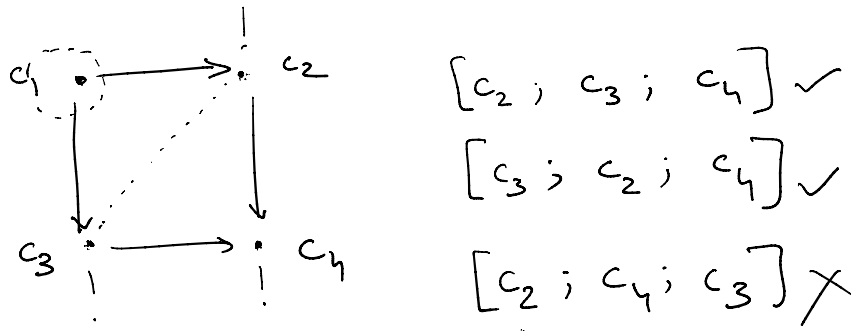


# Homework 3



key  $+3$  : ATTACK AT DAWN ← Plaintext  
 DWDFN DWGDZA ← Cipher text } Caesar cipher

+ key  $+3$   
 ATTACK AT DAWN  
 LEMON LEMONLE  
 LY... } Vigenere cipher

## How to evaluate Relational Algebra

- How to evaluate Cartesian products

To evaluate :  $Q_1 * Q_2$

Let  $T_1 = \text{eval}(Q_1)$

$T_2 = \text{eval}(Q_2)$

Create a new table  $T_3$

For each  $t_1 \in T_1$   $t_2 \in T_2$

Insert  $(t_1 @ t_2)$  into  $T_3$

Return  $T_3$

Employees  $\bowtie_{E.z=U.z}$  USPS

$= \sigma_{E.z=U.z} (\text{Employees} \times \text{USPS})$

9 intermediate tuples.

Name	zip code		zip code	state
John	90089		90089	CA
David	90015		90015	CA
Mary	19104		19104	PA

To evaluate  $\text{Emp} \bowtie_{E.z=U.z}$  USPS

Name	zip		Z	S
J	90089	→	90089	CA
D	90015	→	90015	CA
M	19104	→	19104	PA

With 20000 employees  
20000 zip codes } → 40 billion intermediate tuples

20000 zip codes

Intermediate tuples  
with Cartesian products

~ 40 seconds  
on a modern machine

With targeted joins

Create 20000 tuples.

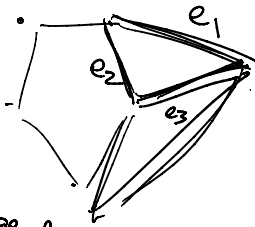
$\frac{20000 * \log(20000)}{\# \text{ of employees}}$  operations to lookup zip code

~ 300 K operations

≤ 1 second.

Question: How do we evaluate joins quickly?

### Triangle Query



select  $n_1$   $n_2$   $n_3$

from edge  $e_1$  edge  $e_2$  edge  $e_3$

where  $e_1 \cdot v_1 = e_2 \cdot v_2$  and

$e_2 \cdot v_1 = e_3 \cdot v_2$  and

$e_3 \cdot v_1 = e_1 \cdot v_2$

$n_1 = e_1 \cdot v_1$ ,  $n_2 = e_2 \cdot v_1$ ,  $n_3 = e_3 \cdot v_1$

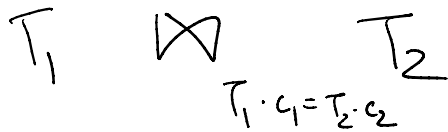
There exists algorithm which takes  $O(n^2)$  time

There exist tables where you output

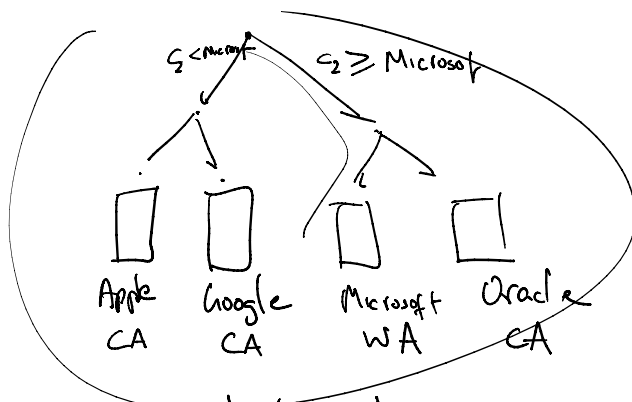
$O(n^{1.7})$  tuples

Loop Nest Train

# Loop Nest Join



- Construct an index on  $T_2.c_2$

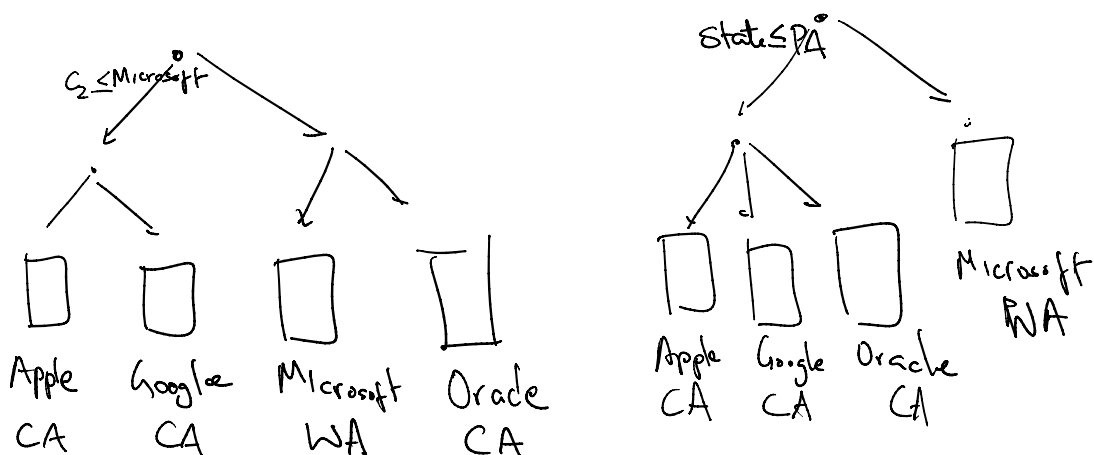


- Loop over all tuples  $t_1 \in T_1$

- Consult the column  $t_1.c_1$

- Find if a corresponding entry exists in  $T_2$

-  $O(|T_1| \cdot \log(|T_2|) + |T_2|)$  ↖ To construct the index





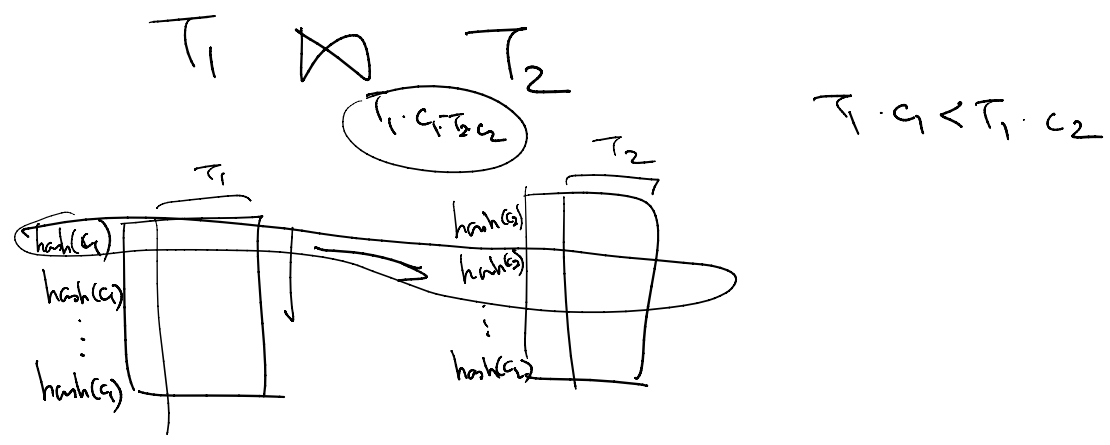
Apple CA    Google CA    Microsoft WA    Oracle CA

CA    CA    CA

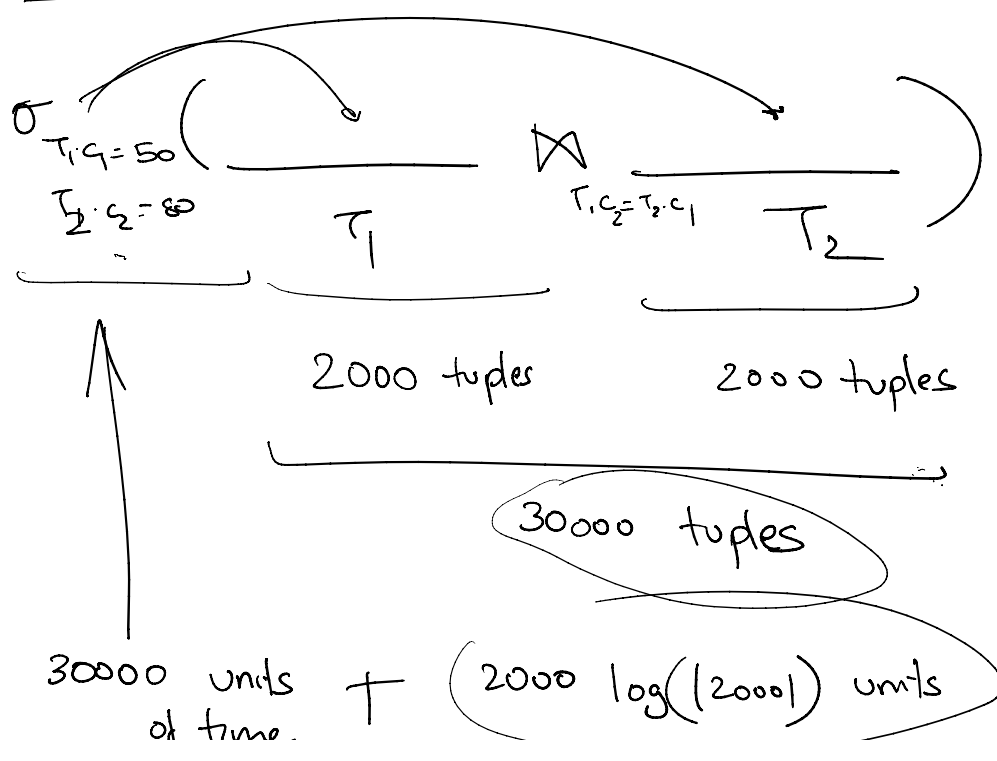
"All companies named Microsoft"

"All companies headquartered in WA"

### Hash join



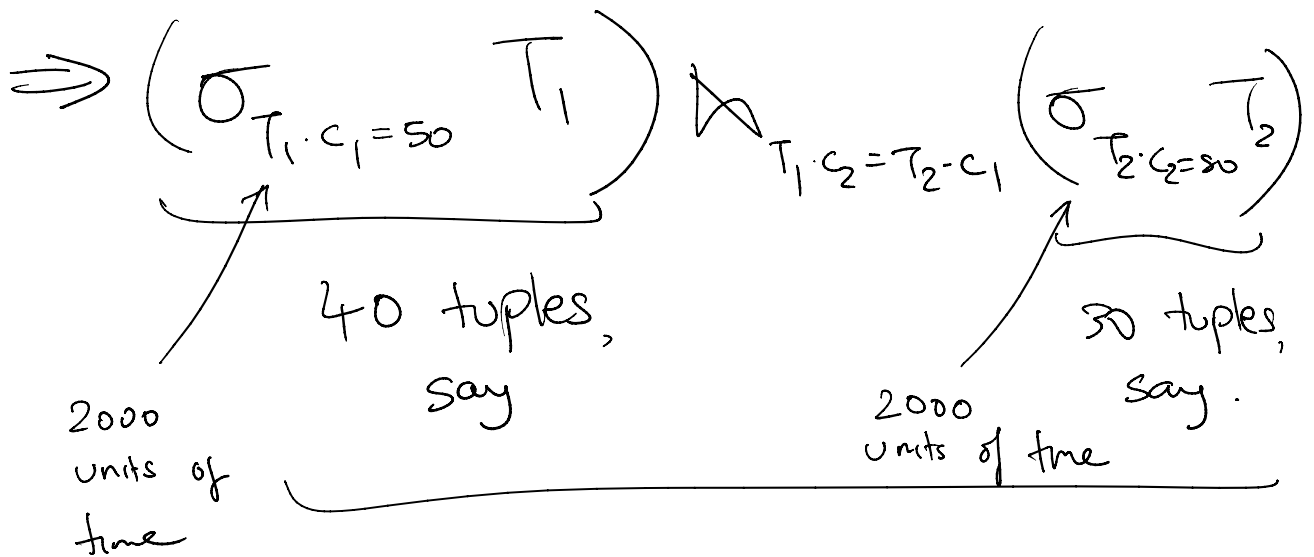
### Two optimizations



30000 units of time +  $(2000 \log(2000))$  units

"Filter them first"

$$\sigma_{\substack{T_1 \cdot c_1 = 50 \\ T_2 \cdot c_2 = 80}} (T_1 \bowtie_{T_1 \cdot c_2 = T_2 \cdot c_1} T_2)$$



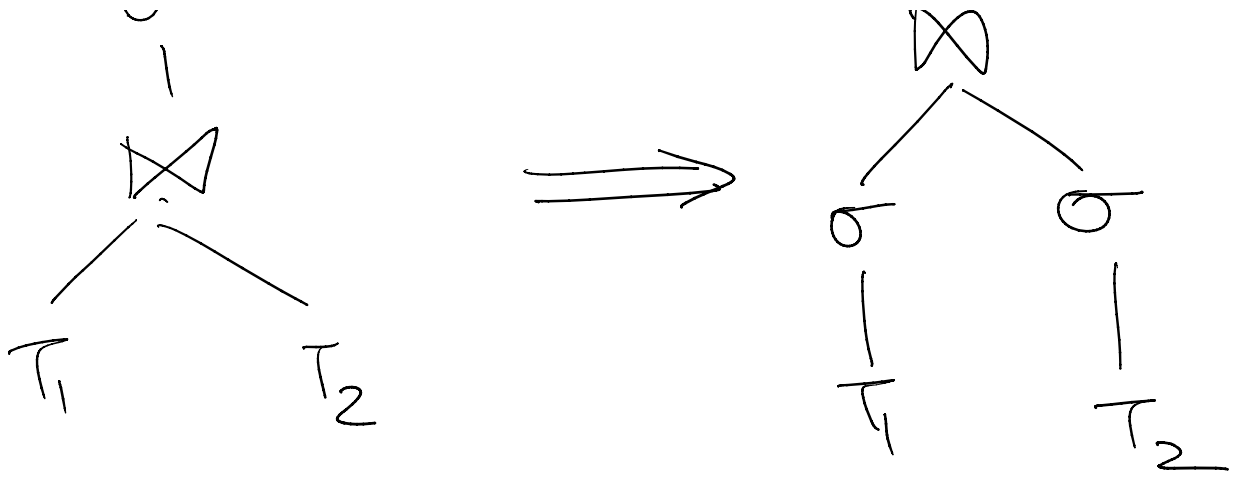
1200  $\approx$  1500 operations.

5500 operations

"Selection push down" / Distributivity.

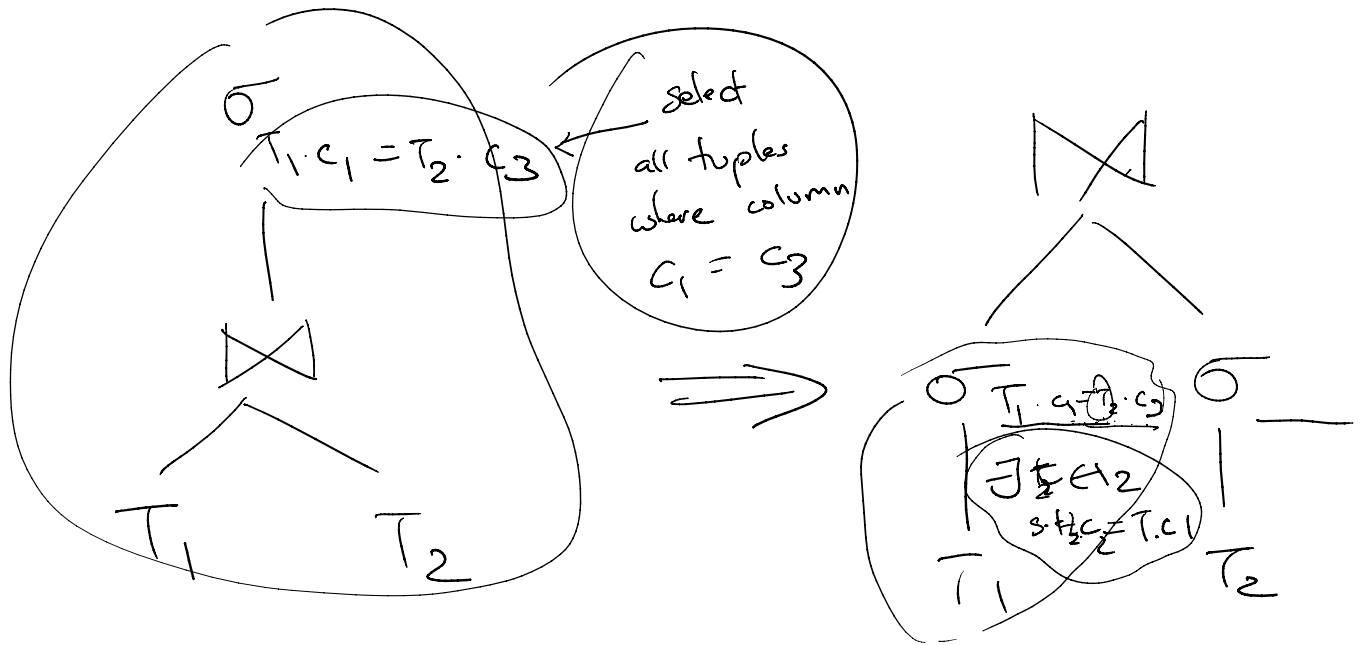
$\sigma$   
|

$\bowtie$   
|

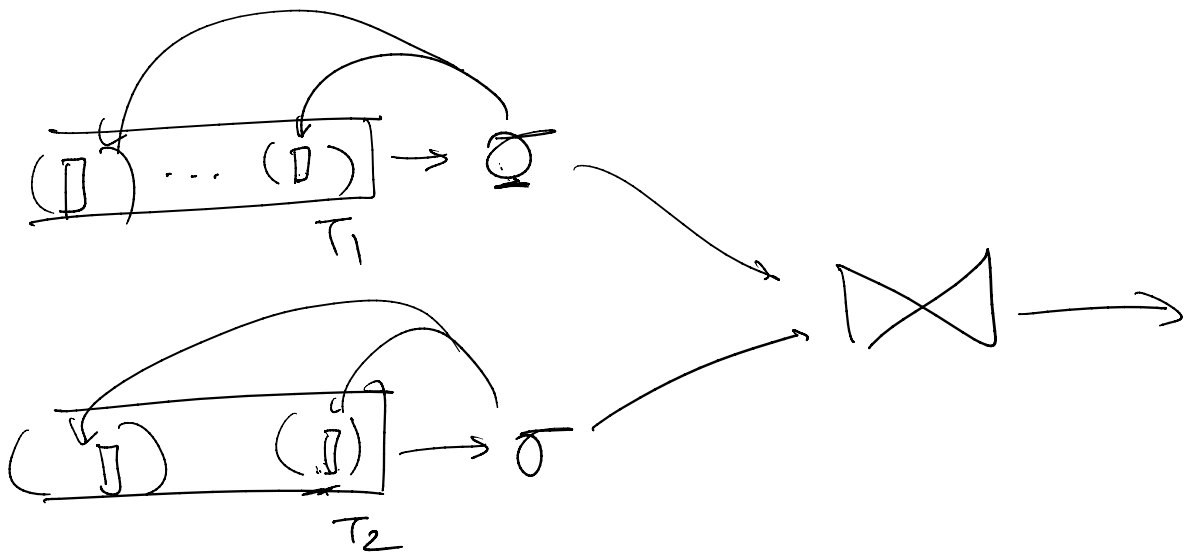
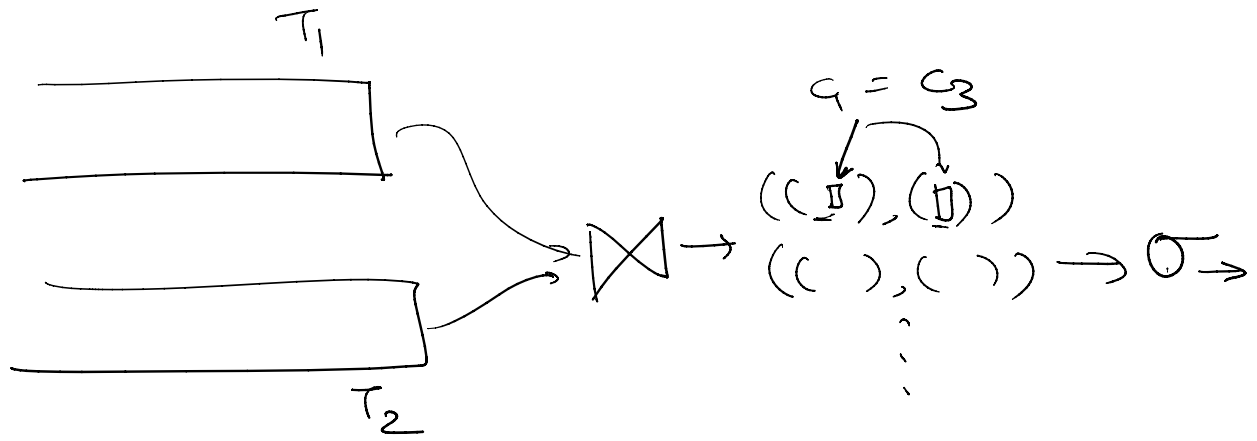


Question 1: Why does this <sup>performance improvement</sup> optimization work?

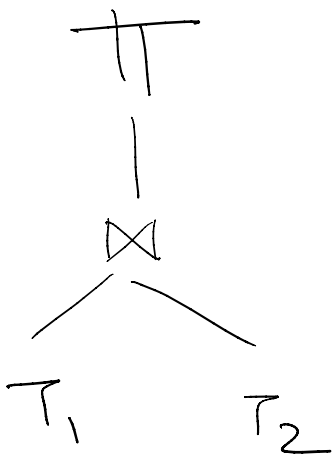
Question 2: Does it always make sense?



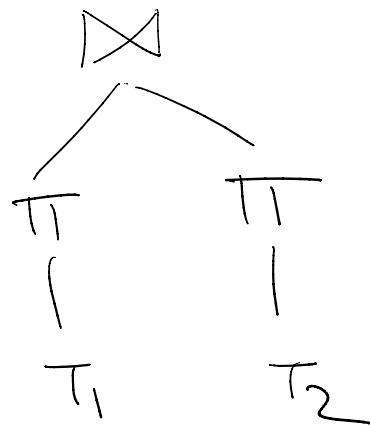
Selection pushdown works only if it is locally filtering the elements of the table.



## Projection pushdowns



$\implies$   
only if  
projection  
preserves all  
" "

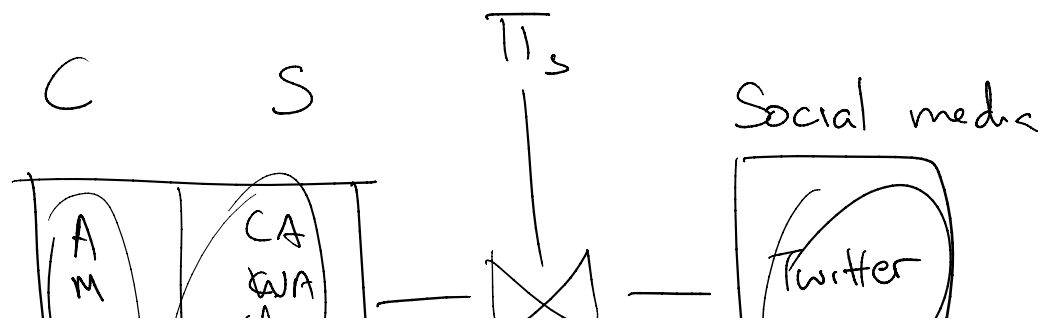
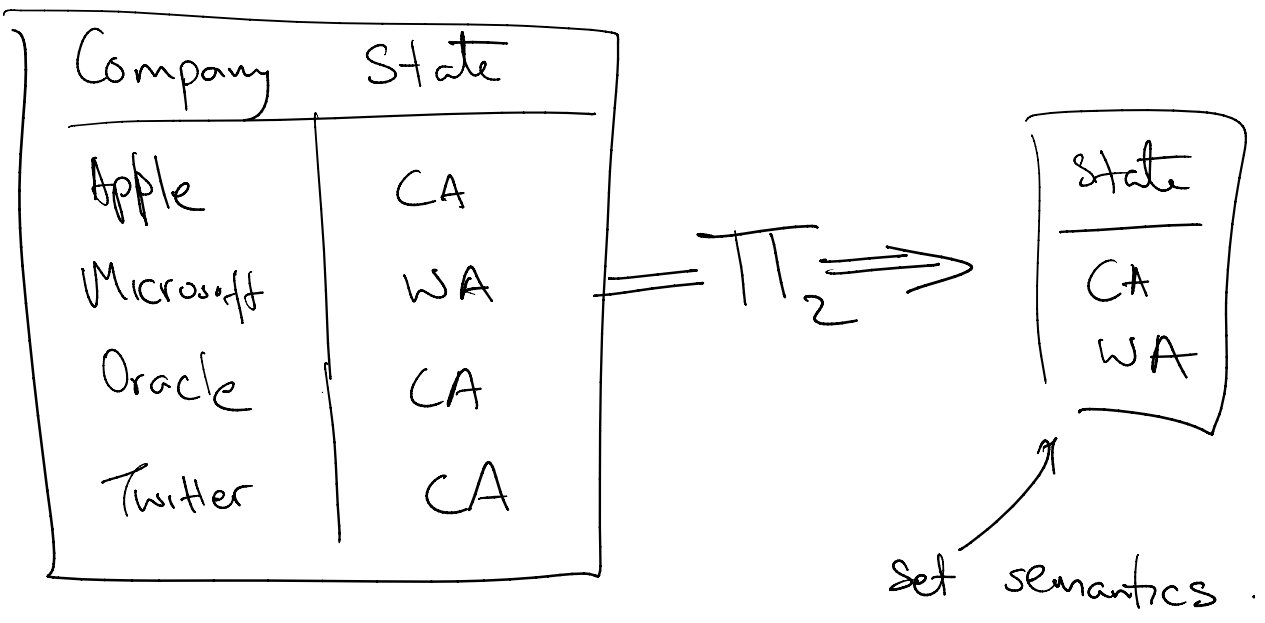


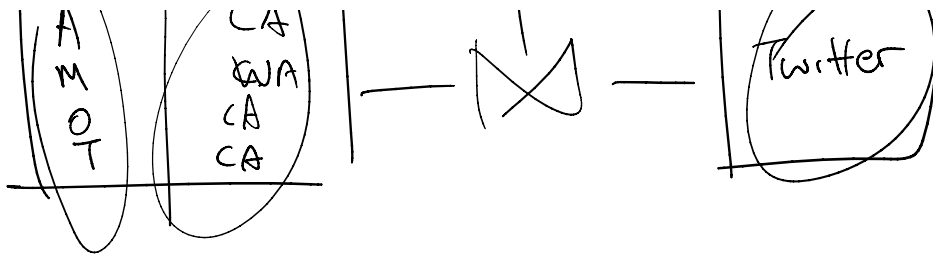
$\tau_1$        $\tau_2$       preserves all columns needed for join.       $\tau_1$        $\tau_2$       reduces memory.

Question 1: Is this a performance improvement?

Question 2: Does it always make sense. (correctness)

Question 1a: Does projection pushdown preserve the number of tuples?





"All states which have the headquarters of a social network."

## Things not discussed

1. Query planning: which operations to execute & in what order

$$(T_1 \bowtie T_2 \bowtie T_3 \dots \bowtie T_k)$$

$$\overbrace{T_1 \bowtie T_2}^{2000} \bowtie T_3$$

2000                      50

2. Rewriting strategies: selection pushdowns  
projection pushdowns.

# "Declarative programming"

---

How do you describe reachability in a graph?  
edge

