

DI, FCT-NOVA

June 29, 2020

Sistemas de Bases de Dados

Exam, 2019/20

Duration: 3 Hours

Group 1

Consider a database for managing an App for tracking proximity contacts of COVID-19 infected patients. This database includes the following tables (where the attributes that make the primary key are underlined in each table):

municipalities(codMun, name)

persons(idP, nameP, sex, age, codMun)

phones(numberPh, model, idP)

tests(idP, dateT, result)

tokens(tokenStr, numberPh, date)

contacts(tokenStr1, tokenStr2, when)

Each of these tables has a B+ tree clustered index over the primary key.

For each person, the database stores her id, name, sex age and the code of the municipality in which she lives. The names of the municipalities are stored in a separate table. Each person can make more than one test, but never more than one in a single day. Each test has a result which can be positive ('+') or negative ('-'). Persons have phones and each phone has a unique number (stored in the *phones* table). From time to time phones are given unique tokens (stored in the *tokens* table). When two phones are close together, they communicate their most recent tokens, and a tuple is stored in the *contacts* table with the most recent tokens of both phones and the date and time (in *when*) of when that occurred.

The database management system (DBMS) uses 4KB blocks and, for the exercises below, consider that the memory can have up to 100 blocks. Tuples of all these tables are of variable size and, on average, a tuple of the municipalities' table has 50 bytes, and a tuple of each of the other tables has 100 bytes.

Furthermore, at a given moment the *municipalities* table has 308 tuples, the *persons* table has 10.000.000 tuples, the *phones* table has 5.000.000 tuples, the *tests* 1.000.000 tuples, the *tokens* 20.000.000 tuples, and *contacts* table has 1.000.000.000 (one billion) tuples. Note that, as such, the contacts table has approximately 100GB.

Note: In this group, whenever an example is asked for, the example **must** be in terms of the database above. Moreover, **all** your answers must include a brief justification.

- 1 a) Show two execution plans for the following query (returning the phone numbers of persons who contacted with other persons who tested positive) and justify which one has the lowest cost.

```
select distinct C.numberPh
from tests natural inner join phones natural inner join tokens T
      contacts, tokens C
where result = '+' and (T.tokenStr = tokenStr1 or T.tokenStr = tokenStr2)
      and (C.tokenStr = tokenStr1 or C.tokenStr = tokenStr2)
```

- 1 b) By using this database over time, most of the tuples in the *contacts* table become irrelevant for warning those who have contacted persons with positive tests. In fact, for that purpose the only relevant contacts are the ones made in the 2 weeks before the positive test. Most DBMS have a mechanism of organisation of data that may significantly improve the efficiency in these situations. What is that mechanism, and how could it be used in this specific case (i.e. for getting the phones of those who, in the last 2 weeks, contacted persons who tested positive)?
- 1 c) The *tests* table has a primary key composed by the concatenation of 2 attributes: *idP* and *dateT*. As such, the clustering index can be created taking into account the concatenation (*idP*, *dateT*) or the concatenation (*dateT*, *idP*). Which one would you prefer for this particular database (given its expected usage), and why?

