

Distributed Data Management – Exam

Winter Term 2019/2020

Matriculation Number: _____

0	1	2	3	4	5	6	7	8	9	Σ
1	7	8	15	6	5	3	7	6	12	70

Important Rules:

- The exam must be solved within **180 minutes** (14:00 – 17:00).
- Fill in your matriculation number (Matrikelnummer) above and **on every page**.
- Answers can be given in English or German.
- Any usage of external resources (such as scripts, prepared pages, electronic devices, or books) is not permitted.
- Make all **calculations transparent and reproducible!**
- Use the **free space under each task for your answers**. If you need more space, continue on the **back of the page**. Use the extra pages at the end of the exam only if necessary. **Provide a pointer to an extra page** if it should be considered for grading. The main purpose of the extra pages is for drafting.
- Please write clearly. **Do not use the color red or pencils.**
- If you have any questions, raise your hand.
- The exam consists of 33 pages including cover page and extra pages.
- For any multiple choice question, more or fewer than one answer might be correct.
- Good luck!

Task 0: Matriculation Number

Fill in your matriculation number (Matrikelnummer) on every page including the cover page and the extra pages (even if you don't use them).

Hint: Do it now!

1 point

Task 1: Distributed Systems

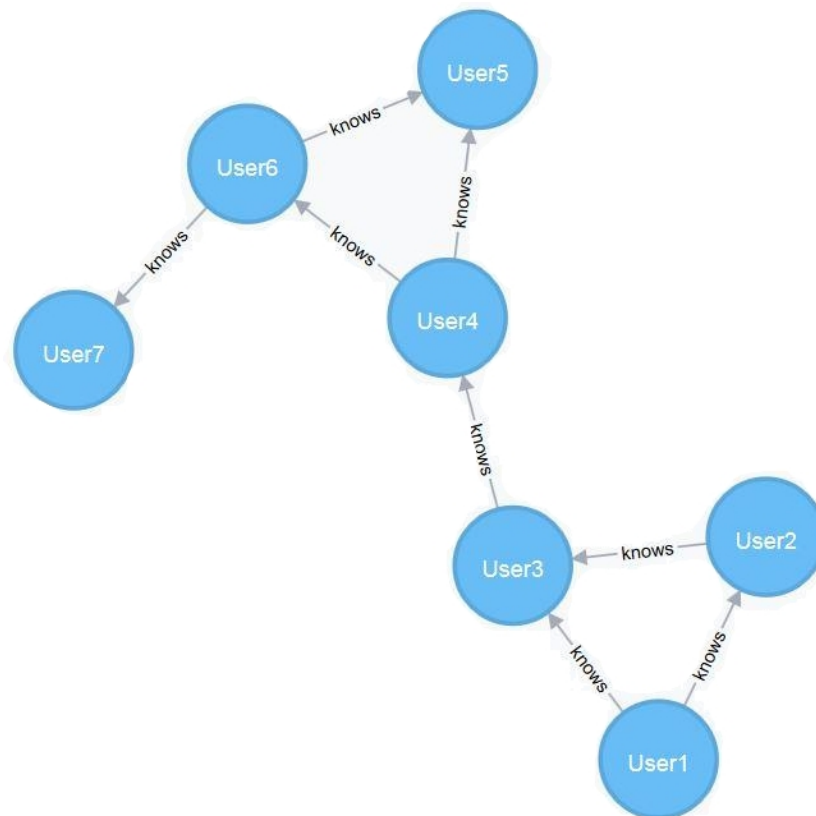
1. Which of the following statements about distributed systems are true? Tick all true statements. **4 points**

- A distributed system is a group of independent compute nodes that communicate and collaborate to solve a common task.
- A distributed system is reliable only if it prevents faults, errors and failures.
- A distributed system can be scaled vertically by adding more nodes into the cluster.
- A distributed system that uses task-parallelism cannot guarantee ACID.

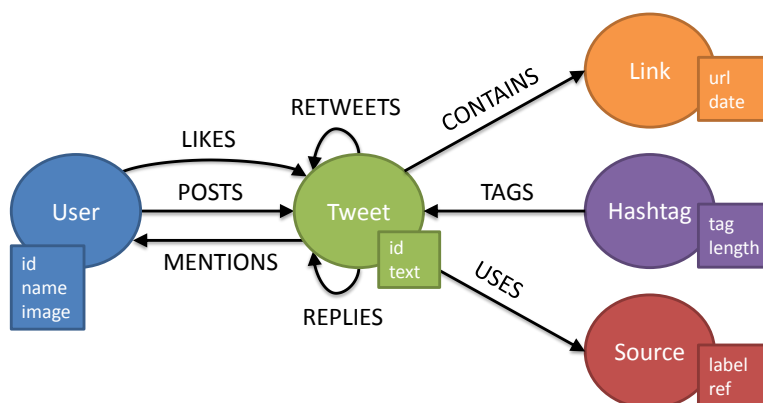
2. Consider the following situation: You are given two algorithms, *mineFast* and *mineDistributed*, that both solve the frequent itemset mining problem. The algorithm *mineFast* is highly efficient, but non-distributable; the algorithm *mineDistributed*, on the contrary, is not quite as efficient, but distributable. While experimenting with different datasets on one machine, you found that *mineFast* is on average twice as fast as *mineDistributed*. You also figured out that 60% of *mineDistributed*'s runtime is distributable and the algorithm scales linearly with the number of available machines. Given that you have five equal machines, which algorithm is the faster approach? What is the expected runtime x of *mineDistributed* for some input dataset on five machines given the runtime y of *mineFast* for that same input dataset on one machine? **3 points**

Task 2: Data Models and Query Languages

1. The graph below depicts a network of seven users. Some users know other users. Assume that a user can ask the users that she knows for their known users so that she then also knows these users – in that way, the **knows** edge becomes transitive. Iteratively calculating all transitive **knows** edges provides us with the transitive closure of the depicted graph that holds an explicit **knows** edge between a user and all other users that it can reach via direct or indirect **knows** edges. If we calculate the transitive closure with a *Bulk Synchronous Parallel* (BSP) transitive closure algorithm, how many steps would that algorithm need? Illustrate your answer or explain it in one or two sentences! **2 points**



2. The following graph describes the schema of a database storing twitter users and their tweets. Each circle describes a possible label for node instances and each edge their possible relationships to other nodes. Assume that the graph/schema is not completely shown here and further nodes and edges might exist (open world assumption). The database was created in Neo4J and you have to use Cypher to query its content.



Write one Cypher query that adds a node with two edges to the depicted graph. The node we want to add has the new label **Website** and an attribute *content* with value “my_CV.pdf”. There should also be a new edge labeled **REFERENCES** that points from the **Link** node with *url* “www.usr42.com” to the new **Website** node and a new edge labeled **DESCRIBES** that points from the new **Website** node to the **User** node with *id* “42”. **3 points**

3. Still consider the graph of twitter users and tweets depicted above. We are now interested in tweets with the **Hashtag** *tag* “#ClimateChange” that have been retweeted at least 1000 times. Write a Cypher query that returns all these popular climate change **Tweet** nodes. **3 points**

Task 3: Actor Programming

1. Which of the following statements on the actor model and actor programming are *true*? Tick the true ones. **6 points**

- An *actor* is a special type of a thread that has its own private state, a mailbox and behavior.
- Actor *messages* need to be serializable, because the actor system needs to serialize and de-serialize every message that is sent from one actor to another.
- The specification of an *actor system* in Akka guarantees only at-most-once message delivery, because Akka uses fire-and-forget messaging and when firing messages with the UDP protocol, these messages might get lost.
- The *actor model* is very flexible: It allows the implementation of both task- and data-parallel applications, it supports pull- and push-based communication protocols and it lets us dynamically spawn and terminate actors at runtime.
- The *reaper pattern* defines a special actor, called the reaper, who sends out PoisonPill messages that terminate all actors in the actor system in a clean way.
- Using a *side channel* for large messages is advisable, because large messages may otherwise block important system messages, such as cluster heartbeats.

2. The Actor model implementation in Akka ensures that every user defined actor is supervised by some other actor. If an actor encounters an error or crashes, its supervisor is notified so that the supervisor can analyze and handle the error. The supervisor then has four different options according to the actors' lifecycles to proceed the program. What are these four options? **2 points**

3. The following two actors should implement the *Network Time Protocol* (NTP) algorithm. Fill out the gaps in the code such that the **Client** actor frequently synchronizes its time according to the **Server** actor's reference time. **7 points**

Hint: Use the space on this page for drafting. Then write your solution into the gaps of the already provided template code. Use the back of the template pages if you need more space.

```
package de.hpi.ddm.ntp;

import java.io.Serializable;
import java.util.concurrent.TimeUnit;
import akka.actor.AbstractLoggingActor;
import akka.actor.ActorRef;
import akka.actor.Cancellable;
import akka.actor.Props;
import scala.concurrent.duration.Duration;

public class Client extends AbstractLoggingActor {

    public static Props props(final ActorRef server) {
        return Props.create(Client.class, () -> new Client(server));
    }

    public Client(final ActorRef server) {
        this.timeSyncProcess = this.getContext().system().scheduler().schedule(
            Duration.create(0, TimeUnit.SECONDS),
            Duration.create(3, TimeUnit.SECONDS),
            server,
            -----
            |
            |
            -----
            this.getContext().dispatcher(), null);
    }

    @Override
    public void postStop() throws Exception {
        this.timeSyncProcess.cancel();
    }

    public static class TimeSyncResponse implements Serializable {
        private static final long serialVersionUID = 1208000708229308005L;
        -----
        |
        |
        -----
    }

    private final Cancellable timeSyncProcess;
    private long offset = 0;
    private float offsetAdjustmentFactor = 0.3f;

    private long getCurrentTime() {
        return System.currentTimeMillis() + this.offset;
    }

    @Override
    public Receive createReceive() {
        return receiveBuilder()
            .match(TimeSyncResponse.class, this::handle)
            .matchAny(object -> this.log().info("Unknown message: \"{}\"", object.toString()))
            .build();
    }

    private void handle(TimeSyncResponse message) {
        -----
        |
        |
        -----
    }
}
}
```

```
package de.hpi.ddm.ntp;

import java.io.Serializable;
import akka.actor.AbstractLoggingActor;
import akka.actor.Props;

public class Server extends AbstractLoggingActor {

    public static Props props() {
        return Props.create(Server.class);
    }

    public static class TimeSyncRequest implements Serializable {
        private static final long serialVersionUID = -3057724412864626584L;
        -----
    }

    private long getCurrentTime() {
        return System.currentTimeMillis();
    }

    @Override
    public Receive createReceive() {
        return receiveBuilder()
            .match(TimeSyncRequest.class, this::handle)
            .matchAny(object -> this.log().info("Unknown message: \"{}\"", object.toString()))
            .build();
    }

    protected void handle(TimeSyncRequest message) {
        -----
    }
}
```


Task 4: Replication and Partitioning

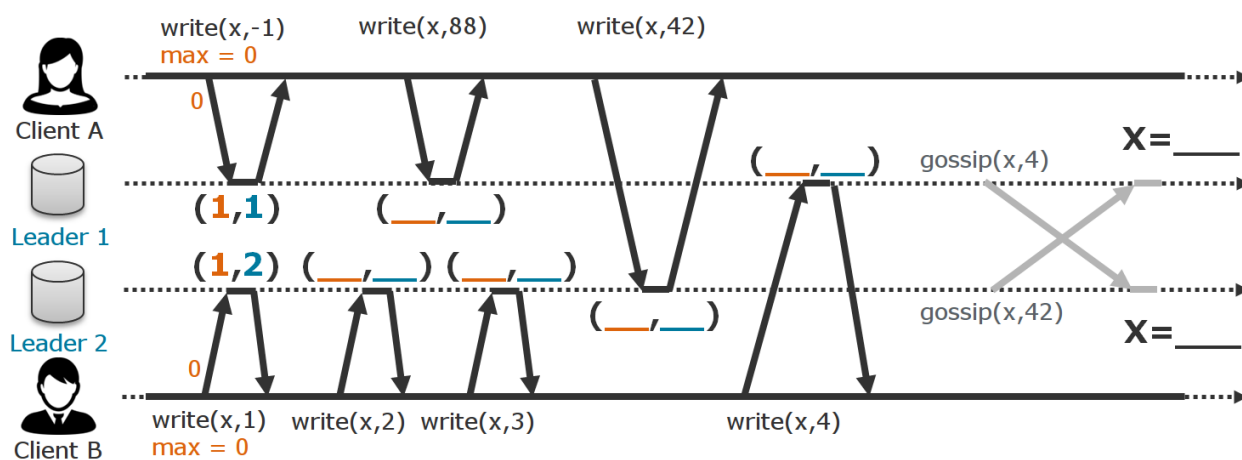
1. Assume you have a cluster of 12 801 nodes. The cluster runs a leaderless replicated, distributed database. Quorum reads and writes are used to ensure consistency and a gossip protocol ensures that all updates will eventually spread to all nodes in the cluster. For this task, assume also that the gossip protocol is perfect, i.e., newly written data is always gossiped to nodes that do not already know it.

Which read (r) and write (w) values do we need to define in our quorum $q(r,w)$ if writes should return as fast as possible and every successfully written value should reach all cluster nodes in not more than seven rounds of gossip? **4 points**

2. What is the disadvantage of the modulo operation when used for partitioning key spaces? What other approach prevents this issue? **2 points**

Task 5: Consistency and Transactions

1. Given a database management system that implements *Lamport Timestamps* for causal write ordering. A lamport timestamp is a pair (c, i) with a write counter c and a node identifier i . Every write operation is associated with such a timestamp. We can use these timestamps for write ordering, because lamport timestamps are comparable: $(c, i) > (c', i')$ iff $(c > c') \vee (c = c' \wedge i > i')$. Add the lamport timestamps in the following example. What is the final value of the field x on leader 1 and on leader 2? **3 points**

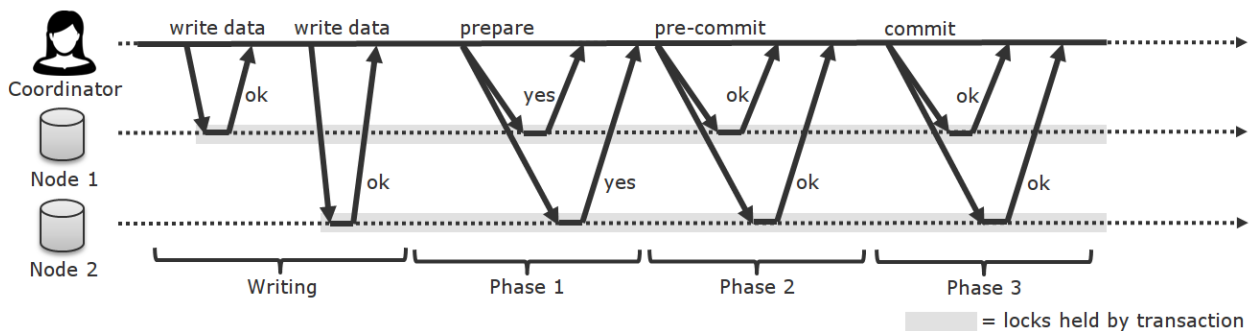


2. The following figure depicts the *Three-Phase Commit* (3PC) protocol. If the coordinator in that protocol dies, another coordinator can take over and finish the transaction by either rolling it back or consistently committing it. Draw two vertical lines into the 3PC process:

I.) The first line should be at exactly that place up until which a new coordinator would definitely roll the transaction back, i.e., coordinator crashes left to that line would cause the transaction to be rolled back.

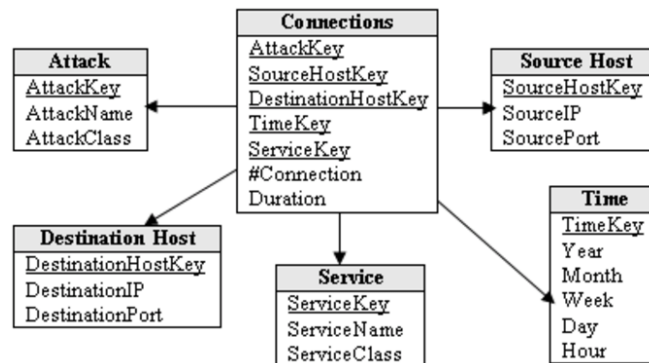
II.) The second line should be at exactly that place from which onwards a new coordinator would definitely commit the transaction, i.e., coordinator crashes right to that line would cause the transaction to be committed.

Coordinator crashed in between these two vertical lines could lead to both rollback and commit situations depending on what information has been lost. **2 points**



Task 6: Data Warehousing

1. Consider the following data warehouse *star schema* and the join query on that schema. Rewrite the join query as a star join and explain why the star join makes sense especially in distributed data warehouses. **2 points**



$((Connections \bowtie Attack) \bowtie Service) \bowtie Time$

2. What kind of schema mapping strategy is expressed by the following schemata and views? Write down the name of the strategy. **1 points**

Global schema
 Students(matrikel, first, last)
 Registrations(matrikel, course)

Local schemata
 ITSE_students(matrikel, first, last)
 DE_students(matrikel, first, last)
 DH_students(matrikel, first, last)
 DS_students(matrikel, last, email)
 HPI_registration(matrikel, course, lp)
 PULS_registration(matrikel, course, accepted)

Views
 V1: (SELECT * FROM ITSE_students, DE_students, DH_students)
 UNION
 (SELECT matrikel, null, last FROM DS_students);
 V2: (SELECT matrikel, course FROM HPI_registration)
 UNION
 (SELECT matrikel, course FROM PULS_registration WHERE accepted=true);

Task 7: Distributed Query Optimization

1. Which of the following statements about distributed query optimization are true? Tick all true statements. **4 points**

- The query optimizer should, if possible, always try to push selection, projection and grouping operations to the nodes that hold the relevant data.
- The query optimizer does not need to push set operations, such as union, intersect and except, to data nodes, because evaluating set operations on data nodes does not reduce the network traffic.
- To minimize the network traffic for join operations, the query optimizer has to run a full reducer semi-join program.
- Bloom filters can be used to approximate semi-joins: They might not remove all non-joining tuples, but all removed tuples are true non-joining tuples.

2. The following picture shows four relations that should be joined. Write down a *reducer program* of semi-joins that reduces R1. Then, write down the result of a full reducer program, i.e., the *reduced relations* R1, R2, R3, and R4. **3 points**



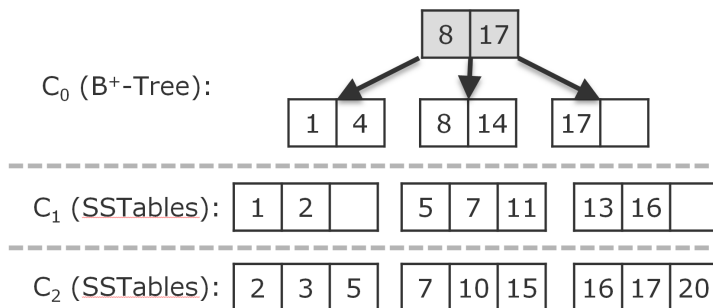
Task 8: SSTables and LSM Trees

The following figures depict LSM tree instances with three levels, which are C_0 , C_1 , and C_2 . In the LSM tree, we show only the keys and not their values. Assume that all values have the same size, so that a B^+ -tree leaf can hold exactly up to two key-value pairs and each SSTable can hold exactly up to three key-value pairs. The maximum depth of the B^+ -tree shall be two, which means that the depicted tree cannot grow any further. We now insert new elements into the LSM tree. Use the free space next to the figures for drafting and write your final results into the LSM tree templates below each instance figure.

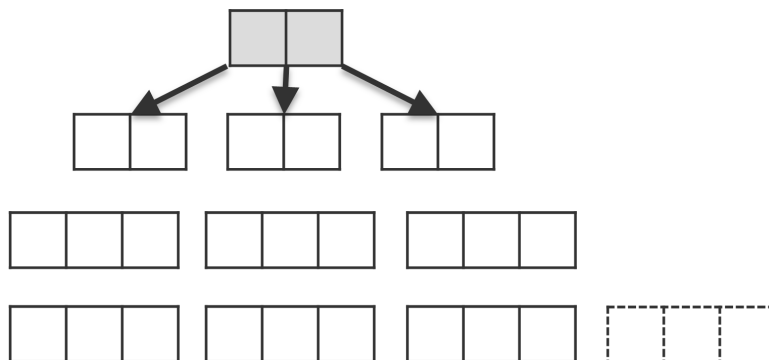
1. Insert the key **18** into the depicted LSM tree instance:

1 points

LSM tree instance:

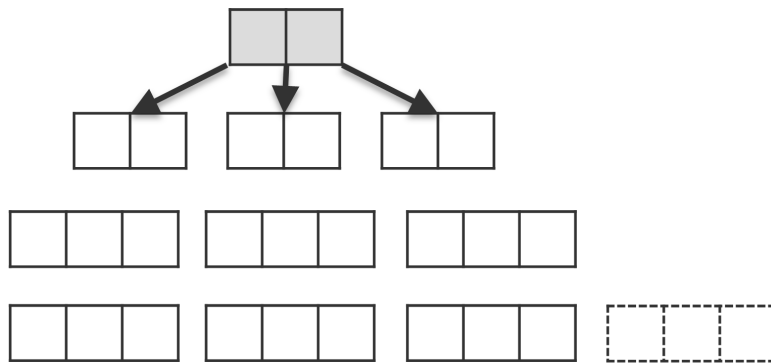
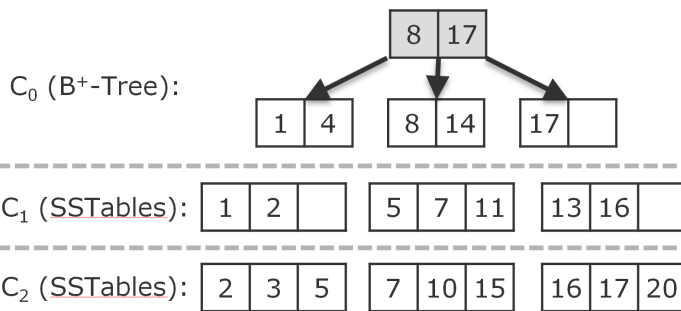


LSM tree template:



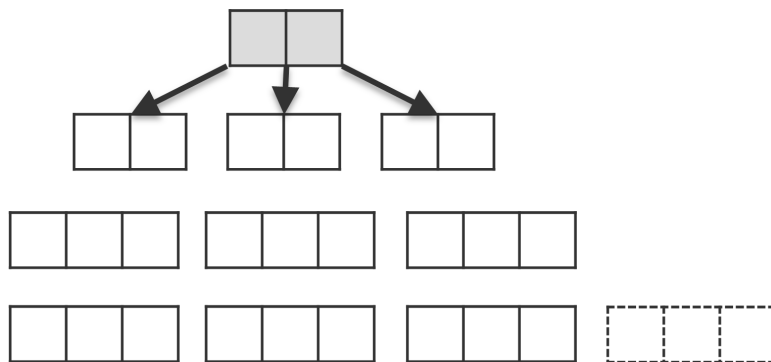
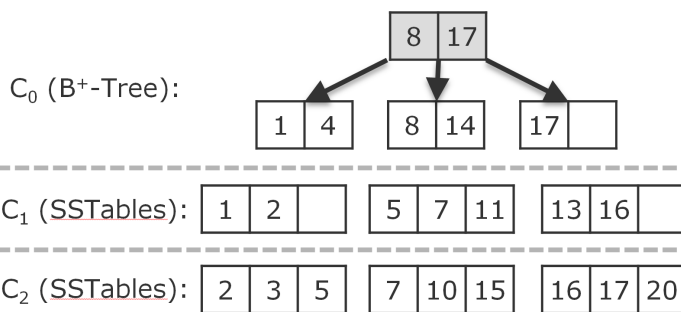
2. Insert the key **3** into the depicted LSM tree instance:

2 points



3. Insert the key **10** into the depicted LSM tree instance:

3 points



Task 9: Batch Processing

Suppose you are given three datasets: A `students` dataset that contains general information about students, a `courses` dataset that describes available courses for the students, and an `enrollment` dataset, which is basically a join table between `students` and their `courses`. The following code snippets read the three datasets into Spark Datasets:

```
val students = spark
  .read
  .option("quote", "\"")
  .option("delimiter", ",")
  .csv(s"data/students.csv")
  .toDF("ID", "Name", "Semester", "Supervisor")
  .as[(String, String, String, String)]
```

```
val enrollments = spark
  .read
  .option("quote", "\"")
  .option("delimiter", ",")
  .csv(s"data/enrollments.csv")
  .toDF("StudentID", "CourseID", "Credits")
  .as[(String, String, String)]
```

```
val courses = spark
  .read
  .option("quote", "\"")
  .option("delimiter", ",")
  .csv(s"data/courses.csv")
  .toDF("ID", "Title", "Teacher", "Topic")
  .as[(String, String, String, String)]
```

Use the three Datasets to solve the following tasks. You may use Spark's *Dataset* and/or *DataFrame* API but no SQL! Also have a look at the *Dataset* API documentation at the end of this task. If you are not sure about how a particular interface, call, or class works, make a good guess and provide a comment on how you *think* it works.

1. Write a Spark transformation pipeline that starts with the `enrollments` Dataset, then removes all those enrollments that grand less than 3 credits, then multiplies the credits by $2/3$ in order to transform them into *SWSs*, and finally displays the results in tabular form with schema (`StudentID`, `CourseID`, `SWS`) on the standard output. **4 points**

2. Translate the following SQL query into a Spark transformation pipeline. **4 points**

```
SELECT Semester, COUNT(ID) AS Students
FROM students
WHERE Semester <= 10
GROUP BY Semester;
```

- Rank all teachers by the number of enrollments that they got for all their courses. You can assume that teachers have unique names in `courses`. `Teacher` and that the `enrollments` file stores all enrollments for all courses ever given. Teacher that never gave a course do not exist in the files and will not appear in the ranking. Report a list of `(Teacher, SumEnrollments)` that is sorted by `SumEnrollments`.

4 points

Typed transformations

- ▶ `def as(alias: String): Dataset[T]`
Returns a new Dataset with an alias set.

- ▶ `def distinct(): Dataset[T]`
Returns a new Dataset that contains only the unique rows from this Dataset.

- ▶ `def except(other: Dataset[T]): Dataset[T]`
Returns a new Dataset containing rows in this Dataset but not in another Dataset.

- ▶ `def filter(func: FilterFunction[T]): Dataset[T]`
(Java-specific) Returns a new Dataset that only contains elements where `func` returns `true`.

- ▶ `def filter(func: (T) => Boolean): Dataset[T]`
(Scala-specific) Returns a new Dataset that only contains elements where `func` returns `true`.

- ▶ `def flatMap[U](f: FlatMapFunction[T, U], encoder: Encoder[U]): Dataset[U]`
(Java-specific) Returns a new Dataset by first applying a function to all elements of this Dataset, and then flattening the results.

- ▶ `def flatMap[U](func: (T) => TraversableOnce[U])(implicit arg0: Encoder[U]): Dataset[U]`
(Scala-specific) Returns a new Dataset by first applying a function to all elements of this Dataset, and then flattening the results.

- ▶ `def groupByKey[K](func: MapFunction[T, K], encoder: Encoder[K]): KeyValueGroupedDataset[K, T]`
(Java-specific) Returns a `KeyValueGroupedDataset` where the data is grouped by the given key `func`.

- ▶ `def groupByKey[K](func: (T) => K)(implicit arg0: Encoder[K]): KeyValueGroupedDataset[K, T]`
(Scala-specific) Returns a `KeyValueGroupedDataset` where the data is grouped by the given key `func`.

- ▶ `def intersect(other: Dataset[T]): Dataset[T]`
Returns a new Dataset containing rows only in both this Dataset and another Dataset.

- ▶ `def joinWith[U](other: Dataset[U], condition: Column): Dataset[(T, U)]`
Using inner equi-join to join this Dataset returning a `Tuple2` for each pair where `condition` evaluates to `true`.

- ▶ `def map[U](func: MapFunction[T, U], encoder: Encoder[U]): Dataset[U]`
(Java-specific) Returns a new Dataset that contains the result of applying `func` to each element.

- ▶ `def map[U](func: (T) => U)(implicit arg0: Encoder[U]): Dataset[U]`
(Scala-specific) Returns a new Dataset that contains the result of applying `func` to each element.

- ▶ `def sort(sortExprs: Column*): Dataset[T]`
Returns a new Dataset sorted by the given expressions.

- ▶ `def sort(sortCol: String, sortCols: String*): Dataset[T]`
Returns a new Dataset sorted by the specified column, all in ascending order.

- ▶ `def union(other: Dataset[T]): Dataset[T]`
Returns a new Dataset containing union of rows in this Dataset and another Dataset.

Untyped transformations

- ▶ `def col(colName: String): Column`
Selects column based on the column name and return it as a `Column`.

Actions

- ▶ `def collect(): Array[T]`
Returns an array that contains all rows in this Dataset.

- ▶ `def show(numRows: Int, truncate: Int): Unit`
Displays the Dataset in a tabular form.

- ▶ `def foreach(f: (T) => Unit): Unit`
Applies a function `f` to all rows.

- ▶ `def reduce(func: ReduceFunction[T]): T`
(Java-specific) Reduces the elements of this Dataset using the specified binary function.

- ▶ `def reduce(func: (T, T) => T): T`
(Scala-specific) Reduces the elements of this Dataset using the specified binary function.

- ▶ `class KeyValueGroupedDataset[K, V] extends Serializable`

- ▶ `def mapGroups[U](f: MapGroupsFunction[K, V, U], encoder: Encoder[U]): Dataset[U]`
(Java-specific) Applies the given function to each group of data.

- ▶ `def mapGroups[U](f: (K, Iterator[V]) => U)(implicit arg0: Encoder[U]): Dataset[U]`
(Scala-specific) Applies the given function to each group of data.

Extra page 1

Extra page 2

Extra page 3