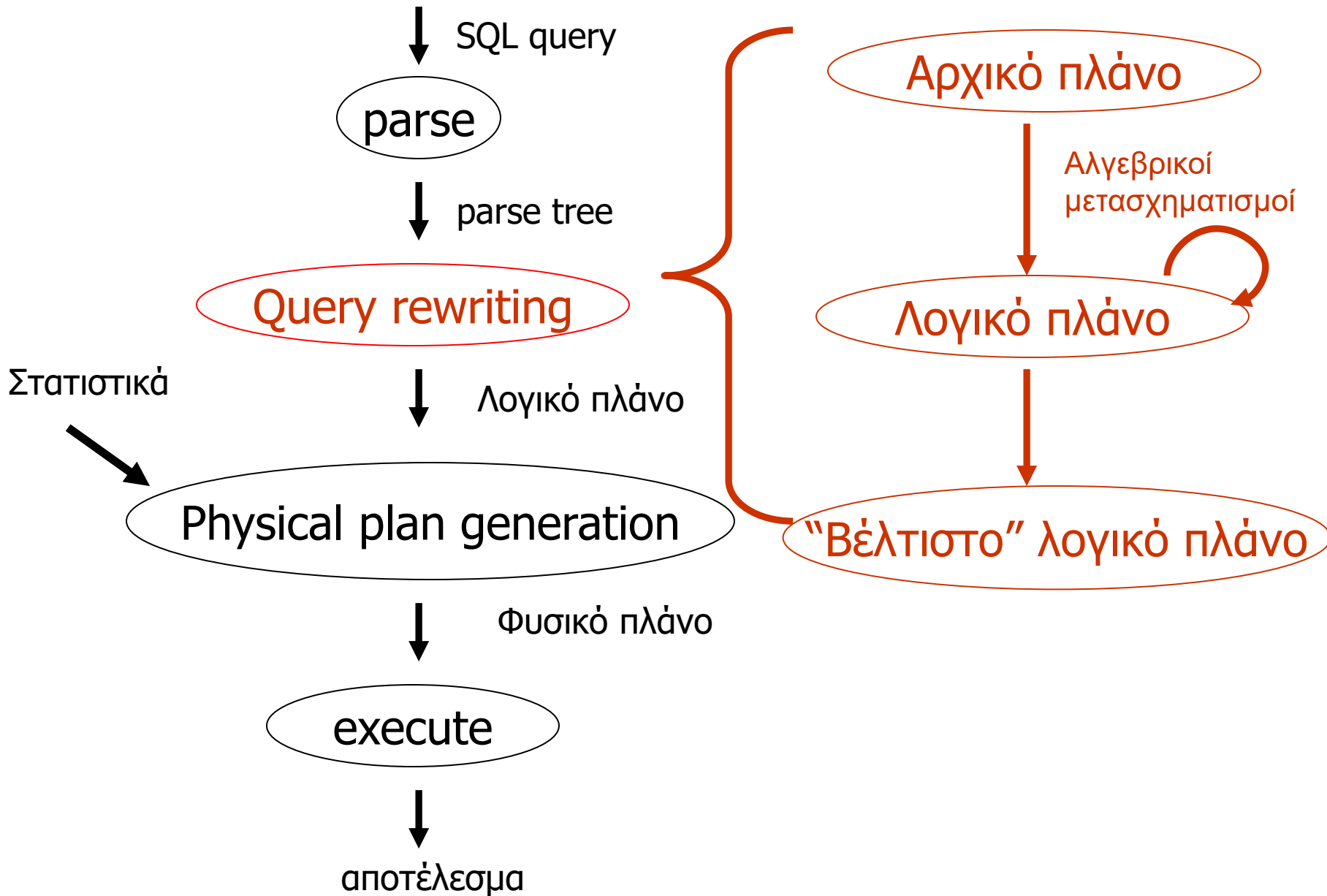
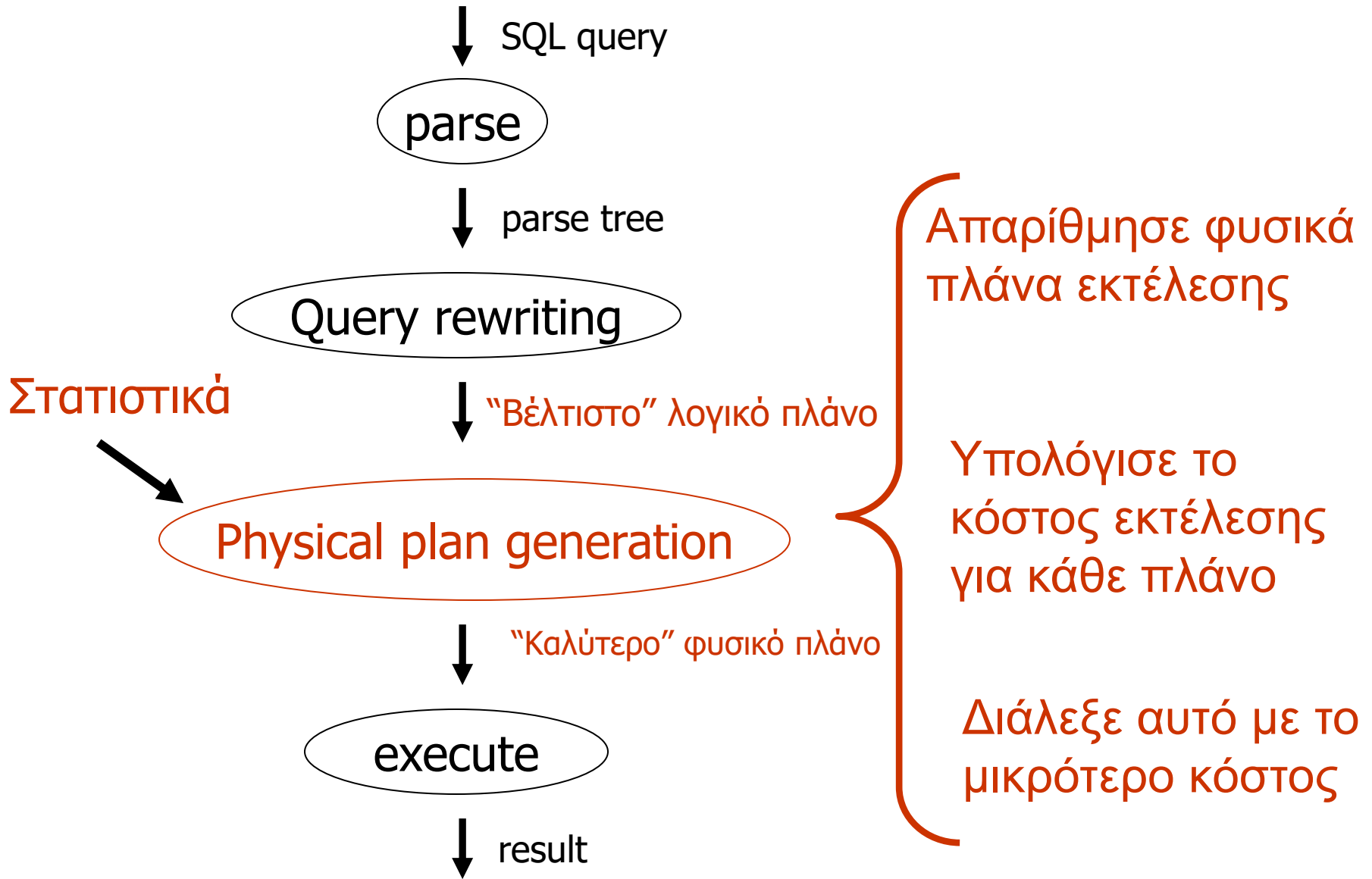




# ΣΤΑΤΙΣΤΙΚΑ

Γιάννης Κωτίδης





# Μοντέλο Κόστους Φυσικού Τελεστή

Κόστος = CPU Time + I/O Time

- Συνήθως CPU Time  $\ll$  I/O Time
- Σκεφτείτε πότε αυτό ενδέχεται να μην ισχύει...

Απλοποίηση

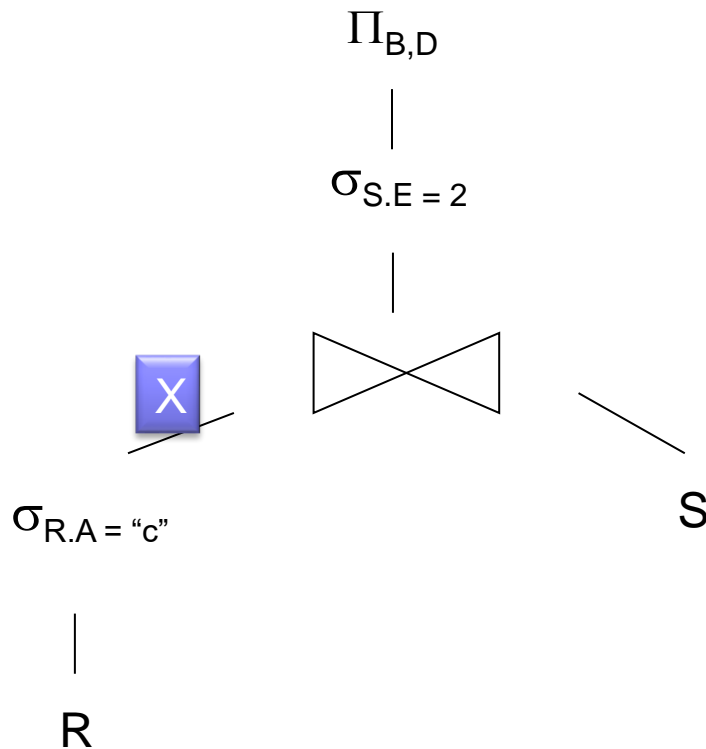
Προσμετρώ μόνο το κόστους του I/O

- Κόστος = Αριθμός σελίδων που διαβάζω από ή γράφω στο δίσκο

# Πρόσθετες απλοποιήσεις για τη συζήτηση που ακολουθεί

- Αγνοώ τον τύπο του I/O (σειριακή, τυχαία προσπέλαση)
- 1 CPU (1 core), 1 Δίσκος
- Δεν προσμετρώ το κόστος για να κάνω output το τελικό αποτέλεσμα
  - Αυτό το κόστος είναι το ίδιο (γιατί?) για όλα τα πλάνα οπότε μπορώ να το αγνοήσω κατά τη φάση της βελτιστοποίησης

# Υπολογισμός Κόστους



- Έστω  $X$  το αποτέλεσμα της επιλογής
- Το κόστος του φυσικού τελεστή που θα υλοποιήσει τη σύζευξη εξαρτάται και από το μέγεθος των 2 σχέσεων  $X, S$  που του δίνονται σαν είσοδο
  - Πόσες εγγραφές περιέχει η σχέση  $X$ ;

# Άρα

- Για να εκτιμήσω το κόστος εκτέλεσης ενός φυσικού τελεστή χρειάζεται να μπορώ να εκτιμώ **κατά προσέγγιση** το μέγεθος του αποτελέσματος των σχεσιακών τελεστών που του παρέχουν είσοδο
  - Το μέγεθος του αποτελέσματος ενός σχεσιακού τελεστή δεν επηρεάζεται από το συγκεκριμένο φυσικό τελεστή που τον υλοποιεί
    - Στο προηγούμενο παράδειγμα ο αριθμός εγγραφών στην X είναι ο ίδιος είτε η επιλογή εκτελεσθεί μέσω ευρετηρίου είτε χωρίς αυτό
      - ...το κόστος αποτίμησης όμως είναι προφανώς διαφορετικό
      - ...η σειρά των αποτελεσμάτων μπορεί να είναι διαφορετική
        - ...και αυτό μπορεί να επηρεάσει το κόστος του τελεστή που υλοποιεί το JOIN, όπως θα δούμε



# Τι χρειάζομαι;

- (1) Τρόπους υπολογισμού του μεγέθους του αποτελέσματος ενός σχεσιακού τελεστή (πριν τον εκτελέσω!)
- (2) Υπολογισμό του κόστους σε I/O

...ας δούμε πρώτα το (1)



# Χρήση απλών στατιστικών

- $T(R)$ : # εγγραφών (**T**uples) στη σχέση  $R$
- $S(R)$ : # μέγεθος (**S**ize) εγγραφής της  $R$  σε bytes
- $B(R)$ : # αριθμός σελίδων (**B**locks) που καταλαμβάνει η  $R$  στο δίσκο
- $V(R, A)$ : # αριθμός **διακριτών** τιμών (**V**alues) του γνωρίσματος  $R.A$

# Παράδειγμα

R

A	B	C	D
cat	1	10	a
cat	1	20	b
dog	1	30	a
dog	1	40	c
bat	1	50	d

A: 20 byte string

B: 4 byte integer

C: 8 byte date

D: 5 byte string

$$T(R) = 5 \quad S(R) = 37$$

$$V(R,A) = 3$$

$$V(R,C) = 5$$

$$V(R,B) = 1$$

$$V(R,D) = 4$$

Εκτίμηση για καρτεσιανό γινόμενο

$$W = R1 \times R2$$

$$T(W) = T(R1) \times T(R2)$$

$$S(W) = S(R1) + S(R2)$$

# Εκτίμηση για επιλογή

$$W = \sigma_{A=a}(R)$$

$$S(W) = S(R)$$

$$T(W) = ?$$

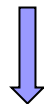
# Πρέπει να κάνω κάποιες **παραδοχές** ελλείψη άλλων στατιστικών

1. Οι τιμή  $A=\alpha$  που ψάχνω με την επιλογή εμφανίζεται μέσα στις εγγραφές της  $R$  [**ψάχνω κάτι που υπάρχει**]
2. Η κατανομή των τιμών του γνωρίσματος  $A$  είναι **ομοιόμορφη** ανάμεσα στις υπάρχουσες τιμές (συνολικά  $V(R,A)$  διακριτές τιμές)

Άρα ψάχνω μία από τις  $V(R,A)$  τιμές που παίρνει το γνώρισμα  $A$  στη σχέση  $R$

Ομοιόμορφη κατανομή  $\rightarrow$  η τιμή  $\alpha$  εμφανίζεται με πιθανότητα  $\frac{1}{V(R,A)}$  σε μία εγγραφή

Έχω  $T(R)$  εγγραφές



Η επιλογή  $\sigma_{A=\alpha}(R)$  επιστρέφει  $T(R) * \frac{1}{V(R,A)} = \frac{T(R)}{V(R,A)}$  εγγραφές

## Παράδειγμα

R	A	B	C	D
cat	1	10	a	
cat	1	20	b	
dog	1	30	a	
dog	1	40	c	
bat	1	50	d	

$$V(R,A)=3$$

$$V(R,B)=1$$

$$V(R,C)=5$$

$$V(R,D)=4$$

$$W = \sigma_{z=val}(R) \quad T(W) = \frac{T(R)}{V(R,Z)}$$

$$\text{Example: } T(\sigma_{A=cat}(R)) = \frac{T(R)}{V(R,A)} = \frac{5}{3}$$

# Αιτιολόγηση

R	A	B	C	D
cat	1	10	a	
cat	1	20	b	
dog	1	30	a	
dog	1	40	c	
bat	1	50	d	

Υπολογίσαμε:  $T(\sigma_{A=\alpha}(R)) = \frac{5}{3}$  [το ίδιο για  $\alpha = \text{cat}, \text{dog}, \text{bat}$ ]

Πόσες εγγραφές βρίσκω με  $A = \text{cat}$  ..... 2

Πόσες εγγραφές βρίσκω με  $A = \text{dog}$  ..... 2

Πόσες εγγραφές βρίσκω με  $A = \text{bat}$  ..... 1

Κατά μέσο όρο βρίσκω  $\frac{2+2+1}{3} = \frac{5}{3} = \frac{T(R)}{V(R,A)}$  εγγραφές

## Παράδειγμα (γνωρίσματα B,C,D)

R	A	B	C	D
	cat	1	10	a
	cat	1	20	b
	dog	1	30	a
	dog	1	40	c
	bat	1	50	d

$$V(R,A)=3$$

$$V(R,B)=1$$

$$V(R,C)=5$$

$$V(R,D)=4$$

$$T(\sigma_{B=1}(R)) = \frac{T(R)}{V(R,B)} = \frac{5}{1} = 5$$

$$T(\sigma_{C=20}(R)) = \frac{T(R)}{V(R,C)} = \frac{5}{5} = 1$$

$$T(\sigma_{D=b}(R)) = \frac{T(R)}{V(R,D)} = \frac{5}{4}$$



# Άσκηση

- Εκτίμηση  $T(\sigma_{A \neq \alpha}(R)) = ?$

# Εκτίμηση $T(\sigma_{A \neq \alpha}(R))$

- Με τις προηγούμενες παραδοχές η πιθανότητα να βρω τη τιμή  $\alpha$  είναι

- $P(A = \alpha) = \frac{1}{V(R,A)}$

- Άρα  $P(A \neq \alpha) = 1 - \frac{1}{V(R,A)} = \frac{V(R,A)-1}{V(R,A)}$

- Επομένως  $T(\sigma_{A \neq \alpha}(R)) = T(R) * \frac{V(R,A)-1}{V(R,A)}$

## Παράδειγμα

R	A	B	C	D
cat	1	10	a	
cat	1	20	b	
dog	1	30	a	
dog	1	40	c	
bat	1	50	d	

$$V(R,A)=3$$

$$V(R,B)=1$$

$$V(R,C)=5$$

$$V(R,D)=4$$

$$\text{Example: } T(\sigma_{A \neq \text{cat}}(R)) = T(R) * \frac{V(R,A)-1}{V(R,A)} = 5 * \frac{3-1}{3} = \frac{10}{3}$$

$$\text{Recall: } T(\sigma_{A=\text{cat}}(R)) = \frac{5}{3}$$

$$\text{Thus: } T(\sigma_{A \neq \text{cat}}(R)) + T(\sigma_{A=\text{cat}}(R)) = \frac{15}{3} = 5 \checkmark$$

Πως υπολογίζω το μέγεθος της σύζευξης:  $W = R1 \bowtie R2$

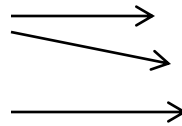
- Ποια είναι η **ελάχιστη** και ποια η **μέγιστη** τιμή του  $T(W)$ ;
- Τι γίνεται στην πράξη;
- Παρατήρηση: συνήθως κάνω join σχέσεις που έχουν διασπασθεί κατά τη διαδικασία της κανονικοποίησης

# Ειδική περίπτωση σχέση (1:N) με περιορισμό κλειδιού

- Έστω οι σχέσεις:  
**Customer(CustId,Name,Address)** και  
**Buys(CustId,ProdId,Date)**
- **Customer**  $\bowtie$  **Buys** γυρίζει τους πελάτες με τις αγορές τους

Customer

<u>CustId</u>	Name	Address
C1	John	Athens
C2	Nick	Athens
C3	Eleni	Piraeus



Buys

CustId	ProdId	Date
C1	P2	Jan 3 2012
C1	P4	Dec 12 2011
C3	P1	Feb 7 2012

# W = Customer Buys

- Ειδική αλλά αρκετά συνηθισμένη περίπτωση
  - Το **CustId** είναι κλειδί στη σχέση **Customer** και ξένο κλειδί στη σχέση **Buys**
  - $V(\text{Customer}, \text{CustId}) =$  συνολικός αριθμός πελατών
  - $V(\text{Buys}, \text{CustId}) =$  αριθμός πελατών με μία τουλάχιστον αγορά

<u>CustId</u>	Name	Address
C1	John	Athens
C2	Nick	Athens
C3	Eleni	Piraeus

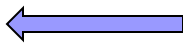


<u>CustId</u>	ProdId	Date
C1	P2	Jan 3 2012
C1	P4	Dec 12 2011
C3	P1	Feb 7 2012

# W = Customer ⋈ Buys

- Κάθε εγγραφή από τη σχέση Buys κάνει join με μία ακριβώς εγγραφή από τη σχέση Customer με βάση την τιμή του κλειδιού CustId (γιατί?)
  - $T(W) = T(\text{Buys})$

CustId	Name	Address
C1	John	Athens
C2	Nick	Athens
C3	Eleni	Piraeus



CustId	ProdId	Date
C1	P2	Jan 3 2012
C1	P4	Dec 12 2011
C3	P1	Feb 7 2012

# Ας δούμε τώρα μία ποιο γενική περίπτωση

- Βασική υπόθεση (informally): κάνουμε join σχέσεις που κάνουν join!!!
- Επίσης όπως και στην επιλογή υποθέτουμε ομοιόμορφη κατανομή



$$\underline{W = R1 \bowtie R2 \quad X \cap Y = A}$$

R1	A	B	C

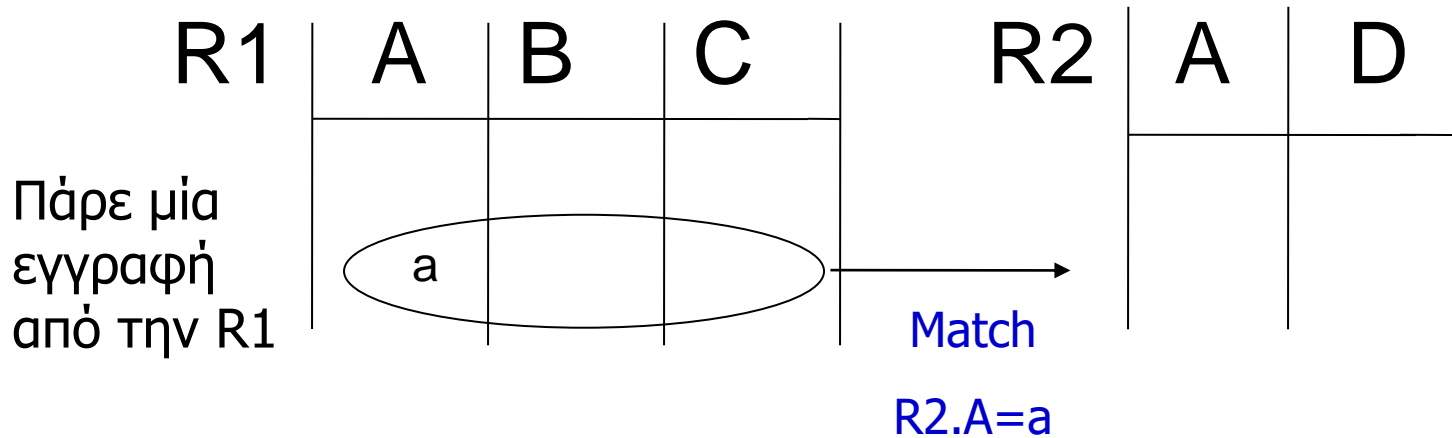
R2	A	D

Έστω ότι  $V(R1,A) \leq V(R2,A)$  :

Θα κάνω την υπόθεση ότι κάθε τιμή του γνωρίσματος R1.A εμφανίζεται στη σχέση R2

(Containment of value sets)

# Υπολογισμός $T(W)$ εφόσον $V(R1,A) \leq V(R2,A)$



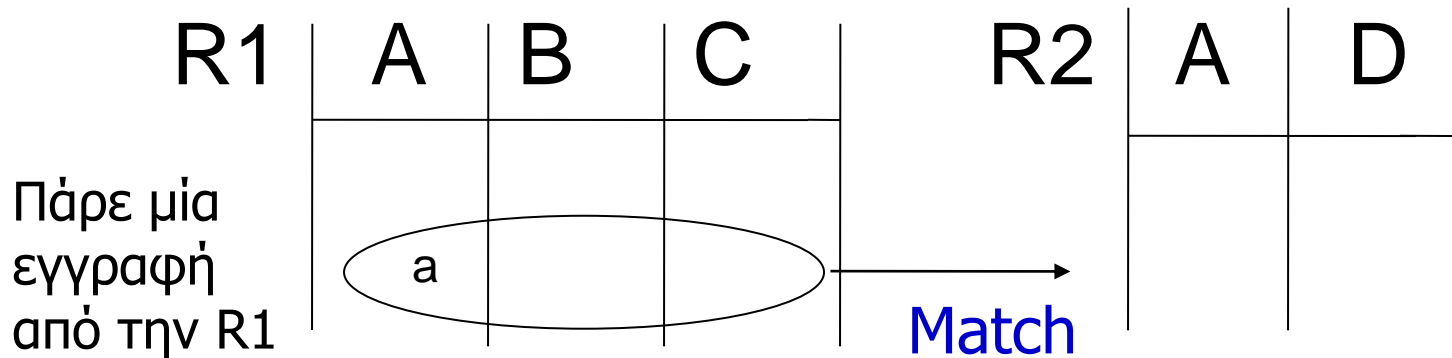
Κάθε εγγραφή της R1 κάνει join με  $\frac{T(R2)}{V(R2,A)}$   
εγγραφές της σχέσης R2

Γιατί?

# Ας σκεφτούμε

- Η εγγραφή που πήρα από την R1 ταιριάζει (κάνει join) με εγγραφές της R2 που έχουν την ίδια τιμή «α» στο γνώρισμα R2.A
- Για να βρω με ποιες ταιριάζει ας κάνω την επιλογή  $\sigma_{R2.A=a}(R2)$
- Έχουμε ήδη βρει ότι:  $T(\sigma_{R2.A=a}(R2)) = \frac{T(R2)}{V(R2,A)}$

# Υπολογισμός $T(W)$ όταν $V(R1,A) \leq V(R2,A)$



Κάθε εγγραφή κάνει join με  $\frac{T(R2)}{V(R2,A)}$  εγγραφές

άρα 
$$T(W) = \frac{T(R1) * T(R2)}{V(R2,A)}$$

$$W = R1 \bowtie R2 \quad X \cap Y = A^*$$

---

- Αν  $V(R1, A) \leq V(R2, A)$ :

$$T(W) = \frac{T(R1) * T(R2)}{V(R2, A)}$$

- Αν  $V(R2, A) \leq V(R1, A)$ :

(προκύπτει με τον ίδιο συλλογισμό αν αντιστρέψετε τη σειρά των σχέσεων)

$$T(W) = \frac{T(R1) * T(R2)}{V(R1, A)}$$

\*A είναι το κοινό γνώρισμα στη φυσική σύζευξη

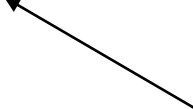
Οπότε (για  $W = R1 \bowtie R2$ )

$$T(W) = \frac{T(R1) * T(R2)}{\max\{V(R1,A), V(R2,A)\}}$$

# Επίσης

$$S(W) = S(R1) + S(R2) - S(A)$$

μέγεθος (bytes) γνωρίσματος A



# Σύνθετο Παράδειγμα

- Έστω ΒΔ με τους εργαζόμενους σε μια εταιρία με πέντε τμήματα
- Για τους υπαλλήλους στο τμήμα έρευνας (Research) θέλω να δώσω bonus ανάλογα με την αρχαιότητα τους (seniority level: 0..9)
- Τα δεδομένα αποθηκεύονται σε δύο πίνακες **Emp** και **Bonus**
  - Κάθε υπάλληλος έχει μια τιμή seniority **level**  $\in [0..9]$
  - Το τμήμα έρευνας έχει υπαλλήλους από όλα τα επίπεδα αρχαιότητας
  - Υπάρχουν συνολικά 80 διαφορετικά είδη bonus στο πίνακα Bonus που αναφέρονται σε 7 από τα 10 seniority levels
    - Υπάρχουν διαφορετικά πιθανά bonuses για κάποια επίπεδα αρχαιότητας, αλλά όχι για όλα!



# ΣΤΑΤΙΣΤΙΚΑ ΣΧΕΣΕΩΝ Emp, Bonus

## ■ Σχέση Emp

- $T(\text{Emp}) = 1000$
- $V(\text{Emp}, \text{dept}) = 5$
- $V(\text{Emp}, \text{level}) = 10$

## ■ Σχέση Bonus

- $T(\text{Bonus}) = 80$
- $V(\text{Bonus}, \text{level}) = 7$

# Δίνεται η παρακάτω SQL επερώτηση

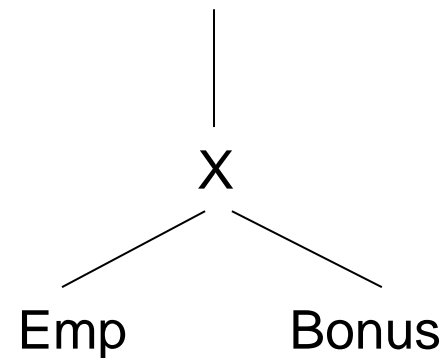
```
SELECT * FROM Emp,Bonus  
WHERE Emp.level = Bonus.level  
AND Emp.dept = "Research"
```

- Σχεδιάστε το αρχικό και τελικό λογικό πλάνο
- Πόσα αποτελέσματα αναμένεται να επιστρέψει η επερώτηση?

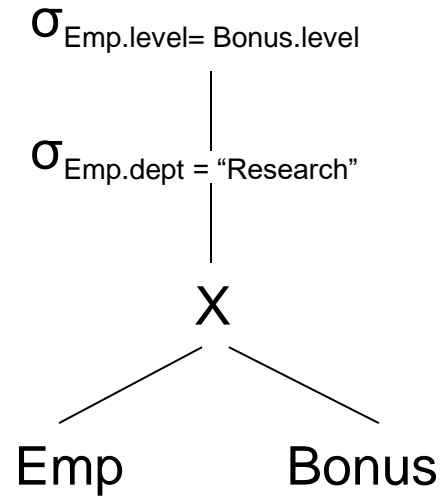
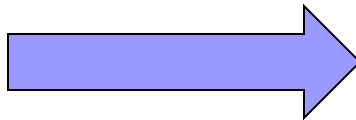
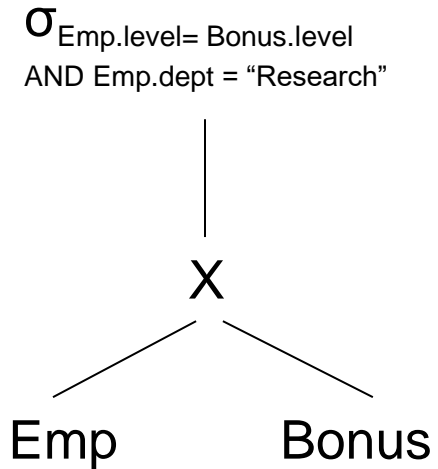
# Αρχικό Λογικό Πλάνο

```
SELECT * FROM Emp,Bonus  
WHERE Emp.level = Bonus.level  
AND Emp.dept = "Research"
```

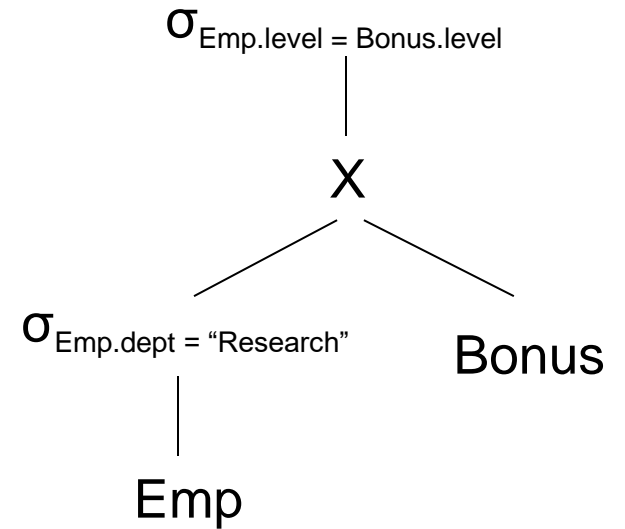
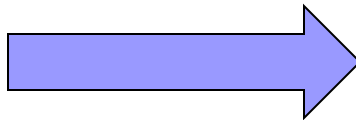
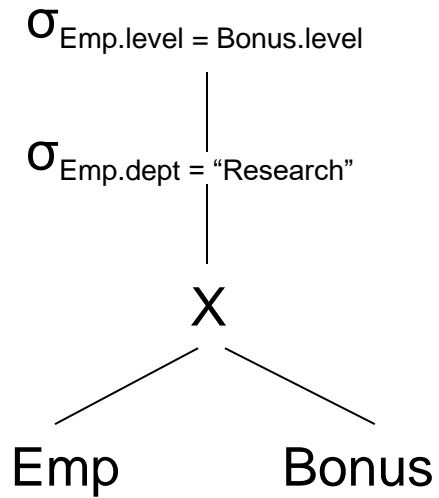
$\sigma_{\text{Emp.level} = \text{Bonus.level}}$   
AND Emp.dept = "Research"



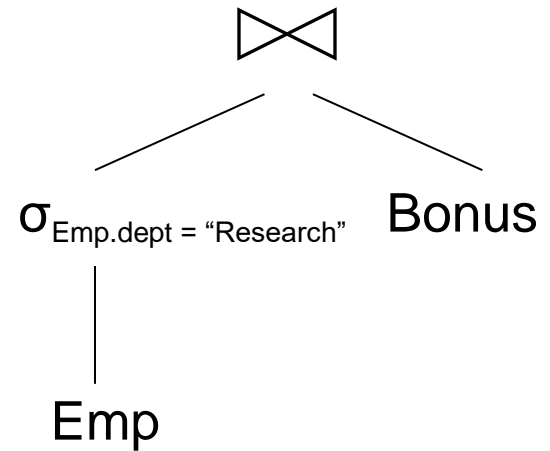
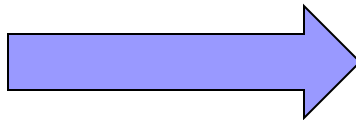
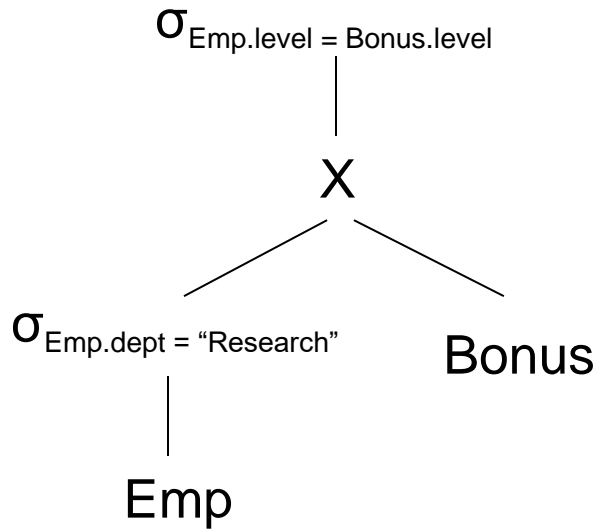
# Rewrites (1)



# Rewrites (2)

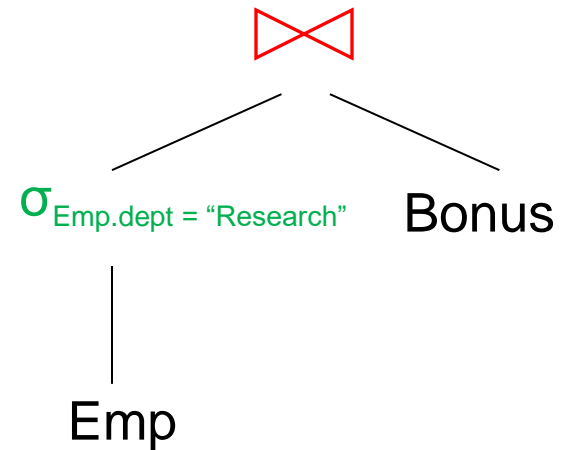


# Rewrites (3)



# Τελικό Λογικό Πλάνο

```
SELECT * FROM Emp,Bonus  
WHERE Emp.level= Bonus.level  
AND    Emp.dept = "Research"
```



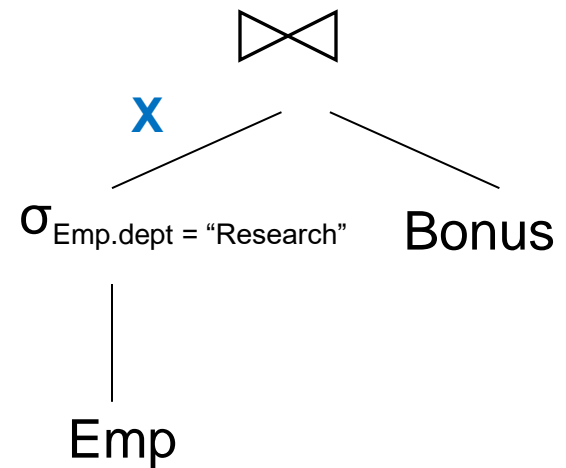
# Υπολογισμός Στατιστικών ( $\sigma$ )

Έστω  $X = \sigma_{\text{Emp.dept} = \text{"Research"}}(\text{Emp})$   
το αποτέλεσμα της επιλογής

Πόσα αποτελέσματα αναμένουμε  
να επιστρέψει το  $X$ ?

Δίνονται:

- $T(\text{Emp}) = 1000$
- $V(\text{Emp}, \text{dept}) = 5$



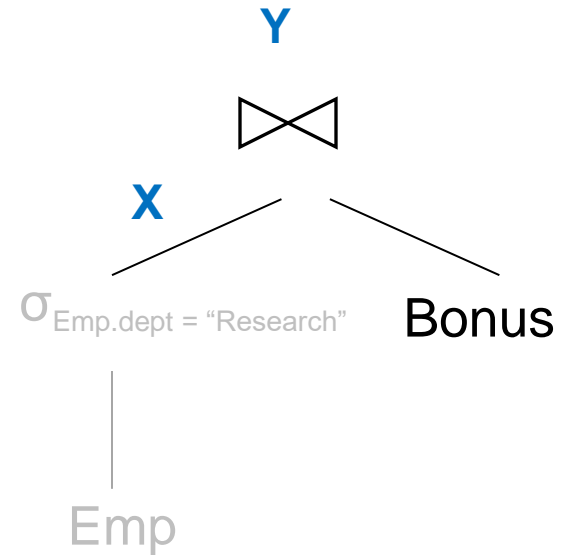
$$\text{Εκτίμηση: } T(X) = \frac{1000}{5} = 200 \text{ εγγραφές}$$



# Υπολογισμός Στατιστικών ( $\bowtie$ )

- $T(X) = 200$
- $V(X, \text{level}) = 10$  (υπόθεση ανεξαρτησίας ανάμεσα σε dept και level)
- $T(\text{Bonus}) = 80$
- $V(\text{Bonus}, \text{level}) = 7$

Έστω  $Y = X \bowtie \text{Bonus}$



$$\text{Εκτίμηση: } T(Y) = \frac{T(X) * T(\text{Bonus})}{\max(V(X, \text{level}), V(\text{Bonus}, \text{level}))} = \frac{200 * 80}{\max(10, 7)} = 1600$$

**Επομένως περιμένω (κατ' εκτίμηση) 1600 αποτελέσματα**



# Preservation of Set Values

- When joining two relations, any attribute **that is not a join** attribute does not lose values from its set of possible values
  - Unless join result is smaller than the number of distinct values
- Assumption useful in-cases we have a chain of joins

# Παράδειγμα

- Έστω  $R(A, B)$  JOIN  $S(B, C)$  JOIN  $W(C, D)$
- Δίνονται
  - $T(R) = 100, V(R, B) = 10$
  - $T(S) = 50, V(S, B) = 5, V(S, C) = 8$
  - $T(W) = 20, V(W, C) = 5$

**$R(A,B) \text{ JOIN } S(B,C)$**  JOIN  $W(C,D)$

■ Δίνονται

□  $T(R) = 100, V(R,B) = 10$

□  $T(S) = 50, V(S,B) = 5, V(S,C) = 8$

□  $T(W) = 20, V(W,C) = 5$

■ Έστω  **$U = R(A,B) \text{ JOIN } S(B,C)$**

■  $T(U) = ?$

# R(A,B) JOIN S(B,C) JOIN W(C,D)

## ■ Δίνονται

□  $T(R) = 100, V(R,B) = 10$

□  $T(S) = 50, V(S,B) = 5, V(S,C) = 8$

□  $T(W) = 20, V(W,C) = 5$

## ■ Έστω $U = R(A,B) \text{ JOIN } S(B,C)$

$$T(U) = \frac{T(R) * T(S)}{\max(V(R,B), V(S,B))} = \frac{100 * 50}{\max(10, 5)} = 500$$

## ■ $V(U,C) = ?$

# R(A,B) JOIN S(B,C) JOIN W(C,D)

- Δίνονται

- $T(R) = 100, V(R,B) = 10$

- $T(S) = 50, V(S,B) = 5, V(S,C) = 8$

- $T(W) = 20, V(W,C) = 5$

- Έστω  $U = R(A,B) \text{ JOIN } S(B,C)$

- $T(U) = 500$

- $V(U,C) = V(S,C) = 8$  (preservation of set values)

U

$R(A,B) \text{ JOIN } S(B,C) \text{ JOIN } W(C,D)$

- Δίνονται
  - $T(W) = 20, V(W,C) = 5$
- Βρήκαμε ( $U=R(A,B) \text{ JOIN } S(B,C)$ )
  - $T(U) = 500$
  - $V(U,C) = 8$
- $T(U \text{ JOIN } W) = ?$

# R(A,B) JOIN S(B,C) JOIN W(C,D)

- Δίνονται

- $T(W) = 20, V(W,C) = 5$

- Βρήκαμε ( $U=R(A,B) JOIN S(B,C)$ )

- $T(U) = 500$

- $V(U,C) = 8$

- $T(U JOIN W) = \frac{500*20}{\max(5,8)} = 1250$

- Επομένως  $T(R JOIN S JOIN W)=1250$