

CS3700 - Introduction to Database Systems  
End Semester Examination - 21st - 22nd Dec

I understand that take-home examination is to test my own understanding of the subject and I will offer an honest attempt to meet this goal. I will not communicate with anyone else while working, answering the examination problems / tasks. I will not post the questions to any public fora and seek answers. I will not help other course mates taking the examination in any manner. I understand that I can consult books, course slides and videos by the teacher while answering the examination. I also understand that I should not copy materials from books, internet sites etc as a part of my answers for the examination questions / tasks. These restrictions apply for whole examination duration (till the deadline for submission).



ARABHI SUBHASH

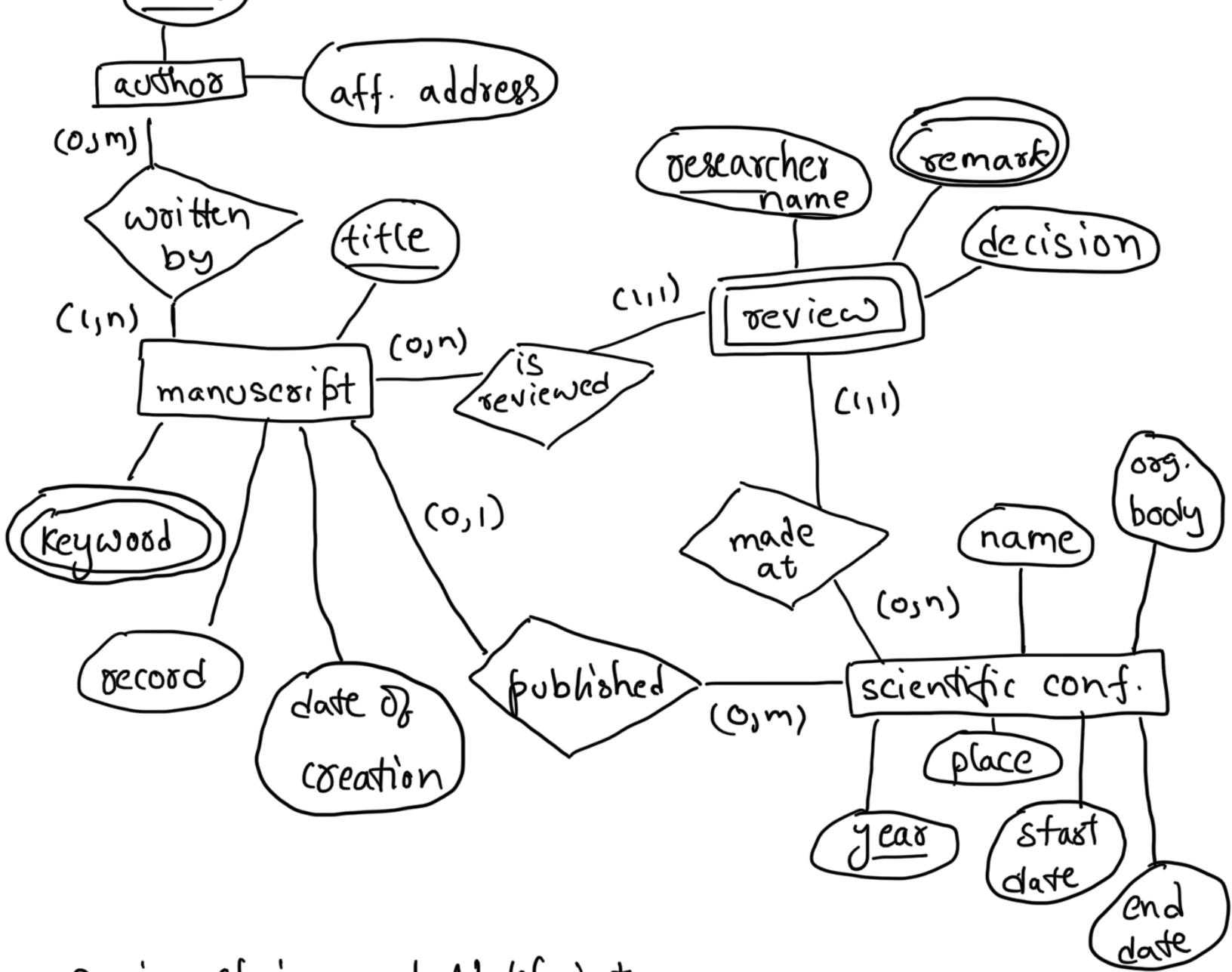
CS17B005

## Task #1

$5 \text{ mod } 3 = 2$ ;  $2+1=3$ ; working on 3rd domain description.

a)

ER-Diagram



### Design choices and Notepoints

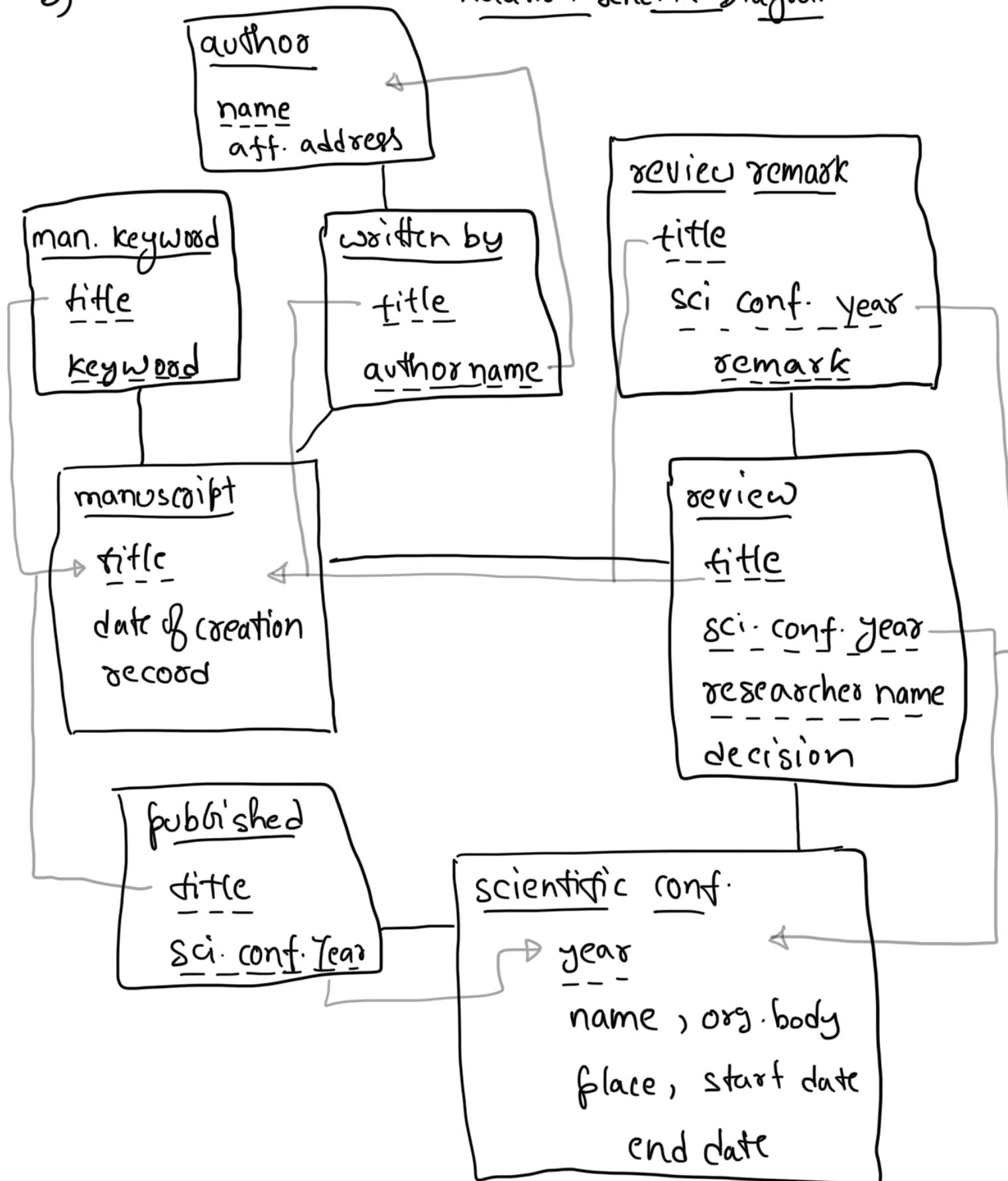
- 1) Didn't maintain separate relation for "Submitted" because that info. can be obtained by checking whether there is a review for that manuscript at that conference (Implicit)
- 2) Made **review** a weak entity whose existence depends on sci. conf. and manuscript. The reason for this is to accomodate multiple reviews

for a pair of conf. & manuscript

- 3) Assumed title is unique for manuscripts or else we can also add an attribute named man. id. and make it a primary key
- 4) Maintained a relation "published" with min-max constraint (0,1) to impose restriction that a manuscript can be published at exactly one conference or not published yet
- 5) keyword, remark are multivalued attributes
- 6) The contents of book "proceedings" of a certain conference is implicitly stored in the relation "published"

b)

### Relation Schema Diagram



" $\rightarrow$ " - foreign key

" $--$ " - primary key

### Task #2

a)

Query: Get the courseID and name of courses which are all male students in odd sem of 2002. Assume every course is offered and has atleast 1 student

enrolled in the given semester and year

Query TRC : { c.courseId, c.name | course(c)  $\wedge$   
 $(\forall s)( (s \text{ student}(s) \wedge \text{enrollment}(e)$   
 $\wedge e.\text{year} = 2002 \wedge e.\text{courseId} = c.\text{courseId}$   
 $\wedge e.\text{sem} = \text{'odd'} \wedge s.\text{rollNo} = e.\text{rollNo}) \rightarrow$   
 $(s.\text{sex} = \text{'male'}) ) \}$

Explanation: Basically we are picking all the students enrolled in the course and checking the gender. If all are male then  $\forall$  condition satisfies and output that tuple

Query Relational Algebra:

- 1) courses-with-female  $\leftarrow \Pi_{\text{courseId}} (\sigma_{\text{sem} = \text{'female'}} (\Pi_{\text{rollNo}, \text{courseId}} (\sigma_{\text{year} = 2002 \wedge \text{sem} = \text{'odd'}} (\text{enrollment})) \bowtie \Pi_{\text{rollNo}, \text{sem}} (\text{student})) ) ;$
- 2) courses-only-male  $\leftarrow (\Pi_{\text{courseId}} (\text{course}) - (\text{courses-with-female})) ;$
- 3) result  $\leftarrow \Pi_{\text{courseId}, \text{name}} (\text{course} \bowtie \text{courses-only-male}) ;$

Explanation: In part 1 collected all the courses in given sem and year with female students. Did that using inner join of student & enrollment. In second part deleted them from all courses to get male only courses. In last part joined it with courses to get course name

into output.

b) query : Get the name and rollNo. of all the students who did a 4 credit course in even semester of 2002

sql query : select s.name, s.rollNo from students, enrollment e, course c where  
e.rollNo = s.rollNo and e.courseId = c.courseId  
and e.sem = 'even' and e.year = 2002  
and c.credits = 4

Relational Algebra expression :

$$\Pi_{s.name, s.rollNo} (\sigma_{\alpha}(s \times e \times c)) \text{ where } s: \text{student}$$

$\alpha = (e.rollNo = s.rollNo \text{ and } e.courseId = c.courseId \text{ and } e.sem = 'even' \text{ and } e.year = 2002 \text{ and } c.credits = 4)$

optimisations using main heuristic rule :

1)  $\Pi_{s.name, s.rollNo} (\sigma_{B_1} (s \times \sigma_{B_2}(e) \times \sigma_{B_3}(c)))$

$B_3 = (c.credits = 4)$

$B_2 = (e.sem = 'even' \text{ and } e.year = 2002)$

$B_1 = (e.rollNo = s.rollNo \text{ and } e.courseId = c.courseId)$

→ pushed independent constraints inside of brackets

2)  $\Pi_{s.name, s.rollNo} (\sigma_{\Theta_1} (s \times \sigma_{\Theta_2} (\sigma_{\Theta_3}(e) \times \sigma_{\Theta_4}(c))))$

$\Theta_4 = (c.credits = 4)$

$\Theta_3 = (e.sem = 'even' \text{ and } e.year = 2002)$

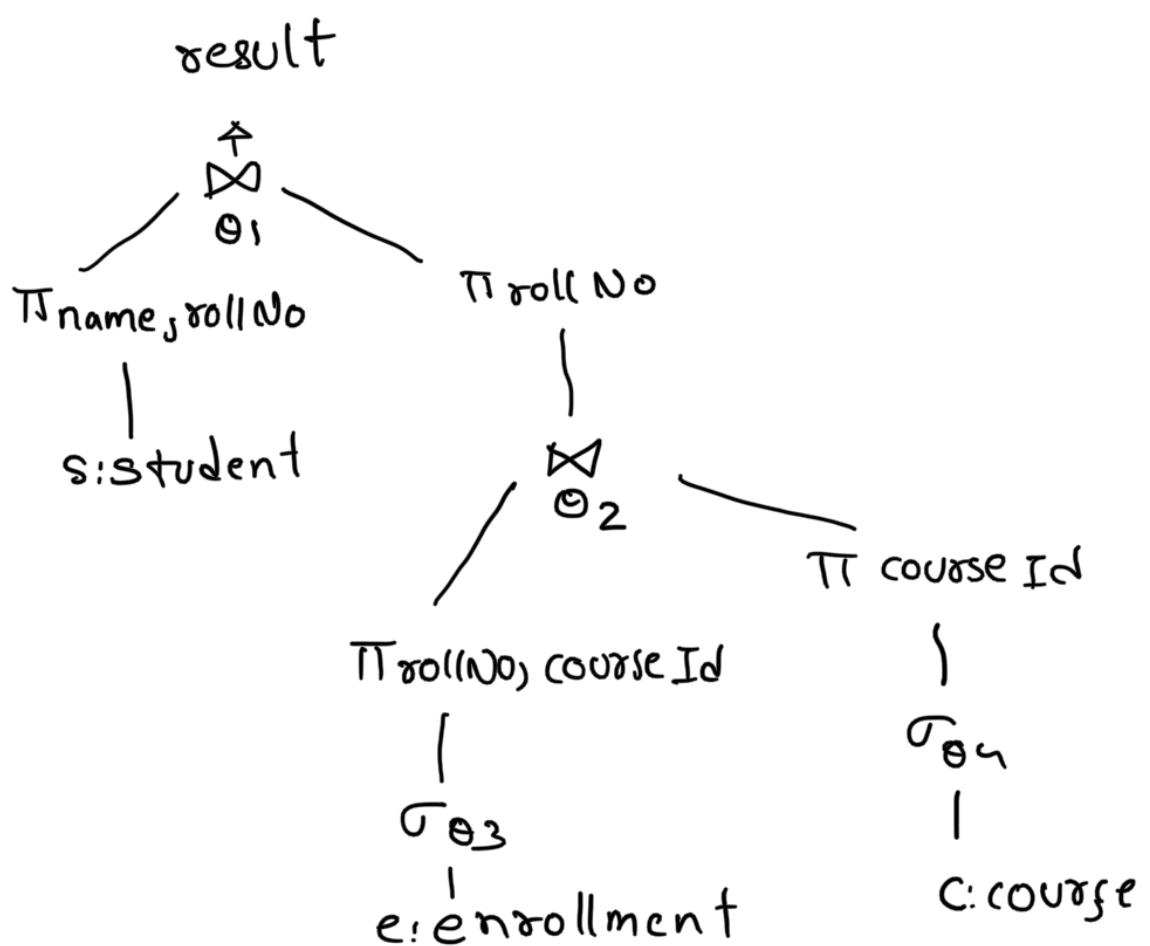
$\Theta_2 = (e.courseId = c.courseId)$

$\Theta_1 = (e.rollNo = s.rollNo)$

→ further pushing of constraints

- 3)  $\Pi_{s.name, s.rollNo} (s \bowtie_{\theta_1} (\sigma_{\theta_3(e)} \bowtie_{\theta_2} \sigma_{\theta_4(c)})$   
 → converted cross product + selection into join.
- 4)  $(\Pi_{name, rollNo(s)} \bowtie_{\theta_1} \Pi_{rollNo} (\Pi_{rollNo, courseId} (\sigma_{\theta_3(e)} \bowtie_{\theta_2} \Pi_{courseId} (\sigma_{\theta_4(c)})))$   
 → Finally, Applied projection to lower levels

Above RA expressions with brackets is equivalent to tree representation. Final expression tree is



c) query: Get max and min of avg credits offered by each department. Assume each dept. offers atleast one course.

SQL query :  
`create view avg-credits-dept as  
 select c.deptNo, avg(c.credits) as avg-credits  
 from course c group by c.deptNo;  
 select max(avg-credits), min(avg-credits) from avg-credits  
 -dept;`

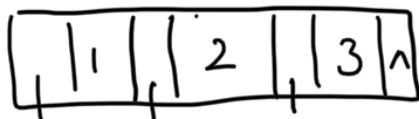
Explanation: In Given query we need to max & min of avg (aggregate over aggregate) So we need views.

In the first part using avg we have created a view avg-credits-dept which stores deptNo and its avg. credits. Finally we applied min\_max on that view.

### Task #3

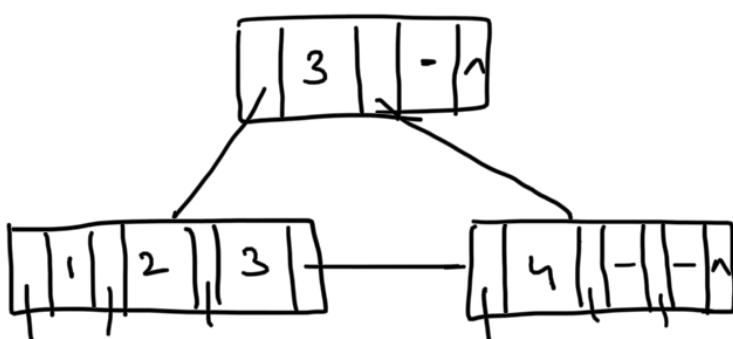
a)  $o = 3$     $o_{leaf} = 3$

i) Insert 1,2,3

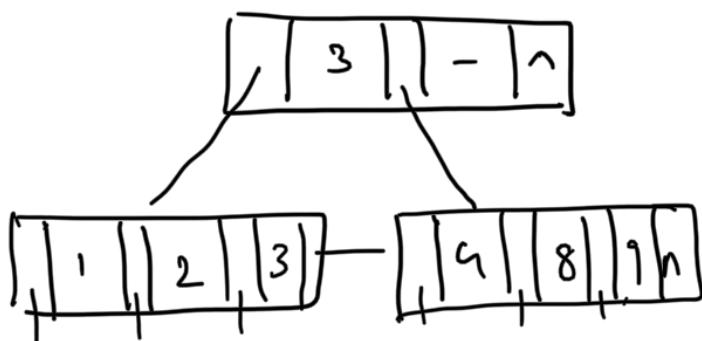


ii) Insert 4

overflow occurs at  $\lceil \frac{3+1}{2} \rceil = 2$ , leafnode

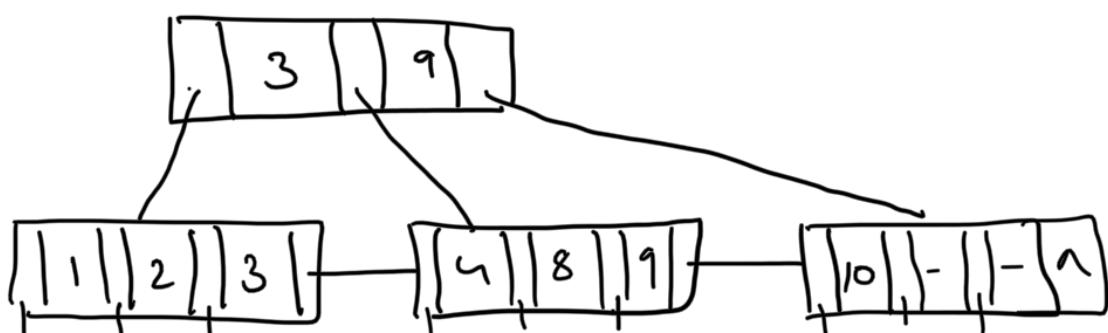


iii) Insert 8,9

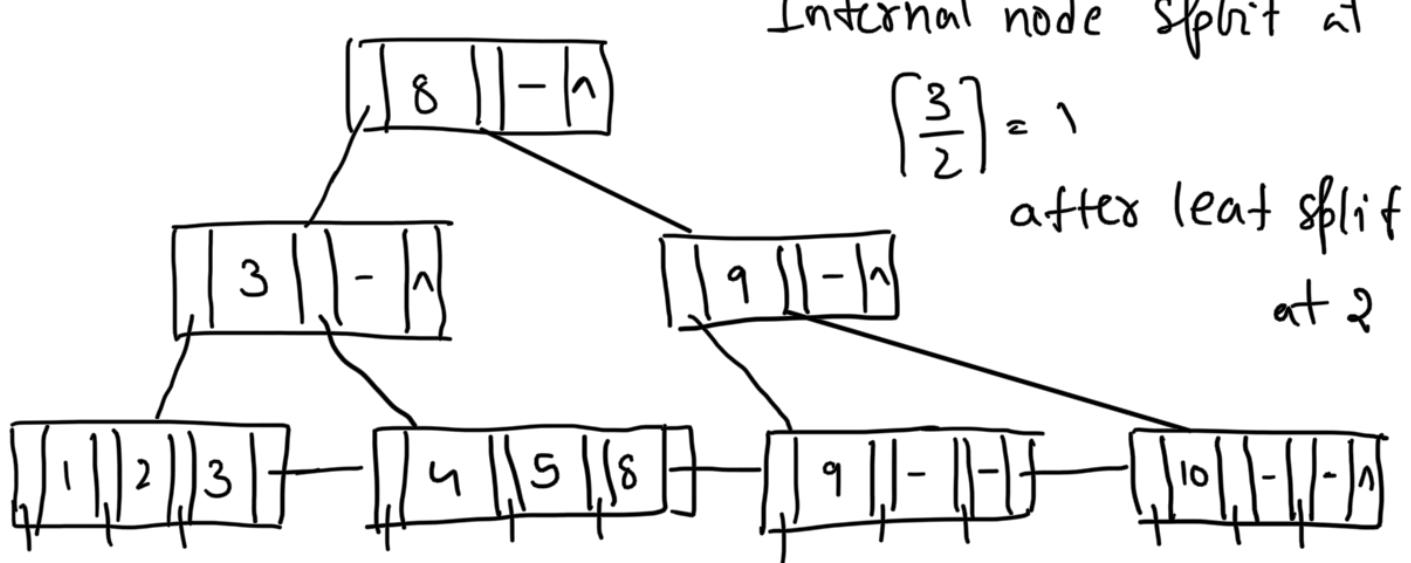


iv) Insert 10

same, overflow at 2



v) Insert 5



After 8 inserts we saw 2 leaf splits and 1 internal node + leaf split. The final tree above still have spaces but we stopped inserting as we have seen the 2 different splits and reached 2 levels.

b) Pseudo code for Qinear hashing insert:-

```

M = 1 # initial length of bucket
N = 0 # split pointer
i = 0 # present set of hash functions

def hash_func(i,x):
    return x%(pow(2,i)*M)

def split_bucket():
    # add bucket (pow(2,i)*M)+N to the set of buckets
    for values: v in bucket N and it's over_flow:
        if hash_func(i+1,v) != N:
            # remove v from bucket N
            # add v to bucket (pow(2,i)*M)+N (Newly added Bucket)

def place_record(b,x):
    # check if bucket b is full
    if not full:
        # add x to bucket b
    else:
        # add x to over flow of bucket b
        split_bucket()
        if N == (pow(2,i)*M)-1:
            i+=1
            N = 0
        else:
            N+=1

def insert(x):
    if hash_func(i,x) < N:
        place_record(hash_func(i+1,x),x)
    else:
        place_record(hash_func(i,x),x)

```

Brief: Above is pseudo code for linear hashing where hashing functions are  $(x \bmod 2^i \cdot M)$ . Simple adding to bucket and adding to overflow are presented as comments

### Task #5

a) Given statement has if and only if (double implication) so proof will be in 2 parts

i) forward implication

→ if a schedule is conflict serializable then precedence graph  $\sigma$  has no cycles

proof by contradiction:

→ consider a schedule  $x$  which is conflict serializable and has a cycle in  $G_x \equiv A \wedge \sim B \equiv \sim(\sim A \vee B) \equiv \sim(A \rightarrow B)$

→ let the cycle be  $T_1, T_2, T_3 \dots T_k$

→ This implies that there must be operations  $OP_1(x) OP_2(x) \dots OP_k(n) OP_1(x)$  in the schedule

→ To make it conflict serializable we must bring the first and last  $OP_1$  operations together - There are 2 ways of doing this

→ one way is to move the first operation forward but this is not possible due to conflict with  $OP_2, OP_3 \dots OP_k$  of lower precedence

- other way is to move last op. back but that too conflicts due to higher precedence of same set of operations
- Hence we can say that our assumption is false and by contradiction the required forward implication is true.

## 2) Backward implication

- If precedence graph is acyclic then schedule is conflict serializable.

proof by induction on number of Transactions:

Basis step:  $n=1$

→ It is trivially true for one transaction

Hypothesis step:  $n=k$

- Assume that for  $k$  transactions the reqd. implication is true

Induction step:  $n=k+1$

→ Given precedence graph is acyclic then there will be at least one node with no in-edge. Let this node be  $T_m$

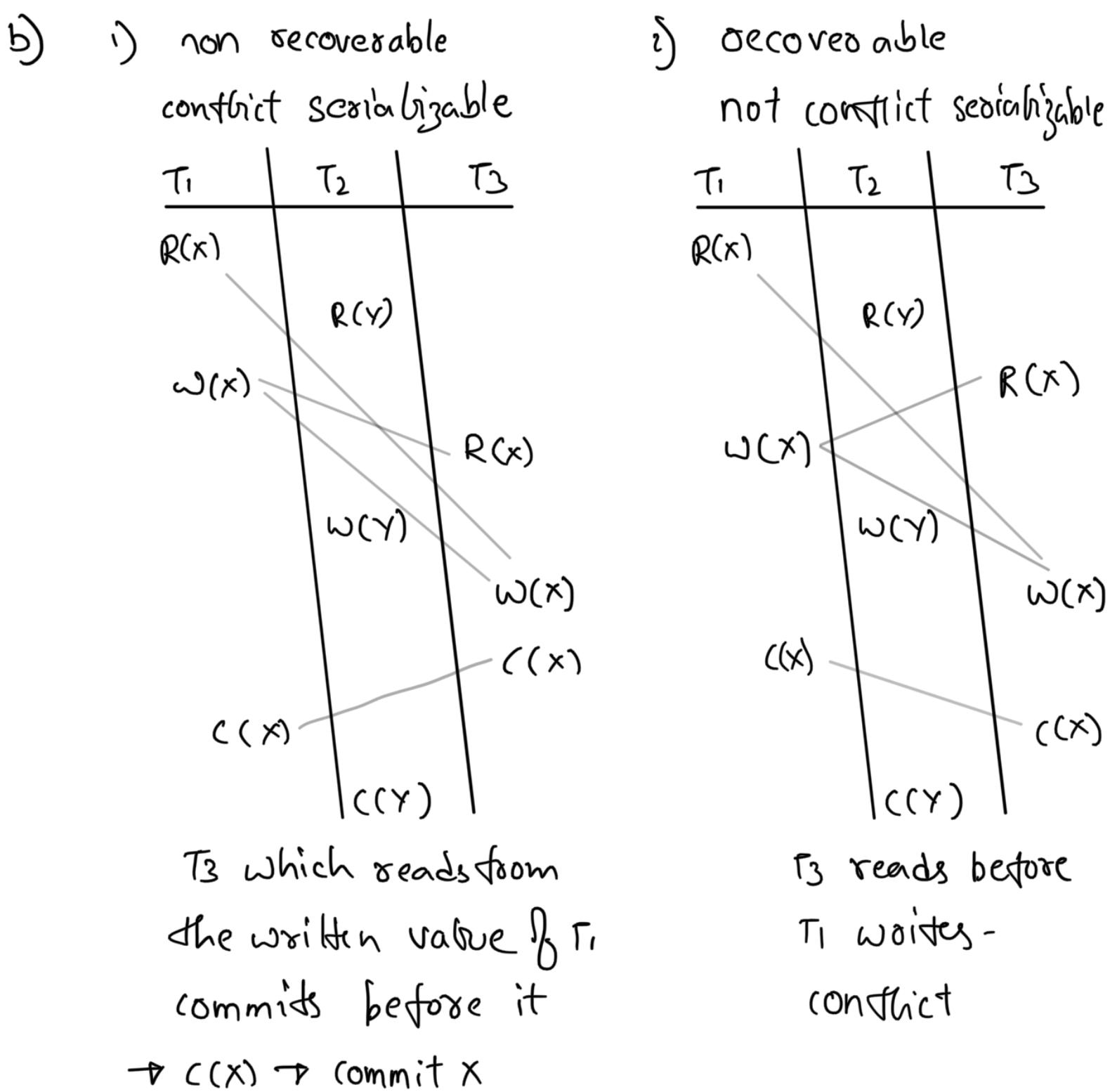
- As there is nothing preceding transaction  $m$ , all the operations corresponding to this transaction can be swapped to start without any conflict

#  $OP_m(x) \dots OP_1(x) \dots OP_m(x) \dots OP_L(x) \dots OP_m(x)$

transformed to

$OP_m(x) \quad OP_m(x) \quad OP_m(x) \dots OP_1(x) \quad OP_L(x) \dots$

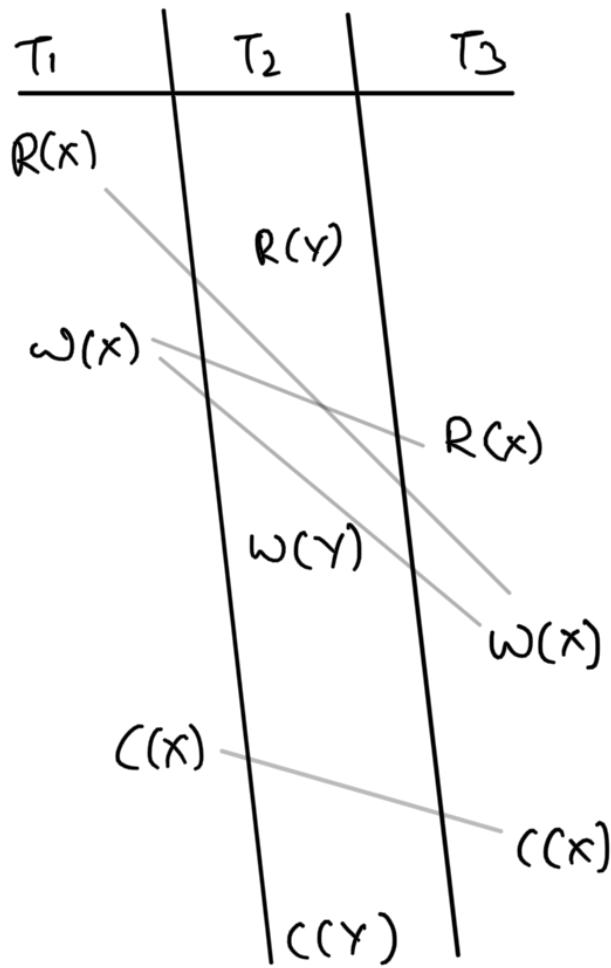
- So the original schedule  $s$  can be written as  $s = T_m ; s'$ . here  $s'$  is the schedule of remaining  $K$  transactions
- From hypothesis step we can say that  $s'$  is conflict serializable. By merging this with transaction  $m$  we can say that our assumption is true for  $n = k + 1$
- From induction we can conclude that our backward implication is true for any number of transactions



3) recoverable

conflict serializable

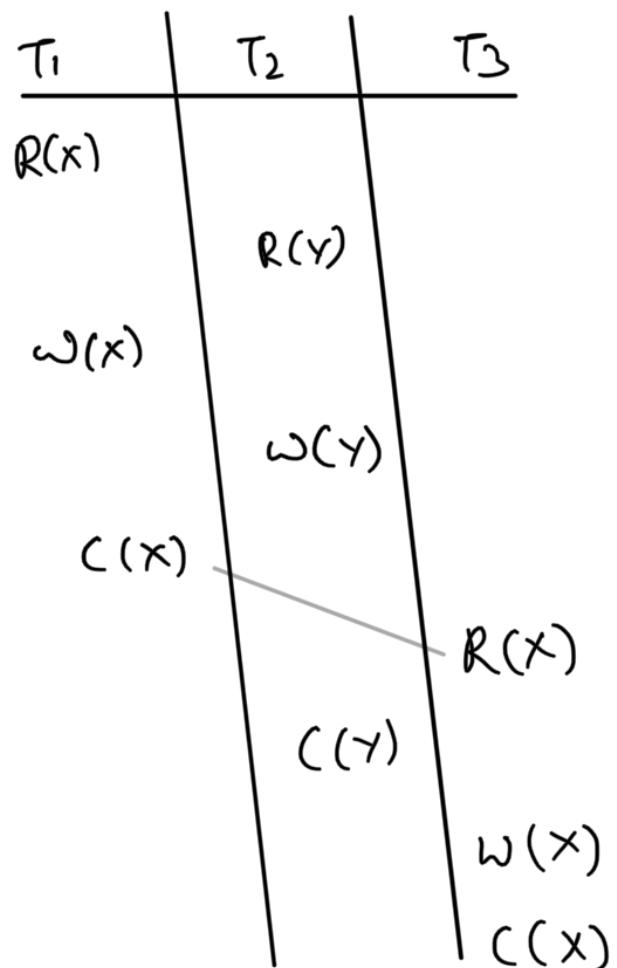
not cascading



4) recoverable

conflict serializable

cascading



T<sub>3</sub> reads from T<sub>1</sub>  
after it writes  
and before it  
commits - cascading  
- no conflicts

T<sub>3</sub> reads after  
T<sub>1</sub> commits its  
written value  
- no cascading  
- no conflicts too

C) RDBMS logging stores the final values of the transaction. As transactions proceed in forward direction so should the reading of log be. Below is a small example to show it

Consider a state where A has 500/- and B has 1000/-. From this state 2 transactions going to happen - i) A spends 100/-

i) A spends 200

### State Table

Step	Action	m	Mem-A	Disk-A	Log
1)					$\langle \text{start } T_1 \rangle$
2)	$R_1(A, m)$	500	500	500	
3)	$m = m - 100$	400	500	500	
4)	$w_1(A, m)$	400	400	500	$\langle T_1, A, 400 \rangle$
5)					$\langle \text{commit}, T_1 \rangle$
6)	Flush Log				
7)	Output(A)	400	400	400	
8)					$\langle \text{start}, T_2 \rangle$
9)	$R_2(A, m)$	400	400	400	
10)	$m = m - 200$	200	400	400	
11)	$w_2(A, m)$	200	200	400	$\langle T_2, A, 200 \rangle$
12)					$\langle \text{commit}, T_2 \rangle$
13)	Flush Log				
14)	Output(A)	200	200	200	

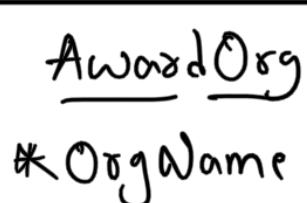
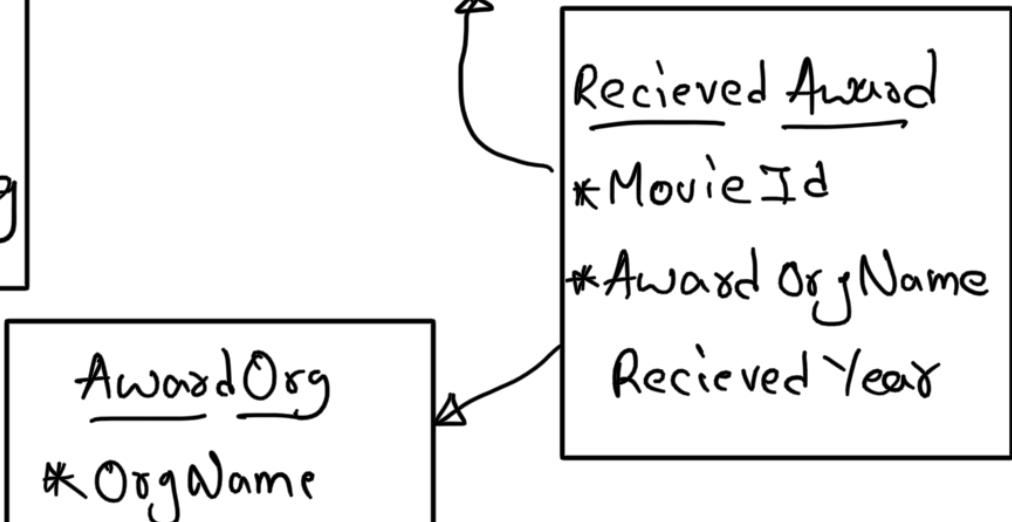
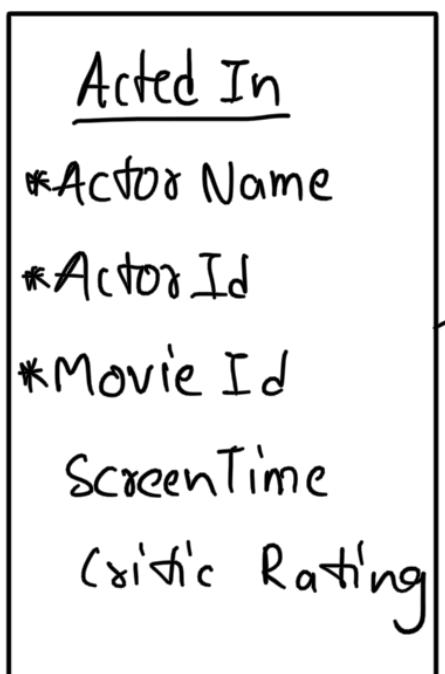
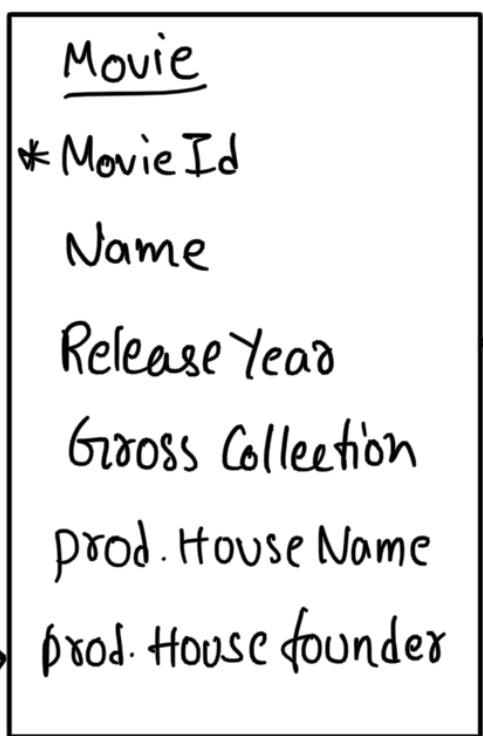
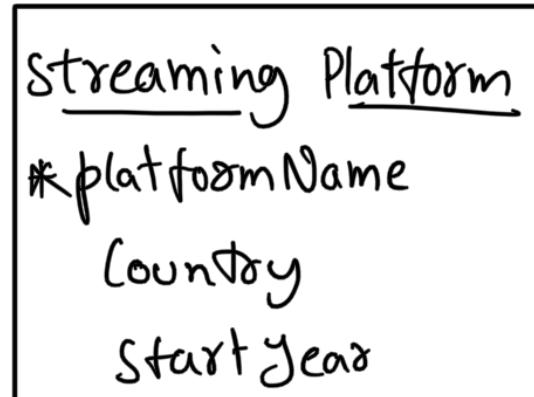
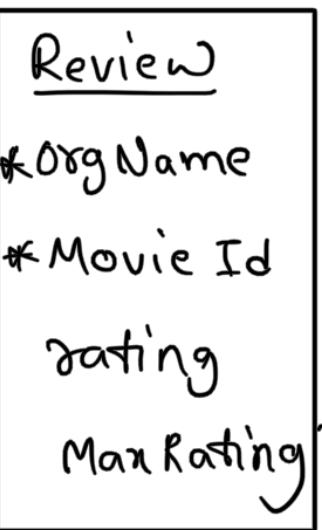
- If crash occurs before 6 then commit isn't present in the disk and the values are also not updated so no action required.
- If it occurs after 6 and before 13 Transaction  $T_1$  is redone using log entries and no action required for  $T_2$ .
- Consider crash occurs after 13 and we don't know whether result of  $T_1$  or  $T_2$  flushed to disk. Here we need redo transactions from start as we have both

commits in disk. Forward reading of log puts a value of 200 to the disk - A where as backward or other readings will write 400.

We can see from the table that forward (forward) reading is only correct as transactions happen in this direction.

### Task #4

- a) Modified DB scheme of database designed in the assignment - i



Country	
Start Year	

\* → key

- Relations are modified so different Normal forms can be observed
- Name of attribute intuitively provides its meaning. Few of them are explicitly defined below

MovieId :- Unique Id given to a movie

orgName, PlatformName :- unique name given to organization or platform respectively

StartYear :- Established Year of the company

streams: Region :- The region / part of world where the movie can be streamed by the platform.

ActorId :- unique Id given to each actor in the movie industry

ActorName :- Name of Actor (Assume it's unique)

MaxRating :- The scale on which rating is given starting from 0 by Revieworg.  
(Ex: 10 → (0-10), 100 → (0-100))

Prod. House :- production House

Screen Time :- The total time actor is present on screen during the movie

Critic Rating :- Avg performance rating given by critics

- Functional Dependencies in each relations and

Their highest normal form are described below

i) Review Key - (OrgName, MovieId)

FD - (OrgName, MovieId)  $\rightarrow$  Rating

\* Each org. gives a single rating to movie (overall rating)

- OrgName  $\rightarrow$  MaxRating

\* Each org. has specific scale on which they rate movies

NF - Not 2NF as non prime attribute (MaxRating) is not fully functionally dependent on key (Orgname, MovieId)

So it's just 1NF

ii) Movie Key - MovieId

FD - MovieId  $\rightarrow$  Name, ReleaseYear, Gross

prod. House Name

prod. House founder

\* MovieId is uniquely given to a movie

- prod. House Name  $\rightarrow$  prod. House founder

\* Assuming naming is unique and a house has only one founder

NF - clearly it's 2NF as every attribute depends on single key (MovieId) but it's not 3NF as there is a transitive relation -

MovieId  $\rightarrow$  prod. House Name  $\rightarrow$  prod. House founder

3) Acted In

Keys - (ActorId, MovieId),

(Actor Name, MovieId)

- FD -  $(ActorId, MovieId) \rightarrow ScreenTime, CriticRating$
- $(ActorName, MovieId) \rightarrow ScreenTime, CriticRating$
  - \* As actor can't have two or more screen times or avg critic ratings for a movie.
  - $ActorId \rightarrow ActorName$
  - $ActorName \rightarrow ActorId$
  - \* Assuming name and Id are unique for an actor in movie industry

NF - clearly its 3NF as its 2NF and no non prime transitive property but It's not BCNF as  $ActorId \rightarrow ActorName$  is nontrivial

FD and  $ActorId$  is not a superkey.

- 5) Streaming platform key - platformName
- FD - platformName  $\rightarrow$  Country, StartYear
- \* trivial, Name is unique for a platform
- NF - BCNF, In every nontrivial FD  $X \rightarrow A$  we can see X is superkey

- 6) Streams key - (platformName, MovieId)
- FD - platformName, MovieId  $\rightarrow$  Region
- \* A company buys streaming rights of a movie for a region
- NF - BCNF

- 7) Award Org key - OrgName
- FD - OrgName  $\rightarrow$  Country, StartYear
- \* trivial, Naming is unique
- NF - BCNF

→ Received Award      key - (MovieId, AwardOrgName)

FD - (MovieId, Award OrgName) → Received Year

\* Trivial, A film is considered for awards by an org. only one time (one year).

NF - BCNF

b) Given : If the matrix operated on by L-test doesn't have a row of all "a" symbols at the end of its running then the given decomposition is lossy wrt given FDs

Proof:

Consider a relation Schema  $R = A_1, A_2 \dots A_n$  ( $A_i$  - attribute) and the decomposition

$S - R \rightarrow R_1, R_2, R_3 \dots R_k$

L-Test : A table of  $k \times n$  is constructed with  $R$  as rows and  $A$  as columns. This table is instantiated with  $b_{ij}$ .

→ Now if  $R_i \rightarrow A_j \dots A_k$  then the values  $b_{ij} \dots b_{ik}$  are replaced with  $a_j \dots a_k$

↔ This step is equivalent to adding an instance  $a_1, a_2 \dots a_n$  to  $R$  (adding a data point  $(a_1, a_2 \dots a_n)$  to initial Relation Schema)

→ After that we use FDs to

change the values of matrix

⇒ Let us consider an FD  $A_j \rightarrow A_k$

and  $R_i$  has elements  $a_j$  and  $b_{ik}$ .

As functional dependency independent  
of decomposition from R we can say  
if  $A_j$  is  $a_j$  then  $A_k$  is  $a_k$ . Therefore  
we can change the value of  $b_{ik}$  to  $a_k$   
 $\rightarrow$  even after robustly applying FDs if  
we couldn't get a row with all "a"s  
means that our decomposition doesn't  
have the instance  $a_1a_2\dots a_n$  wrt given  
FDs

$\rightarrow$  As the instance  $a_1a_2\dots a_n$  present in R  
(lost after decomposition, this is a  
lossy decomposition wrt given FDs)