

# Udiddit, a social news aggregator

## Introduction

Udiddit, a social news aggregation, web content rating, and discussion website, is currently using a risky and unreliable Postgres database schema to store the forum posts, discussions, and votes made by their users about different topics.

The schema allows posts to be created by registered users on certain topics, and can include a URL or a text content. It also allows registered users to cast an upvote (like) or downvote (dislike) for any forum post that has been created. In addition to this, the schema also allows registered users to add comments on posts.

Here is the DDL used to create the schema:

```
CREATE TABLE bad_posts (  
    id SERIAL PRIMARY KEY,  
    topic VARCHAR(50),  
    username VARCHAR(50),  
    title VARCHAR(150),  
    url VARCHAR(4000) DEFAULT NULL,  
    text_content TEXT DEFAULT NULL,  
    upvotes TEXT,  
    downvotes TEXT  
);  
  
CREATE TABLE bad_comments (  
    id SERIAL PRIMARY KEY,  
    username VARCHAR(50),  
    post_id BIGINT,  
    text_content TEXT  
);
```

## Part I: Investigate the existing schema

As a first step, investigate this schema and some of the sample data in the project's SQL workspace. Then, in your own words, outline three (3) specific things that could be improved about this schema. Don't hesitate to outline more if you want to stand out!

1. Foreign keys are currently not established between the tables. It is recommended to define the "post\_id" in the "bad\_comments" table as a foreign key to explicitly indicate the association between "bad\_comments" and "bad\_posts."
2. Further normalization is advisable for this schema. Consider decomposing the columns in the "bad\_posts" table into smaller tables, specifically separating data into tables for 'title' and 'topic.'
3. The "upvotes" and "downvotes" columns in the "bad\_posts" table contain multiple values per row. To adhere to the First Normal Form, it is essential to reorganize this structure.
4. Enhance search performance by creating an index on the "username" column for both tables. This indexing will significantly improve the efficiency of search operations.
5. Employ the unique constraint for the "username" column in the "users" table to prevent the occurrence of duplicate entries. This ensures data integrity and avoids redundancy.

## Part II: Create the DDL for your new schema

Having done this initial investigation and assessment, your next goal is to dive deep into the heart of the problem and create a new schema for Udiddit. Your new schema should at least reflect fixes to the shortcomings you pointed to in the previous exercise. To help you create the new schema, a few guidelines are provided to you:

1. Guideline #1: here is a list of features and specifications that Udiddit needs in order to support its website and administrative interface:
  - a. Allow new users to register:
    - i. Each username has to be unique
    - ii. Usernames can be composed of at most 25 characters
    - iii. Usernames can't be empty
    - iv. We won't worry about user passwords for this project
  - b. Allow registered users to create new topics:
    - i. Topic names have to be unique.
    - ii. The topic's name is at most 30 characters
    - iii. The topic's name can't be empty
    - iv. Topics can have an optional description of at most 500 characters.
  - c. Allow registered users to create new posts on existing topics:
    - i. Posts have a required title of at most 100 characters
    - ii. The title of a post can't be empty.
    - iii. Posts should contain either a URL or a text content, **but not both**.
    - iv. If a topic gets deleted, all the posts associated with it should be automatically deleted too.
    - v. If the user who created the post gets deleted, then the post will remain, but it will become dissociated from that user.
  - d. Allow registered users to comment on existing posts:
    - i. A comment's text content can't be empty.
    - ii. Contrary to the current linear comments, the new structure should allow comment threads at arbitrary levels.
    - iii. If a post gets deleted, all comments associated with it should be automatically deleted too.
    - iv. If the user who created the comment gets deleted, then the comment will remain, but it will become dissociated from that user.
    - v. If a comment gets deleted, then all its descendants in the thread structure should be automatically deleted too.

- e. Make sure that a given user can only vote once on a given post:
  - i. Hint: you can store the (up/down) value of the vote as the values 1 and -1 respectively.
  - ii. If the user who cast a vote gets deleted, then all their votes will remain, but will become dissociated from the user.
  - iii. If a post gets deleted, then all the votes for that post should be automatically deleted too.
- 2. Guideline #2: here is a list of queries that Udiddit needs in order to support its website and administrative interface. Note that you don't need to produce the DQL for those queries: they are only provided to guide the design of your new database schema.
  - a. List all users who haven't logged in in the last year.
  - b. List all users who haven't created any post.
  - c. Find a user by their username.
  - d. List all topics that don't have any posts.
  - e. Find a topic by its name.
  - f. List the latest 20 posts for a given topic.
  - g. List the latest 20 posts made by a given user.
  - h. Find all posts that link to a specific URL, for moderation purposes.
  - i. List all the top-level comments (those that don't have a parent comment) for a given post.
  - j. List all the direct children of a parent comment.
  - k. List the latest 20 comments made by a given user.
  - l. Compute the score of a post, defined as the difference between the number of upvotes and the number of downvotes
- 3. Guideline #3: you'll need to use normalization, various constraints, as well as indexes in your new database schema. You should use named constraints and indexes to make your schema cleaner.
- 4. Guideline #4: your new database schema will be composed of five (5) tables that should have an auto-incrementing id as their primary key.

Once you've taken the time to think about your new schema, write the DDL for it in the space provided here:

```
CREATE TABLE "users" (  
  
  "id" SERIAL PRIMARY KEY,  
  "username" VARCHAR(25) UNIQUE NOT NULL,  
  "login_date" TIMESTAMP,  
  CONSTRAINT "check_username_empty" CHECK(LENGTH(TRIM("username"))>0)  
  
);
```

-- 2.c.Find a user by their username

```
CREATE INDEX "find_username" ON "users" ("username" VARCHAR_PATTERN_OPS);
```

-- 2.a.List all users who haven't logged in in the last year

```
CREATE INDEX "find_login" ON "users" ("login_date");
```

```
CREATE TABLE "topics" (  
  
  "id" SERIAL PRIMARY KEY,  
  "topic" VARCHAR(30) UNIQUE NOT NULL,  
  "description" VARCHAR(500),  
  CONSTRAINT "check_topic_empty" CHECK(LENGTH(TRIM("topic"))>0)  
  
);
```

-- 2.e.Find a topic by its name

```
CREATE INDEX "find_topic" ON "topics" ("topic");
```

```
CREATE TABLE "posts" (  
  
  "id" SERIAL PRIMARY KEY,  
  "title" VARCHAR(100),  
  "url" TEXT,  
  "text_content" TEXT,  
  "post_date" TIMESTAMP,  
  "user_id" INTEGER REFERENCES "users" ON DELETE SET NULL,  
  "topic_id" INTEGER REFERENCES "topics" ON DELETE CASCADE,  
  CONSTRAINT "check_title_empty" CHECK(LENGTH(TRIM("title"))>0),  
  CONSTRAINT "url_or_text" CHECK(((("url") IS NULL AND ("text_content") IS NOT NULL) OR ((("url") IS NOT NULL AND ("text_content") IS NULL))  
  
);
```

-- 2.b.List all users who haven't created any post

```
CREATE INDEX "find_user_id" ON "posts" ("user_id");
```

-- 2.d.List all topics that don't have any posts.

```
CREATE INDEX "find_topic_id" ON "posts" ("topic_id");
```

-- 2.f.List the latest 20 posts for a given topic.

```
CREATE INDEX "find_latest_by_topic" ON "posts" ("id", "post_date", "topic_id");
```

-- 2. g.List the latest 20 posts made by a given user.

```
CREATE INDEX "find_latest_by_user" ON "posts" ("id", "post_date", "user_id");
```

-- 2. h.Find all posts that link to a specific URL, for moderation purposes.

```
CREATE INDEX "find_posts_by_URL" ON "posts" ("id", "url");
```

```
CREATE TABLE "comments" (  
  "id" SERIAL PRIMARY KEY,  
  "comment" TEXT NOT NULL,  
  "parent_id" INTEGER REFERENCES "comments" ON DELETE CASCADE,  
  "user_id" INTEGER REFERENCES "users" ON DELETE SET NULL,  
  "post_id" INTEGER REFERENCES "posts" ON DELETE CASCADE,  
  "comment_date" TIMESTAMP,  
  CONSTRAINT "check_comment_empty" CHECK(LENGTH(TRIM("comment"))>0)  
);
```

-- 2. i.List all the top-level comments (those that don't have a parent comment) for a given post.

```
CREATE INDEX "find_comment_by_post" ON "comments" ("id", "parent_id",  
  "post_id");
```

-- 2. j.List all the direct children of a parent comment.

```
CREATE INDEX "find_direct_children" ON "comments" ("id", "parent_id");
```

-- 2.k.List the latest 20 comments made by a given user.

```
CREATE INDEX "find_latest_comments" ON "comments" ("id", "comment_date",  
"user_id");
```

```
CREATE TABLE "votes" (  
  
"id" SERIAL PRIMARY KEY,  
"vote" INTEGER,  
"user_id" INTEGER REFERENCES "users" ON DELETE SET NULL,  
"post_id" INTEGER REFERENCES "posts" ON DELETE CASCADE,  
CONSTRAINT "check_vote" CHECK ("vote" = 1 OR "vote" = -1)  
  
);
```

2.1.Compute the score of a post, defined as the difference between the number of upvotes and the number of downvotes

```
CREATE INDEX "compute_post_score" ON "votes" ("post_id", "vote");
```

## Part III: Migrate the provided data

Now that your new schema is created, it's time to migrate the data from the provided schema in the project's SQL Workspace to your own schema. This will allow you to review some DML and DQL concepts, as you'll be using INSERT...SELECT queries to do so. Here are a few guidelines to help you in this process:

1. Topic descriptions can all be empty
2. Since the bad\_comments table doesn't have the threading feature, you can migrate all comments as top-level comments, i.e. without a parent
3. You can use the Postgres string function **regexp\_split\_to\_table** to unwind the comma-separated votes values into separate rows
4. Don't forget that some users only vote or comment, and haven't created any posts. You'll have to create those users too.
5. The order of your migrations matter! For example, since posts depend on users and topics, you'll have to migrate the latter first.
6. Tip: You can start by running only SELECTs to fine-tune your queries, and use a LIMIT to avoid large data sets. Once you know you have the correct query, you can then run your full INSERT...SELECT query.
7. **NOTE:** The data in your SQL Workspace contains thousands of posts and comments. The DML queries may take at least 10-15 seconds to run.

Write the DML to migrate the current data in bad\_posts and bad\_comments to your new database schema:

```
-- Migrate usernames into "users"
```

```
WITH temp_usernames AS (  
  SELECT DISTINCT username  
  FROM "bad_posts"  
  UNION ALL  
  SELECT DISTINCT username  
  FROM "bad_comments"  
  UNION ALL  
  SELECT DISTINCT regexp_split_to_table(upvotes, ',') as username  
  FROM "bad_posts"
```



```
        UNION ALL
        SELECT DISTINCT regexp_split_to_table(downvotes, ',') as username
        FROM "bad_posts"
    )
INSERT INTO "users" ("username")
SELECT DISTINCT username
FROM temp_usernames
ORDER BY username ASC;
```

-- Migrate the data into "topics"

```
INSERT INTO "topics" ("topic")
SELECT DISTINCT topic
FROM bad_posts;
```

-- Migrate the data into "posts"

```
INSERT INTO "posts" ("id", "title", "url", "text_content", "user_id",
"topic_id")
SELECT
    bp.id,
    SUBSTRING(bp.title, 1, 100) AS title,
    bp.url,
    bp.text_content,
    u.id AS user_id,
    t.id AS topic_id
FROM "bad_posts" bp
INNER JOIN "users" u
ON bp.username = u.username
INNER JOIN "topics" t
ON bp.topic = t.topic_name;
```

-- Migrate the data into "comments"

```
INSERT INTO "comments" ("comment", "user_id", "post_id")
SELECT
    bc.text_content AS comment,
    u.id AS user_id,
    p.id AS post_id
FROM "bad_comments" bc
INNER JOIN "posts" p
ON bc.post_id = p.id
```

```
INNER JOIN "users" u
ON bc.username = u.username;
```

```
-- Migrate the data into "votes"
```

```
WITH downvote_ids AS (
  SELECT id, regexp_split_to_table(downvotes, ',') AS downvote
  FROM bad_posts
), upvote_ids AS (
  SELECT id, regexp_split_to_table(upvotes, ',') AS upvote
  FROM bad_posts
)
INSERT INTO votes ("vote", "user_id", "post_id")
SELECT
  -1 AS vote
  u.id AS user_id,
  dv.id AS post_id,
FROM downvote_ids dv
INNER JOIN users u
ON u.username = dv.downvote
UNION ALL
SELECT
  1 AS vote
  u.id AS user_id,
  uv.id AS post_id,
FROM upvote_ids uv
INNER JOIN users u
ON u.username = uv.upvote;
```