```
1
    -------
2
    test 1 ausarbeitung
3
    ------
4
5
6
    ***ignoriere was am test steht - ist mit sicherheit grš§teils falsch***
7
8
9
    1) DATA WAREHOUSE CHARACTERISTICS (6 points)
10
    subject-oriented:
11
    - In the DW, data is not stored by operational applications, but by business subjects
12
    - Data is grouped around subjects
13
    - its structure is designed to make querying the data simple.
14
    integrated:
15
    - DW contains consolidated data from several (operational) applications/databases
16
    - Usually the data in the DW is not updated or deleted
17
    - Data inconsistencies are removed
18
    time-variant and nonvolatile:
19
    - The data stored in operational systems contains the current values
20
    - DW: When new current data becomes available, the "old" data is not overwritten
21
    - A data warehouse has to contain historical data to
22
        - Allow for analysis of the past
23
        - Relate information to the present
24
        - Enable forecasts for the future
25
26
27
    2) INFORMATION INTEGRATION APPROACHES (6 points)
28
    - Federation: Everybody talks directly to everyone else.
29
        - issue: n applications / data stores => up to n^2 connections
30
    - Warehouse: Sources are translated from their local schema to a global schema and
    copied to a central DB.
31
        - issue: usually only one-directional data flows supported
32
    - Mediator: Virtual warehouse - turns a user query into a sequence of source queries
    and assembles the results of this queries into an "aggregate" result.
33
        - issue: complex architecture, potentially slow, difficult to maintain
34
35
    3) DWH USE CASES (3 points)
36
37
       List three typical use cases for data warehousing
38
    - Manufacturing: Logistics management
39
    - Resource management
40
    - Finance: Risk management
41
    - Airlines: Route profitability
42
43
44
    4) DWH ARCHITECTURE REFERENCE MODEL (5 points)
45
    picture from left to right:
    Source DBs
46
47
    Landing Area
48
    Staging Area
49
    Data Warehouse
50
    Data Marts
51
52
53
    Source DB: DB of an application which supports types of business transactions.
54
    Landing Area: DB that is able to store a single data extract of a subset of one
    Source DB. Schema 1:1 with source DB scheme.
55
    Staging Area: DB that is able to store matching data extracts from various Landing
    Areas in an integrated format, waiting for the upload to the DW once data from all
    Landing Areas are available. Schema 1:1 with DW scheme.
56
    Data Warehouse: DB containing the history of all complete Staging Areas. Schema: 3NF.
57
    Data Mart: DB - on disk or in main memory - containing data describing the (present
    and past) performance of one or more types of business transactions, taken form the
    DW. Schema: denormalized star.
58
59
60
    5) OLTP VS OLAP (4 points)
61
    Х
62
    Х
63
    O (mit etwas bauchweh: fkeys, indexes -> redundant?)
64
    Х
65
66
```

```
67
      6) NORMALIZATION (3 points)
      A nonkey attribute must provide information about the key, the whole key, and
 68
      nothing but the key (so help me Codd).
 69
      (Anspielung auf 3NF)
 70
 71
 72
      7) FACTS VS DIMENSIONS (4 points)
 73
      Х
 74
      Х
 75
      0
 76
      Х
      (juhu)
 77
 78
 79
      8) SNOWFLAKE SCHEMA (4 points)
 80
      Х
 81
      Х
 82
      0
 83
      0
 84
 85
      9) OLAP OPERATIONS (4 points)
 86
      3
 87
      1
 88
      2
 89
      3
 90
 91
      10) R-OLAP VS M-OLAP VS H-OLAP (4 points)
 92
      Х
 93
      Х
 94
      Х
 95
      0
 96
 97
 98
      11) HORIZONTAL PARTITIONING (4 points)
 99
      Х
100
      0
101
      Х
102
      0
103
104
      12) VERTICAL PARTITIONING AND COLUMN STORES (4 points)
105
      0
106
      0
107
      Х
108
      Х
109
110
      13) BITMAP INDICES (13 points)
111
          a)
112
                       index rating avg 1
                                            index rating avg 2 index rating avg 3
                       index rating avg 4
113
                                                                                        0
               101
                                             0
                                                                  1
                       0
114
                                                                                        0
               102
                                             1
                                                                  0
                       0
115
                                                                                        0
               103
                       1
                                             0
                                                                  0
116
               104
                       0
                                             0
                                                                  0
                                                                                        1
117
          b)
                I) SELECT a.name
118
119
                       FROM article AS a
120
                       WHERE a.rating avg = 4;
121
               Option 2:
122
                   SELECT name
123
                       FROM article
124
                       WHERE index rating avg 4 = 1;
125
126
               II) SELECT count (article id)
127
                       FROM article AS a
128
                       WHERE a.lactoseFree = 'TRUE' AND ( a.rating avg = 3 OR a.rating avg
                       = 4 );
129
               Option 2:
130
                   SELECT count(article id)
131
                       FROM article
132
                       WHERE (index_rating_avg_4 = 1 OR index_rating_avg_3 = 1) AND
                       actoseFree = 1;
133
134
      14) ETL (4 points)
135
      Х
```

136 Х 137 Х 138 0 139 15) KIMBALLS DATA WAREHOUSE DEVELOPMENT LIFECYCLE (4 points) 140 141 142 143 Business Requirements 144 145 Technical Architecture Product Selection & 146 Installation 147 Innovate 148 Plan Dimensional modelling Physical Design Data Integration and (ETL) Grow 149 150 Reporting and Reporting and Analysics Design 151 and Analysics 152 Development 153 154 Project Management 155 156 157 158 16) AGILE BI (6 points) 159 Relevance: 160 There are lot of companies with very dynamic business intelligence applications. 161 This means that there are many on going changes in the business intelligence applications due to changing requirements form business. 162 These requirements are regarding reports, KPIs. 163 164 Classical waterfall model has issues with handling a very dynamic environment like mentioned above.. 165 There is a huge overhead when using a waterfall model which is not efficient for a environment where the requirements change nearly every day. 166 Nevertheless there are companies which do not face such a huge dynamic in their business intelligence applications, so in this case a waterfall model can be prefered. 167 168 Problems: Compliance or Regulation issues in areas like finance can stop you from using Agile BI 169 Risks: Keep the history after changes 170 171 172 17) SCHEMA-ON-WRITE VS SCHEMA-ON-RED (6 points) 173 Comparision: 174 Schema-on-write 175 (1) Schema must be created before any data can be loaded 176 (2) Explicit load operation transforms data to DB internal structure 177 (3) New columns must be added explicitly before new data for such columns can be loaded into DB 178 179 Schema-on-read 180 (1) Data is simply copied to the file store, no transformation 181 (2) Late binding: Serializer/Deserializer extracts required columns at read time 182 (3) New data can start flowing any time and will appear retroactively once the Ser/De is updated to parse it 183 184 Advantages: 185 Schema-on-write 186 + Read is fast 187 + Standards/Governance 188 189 Schema-on-read 190 + Load is fast 191 + Flexibility/Agility 192 193 18) HORIZONTAL VS. VERTICAL SCALING ISSUES (unfortunately unclear) 194 195 - Vertical (i.e., bigger machines): unfavorable economics, fundamental limits - Horizontal (i.e., partitioning): Data Sharding for RDBMS developed as an add-on, 196

inherently difficult 197 - Price vs. performance: 198 "The cost structures of data warehouses are anywhere from 10 to 100 times 199 higher than what they can drive per terabyte on a Hadoop cluster. 200 Now they can store multiple years of data, not just a month or two." 201 202 Storage (Kryder's law) grows much faster than processing power (Moore's law) and network speed (Nielsen's law). 203 204 Solution: Distributed data stores 205 BUT: ACID (Atomicity, Consistency, Isolation, Durability) difficult in a distributed environment 206 207 19) CAP THEOREM (unfortunately unclear) 208 - Consistency: all nodes see the same data at the same time - Availability: every request receives a response about whether it was successful or 209 failed 210 - Partition tolerance: system continues to operate despite arbitrary message loss or failure of part of the system 211 Impossible to achieve all three simultaneously! 212 213 20) MC Frage: HDFS is 214 X a specific file system that runs natively on a local machine ("Sits on top of the native filesystem") 215 0 conceptually based on Google's Big Table 216 X conceptually based on Google's GFS 217 0 stored in HDFS can be changed ("Files are "write once" - can be replaced, but not changed") 218 219 220 21) MC Frage: HBase 221 0 distributed relational database 222 X is based on a column-family oriented data model X is a key-value data store 223 224 0 uses typed columns 225 226 227 22 zusatzfrage) mehr daten sind immer besser als weniger daten - stimmen sie dem zu? (10)228 siehe reading "The Economist: A different game"