

Basi di dati II

Prova parziale — 11 aprile 2012 — Compito A

Rispondere su questo fascicolo.

Tempo a disposizione: un'ora e quindici minuti.

Cognome _____ Nome _____ Matricola _____ Ordin. _____

Domanda 1 (20%)

Considerare un sistema con dischi con $N = 400$ blocchi per traccia

- tempo medio di posizionamento della testina (tempo di seek) $t_S = 5$ msec
- tempo medio di latenza (attesa dovuta alla rotazione) $t_L = 3$ msec
- tempo minimo di lettura di un blocco $t_B = 15$ μ sec

Rispondere alle seguenti domande mostrando formula e valore numerico (N.B. non servono calcolatrici, i calcoli sono semplici; **se si ritengono le informazioni imprecise, dare risposte approssimate, usando il buon senso**).

1. Qual è il tempo medio necessario per leggere un blocco del quale sia dato l'indirizzo fisico?

2. Qual è il tempo medio necessario per la scansione sequenziale di un file costituito da $F = 100$ blocchi contigui, non letti di recente?

3. Qual è il tempo che si può ipotizzare necessario per eseguire un accesso diretto ad un record di un file attraverso un indice che abbia profondità $p = 4$ e fan-out (fattore di blocco dell'indice) $f_I = 100$, usato di recente, ma in modo non molto intenso?

4. Qual è il tempo che si può ipotizzare necessario per eseguire $m = 1000$ accessi diretti in tempi ravvicinati a record di un file (molto grande) attraverso un indice che abbia profondità $p = 4$, fan-out $f_I = 100$, con disponibilità di circa $P = 150$ pagine di buffer?

Domanda 2 (20%)

Si consideri una base di dati con le relazioni

- R1(A,B,C,D) con
 - $N_1=2.500.000$ ennuple di lunghezza $l_1 = 20$ Byte
 - $b=5$ valori diversi sull'attributo B (tutti gli interi compresi fra 1 e b)
 - $c=25.000$ valori diversi sull'attributo C (tutti gli interi compresi fra 1 e c)
 - una struttura disordinata, un indice sulla chiave primaria A e un altro sull'attributo C;
- R2(E,F,G) con
 - $N_2=1.000.000$ ennuple di lunghezza $l_2 = 100$ Byte
 - una struttura disordinata, un indice sulla chiave primaria E

Supporre che:

- i blocchi abbiano dimensione $P = 2$ KByte (approssimabile a 2000 Byte);
- ogni operazione possa contare su un numero di pagine di buffer pari a circa $q=300$;
- gli indici abbiano tutti $p=4$ livelli (contando anche radice e foglie) e fattore di blocco massimo $f_i=100$;
- il sistema esegua i join come nested loop oppure come hash-join (si ricorda che questi ultimi hanno un costo pari a circa tre volte la somma dei blocchi dei due file, a condizione che il quadrato del numero di pagine di buffer disponibili sia maggiore del numero di blocchi del più piccolo dei due file).

Valutare il costo, indicandolo in modo sia simbolico sia numerico e specificando quale algoritmo si utilizza per il join di ciascuna delle interrogazioni seguenti (NB: **al di là del dettaglio, la cui precisione non è forse nemmeno possibile, è essenziale mostrare una comprensione degli ordini di grandezza**):

```
select *
from R1 join R2 on D=E
```

```
select *
from R1 join R2 on D=E
where B=5
```

```
select *
from R1 join R2 on D=E
where C=502
```

Domanda 3 (20%)

Per ciascuno degli schedule sotto riportati, indicare, scrivendo **si** o **no** in ciascuna casella, a quali classi appartiene: S (seriale, rispetto a letture e scritture, ignorare commit e abort), CSR (conflict-serializzabile), S2PL (cioè generabile da uno scheduler basato su 2PL stretto), MV (cioè generabile da uno scheduler multiversion con controllo di serializzabilità: “a serializable transaction cannot modify or lock rows changed by other transactions after the serializable transaction began”). Negli schedule, s_i indica l’inizio della transazione i e c_i il suo commit.

	S	CSR	S2PL	MV
$s_2, s_1, r_2(x), w_2(x), c_2, r_1(x), w_1(x), c_1$				
$s_2, r_2(x), w_2(x), c_2, s_1, r_1(x), w_1(x), c_1$				
$s_1, r_1(x), s_2, r_2(x), w_1(x), c_1, w_2(x), c_2$				
$s_1, r_1(x), s_2, r_2(x), w_2(x), r_2(y), w_2(y), c_2, r_1(y), c_1$				

Domanda 4 (20%)

Una catena di supermercati ha una base di dati dei propri clienti che dispongono di una “tessera fedeltà,” con varie informazioni su ciascun cliente, fra cui (a) il totale dei punti acquisiti attraverso l’uso della tessera e (b) il negozio della catena cui fa riferimento (ad esempio, quello presso cui ha inizialmente richiesto la tessera). Si vuole eseguire su di essa l’interrogazione che calcola, per ciascun negozio, il numero dei clienti, la somma dei punti fedeltà dei clienti e la relativa media per cliente. Indicare quale livello di isolamento (**READ UNCOMMITTED**, **READ COMMITTED**, **REPEATABLE READ** o **SERIALIZABLE**) si potrebbe scegliere in ciascuno dei seguenti casi (si supponga che, in generale, sia stato rilevato che, nel corso degli inserimenti e delle modifiche, vengono inseriti valori sbagliati anche di vari ordini di grandezza, che sono poi corretti prima del commit):

1. L’operazione è eseguita mentre vengono inseriti alcuni nuovi clienti (per ciascun negozio pochi rispetto a quelli già presenti), con la finalità di acquisire informazioni approssimate ma ragionevolmente indicative sugli andamenti complessivi.
2. L’operazione è eseguita mentre vengono inseriti molti nuovi clienti, con la finalità di acquisire informazioni approssimate ma ragionevolmente indicative sugli andamenti complessivi.
3. L’operazione è eseguita mentre vengono inseriti alcuni nuovi clienti (per ciascun negozio pochi rispetto a quelli già presenti), con la finalità di individuare i primi tre negozi da premiare in una campagna promozionale sulla base dei punti acquisiti dai rispettivi clienti.
4. L’operazione è eseguita mentre vengono modificati i valori dei punti fedeltà di tutti i clienti (a seguito di una ridefinizione dei criteri di assegnazione dei punti stessi), con la finalità di individuare i primi tre negozi da premiare in una campagna promozionale sulla base dei punti acquisiti dai rispettivi clienti.
5. L’operazione è eseguita in un momento in cui non ci sono aggiornamenti di alcun genere, con la finalità di acquisire informazioni approssimate ma ragionevolmente indicative sugli andamenti complessivi.

1.	2.	3.	4.	5.

Domanda 5 (20%)

Si considerino due relazioni $R_1(\underline{A}, C)$, $R_2(\underline{A}, D, E, F)$, in cui gli attributi hanno tutti la stessa dimensione $L = 4\text{Byte}$, molto più piccola della dimensione del blocco pari a $B = 4000\text{ Byte}$. Si supponga che le relazioni abbiano entrambe $N = 1.000.000$ ennuple, con gli stessi valori su A , e che le operazioni più frequenti su di essa siano le seguenti:

- o_1 selezione di una ennupla del join di R_1 e R_2 (sulla base del valore di A), con frequenza $f_1 = 10.000$;
- o_2 scansione dell'intera relazione R_1 , con frequenza $f_2 = 1$

Valutare le due seguenti alternative di memorizzazione, calcolando il costo complessivo (riportare la formula che indica il numero di accessi nell'unità di tempo, in base alle variabili sopra citate):

- (i) memorizzazione separata delle due relazioni, entrambe ordinate su A e con indice primario su A con profondità $p = 4$, con 2 livelli mediamente disponibili nel buffer.

costo unitario di o_1 :	costo unitario di o_2 :
costo complessivo:	

- (ii) memorizzazione in un cluster delle due relazioni pure entrambe ordinate su A e con indice primario su A con profondità sempre $p = 4$, con 2 livelli mediamente disponibili nel buffer.

costo unitario di o_1 :	costo unitario di o_2 :
costo complessivo:	

In conclusione, conviene quindi la memorizzazione in un cluster? (Sì o No)

Basi di dati II

Prova parziale — 11 aprile 2012 — Compito B

Rispondere su questo fascicolo.

Tempo a disposizione: un'ora e quindici minuti.

Cognome _____ Nome _____ Matricola _____ Ordin. _____

Domanda 1 (20%)

Considerare un sistema con dischi con $N = 400$ blocchi per traccia

- tempo medio di posizionamento della testina (tempo di seek) $t_S = 4$ msec
- tempo medio di latenza (attesa dovuta alla rotazione) $t_L = 4$ msec
- tempo minimo di lettura di un blocco $t_B = 15 \mu\text{sec}$

Rispondere alle seguenti domande mostrando formula e valore numerico (N.B. non servono calcolatrici, i calcoli sono semplici; **se si ritengono le informazioni imprecise, dare risposte approssimate, usando il buon senso**).

1. Qual è il tempo medio necessario per leggere un blocco del quale sia dato l'indirizzo fisico?

2. Qual è il tempo medio necessario per la scansione sequenziale di un file costituito da $F = 100$ blocchi contigui, non letti di recente?

3. Qual è il tempo che si può ipotizzare necessario per eseguire un accesso diretto ad un record di un file attraverso un indice che abbia profondità $p = 4$ e fan-out (fattore di blocco dell'indice) $f_I = 100$, usato di recente, ma in modo non molto intenso?

4. Qual è il tempo che si può ipotizzare necessario per eseguire $m = 2000$ accessi diretti in tempi ravvicinati a record di un file (molto grande) attraverso un indice che abbia profondità $p = 4$, fan-out $f_I = 100$, con disponibilità di circa $P = 150$ pagine di buffer?

Domanda 2 (20%)

Si consideri una base di dati con le relazioni

- R1(A,B,C,D) con
 - $L_1=1.500.000$ ennuple di lunghezza $l_1 = 40$ Byte
 - $b=5$ valori diversi sull'attributo B (tutti gli interi compresi fra 1 e b)
 - $c=15.000$ valori diversi sull'attributo C (tutti gli interi compresi fra 1 e c)
 - una struttura disordinata, un indice sulla chiave primaria A e un altro sull'attributo C;
- R2(E,F,G) con
 - $L_2=2.000.000$ ennuple di lunghezza $l_2 = 200$ Byte
 - una struttura disordinata, un indice sulla chiave primaria E

Supporre che:

- i blocchi abbiano dimensione $P = 4$ KByte (approssimabile a 4000 Byte);
- ogni operazione possa contare su un numero di pagine di buffer pari a circa $q=300$;
- gli indici abbiano tutti $i=4$ livelli (contando anche radice e foglie) e fattore di blocco massimo $f_i=100$;
- il sistema esegua i join come nested loop oppure come hash-join (si ricorda che questi ultimi hanno un costo pari a circa tre volte la somma dei blocchi dei due file, a condizione che il quadrato del numero di pagine di buffer disponibili sia maggiore del numero di blocchi del più piccolo dei due file).

Valutare il costo, indicandolo in modo sia simbolico sia numerico e specificando quale algoritmo si utilizza per il join di ciascuna delle interrogazioni seguenti (NB: **al di là del dettaglio, la cui precisione non è forse nemmeno possibile, è essenziale mostrare una comprensione degli ordini di grandezza**):

```
select *
from R1 join R2 on D=E
```

```
select *
from R1 join R2 on D=E
where B=5
```

```
select *
from R1 join R2 on D=E
where C=502
```

Domanda 3 (20%)

Per ciascuno degli schedule sotto riportati, indicare, scrivendo **si** o **no** in ciascuna casella, a quali classi appartiene: S (seriale, rispetto a letture e scritture, ignorare commit e abort), CSR (conflict-serializzabile), S2PL (cioè generabile da uno scheduler basato su 2PL stretto), MV (cioè generabile da uno scheduler multiversion con controllo di serializzabilità: “a serializable transaction cannot modify or lock rows changed by other transactions after the serializable transaction began”). Negli schedule, s_i indica l’inizio della transazione i e c_i il suo commit.

	S	CSR	S2PL	MV
$s_2, r_2(x), w_2(x), c_2, s_1, r_1(x), w_1(x), c_1$				
$s_1, r_1(x), s_2, r_2(x), w_1(x), c_1, w_2(x), c_2$				
$s_2, s_1, r_2(x), w_2(x), c_2, r_1(x), w_1(x), c_1$				
$s_1, r_1(x), s_2, r_2(x), w_2(x), r_2(y), w_2(y), c_2, r_1(y), c_1$				

Domanda 4 (20%)

Una catena di supermercati ha una base di dati dei propri clienti che dispongono di una “tessera fedeltà,” con varie informazioni su ciascun cliente, fra cui (a) il totale dei punti acquisiti attraverso l’uso della tessera e (b) il negozio della catena cui fa riferimento (ad esempio, quello presso cui ha inizialmente richiesto la tessera). Si vuole eseguire su di essa l’interrogazione che calcola, per ciascun negozio, il numero dei clienti, la somma dei punti fedeltà dei clienti e la relativa media per cliente. Indicare quale livello di isolamento (**READ UNCOMMITTED**, **READ COMMITTED**, **REPEATABLE READ** o **SERIALIZABLE**) si potrebbe scegliere in ciascuno dei seguenti casi (si supponga che, in generale, sia stato rilevato che, nel corso degli inserimenti e delle modifiche, vengono inseriti valori sbagliati anche di vari ordini di grandezza, che sono poi corretti prima del commit):

1. L’operazione è eseguita mentre vengono inseriti molti nuovi clienti, con la finalità di acquisire informazioni approssimate ma ragionevolmente indicative sugli andamenti complessivi.
2. L’operazione è eseguita mentre vengono modificati i valori dei punti fedeltà di alcuni clienti (in ciascun negozio una piccola percentuale), con la finalità di acquisire informazioni approssimate ma ragionevolmente indicative sugli andamenti complessivi.
3. L’operazione è eseguita in un momento in cui non ci sono aggiornamenti di alcun genere, con la finalità di individuare i primi tre negozi da premiare in una campagna promozionale sulla base dei punti acquisiti dai rispettivi clienti.
4. L’operazione è eseguita mentre vengono inseriti molti nuovi clienti, con la finalità di individuare i primi tre negozi da premiare in una campagna promozionale sulla base dei punti acquisiti dai rispettivi clienti.
5. L’operazione è eseguita mentre vengono modificati i valori dei punti fedeltà di alcuni clienti (in ciascun negozio una piccola percentuale), con la finalità di individuare i primi tre negozi da premiare in una campagna promozionale sulla base dei punti acquisiti dai rispettivi clienti.

1.	2.	3.	4.	5.

Domanda 5 (20%)

Si considerino due relazioni $R_1(\underline{A}, C)$, $R_2(\underline{A}, D, E, F)$, in cui gli attributi hanno tutti la stessa dimensione $a = 4\text{Byte}$, molto più piccola della dimensione del blocco pari a $B = 4000\text{ Byte}$. Si supponga che le relazioni abbiano entrambe $L = 1.000.000$ ennuple, con gli stessi valori su A , e che le operazioni più frequenti su di essa siano le seguenti:

- o_1 selezione di una ennupla del join di R_1 e R_2 (sulla base del valore di A), con frequenza $f_1 = 1000$;
- o_2 scansione dell'intera relazione R_1 , con frequenza $f_2 = 10$

Valutare le due seguenti alternative di memorizzazione, calcolando il costo complessivo (riportare la formula che indica il numero di accessi nell'unità di tempo, in base alle variabili sopra citate):

- (i) memorizzazione separata delle due relazioni, entrambe ordinate su A e con indice primario su A con profondità $p = 4$, con 2 livelli mediamente disponibili nel buffer.

costo unitario di o_1 :	costo unitario di o_2 :
costo complessivo:	

- (ii) memorizzazione in un cluster delle due relazioni pure entrambe ordinate su A e con indice primario su A con profondità sempre $p = 4$, con 2 livelli mediamente disponibili nel buffer.

costo unitario di o_1 :	costo unitario di o_2 :
costo complessivo:	

In conclusione, conviene quindi la memorizzazione in un cluster? (Sì o No)

Basi di dati II

Prova parziale — 11 aprile 2012 — Compito C

Rispondere su questo fascicolo.

Tempo a disposizione: un'ora e quindici minuti.

Cognome _____ Nome _____ Matricola _____ Ordin. _____

Domanda 1 (20%)

Considerare un sistema con dischi con $N = 400$ blocchi per traccia

- tempo medio di posizionamento della testina (tempo di seek) $t_S = 4$ msec
- tempo medio di latenza (attesa dovuta alla rotazione) $t_L = 4$ msec
- tempo minimo di lettura di un blocco $t_B = 15$ μ sec

Rispondere alle seguenti domande mostrando formula e valore numerico (N.B. non servono calcolatrici, i calcoli sono semplici; **se si ritengono le informazioni imprecise, dare risposte approssimate, usando il buon senso**).

1. Qual è il tempo medio necessario per leggere un blocco del quale sia dato l'indirizzo fisico?

2. Qual è il tempo medio necessario per la scansione sequenziale di un file costituito da $F = 100$ blocchi contigui, non letti di recente?

3. Qual è il tempo che si può ipotizzare necessario per eseguire un accesso diretto ad un record di un file attraverso un indice che abbia profondità $p = 4$ e fan-out (fattore di blocco dell'indice) $f_I = 100$, usato di recente, ma in modo non molto intenso?

4. Qual è il tempo che si può ipotizzare necessario per eseguire $m = 1000$ accessi diretti in tempi ravvicinati a record di un file (molto grande) attraverso un indice che abbia profondità $p = 4$, fan-out $f_I = 100$, con disponibilità di circa $P = 150$ pagine di buffer?

Domanda 2 (20%)

Si consideri una base di dati con le relazioni

- R1(A,B,C,D) con
 - $N_1=2.500.000$ ennuple di lunghezza $l_1 = 20$ Byte
 - $b=5$ valori diversi sull'attributo B (tutti gli interi compresi fra 1 e b)
 - $c=25.000$ valori diversi sull'attributo C (tutti gli interi compresi fra 1 e c)
 - una struttura disordinata, un indice sulla chiave primaria A e un altro sull'attributo C;
- R2(E,F,G) con
 - $N_2=1.000.000$ ennuple di lunghezza $l_2 = 100$ Byte
 - una struttura disordinata, un indice sulla chiave primaria E

Supporre che:

- i blocchi abbiano dimensione $P = 2$ KByte (approssimabile a 2000 Byte);
- ogni operazione possa contare su un numero di pagine di buffer pari a circa $q=300$;
- gli indici abbiano tutti $i=4$ livelli (contando anche radice e foglie) e fattore di blocco massimo $f_i=100$;
- il sistema esegua i join come nested loop oppure come hash-join (si ricorda che questi ultimi hanno un costo pari a circa tre volte la somma dei blocchi dei due file, a condizione che il quadrato del numero di pagine di buffer disponibili sia maggiore del numero di blocchi del più piccolo dei due file).

Valutare il costo, indicandolo in modo sia simbolico sia numerico e specificando quale algoritmo si utilizza per il join di ciascuna delle interrogazioni seguenti (NB: **al di là del dettaglio, la cui precisione non è forse nemmeno possibile, è essenziale mostrare una comprensione degli ordini di grandezza**):

```
select *
from R1 join R2 on D=E
```

```
select *
from R1 join R2 on D=E
where B=5
```

```
select *
from R1 join R2 on D=E
where C=502
```

Domanda 3 (20%)

Per ciascuno degli schedule sotto riportati, indicare, scrivendo **sì** o **no** in ciascuna casella, a quali classi appartiene: S (seriale, rispetto a letture e scritture, ignorare commit e abort), CSR (conflict-serializzabile), S2PL (cioè generabile da uno scheduler basato su 2PL stretto), MV (cioè generabile da uno scheduler multiversion con controllo di serializzabilità: “a serializable transaction cannot modify or lock rows changed by other transactions after the serializable transaction began”). Negli schedule, s_i indica l’inizio della transazione i e c_i il suo commit.

	S	CSR	S2PL	MV
$s_2, s_1, r_2(x), w_2(x), c_2, r_1(x), w_1(x), c_1$				
$s_2, r_2(x), w_2(x), c_2, s_1, r_1(x), w_1(x), c_1$				
$s_1, r_1(x), s_2, r_2(x), w_2(x), r_2(y), w_2(y), c_2, r_1(y), c_1$				
$s_1, r_1(x), s_2, r_2(x), w_1(x), c_1, w_2(x), c_2$				

Domanda 4 (20%)

Una catena di supermercati ha una base di dati dei propri clienti che dispongono di una “tessera fedeltà,” con varie informazioni su ciascun cliente, fra cui (a) il totale dei punti acquisiti attraverso l’uso della tessera e (b) il negozio della catena cui fa riferimento (ad esempio, quello presso cui ha inizialmente richiesto la tessera). Si vuole eseguire su di essa l’interrogazione che calcola, per ciascun negozio, il numero dei clienti, la somma dei punti fedeltà dei clienti e la relativa media per cliente. Indicare quale livello di isolamento (**READ UNCOMMITTED**, **READ COMMITTED**, **REPEATABLE READ** o **SERIALIZABLE**) si potrebbe scegliere in ciascuno dei seguenti casi (si supponga che, in generale, sia stato rilevato che, nel corso degli inserimenti e delle modifiche, vengono inseriti valori sbagliati anche di vari ordini di grandezza, che sono poi corretti prima del commit):

1. L’operazione è eseguita mentre vengono inseriti alcuni nuovi clienti (per ciascun negozio pochi rispetto a quelli già presenti), con la finalità di acquisire informazioni approssimate ma ragionevolmente indicative sugli andamenti complessivi.
2. L’operazione è eseguita mentre vengono modificati i valori dei punti fedeltà di tutti i clienti (a seguito di una ridefinizione dei criteri di assegnazione dei punti stessi), con la finalità di acquisire informazioni approssimate ma ragionevolmente indicative sugli andamenti complessivi.
3. L’operazione è eseguita in un momento in cui non ci sono aggiornamenti di alcun genere, con la finalità di individuare i primi tre negozi da premiare in una campagna promozionale sulla base dei punti acquisiti dai rispettivi clienti.
4. L’operazione è eseguita mentre vengono inseriti alcuni nuovi clienti (per ciascun negozio pochi rispetto a quelli già presenti), con la finalità di individuare i primi tre negozi da premiare in una campagna promozionale sulla base dei punti acquisiti dai rispettivi clienti.
5. L’operazione è eseguita mentre vengono modificati i valori dei punti fedeltà di tutti i clienti (a seguito di una ridefinizione dei criteri di assegnazione dei punti stessi), con la finalità di individuare i primi tre negozi da premiare in una campagna promozionale sulla base dei punti acquisiti dai rispettivi clienti.

1.	2.	3.	4.	5.

Domanda 5 (20%)

Si considerino due relazioni $R_1(\underline{A}, C)$, $R_2(\underline{A}, D, E, F)$, in cui gli attributi hanno tutti la stessa dimensione $a = 4\text{Byte}$, molto più piccola della dimensione del blocco pari a $B = 4000\text{ Byte}$. Si supponga che le relazioni abbiano entrambe $L = 1.000.000$ ennuple, con gli stessi valori su A , e che le operazioni più frequenti su di essa siano le seguenti:

- o_1 selezione di una ennupla del join di R_1 e R_2 (sulla base del valore di A), con frequenza $f_1 = 10.000$;
- o_2 scansione dell'intera relazione R_1 , con frequenza $f_2 = 1$

Valutare le due seguenti alternative di memorizzazione, calcolando il costo complessivo (riportare la formula che indica il numero di accessi nell'unità di tempo, in base alle variabili sopra citate):

- (i) memorizzazione separata delle due relazioni, entrambe ordinate su A e con indice primario su A con profondità $p = 4$, con 2 livelli mediamente disponibili nel buffer.

costo unitario di o_1 :	costo unitario di o_2 :
costo complessivo:	

- (ii) memorizzazione in un cluster delle due relazioni pure entrambe ordinate su A e con indice primario su A con profondità sempre $p = 4$, con 2 livelli mediamente disponibili nel buffer.

costo unitario di o_1 :	costo unitario di o_2 :
costo complessivo:	

In conclusione, conviene quindi la memorizzazione in un cluster? (Sì o No)

Basi di dati II

Prova parziale — 11 aprile 2012 — Compito D

Rispondere su questo fascicolo.

Tempo a disposizione: un'ora e quindici minuti.

Cognome _____ Nome _____ Matricola _____ Ordin. _____

Domanda 1 (20%)

Considerare un sistema con dischi con $N = 400$ blocchi per traccia

- tempo medio di posizionamento della testina (tempo di seek) $t_S = 5$ msec
- tempo medio di latenza (attesa dovuta alla rotazione) $t_L = 3$ msec
- tempo minimo di lettura di un blocco $t_B = 15$ μ sec

Rispondere alle seguenti domande mostrando formula e valore numerico (N.B. non servono calcolatrici, i calcoli sono semplici; **se si ritengono le informazioni imprecise, dare risposte approssimate, usando il buon senso**).

1. Qual è il tempo medio necessario per leggere un blocco del quale sia dato l'indirizzo fisico?

2. Qual è il tempo medio necessario per la scansione sequenziale di un file costituito da $F = 100$ blocchi contigui, non letti di recente?

3. Qual è il tempo che si può ipotizzare necessario per eseguire un accesso diretto ad un record di un file attraverso un indice che abbia profondità $p = 4$ e fan-out (fattore di blocco dell'indice) $f_I = 100$, usato di recente, ma in modo non molto intenso?

4. Qual è il tempo che si può ipotizzare necessario per eseguire $m = 2000$ accessi diretti in tempi ravvicinati a record di un file (molto grande) attraverso un indice che abbia profondità $p = 4$, fan-out $f_I = 100$, con disponibilità di circa $P = 150$ pagine di buffer?

Domanda 2 (20%)

Si consideri una base di dati con le relazioni

- R1(A,B,C,D) con
 - $L_1=1.500.000$ ennuple di lunghezza $l_1 = 40$ Byte
 - $b=5$ valori diversi sull'attributo B (tutti gli interi compresi fra 1 e b)
 - $c=15.000$ valori diversi sull'attributo C (tutti gli interi compresi fra 1 e c)
 - una struttura disordinata, un indice sulla chiave primaria A e un altro sull'attributo C;
- R2(E,F,G) con
 - $L_2=2.000.000$ ennuple di lunghezza $l_2 = 200$ Byte
 - una struttura disordinata, un indice sulla chiave primaria E

Supporre che:

- i blocchi abbiano dimensione $P = 4$ KByte (approssimabile a 4000 Byte);
- ogni operazione possa contare su un numero di pagine di buffer pari a circa $q=300$;
- gli indici abbiano tutti $p=4$ livelli (contando anche radice e foglie) e fattore di blocco massimo $f_i=100$;
- il sistema esegua i join come nested loop oppure come hash-join (si ricorda che questi ultimi hanno un costo pari a circa tre volte la somma dei blocchi dei due file, a condizione che il quadrato del numero di pagine di buffer disponibili sia maggiore del numero di blocchi del più piccolo dei due file).

Valutare il costo, indicandolo in modo sia simbolico sia numerico e specificando quale algoritmo si utilizza per il join di ciascuna delle interrogazioni seguenti (NB: **al di là del dettaglio, la cui precisione non è forse nemmeno possibile, è essenziale mostrare una comprensione degli ordini di grandezza**):

```
select *
from R1 join R2 on D=E
```

```
select *
from R1 join R2 on D=E
where B=5
```

```
select *
from R1 join R2 on D=E
where C=502
```

Domanda 3 (20%)

Per ciascuno degli schedule sotto riportati, indicare, scrivendo **si** o **no** in ciascuna casella, a quali classi appartiene: S (seriale, rispetto a letture e scritture, ignorare commit e abort), CSR (conflict-serializzabile), S2PL (cioè generabile da uno scheduler basato su 2PL stretto), MV (cioè generabile da uno scheduler multiversion con controllo di serializzabilità: “a serializable transaction cannot modify or lock rows changed by other transactions after the serializable transaction began”). Negli schedule, s_i indica l’inizio della transazione i e c_i il suo commit.

	S	CSR	S2PL	MV
$s_1, r_1(x), s_2, r_2(x), w_2(x), r_2(y), w_2(y), c_2, r_1(y), c_1$				
$s_1, r_1(x), s_2, r_2(x), w_1(x), c_1, w_2(x), c_2$				
$s_2, s_1, r_2(x), w_2(x), c_2, r_1(x), w_1(x), c_1$				
$s_2, r_2(x), w_2(x), c_2, s_1, r_1(x), w_1(x), c_1$				

Domanda 4 (20%)

Una catena di supermercati ha una base di dati dei propri clienti che dispongono di una “tessera fedeltà,” con varie informazioni su ciascun cliente, fra cui (a) il totale dei punti acquisiti attraverso l’uso della tessera e (b) il negozio della catena cui fa riferimento (ad esempio, quello presso cui ha inizialmente richiesto la tessera). Si vuole eseguire su di essa l’interrogazione che calcola, per ciascun negozio, il numero dei clienti, la somma dei punti fedeltà dei clienti e la relativa media per cliente. Indicare quale livello di isolamento (**READ UNCOMMITTED**, **READ COMMITTED**, **REPEATABLE READ** o **SERIALIZABLE**) si potrebbe scegliere in ciascuno dei seguenti casi (si supponga che, in generale, sia stato rilevato che, nel corso degli inserimenti e delle modifiche, vengono inseriti valori sbagliati anche di vari ordini di grandezza, che sono poi corretti prima del commit):

1. L’operazione è eseguita mentre vengono modificati i valori dei punti fedeltà di tutti i clienti (a seguito di una ridefinizione dei criteri di assegnazione dei punti stessi), con la finalità di acquisire informazioni approssimate ma ragionevolmente indicative sugli andamenti complessivi.
2. L’operazione è eseguita mentre vengono inseriti molti nuovi clienti, con la finalità di individuare i primi tre negozi da premiare in una campagna promozionale sulla base dei punti acquisiti dai rispettivi clienti.
3. L’operazione è eseguita mentre vengono modificati i valori dei punti fedeltà di alcuni clienti (in ciascun negozio una piccola percentuale), con la finalità di acquisire informazioni approssimate ma ragionevolmente indicative sugli andamenti complessivi.
4. L’operazione è eseguita mentre vengono modificati i valori dei punti fedeltà di alcuni clienti (in ciascun negozio una piccola percentuale), con la finalità di individuare i primi tre negozi da premiare in una campagna promozionale sulla base dei punti acquisiti dai rispettivi clienti.
5. L’operazione è eseguita in un momento in cui non ci sono aggiornamenti di alcun genere, con la finalità di individuare i primi tre negozi da premiare in una campagna promozionale sulla base dei punti acquisiti dai rispettivi clienti.

1.	2.	3.	4.	5.

Domanda 5 (20%)

Si considerino due relazioni $R_1(\underline{A}, C)$, $R_2(\underline{A}, D, E, F)$, in cui gli attributi hanno tutti la stessa dimensione $L = 4\text{Byte}$, molto più piccola della dimensione del blocco pari a $B = 4000\text{ Byte}$. Si supponga che le relazioni abbiano entrambe $N = 1.000.000$ ennuple, con gli stessi valori su A , e che le operazioni più frequenti su di essa siano le seguenti:

- o_1 selezione di una ennupla del join di R_1 e R_2 (sulla base del valore di A), con frequenza $f_1 = 1000$;
- o_2 scansione dell'intera relazione R_1 , con frequenza $f_2 = 10$

Valutare le due seguenti alternative di memorizzazione, calcolando il costo complessivo (riportare la formula che indica il numero di accessi nell'unità di tempo, in base alle variabili sopra citate):

- (i) memorizzazione separata delle due relazioni, entrambe ordinate su A e con indice primario su A con profondità $p = 4$, con 2 livelli mediamente disponibili nel buffer.

costo unitario di o_1 :	costo unitario di o_2 :
costo complessivo:	

- (ii) memorizzazione in un cluster delle due relazioni pure entrambe ordinate su A e con indice primario su A con profondità sempre $p = 4$, con 2 livelli mediamente disponibili nel buffer.

costo unitario di o_1 :	costo unitario di o_2 :
costo complessivo:	

In conclusione, conviene quindi la memorizzazione in un cluster? (Sì o No)