

Επεξεργασία Δοσοληψιών (συνέχεια)

Πρόβλημα

«Σωστή» εκτέλεση προγραμμάτων όταν επιτρέπουμε *ταυτοχρονισμό* και ακόμα και αν υπάρχουν *αποτυχίες*

Δοσοληψία (transaction)

Η εκτέλεση ενός προγράμματος σε μια βάση δεδομένων. Μια ακολουθία από πράξεις R, W που τελειώνει ή με Commit ή με Abort

- **BEGIN**
- **R(X) W(X)**
- **END**
- **COMMIT** (επικύρωση) - επιτυχία - όλες οι τροποποιήσεις επικυρώνονται και δεν μπορούν να αναιρεθούν
- **ABORT** (ακύρωση ή ανάκληση) - αποτυχία - όλες οι τροποποιήσεις πρέπει να αναιρεθούν

Ιδιότητες Δοσοληψιών

- **Atomicity** (ατομικότητα) - είτε όλες οι πράξεις είτε καμία
- **Consistency** (συνέπεια) - διατήρηση συνέπειας της ΒΔ
- **Isolation** (απομόνωση) - δεν αποκαλύπτει ενδιάμεσα αποτελέσματα
- **Durability** (μονιμότητα ή διάρκεια) - μετά την επικύρωση μιας δοσοληψίας οι αλλαγές δεν είναι δυνατόν να χαθούν

- **Atomicity** (ατομικότητα) → ΤΕΧΝΙΚΕΣ ΑΝΑΚΑΜΨΕΙΣ
- **Consistency** (συνέπεια) → ΥΠΕΥΘΥΝΟΤΗΤΑ ΤΟΥ ΠΡΟΓΡΑΜΜΑΤΙΣΤΗ
- **Isolation** (απομόνωση) → ΕΛΕΓΧΟΣ ΣΥΝΔΡΟΜΙΚΟΤΗΤΑΣ
- **Durability** (μονιμότητα ή διάρκεια) → ΤΕΧΝΙΚΕΣ ΑΝΑΚΑΜΨΕΙΣ

Ισοδύναμα Χρονοπρογράμματα βάσει Συγκρούσεων:

Δυο χρονοπρογράμματα είναι ισοδύναμα βάσει συγκρούσεων αν η διάταξη κάθε ζεύγους συγκρουόμενων πράξεων είναι ίδια και στα δυο χρονοπρογράμματα.

Ένα χρονοπρόγραμμα είναι σειριοποιήσιμο (βάσει συγκρούσεων) αν και μόνο αν ο γράφος προήγησής του είναι ακυκλικός.

Επεξεργασία Δοσοληψιών

Άσκηση 1

Ο γράφος προήγησης για ένα συγκεκριμένο χρονοπρόγραμμα περιέχει μόνο τις παρακάτω ακμές και κόμβους. Πόσα διαφορετικά σειριακά χρονοπρόγραμμα είναι ισοδύναμα με αυτό το χρονοπρόγραμμα;



ΠΑΡΑΤΗΡΗΣΗ: Διαφορετικά χρονοπρόγραμμα: Έστω T_i με m_i πράξεις
 $(m_1 + m_2 + \dots + m_N)!$ / $(m_1! m_2! \dots m_N!)$

Επανάληψη: Χρονοπρόγραμμα και Δυνατότητα Ανάκαμψης

• Χρονοπρόγραμμα με δυνατότητα ανάκαμψης (recoverable)

αν καμία δοσοληψία T στο S δεν επικυρώνεται έως ότου επικυρωθούν όλες οι δοσοληψίες οι οποίες τροποποίησαν ένα δεδομένο που διαβάζει η T

• Χρονοπρόγραμμα χωρίς διάδοση ανακλήσεων (avoids cascading aborts)

αν κάθε δοσοληψία T στο S διαβάζει μόνο στοιχεία που έχουν γραφεί από επικυρωμένες δοσοληψίες

• Αυστηρά Χρονοπρόγραμμα (strict)

οι δοσοληψίες δεν μπορούν ούτε να διαβάσουν ούτε να γράψουν ένα στοιχείο X έως ότου επικυρωθεί η δοσοληψία που έγραψε το X

Τεχνικές Ελέγχου Συνδρομικότητας

- Ο χρήστης δεν ασχολείται με τη συνδρομικότητα
- Το ΣΔΒΔ εξασφαλίζει «ασυτή συνδρομικότητα»: *δρομολογεί* τις πράξεις των δοσοληψιών ώστε να προκύπτουν χρονοπρόγραμμα σειριοποιησίμα βάσει συγκρούσεων
- Μέσω τεχνικών ελέγχου συνδρομικότητας

Τεχνικές Ελέγχου Συνδρομικότητας

Τεχνικές

1. **Κλειδώματος (locking)** για να αποτρέψουν τη συνδρομική (ταυτόχρονη) προσπέλαση των δεδομένων από πολλές δοσοληψίες
2. **Διάταξης χρονοσημάτων (timestamps)**
3. **Πιστοποίησης (validation)** μιας δοσοληψίας (αισιόδοξα πρωτόκολλα)

Τεχνικές Κλειδώματος

Δύο ειδών κλειδιά:

- **διαμοιραζόμενο (shared)** κλειδί ή κλειδί ανάγνωσης
- **αποκλειστικό (exclusive)** κλειδί ή κλειδί εγγραφής

	S-Lock(X)	X-Lock(X)
S-Lock(X)	✓	
X-Lock(X)		

Πρωτόκολλο κλειδώματος δυο φάσεων (Two-Phase Locking 2PL)

Όλες οι πράξεις (αιτήσεις) κλειδώματος μιας δοσοληψίας προηγούνται της πρώτης πράξης (αίτησης) άρσης κλειδώματος της διαδικασίας

- Δυνατή η «αναβάθμιση» κλειδιών

Τεχνικές Κλειδώματος

- Κάθε δοσοληψία δυο φάσεις
 - μια φάση επέκτασης ή εξάπλωσης
 - μια φάση συρρίκνωσης

• Αποδεικνύεται ότι είναι σωστό: παράνε χρονοπρόγραμμα σειριοποιησίμα με συγκρούσεις. Υπάρχουν σειριοποιησίμα με συγκρούσεις χρονοπρόγραμμα που δεν παράγονται με κλειδίμα δυο φάσεων;

• Παράγεται αυστηρό χρονοπρόγραμμα: Όχι.

Παραλλαγές:

Αυστηρό Κλειδίμα Δύο Φάσεων: μια δοσοληψία δεν απελευθερώνει κανένα αποκλειστικό κλειδί πριν επικυρωθεί ή ακυρωθεί

Υπερ-αυστηρό Κλειδίμα Δύο Φάσεων: μια δοσοληψία δεν απελευθερώνει κανένα (αποκλειστικό και διαμοιραζόμενο) κλειδί πριν επικυρωθεί ή ακυρωθεί

Τα πιο δημοφιλή πρωτόκολλα - ευρεία χρήση σε εμπορικά συστήματα

Τεχνικές Κλειδώματος

- Πρόβλημα με τη δημιουργία αδιεξόδων
- Κόστος διαχείρισης κλειδιών

Δυο τεχνικές:

• **Πρωτόκολλα Πρόληψης Αδιεξόδων** (Deadlock Prevention): Αποφυγή δημιουργίας αδιεξόδου

- όλα τα κλειδιά δίνονται μαζί
- διάταξη των δοσοληπιών (με χρήση χρονοσήμων)

• **Πρωτόκολλα Ανίχνευσης Αδιεξόδου** (Deadlock Detection): Ελέγχουμε περιοδικά αν το σύστημα βρίσκεται σε κατάσταση αδιεξόδου (π.χ με χρήση γράφου αναμονής)

Κλειδωμα Δυο Φάσεων: Άσκηση 18.17

Άσκηση

Αποδείξτε ότι το βασικό πρωτόκολλο κλειδώματος δύο φάσεων εξασφαλίζει σειριοποιησιμότητα βάσει συγκρούσεων

Με απαγωγή στο άτοπο

Έστω ότι δεν το εξασφαλίζει, δηλαδή κλειδωμα δυο φάσεων και μη σειριοποιησιμότητα βάσει συγκρούσεων \Rightarrow

Κύκλος στο γράφο σειριοποιησιμότητας

Κλειδωμα Δυο Φάσεων: Άσκηση

Άσκηση

Θεωρείστε το παρακάτω χρονοπρόγραμμα των δοσοληπιών T_1 και T_2 :

S: b $W_1(X)$ $R_2(X)$ $W_2(Y)$ C_2 $W_1(Z)$ C_1 e

b αρχή και e τέλος του χρονοπρογράμματος. Υποθέστε 2PL.

Πότε -- μεταξύ ποιών γεγονότων μπορούν να συμβούν τα παρακάτω

π.χ., T_1 μπορεί να κλειδώσει το X μεταξύ των b και $W_1(X)$

Αντιμετώπιση Αδιεξόδων

Δυο τεχνικές:

• **Πρωτόκολλα Πρόληψης Αδιεξόδων** (Deadlock Prevention): Αποφυγή δημιουργίας αδιεξόδου - *χρονόσημα (προτεραιότητες) σε κάθε δοσοληψία, επανεκκίνηση δοσοληψίας με το ίδιο χρονόσημα*

• *Αναμονή - θανάτωση*

• *Τραυματισμός αναμονής*

• **Πρωτόκολλα Ανίχνευσης Αδιεξόδου** (Deadlock Detection): Ελέγχουμε περιοδικά αν το σύστημα βρίσκεται σε κατάσταση αδιεξόδου *γράφος αναμονής*

Κλειδωμα Δυο Φάσεων: Άσκηση

S: b $W_1(X)$ $R_2(X)$ $W_2(Y)$ C_2 $W_1(Z)$ C_1 e

Πότε -- μεταξύ ποιών γεγονότων μπορούν να συμβούν τα παρακάτω

(α) Πότε μπορεί η T_1 να κάνει unlock(X)

(β) Πότε μπορεί η T_1 να κάνει lock(Z)

Αντιμετώπιση Αδιεξόδων: Άσκηση

Άσκηση

Θεωρείστε ότι οι πράξεις υποβάλλονται με την παρακάτω σειρά. Θεωρείστε ότι οι εντολές έρχονται στο scheduler με την παρακάτω σειρά. Περιγράψτε για το καθένα από τα παρακάτω πρωτόκολλα πως χειρίζεται τις πράξεις ο μηχανισμός ελέγχου ταυτοχρονισμού

$R_1(X)$ $W_2(X)$ $W_2(Y)$ $W_3(Y)$ $W_1(Y)$ C_1 C_2 C_3

$R_1(X)$ $W_2(Y)$ $W_2(X)$ $W_3(Y)$ $W_1(Y)$ C_1 C_2 C_3

(α) Αυστηρό 2PL με χρήση χρονοσημάτων για αποφυγή (πρόληψη) αδιεξόδου και αναμονή-θανάτωση (το κλειδί το έχει με μικρότερο χρονόσημα (παλιότερη) - ακυρώνεται, αλλιώς περιμένει)

(β) Αυστηρό 2PL με χρήση γράφου αναμονής για ανίχνευση αδιεξόδου

Διάταξη Χρονοσημάτων

Το χρονοσημα δημιουργείται από το ΣΔΒΔ και προσδιορίζει μοναδικά μια δοσοληψία

Ιδέα: διάταξη των δοσοληψιών με βάση το χρονοσημα τους (δηλαδή, χρονοπρόγραμμα ισοδύναμο με σειριακό στο οποίο οι δοσοληψίες εμφανίζονται διατεταγμένες με βάση τις τιμές των χρονοσημάτων)

⇒ άρα η σειρά προπέλασης στα δεδομένα πρέπει να μη παραβιάζει τη σειριοποιησιμότητα

Δηλαδή: αν μια πράξη a_i μιας δοσοληψίας T_i συγκρούεται με μια πράξη a_j μιας δοσοληψίας T_j και $TS(T_i) < TS(T_j)$, τότε η a_i πρέπει να προηγείται της a_j . Αλλιώς, restart τη δοσοληψία.

Διάταξη Χρονοσημάτων

Κάθε δεδομένο X έχει δύο τιμές χρονοσημάτων:

$ΧΣΑ(X)$ (χρονοσημα ανάγνωσης) το μεγαλύτερο μεταξύ όλων των χρονοσημάτων των δοσοληψιών που διάβασαν το X (διαισθητικά, η πιο πρόσφατη που το διάβασε)

$ΧΣΕ(X)$ (χρονοσημα εγγραφής) το μεγαλύτερο μεταξύ όλων των χρονοσημάτων των δοσοληψιών που έγραψαν το X

Κάθε πράξη ανάγνωσης και εγγραφής μιας δοσοληψίας συνοδεύεται από κατάλληλους ελέγχους του χρονοσημάτος της και των χρονοσημάτων του αντίστοιχου δεδομένου.

Αν ο έλεγχος αποτύχει, η δοσοληψία απορρίπτεται και αρχίζει πάλι

Διάταξη Χρονοσημάτων

- Παράγει σωστά χρονοπρογράμματα
- Δε δημιουργούνται αδιέξοδα (αφού καμία δοσοληψία δεν περιμένει)
- Μπορεί να παραχθούν χρονοπρογράμματα που δεν έχουν δυνατότητα ανάκαμψης (not recoverable)

Παραλλαγή για χρονοπρόγραμμα χωρίς διάδοση ανακλήσεων (ιδέα: commit bit με κάθε δεδομένο - καθυστέρηση αναγνώσεων)

Παραλλαγή για ασυζητά (αnamονή αναγνώσεων και εγγραφών)

Άλλη ιδέα ($R(X)$ ή $W(X)$ από T με $ΧΣ(T) > ΧΣΕ(X)$) (αν μικρότερο ακυρώνεται) περιμένει μέχρι να επικυρωθεί ή να ακυρωθεί η δοσοληψία που έγραψε το X (όχι αδιέξοδα)

Επιβάρυνση διαχείρισης χρονοσημάτων

Απόρριψη και επανεκκίνηση δοσοληψιών

Διάταξη Χρονοσημάτων: Άσκηση 7

Άσκηση

Εφαρμόστε τον αλγόριθμο διάταξης χρονοσημάτων στα χρονοπρογράμματα (α), (β), και (γ)

(α)

$R_1(X) W_2(Y) W_2(X) W_3(X) W_1(Y) C_1 C_2 C_3$

Διάταξη Χρονοσημάτων: Άσκηση 7(β)

T1	T2	T3
	$R_2(Z)$	
	$R_2(Y)$	
	$W_2(Y)$	
$R_1(X)$		$R_3(Y)$
$W_1(X)$		$R_3(Z)$
		$W_3(Y)$
		$W_3(Z)$
$R_1(Y)$	$R_2(X)$	
$W_1(Y)$		
	$W_2(X)$	
		• σειριοποίησησημο;

Διάταξη Χρονοσημάτων: Άσκηση 7(γ)

T1	T2	T3
		$R_3(Y)$
		$R_3(Z)$
$R_1(X)$		
$W_1(X)$		
	$R_2(Z)$	$W_3(Y)$
		$W_3(Z)$
$R_1(Y)$		
$W_1(Y)$		
	$R_2(Y)$	
	$W_2(Y)$	
	$R_2(X)$	
	$W_2(X)$	
		• σειριοποίησησημο;

Οι τεχνικές κλειδώματος είναι συντηρητικές (αποφεύγονται οι συγκρούσεις)

Μειονεκτήματα

- επιβάρυνση (overhead) χειρισμού κλειδώματος
- αποφυγή/ανίχνευση αδιεξόδων
- lock contention για τα δεδομένα που χρησιμοποιούνται συχνά

Αν οι συγκρούσεις είναι σπάνιες, μεγαλύτερη συγχρονικότητα, αν αντί για κλείδωμα, έλεγχος για συγκρούσεις όταν μια δοσοληψία επικυρώνεται (commits)

Κάθε δοσοληψία έχει τρεις φάσεις

- **ΑΝΑΓΝΩΣΗ:** η δοσοληψία διαβάζει από τη βδ, αλλά τροποποιεί προσωπικά αντίγραφα των δεδομένων
- **ΠΙΣΤΟΠΟΙΗΣΗ:** έλεγχος για συγκρούσεις
- **ΕΓΓΡΑΦΗ:** γράφει τα τοπικά αντίγραφα στη βδ

Κατά την πιστοποίηση μιας δοσοληψίας ελέγχουμε τα σύνολα

ReadSet(T): το σύνολο των δεδομένων που διάβασε η T

WriteSet(T): το σύνολο των δεδομένων που διάβασε η T

με τα αντίστοιχα σύνολα όλων των ταυτόχρονα εκτελούμενων δοσοληψιών του συστήματος

Έλεγχος συνηκών που είναι ικανές για να εξασφαλίσουν ότι δεν υπήρχαν συγκρούσεις

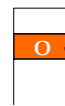
- Κάθε δοσοληψία T έχει ένα μοναδικό αριθμό TID (χρονόσημα)
- Το TID ανατίθεται στο τέλος της φάσης ΑΝΑΓΝΩΣΗΣ (ακριβώς πριν αρχίσει η πιστοποίηση)
- Με κάθε δοσοληψία

ReadSet(T): το σύνολο των δεδομένων που διάβασε η T

WriteSet(T): το σύνολο των δεδομένων που διάβασε η T

Κάθε εγγραφή παράγει ένα καινούργιο αντίγραφο ενώ οι αναγνώσεις διαβάζουν ένα κατάλληλο παλιό

ΚΥΡΙΟ ΤΜΗΜΑ (Τρέχουσες εκδόσεις των αντικειμένων της ΒΔ)



ΣΥΛΛΟΓΗ ΕΚΔΟΣΕΩΝ VERSION POOL (Παλιές εκδόσεις που μπορεί να είναι χρήσιμες σε κάποιους ενεργούς αναγνώστες)



Οι αναγνώστες μπορούν πάντα να προχωρήσουν -- αλλά μπορεί να χρειαστεί να περιμένουν μέχρι την επικύρωση των εγγραφών

Διατήρηση παλαιότερων εκδοχών, χρήσιμη ως ιστορικό εξέλιξης των τιμών (χρονικές (temporal) βάσεις δεδομένων)

Βασική ιδέα:

οι αναγνώσεις διαβάζουν κατάλληλες εκδόσεις ώστε να διατηρήσουν τη σειριοποιησιμότητα

Παραλλαγές των τεχνικών ελέγχου συγχρονικότητας

Θα δούμε με χρονοσήματα

- **Κάθε έκδοση** ενός αντικειμένου έχει ως ΧΣΕ το ΧΣ της T που την έγραψε και ως ΧΣΑ το ΧΣ της T που διάβασε αυτήν την έκδοση πιο πρόσφατα

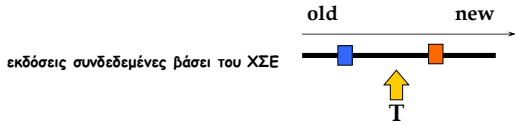
Μια πράξη εγγραφής μια δοσοληψίας T δημιουργεί μια καινούργια έκδοση του δεδομένου, με ΧΣΑ και ΧΣΕ ίσο με ΧΣ(T)

- Οι εκδόσεις συνδέονται (chained backward) μπορούμε να σηθίσουμε εκδόσεις που είναι πολύ «παλιές για να έχουν ενδιαφέρον»

- Κάθε δοσοληψία: **Reader** ή **Writer**.
 - **Writer** μπορεί να γράφει κάποιο αντικείμενο; **Reader** μόνο διαβάσει.
 - Κάθε δοσοληψία δηλώνει αν είναι **Reader** όταν ξεκινά. **Στόχος είναι η Readers να μην αποτυγχάνουν**

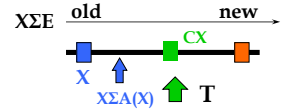
Reader T

- Για την ανάγνωση ενός αντικειμένου:
 - Βρες την **πιο νεώτερη έκδοση με $ΧΣΕ < ΧΣ(T)$** (Η αναζήτηση ξεκινά από την πιο πρόσφατη - κοιτά σχήμα)
 - Αν υποθέσουμε ότι όλες οι εκδόσεις υπάρχουν, μια Reader T δε χρειάζεται να επανικηθεί.



Writer T

- Για την ανάγνωση, ακολούθησε το πρωτόκολλο του Reader
- Για να γράφεις ένα αντικείμενο X:
 - Βρες τη νεώτερη έκδοση με $ΧΣΕ(X) < ΧΣΕ(T)$
 - Αν $ΧΣΑ(X) < ΧΣ(T)$, η T δημιουργεί ένα αντίγραφο CX του X, με ένα δείκτη στο X, και $ΧΣΕ(CX) = ΧΣ(T)$, $ΧΣΑ(CX) = ΧΣ(T)$. (Οι εγγραφές γίνονται buffered μέχρι την επικύρωση της T, οι άλλες δοσοληψίες μπορούν να δουν τις τιμές των ΧΣ αλλά δεν μπορούν να διαβάσουν την έκδοση CX μέχρι η T να επικυρωθεί.)
 - Αλλιώς, απόρριψη.



Δυναμικές Βάσεις Δεδομένων

Αν επιτρέψουμε εγγραφές και διαγραφές στοιχείων, τα πρωτόκολλα δε δουλεύουν σωστά (ούτε το αυστηρό 2PL)

Sailors(sid, sname, rating, age)
Boats(bid, bname, color)
Reserves(sid, bid, day)

T1: βρίσκει το **γνηραιότερο** ναύτη με rating 1 και 2
T2: εισάγει ναύτη με rating 1 και age 96 και διαγράφει τον **πιο ηλικιωμένο** ναύτη με rating 2

- T1 κλειδώνει όλες τις σελίδες που περιέχουν εγγραφές sailor με rating = 1, και βρίσκει τον **πιο ηλικιωμένο** (έστω, age = 71).
 - Μετά, η T2 εισάγει ένα νέο sailor; rating = 1, age = 96.
 - T2 επίσης διαγράφει τον πιο ηλικιωμένο ναυτικό με rating = 2 (έστω, age = 80), και επικυρώνεται.
 - T1 τώρα κλειδώνει όλες τις σελίδες που περιέχουν εγγραφές sailor με rating = 2, βρίσκει τον **πιο ηλικιωμένο** (έστω, age = 63).
- Αποτέλεσμα: 71 63
- Δεν υπάρχει ισοδύναμη με την παραπάνω σειριακή εκτέλεση

Δυναμικές Βάσεις Δεδομένων

Sailors(sid, sname, rating, age)
Boats(bid, bname, color)
Reserves(sid, bid, day)

T1: βρίσκει το **γνηραιότερο** ναύτη με rating 1 και 2
T2: εισάγει ναύτη με rating 1 και age 96 και διαγράφει τον **πιο ηλικιωμένο** ναύτη με rating 2

- T1 κλειδώνει όλες τις σελίδες που περιέχουν εγγραφές sailor με rating = 1, και βρίσκει τον **πιο ηλικιωμένο** (έστω, age = 71).
- Μετά, η T2 εισάγει ένα νέο sailor; rating = 1, age = 96.
- T2 επίσης διαγράφει τον πιο ηλικιωμένο ναυτικό με rating = 2 (έστω, age = 80), και επικυρώνεται.
- T1 τώρα κλειδώνει όλες τις σελίδες που περιέχουν εγγραφές sailor με rating = 2, βρίσκει τον **πιο ηλικιωμένο** (έστω, age = 63).

Αποτέλεσμα 71 63

- Δεν υπάρχει ισοδύναμη με την παραπάνω σειριακή εκτέλεση

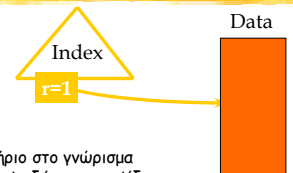
T1 T2: αποτέλεσμα 71 80
T2 T1: αποτέλεσμα 96 63

Δυναμικές Βάσεις Δεδομένων

- Η T1 έμμεσα υποθέτει ότι έχει κλειδώσει το σύνολο όλων των εγγραφών με rating = 1. (εγγραφές **φαντάσματα** σε άλλες σελίδες)
 - Η υπόθεση ισχύει μόνο αν δεν προστεθούν εγγραφές ενώ εκτελείται η T1
 - Χρειάζεται κάποιος μηχανισμός για να το επιβάλει: **κλειδωμά ευρετηρίου** (index locking) και **κλειδωμά συνθηκών** (predicate locking.)

Τα παράδειγμα δείχνει ότι η **σειριοποιησιμότητα** βάσει συγκρούσεων δίνει **σειριοποιησιμότητα** μόνο όταν τα αντικείμενα στη βδ είναι σταθερά

Δυναμικές Βάσεις Δεδομένων



Κλειδωμά Ευρετηρίου

- Αν υπάρχει ένα πυκνό ευρετήριο στο γνώρισμα rating, η T1 θα μπορούσε να κλειδώσει τη σελίδα του ευρετηρίου που περιέχει τις εγγραφές με rating = 1.
 - Αν δεν υπάρχει καμία εγγραφή με rating = 1

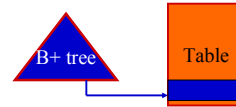
Κλειδωμα Συνθήκης

- Κλειδώσε όλες τις εγγραφές που ικανοποιούν κάποια συνθήκη λογικής, π.χ., $age > 2 * salary$.
- Το κλειδωμα ευρετηρίου είναι μια ειδική μορφή κλειδώματος συνθήκης όπου το ευρετήριο επιτρέπει αποδοτική υλοποίηση του κλειδώματος συνθήκης

ΑΡΑ:

Αν είχαμε ένα index (π.χ., B+ tree) πάνω στο πεδίο rating, θα μπορούσαμε να κλειδώνουμε γρήγορα όλες τις εγγραφές με $RATING=1$.

- ☑ Ακόμα κι αν δεν υπάρχουν τέτοιες εγγραφές, κλειδώνοντας τον index, μπορούμε να απαγορεύσουμε να δημιουργηθούν από άλλους κατά τη διάρκεια της δοσοληψίας
- ☑ Ειδική περίπτωση του κλειδώματος κατηγορήματος



Το κλειδωμα του ευρετηρίου είναι μια ειδική μορφή δεδομένων, πως θα γίνει το κλειδωμα:

- ☞ Ας υποθέσουμε ότι θέλουμε να εισάγουμε μια εγγραφή με $RATING = 1$ στη διάρκεια της T1
- ☞ Κατεβαίνοντας από τη ρίζα προς τα φύλλα του B+ δέντρου, μπορούμε να κλειδώνουμε κάθε κόμβο με X-lock [του φύλλου συμπεριλαμβανομένου]

ΠΡΟΒΛΗΜΑ: ήδη από το πρώτο βήμα έχουμε κλειδώσει τη ρίζα, δηλαδή, **ΟΛΟ** το δέντρο.

ΠΑΡΑΤΗΡΗΣΗ

Στα B+ δέντρα, οι εσωτερικοί κόμβοι χρησιμεύουν **ΜΟΝΟ** ως μονοπάτια για τα φύλλα, όλα τα πραγματικά δεδομένα είναι στα φύλλα [βασική διαφορά με τα B δέντρα]

Δηλαδή, «συνήθως» οι εσωτερικοί κόμβοι δεν τροποποιούνται, εκτός ΑΝ;

- ☞ Αν κλειδωνα με S-lock τους ενδιάμεσους κόμβους και με X-lock τα φύλλα?

Εκτός, αν χρειαστεί να διασπάσω ή να συνενώσω κάποιο εσωτερικό κόμβο

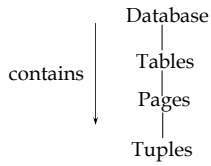
Πότε;

Αν overflow/underflow του παιδιού του!

- ☞ **Αναζητήσεις:** Ξεκινώντας από τη ρίζα κλειδώνουμε με S-lock κάθε κόμβο. Μετά, ξεκλειδώνουμε τον γονέα του.
- ☞ **Εισαγωγές:** Ξεκινώντας από τη ρίζα κλειδώνουμε με X-locks κατά περίπτωση. Μετά ελέγχουμε αν ο κλειδωθείς κόμβος είναι **ασφαλής**:
 - ☑ Αν ναι, ξεκλειδώνουμε τους προγόνους του
- ☞ **Ασφαλής κόμβος:** οι όποιες αλλαγές δεν θα διαδοθούν προς τα πάνω από τον κόμβο αυτό
 - ☑ INS: ο κόμβος δεν είναι πλήρης (όχι overflow)
 - ☑ DEL: ο κόμβος δεν είναι άδειος κατά το ήμισυ. (όχι underflow)

Διακριτότητα Κλειδώματος

- Μέχρι τώρα θεωρούσαμε **εγγραφές**
- Στην πράξη, συνήθως μιλάμε για **σελίδες**
- Μπορούμε να κλειδώσουμε τη βάση σε πολλά «επίπεδα» που ορίζουν μια «**ιεραρχία διακριτότητας**»



Ορισμοί

- Διακριτότητα** (granularity): πόσο μικρό είναι ένα αντικείμενο υπό παρατήρηση
- Μικρή διακριτότητα = μεγάλο μέγεθος στοιχείου (π.χ., πίνακας ή ΒΔ)
- Μεγάλη διακριτότητα = μικρό μέγεθος στοιχείου (π.χ., σελίδα ή εγγραφή)
- Πρόγονος ενός αντικειμένου**: οποιοδήποτε αντικείμενο στην ιεραρχία διακριτότητας το περιλαμβάνει (π.χ., η ΒΔ είναι πρόγονος μιας σελίδας)

Κλειδωμα σε διαφορετικά επίπεδα διακριτότητας

- Κλειδωμα σε **μικρή διακριτότητα**: αποκλεισμός και καθυστέρηση των δοσοληψιών που τρέχουν εν παραλλήλω (π.χ., τι θα γίνει αν κλειδώσω όλη τη ΒΔ?)
- Κλειδωμα σε **μεγάλη διακριτότητα**: μεγάλος αριθμός κλειδωμάτων (χώρος μνήμης, καθυστέρηση διαχείρισης)

Εκμετάλλευση της ιεραρχίας

- Η ιδέα είναι να κλειδώνουμε με κατάλληλο τρόπο σε διαφορετικά επίπεδα διακριτότητας
- Χρειαζόμαστε ένα νέο είδος κλειδώματος
- Κλειδωμα πρόθεσης (intention lock)**: εφαρμόζεται σε ένα πρόγονο, χαρακτηρίζοντας τι πράξη (για την ακρίβεια: κλειδωμα) θα εφαρμοστεί στον απόγονό του
- Δύο ειδών κλειδωματα πρόθεσης
 - Πρόθεση αποκλειστικού κλειδιού
 - Πρόθεση διαμοιραζόμενου κλειδιού

Κλειδωμα Πολλαπλής Κλιμάκωσης

Ιεραρχική δομή

Νέα είδη κλειδίων **IS** (πρόθεση διαμοιραζόμενου) **IX** (πρόθεση αποκλειστικό)

Για να κλειδώσει με **S** (**X**) πρέπει να κλειδώσει όλους τους προγόνους με **IS** (**SX**) αρχίζοντας από τη ρίζα

	IS	IX	S	X
IS	T	T	T	
IX	T	T		
S	T		T	
X				

Κλειδωμα Πολλαπλής Κλιμάκωσης

Συχνά **IX** και **S** (διάβασμα όλου, τροποποίηση ενός) είδος κλειδιού **SIX**

	IS	IX	S	SIX	X
IS	T	T	T	T	
IX	T	T			
S	T		T		
SIX	T				
X					

Αλγόριθμος Κλειδωμάτων

- ☞ Το κλειδωμα ξεκινά από την κορυφή της ιεραρχίας
- ☞ Για να πάρω ένα κλειδωμα S- (ή IS-) σε ένα αντικείμενο, πρέπει να έχω IS- (ή IX-) στον πρόγονό του [αλλιώς, περιμένω]
- ☞ Για να πάρω ένα κλειδωμα X- (ή IX- ή SIX) σε ένα αντικείμενο, πρέπει να έχω IX- (ή SIX-) στον πρόγονό του [αλλιώς, περιμένω]
- ☞ Η αποδέσμευση των κλειδιών πρέπει να γίνει με την αντίστροφη σειρά (bottom-up)

Κλειδωμα Πολλαπλής Κλιμάκωσης

Δύο φάσεις
Σειρά αποδέσμευσης κλειδιών:

Υποστήριξη δοσοληπιών στην SQL-92

Δοσοληπία αυτόματα με προτάσεις που τροποποιούν είτε τη βάση δεδομένων είτε το σχήμα (πχ select, update, create table)

Commit ή Rollback (για abort)

▪ Μέγεθος διαγνωστικών

▪ Μέθοδος πρόσβασης (ή μέθοδος προσπέλασης):

READ ONLY ή

READ WRITE

▪ Επίπεδο απομόνωσης READ UNCOMMITTED, READ COMMITTED, REPEATABLE READ, SERIALIZABLE (αυστηρό 2PL)

SET TRANSACTION ISOLATION LEVEL SERIALIZABLE READ ONLY

default: SERIALIZABLE, READ WRITE

Υποστήριξη δοσοληπιών στην SQL-92

Isolation Level	Dirty Read	Unrepeatable Read	Phantom Problem
Read Uncommitted	Maybe	Maybe	Maybe
Read Committed	No	Maybe	Maybe
Repeatable Reads	No	No	Maybe
Serializable	No	No	No

Αναλυτικά στο βιβλίο

Υποστήριξη δοσοληπιών

IBM DB2, Informix, Microsoft SQL server, Sybase Αυστηρό 2PL (με παραλλαγές για επίπεδα απομόνωσης λιγότερα αυστηρά του SERIALIZABLE)

Microsoft SQL server: χρονοσήμα εγγραφής

Oracle 8: πολλαπλές εκδόσεις, ώστε οι αναγνώσεις δεν περιμένουν ποτέ

Κλειδωμα πολλαπλής κλιμάκωσης σε επίπεδο πίνακα, σελίδας και εγγραφής

Ανάκαμψη από Αποτυχίες

Είδη Αποτυχιών

Δυο κατηγορίες: καταστροφή ή όχι της μόνιμης αποθήκευσης (δίσκου)

Στόχος της ανάκαμψης: επαναφορά στην πιο πρόσφατη συνεχή κατάσταση ακριβώς πριν τη στιγμή της αποτυχίας

☞ Ατομικότητα:

☒ όχι επικυρωμένες -- undo -- αναίρεση

☞ Διάρκεια:

☒ επικυρωμένες -- redo -- επανάληψη

Στην πιο εύκολη περίπτωση

- όλες οι ενημερώσεις μιας δοσοληψίας γράφονται στον τοπικό χώρο εργασίας της δοσοληψίας
- η βδ ενημερώνεται *μόνον* αφού μια δοσοληψία φτάσει στο σημείο επικύρωσής της

Στην πραγματικότητα --Steal (flush) ::

- κάποιες ενημερώσεις μιας δοσοληψίας γράφονται στον δίσκο (βδ) πριν η δοσοληψία επικυρωθεί
- Γιατί: βελτίωση throughput
- Πρόβλημα: **Ατομικότητα** (ανάρρηση (undo) στην περίπτωση που η δοσοληψία αποτύχει)
- Χρειαζόμαστε την παλιά τιμή (**BeFore Image**)

Το αντίθετο

ετραχρονισμός των ενημερώσεων -- deferred updates (no steal) όλες οι ενημερώσεις μιας δοσοληψίας γράφονται στον τοπικό χώρο εργασίας της δοσοληψίας -- η βδ ενημερώνεται μόνον αφού μια δοσοληψία φτάσει στο σημείο επικύρωσής της

- χωρίς ανάρρηση/ με επανάληψη (no undo/redo)

Στην πραγματικότητα -- No Force ::

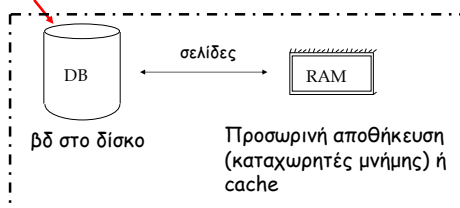
- κάποιες ενημερώσεις μιας δοσοληψίας δε γράφονται στον δίσκο (βδ) ακόμα και αφού η δοσοληψία επικυρωθεί
- Γιατί: βελτίωση χρόνου απόκρισης
- Πρόβλημα: **Διάρκεια** (επανάληψη (redo) στην περίπτωση που το σύστημα αποτύχει)
- Χρειαζόμαστε την νέα τιμή (**AFter Image**)

Ημερολόγιο Συστήματος (Log)

- Για να είναι δυνατή η ανάκαμψη από αποτυχίες, καταχωρούνται πληροφορίες για τις πράξεις των δοσοληψιών
 - Αποθηκεύονται στο δίσκο
 - Τύποι πληροφορίας: έναρξη δοσοληψίας
εγγραφή στοιχείου (παλιά, νέα τιμή)
ανάγνωση στοιχείου
επικύρωση/ακύρωση
- Αυξητικό ημερολόγιο

Ημερολόγιο Συστήματος (system log)
Πληροφορίες για κάθε δοσοληψία

το μοντέλο του συστήματος



Σημείο Ελέγχου

Γενικότερη έννοια από ότι στην ARIES

Εγγραφή όλων των τροποποιημένων σελίδων στο δίσκο

Ασαφής ή όχι (αναστολή εκτέλεσης των δοσοληψιών)

Δε χρειάζεται η επανάληψη (redo) όλων των πράξεων εγγραφής κατά την ανάκαμψη

Το πρωτόκολλο προεγγραφής ημερολογίου (Write Ahead Log)

- (1) Η καταχώρηση στο log για μια εγγραφή γράφεται στο δίσκο πριν οι αντίστοιχες σελίδες να γραφούν στο δίσκο
- (2) Πριν την επικύρωση (commit) μιας δοσοληψίας όλες οι καταχωρήσεις του log που την αφορούν γράφονται στο δίσκο

Το (1) δίνει ατομικότητα
 Το (2) δίνει διάρκεια

Σειριακή εγγραφή στο log

Για κάθε εγγραφή στοιχείου, στο log μια καταχώρηση:
 <XID, pageID, offset, length, old data, new data>



Πεδία LogRecord:

Προηγούμενο LSN για τη δοσοληψία

- prevLSN
- XID
- type
- pageID
- length
- offset
- before-image
- after-image

Μόνο για εγγραφές

Τύποι καταχωρήσεων

- Update
- Commit
- Abort
- End
- Compensation Log Records (CLRs)
 - για πράξεις UNDO



LogRecords

- prevLSN
- XID
- type
- pageID
- length
- offset
- before-image
- after-image



Σελίδες Δεδομένων

σε κάθε σελίδα
 pageLSN



Πίνακας Δοσοληψιών

lastLSN
 status

Πίνακας Τροποποιημένων Σελίδων

recLSN

flushedLSN

Τρεις Φάσεις

ΑΝΑΛΥΣΗ -- ανακατασκευή τον πινάκων του συστήματος: ποιες σελίδες είναι τροποποιημένες, ποια είναι η κατάσταση των δοσοληψιών

ΕΠΑΝΑΛΗΨΗ (REDO) -- επανάληψη όλων των δοσοληψιών

ΑΚΥΡΩΣΗ (UNDO) -- αναίρεση των δοσοληψιών που δεν ήταν επικυρωμένες

Παλιότερη καταχώρηση log δοσοληψίας ενεργής κατά την αποτυχία

Μικρότερο recLSN στον Πίνακα Τροποποιημένων Σελίδων μετά την Ανάλυση

Τελευταίο chkpt

ΑΠΟΤΥΧΙΑ



❖ Άρχισε από ένα σημείο ελέγχου (που βρίσκεται μέσω του master record).

❖ Τρεις Φάσεις.

- Εύρεση ποιες δοσοληψίες μετά το σημείο ελέγχου επικυρώθηκαν, ποιες ακυρώθηκαν (**ΑΝΑΛΥΣΗ**).
- **REDO** όλες τις δοσοληψίες.
 - ♦ (repeat history)
- **UNDO** το αποτέλεσμα των αποτυχημένων δοσοληψιών

Τεχνικές Ανάκαμψης από Σφάλματα

Άσκηση 8

Φάση της Ανάλυσης

00 begin_checkpoint	T1 abort 70	P5 20
10 end_checkpoint	T3 running 60	P3 60
20 update: T1 writes P5	Πίνακας Δοσοληψιών	Πίνακας Τροποποιημένων Σελίδων
40 T2 Commit	Φάση της Επανάληψης (REDO)	
50 T2 end	Ξεκινάμε από το 20, γιατί:	
60 update: T3 writes P3	Φάση της Αναιρέσεως (UNDO)	
70 T1 abort	ToUndo = {70, 60}	

Τεχνικές Ανάκαμψης από Σφάλματα

Στο ARIES υποθέτουμε αυστηρό 2PL, άρα όχι διάδοση ανακλήσεων, εγγραφές

Γενικά μπορεί να χρειαστεί ανάκληση

(άρα στο log εκτός των εντολών εγγραφής και εντολής ανάγνωσης, για να ελέγξουμε ποιες πρέπει να ανακληθούν)