
Project Portfolio 1

RAWDATA 2017

Project Report

Group Name/ Group number 8

Roskilde University

Computer Science Department RAWDATA 2017



Computer Science Department
Roskilde University
<http://www.ruc.dk>

Title:

Project Portfolio 1

Project Period:

Fall Semester 2017

Project Group:

8

Participant(s):

Eugenio Maria Capuani
Kidane Mahari Tesfai
Kristoffer Carl Schou
Rune Bastian Barrett

Page Numbers: 17

Date of Completion:

October 9, 2017

Abstract:

In this report we go over the design of a database to support a Stack Overflow viewing application (for later implementation). We will also go through the functionality of that application and the normalization of the database.

Contents

1	Design	1
1.1	The Post table	1
1.2	Schema structure	2
1.3	Design	5
2	Functionality	7
2.1	Creation of the database	7
2.1.1	Splitting up non atomic tags	7
2.2	API for querying the database	8
2.2.1	Searching	9
2.2.2	History, Marking and Listing Data	10
3	Normalization of the database	12
3.1	The data sample	12
3.2	Normal forms	13
3.2.1	First normal form	13
3.2.2	Second normal form	13
3.2.3	Third normal form	13
3.2.4	Boyce-Codd normal form	14
3.3	Entity-relationship model	15
	Bibliography	16
A	Initial design document	17
B	stack_overflow_normalized.sql	19

Chapter 1

Design

The initial schema contained two tables, namely *Posts* 1.1 and *Comments*. It was clear that the tables were badly designed, and as a result in order to avoid redundancy, this schema had to be normalized. This first section contains a break-down of our normalized database schema, followed by some design considerations.

1.1 The Post table

As it can be clearly seen from the figure, the post table 1.1 contains information about the the post and the user: this is not an appropriate way of storing data in the Database, as it lead to unnecessary data redundancy. As a result, the table has been decomposed in to the following two smaller tables: The *Post* table, which stores information only about the post, and the *User* tables stores information only about the user. In addition to that the two tables are linked using a foreign key, namely *owner_user_id*. The *Comments* table has been implemented in a similar way to how it looks in the original schema, but it now references our super-table *Post* through *post_id* and *user_id*.

post_id	creation_date	score	body	title	owner_user_id	type_id
19	2008-08-01 05:21:22	164	Solut...	13	1	
71	2008-08-01 13:38:00	43	<p>Her...		49	2
531	2008-08-02 18:22:53	110	<p>The ...		157	2
709	2008-08-03 14:53:54	28	<p>Il... .NET Testi		111	1
712	2008-08-03 14:57:45	14	<p>Scrw.		91	2
713	2008-08-03 14:59:21	32	<p>Int...		34	2
718	2008-08-03 15:07:20	5	<p><n...		27	2
798	2008-08-03 19:36:48	3	<code...		230	2
1053	2008-08-04 06:21:38	7	<part...	13	1	
1054	2008-08-04 06:22:19	5	<Sche...		13	2

Figure 1.1: The Posts table

user_id	user_name	user_creation_date	user_location	user_age
1	Jeff Atwood	2008-07-31 14:22:31	El Cerrito, CA	45
3	Jarrod Dixon	2008-07-31 14:22:31	New York, NY	36
4	Joel Spolsky	2008-07-31 14:22:31	New York, NY	
5	Jon Galloway	2008-07-31 14:22:31	San Diego, CA	45
13	Chris Jester-Young	2008-08-01 04:18:05	Raleigh, NC	35
17	Nick Berardi	2008-08-01 12:02:22	Collegeville, PA	35
20	Tom	2008-08-01 12:09:11	Atlanta, GA	
25	CodingWithoutComments	2008-08-01 12:15:23	Seattle, WA	35
26	Shawn	2008-08-01 12:18:15	San Francisco, CA	31

Figure 1.2: The User table

1.2 Schema structure

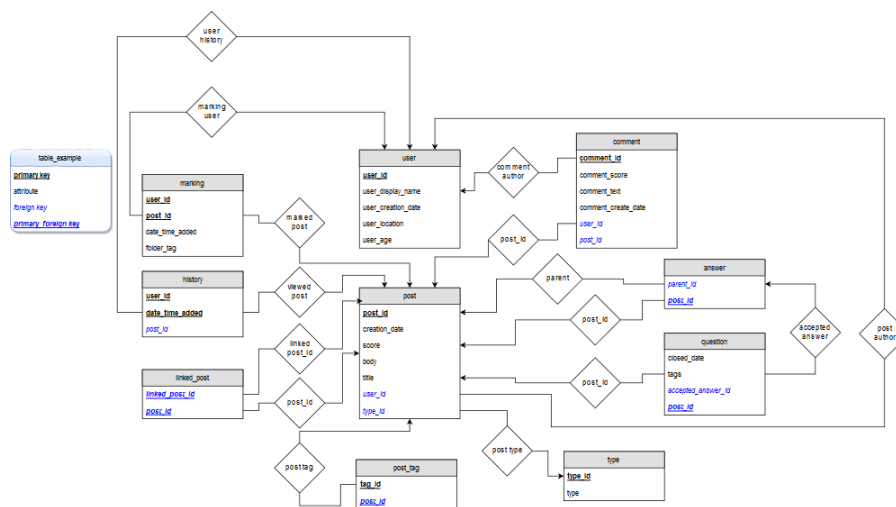


Figure 1.3: The ER model of our normalized schema

- *Post*

The *Post* table is a super-table: it contains many of the elements common to Questions and Answers such as a *post_id*, *body* and *title*. The *post_id* attribute is used to link information from other tables to a given post.

- *User*

The *User* table handles user data, It is referenced by other tables by the attribute *user_id*, that identifies the author of a post.

- *Question*

The *Question* table is a specialization of the *Post* table, with the notable characteristic that it can have an accepted answer, in the form of a foreign key referencing a *post_id*.

- *Answer*

The *Answer* table is also derivate from the *Post* table, and is only being referenced by *post_id*. The one distinctive attribute of the *Answer* Table is the *parent_id* attribute, which is a *post_id* that defines what question any given answer is attached to.

- *Type*

The *Type* table will be discussed in more detail in the Design section1.3, as it ties into one of the main design decision we did when refactoring the database from the original schema. For now, it should suffice to say that this class defines what kind of post there can be, either a question or an answer; every type has an associated *type_id*.

- *comment*

The comments table has been implemented such that it now references the super-table *Post* through *post_id* and *user_id*.

- *post_tags*

The *Post_tag* table contains tags, which can be attached to a post; this table is very simple in its structure, but its important because in the original schema of the Stack Overflow database, tags were stored as a string containing multiple tags separated by ":", as shown in the following picture.

tags
vb.net::vb.net-2010
c++::c::boost::makefile::cmake
c#::html::asp.net-mvc-4
java::class::casting::classloader
java::class::casting::classloader
ios::objective-c::constraints::autolayout
visual-studio-2010::sharepoint::.net-framework...
visual-studio-2010::sharepoint::.net-framework...
visual-studio-2010::sharepoint::.net-framework...
javascript::jquery::html::jquery-selectors
html::css
android
email::maven::maven-plugin
powershell
php::javascript::jquery::session::uploadify
php::javascript::jquery::session::uploadify
php::javascript::jquery::session::uploadify
javascript::python
c++::c++11::unordered-map::unordered-set::...
c++::c++11::unordered-map::unordered-set::...
c++::c++11::unordered-map::unordered-set::...
c++::c++11::unordered-map::unordered-set::...
python::class::methods::couchdb::decorator

Figure 1.4: Figure of how tags are stored in the sample database schema, from SQL workbench output

Since we needed to be able to include specific tags in our search, this data structure didn't suit our needs. In order to obviate this problem, our database file contains a procedure that splits each row in the *tag* column of the original schema into substrings, divided at the "::", that are subsequently inserted into our table. This procedure is described in more detail in chapter 2.

- *History*

The *History* table is the first of two tables that are needed to provide certain features in the software, in this case a history of the posts the user has viewed; This will most likely be a fixed-length list, and it will probably work as a stack: when the set limit is reached, the oldest post is deleted and the new one is added at the top of the list. The specific implementation of this feature will be done in software later on, but it will be supported by the database either via a special procedure or a trigger.

- *Marking*

The *Marking* table supports marking functionality in the software, similar to bookmarks. The idea is for the user to be able to bookmark certain questions, and furthermore he should be able to organize his bookmarks and group them

through the use of labels. To that end, the *label_tag* attribute stores the label of a given bookmark. As it can be seen from the table and schema, this attribute can only have one value for every post, we consider this a constraint of our design, the reason behind which is explained in 1.3

- *Linked_post*

The *Linked_post* table contains linked posts, which are feature of the Stack Overflow website and dataset. In brief, users can put links to other posts in a question or an answer. The *Linked_post* table contains the *post_id* of the post the link is posted in, and it thus has a similar function to a *parent_id*, while the attribute *linked_post_id* contains the id of the post the link is pointing to.

1.3 Design

Our database design is centered around the super-table *Post*: As it can be seen from the diagram above, all other tables reference the *Post* table by *post_id*. During our design process we took an object oriented approach, and implemented questions, answers and comments as different kinds of posts.

Our schema is made up of 10 tables: two of them (*History* and *Marking*) support features to be implemented in software, while the others try to adapt the existing data-set from Stack Overflow. We tried to ensure compatibility with the existing data and the schema of the Stack Overflow database wherever possible, even if that meant adding complexity to our schema. An example of this can be seen in our implementation of the *Type* table and the associated *type_id*.

A *type_id* is an attribute that identifies a post as an either a question or an answer: this was necessary in the original database schema because there was no distinction between the two; Such a distinction is present in our schema however, and thus *type_id* attributes can be seen as redundant data: *type* could be reasonably be inferred by *parent_id*, since answers have a parent questions while questions do not.

Our schema could have been simplified by removing them, but we decided against it because the removal of *type* may have led to unintended consequences or loss of data, especially when one takes bigger and more complex data samples into account. we wanted to ensure that our schema was compatible with the data, as well as make our transition process easier. Another question that came up during the design process was how to handle labels for the bookmarking functionality: as explained in the previous section, we wanted to give the user the possibility to customize and

organize his bookmarks; in our initial design, this feature would have taken the form of a hierarchical structure of labels or "folders". However in the end we decided against this kind of hierarchical data-structure, because it would have made the table non-atomic. As a result, at present there can only be a single label associated with any given bookmark.

Chapter 2

Functionality

This chapter looks at the functions and procedures implemented in the database.

2.1 Creation of the database

Due to the way the original posts table was structured, the field *tags* would violate the first normal form by having multiple values together. This led us to write a procedure that splits the tag field into atomic values, and inserts these into a table called *post_tags*. (The structure is described in chapter 1).

The functions and procedures used for this are described in section 2.1.1.

2.1.1 Splitting up non atomic tags

num_values_in_delimited_str(tag)

This is a helper function that takes a string delimited by ':' and returns the number of values in it. Eg. "*value_a::value_b::value_c*" returns 3. This is used by the procedure seen in listing 2.3.

Listing 2.1: num_values_in_delimited_str(tag)

```
RETURN ROUND((LENGTH(tag) - LENGTH(REPLACE(tag, '::', ''))) /  
    LENGTH('::')) + 1;
```

string_at_delimited_pos()

Another helper function is needed for extracting the value in a delimited string. This is used by the procedure seen in listing 2.3.

Listing 2.2: string_at_delimited_pos

```
RETURN REPLACE(SUBSTRING(SUBSTRING_INDEX(str, delim, pos),
    LENGTH(SUBSTRING_INDEX(str, delim, pos - 1)) + 1),
    delim, '');
```

split_insert_into_tags()

This is the procedure that splits up tags. A cursor containing *post_id* and *tag* is created and looped over using a read loop. Each value in the currently read tag is then extracted by looping over the tag using the helper functions 2.1 and 2.2, and inserting the returned value into the table *post_tag*.

Listing 2.3: split_insert_into_tags()

```
open tags_cur;
read_loop: loop
    fetch tags_cur into post_id, tag;
    set i_max = num_values_in_delimited_str(tag);
    set i = 1;
    if done then
        leave read_loop;
    end if;
    start transaction;
    while i <= i_max do
        insert ignore into post_tags(post_id, tag_name
            ) values (post_id, (string_at_delimited_pos
            (tag, ":", i)));
        set i=i+1;
    end while;
    commit;
end loop;
close tags_cur;
```

2.2 API for querying the database

We have introduced various functions that we expect will be useful as an API when developing an application using this database. Mainly these consists of search queries, and queries for listing data from different tables.

2.2.1 Searching

We currently have a few ways of searching through the posts, whose purpose and implementation will be described below.

`search_questions_by_tag`

This procedure simply queries all questions that has the given tag, and orders them by score. It seems likely that such a procedure could be useful when developing an API.

Listing 2.4: `search_questions_by_tag`

```
BEGIN
    select post.post_id, title, body, score, creation_date,
        closed_date
    from post, question, post_tags
    where post.post_id = question.post_id and post.post_id =
        post_tags.post_id and post_tags.tag_name = tag
    order by post.score desc limit lim;
END //
```

`fulltext_search`

This functions performs a fulltext search in question.title and question.body with the given string. It makes use of mysql's implementation of natural language fulltext search.

Listing 2.5: `fulltext_search`

```
BEGIN
    SELECT post_id, title, body, match (title,body) AGAINST
        (search_str IN NATURAL LANGUAGE MODE) AS score
    FROM post WHERE post.type_id = post_type and MATCH (title,body
        ) AGAINST
        (search_str IN NATURAL LANGUAGE MODE);
END //
```

Searching_Question

Finds questions that matches the last meaningful word of the input string, where that word must be in the title.

Listing 2.6: search_question

```
begin
  declare done int default false;
  declare a char(200);

  declare curl cursor for (
    select tag_name
    from post_tags
    where match(tag_name) against(inpute IN BOOLEAN MODE));
  declare continue handler for not found set done = true;

  open curl;
  read_loop: loop
    fetch curl into a;
    leave read_loop;
  end loop;
  close curl;

  select post_id, type_id, title, body, score from
  (select * from post_tags where match(tag_name) against(a IN
    BOOLEAN MODE)) as t
  natural join post
  where match(post.title) against(+a IN BOOLEAN MODE)
  and match(post.body) against(inpute in natural language mode)
  order by post.score desc;

end; //
```

2.2.2 History, Marking and Listing Data

add_marking

Inserts a marking to a given post for a given user. We have included a 'marking_label' here as well, with the purpose of using that to organize markings using labels.

Listing 2.7: add_marking

```
BEGIN
  insert into marking values (user_id, post_id, now(),
    marking_label);
END //
```

retrieve_answers

Retrieves answers to a given question.

Listing 2.8: retrieve_answers

```
BEGIN
  select /*post.post_id, title, body, score, creation_date*/
    from post, answer
  where post.post_id = answer.post_id
  order by post.score desc limit lim;
END //
```

Chapter 3

Normalization of the database

In this chapter we will go through the normalization of the data from the Stack Overflow sample. This will include going through first, second, third and Boyce-Codd normal form, as well as showing the Entity-relationship model 3.2 of the normalized dataset.

3.1 The data sample

The original structure of the Stack Overflow sample stored all the data in to tables (as seen below) with a lot of redundant entries, like how all posts and comments have all data about the user that made them.

- *Comments(commentid, postid, commentscore, commenttext, commentcreate-date, userid, userdisplayname, usercreationdate, userlocation, usage)*
- *Posts(id, posttypeid, parentid, acceptedanswerid, creationdate, score, body, closed-date, title, tags, owneruserid, owneruserdisplayname, ownerusercreationdate, owneruserlocation, ownerusage, linkpostid)* [1]

In order to facilitate the history and marking functionality of the Stack Overflow Viewer application we added two more tables (shown below).

- *History(user_id, date_time_added, post_id)*
- *Marking(user_id, post_id, date_time_added, folder_tag)*

To avoid having redundant entries in the data we normalized the structure of the database, to fit first, second, third and Boyce-Codd normal form.

3.2 Normal forms

3.2.1 First normal form

A table is in first normal form if it's atomic. Meaning that it doesn't contain any repeating attributes or groups of attributes.

The data set was already close to being compliant with first normal form and only needed to have the tags split into a separate table.

3.2.2 Second normal form

A table is in second normal form when there's no attributes in it that depend on part of a concatenated key. Meaning that if the table has a primary key that contains two attributes no attribute in the table depends on only one attribute of the key.

The data was already compliant with the second normal form, without furtherer work, as there were no partial dependencies.

3.2.3 Third normal form

A table is in third normal form when no attributes depend on a non-key attribute.

To be compliant with third normal form we removed the transitive dependencies, like how the user attributes are dependent on *user_id* which in turn is dependent on *post_id*. This was done by extracting the user information from *Posts* and *Comments* and storing this information in a new table called *User* and splitting the *Posts* table into three separate tables: *Post*, as a super-table, and *Answer* and *Question*, as sub-tables.

- *User*(*user_id*, *user_name*, *user_creation_date*, *user_location*, *user_age*)
- *Post*(*post_id*, *creation_date*, *score*, *body*, *title*, *owner_user_id*, *type_id*)
- *Question*(*accepted_answer_id*, *closed_date*, *post_id*)
- *Answer*(*parent_id*, *post_id*)
- *Comment*(*comment_id*, *comment_score*, *comment_text*, *comment_create_date*, *user_id*, *post_id*)
- *Post_tags*(*post_id*, *tag*)
- *Linked_posts*(*link_post_id*, *post_id*)
- *History*(*user_id*, *date_time_added*, *post_id*)
- *Marking*(*user_id*, *post_id*, *date_time_added*, *folder_tag*)

3.2.4 Boyce-Codd normal form

Boyce-Codd normal form is a stricter version of third normal form, that ensures that all redundancy from functional dependencies is removed. This is done by requiring that all dependences are either trivial functional dependencies, meaning that one side of the equation is a subset of the other side, or that the dependency is on the primary key.

We used the Boyce-Codd Decomposition Algorithm (seen below) to ensure that there was no redundancy left from the functional dependencies.

```

result := {R};
done := false;
compute F+;
while (not done) do
  if (there is a schema Ri in result that is not in BCNF)
    then begin
      let  $\alpha \rightarrow \beta$  be a nontrivial functional dependency that
        holds on Ri such that  $\alpha \rightarrow R_i$  is not in F+,
        and  $\alpha \cap \beta = \emptyset$ ;
      result := (result − Ri) ∪ (Ri −  $\beta$ ) ∪ ( $\alpha, \beta$ );
    end
  else done := true;

```

Figure 3.1: The Boyce-Codd Decomposition Algorithm

This resulted in a database structure like this:

- *User*(*user_id*, *user_name*, *user_creation_date*, *user_location*, *user_age*)
- *Post*(*post_id*, *creation_date*, *score*, *body*, *title*, *owner_user_id*, *type_id*)
- *Question*(*accepted_answer_id*, *closed_date*, *post_id*)
- *Answer*(*parent_id*, *post_id*)
- *Comment*(*comment_id*, *comment_score*, *comment_text*, *comment_create_date*, *user_id*, *post_id*)
- *Post_tags*(*post_id*, *tag*)
- *Linked_posts*(*link_post_id*, *post_id*)
- *History*(*user_id*, *date_time_added*, *post_id*)
- *Marking*(*user_id*, *post_id*, *date_time_added*, *folder_tag*)

3.3 Entity-relationship model

After normalization of the dataset we visualized it with an ER-model 3.2, both for later reference and to help keep an overview of the dataset.

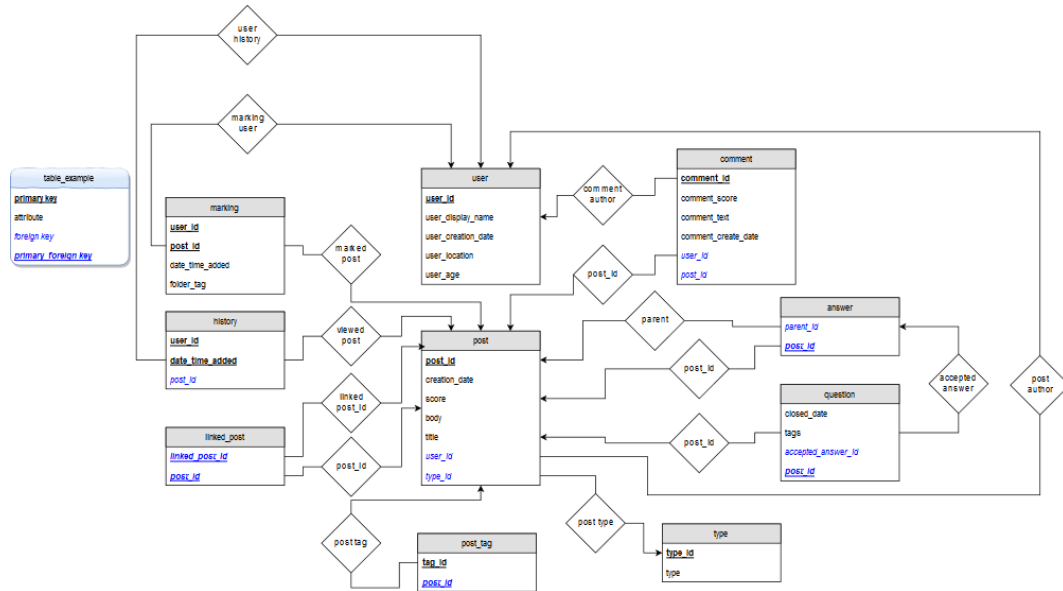


Figure 3.2: Entity-relationship model

Bibliography

- [1] Troels Andreasen & Henrik Bulskov. *RAWDATA Portfolio Subproject 1 Requirements*. <https://moodle.ruc.dk/>. 2017.

Appendix A

Initial design document

17

During design and normalization we worked on the tables in the following text document.

```
user (user_id(PK), user_name, user_creation_date, user_location, user_age)

post(post_id(PK), creation_date, score, body, title, owner_user_id(FK), type_id) /* type id
    somewhat redundant */

question(accepted_answer_id(FK), closed_date, tags, _post_id(PK,FK))
answer(parent_id, post_id(PK,FK))

comment(comment_id(PK), comment_score, comment_text, comment_create_date, user_id(FK), post_id(FK)
)

post_tags((post_id(FK), tag)(PK))//tag_id
//tags(TAG_ID, tag) <split tags on '::'>
```

```
//post_type(type_id, type); /* redundant? a post with a parentid is an answer, other posts are
    questions */ /* could remove and store type directly in the post table*/

history((USER_ID, DATE_TIME_ADDED) (PK), post_id(FK)); /*Index by users for search*/

marking(USER_ID, POST_ID, date_time_added, folder_tag);

linked_posts(LINK_POST_ID, POST_ID)

-----

//In posts(), index questions aside from answers
```

Appendix B

stack_overflow_normalized.sql

19

During design and normalization we worked on the tables in the following text document.

```
DROP DATABASE if exists stack_overflow_normalized;
CREATE DATABASE stack_overflow_normalized;
USE stack_overflow_normalized;

-- ----- DATABASE CREATION AND DATA INSERTION -----

-- user (user_id(PK), user_name, user_creation_date, user_location, user_age)
CREATE TABLE user (
    user_id INT UNSIGNED PRIMARY KEY,
    user_name VARCHAR(30) NOT NULL,
    user_creation_date DATETIME,
    user_location VARCHAR(200),
    user_age INT UNSIGNED
);
```

```

insert into user (user_id, user_name, user_creation_date, user_location, user_age)
select distinct owneruserid, owneruserdisplayname, ownerusercreationdate, owneruserlocation,
    owneruserage
from stackoverflow_sample_universal.posts
union
select distinct userid, userdisplayname, usercreationdate, userlocation, userage
from stackoverflow_sample_universal.comments;

-- insert ignore into user (user_id, user_name, user_creation_date, user_location, user_age)
-- select userid, userdisplayname, usercreationdate, userlocation, userage
-- from stackoverflow_sample_universal.comments;

-- post(post_id(PK), creation_date, score, body, title, owner_user_id(FK), type_id) /* type id
    somewhat redundant */
drop table if exists post;
CREATE TABLE post (
    post_id INT UNSIGNED PRIMARY KEY,
    creation_date DATETIME,
    score BIGINT,
    body TEXT,
    title VARCHAR(300),
    owner_user_id INT UNSIGNED NOT NULL REFERENCES user (user_id),
    type_id INT UNSIGNED,
    FULLTEXT (title,body), -- used by: fulltext_search
    FULLTEXT (body), -- used by: Searching_Questions
    FULLTEXT (title)-- used by: Searching_Questions
);

insert into post (post_id, creation_date, score, body, title, owner_user_id, type_id)
select distinct id, creationdate, score, body, title, owneruserid, posttypeid
from stackoverflow_sample_universal.posts;

```

```

-- question(post_id(PK,FK), accepted_answer_id(FK), closed_date, tags) /* no need for tags here
right? */
CREATE TABLE question (
    post_id INT PRIMARY KEY REFERENCES post (post_id),
    accepted_answer_id INT REFERENCES post (post_id),
    closed_date DATETIME
);

insert into question (post_id, accepted_answer_id, closed_date)
select distinct id, acceptedanswerid, closeddate
from stackoverflow_sample_universal.posts as p where p.posttypeid = 1;

-- answer(parent_id, post_id(PK,FK))
CREATE TABLE answer (
    post_id INT PRIMARY KEY REFERENCES post (post_id),
    parent_id INT
);

insert into answer (post_id, parent_id)
select distinct id, parentid
from stackoverflow_sample_universal.posts as p where p.posttypeid = 2;

-- comment(comment_id(PK), comment_score, comment_text, comment_create_date, user_id(FK), post_id(
FK))
CREATE TABLE comment (
    comment_id INT PRIMARY KEY,
    comment_score INT,
    comment_text MEDIUMTEXT,
    comment_create_date DATETIME,

```



```

        user_id INT REFERENCES user(user_id),
        post_id INT references post(post_id)
    );
insert into comment (comment_id, comment_score, comment_text, comment_create_date, user_id,
    post_id)
select commentid, commentscore, commenttext, commentcreatedate, userid, postid
from stackoverflow_sample_universal.comments;

-- post_tags((post_id(FK), tag) (PK))//tag_id
drop table if exists post_tags;
CREATE TABLE post_tags (
    post_id INT REFERENCES post (post_id),
    tag_name VARCHAR(50),
    primary key(post_id, tag_name),
    fulltext(tag_name) -- used in Searching_Questions
);

-- tags(TAG_ID, tag) <split tags on '::'>
/*drop table if exists tags;
CREATE TABLE tags (
    tag_id INT AUTO_INCREMENT PRIMARY KEY,
    tag_name VARCHAR(50)
);*/

-- string_at_delimited_pos(str, delim, pos) /* returns the value at n'th position eg.
    string_at_delimited_pos("a::b::c::d", "::", 3) returns "c" */
-- drop function if exists string_at_delimited_pos;
CREATE FUNCTION string_at_delimited_pos(str VARCHAR(255), delim VARCHAR(12), pos INT)
RETURNS VARCHAR(255)
RETURN REPLACE(SUBSTRING(SUBSTRING_INDEX(str, delim, pos),
    LENGTH(SUBSTRING_INDEX(str, delim, pos - 1)) + 1),

```

```

        delim, '');
-- select string_at_delimited_pos("qwer::rtyu::khhj", "::", 2);

-- num_values_in_delimited_str(str) /* Returns the number of values such that a::b::c returns 3*/
DELIMITER //
CREATE FUNCTION num_values_in_delimited_str ( tag varchar(200) )
RETURNS INT
BEGIN
    RETURN ROUND((LENGTH(tag) - LENGTH(REPLACE(tag, '::', ''))) / LENGTH('::')) + 1;
END; //
DELIMITER ;

/* This procedure creates a cursor containing id, tag from the original 'posts' table,
   iterates over that splitting up each tag inside the original 'tags' string, and inserts them
   into the post_tag table */
-- drop procedure if exists split_insert_into_tags;
delimiter //
create procedure split_insert_into_tags()
begin
    declare done int default false;
    declare tag varchar(90);
    declare post_id int;
    declare i_max int;
    declare i int;

    declare tags_cur cursor for select id, tags from stackoverflow_sample_universal.posts;
    declare continue handler for not found set done = true;

    open tags_cur;
    read_loop: loop
        fetch tags_cur into post_id, tag;

```

```

        set i_max = num_values_in_delimited_str(tag);
        set i = 1;
        if done then
            leave read_loop;
        end if;
        start transaction;
            while i <= i_max do
                insert ignore into post_tags(post_id, tag_name) values (post_id, (
                    string_at_delimited_pos(tag, ":", i)));
                set i=i+1;
            end while;
        commit;
    end loop;
    close tags_cur;
end; //
delimiter ;

call split_insert_into_tags();

-- linked_posts(LINK_POST_ID, POST_ID)
CREATE TABLE linked_posts (
    post_id INT REFERENCES post (post_id),
    link_post_id INT REFERENCES post (post_id),
    PRIMARY KEY (link_post_id , post_id)
);

insert into linked_posts (post_id, link_post_id)
select id, linkpostid
from stackoverflow_sample_universal.posts where linkpostid > "";

-- history((USER_ID, DATE_TIME_ADDED) (PK), post_id(FK)); /*Index by users for search*/

```

```

CREATE TABLE history (
    user_id INT REFERENCES user(user_id),
    datetime_added DATETIME,
    link_post_id INT REFERENCES post (post_id),
    PRIMARY KEY (user_id, datetime_added)
);
-- marking(USER_ID, POST_ID, date_time_added, folder_tag);
CREATE TABLE marking (
    user_id INT REFERENCES user(user_id),
    post_id INT REFERENCES post(post_id),
    datetime_added DATETIME,
    folder_tag varchar(200),
    PRIMARY KEY (user_id, post_id) /* added post_id to the PK instead of datetime_added. no need
        for multiple marks to one post */
);
-- //post_type(type_id, type); /* somewhat redundant. a post with a parentid is an answer, posts
    without are questions */ /* could also remove and store type directly in the post table*/

-- ----- FUNCTIONALITY / API -----

-- search_questions_by_tag(tag, lim) /* Basic tag search query. Returns relevant data about
    questions that contains the <tag>, ordered by score, limited to <lim> */
-- drop procedure if exists search_questions_by_tag;
DELIMITER //
CREATE PROCEDURE search_questions_by_tag (IN tag varchar(200), lim int)
BEGIN
    select post.post_id, title, body, score, creation_date, closed_date
    from post, question, post_tags
    where post.post_id = question.post_id and post.post_id = post_tags.post_id and post_tags.
        tag_name = tag
    order by post.score desc limit lim;

```

```

END //
DELIMITER ;
-- call search_questions_by_tag("c#", 50);

-- retrieve_answers(question_id) /* Retrieves the answers to a given question */
-- drop procedure if exists retrieve_answers;
DELIMITER //
CREATE PROCEDURE retrieve_answers (IN question_id int, lim int)
BEGIN
    select /*post.post_id, title, body, score, creation_date*/ from post, answer
    where post.post_id = answer.post_id
    order by post.score desc limit lim;
END //
DELIMITER ;
-- call retrieve_answers(9033, 50);

-- fulltext_search(search_str) /* Procedure that finds questions using mysql's built in fulltext
search (ignoring useless words, using multiword strings), searching in both title and body */
drop procedure if exists fulltext_search;
DELIMITER //
CREATE PROCEDURE fulltext_search (in search_str varchar(400), post_type int)
BEGIN
    SELECT post_id, title, body, match (title,body) AGAINST
    (search_str IN NATURAL LANGUAGE MODE) AS score
    FROM post WHERE post.type_id = post_type and MATCH (title,body) AGAINST
    (search_str IN NATURAL LANGUAGE MODE);
END //
DELIMITER ;
-- call fulltext_search('mysql', 1);
-- call fulltext_search('mysql', 2);
-- call fulltext_search('javascript tutorial',1);

```

```

-- call fulltext_search('machine learning',1);
-- call fulltext_search('how to python good',1);
-- call fulltext_search('Hi database teach me to be the bestest at searching thank you bye bye',1)
;

-- add_marking(user_id, post_id, marking_label) /* Inserts a marking to at given post <post_id>,
    for a given user <user_id>, with a given folder name <marking_label> and a <now()> timestamp*/
-- drop procedure if exists add_marking;
DELIMITER //
CREATE PROCEDURE add_marking (IN user_id int, post_id int, marking_label varchar(200))
BEGIN
    insert into marking values (user_id, post_id, now(), marking_label);
END //
DELIMITER ; /* todo: handle the user inserting multiple identical marks. gives duplicate error now
    . insert ignore would do it but its probably bad design*/
-- call add_marking(1185, 9033, 'MyFolder');
-- select * from marking;

-- Searching_Questions /* Finds questions that matches the last meaningful word of the input
    string, where that word must be in the title. */
drop procedure if exists Searching_Questions;
delimiter //
create procedure Searching_Questions( in inpute char (200))
begin
    declare done int default false;
    declare a char(200);

    declare curl cursor for (
        select tag_name
        from post_tags
        where match(tag_name) against(inpute IN BOOLEAN MODE));

```

```

declare continue handler for not found set done = true;

open curl;
  read_loop: loop
    fetch curl into a;
    leave read_loop;
  end loop;
close curl;

select post_id, type_id, title, body, score from
  (select * from post_tags where match(tag_name) against(a IN BOOLEAN MODE)) as t
natural join post
where match(post.title) against(+a IN BOOLEAN MODE)
and match(post.body) against(inpute in natural language mode)
order by post.score desc;

end;//
delimiter ;

-- Search by %word% in title
-- Search by tag in body

```

