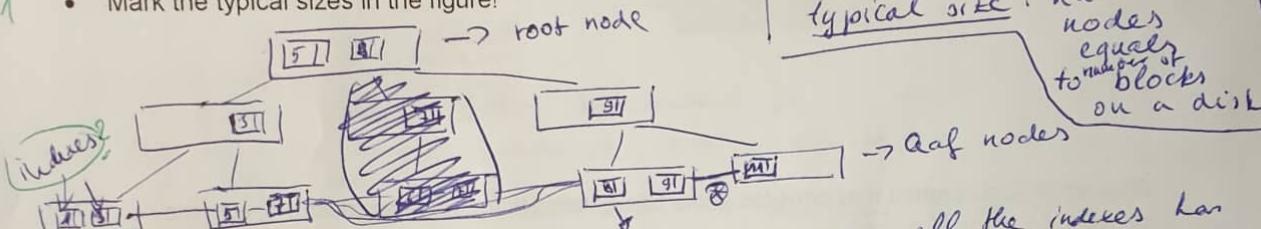


(10+5+5+10+3+3+3+6= 48points)

6
2.5

A. Indexes (10 points)

- Draw a figure showing the structure of a B+ index and the corresponding table!
- Mark the important components in the figure!
- Mark the typical sizes in the figure!



The leafnodes contains all the indexes. As we can see all the indexes has a pointer that shows us the next item (next index). It points towards the next leaf node too.

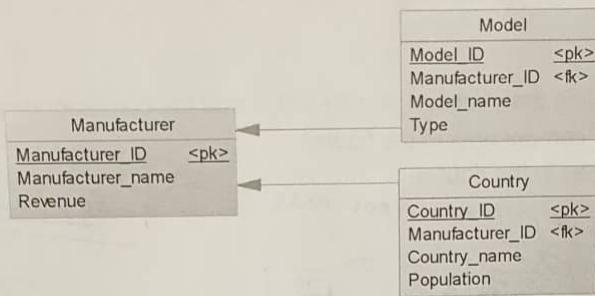
1 Which size of a B+ tree determines the cost of a query (let us consider a simple case, when the query is very selective, returns only a few records)?

The height of the B+ tree determines the cost of a query, that is how fast the tree can return an index for a query.

1.5 What can we say about index usage in case of a query with bad selectivity (high ratio of rows returned)? Please explain your answer!

In case of bad selectivity, index usage doesn't make the query cost less. However if the index is composed for the selected query it can return the records much faster.

- B. The following relational schema diagram is given, describing bike manufacturers, their bike models and the countries in which they have representatives.

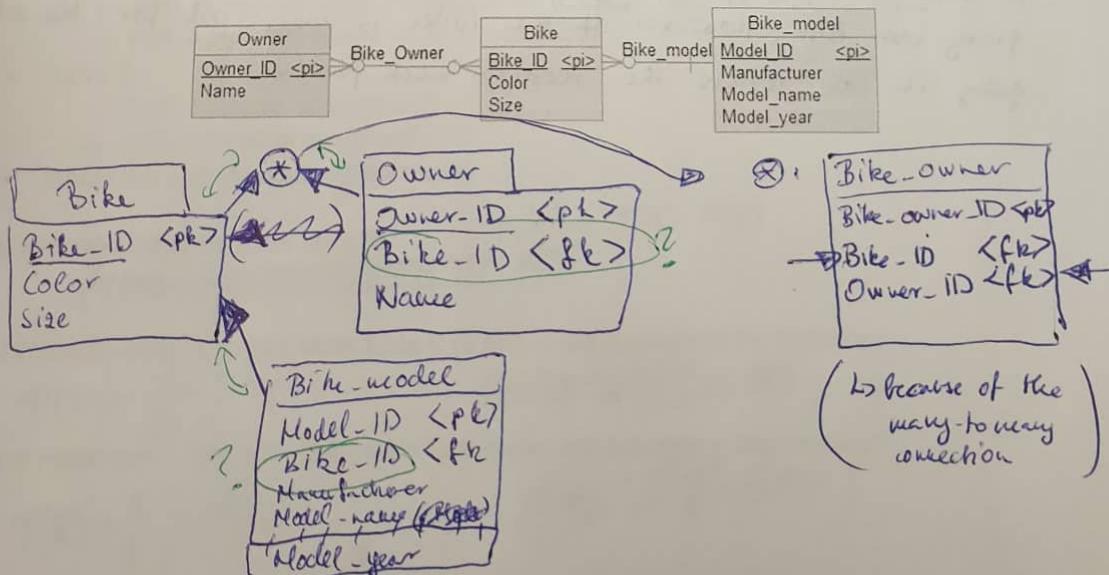


Write an SQL statement that provides some basic statistics for the manufacturers, with the following columns (5 points):

Manufacturer_name	Sum of population (sum of population of countries in which the manufacturer has a representative)	Number of models (number of models the manufacturer produces)

Joining the tables {
 The query {
 FROM (Manufacturer LEFT JOIN Model ON Manufacturer.ID = Model.Manufacturer_ID) AS (Table) T1
 LEFT JOIN Country ON T1.Manufacturer_ID = Country.Manufacturer_ID
 SELECT Manufacturer_name, SUM(Population), COUNT(Model_ID)

- C. Create a relational diagram for the following ER diagram. Mark primary and foreign keys clearly! Use the relational schema diagram notation of task B! (5 points)



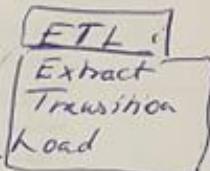
5.5

D. Explain major concepts of data warehousing! (10 points)

Motivation, definition: Data warehousing is (for) used for collect data from relational databases and from other good sources to examine and structure the data for enterprise ~~usage~~ usage, decision making, statistics etc.

Architecture, major elements:

- Data source: Relational DBs, Data仓库, etc. $\leftarrow E$
- ~~(Process)~~ Extract the data, clean it, and reorder it to a well-structured form $\leftarrow T$
- End-users: Who are making queries from the Warehouse database



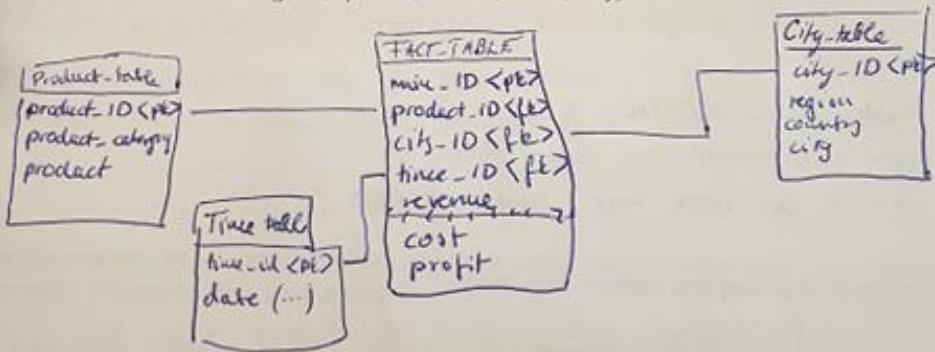
Challenges (some): and can even change it.

~~Eventual consistency~~ (all the infos are the same for all the users after time after every query and change.)

~~CAP theorem~~: Capacity, Availability, Partition tolerance \rightarrow only 2 can be guaranteed.

Draw a star schema for the following case! (Use the relational schema diagram notation in task ...!)

- 3
- Measures: revenue, cost, profit
 - Dimension 1: region, country, city
 - Dimension 2: product category, product
 - Dimension 3: date (year, quarter, month, week, day)



O

Possibilities to improve the performance/optimize a data warehouse (in a wide sense)?

- decentralize the system
- using indexes (for queries too)
- user friendly interface to communicate with the DW
- "Clean" the data from uninteresting information
- (Rebuild the db)
- optimise the structure after time, for example reorder or restructure after many changes.

E. What advantages/disadvantages has Hibernate over JDBC? (Mark advantages with "+" sign and the disadvantages with "-" sign!) (3 points)

(less) JDBC has (~~better~~) an Object Oriented background with inheritance

F. Compare object-relational mappers and objectrelational databases! What are the typical scenarios for their usage? (3 points)

Object-relational mappers are (~~structured for (read/write) querying~~) used for mapping huge data with many inconsistent data and info.
Objectrelational database is used for containing huge amount of data from many well-structured tables.

G. What advantages/disadvantages has a relational database over the Python Pandas package? (Mark the advantages with "+" sign and the disadvantages with "-" sign!) (3 points)

- Rel. Dat. doesn't have a toolkit for displaying the data and graphs like Pandas
- Performs better with queries (write with selectivity)
- Python Pandas can be used for machine learning, deep learning, and to make neural networks

H. Compare NoSQL databases (key-value/document stores, MongoDB), to relational ("SQL") databases! Mark the advantages of NoSQL systems with a "+" sign, the disadvantages with a "-" sign. (3 points)

- Faster data retrieval
- Better performance on bigger data
- with SQL we can make more complex queries

I. What is the cost of joining two tables? (6 points)

- Table R
 - $T(R)$: number of tuples (records) in R
 - $P(R)$ – number of pages in R
- Table S
 - $T(S)$ – number of tuples in S
 - $P(S)$ number of pages in S
- Ignore the cost of outputting the result!

Cost of nested loop join, outer table R:

$$2 \quad P(R) \cdot [P(S) \cdot [T(S) \cdot T(R)]]$$

$$P(R) + T(R) \cdot P(S)$$

Cost of block nested loop join, outer table R (B pages are available in the memory):

$$2 \quad P(R) + \frac{P(R)}{B+1} \cdot P(S) \quad (\cancel{P(R) + P(S) \cdot \cancel{B+1}})$$

Cost of block nested loop join, outer table S (B pages are available in the memory):

$$2 \quad P(S) + \frac{P(S)}{B-1} \cdot P(R) \quad (\cancel{P(S) + P(R) \cdot \cancel{B-1}})$$