

Vypracované otázky **PB154 - Základy databázových systémů** (částečně recyklované) z fi.muny.cz.

Pokud se vám hodily už vypracované otázky a považujete za cenné, když jsou dostupné otázky z minulých let, pokuste se zapamatovat si ty své a následně je dát na fi.muny.cz. Budoucí generace to ocení. :-)

Popište, co znamená?

Cizí klíč: v prostředí relačních databází definuje vztah mezi dvěma tabulkami takový, že hodnota v určeném sloupci musí existovat v jiné (primární) tabulce. Tím je definováno integritní omezení, které do tabulky položky umožní vložit jen povolené hodnoty. Je tím vlastně vytvořeno spojení jednoho nebo více sloupců se sloupcem nebo více sloupci jiné („cizí“) tabulky. Tomu se též říká reference nebo odkaz.

Cizí klíč umožňuje definovat akce, které mají nastat při pokusu o změnu nebo mazání záznamů v cizí tabulce.

Superklíč: Určuje každý řádek tabulky jednoznačně, ale není nutně minimální.

Kandidátní klíč: v relačním modelování označuje sloupec nebo kombinaci sloupců, ve kterých mají všechny řádky tabulky své hodnoty unikátní. Každý kandidátní klíč tak umožňuje jednoznačně identifikovat každý řádek tabulky. Kandidátní klíč je nejmenší možná taková množina sloupců.

Primární klíč: je pole nebo kombinace polí (jeden vybraný z kandidátních), jednoznačně identifikující každý záznam v databázové tabulce. Žádné pole, které je součástí primárního klíče, nesmí obsahovat hodnotu NULL. Každá tabulka má mít definovaný právě jeden primární klíč.

Integritní omezení se mohou týkat jednotlivých hodnot vkládaných do polí databáze (například známka z předmětu musí být v rozsahu 1 až 5), či může jít o podmínku na kombinaci hodnot v některých polích jednoho záznamu (například datum narození nesmí být pozdější než datum úmrtí). Integritní omezení se může týkat i celé množiny záznamů daného typu – může jít o požadavek na unikátnost hodnot daného pole či kombinace polí v rámci celé množiny záznamů daného typu, které se v databázi vyskytují (například číslo průkazu v záznamech o osobách).

Velmi často používaným integritním omezením v relačních databázích je tzv. referenční integrita. Jedná se o požadavek, aby pro pole záznamu, jež má obsahovat odkaz na jiný záznam někde v databázi, takový odkazovaný záznam skutečně existoval.

Trigger v databázi definuje činnosti, které se mají provést v případě definované události nad databázovou tabulkou. Definovanou událostí může být například vložení nebo smazání dat.

Běžný rotační pevný disk?

Popište 3 základní složky časových nákladů potřebných ke čtení jednoho bloku

Zařaďte disk do hierarchie pamětí.

Čtení bloku:

1. Nastavení hlav nad správnou stopu.
2. Vyčkání, až bude správný sektor pod hlavou.
3. Přečtení a odeslání informace.
- 4-11 milisekund

Klasifikace paměťových médií:

Rychlost přístupu

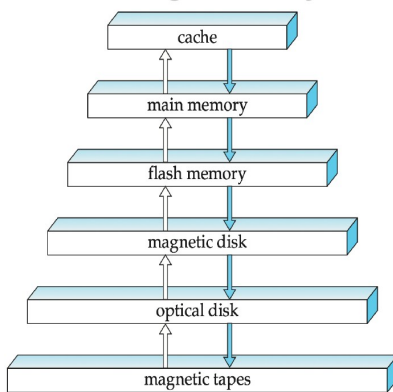
Cena za jednotku paměti, typicky jeden bit.

Spolehlivost: Ztráta dat při výpadku energie. Frekvence HW poruch.

Volatilní: Data jsou ztracena při odpojení napájení.

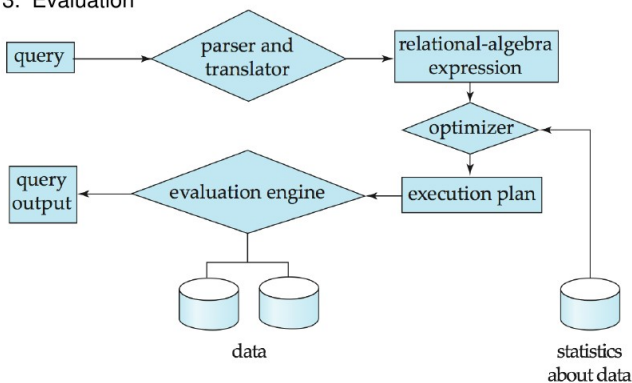
Nevolatilní: Data zůstávají i po odpojení napájení.

Hierarchie paměti:



Nakreslete schéma kroků zpracování dotazu a popište?

1. Parsing and translation
2. Optimization
3. Evaluation



1. rozbor a překlad (parsing and translation)
2. překlad dotazu do vnitřní formy a následně relační algebry, kontrola syntaxe a verifikace relací
3. optimalizace (optimization) zkoušení nalézt nejlevnější plán (posloupnost kroků) pro vyhodnocení dotazu (výrazy relační algebry lze totiž vyjádřit více ekvivalentními výrazy a každý výraz relační algebry může být různě složitý na provedení).
4. vyhodnocení (evaluation) převzetí a provedení plánu a vrácení odpovědi na dotaz.

Pro optimalizaci se využívají známa statistická data o databázi (počty a velikosti sloupců a hodnot). Hodnocení je založeno na odhadu doby, kterou bude asi vyhodnocení trvat. Ta může být ovlivněna i rychlostí procesoru, disků, ... počtem sektorů k přečtení a zapsání, počtem přístupů na disk, ...

**Popište, co je sekvenční soubor a uveďte příklad souboru s vyhledávacím klíčem UČO
Zkonstruuje nad souborem index pro jméno a uveďte typ indexu?**

Sekvenční soubor ukládá data uspořádaně za sebou s ohledem na hodnotu vyhledávacího klíče pro každý záznam. Používá se přetoková oblast pro vkládání záznamů, na které není místo. Musí se aktualizovat ukazatele a čas od času provést náročná reorganizace.
Sobor bude obsahovat řádky uspořádané podle UČO.

Hustý index – indexový záznam je uložený pouze pro každou hodnotu vyhledávacího klíče (ale stejné hodnoty klíče se v indexu neopakují)
(Každé různé jméno by mělo vlastní záznam v indexu a ukazatel na něj v indexované tabulce.)

Řídký index – indexové záznamy jsou pouze pro některé hodnoty vyhledávacího klíče.
Pro nalezení záznamu s vyhledávacím klíčem K musíme:
Nalézt indexový záznam s největším vyhledávacím klíčem menším než K.
Prohledat sekvenčně soubor od tohoto záznamu

Méně prostoru pro uložení indexu a méně udržujících operací při vkládání a mazání záznamu.
Obecně je ale při vyhledávání pomalejší než hustý index.
Vhodné řešení je řídký indexový soubor pro každý blok v souboru.
(Vybrat jenom některé položky souboru, tu mají záznam s ukazatelem v indexu.)

Jiné způsoby ukládání do souboru:

Halda: Sama se uspořádává.

Hash: Funkce počítá podle obsahu záznamu, kam má být uložen. Může to špatně dopadnout ve chvíli, kdy je mnoho záznamů mapovaných do stejného místa. Používají se přetokové oblasti, které se prohledávají sekvenčně. Hashovací funkce musí být dobře navrhnutá, aby rozmísťovala data rovnoměrně.

Co je pohled View v SQL a jakým způsobem je definován + příklad relace a pohledu

Uveďte dva způsoby používání pohledů?

V některých případech není žádoucí, aby uživatel viděl všechna data z databáze. Pohledy slouží k zakrytí některých dat pro různé uživatele.

create view POHLED as < SQL výraz >

Pohled uloží výraz, který se následně substituuje za název relace.

create view faculty as select ID, name, dept_name from instructor

Co je slotted page, jakým způsobem řeší fragmentaci volného místa?

Způsob ukládání záznamů v souboru.

Hlavička souboru obsahuje údaj o počtu záznamů, ukazatel na začátek volného místa a ukazatele na jednotlivé záznamy. Záznamy mohou být v souboru přeskupovány tak, aby mezi nimi nevznikalo nevyužitelné prázdné místo.

Záznamy mohou být proměnné délky.

Co je RAID, výhody použití, popis RAIDů?

Využití mnoha disků, které pracují paralelně, čímž se zvedá rychlost a kapacita. Redundance dat vytváří odolnost vůči poruchám.

Ukládají se informace, ze kterých je možné případná poškozená data zrekonstruovat.

Mirroring (shadowing):

Každý logický disk se skládá z dvou disků. Data jsou kládána paralelně stejná na oba disky. Pokud jeden selže, data zůstávají na druhém. Ke ztrátě dat by došlo pouze pokud by selhaly oba disky nebo druhý před nahrazením prvního.

Bit-level striping:

Každý bit každého bajtu ukládaných dat se uloží na jiný disk. V poli osmi disků, by každý bajt byl po bitech uložený na všech osmi discích. Rychlejší čtení než z jednoho disku, ale horší doba vyhledání záznamu.

Není dnes příliš používáno.

Block-level striping:

Každý blok ukládaného souboru se uloží na jiný disk. Dotazy na poskytnutí bloků mohou fungovat paralelně, pokud jsou bloky na různých discích.

RAID 0:

Používá block striping, nemá redundantní data. Využití v aplikacích vyžadujících rychlost, kde nezáleží na ztrátě dat.

RAID 1:

Disky zrcadleny, využívá block striping. Velká rychlost zápisu. Použití například pro ukládání logů.

RAID 2:

Bit striping kombinovaný s automatickou opravou chyb pomocí opravujícího se kódování.

RAID 3:

Využívá disky pro ukládání paritních bitů, které slouží v kombinaci s ostatními disky k rekonstrukci dat. (XOR s daty na ostatních discích).

RAID 4:

Využívá block striping a paritní disk na ukládání celých paritních bloků. Protože paritní blok se ukládá pro každý zápis, paritní disk je slabé místo, které může zpomalovat.

RAID 5:

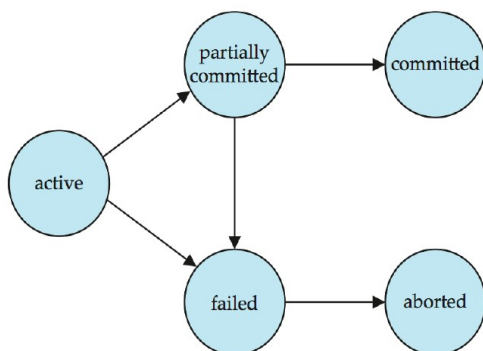
Paritní bloky se ukládají přímo na discích raidu na dalším disku v pořadí modulo počet disků + 1, ne na zvláštním disku mimo. Rychlejší než RAID 4 a odstraňuje problém s jedním vytíženým paritním diskem.

RAID 6:

Stejný princip jako RAID 5, ale navíc se ukládají redundantní data pro případ selhání.

RAID může být softwarový nebo hardwarový (složitější, více možností k poruše, používá nevolatilní RAM paměť)

**Co je plán zpracování transakcí a co musí splňovat, co je sériový plán (příklad)?
Nakreslit stavový diagram?**



Aktivní: Počáteční stav transakce

Částečně provedená: Po vykonání posledního příkazu transakce.

Provedená: Po úspěšném dokončení a zapsání.

Selhala: Při chybě provádění.

Zrušená: Ve chvíli, kdy transakce selhala a databáze se vrátila do původního stavu před provedením. Transakce může být znovu spuštěna nebo zahozena.

Transakce (transaction) je posloupnost jednotlivých příkazů, která má určitý význam vlastnosti zaručující integritu:

- **atomicita:** transakce musí být buď provedena celá nebo vůbec, nesmí dojít ke ztrátě dat při selhání systému v průběhu.
- **konzistence:** databáze musí být před prováděním transakce v konzistentním stavu a musí v něm být i po provedení. Například součty hodnot musí odpovídat před i po.
- **izolovanost:** každá transakce z více najednou prováděných transakcí nesmí vědět o jiných a její mezivýsledky musí být schovány před ostatními.
- **trvanlivost:** po úspěšném dokončení transakce, změny musí v databázi přetrvat (i v případě selhání systému).

T_1	T_2
read (A) $A := A - 50$ write (A) read (B) $B := B + 50$ write (B) commit	read (A) $temp := A * 0.1$ $A := A - temp$ write (A) read (B) $B := B + temp$ write (B) commit

Serialize je způsob provádění transakcí jednu po druhé, což udržuje databázi v konzistentním stavu, protože transakce vždy přistupují k databázi po všech předchozích změnách. Serializovatelnost současně běžících transakcí znamená, že by mohly být provedeny seriově a tudíž neporuší konzistenci databáze.

Porovnejte, popište od nejlepší k nejhorší hashování, B strom a B+ strom?

Hledání prvků s konkrétní hodnotou: H, B, B+

Hledání prvků z intervalu: B+, B, H

Vkládání prvku: H, B+, B

Princip algoritmu pro spojování relací pomocí hashování na příkladu?

Hashovací funkce rozdělí tabulky do odpovídajících si bloků, které mají stejnou hodnotu hashovací funkce na spojovacích attributech. Bloky je potom potřeba porovnat a spojit jen odpovídající si mezi sebou, nikoli všechny se všemi. Vhodné pro natural join.

1.) Co jsou složky ER-modelu?

a) Entitní množiny

b) Množiny vztahu

c) Atributy

a) Database může být postavena na modelu množiny entit nebo vztahu mezi entitami.

ENTITA - je hmotný nebo abstraktní objekt, který existuje v reálném světě a je jednoznačně odlišitelný od ostatních objektů (osoba nebo událost).

ENTITNÍ MNOŽINA - množina entit stejného typu, které sdílejí stejné vlastnosti (množina studentů MU)

b) VZTAH - je spojitost mezi několika entitami (vkladatel: zákazník-účet)

MNOŽINA VZTAHU - je matematická relace mezi minimálně dvěma entitami, z nichž každá prochází z určité množiny entit

- může mít také své atributy (vkladatel: datum přístupu)

STUPĚN MNOŽINY VZTAHU - je počet entitních množin, které se sdílejí na množině vztahu

- většina je binární, případně ternární

c) entita je představována množinou atributů, vlastností, jímž jsou obdareni všichni členové dané entitní množiny (student jméno a číslo)

DOMÉNA - je množina všech možných hodnot atributů

TPY - jednoduše a složené (adresa), jednohodnotové (ÚČTO), vícehodnotové {tel., te., ...}

nulové, odvozené (věk z data narození)

6.) Jakým způsobem mohou diskové pole urychlit zpracování přístupu k datům?

a) malých objektů

b) velkých objektů

a) Větší přístupová rychlost

b) Velká propustnost dat

9.) Charakterizujte stručně jazyk pro definici dat (DDL) a jazyk pro manipulaci s daty (DML)?

DDL - jazyk pro definici databázového schématu umožňuje specifikaci

- schématu relace

- množiny atributů

- doménu hodnot spojenou s každým atributem

- integritní omezení

- množinu indexů, které budou udržovány pro každou relaci,

- bezpečnostní opatření

a informace o oprávnění přístupu pro každou relaci, fyzickou strukturu ukládání na disk pro každou relaci.

DML - jazyk pro manipulaci s daty

Umožňuje pokládat dotazy do databáze a provádět modifikaci dat (insert, delete, update)

- NEPROCEDURALNI - uživatel specifikuje, jaká data požaduje

- PROCEDURALNI - uživatel specifikuje navíc i způsob, jak uskutečnit požadavek

10.) Definujte SELECT, PROJECT?

SELECT – odpovídá operaci projekce v relační algebře

- používá se k vypisování požadovaných atributů ve výsledku dotazu
- hvězdička(*, asterisk) v klauzuli značí všechny atributy
- SQL umožňuje vyskyt duplicit
- klíčové slovo DISTINCT uvedené bezprostředně za slovem SELECT způsobí odstranění duplicit
- klíčové slovo ALL uvedené hned za slovem SELECT způsobí ponechání duplicit
- může obsahovat aritmetické výrazy(+, -, *, /) a operace na konstantách nebo atributech n-tic

14.) Napiš 3 metody implementace JOIN?

Hash-join, merge-join, nested-loop join, block nested-loop join, indexed nested-loop join

16.) Definujte fyzickou a datovou nezávislost?

Fyzická úroveň - popisuje způsob uložení dat.

Logická úroveň - popisuje data uložená v databázi a vztahy mezi nimi.

Nezávislost dat - modifikace schématu nebo prostředí jedné úrovně nemá vliv na schéma vyšší úrovně.

- rozhraní mezi různými stupni a komponentami by měla být dobře definována, aby změny v některých částech neměly významný vliv na jiné části

Fyzická datová nezávislost: aplikační programy jsou nezávislé na fyzické datové struktuře.

Logická datová nezávislost: aplikační programy jsou nezávislé na změnách v logické struktuře databázového souboru.

18.) Co je outer join a jaké jsou jeho formy?

Rozšíření operace přirozené spojení, které zabraňuje ztrátě informací.

Spočítá operaci spojení a přidá n-tice z jedné relace, které neodpovídají n-ticím v druhé relaci k výsledkům operace spojení, řádky, které není s čím spojit jsou do výsledku zkopírovány a doplněny speciální hodnotou null, která vyjadřuje absenci hodnoty (hodnota je neznámá). Null hodnoty mají logickou hodnotu false.

20.) Kdy dojde k přetečení datové oblasti při hashování a jak tomu předjet?

Pokud hashovací funkce namapuje více záznamů do jednoho místa (více záznamů má stejný vyhledávací klíč, špatně navržená funkce). Dojde k přetečení bucketu a je vytvořen další zřetězený. Přetokové oblasti se potom prohledávají sekvenčně.

Použití rozšiřitelného hashování, které zvětšuje nebo zmenšuje slovník podle potřeby.

23.) Rozdíl mezi relačním schématem a relací?

Relace je množina uspořádaných n-tic a odpovídá tabulce.

Uspořádaná n-tice odpovídá řádku tabulky. Je prvkem relace.

Relační schéma je uspořádaná n-tice, nesmí být prázdné a odpovídá záhlaví tabulky.

25.) Rozdíl mezi kartézským součinem a natural-join?

Kartézský součin spojí všechny řádky první tabulky se všemi řádky druhé tabulky, vzniknou všechny kombinace.

Natural join spojuje jen řádky tabulky, které mají stejnou hodnotu jednoho atributu.

32.) Co je to víceúrovňový index a kdy se používá?

Primární index je seřazený podle vyhledávacího klíče.

Sekundární index je seřazený jinak než sekvencí uspořádání souboru.

Víceúrovňový index: Pokud se index nevejde do paměti, je rozdělen řídkým indexem, který specifikuje kde hledat. Pokud se ani řídký index nevejde do paměti, je přidána další úroveň. Všechny indexy je potřeba updatovat v případě změny dat.

33.) Co je to datový model a jaké jsou jeho základní formy?

Soubor nástrojů pro popis dat. Objektově logický model - OO model, funkcionální model...

Záznamově logické modely - relační, síťové, hierarchické