

Transaksjonshåndtering

Del 2

Ragnar Normann

Noen figurer er basert på en original laget av
Hector Garcia-Molina

En ny type serialiseringsprotokoll

- Hittil har vi bare sett på 2PL-baserte protokoller
- Alle slike protokoller har følgende struktur:
 - først setter vi låser på de dataelementene vi er interessert i
 - så leser og (eventuelt) skriver vi disse dataelementene
 - til slutt frigir vi låsene og slipper andre transaksjoner til
- Slike protokoller egner seg godt når dataene er lagret på tabellform (som arrayer og hashtabeller)
- Den andre hovedmåten å lagre data på er å organisere dem som trær (som oftest B-trær eller B⁺-trær)
- For slike data finnes en mer hensiktsmessig protokoll, kalt **treprotokollen**

Treprotokollen

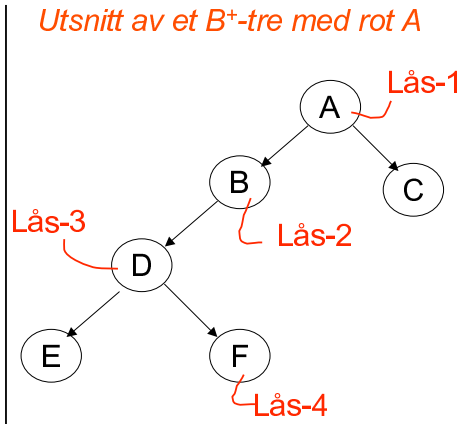
- Treprotokollen er formulert ut fra at vi bare har én låstype, men den fungerer like bra med flere:
 1. En transaksjon kan sette sin første lås på en vilkårlig node i treet
 2. Senere kan transaksjonen bare sette lås på en node hvis den har låst foreldrenoden
 3. En transaksjon kan når som helst slette en lås den har på en node
 4. En transaksjon kan ikke sette lås på en node den tidligere har frigitt (selv om den fortsatt har lås på foreldrenoden)

Treprotokollen på B-trær (B⁺-trær)

- I et B-tre er det alltid roten som låses først
- I teorien må transaksjonen holde en U-lås (evt. X-lås) på roten til den er ferdig fordi både insert og delete kan føre til at roten blir endret
- Det ville ha medført at bare én skrivetransaksjon om gangen kunne aksessere B-treet
- Imidlertid vil det normale være at når man aksesserer neste nivå i treet, så kan man straks fastslå at foreldrenoden ikke vil bli berørt, og den kan da øyeblikkelig frigies
- I praksis gir dette høy grad av samtidighet

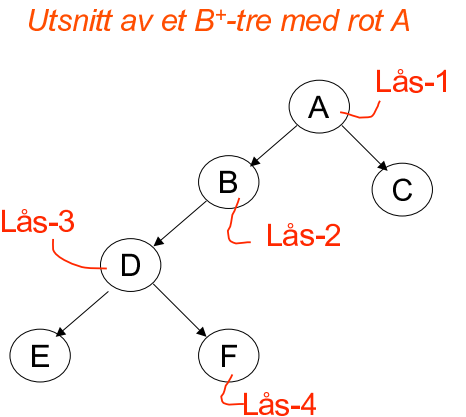
Eksempel på treprotokollen

1. T låser A (lås-1)
2. T låser B (lås-2)
3. T ser at A ikke blir endret og hever lås-1
4. T låser D (lås-3)
5. T ser at B ikke blir endret og hever lås-2
6. T låser F (lås-4)
7. T ser at D vil bli endret
8. T skriver F og hever lås-4
9. T skriver D og hever lås-3



Treprotokollen sikrer serialiserbarhet

- Anta at T_1 er som før, at T_2 også vil skrive, og at T_1 låser A først
- T_2 prøver å låse A og blir lagt i A-køen til T_1 slipper lås-1
- T_2 får låse A og kan gå videre til C, men ikke til B før T_1 har sluppet lås-2
- T_2 kan aldri «passere» T_1 , så rekkefølgen ved roten bestemmer en serialiseringsrekkefølge



Tidsstempling

- Tidsstempling er grunnlag for en familie av serialiserbarhetsprotokoller som ikke bruker låser
- Tidsstempling gir optimistisk samtidighetskontroll hvor transaksjonene får gå uhindret til man evt. oppdager at noe gikk galt slik at transaksjonen må aborteres (2PL er pessimistisk og prøver å hindre feil på forhånd)
- Når en transaksjon T starter, får den et **tidsstempel**, $TS(T)$
- To transaksjoner kan ikke ha samme tidsstempel, og hvis T_1 starter før T_2 , så skal vi ha $TS(T_1) < TS(T_2)$
- Serialiseringsrekkefølgen er bestemt av tidsstemplene
- De to vanligste formene for tidsstempler er:
 - tidspunktet da T startet (verdien av systemklokken)
 - et løpenummer (transaksjonsnr i systemets levetid)

Tidsstempler på data

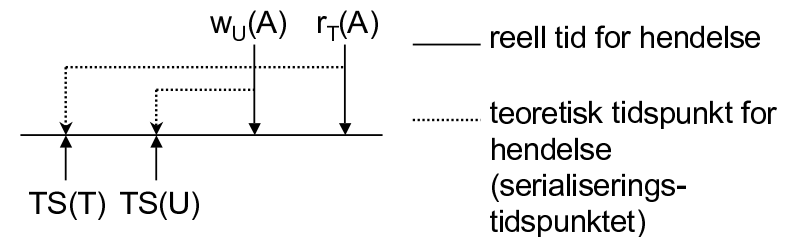
- Planleggeren (The Scheduler) må vedlikeholde en tabell over aktive transaksjoner og deres tidsstempler
- For å bruke tidsstempling til samtidighetskontroll må det knyttes to tidsstempler og en boolsk variabel til hvert eneste dataelement A i databasen:
 - $RT(A)$ – **lesetiden til A**: Det høyeste tidsstempelet til noen transaksjon som har lest A
 - $WT(A)$ – **skrivetiden til A**: Det høyeste tidsstempelet til noen transaksjon som har skrevet A
 - $C(A)$ – **commit-flagget til A**: Sant hvis, og bare hvis, den siste transaksjonen som skrev A, er committed

Tidsstempelserialisering

- Protokollen for serialisering ved hjelp av tidsstempler har som utgangspunkt at serialiseringsrekkefølgen er bestemt av tidsstemplene
- Protokollen sikrer at vi får en eksekveringsplan som er konfliktekvivalent med den serielle planen vi hadde fått om alle transaksjonene hadde gjort alle sine lese- og skriveoperasjoner ved tidsstempelet sitt (det vil si dersom transaksjonene var momentane og ikke brukte noe tid fra de startet til de var ferdige)
- Det er to problemer som kan oppstå:
 - Transaksjonen leser for sent
 - Transaksjonen skriver for sent

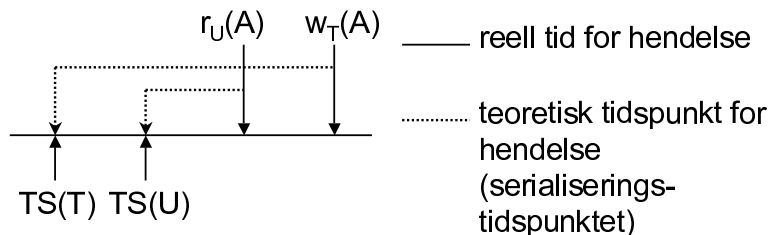
Transaksjonen leser for sent

- Transaksjon T vil lese dataelement A
- Verdien av A er skrevet av en transaksjon som startet etter T, dvs at $WT(A) > TS(T)$
- T ville komme til å lese «gal» verdi av A
- Konsekvens: T må abortere



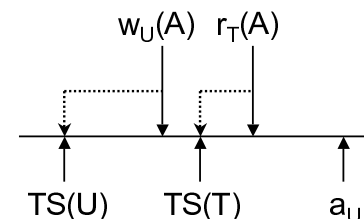
Transaksjonen skriver for sent

- T vil skrive dataelement A, men A er allerede lest av U som startet etter T, dvs at $RT(A) > TS(T)$
- Hvis $WT(A) > TS(T)$, skal ikke T skrive A (alt er OK)
- Hvis ikke, skulle U ha lest den verdien av A som T nå skal skrive, og T må abortere



Skitne data

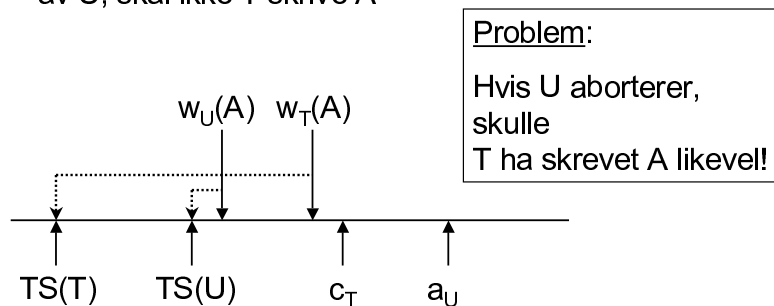
- Transaksjon T leser dataelement A som er skrevet av en annen transaksjon U
- Hvis U aborterer etter at T har lest A, har T lest en verdi av A som aldri skulle ha vært i DB
- Vi sier at T har lest en skitten verdi av A
- Dette er et brudd på isolasjonskravet, så T bør vente med å lese A til commit-flagget til A er sant



Merk!
Vi skriver
 a_T for at T aborterer og
 c_T for at T gjør commit

Thomas' skriverregel

- Transaksjon T ønsker å skrive dataelement A som allerede er skrevet av en annen transaksjon U med $TS(U) > TS(T)$
- Siden senere transaksjoner skal lese A-verdien skrevet av U, skal ikke T skrive A



INF3100 – 9.5.2006 – Ragnar Normann

13

Tidsstempelserialiseringsprotokollen–1

- 1 T ønsker å lese A
 - a) Hvis $TS(T) \geq WT(A)$, kan T få lov til å lese A:
 - I. Hvis $C(A)$, så utfør $r_T(A)$ og sett $RT(A) = \max(RT(A), TS(T))$
 - II. Hvis ikke $C(A)$, la T vente til $C(A)$ blir sann (evt. til transaksjonen som skrev A, aborterer)
 - b) Hvis $TS(T) < WT(A)$, er det fysisk umulig for T å lese (riktig verdi av) A
T må ruller tilbake (vi må utføre $rollback(T)$):
 - I. Utfør a_T
 - II. Start T på nytt med et høyere tidsstempel

INF3100 – 9.5.2006 – Ragnar Normann

14

Tidsstempelserialiseringsprotokollen–2

- 2 T ønsker å skrive A
 - a) Hvis $TS(T) \geq RT(A)$ og $TS(T) \geq WT(A)$, får T skrive:
 - I. Utfør $w_T(A)$
 - II. Sett $WT(A) = TS(T)$
 - III. Sett $C(A) = \text{falsk}$
 - b) Hvis $TS(T) \geq RT(A)$ og $TS(T) < WT(A)$, er A skrevet av en nyere transaksjon enn T
Hvis $C(A)$, ignorer ønsket (Thomas' skriverregel)
Hvis ikke $C(A)$, la T vente til $C(A)$ blir sann (evt. til transaksjonen som sist skrev A, aborterer)
 - c) Hvis $TS(T) < RT(A)$, er A lest av en nyere transaksjon enn T, og T må ruller tilbake

INF3100 – 9.5.2006 – Ragnar Normann

15

Tidsstempelserialiseringsprotokollen–3

- 3 T ønsker å committe
 - a) Utfør c_T (dvs. skriv $commit(T)$ i loggen)
 - b) For hver A som T har skrevet, sett $C(A) = \text{sann}$
 - c) La transaksjoner som venter på å lese eller skrive en slik A, få fortsette
- 4 T blir abortert (eller ønsker selv å abortere)
 - a) Bruk loggen til å gjøre Undo på alle skrivinger foretatt av T og skriv $abort(T)$ i loggen
 - b) La alle transaksjoner som venter på å lese eller skrive et dataelement skrevet av T, forsøke på nytt

INF3100 – 9.5.2006 – Ragnar Normann

16

Versjonering

- Metode for å unngå abort pga at T vil lese en A skrevet av en nyere transaksjon
- Når en skriveoperasjon $w_T(A)$ utføres, blir det laget en ny versjon A_t av A, der $t = TS(T)$
- Når en leseoperasjon $r_T(A)$ utføres, leses den versjonen A_t som har den største $t \leq TS(T)$
- Foreldede versjoner av A kan, og bør, slettes
- For at A_t skal kunne slettes, må det finnes en versjon $A_{t'}$ der $t < t'$ og der $t' \leq TS(T)$ for alle aktive transaksjoner T
- En effektiv implementasjon av versjonering krever i praksis at dataelementene er blokker

Tidsstemping vs låsing

- Tidsstemping er best hvis de fleste transaksjonene bare leser, eller hvis konflikter er sjeldne
 - Hvis konflikter er vanlige, vil tidsstemping føre til mange tilbakerullinger, og låsing er bedre
 - Noen kommersielle systemer tilbyr et kompromiss:
 - Lese/skrive-transaksjoner bruker 2PL
 - Rene lesetransaksjoner bruker multiversjon tidsstemping
- Dette gjør at lesetransaksjoner aldri må abortere og sjelden må vente

Validering

- Validering er en optimistisk serialiseringsstrategi som baserer seg på tidsstemping
- Den skiller seg fra vanlig tidsstemping ved at man ikke lagrer lese- og skrive-tidsstempel for alle dataelementene i databasen
- For hver aktiv transaksjon T lagres to mengder
 - **lesemengden** til T, $RS(T)$
 - **skrivemengden** til T, $WS(T)$som inneholder alle dataelementer som T henholdsvis leser og skriver

Validering (forts.)

Utførelsen av en transaksjon T deles inn i tre faser:

1. **Lesefasen**

All lesing og beregning gjøres her
 $RS(T)$ og $WS(T)$ bygges opp i Ts adresserom
Starttidspunktet for lesefasen kalles $Start(T)$

2. **Valideringsfasen**

T valideres ved å sammenligne $RS(T)$ og $WS(T)$ med lese- og skrivemengdene til andre transaksjoner (detaljene kommer snart)
Hvis valideringen feiler, rulles T tilbake
Sluttidspunktet for valideringsfasen kalles $Val(T)$

3. **Skrivefasen**

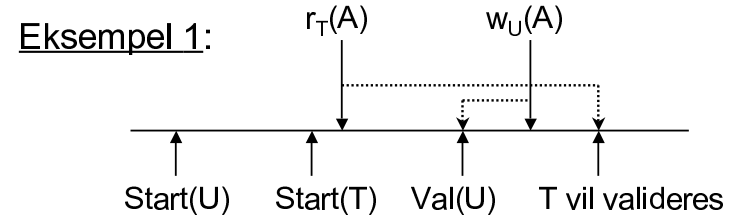
Her skriver T verdiene i $WS(T)$ til databasen
Sluttidspunktet for Skrivefasen kalles $Fin(T)$

Verdien av $Val(T)$ bestemmer serialiseringsrekkefølgen

Validering (forts.)

- Planleggeren vedlikeholder tre mengder med transaksjoner:
- START**
De som har startet, men ennå ikke har avsluttet valideringsfasen
For hver T i START lagres Start(T)
 - VAL**
De som er validert, men ikke har avsluttet skrivefasen
For hver T i VAL lagres Start(T) og Val(T)
 - FIN**
De som (nylig) har avsluttet skrivefasen
For hver T i VAL lagres Start(T), Val(T) og Fin(T)
T kan fjernes fra FIN når vi for alle $U \in \text{START} \cup \text{VAL}$ har at $\text{Start}(U) > \text{Fin}(T)$

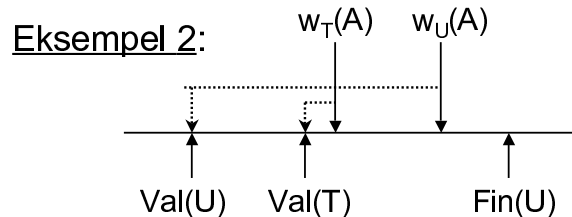
Valideringsfasen



Anta at idet T vil valideres, så finnes en U slik at

- $U \in \text{VAL} \cup \text{FIN}$ (dvs at U er validert)
 - $\text{Fin}(U) > \text{Start}(T)$ (dvs at U ikke var avsluttet da T startet)
 - $\text{RS}(T) \cap \text{WS}(U) \neq \emptyset$ (på figuren er $A \in \text{RS}(T) \cap \text{WS}(U)$)
- T må ruller tilbake fordi T kan ha lest en «gal» verdi av A

Valideringsfasen (forts.)



Anta at i det T vil valideres, så finnes en U slik at

- $U \in \text{VAL}$ (dvs at U er validert)
- $\text{Fin}(U) > \text{Val}(T)$ (dvs at U ikke er avsluttet ennå)
- $\text{WS}(T) \cap \text{WS}(U) \neq \emptyset$ (på figuren er $A \in \text{WS}(T) \cap \text{WS}(U)$)

T må ruller tilbake da T kan komme til å skrive A før U gjør det

Valideringstesten

- De to foregående eksemplene dekker alle mulige feilsituasjoner
- Dermed består valideringsfasen av følgende to tester:
 - Sjekk at $\text{RS}(T) \cap \text{WS}(U) = \emptyset$ for alle validerte U med $\text{Fin}(U) > \text{Start}(T)$
 - Sjekk at $\text{WS}(T) \cap \text{WS}(U) = \emptyset$ for alle validerte U som ennå ikke er avsluttet
- Hvis T består begge testene, er T validert og går inn i skrivefasen
- Hvis ikke, må T ruller tilbake og starte på nytt

Kaskadetilbakerulling

T ₁	T ₂	A	B
I ₁ (A); r ₁ (A); A←A+100; w ₁ (A); I ₁ (B); u ₁ (A);		25	25
	I ₂ (A); r ₂ (A); A←Ax2; w ₂ (A); I ₂ (B); Avslått	125	
r ₁ (B); a₁ ; u ₁ (B);		250	
	I ₂ (B); u ₂ (A); r ₂ (B); B←Bx2; w ₂ (B); u ₂ (B);		50
		250	50

Når T₁ aborterer, sletter planleggeren alle låsene T₁ har
Hvis T₂ får fortsette, vil T₂ lage en inkonsistent tilstand
Altså må T₂ ruller tilbake (fordi T₂ har lest en skitten A)

Kaskadetilbakerulling (forts.)

- Kaskadetilbakerulling kan være rekursiv:
Abort av T₁ kan føre til abort av T₂ som kan føre til abort av T₃ osv
- Kaskadetilbakerulling kan omfatte committede transaksjoner, noe som er i strid med D i ACID
- Kaskadetilbakerulling bør derfor unngås
- Tidsstempelprotokoller med commit-flagg sikrer mot kaskadetilbakerulling
- Verifisering sikrer også mot kaskadetilbakerulling (ingen skrivning gjøres før man vet at det ikke blir noen abort)

Gjenopprettbare eksekveringsplaner

- En eksekveringsplan er gjenopprettbar hvis ingen transaksjon T gjør commit før alle transaksjoner som har skrevet data som T har lest, har gjort commit
- Eksempel 1: En gjenopprettbar serialiserbar (seriell) plan:
S₁: w₁(A) w₁(B) w₂(A) r₂(B) c₁ c₂
- Eksempel 2: En gjenopprettbar ikke-serialiserbar plan:
S₂: w₂(A) w₁(B) w₁(A) r₂(B) c₁ c₂
- Eksempel 3: En serialiserbar, ikke gjenopprettbar, plan:
S₃: w₁(A) w₁(B) w₂(A) r₂(B) c₂ c₁
- For at en eksekveringsplan skal være gjenopprettbar både for undo-, redo- og undo/redo-logging, må loggens commit-poster skrives til disk i samme rekkefølge som de skrives i loggen (flere loggposter kan ligge i samme blokk og bli skrevet samtidig)

ACR-planer

- En eksekveringsplan unngår kaskadetilbakerullinger hvis transaksjonene bare kan lese data skrevet av committede transaksjoner
- Slike planer kalles ACR-planer (ACR = Avoid Cascade Rollback)
- Alle ACR-planer er gjenopprettbare

Bevis:

Anta at T₂ leser en verdi skrevet av T₁ etter at T₁ har gjort commit

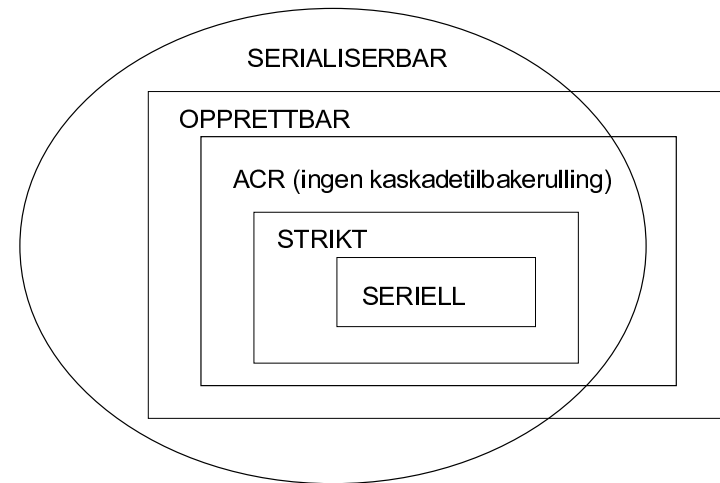
Siden T₂ hverken har gjort commit eller abort før den leser, må T₂ gjøre sin commit eller abort etter at T₁ gjorde sin commit

qed

Strikte eksekveringsplaner

- En eksekveringsplan kalles **strikt** hvis den er basert på låser og følger følgende regel:
- **Strikt låseregel:** En transaksjon kan ikke friggi noen skrive-lås før den har gjort commit eller abort, og commit- eller abort-loggposten er skrevet til disk
- En strikt eksekveringsplan er en ACR-plan
Bevis: Fordi skrive-låsene ikke blir frigitt før etter at transaksjonene har gjort commit, kan ingen lese data skrevet av en transaksjon som ikke har gjort commit
- En strikt eksekveringsplan er serialiserbar
Bevis: En strikt eksekveringsplan er åpenbart konflikt-ekvivalent med den serielle planen vi får ved å la hver transaksjon bli utført ved committidspunktet

Eksekveringsplantyper



Tilbakerulling ved bruk av låser

- Hvis dataelementene er blokker, er alt enkelt:
All skriving gjøres i bufferet;
intet skrives til disk før commit
Ved abort frigies blokken som blir ubrukt bufferareale
Samme teknikk virker ved versjonering;
blokken med «abortert versjon» frigies
- Hvis det er flere dataelementer i hver blokk, er det tre måter å restaurere data på etter en abort
 1. Originalen kan leses fra databasen på disk
 2. Med en undo- eller undo/redo-logg kan originalen hentes fra loggen
 3. Hver aktiv transaksjon kan ha sin egen logg over sine endringer i hukommelsen

Logisk logging

- Dette er en loggtype som benytter seg av transaksjons-logikken til å foreta tilbakerullinger
- Typiske logiske loggposter består av fire felt
 - L – et loggpostløpenummer
 - T – transaksjonen
 - A – aksjon (operasjon) utført (f.eks. insert tuppel t)
 - B – blokken hvor A ble utført
- For hver aksjon finnes en kompensierende aksjon som opphever virkningen (f.eks. delete for insert)
- Hvis T aborterer, blir alle Ts aksjoner kompensert, og kompensasjonene blir loggført
- Hver blokk har loggpostnummeret på siste aksjon som berørte den