



Magdeburg, February 25th, 2016

Written Exam "Distributed Data Management"

Name:	Firstname:
Immatriculation Number:	
Course of Studies: DKE, DigiEng, WIF, CV, INF, IngIf (underline your course)	
Number of additionally submitted sheets :	
Student signature:	Examinant signature:

Filled out by the examiners:

Part 1: Foundations	Part 2: Design	Part 3: Query Processing	Part 4: Transactions	Part 5: Single Choice, True False	Overall
(of 10)	(of 25)	(of 25)	(of 20)	(of 20)	(of 100)

Hints for the exam:

- Please check the completeness of the exam sheets (10 sheets, 20 pages).
- Fill out the first page with your personal data. Write your name in the according fields in the header of all of the following sheets.
- For the answers, use the space left between the questions. Use additional free pages at the end for further answers (do not forget to reference the number of the question). If additional paper is required, please contact the examiners. Do not use your paper.
- Do not use red or green pens. Do not use a pencil. Black and blue ballpoint or fountain pens preferred.
- Put the required things on the table, i.e. pens, student ID card, food/beverages.
- Turn off your mobile phone!
- Do not use any impermissible help (containing content of the lecture, allowing outside communication), especially books, slides, notes, printouts, mobile phones, etc.
- Using impermissible help will be considered an attempt of fraud and will cause the rejection from the exam. The exam will be graded as 5.0 ("not passed").
- Write clearly. Unreadable parts will not be considered.
- Answer only the posed questions. Additional text not referring to the questions will not be considered for extra points.

Good luck!

Name: _____

- d. Explain the term *Allocation Transparency*! (2 Points)
- e. Name the levels of the *5-Level-Schema-Architecture*! Where is information about fragmentation and allocation stored? (3 Points)

Name: _____

Part 2: Distributed Database Design (25 Points)

2. Fragmentation and Allocation (4 Points)

a. Explain the terms *Fragmentation* and *Allocation* and their relation! **(2Points)**

b. What are the necessary preconditions for a *Derived Horizontal Fragmentation*? **(2 Points)**

Name: _____

3. Fragmentation Methods (16 Points)

Consider an example relation

STUDENT (*MATRNO*, *FIRSTNAME*, *LASTNAME*, *COURSE*, *GRADE*)

containing 1.000 rows, where *MATRNO* is the primary key ranging from 1 to 1.000.

- a. Create 3 fragments of approximately equal size by *Horizontal Fragmentation*. Describe the three fragments by using the appropriate relational algebra operation! **(6 Points)**

- b. Name the *3 Rules for Fragmentation* and explain them based on the horizontal fragmentation created in 3.a! **(6 Points)**

Name: _____

- c. Which of the three rules is always violated for *Vertical Fragmentation*? Why? (2 Points)

- d. Given two fragments

*STUDENT*₁ (*MATRNO*, *FIRSTNAME*, *LASTNAME*)

*STUDENT*₂ (*MATRNO*, *GRADE*)

created by vertical fragmentation, which rule would be violated besides the one discussed in 3.c? Explain! (2 Points)

Name: _____

4. Allocation (4 Points)

In a network of three nodes the relations A, B and C should be stored on nodes N_1 , N_2 and N_3 . The relations have the following sizes:

- A is 10 Terabyte
- B is 15 Terabyte
- C is 5 Terabyte

The storage capacities of the each node are limited to 25 Terabyte

- a.** Give two example allocations of the three tables, one as *allocation without replication* and one as *allocation with replication*! **(2 Points)**

- b.** Name one possible advantage and one possible disadvantage of the *allocation of relations with replication*! **(2 Points)**

Name: _____

Part 3: Distributed and Parallel Query Processing (25 Points)**5. Phases of Query Processing (5 Points)**

Given an initial query plan $\sigma_{ID < 100}(R)$ where R is defined as $R(\underline{ID}, A, B, C)$ and, R is horizontally fragmented to R_1 and R_2 , such that R_1 contains all tuples with $ID \leq 3000$ and R_2 contains all tuples with $ID > 3000$.

a. How is the query plan changed by the step of *Data Localization*? (2 Points)

b. How can the resulting query be optimized by applying the rule to *Push Down Selections*? (1 Point)

c. How can the query processing furthermore be optimized by applying *Quantified Relations*? (2 Points)

Name: _____

6. Distributed Join Processing (10 Points)

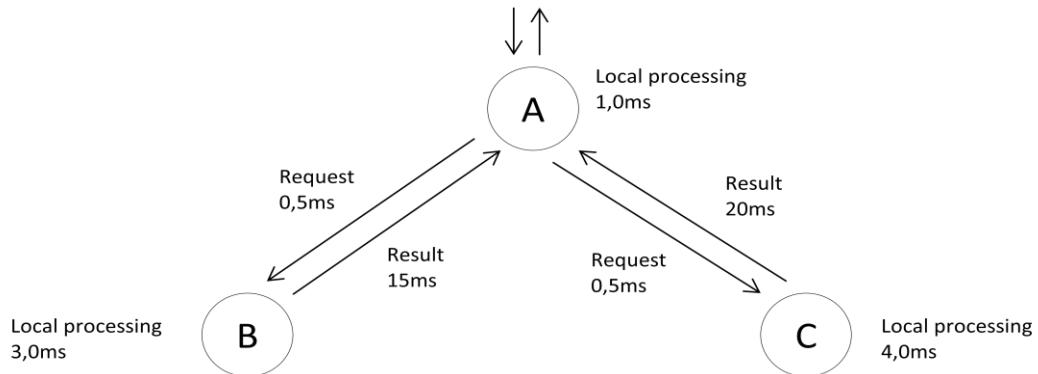
For two relations R and S stored on two different Nodes N_R and N_S , respectively, the following information is known: R has 5 attributes and 200 tuples, S has 3 attributes and 1000 tuples, and both are connected via an attribute defined as a primary key in R and as a foreign key in S . On average, there are 5 join partners for each tuple of R in S . To compute the join $R \bowtie S$, how many messages and how many attribute values need to be transferred for the following shipping strategies? **(2 Points each)**

- *Ship Whole* on node N_R
- *Ship Whole* on node N_S
- *Fetch As Needed* on node N_R
- *Fetch As Needed* on node N_S
- *Semi-Join* on node N_R

Name: _____

7. Optimization of Distributed Queries (10 Points)

As depicted in the figure below, a query is initiated at a node A. Parts of the query have to be independently computed on nodes B and C. The times required for local processing and network communication (request and result transfer) are outlined in the figure below.



- What is the *Total Time* of this execution? (2 Points)
- What is the *Response Time* of this execution? (2 Points)
- Why can the *Total Time* and the *Response Time* be different? (2 Points)

Name: _____

A relation containing data about master students has among other attributes the gender (female or male) and the course of studies (with six possible values: DKE, DigiEng, WIF, IF, IIF, CV).

- d. Assuming an equal distribution of values for all named attributes, what would be an estimation for the *Selectivity Factor* of a query searching for all female DKE master students? **(2 Points)**

- e. Why is the estimation of the *Selectivity Factor* of great importance for Distributed Database Systems? **(2 Points)**

Name: _____

- c. 5 nodes are participating in a transaction. How many messages are necessary to carry out the distributed commit in case of a successful transaction? Explain! **(2 Points)**

- d. Given the same scenario, how many messages have to be transferred, if 1 of the sub-transactions decides for an abort? Explain! **(2 Points)**

Name: _____

9. Transactional Replication (10 Points)

Five nodes in a distributed database hold a replica of a table. Read and write operations are performed on all nodes.

- a. Using the *Read One Write All*-approach (ROWA), how are read and write accesses synchronized? **(2 Points)**

- b. Name and explain one disadvantage of the ROWA-Approach! **(2 Points)**

Name: _____

c. Using the *Primary Copy*-approach, how are read and write accesses synchronized? **(2 Points)**

d. Compared to the ROWA-approach, name and explain one advantage and one disadvantage of the *Primary Copy*-approach! **(4 Points)**

Name: _____

Part 5: Single-Choice /True-False-Questions (20 Points)**1. Single Choice Questions**

For the following statements mark **exactly one** headword to continue the statement correctly. A correct mark scores one point for each statement. For an incorrect mark 1 point is subtracted. No mark is 0 points. Accordingly, only mark answers where you are very certain about the correctness. The overall score cannot be less than 0.

a. Multiple processors may access the same DBMS buffer (main memory) in the

- Shared-Nothing-Architecture
- Shared-Disk-Architecture
- Shared-Everything-Architecture

b. The fact that redundant copies of tables/fragments are visible to the user as only one table/fragment is referred to as

- Fault-Tolerance
- Replication Transparency
- One-Copy-Serializability

c. Horizontal fragmentation decomposes a relation

- by using the selection operation σ into sets of tuples
- by using the intersection operation \cap into overlapping subsets
- by using the projection operation π into sets of columns

d. The usage of several fragmentation methods on one relation is called

- Transactional Fragmentation
- Derived Fragmentation
- Hybrid Fragmentation

e. During query optimization, using a description of the content of a relation by a logical condition is referred to as

- Qualified Relations
- Relational Algebra
- Federated Databases

Name: _____

- f.** Given the relations $R(A, B, C)$, $S(C, D, E)$ and $T(E, F, G)$ an undesirable join order containing a Cartesian Product is
- $R \bowtie S \bowtie T$
 - $T \bowtie S \bowtie R$
 - $R \bowtie T \bowtie S$
- g.** The Selectivity Factor of an equivalence condition $A = v$ can be estimated as
- $SF(A = v) = \frac{v}{A_{max} - A_{min}}$
 - $SF(A = v) = \frac{1}{val(A)}$, where $val(A)$ is the number of distinct values of A
 - $SF(A = v) = \frac{v - A_{min}}{A_{max} - A_{min}}$
- h.** Full One-Copy-Serializability is granted for replicated data by
- the Primary Copy-approach
 - the Read One Write All-approach
 - the 5-Level-Schema-Architecture
- i.** The term “Snapshot Semantics” used in Primary Copy-replication refers to the fact, that
- disk images are used to check for possible inconsistencies
 - the content of multi-media data cannot be verified
 - secondary copies may be out of date with the primary copy
- j.** To compute a join result, the Ship Whole-strategy for distributed joins
- first transfers the join column of one relation
 - transfers a bit vector indicating the join partners
 - always transfers one relation in its entirety

Name: _____

2. True or false? (10 Points)

In the following table please mark whether the statement is true or false. A correct mark scores one point for each statement. For an incorrect mark 1 point is subtracted. No mark is 0 points. Accordingly, only mark answers where you are very certain about the correctness. The overall score cannot be less than 0.

	True	False
1. Federated Database Systems are a solution to integrate previously existing and possibly distributed databases.		
2. The reconstruction operation for fragments created by vertical fragmentation is the Semi-Join.		
3. An optimal horizontal fragmentation can be derived from the analysis of typical queries and their frequencies.		
4. A disadvantage of the Bit Vector Join is that tuples not required for a join are possibly transferred over the network.		
5. The Ship Whole-strategy is always preferable to the Fetch as Needed-strategy.		
6. The Fetch as Needed-strategy is always preferable to the Ship Whole-strategy.		
7. Deadlocks are problematic in a distributed database, because according wait relationships cannot always be detected locally.		
8. For Timestamp-based synchronization transactions are assigned a unique timestamp at the beginning.		
9. The 2-Phase-Commit Protocol grants Availability for distributed transactions.		
10. For Cost-based Optimization less possible plans have to be considered in Distributed Databases, e.g. because Replication leads to less redundant plans.		

Name: _____

Name: _____