Data Management – exam of 13/07/2022

Problem 1

In a schedule S on transactions $\{T_1, \ldots, T_n\}$ we say that two transactions T_i, T_j share the element X of the database if there exist actions $\alpha(X)$ in T_i and $\beta(X)$ in T_j such that α is either r_i or w_i and β is either r_j or w_j . Moreover, S is called "chary" if (i) no transaction in S uses the same element twice, and (ii) for every $T_i, T_j \in \{T_1, \ldots, T_n\}$, T_i and T_j share at most one element. Prove or disprove the following claims:

- 1.1 Every chary schedule is view-serializable.
- 1.2 Every chary schedule on two transactions is conflict-serializable.
- 1.3 Every chary schedule on two transactions is a 2PL schedule with exclusive and shared locks.

Problem 2

Let S be the schedule: $r_1(Z) w_3(Y) w_3(V) r_1(Y) r_2(V) w_2(Y) w_3(X) r_2(X) r_2(Z) r_3(Z) w_4(Z) w_4(X) w_2(X)$

- 2.1 Tell whether S is accepted by the 2PL scheduler with exclusive and shared locks. If the answer is yes, then specify the 2PL schedule obtained from S by adding suitable lock and unlock commands. If the answer is no, then explain the answer.
- 2.2 Tell whether S is view-serializable. If the answer is yes, then illustrate a serial schedule which is view-equivalent to S. If the answer is no, then explain the answer.
- 2.3 Answer all the following questions, motivating the answers: (i) Is S recoverable? (ii) Is S ACR? (iii) Is S strict?

Problem 3

Let $R(\underline{A},B,C)$, $S(\underline{A},D,E)$, T(A,B,C) be three tables (where T is a bag) and let τ indicate the ternary operator such that $\tau(R,S,T) = \delta(T \cup_b \pi_{A,B,C}(R \bowtie S))$, where δ denotes duplicate elimination, \cup_b denotes bag union and \bowtie denotes natural join.

- 3.1 Design and describe in detail a one pass algorithm that, given R,S,T as above, each one stored as a heap, computes $\tau(R,S,T)$.
- 3.2 Tell what is the weakest condition under which the algorithm can be used and illustrate the cost of the algorithm in terms of number of page accesses.
- 3.3 Tell what does it change if all the tables have A as key and are stored as sorted file with search key A.

Problem 4

Consider the relations Flight(code,company,type) with 1.000 pages and 10.000 tuples, and Ticket(number,code,company,type) with 2.000 pages and an associated index on Ticket with search key (company,type), for which we know that the cost of retrieving the records with a given value of attribute company is 3 page accesses. Assume a buffer with 50 frames, and consider the two queries shown below.

```
Query Q_1: select code, compa
```

select code, company from Flight except all —— not removing duplicates select code, company from Ticket

Query Q_2 :

select company, type from Flight except all -- not removing duplicates select company, type from Ticket

where "except all" denotes bag difference. For both queries Q_1 and Q_2 , tell (i) whether it is possible to process the query by using a block-nested loop algorithm, and (ii) whether it is possible to process the query by using an index-based algorithm. In all four cases, if the answer is positive, then describe the algorithm and tell which is its cost in terms of number of page accesses. If the answer is negative, then motivate the answer in detail.

Problem 5 (A.Y. 2021/22)

Describe in detail the notion of "star schema" in data warehousing and illustrate the difference between such a notion and the notion of "snowflake schema".

Problem 5 (A.Y. before 2021/22)

Consider the relations R(A,B,C) (with 1.500 pages and with one duplicate, in the average, for each tuple) and $S(\underline{D},E)$ (with 7.000 pages), and the query select distinct B,E from R, S where A <> 3.

Assume that the buffer contains 260 frames and show the logical query plan associated to the query, as well as the logical query plan and the physical query plan you would choose for executing the query as efficiently as possible. Also, tell which is the cost (in terms of number of page accesses) of executing the query according to the chosen physical query plan.

2)

In a schedule S on transactions $\{T_1, \ldots, T_n\}$ we say that two transactions T_i, T_j share the element X of the database if there exist actions $\alpha(X)$ in T_i and $\beta(X)$ in T_j such that α is either r_i or w_i and β is either r_j or w_j . Moreover, S is called "chary" if (i) no transaction in S uses the same element twice, and (ii) for every $T_i, T_j \in \{T_1, \ldots, T_n\}$, T_i and T_j share at most one element. Prove or disprove the following claims:

- 1.1 Every chary schedule is view-serializable.
- 1.2 Every chary schedule on two transactions is conflict-serializable.
- 1.3 Every chary schedule on two transactions is a 2PL schedule with exclusive and shared locks.

ı)	T	:	V	N,	(x))	r,	(y))				3	5 :	W	, (×	()	W	<u> (</u> (y)	r ₂	(×)	W	(3)	r	, (y	·)	W	3 (F	2)			
	T	2 :	1	W	(5	-)	rz	(x))																								
	T	3:		W;	3 (7	E)	W	3 (/)																7	3	->	·T	1				
			RI	EΑ	۵	-F	Ro	: דנ		(h.	ر)	٥,	W,	(x) >,	, د	r,	(4)	, ,	N3	(y)>			7	آ ء	-	>	Tз				
			FI	NA)L-	w	RI.	TE	:	W	,(;	د) ,	W	J3 (y)	, v	Nz	(3)						•	Ta	_	-7	Tz	2			
	5	•	ıs	1	VIE	W	-8	ER		ıF	т	HE	RE		E	XIZ	STS	5	A	S	CH	ιEι	DUI	E	3	`	se	RI	AL	S	UCH	TA	IAT
									f	₹E	AC) - (FR	רנס	1 (2	s):	2	RE	AI	٠ د	F	Ro	77 ((s,)									
									F	ĪW	AL	\	NR	١T	E ((s)	=	FIA	υA	L-	W	RI-	ΤE	(s,)								

- IF TWO TRANSACTIONS SHARE ONLY ONE ELEMENT, HEARS THAT THERE IS ONLY ONE CONFLICT IN A DIRECTION.

 ALSO BECAUSE WE CAN'T USE THE SAME VARIABLE TWICE.

 THERE FORE THE PRECEDENCE GRAPH IS ACYCLIC -> CONFLICT SER
- 3) T: w, (x) r, (y)
 T: w₂(x)

S: W, (x) W2(x) r, (y)

S: x2,(x) w,(x) S.L,(y) U,(x) x L2(x) W2(x) U2(x) r,(y) U,(y)

SINCE ONLY ONE VARIABLE IS SHARED, EVEN IF WE HAVE W, (x) W2 (x), WE CAN ANTICIPATE ALL LOCKS IN A BEFORE W2(x), AND UNLOCK U, (x)

Let S be the schedule: $r_1(Z) w_3(Y) w_3(V) r_1(Y) r_2(V) w_2(Y) w_3(X) r_2(X) r_2(Z) r_3(Z) w_4(Z) w_4(X) w_2(X)$

- 2.1 Tell whether S is accepted by the 2PL scheduler with exclusive and shared locks. If the answer is yes, then specify the 2PL schedule obtained from S by adding suitable lock and unlock commands. If the answer is no, then explain the answer.
- 2.2 Tell whether S is view-serializable. If the answer is yes, then illustrate a serial schedule which is view-equivalent to S. If the answer is no, then explain the answer.
- 2.3 Answer all the following questions, motivating the answers: (i) Is S recoverable? (ii) Is S ACR? (iii) Is S strict?

```
r_1(Z) w_3(Y) w_3(V) r_1(Y) r_2(V) w_2(Y) w_3(X) r_2(X) r_2(Z) r_3(Z) w_4(Z) w_4(X) w_2(X)
5: sl,(2) r(2) xl3(4) w3(4) xl3(6) m3(4) xl3(x) sl3(8) 03(4)
     SL, (Y) r, (Y) U3 (V) SL2 (V) r2 (V) U, (Y) xL2 (Y) W2 (Y) W3 (X) U3 (X)
     sl<sub>2</sub>(x) r<sub>2</sub>(x) sl<sub>2</sub>(z) r<sub>2</sub>(z) r<sub>3</sub>(z) U<sub>1</sub>(z)...
                                                   r2 (x) .. W4(x) ... W2 (x), 50
                      BECAUSE WE HAVE
IT'S NOT 2PL
                  ANTICIPATE
r_1(Z)\,w_3(Y)\,w_3(V)\,r_1(Y)\,r_2(V)\,w_2(Y)\,w_3(X)\,r_2(X)\,r_2(Z)\,r_3(Z)\,w_4(Z)\,w_4(X)\,w_2(X)
   READ-FROM: < 1, (4), W3 (4) > , < r2 (4), W3 (x) > , < r2 (x), W3 (x) >
                                                                        T_3 \rightarrow T_1 T_2 \rightarrow T_2
   FINAL - WRITE : W2 (x), W4 (2), W2 (y), W3 (V)
        VIEW-SER IF THERE EXISTS A SCHEDULE S' SERIAL SUCH THAT
                      READ - FROM (S) = READ - FROM (S')
                      FINAL-WRITE (S) = FINAL-WRITE (S')
     T4 CHANGE FINAL WRITE FOR X
                                                        NOT
                                                                VIEW
 Ty CHANGE READ FROM FOR 2
r_1(Z) w_3(Y) w_3(V) r_1(Y) r_2(V) w_2(Y) w_3(X) r_2(X) r_2(Z) r_3(Z) w_4(Z) w_4(X) w_2(X)
```

READ-FROM: < 1, (Y), W3 (Y)>, < r2 (V), W3 (V)>, < r2 (x), W3 (x)>

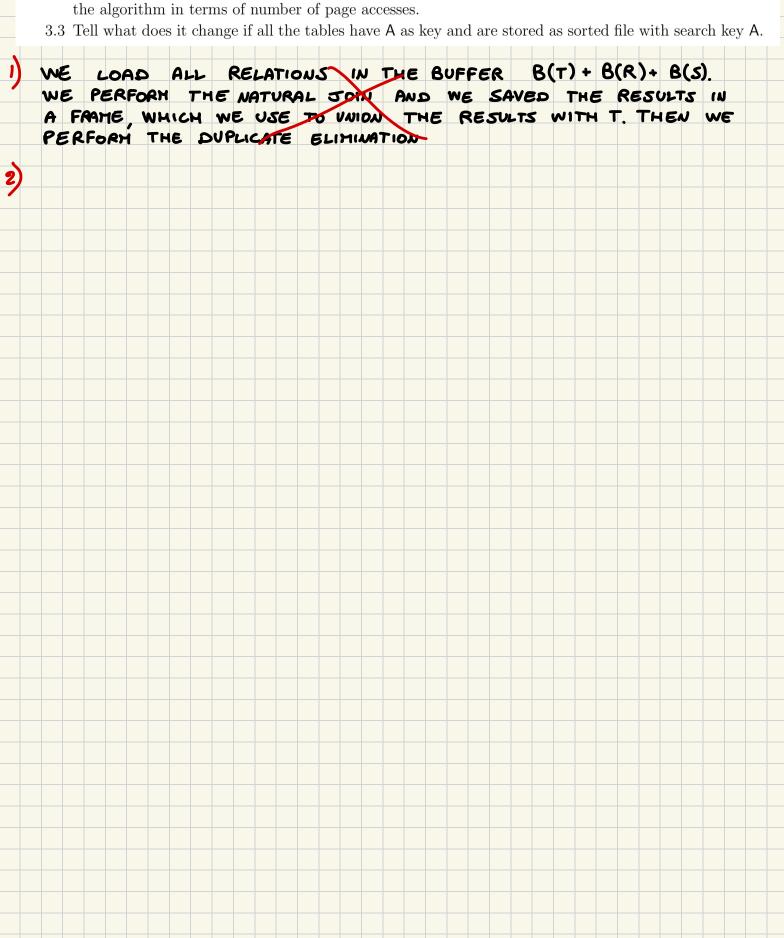
63 BEFORE C, 63 BEFORE 62 RECOVERABLE

NOT ALR BELAUSE SOME T READ BEFORE COMMITS

NOT ACR => NOT STRICT

Let $R(\underline{A},B,C)$, $S(\underline{A},D,E)$, T(A,B,C) be three tables (where T is a bag) and let τ indicate the ternary operator such that $\tau(R,S,T) = \delta(T \cup_b \pi_{A,B,C}(R \bowtie S))$, where δ denotes duplicate elimination, \cup_b denotes bag union and \bowtie denotes natural join.

- 3.1 Design and describe in detail a one pass algorithm that, given R,S,T as above, each one stored as a heap, computes $\tau(R,S,T)$.
- 3.2 Tell what is the weakest condition under which the algorithm can be used and illustrate the cost of the algorithm in terms of number of page accesses.



Consider the relations Flight(code,company,type) with 1.000 pages and 10.000 tuples, and Ticket(number,code,company,type) with 2.000 pages and an associated index on Ticket with search key (company,type), for which we know that the cost of retrieving the records with a given value of attribute company is 3 page accesses. Assume a buffer with 50 frames, and consider the two queries shown below.

Query Q_1 :

select code, company from Flight except all -- not removing duplicates select code, company from Ticket

Query Q_2 :

select company, type from Flight except all $\,--\,not\,\,removing\,\,duplicates$ select company, type from Ticket

where "except all" denotes bag difference. For both queries Q_1 and Q_2 , tell (i) whether it is possible to process the query by using a block-nested loop algorithm, and (ii) whether it is possible to process the query by using an index-based algorithm. In all four cases, if the answer is positive, then describe the algorithm and tell which is its cost in terms of number of page accesses. If the answer is negative, then motivate the answer in detail.

