>marr-Lic</i></u>	<i>age</i>
Simpson	Homer	777	36
Lovejoy	Timothy	234	36
Simpson	Bart	<i>null</i>	9

$r \bowtie_{r.age < s.age \text{ AND } r.Marr-Lic = s.Marr-Lic} S$

<i>r.l-name</i>	<i>r.f-name</i>	<u><i>r.Marr-Lic</i></u>	<i>r.age</i>	<i>s.l-name</i>	<i>s.f-name</i>	<u><i>s.marr-Lic</i></u>	<i>s.age</i>
Simpson	Marge	777	35	Simpson	Homer	777	36

Note we have removed ambiguity of attribute names by using “dot” notation

Also note the redundant information in the *marr-lic* attributes

# Equi-Join

- Equi-Join: Special case of conditional join where the conditions consist only of equalities.
- Natural Join: Special case of equi-join in which equalities are specified on ALL fields having the same names in both relations.

# Equi-Join

*r*

<i>l-name</i>	<i>f-name</i>	<u><i>marr-Lic</i></u>	<i>age</i>
Simpson	Marge	777	35
Lovejoy	Helen	234	38
Flanders	Maude	555	24
Krabappel	Edna	978	40

*S*

<i>l-name</i>	<i>f-name</i>	<u><i>marr-Lic</i></u>	<i>age</i>
Simpson	Homer	777	36
Lovejoy	Timothy	234	36
Simpson	Bart	<i>null</i>	9

$$r \bowtie_{r.Marr-Lic = s.Marr-Lic} S$$

<i>r.l-name</i>	<i>r.f-name</i>	<u><i>Marr-Lic</i></u>	<i>r.age</i>	<i>s.l-name</i>	<i>s.f-name</i>	<i>s.age</i>
Simpson	Marge	777	35	Simpson	Homer	36
Lovejoy	Helen	234	38	Lovejoy	Timothy	36

# Review on Joins

- All joins combine a cross product and a selection into one operation.
- Conditional Join
  - the selection condition can be of any predicate (e.g.  $\text{rating1} > \text{rating2}$ )
- Equi-Join:
  - Special case of conditional join where the conditions consist only of equalities.
- Natural Join
  - Special case of equi-join in which equalities are specified on ALL fields having the same names in both relations.

# A Note on Precedence

- Unary operators have the highest precedence:  $[\sigma, \pi, \rho]$
- Then “multiplicative” operators:  $[\times, \otimes]$
- Then “additive” operators:  $[\cap, \cup, -]$



# Banking Examples

*branch (branch-id, branch-city, assets)*

*customer (customer-id, customer-name, customer-city)*

*account (account-number, branch-id, balance)*

*loan (loan-number, branch-id, amount)*

*depositor (customer-id, account-number)*

*borrower (customer-id, loan-number)*

# Example Queries 1

- Find all loans over \$1200

“select from the relation *loan*, only the rows which have a *amount* greater than 1200”

*loan*

<i>loan-number</i>	<i>branch-id</i>	<i>amount</i>
1234	001	1,923.03
3421	002	123.00
2342	004	56.25
4531	005	120.03

$\sigma_{amount > 1200} (loan)$

1234	001	1,923.03
------	-----	----------

# Example Queries 2

- Find the loan number for each loan of an amount greater than \$1200

“select from the relation *loan*, only the rows which have a *amount* greater than 1200, then project out just the *loan\_number*”

*loan*

<i>loan-number</i>	<i>branch-id</i>	<i>amount</i>
1234	001	1,923.03
3421	002	123.00
2342	004	56.25
4531	005	120.03

$\sigma_{amount > 1200} (loan)$

1234	001	1,923.03
------	-----	----------

$\pi_{loan-number} (\sigma_{amount > 1200} (loan))$

1234
------

# Example Queries 3

- Find all loans greater than \$1200 or less than \$75

“select from the relation *loan*, only the rows which have a *amount* greater than 1200 or an *amount* less than 75

*loan*

<i>loan-number</i>	<i>branch-id</i>	<i>amount</i>
1234	001	1,923.03
3421	002	123.00
2342	004	56.25
4531	005	120.03

$\sigma_{amount > 1200 \vee amount < 75}(loan)$

1234	001	1,923.03
2342	004	56.25

# Example Queries 4

- Find the IDs of all customers who have a loan, an account, or both, from the bank

*borrower*

<i>customer-id</i>	<i>loan-number</i>
201	1234
304	3421
425	2342
109	4531

*depositor*

<i>customer-id</i>	<i>account-number</i>
333	3467
304	2312
201	9999
492	3423

$\pi_{customer-id}(borrower)$

201
304
425
109

201
304
425
109
333
492

$\pi_{customer-id}(depositor)$

333
304
201
492

$$\pi_{customer-id}(borrower) \cup \pi_{customer-id}(depositor)$$

# Example Queries 5

Note this example is split over two slides!

Find the IDs of all customers who have a loan at branch 001.

*borrower*

<i>customer-id</i>	<i>loan-number</i>
201	1234
304	3421

We retrieve  
*borrower* and  
*loan*...

*loan*

<i>loan-number</i>	<i>branch-id</i>	<i>amount</i>
1234	001	1,923.03
3421	002	123.00

...we  
calculate  
their cross  
product...

<i>customer-id</i>	<i>borrower.loan-number</i>	<i>loan.loan-number</i>	<i>branch-id</i>	<i>amount</i>
201	1234	1234	001	1,923.03
201	1234	3421	002	123.00
304	3421	1234	001	1,923.03
304	3421	3421	002	123.00

...we calculate their cross product...

<i>customer-id</i>	<i>borrower.loan-number</i>	<i>loan.loan-number</i>	<i>branch-id</i>	<i>amount</i>
201	1234	1234	001	1,923.03
201	1234	3421	002	123.00
304	3421	1234	001	1,923.03
304	3421	3421	002	123.00

...we select the rows where *borrower.loan-number* is equal to *loan.loan-number*...

<i>customer-id</i>	<i>borrower.loan-number</i>	<i>loan.loan-number</i>	<i>branch-id</i>	<i>amount</i>
201	1234	1234	001	1,923.03
304	3421	3421	002	123.00

...we select the rows where *branch-id* is equal to "001"

<i>customer-id</i>	<i>borrower.loan-number</i>	<i>loan.loan-number</i>	<i>branch-id</i>	<i>amount</i>
201	1234	1234	001	1,923.03

...we project out the *customer-id*.

201

$$\pi_{customer-id} (\sigma_{branch-id='001'} (\sigma_{borrower.loan-number = loan.loan-number}(borrower \times loan)))$$

# Now Using Natural Join

Find the IDs of all customers who have a loan at branch 001.

We retrieve *borrower*  
and *loan*...

*borrower*

<i>customer-id</i>	<i>loan-number</i>
201	1234
304	3421

*loan*

<i>loan-number</i>	<i>branch-id</i>	<i>amount</i>
1234	001	1,923.03
3421	002	123.00

1234 in *borrower* is  
matched with 1234 in  
*loan*...

<i>customer-id</i>	<i>loan-number</i>	<i>branch-id</i>	<i>amount</i>
201	1234	001	1,923.03
304	3421	002	123.00

3421 in *borrower* is  
matched with 3421 in  
*loan*...

The rest is the same.

<i>customer-id</i>	<i>loan-number</i>	<i>branch-id</i>	<i>amount</i>
201	1234	001	1,923.03

$$\pi_{customer-id} (\sigma_{branch-id='001'} (\sigma_{borrower.loan-number = loan.loan-number} (borrower \times loan)))$$

$$= \pi_{customer-id} (\sigma_{branch-id='001'} (borrower \bowtie loan))$$



# Example Queries 7

Note this example is split over two slides!

- Find the *names* of all customers who have a loan, an account, or both, from the bank

$\pi_{customer-id} (borrower) \cup \pi_{customer-id} (depositor)$       customer

<i>customer-id</i>	<i>customer-name</i>	<i>customer-city</i>
101	Carol	Fairfax
109	David	Fairfax
201	John	Vienna
304	Mary	McLean
333	Ben	Chantilly
425	David	Manassas
492	Jason	Fairfax
501	Adam	Burke

# Example Queries

- Find the *names* of all customers who have a loan, an account, or both, from the bank

<i>customer-id</i>	<i>customer-name</i>	<i>customer-city</i>
109	David	Fairfax
201	John	Vienna
304	Mary	McLean
333	Ben	Chantilly
425	David	Manassas
492	Jason	Fairfax

<i>customer-name</i>
David
John
Mary
Ben
David
Jason

<i>customer-name</i>
David
John
Mary
Ben
Jason

$$\pi_{customer-name} \left( \left( \pi_{customer-id} (borrower) \cup \pi_{customer-id} (depositor) \right) \bowtie customer \right)$$

# Example Queries 8

Note this example is split over three slides!

Find the largest account balance

*account*

<i>account-number</i>	<i>branch-id</i>	<i>balance</i>
7777	001	100.30
8888	003	12.34
6666	004	45.34

We do a rename to get a “copy” of the balance column from *account*. We call this copy *d*...

*d*

<i>balance</i>
100.30
12.34
45.34

... next we will do a cross product...

... do a cross product...

...select out all rows where *account.balance* is less than *d.balance*...

<i>account-number</i>	<i>branch-id</i>	<i>account.balance</i>	<i>d.balance</i>
7777	001	100.30	100.30
7777	001	100.30	12.34
7777	001	100.30	45.34
8888	003	12.34	100.30
8888	003	12.34	12.34
8888	003	12.34	45.34
6666	004	45.34	100.30
6666	004	45.34	12.34
6666	004	45.34	45.34

<i>account-number</i>	<i>branch-id</i>	<i>account.balance</i>	<i>d.balance</i>
8888	003	12.34	100.30
8888	003	12.34	45.34
6666	004	45.34	100.30

.. next we project out *account.balance*...

...then we do a set difference between it and the original *account.balance* from the account relation...

... the set difference leaves us with one number, the largest value!

<i>account-number</i>	<i>branch-id</i>	<i>account.balance</i>	<i>d.balance</i>
8888	003	12.34	100.30
8888	003	12.34	45.34
6666	004	45.34	100.30

*balance from account*

<i>balance</i>
100.30
12.34
45.34

—

<i>account.balance</i>
12.34
45.34

100.30
--------

$$\pi_{balance}(account) - \pi_{account.balance}(\sigma_{account.balance < d.balance} (account \times \rho(d, \pi_{balance}(account))))$$

# Now Using Conditional Join

Find the largest account balance

$\pi_{balance}(account) - \pi_{account.balance}(\sigma_{account.balance < d.balance} (account \times \rho(d, \pi_{balance}(account))))$

$\rho(d, \pi_{account.balance}(account))$

$\pi_{balance}(account) - \pi_{account.balance}(account \bowtie_{account.balance < d.balance} d)$

# More Examples on Sailors Relations

Sailors(sid, sname, rating, age)

Boats(bid, bname, color)

Reserves(sid, bid, day)

# Find names of sailors who've reserved boat #103

- Solution 1: Find those who reserved boat 103, join with Sailors to find the names, and project out the names

$$\pi_{sname}((\sigma_{bid=103}(\text{Reserves})) \bowtie \text{Sailors})$$

- Solution 2: Join Reserves and Sailors to get all information, and find those who reserved boat 103. Project out the names.

$$\pi_{sname}(\sigma_{bid=103}(\text{Reserves} \bowtie \text{Sailors}))$$

Which one is more efficient?



## Find names of sailors who've reserved a red boat

- Information about boat color only available in Boats; so need an extra join:

- A more efficient solution: Find the bids of red boats first before doing the join.

☞ *A query optimizer can find this given the first solution!*

Find sailors who've reserved a **red** or a **green**  
boat

- Can identify all red or green boats, then find sailors who've reserved one of these boats:

# Find sailors who've reserved a **red** or a **green** boat

- Can identify all red or green boats, then find sailors who've reserved one of these boats:

- Can also define Tempboats using union:

- What happens if “or” is replaced by “and”?

---

Find sailors who've reserved a **red** and a  
**green** boat

Find sailors who've reserved a **red** and a **green** boat

- Previous first approach won't work! (Why not?) Must use intersection.

$$Tempred = \pi_{sid}((\sigma_{color='red'}(Boats)) \bowtie Reserves)$$
$$Tempgreen = \pi_{sid}((\sigma_{color='green'}(Boats)) \bowtie Reserves)$$
$$Result = \pi_{sname}((Tempred \cap Tempgreen) \bowtie Sailors)$$

# Consider yet another query

- Find the sailor(s) who reserved all the red boats.

*R1*

<i>sid</i>	<i>bid</i>	<i>day</i>
22	101	10/10/96
22	103	10/11/96
56	102	11/12/96

*B*

<i>bid</i>	<i>color</i>
101	Red
102	Green
103	Red

# Start an attempt

- Who reserved what boat:

$$S1 = \pi_{sid, bid}(R1) =$$

<i>sid</i>	<i>bid</i>
22	101
22	103
56	102

- All the red boats:

$$S2 = \pi_{bid}(\sigma_{color=red}(B)) =$$

<i>bid</i>
101
103

Now what?

Find the sailor(s) who reserved all the red boats.

- We will solve the problem the “hard” way, and then will introduce an operator specifically for this kind of problem.
- *Idea*: Compute the sids of sailors who *didn't* reserve all red boats.
  1. Find all possible reservations that could be made on red boats.
  2. Find *actual* reservations on red boats
  3. Find the possible reservations on red boats that were not actually made (#1 – #2) <- that is a minus sign.
  4. Project out the sids from 3 – these are the sailors who didn't have reservation on some red boat(s).



# Find the sailor(s) who reserved all the red boats.

- *Idea:* Compute the sids of sailors who *didn't* reserve all red boats (then find the difference between this set and set of all sailors).

1. Find all possible reservations that could be made on red boats.

$$\text{AllPossibleRes} = \pi_{\text{sid}} (\text{R1}) \times \pi_{\text{bid}} \sigma_{\text{color}=\text{"red"}} (\text{B})$$

2. Find *actual* reservations on red boats

$$\text{AllRedRes} = \pi_{\text{sid,bid}} (\text{R1}) \bowtie \pi_{\text{bid}} \sigma_{\text{color}=\text{"red"}} (\text{B})$$

3. 4. Find the possible reservations on red boats that were not actually made, and project out the sids.

$$\pi_{\text{sid}} (\text{AllPossibleRes} - \text{AllRedRes})$$

5. Find sids that are not in the result from above (sailors such that there is no red boat that's not reserved by him/her)

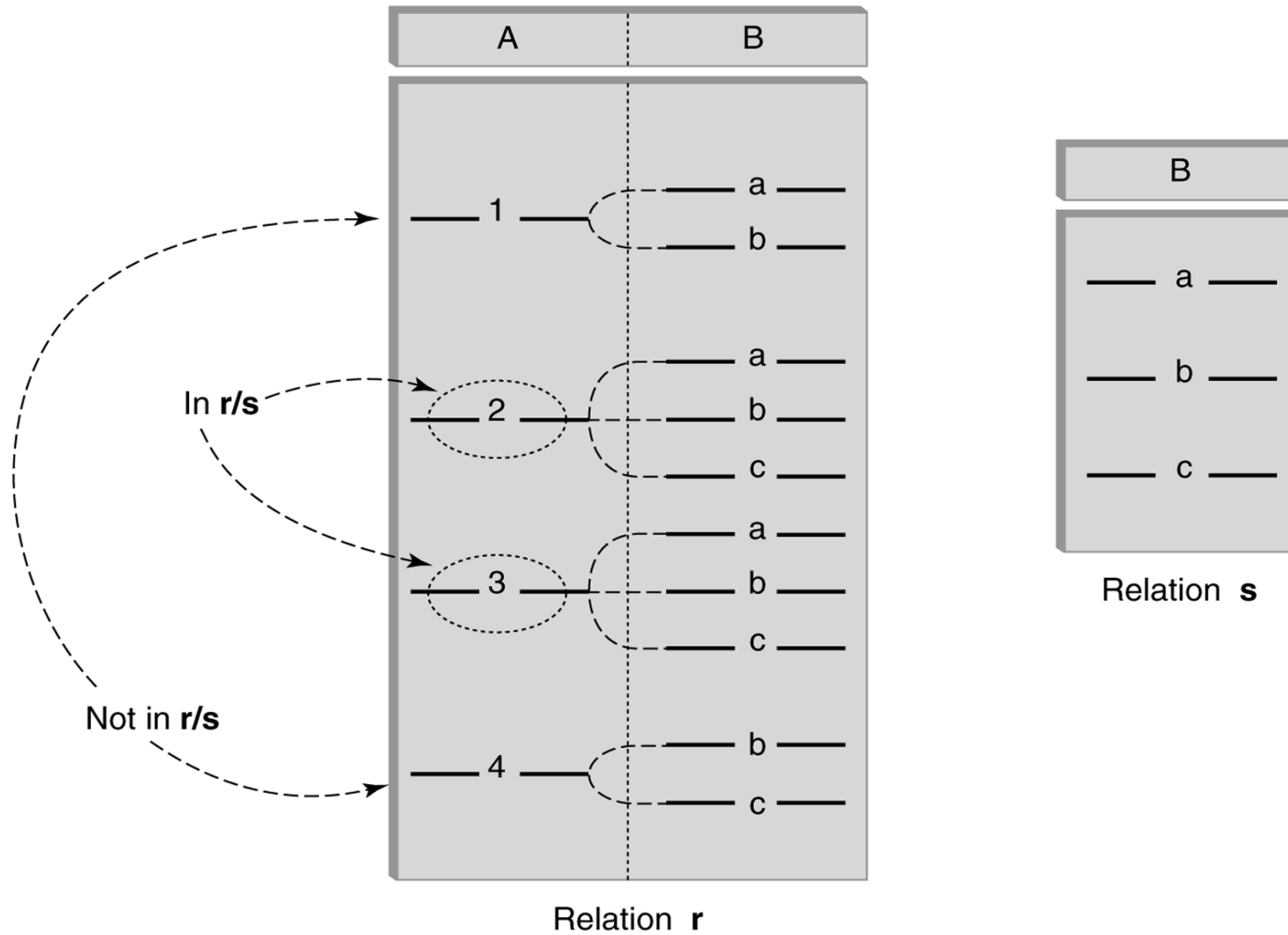
$$\pi_{\text{sid}} (\text{R1}) - \pi_{\text{sid}} (\text{AllPossibleRes} - \text{AllRedRes})$$

# Division Operation

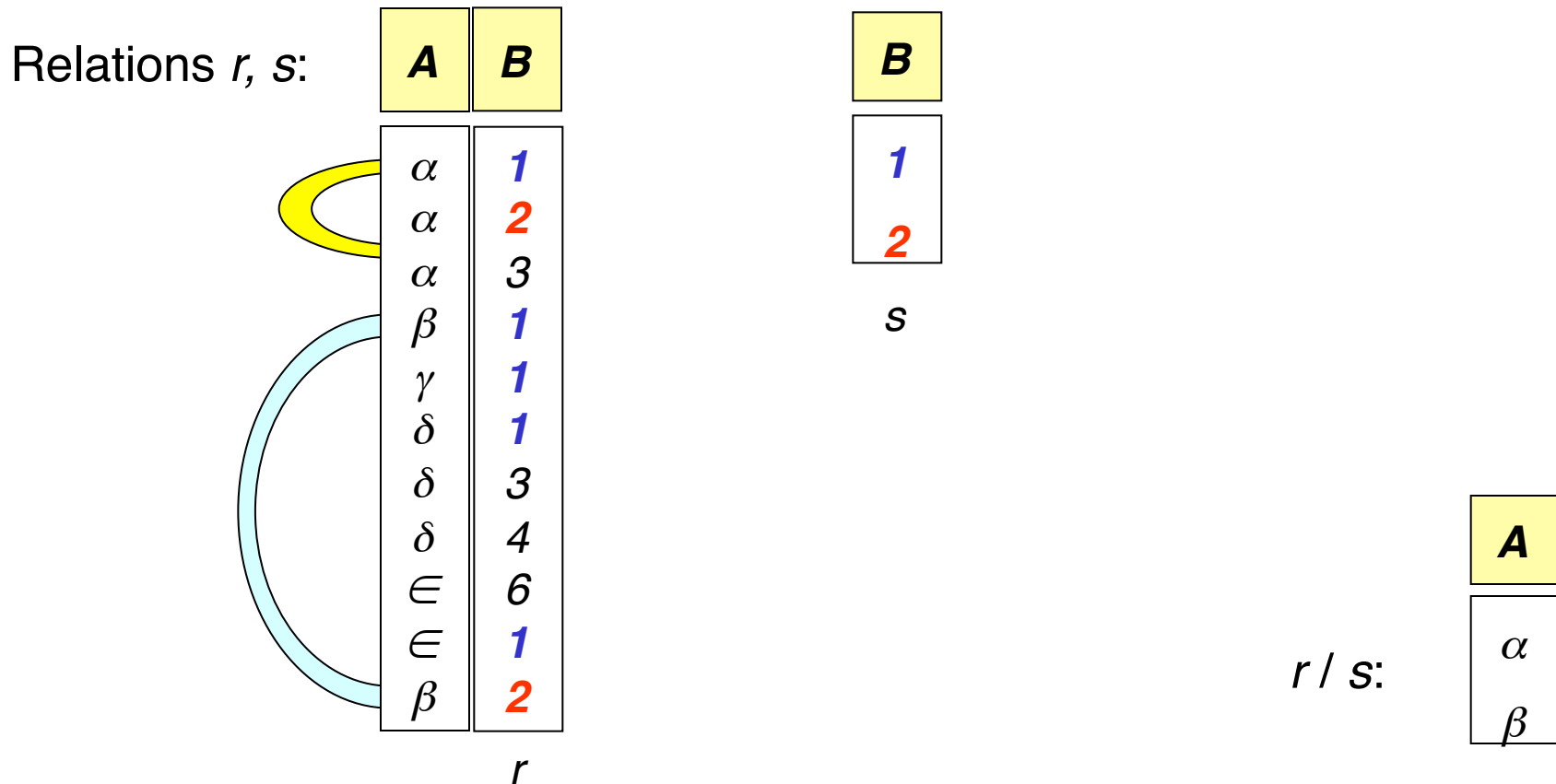
$r / s$

- Suited to queries that include the phrase “for all”, e.g. *Find sailors who have reserved all red boats.*
- Produce the tuples in one relation,  $r$ , that match *all* tuples in another relation,  $s$
- Let  $S1$  have 2 fields,  $x$  and  $y$ ;  $S2$  have only field  $y$ :
  - $S1/S2 = \{ \langle x \rangle \mid \forall \langle y \rangle \text{ in } S2 (\exists \langle x, y \rangle \text{ in } S1) \}$
  - i.e.,  **$S1/S2$  contains all  $x$  tuples (sailors) such that for every  $y$  tuple (redboat) in  $S2$ , there is an  $xy$  tuple in  $S1$  (i.e,  $x$  reserved  $y$ ).**
- In general,  $x$  and  $y$  can be any lists of fields;  $y$  is the list of fields in  $S2$ , and  $x \cup y$  is the list of fields of  $S1$ .
- Let  $r$  and  $s$  be relations on schemas  $R$  and  $S$  respectively where
  - $R = (A_1, \dots, A_m, B_1, \dots, B_n)$ ,
  - $S = (B_1, \dots, B_n)$ ,The result of  $r / s$  is a relation on schema  
 $R - S = (A_1, \dots, A_m)$

# Division (cont'd)



# Division Operation – Example



$\alpha$  occurs in the presence of both **1** and **2**, so it is returned.

$\beta$  occurs in the presence of both **1** and **2**, so it is returned.

$\gamma$  does not occur in the presence of both **1** and **2**, so is ignored.

...

# Another Division Example

Relations  $r, s$ :

A	B	C	D	E
$\alpha$	a	$\alpha$	a	1
$\alpha$	a	$\gamma$	a	1
$\alpha$	a	$\gamma$	b	1
$\beta$	a	$\gamma$	a	1
$\beta$	a	$\gamma$	b	3
$\gamma$	a	$\gamma$	a	1
$\gamma$	a	$\gamma$	b	1
$\gamma$	a	$\beta$	b	1

$r$

D	E
a	1
b	1

$s$

$r/s$ :

A	B	C
$\alpha$	a	$\gamma$
$\gamma$	a	$\gamma$

$\langle \alpha, a, \gamma \rangle$  occurs in the presence of both  $\langle a, 1 \rangle$  and  $\langle b, 1 \rangle$ , so it is returned.

$\langle \gamma, a, \gamma \rangle$  occurs in the presence of both  $\langle a, 1 \rangle$  and  $\langle b, 1 \rangle$ , so it is returned.

$\langle \beta, a, \gamma \rangle$  does not occur in the presence of both  $\langle a, 1 \rangle$  and  $\langle b, 1 \rangle$ , so it is ignored.

# More Division Examples: A/B

sno	pno
s1	p1
s1	p2
s1	p3
s1	p4
s2	p1
s2	p2
s3	p2
s4	p2
s4	p4

*A*

pno
p2

*B1*

sno
s1
s2
s3
s4

*A/B1*

pno
p2
p4

*B2*

sno
s1
s4

*A/B2*

pno
p1
p2
p4

*B3*

sno
s1

*A/B3*

# Find the sailor(s) who reserved ALL red boats

- who reserved what boat:

$$S1 = \pi_{sid, bid}(R1) =$$

<u>sid</u>	<u>bid</u>
22	101
22	103
58	102

- All the red boats:

$$S2 = \pi_{bid}(\sigma_{color=red}(B)) =$$

<u>bid</u>
101
103

**=> S1/S2**

Find the names of sailors who've reserved all boats

- Uses division; schemas of the input relations to “divide” must be carefully chosen:

$$\text{Tempsids} = (\pi_{sid, bid}(\text{Reserves})) / (\pi_{bid}(\text{Boats}))$$

$$\text{Result} = \pi_{sname}(\text{Tempsids} \bowtie \text{Sailors})$$



- SALES(supId, prodId);
- PRODUCTS(prodId);
- SALES/PRODUCTS = ?

# Expressing A/B Using Basic Operators

- Division is not essential op; just a useful shorthand.
  - (Also true of joins, but joins are so common that systems implement joins specially. Division is NOT implemented in SQL).
- *Idea:* For *SALES/PRODUCTS*, compute the IDs of suppliers that don't supply all products.

$$A = \pi_{sid}((\pi_{sid}(Sales) \times Products) - Sales)$$

The answer is  $\pi_{sid}(Sales) - A$

# Additional Operator: Outer Join

- An extension of the join operation that avoids loss of information.
- Computes the join and then adds tuples from one relation that does not match tuples in the other relation to the result of the join.
- Uses *null* values:
  - *null* signifies that the value is unknown or does not exist
  - All comparisons involving *null* are (roughly speaking) **false** by definition.
    - Will study precise meaning of comparisons with nulls later

# Outer Join – Example

Relation *loan*

<i>loan-number</i>	<i>branch-name</i>	<i>amount</i>
L-170	Springfield	3000
L-230	Shelbyville	4000
L-260	Dublin	1700

Relation *borrower*

<i>customer-name</i>	<i>loan-number</i>
Simpson	L-170
Wiggum	L-230
Flanders	L-155

# Outer Join – Example

- **Inner Join**

<i>loan-number</i>	<i>branch-name</i>	<i>amount</i>	<i>customer-name</i>
L-170	Springfield	3000	Simpson
L-230	Shelbyville	4000	Wiggum

*loan* ⋈ *Borrower*

<i>loan-number</i>	<i>branch-name</i>	<i>amount</i>
L-170	Springfield	3000
L-230	Shelbyville	4000
L-260	Dublin	1700

<i>customer-name</i>	<i>loan-number</i>
Simpson	L-170
Wiggum	L-230
Flanders	L-155

- **Left Outer Join**

*loan* ⋈<sub>l</sub> *borrower*

<i>loan-number</i>	<i>branch-name</i>	<i>amount</i>	<i>customer-name</i>
L-170	Springfield	3000	Simpson
L-230	Shelbyville	4000	Wiggum
L-260	Dublin	1700	<i>null</i>

# Outer Join – Example

## Right Outer Join

<i>loan-number</i>	<i>branch-name</i>	<i>amount</i>	<i>customer-name</i>
L-170	Springfield	3000	Simpson
L-230	Shelbyville	4000	Wiggum
L-155	<i>null</i>	<i>null</i>	Flanders

*loan* ⋈ *borrower*

<i>loan-number</i>	<i>branch-name</i>	<i>amount</i>
L-170	Springfield	3000
L-230	Shelbyville	4000
L-260	Dublin	1700

<i>customer-name</i>	<i>loan-number</i>
Simpson	L-170
Wiggum	L-230
Flanders	L-155

## Full Outer Join

*loan* ⋈ *borrower*

<i>loan-number</i>	<i>branch-name</i>	<i>amount</i>	<i>customer-name</i>
L-170	Springfield	3000	Simpson
L-230	Shelbyville	4000	Wiggum
L-260	Dublin	1700	<i>null</i>
L-155	<i>null</i>	<i>null</i>	Flanders

# Null Values

- It is possible for tuples to have a null value, denoted by *null*, for some of their attributes
- *null* signifies an unknown value or that a value does not exist.
- The result of any arithmetic expression involving *null* is *null*.
- Aggregate functions simply ignore null values
  - Is an arbitrary decision. Could have returned null as result instead.
  - We follow the semantics of SQL in its handling of null values
- For duplicate elimination and grouping, null is treated like any other value, and two nulls are assumed to be the same
  - Alternative: assume each null is different from each other
  - Both are arbitrary decisions, so we simply follow SQL

# Null Values

- Comparisons with null values return the special truth value *unknown*
- Three-valued logic using the truth value *unknown*:
  - OR: (*unknown* **or** *true*) = *true*,  
(*unknown* **or** *false*) = *unknown*  
(*unknown* **or** *unknown*) = *unknown*
  - AND: (*true* **and** *unknown*) = *unknown*,  
(*false* **and** *unknown*) = *false*,  
(*unknown* **and** *unknown*) = *unknown*
  - NOT: (**not** *unknown*) = *unknown*
  - In SQL “*P* is **unknown**” evaluates to true if predicate *P* evaluates to *unknown*
- Result of select predicate is treated as *false* if it evaluates to *unknown*



# Summary

- The relational model has rigorously defined query languages that are simple and powerful.
- Relational algebra is more operational; useful as internal representation for query evaluation plans.
- Several ways of expressing a given query; a query optimizer should choose the most efficient version.
- Operations covered: 5 basic operations (selection, projection, union, set difference, cross product), rename, joins (natural join, equi-join, conditional join, outer joins), division