

Set Joins and Semijoins in Relational Algebra

Dirk Van Gucht¹

¹Indiana University

Joins and semijoins

- Two types of joins
 - **Regular joins**: these are joins that compute relationships between objects (o_1, o_2) based on **tuple component comparisons** (see the previous lecture on regular joins.)
 - **Set joins**: these are joins that compute relationships between objects (o_1, o_2) based on **set comparisons between sets** $S(o_1)$ and $S(o_2)$ that are associated with objects o_1 and o_2 , respectively.
 - Think for example of finding pairs (s, p) such that student s only takes courses taught by professors.
 - If $S(s)$ denotes the set of courses taken by s and $S(p)$ denotes the set of courses taught by p , then we need to check if $C(s) \subseteq C(p)$.
- Analogously, there are also two types of semijoins: regular semijoins and set semijoins

Set joins (Motivation)

- Consider the relations **Enroll**(sid,cno) and **TaughtBy**(cno,pid).
- (s, c) is in **Enroll** when student s is enrolled in course c
- (c, p) is in **TaughtBy** when course c is taught by professor p
- We are interested in student-professor relationships (s, p) based on some complex join condition.
- "Find the student-professor pairs (s, p) such that student s takes *[some|not only|not all|no|only|all|at least 2|...]* course(s) taught by professor p ."

Set joins (Motivation)

- Consider the relations **Patient**(pid, symptom) and **Disease**(symptom,disease).
- (p, s) is in **Patient** when patient p exhibits health symptom s
- (s, d) is in **Disease** when s is a symptom of disease d
- "Find the patient-disease pairs (p, d) such that patient p has [some|not only|not all|no|only|all|at least 2|...] symptoms associated with disease d ."

Set semi joins (Motivation)

- Consider the relations *Patient(pid,sympton)* and *Flu(sympton)*.
- We are interested in finding patient that satisfy some complex condition.
- "Find each patient *p* such that *p* has *[some|not only|not all|no|only|all|at least 2|...]* flu symptoms."

Set **semi**joins (Advanced Google Search Engine)

- Consider the relations **Document**(doc,word) and the set of words {*Advanced, Database, Concept*}.
- A pair $(d, w) \in \text{Document}$ when document d contains the word w .
- Advanced Google Search supports the following queries:
"Find each document d such that d contains
[*some(any)|no|all*] words in
{*Advanced, Database, Concepts*}.
- You can try this out for the *all* query and you will see documents that refer to B561.
- [https : // www.google.com/advanced_search](https://www.google.com/advanced_search)
- So set semijoins are frequently used in real search applications

Set joins and semijoin in RA (Objective)

- Derive relational algebra expressions for set joins and semijoins that supports queries with quantifiers.
- The technique we use for this follows the Venn-diagram method used in the lecture on "Queries with Quantifiers" to derive SQL queries for such queries.
- We will discuss how the derived RA expressions provide deep insight into the time and space complexities for different set joins
- Note: this is a very technical lecture that relies on understanding of logic and set theory

Set joins in RA (Introduction-Terminology)

We begin by introducing some terminology and state our assumptions.

- Denote by \mathcal{S} the set of all student sids
- Denote by \mathcal{P} the set of all professor pids
- Denote by \mathcal{C} the set of all course cnos
- Denote by E the relation *Enroll*
- Denote by T the relation *TaughtBy*
- We will assume that the natural foreign key constraints hold:

$$\pi_{sid}(E) \subseteq \mathcal{S}$$

$$\pi_{cno}(E) \subseteq \mathcal{C}$$

$$\pi_{cno}(T) \subseteq \mathcal{C}$$

$$\pi_{pid}(T) \subseteq \mathcal{P}$$

Set joins in RA (Parameterized sets)

- Let s be a student sid (i.e., $s \in \mathcal{S}$) and let p be a professor pid (i.e., $p \in \mathcal{P}$).
- We will consider **sets parameterized** by s and p .
- Both these sets are subsets of \mathcal{C} , i.e., the set of cnos of courses.
 - $E(s)$ is the set of cnos of courses in which student s is enrolled.¹

$$E(s) = \{c \mid (s, c) \text{ is a tuple in the Enroll relation}\}$$

- $T(p)$ is the set of cnos of courses taught by professor p .²

$$T(p) = \{c \mid (c, p) \text{ is a tuple in the TaughtBy relation}\}$$

¹If s is not enrolled in any course then $E(s) = \emptyset$.

²If p does not teach any course then $T(p) = \emptyset$.

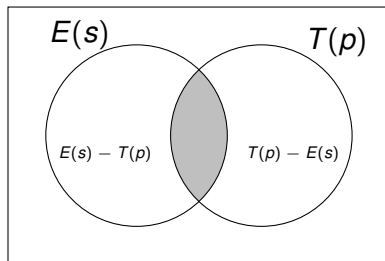
Set joins (Regions)

- For each student-professor pair (s, p) , the sets $E(s)$ and $T(p)$ define 4 pairwise disjoint "regions" in the Venn diagram associated with these sets.³
- These regions partition the set of cno of courses \mathcal{C} :

$E(s) \cap T(p)$	=	the courses enrolling student s and taught by professor p
$E(s) - T(p)$	=	the courses enrolling student s but not taught by professor p
$T(p) - E(s)$	=	the courses not enrolling student s but taught by professor p
$\mathcal{C} - (E(s) \cup T(p))$	=	the courses not enrolling student s and not taught by professor p

³It is recommended that you revisit the lecture of "Queries with Quantifiers" to further understand this Venn diagram.

Venn diagram and regions for $E(s)$ and $T(p)$



The gray region is $E(s) \cap T(p)$

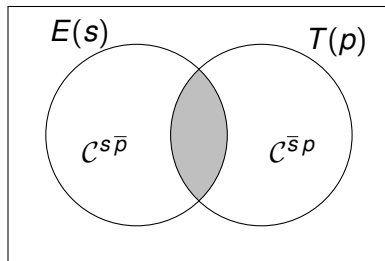
The region outside both $E(s)$ and $T(p)$ is $\mathcal{C} - (E(s) \cup T(p))$

Set joins (Regions)

- For each student-professor pair (s, p) , the sets $E(s)$ and $T(p)$ define 4 pairwise disjoint "regions" that partition the set of cnos of courses \mathcal{C} :
- We use the following notations to denote these regions:

Region	Notation
$E(s) \cap T(p)$	\mathcal{C}^{sp}
$E(s) - T(p)$	$\mathcal{C}^{s\bar{p}}$
$T(p) - E(s)$	$\mathcal{C}^{\bar{s}p}$
$\mathcal{C} - (E(s) \cup T(p))$	$\mathcal{C}^{\bar{s}\bar{p}}$

Venn diagram and regions for $E(s)$ and $T(p)$



The gray region is C^{sp}

The region outside both $E(s)$ and $T(p)$ is $C^{\bar{s}\bar{p}}$. (From now on, we will ignore this region in the lecture.)

Set joins (Region analysis)

- Consider again the regions \mathcal{C}^{sp} , $\mathcal{C}^{s\bar{p}}$, and $\mathcal{C}^{\bar{s}p}$.
- These regions can be either not-empty or empty.
- Depending on this, we can distinguish 6 set joins. They correspond to the **some**, **no**, **not only**, **only**, **not all**, and **all** quantifiers we considered earlier:

$$\mathcal{C}^{sp} \neq \emptyset \quad \leftrightarrow \quad \text{some}$$

$$\mathcal{C}^{sp} = \emptyset \quad \leftrightarrow \quad \text{no}$$

$$\mathcal{C}^{s\bar{p}} \neq \emptyset \quad \leftrightarrow \quad \text{not only}$$

$$\mathcal{C}^{s\bar{p}} = \emptyset \quad \leftrightarrow \quad \text{only}$$

$$\mathcal{C}^{\bar{s}p} \neq \emptyset \quad \leftrightarrow \quad \text{not all}$$

$$\mathcal{C}^{\bar{s}p} = \emptyset \quad \leftrightarrow \quad \text{all}$$

- Once the analysis of these cases is understood, it becomes straightforward to study a very large set of other interesting (and practical) set joins.

Set join (not-empty cases)

- We begin with the set joins for **some**, **not only**, and **not all**. I.e., those corresponding to **non-empty** regions.
- To perform our analysis, it helps to formulate precisely the queries we want to express as RA expressions:

some	$\leftrightarrow \{(s, p) \exists c : c \in \mathcal{C}^{sp}\}$
	$\leftrightarrow \{(s, p) \exists c : c \in E(s) \cap T(p)\}$
	$\leftrightarrow \{(s, p) \exists c : c \in E(s) \wedge c \in T(p)\}$
not only	$\leftrightarrow \{(s, p) \exists c : c \in \mathcal{C}^{s\bar{p}}\}$
	$\leftrightarrow \{(s, p) \exists c : c \in (E(s) - T(p))\}$
	$\leftrightarrow \{(s, p) \exists c : c \in E(s) \wedge c \notin T(p)\}$
not all	$\leftrightarrow \{(s, p) \exists c : c \in \mathcal{C}^{\bar{s}p}\}$
	$\leftrightarrow \{(s, p) \exists c : c \in (T(p) - E(s))\}$
	$\leftrightarrow \{(s, p) \exists c : c \in T(p) \wedge c \notin E(s)\}$

Set join (**some** case)

- Recall

$$\begin{aligned}\text{some} &\leftrightarrow \{(s, p) \mid \exists c : c \in \mathcal{C}^{sp}\} \\ &\leftrightarrow \{(s, p) \mid \exists c : c \in E(s) \wedge c \in T(p)\} \\ &\leftrightarrow \{(s, p) \mid \exists c : E(s, c) \wedge T(c, p)\}\end{aligned}$$

- We first need to find an RA expression for the set

$$\{(s, c, p) \mid E(s, c) \wedge T(c, p)\}.$$

- But this is easy since, by the definition of the natural join \bowtie

$$\{(s, c, p) \mid E(s, c) \wedge T(c, p)\} = E \bowtie T.$$

Consequently, the RA expression for **some** is

$$\pi_{sid, pid}(E \bowtie T)$$

Set join (**not only** case)

- Recall

$$\begin{aligned}\text{not only} &\Leftrightarrow \{(s, p) \mid \exists c : c \in \mathcal{C}^{s\bar{p}}\} \\ &\Leftrightarrow \{(s, p) \mid \exists c : c \in E(s) \wedge \neg(c \in T(p))\} \\ &\Leftrightarrow \{(s, p) \mid \exists c : E(s, c) \wedge \neg T(c, p)\}\end{aligned}$$

- We first need to find an RA expression for the set

$$\{(s, c, p) \mid E(s, c) \wedge \neg T(c, p)\}.$$

$$\begin{aligned}\{(s, c, p) \mid E(s, c) \wedge \neg T(c, p)\} &= \{(s, c, p) \mid (E(s, c) \wedge p \in \mathcal{P}) \wedge \neg(E(s, c) \wedge T(c, p))\} \\ &= \{(s, c, p) \mid (s, c, p) \in (E \times \mathcal{P}) \wedge \neg((s, c, p) \in E \bowtie T)\} \\ &= (E \times \mathcal{P}) - (E \bowtie T)\end{aligned} \tag{1}$$

The equality (1) will be shown on the next slides.
Consequently, the RA expression for **not only** is

$$\pi_{sid, pid}((E \times \mathcal{P}) - (E \bowtie T))$$

Set join (**not only** case using natural join)

- To prove equality (1), we need to show that the following statements are logically equivalent

$$E(s, c) \wedge \neg T(c, p)$$

$$(E(s, c) \wedge p \in \mathcal{P}) \wedge \neg(E(s, c) \wedge T(c, p))$$

Since, clearly $p \in \mathcal{P}$ is true, we have to show that the following statements are equivalent:⁴

$$E(s, c) \wedge \neg T(c, p)$$

$$E(s, c) \wedge \neg(E(s, c) \wedge T(c, p))$$

⁴We are assuming that there is at least one professor.

Set join (**not only** case using natural join)

We need to show that the following statements are equivalent:

$$E(s, c) \wedge \neg T(c, p) \quad (I)$$

$$E(s, c) \wedge \neg(E(s, c) \wedge T(c, p)) \quad (II)$$

To see this, you need to determine if statements (I) and (II) have the same truth values. This is the case as shown in the following truth table:

$E(s, c)$	$T(c, p)$	Statement (I)	Statement (II)
true	true	false	false
true	false	true	true
false	true	false	false
false	false	false	false

Set join (**not all** case)

- Recall

$$\begin{aligned}\text{not all} &\leftrightarrow \{(s, p) \mid \exists c : c \in \bar{c}^s p\} \\ &\leftrightarrow \{(s, p) \mid \exists c : c \in T(p) \wedge \neg(c \in E(s))\} \\ &\leftrightarrow \{(s, p) \mid \exists c : T(c, p) \wedge \neg E(s, c)\}\end{aligned}$$

- We first need to find an RA expression for the set

$$\{(s, c, p) \mid T(c, p) \wedge \neg E(s, c)\}.$$

$$\begin{aligned}\{(s, c, p) \mid T(c, p) \wedge \neg E(s, c)\} &= \{(s, c, p) \mid (s \in S \wedge T(c, p)) \wedge \neg(E(s, c) \wedge T(c, p))\} \\ &= \{(s, c, p) \mid (s, c, p) \in (S \times T) \wedge \neg((s, c, p) \in E \bowtie T)\} \\ &= (S \times T) - (E \bowtie T)\end{aligned} \tag{1}$$

The equality (1) will be shown on the next slides.
Consequently, the RA expression for **not all** is

$$\pi_{sid,pid}((S \times T) - (E \bowtie T))$$

Set join (**not all** case using natural join)

- To prove equality (1), we need to show that the following statements are logically equivalent

$$T(c, p) \wedge \neg E(s, c)$$

$$(s \in \mathcal{S} \wedge T(c, p)) \wedge \neg(E(s, c) \wedge T(c, p))$$

Since, clearly, $s \in \mathcal{S}$ is true, we need to show that the following statements are equivalent:⁵

$$T(c, p) \wedge \neg E(s, c)$$

$$T(c, p) \wedge \neg(E(s, c) \wedge T(c, p))$$

⁵We assume that there is at least one student.

Set join (**not all** case using natural join)

We need to show that the following statements are equivalent:

$$T(c, p) \wedge \neg E(s, c) \quad (I)$$

$$T(c, p) \wedge \neg(E(s, c) \wedge T(c, p)) \quad (II)$$

To see this, you need to determine if statements (I) and (II) have the same truth values. This is the case as shown in the following truth table:

$E(s, c)$	$T(c, p)$	Statement (I)	Statement (II)
true	true	false	false
true	false	false	false
false	true	true	true
false	false	false	false

Set joins (not-empty cases)

- We have now developed the RA expressions for the following set joins:

$$\text{some} = \pi_{sid,pid}(E \bowtie T)$$

$$\text{not only} = \pi_{sid,pid}((E \times \mathcal{P}) - (E \bowtie T))$$

$$\text{not all} = \pi_{sid,pid}((S \times T) - (E \bowtie T))$$

- Notice that the time and space complexities for *not only* and *not all* are at least $O(|E| * |\mathcal{P}|)$ and $O(|T| * |S|)$, respectively. (At IU, $160000 * 2000$ and $6000 * 40000$.)
- These set joins can be very expensive to evaluate.
- Observe that the time and space complexity for $E \bowtie S$ is never more than $O(|E| + |S| + |S| * |\mathcal{P}|)^6$

⁶If we use the hash join algorithm.

Set joins (generalized intersection and difference)

A more compact way to remember these expressions is by using the following notations:

$$\begin{aligned}E \sqcap T &\leftrightarrow E \bowtie T \\E \ominus T &\leftrightarrow (E \times \mathcal{P}) - (E \bowtie T) \\T \ominus E &\leftrightarrow (\mathcal{S} \times T) - (E \bowtie T)\end{aligned}$$

With these notations, we get

$$\begin{aligned}\text{some} &= \pi_{sid,pid}(E \sqcap T) \\ \text{not only} &= \pi_{sid,pid}(E \ominus T) \\ \text{not all} &= \pi_{sid,pid}(T \ominus E)\end{aligned}$$

We call \sqcap the **generalized intersection** and \ominus the **generalized difference**.

Set joins (empty cases)

- Finding RA expressions for the empty cases is now easy: they consist of subtracting the corresponding not-empty cases from the set of all possible student-professor pairs (s, p) , i.e., the set $\mathcal{S} \times \mathcal{P}$.
- So, we get the following RA expressions for the set joins associated with **no**, **not only**, and **all**:

$$\begin{aligned}\text{no} &= \mathcal{S} \times \mathcal{P} - (\pi_{sid,pid}(E \sqcap T)) \\ \text{only} &= \mathcal{S} \times \mathcal{P} - (\pi_{sid,pid}(E \ominus T)) \\ \text{all} &= \mathcal{S} \times \mathcal{P} - (\pi_{sid,pid}(T \ominus E))\end{aligned}$$

Set joins (Summary)

The following are then the RA expressions for all of the cases considered:

some	=	$\pi_{sid,pid}(E \sqcap T)$
not only	=	$\pi_{sid,pid}(E \ominus T)$
not all	=	$\pi_{sid,pid}(T \ominus E)$
no	=	$\mathcal{S} \times \mathcal{P} - (\pi_{sid,pid}(E \sqcap T))$
only	=	$\mathcal{S} \times \mathcal{P} - (\pi_{sid,pid}(E \ominus T))$
all	=	$\mathcal{S} \times \mathcal{P} - (\pi_{sid,pid}(T \ominus E))$

Set joins **some** and **no** in SQL

- **some**

```
SELECT DISTINCT sid, pid
FROM Enroll NATURAL JOIN TaughtBy
```

- **no**

```
SELECT sid, pid
FROM Student CROSS JOIN Professor
EXCEPT
SELECT DISTINCT sid, pid
FROM Enroll NATURAL JOIN TaughtBy
```

Set joins **not only** and **only** in SQL

not only

```
SELECT  sid, pid
FROM    (SELECT sid, cno, pid
        FROM Enroll CROSS JOIN Professor
        EXCEPT
        SELECT sid, cno, pid
        FROM Enroll NATURAL JOIN TaughtBy) q
```

only

```
SELECT  sid, pid
FROM    Student      CROSS JOIN Professor
EXCEPT
SELECT  sid, pid
FROM    (SELECT  sid, cno, pid
        FROM    Enroll CROSS JOIN Professor
        EXCEPT
        SELECT  sid, cno, pid
        FROM    Enroll NATURAL JOIN TaughtBy) q
```

Set joins **not all** and **all** in SQL

not all

```
SELECT  sid, pid
FROM    (SELECT sid, cno, pid
        FROM Student CROSS JOIN TaughtBy
        EXCEPT
        SELECT sid, cno, pid
        FROM Enroll NATURAL JOIN TaughtBy) q
```

all

```
SELECT  sid, pid
FROM    Student      CROSS JOIN Professor
EXCEPT
SELECT  sid, pid
FROM    (SELECT  sid, cno, pid
        FROM    Student CROSS JOIN TaughtBy
        EXCEPT
        SELECT  sid, cno, pid
        FROM    Enroll NATURAL JOIN TaughtBy) q
```

The queries run in $O(|Student| * |TaughtBy| + |Enroll| * |TaughtBy|)$

Set joins **all** in Pure SQL

```
SELECT  sid, pid
FROM    Student
EXCEPT
SELECT  sid, pid
FROM    (SELECT
          FROM    Student CROSS JOIN TaughtBy
        EXCEPT
          SELECT
          FROM    Enroll NATURAL JOIN TaughtBy) q
```

Contrast this with the **all** set join expressed in Pure SQL

```
SELECT  sid, pid
FROM    Student s, Professor p
WHERE   NOT EXISTS (SELECT 1
                    FROM    TaughtBy t
                    WHERE t.pid = p.pid AND
                          t.cno NOT IN (SELECT e.cno
                                       FROM    Enroll e
                                       WHERE e.sid = s.sid))
```

This query runs in $O(|Student| * |Professor| * |TaughtBy| * |Enroll|)$ and is thus orders of magnitude slower than the above RA SQL query which runs in $O(|Student| * |TaughtBy| + |Enroll| * |TaughtBy|)$.

Set joins (Example)

Consider the following queries:

"Find the student sid pairs (s_1, s_2) such that student s_1 takes [some|not only|not all|no|only|all] course(s) that student s_2 takes."

The queries are very similar to those we have been considering.

The only difference is that we now have to involve the *Enroll* relation twice: ones for student s_1 and the other for student s_2 .

Set joins (Example)

Let us denote by \tilde{E} the RA expression $\rho_{sid \rightarrow \tilde{sid}}(\pi_{cno, sid}(E))$. So \tilde{E} is the same as E , except that its schema is (cno, \tilde{sid}) instead of (sid, cno) . We have renamed sid by \tilde{sid} .

For this example, the **search space** is the set of all student pairs $\mathcal{S} \times \mathcal{S}$ and we assume its schema is (sid, \tilde{sid}) .

Then the RA expressions for the above six queries are

some	=	$\pi_{sid, \tilde{sid}}(E \bowtie \tilde{E})$
not only	=	$\pi_{sid, \tilde{sid}}(E \ominus \tilde{E})$
not all	=	$\pi_{sid, \tilde{sid}}(\tilde{E} \ominus E)$
no	=	$\mathcal{S} \times \mathcal{S} - (\pi_{sid, \tilde{sid}}(E \bowtie \tilde{E}))$
only	=	$\mathcal{S} \times \mathcal{S} - (\pi_{sid, \tilde{sid}}(E \ominus \tilde{E}))$
all	=	$\mathcal{S} \times \mathcal{S} - (\pi_{sid, \tilde{sid}}(\tilde{E} \ominus E))$

More set joins

- The RA region expressions $E \sqcap T$, $E \ominus T$, and $T \ominus E$ permit us to express many more set joins. We will consider one such example next.
- "Find the student-professor pairs (s, p) such that student s takes at least 2 courses not taught by professor p ."
- More formally, we need to find a RA expression for the query

$$\{(s, p) \mid \exists c_1, c_2 : c_1 \neq c_2 \wedge c_1 \in (E(s) - T(p)) \wedge c_2 \in (E(s) - T(p))\}$$

- Alternatively,

$$\{(s, p) \mid \exists c_1, c_2 : c_1 \neq c_2 \wedge \{c_1, c_2\} \subseteq (E(s) - T(p))\}$$

More set joins

- "Find the student-professor pairs (s, p) such that student s takes at least 2 courses not taught by professor p ."

$$\{(s, p) \mid \exists c_1, c_2 : c_1 \neq c_2 \wedge \{c_1, c_2\} \subseteq (E(s) - T(p))\}$$

- From this formulation, it is clear that we need to involve the region expression $E \ominus T$, but we need two copies of it, one to reason about c_1 and the other to reason about c_2 .
- Let these two copies be $E_1 \ominus T_1$ with schema (sid, cno_1, pid) and $E_2 \ominus T_2$ with schema (sid, cno_2, pid) .
- Then the RA expression for the query is as follows:

$$\pi_{sid, pid}(\sigma_{cno_1 \neq cno_2}((E_1 \ominus T_1) \bowtie (E_2 \ominus T_2)))$$

Set semijoins (Motivation)

- Consider the relations $\text{Enroll}(\text{sid}, \text{cno})$ and the relation of CS courses $\text{CS}(\text{cno})$.
- We are interested in students satisfying some property.
- "Find each student sid s such that student s takes *[some|not only|not all|no|only|all]* CS course(s)."

Set semijoins (Solution)

An analysis very similar to that for set joins gives the following RA queries for these queries.

First observe that $E \bowtie CS = E \times CS$.

some	=	$\pi_{sid}(E \bowtie CS)$	=	$\pi_{sid}(E \times CS)$
not only	=	$\pi_{sid}(E \ominus CS)$		
not all	=	$\pi_{sid}(CS \ominus E)$		
no	=	$S - (\pi_{sid}(E \times CS))$		
only	=	$S - (\pi_{sid}(E \ominus CS))$		
all	=	$S - (\pi_{sid}(CS \ominus E))$		

Where

$E \ominus CS$	=	$E - (E \times CS)$	=	$E \overline{\times} CS$
$CS \ominus E$	=	$(S \times CS) - (E \times CS)$		

Set semijoins (Efficiency consideration)

Reconsider the 6 semijoins:

some	=	$\pi_{sid}(E \bowtie CS)$
not only	=	$\pi_{sid}(E \overline{\bowtie} CS)$
not all	=	$\pi_{sid}((\pi_{sid}(S) \times CS) - (E \bowtie CS))$
no	=	$\pi_{sid}(S) - (\pi_{sid}(E \bowtie CS))$
only	=	$\pi_{sid}(S) - \pi_{sid}(E \overline{\bowtie} CS)$
all	=	$\pi_{sid}(S) - \pi_{sid}((\pi_{sid}(S) \times CS) - (E \bowtie CS))$

Let $k = |S|$, $l = |CS|$, and $n = |E|$.

- Then, using hashing methods, the **some**, **no**, **not only**, **only** set semijoins can be implemented in **linear time and space** $O(|k| + |l| + |n|)$. (I.e., extremely fast.)
- However, the **not all** and **all** set semijoins require **quadratic time and space** $O((k * l) + n)$ because of the need to compute $\pi_{sid}(S) \times CS$. (I.e., can be extremely expensive.)

Set semijoin **some** and **no** in SQL

- **some**

```
SELECT DISTINCT sid
FROM Enroll NATURAL JOIN CS
```

- **no**

```
SELECT sid
FROM Student
EXCEPT
SELECT DISTINCT sid
FROM Enroll NATURAL JOIN CS
```

Set semijoin **not only** in SQL

- **not only**

```
SELECT DISTINCT sid
FROM (SELECT sid, cno
      FROM Enroll
      EXCEPT
      SELECT sid, cno
      FROM Enroll NATURAL JOIN CS) q
```

Set semijoin **only** SQL

- **only**

```
SELECT  sid
FROM    Student
EXCEPT
SELECT  DISTINCT sid
FROM    (SELECT sid, cno
        FROM    Enroll
        EXCEPT
        SELECT sid, cno
        FROM    Enroll NATURAL JOIN CS) q
```


Set semijoin **not all** in SQL

- **not all**

```
SELECT DISTINCT sid
FROM (SELECT sid, cno
      FROM (SELECT sid from Student) s CROSS JOIN CS
      EXCEPT
      SELECT sid, cno
      FROM Enroll NATURAL JOIN CS) q
```

Set semijoin **all** in SQL

- **all**

```
SELECT  sid
FROM    Student s
EXCEPT
SELECT  DISTINCT sid
FROM    (SELECT sid, cno
        FROM  (SELECT sid from Student) s CROSS JOIN CS
        EXCEPT
        SELECT sid, cno
        FROM  Enroll NATURAL JOIN CS) q
```

Set joins (General case)

- Let $E_1(A_1, \dots, A_m, C_1, \dots, C_k)$ and $E_2(C_1, \dots, C_k, B_1, \dots, B_n)$ be RA expressions.
- We assume that each attribute has an associated domain: these are denoted $\mathcal{A}_1, \dots, \mathcal{A}_m, \mathcal{C}_1, \dots, \mathcal{C}_k, \mathcal{B}_1, \dots, \mathcal{B}_n$.
- We can then build RA expressions $E_1 \sqcap E_2$, $E_1 \ominus E_2$, and $E_2 \ominus E_1$ as follows:

$$E_1 \sqcap E_2 = E_1 \bowtie E_2$$

$$E_1 \ominus E_2 = (E_1 \times \mathcal{B}_1 \times \dots \times \mathcal{B}_n) - (E_1 \bowtie E_2)$$

$$E_2 \ominus E_1 = (\mathcal{A}_1 \times \dots \times \mathcal{A}_m \times E_2) - (E_1 \bowtie E_2)$$

- Using these region expressions we can build numerous other set-joins.

Set joins (From general case to special cases)

$$E_1 \sqcap E_2 = E_1 \bowtie E_2$$

$$E_1 \ominus E_2 = (E_1 \times \mathcal{B}_1 \times \cdots \times \mathcal{B}_n) - (E_1 \bowtie E_2)$$

$$E_2 \ominus E_1 = (\mathcal{A}_1 \times \cdots \times \mathcal{A}_m \times E_2) - (E_1 \bowtie E_2)$$

We can now consider special cases:

- When $m = k = n = 1$ we have the region expressions for the set joins we have seen above.
- When $m = k = 1$ and $n = 0$ we have the region expressions for semi set joins we have seen above.
- When $m = n = 0$, then

$$E_1 \sqcap E_2 = E_1 \cap E_2$$

$$E_1 \ominus E_2 = E_1 - E_2$$

$$E_2 \ominus E_1 = E_2 - E_1$$