

## Zadanie 2 – vyhľadávanie a indexovanie

Odovzdanie do 21.10.2021 23:59 – máte na to 2 týždne – dostanete za to 7,5 boda.

Otázky 1-15 sú dokopy za 6 bodov (každá +/- rovnako keďže sa to nedá deliť). Zadanie 16 je za 1,5 boda. Zadania prosím neopisujte jednoslovnou ale zmysluplnou vetou (nie slohová práca - teda vecne), no je dôležité, aby sme z vašej odpovede pochopili či viete o čom píšete. Ak nebudeme mať tento pocit tak za otázku nedostanete ohodnotenie. Zároveň vždy priložte screenshot z výsledku a z explain analyse, ktorý podporí vaše vysvetlenie.

1. Vyhľadajte v authors username s presnou hodnotou 'mfa\_russia' a analyzujte daný select. Akú metódu vám vybral plánovač a prečo - odôvodnite prečo sa rozhodol tak ako sa rozhodol?
2. Koľko workerov pracovalo na danom selecte a na čo slúžia? Zdvihnite počet workerov a povedzte ako to ovplyvňuje čas. Je tam nejaký strop? Ak áno, prečo? Od čoho to závisí (napíšte a popíšte všetky parametre)?
3. Vytvorte btree index nad username a pozrite ako sa zmenil čas a porovnajte výstup oproti požiadavke bez indexu. Potrebuje plánovač v tejto požiadavke viac workerov? Čo ovplyvnilo zásadnú zmenu času?
4. Vyberte používateľov, ktorý majú followers\_count väčší, rovný ako 100 a zároveň menší, rovný 200. Potom zmeňte rozsah na väčší, rovný ako 100 a zároveň menší, rovný 120. Je tam rozdiel, ak áno prečo?
5. Vytvorte index nad 4 úlohou a v oboch podmienkach popíšte prácu s indexom. Čo je to Bitmap Index Scan a prečo je tam Bitmap Heap Scan? Prečo je tam recheck condition? Použil sa vždy index?
6. Vytvorte ďalšie 3 btree indexy na name, followers\_count, a description a insertnite si svojho používateľa (to je jedno aké dáta) do authors. Koľko to trvalo? Dropnite indexy a spravte to ešte raz. Prečo je tu rozdiel?
7. Vytvorte btree index nad conversations pre retweet\_count a pre content. Porovnajte ich dĺžku vytvárania. Prečo je tu taký rozdiel? Čím je ovplyvnená dĺžka vytvárania indexu a prečo?
8. Porovnajte indexy pre retweet\_count, content, followers\_count, name,... v čom sa líšia pre nasledovné parametre: počet root nódov, level stromu, a priemerná veľkosť itemu. Vysvetlite.
9. Vyhľadajte v conversations content meno „Gates“ na ľubovoľnom mieste a porovnajte výsledok po tom, ako content naindexujete pomocou btree. V čom je rozdiel a prečo?

10. Vyhľadajte tweet, ktorý začína "There are no excuses" a zároveň je obsah potenciálne senzitívny (possibly\_sensitive). Použil sa index? Prečo? Ako query zefektívniť?
11. Vytvorte nový btree index, tak aby ste pomocou neho vedeli vyhľadať tweet, ktorý končí reťazcom „https://t.co/pkFwLXZIEm“ kde nezáleží na tom ako to napíšete. Popíšte čo jednotlivé funkcie robia.
12. Nájdite conversations, ktoré majú reply\_count väčší ako 150, retweet\_count väčší rovný ako 5000 a výsledok zoradte podľa quote\_count. Následne spravte jednoduché indexy a popíšte ktoré má a ktoré nemá zmysel robiť a prečo. Popíšte a vysvetlite query plan, ktorý sa aplikuje v prípade použitia jednoduchých indexov.
13. Na predošlú query spravte zložený index a porovnajte výsledok s tým, keď je sú indexy separátne. Výsledok zdôvodnite. Popíšte použitý query plan. Aký je v nich rozdiel?
14. Napíšte dotaz tak, aby sa v obsahu konverzácie našlo slovo „Putin“ a zároveň spojenie „New World Order“, kde slová idú po sebe a zároveň obsah je senzitívny. Vyhľadávanie má byť indexe. Popíšte použitý query plan pre GiST aj pre GIN. Ktorý je efektívnejší?
15. Vytvorte vhodný index pre vyhľadávanie v links.url tak aby ste našli kampane z 'darujme.sk'. Ukážte dotaz a použitý query plan. Vysvetlite prečo sa použil tento index.
16. Vytvorte query pre slová "Володимир" a "Президент" pomocou FTS (tsvector a tsquery) v angličtine v stĺpcoch conversations.content, authors.decription a authors.username, kde slová sa môžu nachádzať v prvom, druhom alebo treťom stĺpci. Teda vyhovujúci záznam je ak aspoň jeden stĺpec má „match“. Výsledky zoradíte podľa retweet\_count zostupne. Pre túto query vytvorte vhodné indexy tak, aby sa nepoužil ani raz sekvenčný scan (správna query dobehne rádovo v milisekundách, max sekundách na super starých PC). Zdôvodnite čo je problém s OR podmienkou a prečo AND je v poriadku pri joine.