# CSE 4125: Distributed Database Systems
# Chapter – 6
# (Part – D)

Optimization of Access Strategies.

# Semi-join Programs (Join using Semi-join)
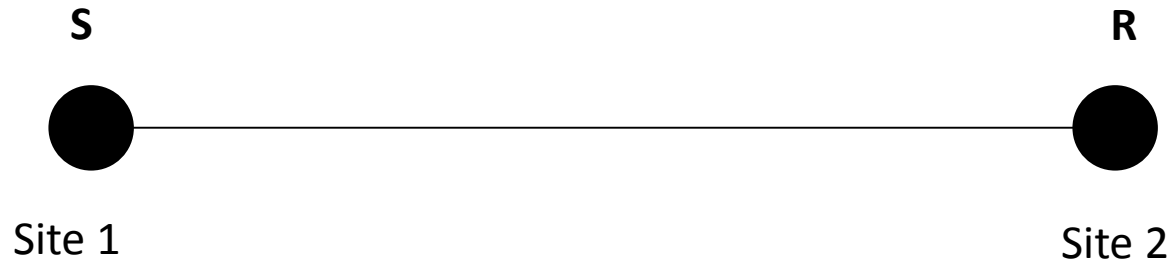
# Semi-join Programs

$$R \ \text{JN}_{C=A} \ S \leftrightarrow (R \ \text{SJ}_{C=A} \ \text{PJ}_A \ S) \ \text{JN}_{C=A} \ S$$

| A | B |
|---|---|
| 1 | 4 |
| 2 | 5 |
| 3 | 6 |
| 3 | 7 |

| C | D |
|---|---|
| 1 | 6 |
| 2 | 7 |
| 4 | 8 |
| 5 | 9 |

*Assume, $C_0 = 0$, $C_1 = 1$, $\rho = 0.2$*
*Size(A) = Size(B) = Size(C) = Size(D) = 2 bytes*
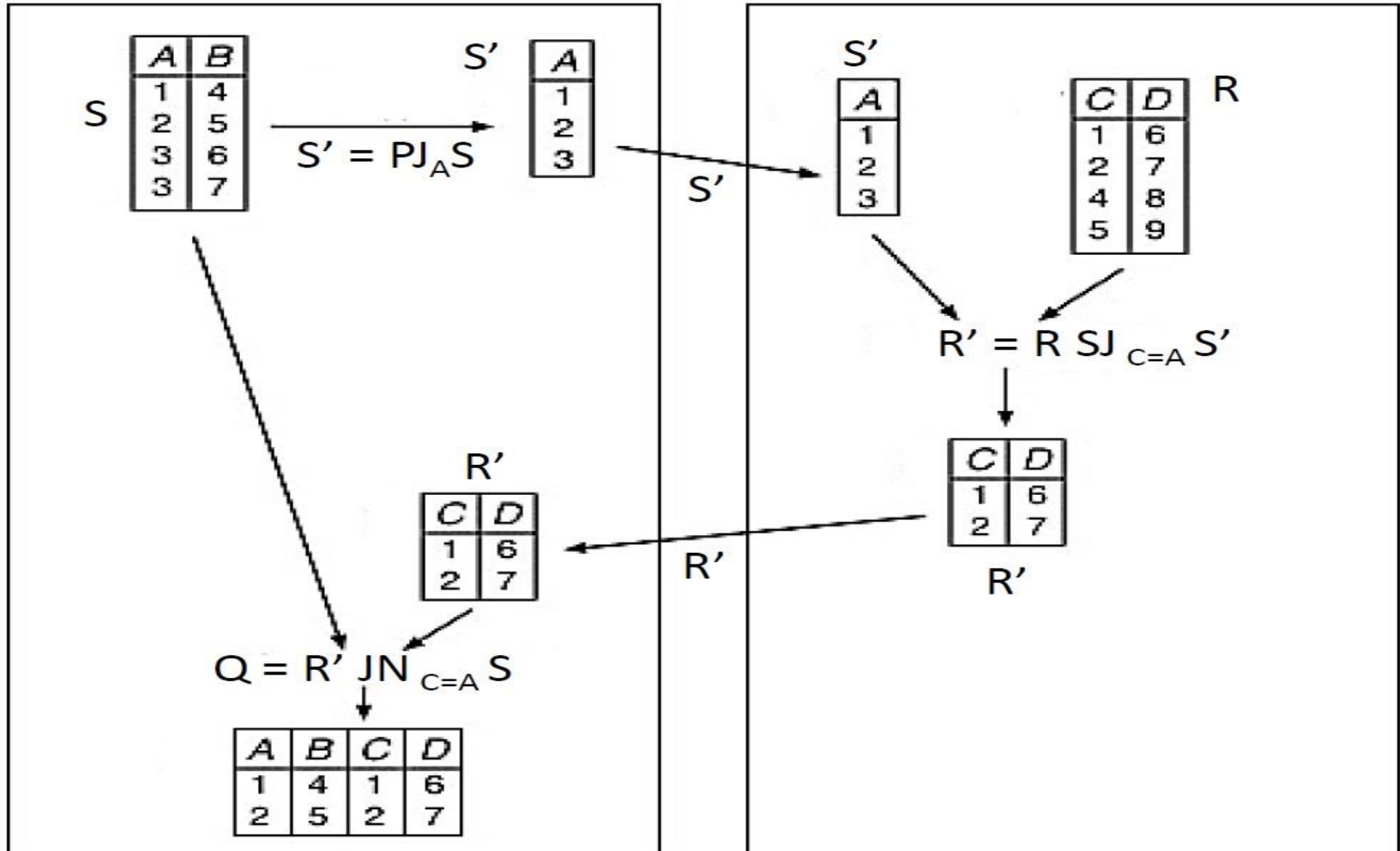
**S**

**R**

Site 1

Site 2

We want to perform $R \ \text{JN}_{C=A} \ S$ at site − 1 using Semi-Join Program.

# Semi-join Programs (cont.)

$$R \text{ JN}_{C=A} S \leftrightarrow (R \text{ SJ}_{C=A} \text{PJ}_A S) \text{ JN}_{C=A} S$$
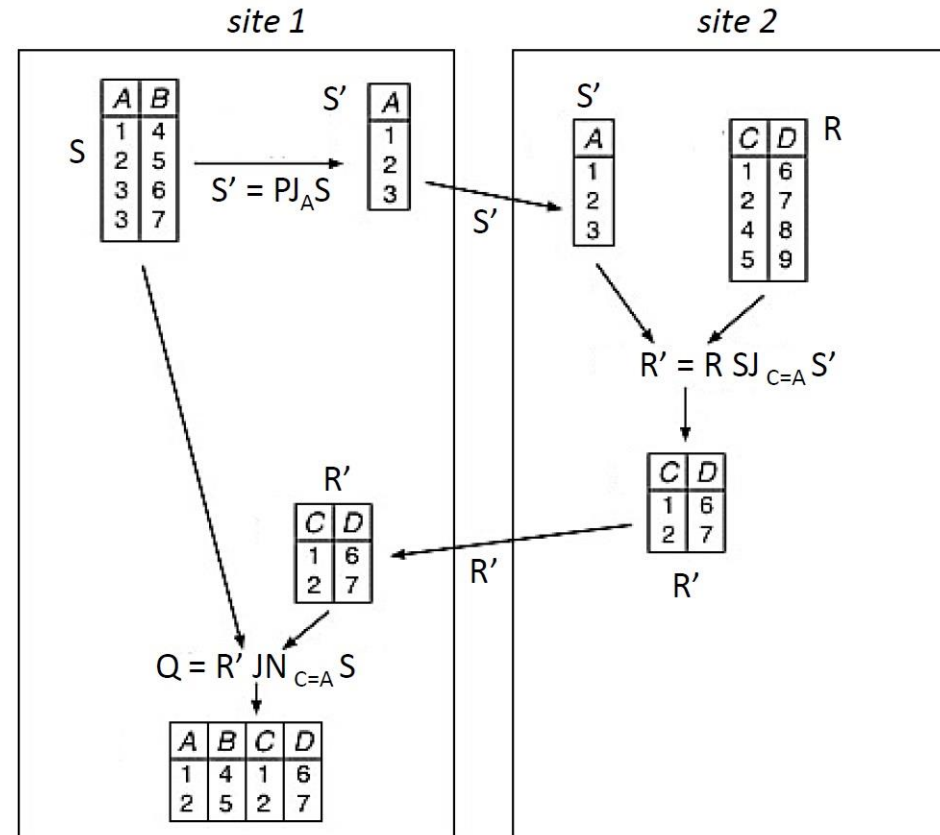
# Semi-join Programs (cont.)

## Steps of Semi-join program

1. Send $S' = PJ_A(S)$ to **site-2**
2. Compute $R' = R \ SJ_{C=A} \ S'$ at **site-2**
3. Send $R'$ to **site-1**
4. Compute $Q = R' \ JN_{C=A} \ S$ at **site-1**

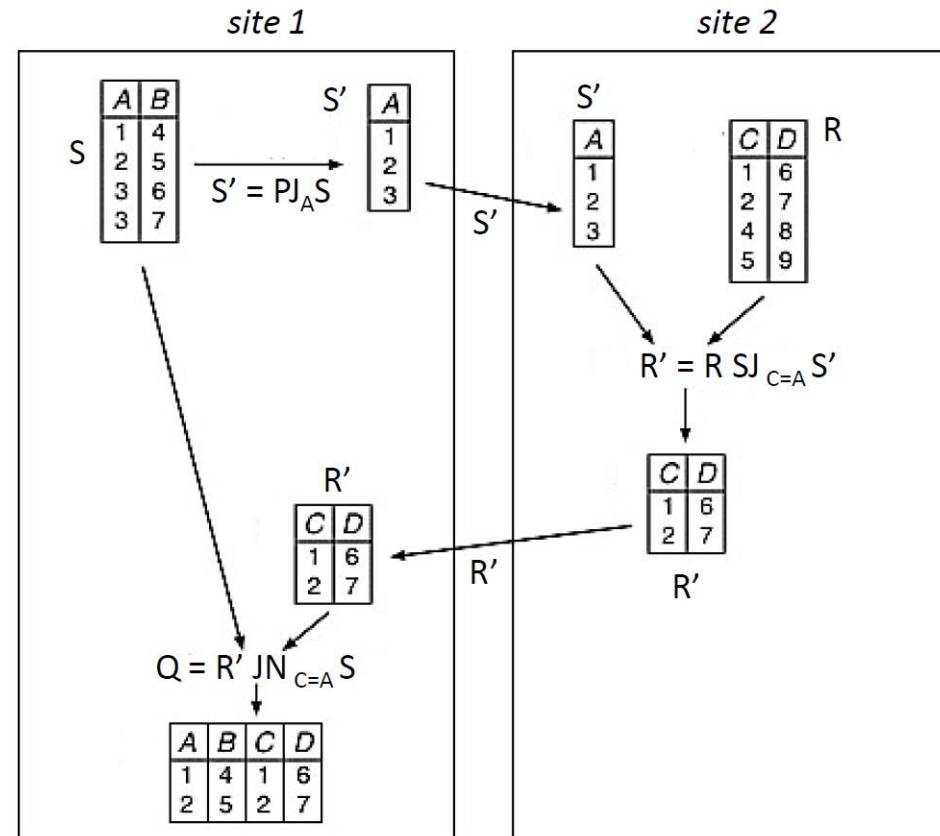*Task*: What will be the steps if we want to perform the join at site 2?

# Semi-join Programs (cont.)

## **Cost of Semi-join program**

Step 1: Send $S' = PJ_A (S)$ to **site-2**

$$TC_1 = C_0 + C_1 * x$$
$$= C_0 + C_1 * Card(S') * size(S')$$
$$= C_0 + C_1 * val(A[S]) * size(A)$$
$$= 0 + 1 * 3 * 2 \text{ bytes}$$
$$= 6 \text{ bytes}$$
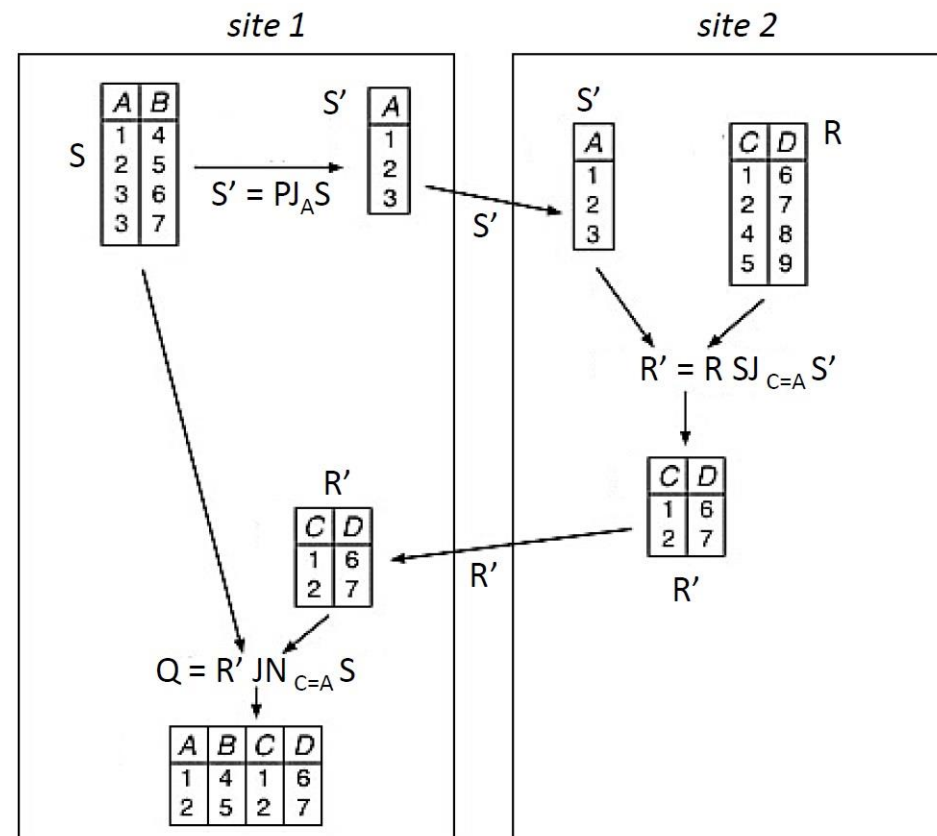$$= 6 * 8 \text{ bits}$$
$$= 48 \text{ bits}$$

# Semi-join Programs (cont.)

## **Cost of Semi-join program**

Step 2: Compute **R' = R JN$_{C=A}$ S'** at **site-2**
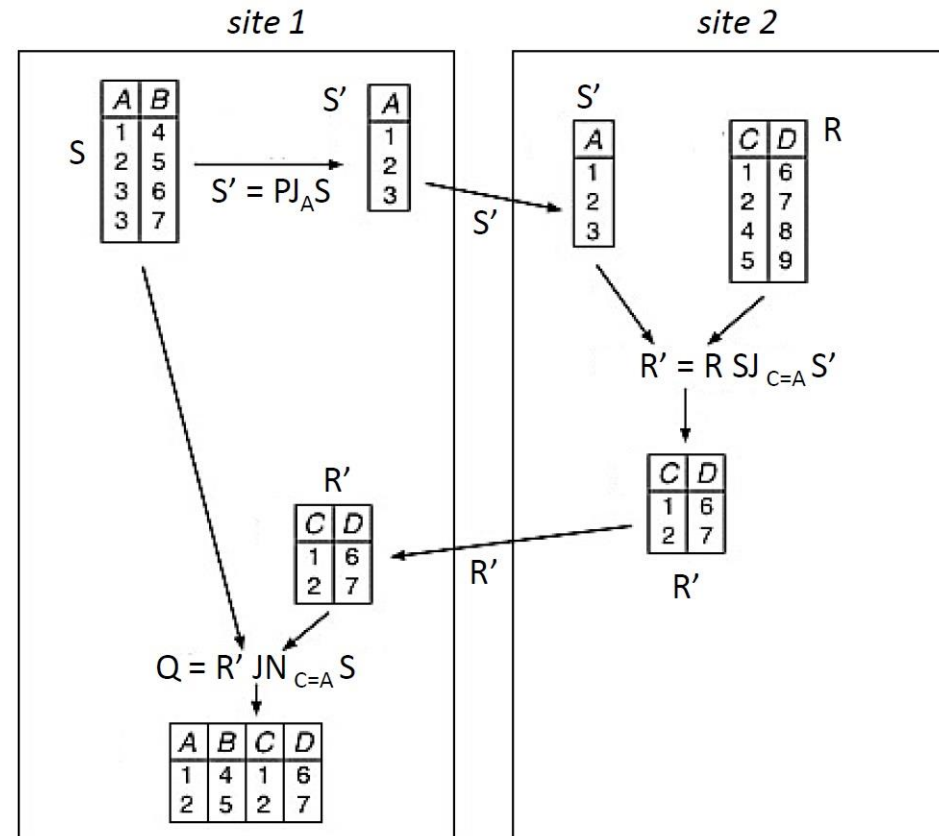
$TC_2 = 0$

# Semi-join Programs (cont.)

## Cost of Semi-join program

Step 3: Send **R'** to **site-1**

$\text{TC}_3 = C_0 + C_1 * x$

$\qquad = C_0 + C_1 * Card(R') * size(R')$

$\qquad = C_0 + C_1 * \rho * Card(R) * size(R)$

$\qquad = 0 + 1 * 0.2 * 4 * 4 \ bytes$

$\qquad = 3.2 \ \text{bytes}$

$\qquad = 3.2 * 8 \ \text{bits}$

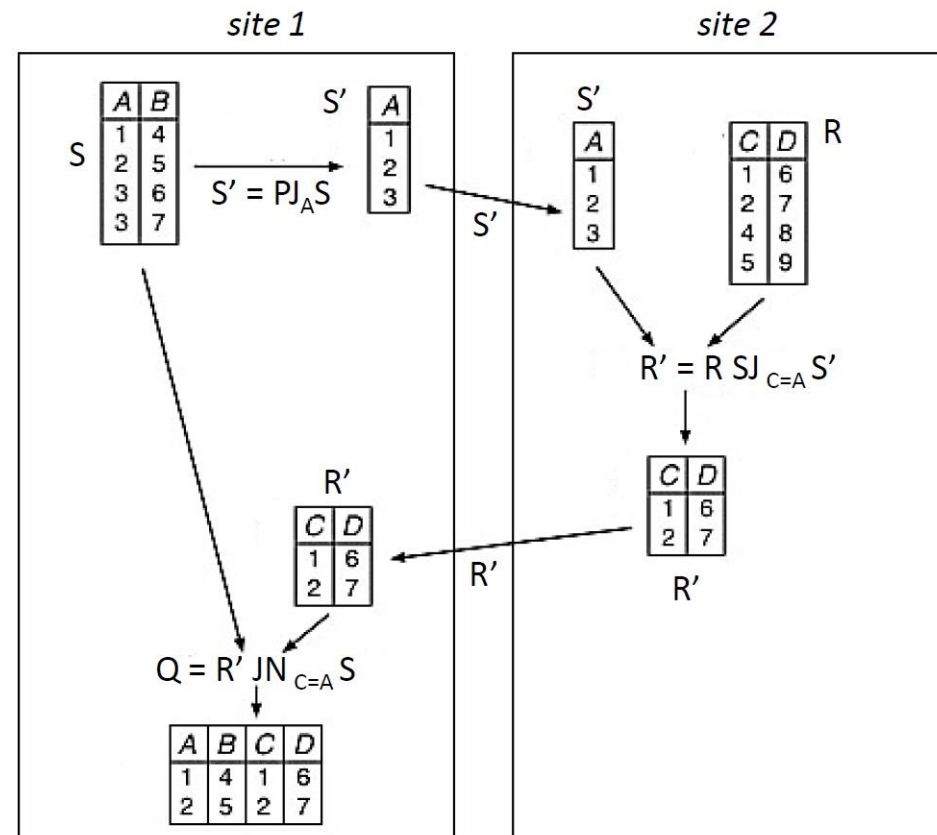$\qquad = 25.6 \ \text{bits}$

$\qquad \approx 26 \ bits$

# Semi-join Programs (cont.)

**<u>Cost of Semi-join program</u>**

Step 4: Compute $\mathbf{Q = R' \ JN_{C=A} \ S}$ at **site-1**

$TC_4 = 0$

# Semi-join Programs (cont.)
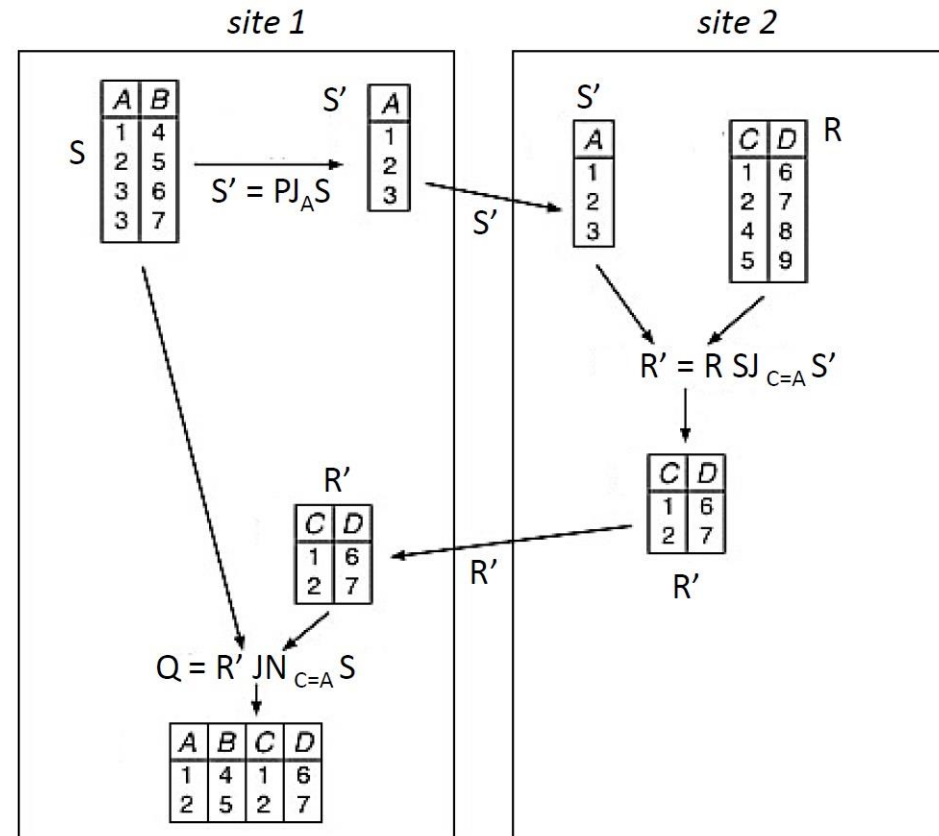
**Cost of Semi-join program**

Total cost

$$TC_{SJ} = TC_1 + TC_2 + TC_3 + TC_4 = 48 + 0 + 26 + 0 = 74 \text{ bits}$$

If $TC_{SJ} < TC_{JN}$ then semi-join program is profitable.

# Semi-join Programs (cont.)

## Cost **without** Semi-join program

$$TC_{JN} = C_0 + C_1 * x$$
$$= C_0 + C_1 * Card(R) * size(R)$$
$$= 0 + 1 * 4 * 4 \text{ bytes}$$
$$= 16 \text{ bytes}$$
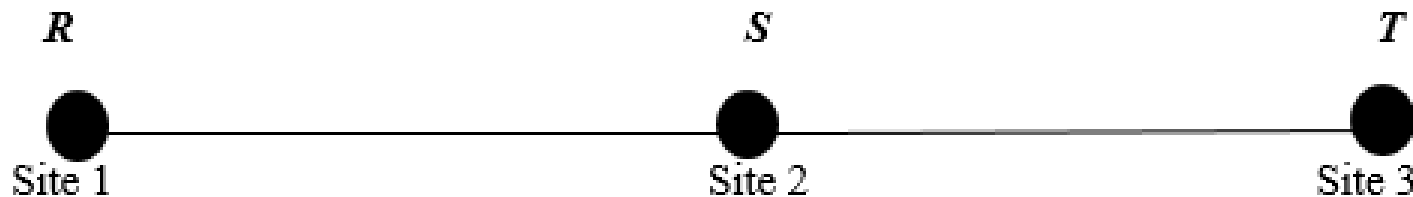$$= 16 * 8 \text{ bits}$$
$$= 128 \text{ bits}$$

# Other Applications of Semi-join Programs

❑ Semi-join programs can be used as fragment reducers (operations that can reduce cardinality of a relation).
– Similarly to unary operations.

❑ Full reducer:
– Chain of semi-joins.

# Exercise

Consider the following distributed database with relations **R**, **S** and **T** over a network of site 1, 2 and 3.

$$R \qquad\qquad\qquad\qquad S \qquad\qquad\qquad\qquad T$$

Site 1 ———————————— Site 2 ———————————— Site 3

Assume the following specifications are given.

$$C_0^{12} = C_0^{23} = C_1^{21} = C_1^{32} = 0 \; unit$$

$$C_1^{12} = C_1^{23} = C_0^{21} = C_0^{32} = 1 \; unit$$

**size (R)** = 20 bytes, **size(T)** = 20 bytes, **size(S)** = 40 bytes, **size(a)** = **size(b)** = 1 byte

**card(R)** = 100, **card(S)** = 50, **card(T)** = 50

val(a[R]) = val(b[S]) = val(a[T]) = 50

$$R \; SJ_{a \, = \, b}S \; has \; selectivity \; \rho = \mathbf{0.1}$$

$S\ SJ_{b\,=\,a}R$ *has selectivity* $\rho = 0.9$

$T\ SJ_{a\,=\,b}S$ *has selectivity* $\rho = 0.5$

$S\ SJ_{b\,=\,a}T$ *has selectivity* $\rho = 0.5$

Determine the total transmission cost of performing $(R\ JN_{a\,=\,b}\ S)\ DF\ (T\ JN_{a\,=\,b}\ S)$ at site **2** using *semi-join program only.* $[\ C^{xy}\ means\ transimission\ cost\ from\ site\ x\ to\ site\ y\ ]$