

# Διάλεξη 15: Ατομική ΚΚΜ Εγγραφής/Ανάγνωσης με Γρήγορες Λειτουργίες

ΕΠΛ 432: Κατανεμημένοι Αλγόριθμοι



# Τι θα δούμε σήμερα

- Γρήγορες Λειτουργίες
- Συστήματα Απαρτίας
- Αλγόριθμος SLIQ – Χρήση Quorum Views

# Γρήγορες Λειτουργίες

- Μέχρι τώρα είδαμε αλγορίθμους όπου
  - Ανάγνωση χρειάζεται **2 φάσεις**
  - Εγγραφή χρειάζεται **1 ή 2 φάσεις**
- **Γρήγορη Λειτουργία:** Είναι η λειτουργία εγγραφής ή ανάγνωσης η οποία χρειάζεται **μόνο 1 φάση ή ένα επικοινωνιακό γύρο** για να ολοκληρωθεί.
  - Π.χ. Όλες οι εγγραφές στον αλγόριθμο ABD είναι γρήγορες
- **Γρήγορος Αλγόριθμος:** Είναι αυτός όπου σε **κάθε εκτέλεσή** του **όλες** οι λειτουργίες είναι γρήγορες

# Συστήματα Απαρτίας

- Παρατήρηση:
  - Η επικοινωνία με την πλειοψηφία των αντιγράφων διασφαλίζει ότι για **κάθε ζεύγος λειτουργιών** υπάρχει **τουλάχιστον ένας διαχειριστής αντιγράφου** που λαμβάνει μήνυμα και από τις δύο λειτουργίες
  - Έπεται από το γεγονός ότι η τομή δύο συνόλων πλειοψηφίας δεν είναι κενή
- Ιδέα Συστημάτων Απαρτίας:
  - Οργάνωσε τους διαχειριστές αντιγράφων σε σύνολα (όχι απαραίτητα σύνολα πλειοψηφίας) έτσι ώστε **κάθε δύο τέτοια σύνολα να έχουν μη κενή τομή**
  - Τα σύνολα ονομάζονται **απαρτίες**

# Ορισμός Συστημάτων Απαρτίας

- Quorum System :  $\mathbf{Q}$

$$\mathbf{Q} = \{Q : Q \subseteq S\} \text{ s.t. } \forall Q_i, Q_j \in \mathbf{Q} : Q_i \cap Q_j \neq \emptyset$$

- **Απενεργοποιημένη Απαρτία:** Μια απαρτία  $Q$  είναι απενεργοποιημένη σε μια εκτέλεση  $\alpha$  εαν

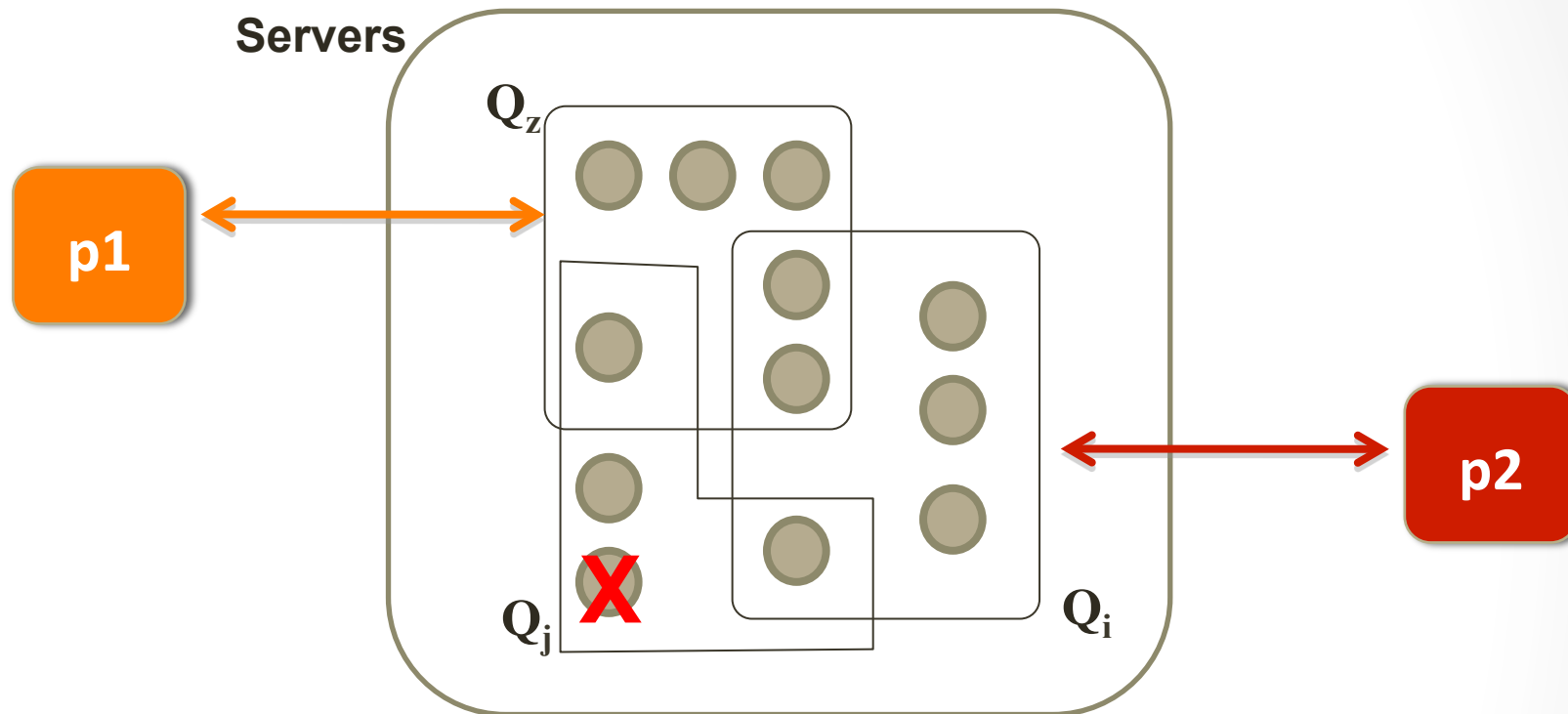
$$\exists s \in Q \text{ s.t. } fail_s \in \alpha$$

- **Απενεργοποιημένο Σύστημα Απαρτίας:** Ένα σύστημα απαρτίας  $\mathbf{Q}$  είναι απενεργοποιημένο εαν

$$\forall Q \in \mathbf{Q}, Q \text{ is faulty}$$

- **Υπόθεση Σφάλματος:** Μια απαρτία παραμένει ενεργοποιημένη

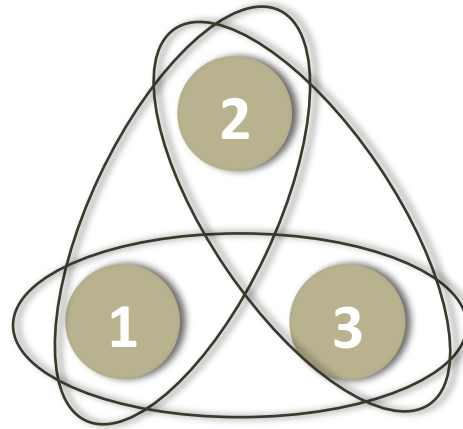
# Συστήματα Απαρτίας Σχηματικά



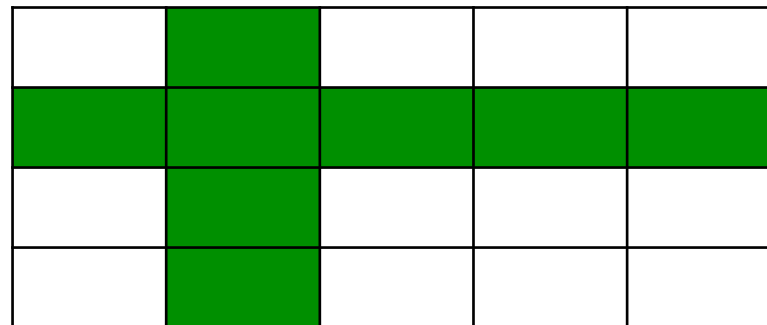
- $Q_i, Q_j, Q_z$  είναι **απαρτίες**
- **Σύστημα Απαρτίας** είναι το σύνολο  $\{Q_i, Q_j, Q_z\}$ 
  - Ιδιότητα: κάθε ζεύγος απαρτιών τέμνονται
- Κάθε λειτουργία Εγγραφής/Ανάγνωσης επικοινωνεί με μια απαρτία
- **Εσφαλμένη Απαρτία**: Αυτή που περιέχει μια εσφαλμένη διεργασία

# Παραδείγματα Συστημάτων Απαρτίας

- Πλειοψηφία

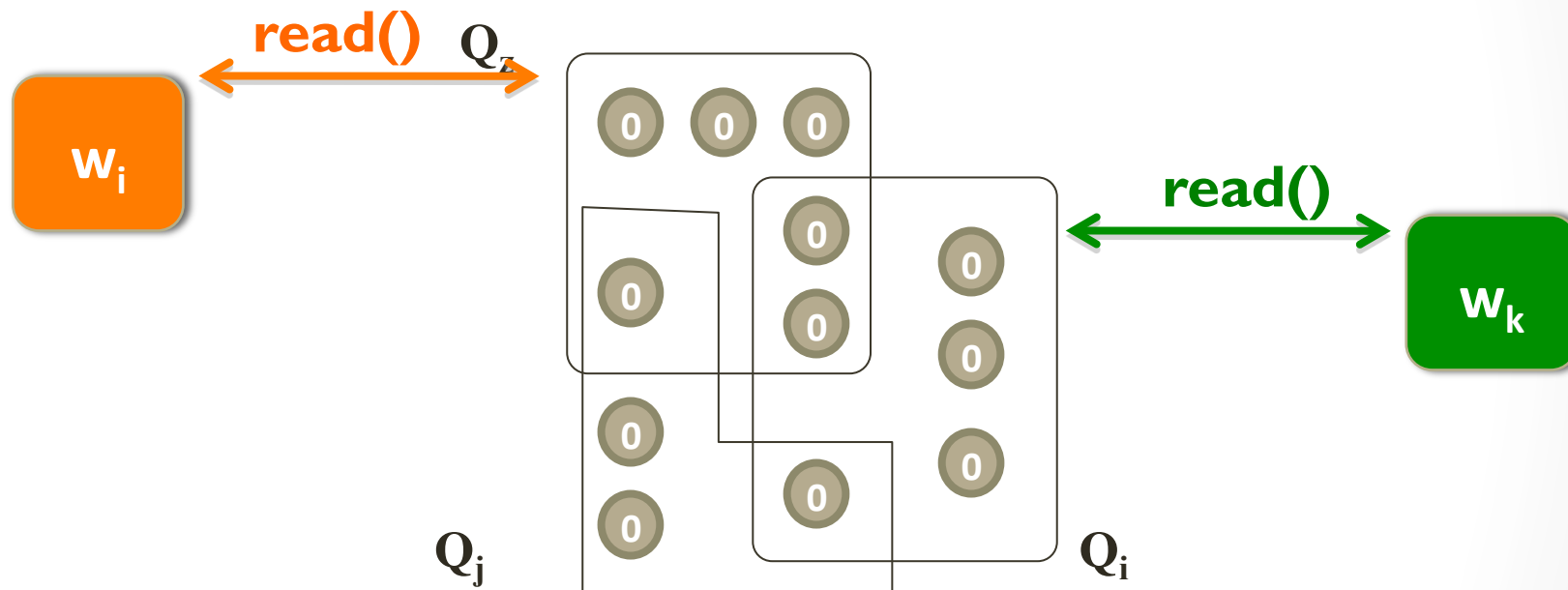


- Απαρτίες σε Πλέγμα: Κάθε απαρτία μια γραμμή και μια στήλη



# Αλγόριθμος MWMR με Απαρτίες

- ▶ Υποθέτουμε ότι  $w_i > w_k$

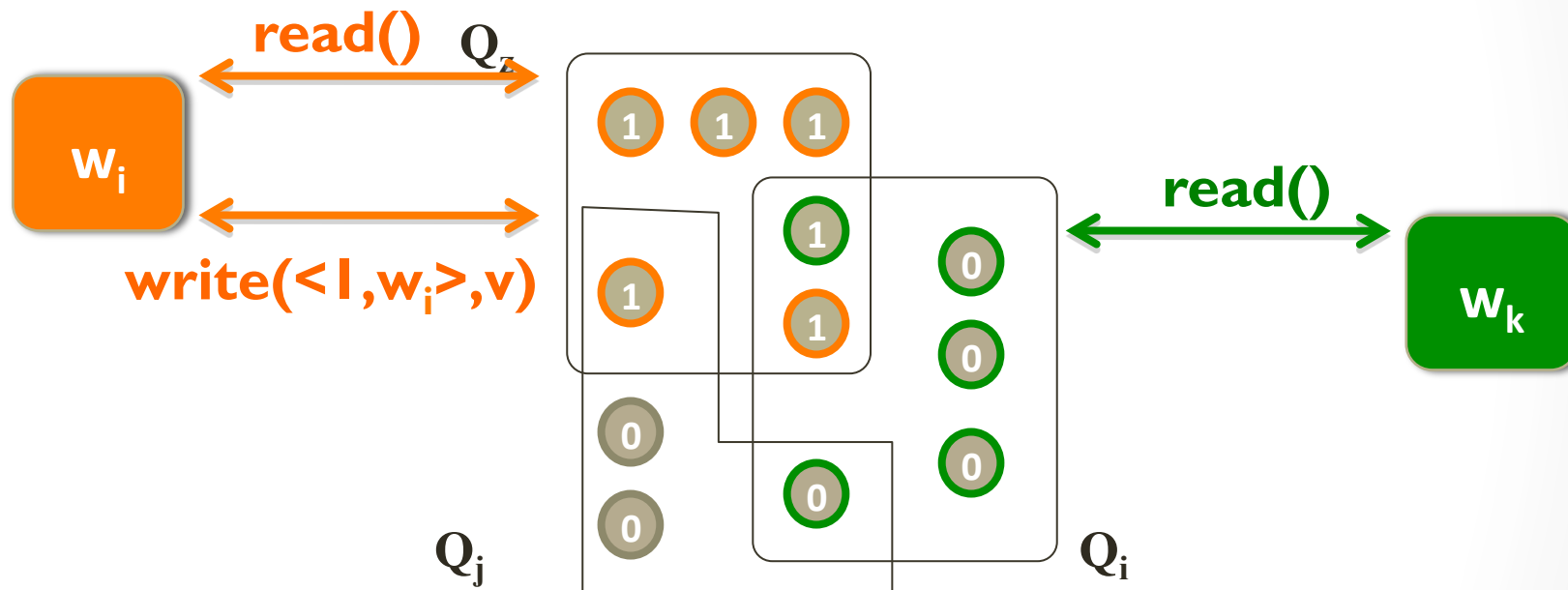


11/9/12



# Αλγόριθμος MWMR με Απαρτίες

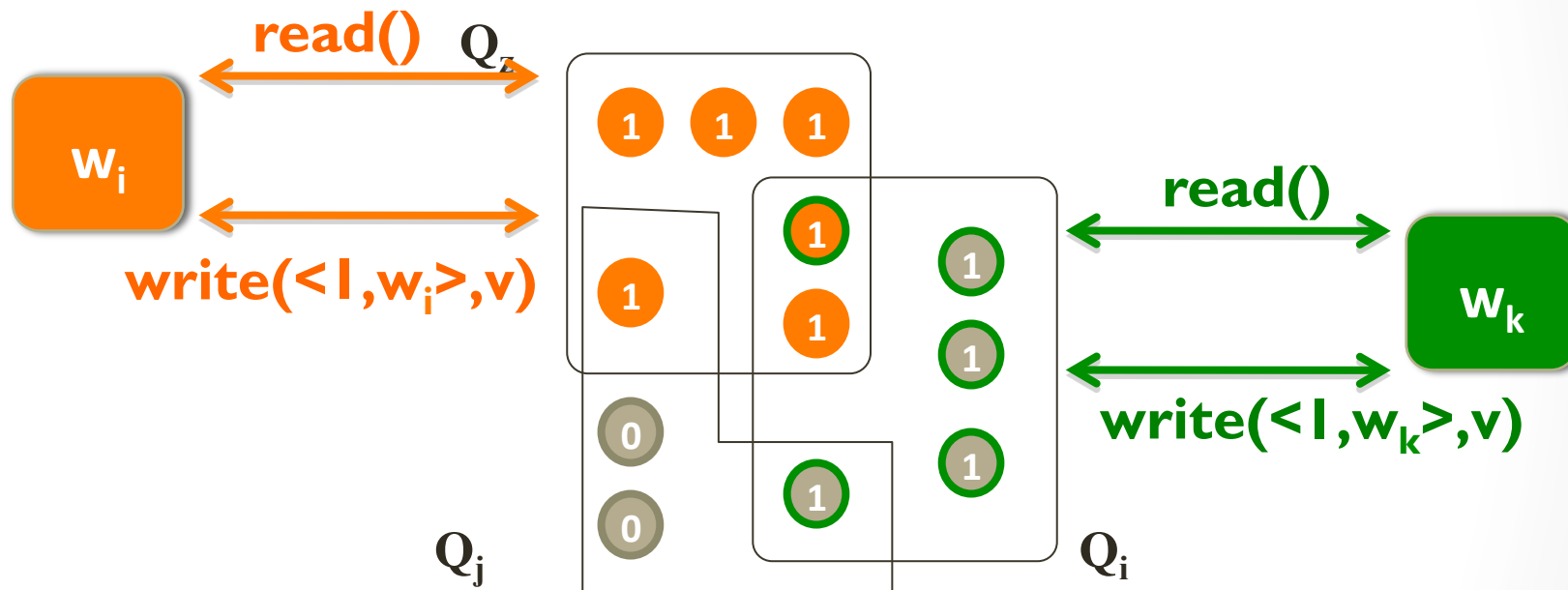
- Υποθέτουμε ότι  $w_i > w_k$



11/9/12

# Αλγόριθμος MWMR με Απαρτίες

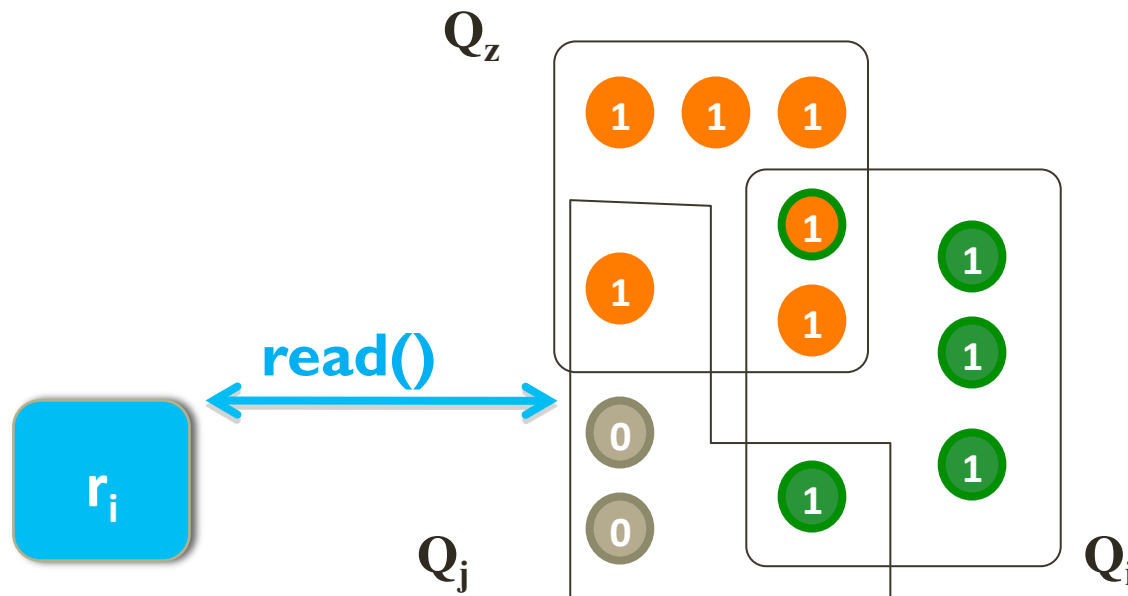
- Υποθέτουμε ότι  $w_i > w_k$



11/9/12

# Αλγόριθμος MWMR με Απαρτίες

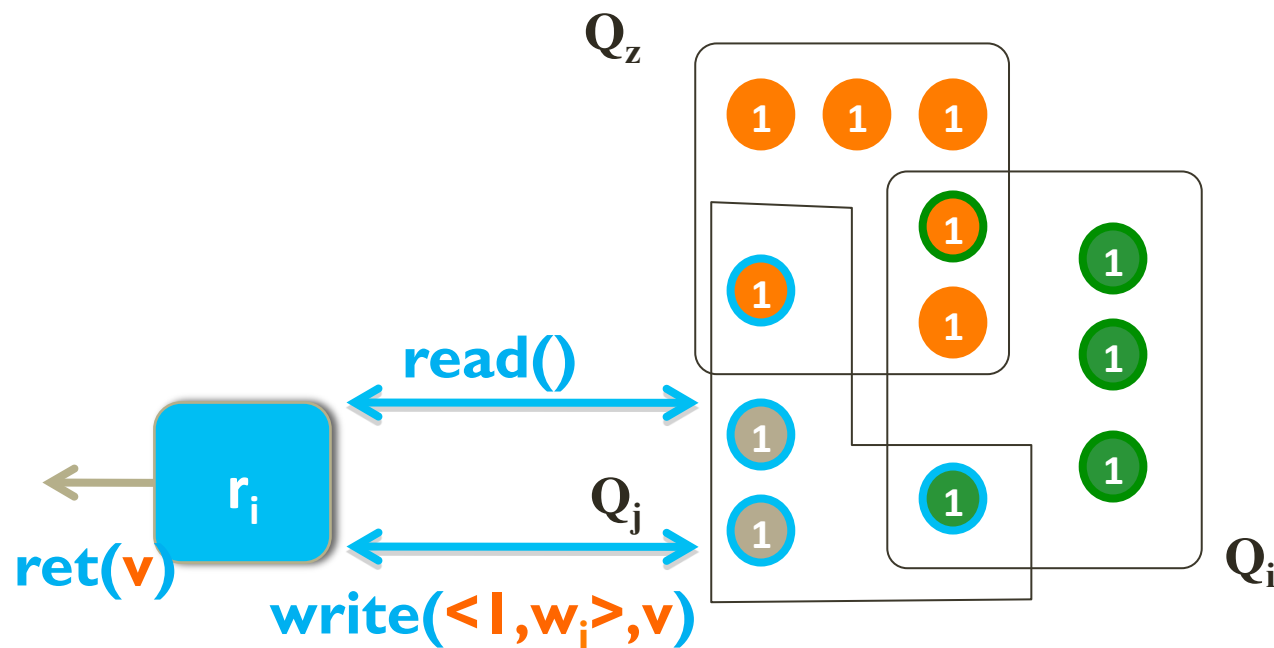
- Υποθέτουμε ότι  $w_i > w_k$



11/9/12

# Αλγόριθμος MWMR με Απαρτίες

- Υποθέτουμε ότι  $w_i > w_k$



11/9/12

Operation Ordering:  $w_k \rightarrow w_i \rightarrow r_i$

# Quorum Views

Ιδέα:

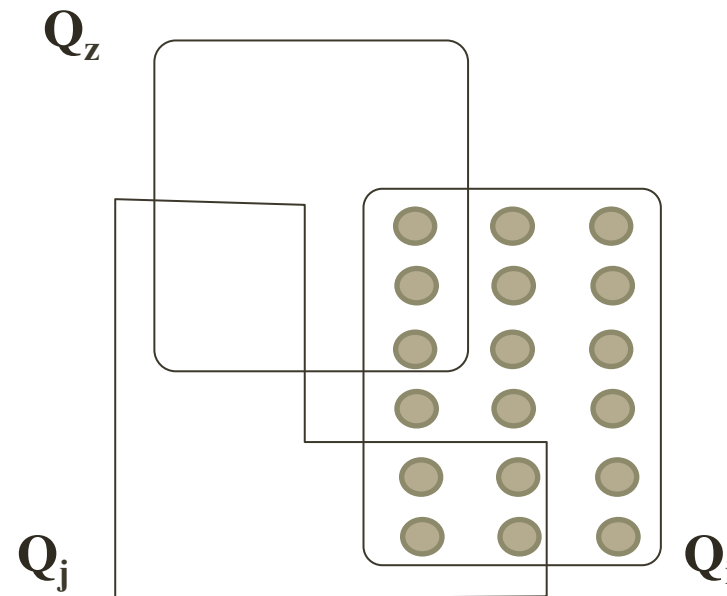
- Προσπαθούμε να προσδιορίσουμε την κατάσταση της τελευταίας λειτουργίας εγγραφής
- Αν η κατάσταση της εγγραφής κατά τον πρώτο γύρο της ανάγνωσης μπορεί να
  - Προσδιοριστεί => Η Ανάγνωση είναι **Γρήγορη**
  - Δεν προσδιορίζεται => Η Ανάγνωση είναι **Αργή**

17/3/2011

# Προσδιοριστέα Εγγραφή - Qview(1)

- Όλα οι εξυπηρετητές της απαρτίας απαντούν με την ίδια χρονοσφραγίδα (ίση με την μέγιστη χρονοσφραγίδα που παρατηρεί ο Αναγνώστης)

$$[qView(1)] : \forall s \in Q_i : s.ts = \max TS$$

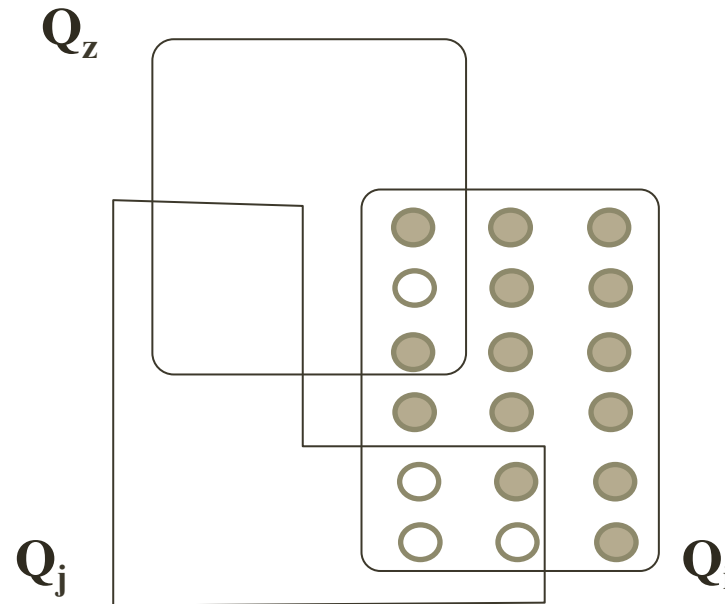


**(Πιθανότητα) Η Εγγραφή έχει Τερματιστεί**

# Προσδιοριστέα Εγγραφή - Qview(2)

- Κάθε τομή περιέχει ένα μέλος με χρονοσφραγίδα μικρότερη από τη μέγιστη

$[qView(2)]: \forall j \neq i, \exists A \subseteq Q_i \cap Q_j \text{ s.t. } A \neq \emptyset \text{ and } \forall s \in A: s.ts < \max TS$

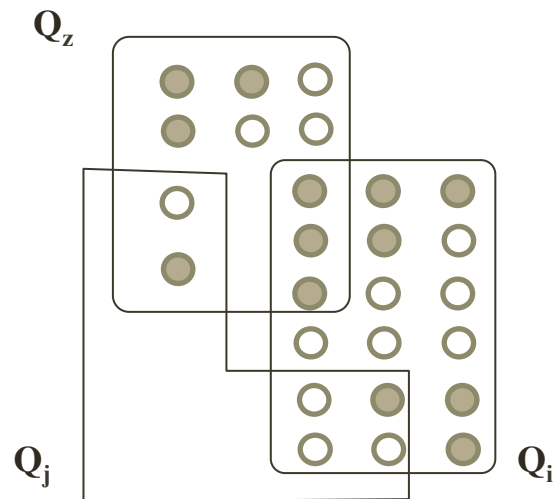


**(Σίγουρα) Η Εγγραφή  $\langle \max Tag, v \rangle$  ΔΕΝ τερματίστηκε**

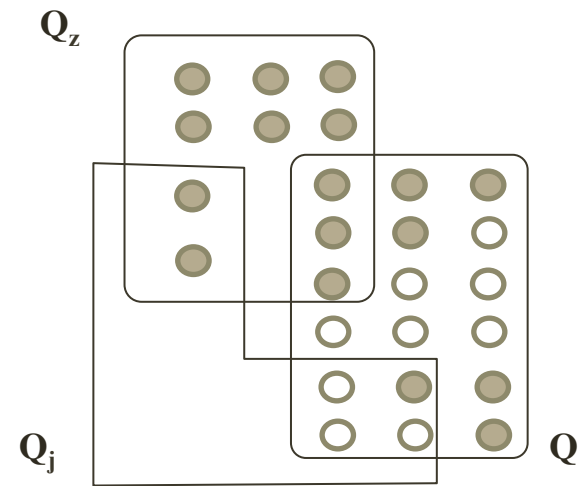
# Απροσδιόριστη Εγγραφή - Qview(3)

- Όλα τα μέλη της τομής της απαρτίας που μας απάντησε και κάποιας άλλης απαρτίας στο σύστημα περιέχουν τη μέγιστη χρονοσφραγίδα.

$[qView(3)]: \exists s \in Q_i \text{ s.t. } s.ts < \max Ts \text{ and } \exists j \neq i \text{ s.t. } \forall s \in Q_i \cap Q_j : s.ts = \max TS$



qV(3) και Μη Τερματισμένη Εγγραφή



qV(3) και Τερματισμένη Εγγραφή

Απροσδιόριστη => Δύο Φάσεις



# Αλγόριθμος SLIQ

## Write Protocol: one round

- P1: Αύξησε την χρονοσφραγίδα σου  $ts$  και στείλε μήνυμα  $write(\langle ts, v \rangle)$  σε όλους και περίμενε απαντήσεις από μια απαρτία
- Μόλις λάβεις απαντήσεις επέστρεψε  $ack$  και τερμάτισε

## Read Protocol: one or two rounds

- P1: Στείλε μήνυμα  $read(\langle ts, v \rangle)$  όπου  $\langle ts, v \rangle$  η τελευταία τιμή που επέστρεψε ο αναγνώστης σε όλα τα αντίγραφα και περίμενε απάντηση από μια απαρτία  $Q$
- $QView_Q(1)$  – **Γρήγορη** και επέστρεψε την τιμή του  $maxTS$
- $QView_Q(2)$  – **Γρήγορη** και επέστρεψε την τιμή του  $maxTS-1$
- $QView_Q(3)$  – **Αργή** προχώρα στη φάση P2 και ακολούθως επεστρεψε  $maxTS$
- P2: προώθησε  $\langle maxTS, v \rangle$  σε μια απαρτία και μετά επέστρεψε  $v$

## Server Protocol: passive role

- Παρέλαβε αιτήσεις  $read(\langle ts, v \rangle)$  και  $write(\langle ts, v \rangle)$
- Εάν  $local.ts < msg.ts$  ενημέρωσε το τοπικό σου αντίγραφο
- Απάντησε με  $reply(\langle ts, v \rangle)$

17/3/2011

# Ορθότητα Αλγορίθμου

- **Ζωτικότητα:** Σύμφωνα με το μοντέλο σφαλμάτων πάντα υπάρχει μια απαρτία να απαντήσει σε κάθε λειτουργία
- **Ατομικότητα:** Πρέπει να δείξουμε ότι
  - Κάθε Ανάγνωση επιστρέφει τουλάχιστον την τιμή που γράφτηκε από την τελευταία Εγγραφή (Νόμιμη Ακολουθία Λειτουργιών)
  - Αν μία Ανάγνωση  $\rho_1$  προηγείται μιας ανάγνωσης  $\rho_2$ , τότε η  $\rho_2$  επιστρέφει την ίδια ή νεότερη τιμή από αυτή που επιστρέφει η  $\rho_1$ .

# Μονοτονικότητα Χρονοσφραγίδας

- **Λήμμα 1:** Αν κάποιος εξυπηρετητής λάβει μήνυμα με  $\langle ts, v \rangle$  τότε κάθε μεταγενέστερη απάντησή του περιλαμβάνει ζεύγος  $\langle ts', v' \rangle$  τ.ω.  $ts' \geq ts$ .
- Απόδειξη:
  - Κάθε φορά που ένας εξυπηρετητής παραλάβει μήνυμα  $read(\langle ts, v \rangle)$  και  $write(\langle ts, v \rangle)$  κάνει τον εξής έλεγχο
    - Αν  $local.ts < msg.ts$  τότε  $local.ts = msg.ts$  και  $local.v = msg.v$
    - Αλλιώς δεν κάνει τίποτα
  - Άρα μετά την παραλαβή μηνύματος  $local.ts \geq msg.ts$
  - Επομένως σε κάθε  $reply(\langle ts', v' \rangle)$  που έπεται οποιασδήποτε από τις πιο πάνω παραλαβές ισχύει ότι  $ts' = local.ts \geq msg.ts = ts$

# Ορθότητα Αλγορίθμου

- **Λήμμα 2:** Αν μια Ανάγνωση  $\rho_1$  έπεται μιας Εγγραφής τότε επιστρέφει τουλάχιστον την τιμή που γράφτηκε
- Απόδειξη:
  - Κάθε εγγραφή επικοινωνεί με μια απαρτία πριν ολοκληρωθεί
  - Έστω  $w$  μια εγγραφή που έλαβε μηνύματα από την απαρτία  $Q_i$  και έγραψε την τιμή  $\langle ts, v \rangle$
  - Έστω  $\rho$  είναι μια ανάγνωση που έπεται της πιο πάνω εγγραφής και λαμβάνει απαντήσεις από την απαρτία  $Q_j$
  - Έχουμε δύο περιπτώσεις: 1)  $Q_i = Q_j$ , και 2)  $Q_i \neq Q_j$

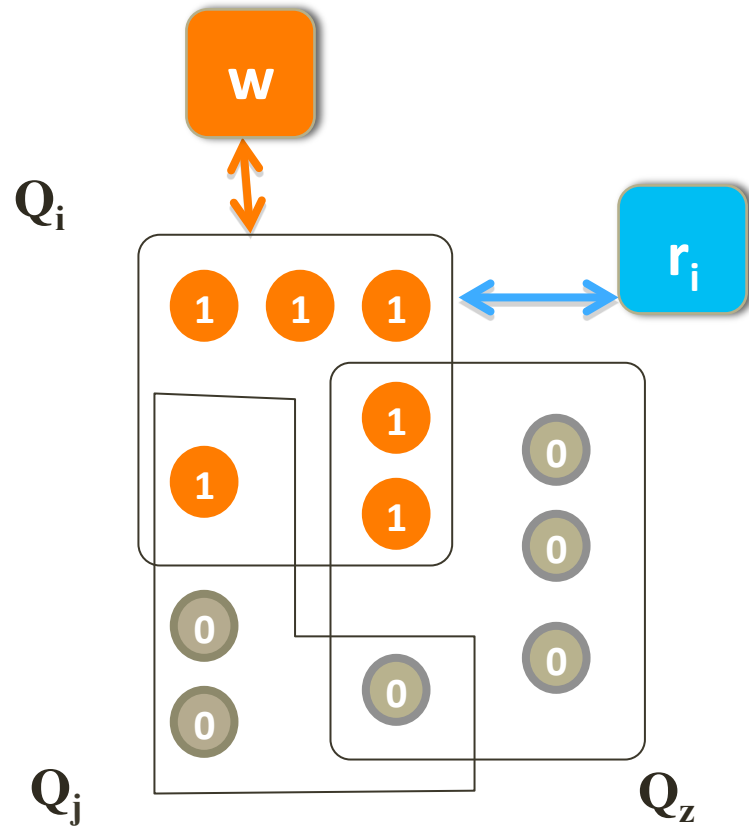
## Λήμμα 2: Περίπτωση 1 ( $Q_i = Q_j$ )

- Όλοι οι εξυπηρετητές στο  $Q_i$  έλαβαν μήνυμα  $\text{write}(\langle ts, v \rangle)$  από τον  $w$  πριν απαντήσουν στην ανάγνωση  $\rho$
- Από το Λήμμα 1 θα στείλουν  $\text{reply}(\langle ts', v' \rangle)$  στον  $\rho$  τ.ω.  $ts' \geq ts$
- Αν  $ts' = ts$  και  $\text{maxTS} = ts$  τότε ο  $\rho$  θα παρατηρήσει  $q\text{View}(1)$  και θα επιστρέψει  $v$ 
  - Όλοι απαντούν με την ίδια χρονοσφραγίδα  $ts$
- Αν  $ts' > ts$  τότε  $\text{maxTS} = ts'$  και ο  $\rho$  θα επιστρέψει τουλάχιστον την τιμή που αντιστοιχεί στην χρονοσφραγίδα  $\text{maxTS}-1 \geq ts$ 
  - Αρα νεότερη ή ίση με  $v$

## Λήμμα 2: Περίπτωση 2 ( $Q_i \neq Q_j$ )

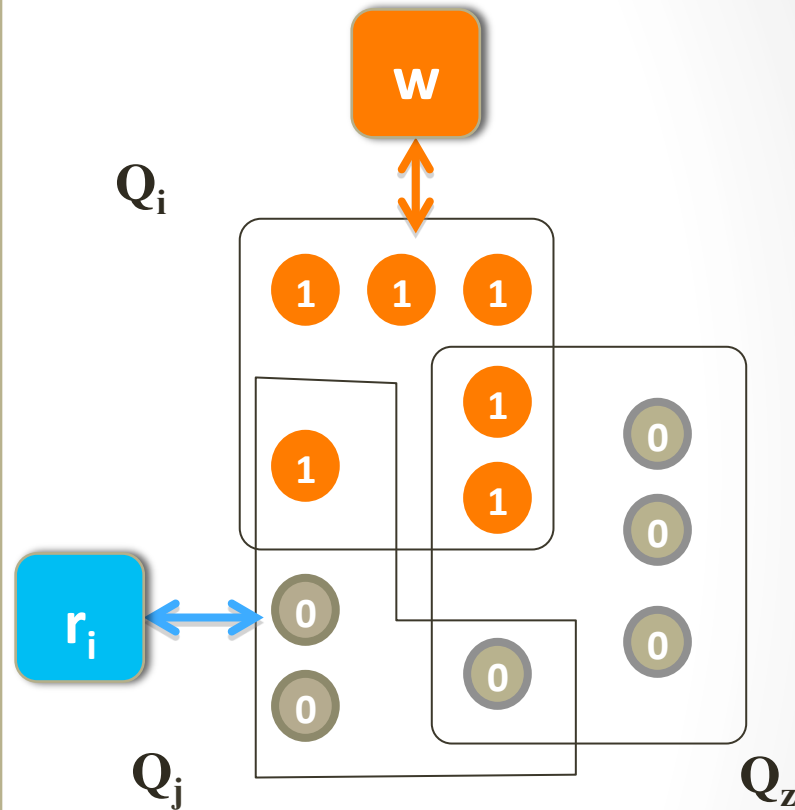
- Όλοι οι εξυπηρετητές στο  $Q_i \cap Q_j$  έλαβαν μήνυμα  $\text{write}(\langle ts, v \rangle)$  από τον  $w$  πριν απαντήσουν στην ανάγνωση  $\rho$
- Από το Λήμμα 1 θα στείλουν  $\text{reply}(\langle ts', v' \rangle)$  στον  $\rho$  τ.ω.  $ts' \geq ts$
- Αν  $ts' = ts$  και  $\text{maxTS} = ts$  τότε ο  $\rho$  θα παρατηρήσει  $\text{qView}(3)$  και θα επιστρέψει  $v$  αφού πρώτα εκτελέσει 2<sup>η</sup> φάση
  - Όλοι σε μια τομή απαντούν με την ίδια χρονοσφραγίδα  $ts$
- Αν  $ts' > ts$  τότε  $\text{maxTS} = ts'$  και ο  $\rho$  θα επιστρέψει τουλάχιστον την τιμή που αντιστοιχεί στην χρονοσφραγίδα  $\text{maxTS}-1 \geq ts$ 
  - Άρα νεότερη ή ίση με  $v$

# Σχηματικά



$Q_i = Q_j \Rightarrow qView(1)$

Returns maxTS=1 in **one round**



$Q_i \neq Q_j \Rightarrow qView(3)$

Returns maxTS=1 in **two rounds**

# Ορθότητα Αλγορίθμου

- **Λήμμα 3:** Αν μία Ανάγνωση  $\rho_1$  προηγείται μιας ανάγνωσης  $\rho_2$ , τότε η  $\rho_2$  επιστρέφει την ίδια ή νεότερη τιμή από αυτή που επιστρέφει η  $\rho_1$
- Απόδειξη: Έστω ότι
  - $\rho_1$  διάβασε από  $Q_i$  και επέστρεψε  $\langle ts, v \rangle$  και
  - $\rho_2$  διάβασε από  $Q_j$  και επέστρεψε  $\langle ts', v' \rangle$
  - Και οι δύο αναγνώσεις προέρχονται από τον ίδιο αναγνώστη
    - Τότε  $ts' \geq ts$ . (Γιατί;)
  - Αν οι δύο αναγνώσεις προέρχονται από 2 διαφορετικούς αναγνώστες πρέπει να μελετήσουμε τις περιπτώσεις όπου: 1)  $\rho_1$  αργή, και 2)  $\rho_1$  γρήγορη



## Λήμμα 3: Περίπτωση 1 (ρ1 αργή)

- Αν ρ1 είναι αργή σημαίνει ότι πριν ολοκληρωθεί στέλνει  $\text{read}(\langle ts, v \rangle)$  στα μέλη μιας απαρτίας (έστω  $Q_z$ )
- Από το Λήμμα 1 όλοι οι εξυπηρετητές στη  $Q_j \cap Q_z$  θα στείλουν  $\text{reply}(\langle ts', v' \rangle)$  στον ρ τ.ω.  $ts' \geq ts$
- Άρα ο ρ2 θα παρατηρήσει  $\text{maxTS} \geq ts'$  με το τέλος του πρώτου γύρου
- Αν  $\text{maxTS} > ts$  τότε στην χειρότερη περίπτωση ο ρ2 θα επιστρέψει την τιμή που αντιστοιχεί στο  $\text{maxTS} - 1 \geq ts$
- Αν  $\text{maxTS} = ts$  τότε ο ρ2 θα επιστρέψει  $\langle \text{maxTS}, v \rangle = \langle ts, v \rangle$  αφού θα παρατηρήσει είτε
  - $\text{qView}(1)$  if  $Q_j = Q_z$ , ή
  - $\text{qView}(3)$  if  $Q_j \neq Q_z$  αφού όλοι οι εξυπηρετητές στο  $Q_j \cap Q_z$  στέλνουν  $\text{reply}(\langle ts, v \rangle)$

## Λήμμα 3: Περίπτωση 2 (ρ1 γρήγορη)

- Αν ρ1 είναι γρήγορη σημαίνει στον πρώτο της γύρο προσεξε
  - 2.1) qView(1) και επέστρεψε  $\max TS = ts$
  - 2.2) qView(2) και επέστρεψε  $\max TS - 1 = ts$
- Περίπτωση 2.1
  - Παρατήρησε qView(1)  $\Rightarrow$  όλα τα μέλη του  $Q_i$  απάντησαν με  $\langle ts, v \rangle$
  - Επομένως από το Λήμμα 1 όλα τα μέλη της  $Q_i \cap Q_j$  θα στείλουν  $\text{reply}(\langle ts', v' \rangle)$  στον ρ τ.ω.  $ts' \geq ts$
  - Άρα ο ρ2 θα παρατηρήσει  $\max TS \geq ts'$  με το τέλος του πρώτου γύρου
  - Με την ίδια λογική όπως και στην περίπτωση 1 μπορούμε να δείξουμε ότι ο ρ2 επιστρέφει  $\langle ts', v' \rangle$  s.t.  $ts' \geq ts$

# Λήμμα 3: Περίπτωση 2 (ρ1 γρήγορη)

- Περίπτωση 2.2
  - Παρατήρησε  $qView(3) \Rightarrow$  ο ρ1 επέστρεψε  $maxTS-1=ts$
  - Κάθε διεργασία τελειώνει μια λειτουργία πριν μεταβεί στην επόμενη.
  - Επομένως αφού
    - ο ρ1 παρατήρησε την χρονοσφραγίδα  $maxTS$  στο σύστημα
    - ο εγγραφέας χρησιμοποιεί μια χρονοσφραγίδα ανα εγγραφή,
    - Σε κάθε εγγραφή αυξάνεται η χρονοσφραγίδα κατά 1, και
    - έχουμε μόνο ένα εγγραφέα

έπεται ότι η εγγραφή που χρησιμοποίησε τη χρονοσφραγίδα  $maxTS-1$  έχει ήδη ολοκληρωθεί πριν ή κατά τη διάρκεια του ρ1.

  - Αφού το ρ1 προηγείται της ρ2 επομένως και η εγγραφή με  $maxTS-1$  προηγείται της ρ2
  - Επομένως από το Λήμμα 2 το ρ2 θα επιστρέψει  $ts' \geq maxTS-1 \Rightarrow ts' \geq ts$ .

# Ερωτήσεις;

