# COL362 - Project - Milestone 1

DataBoys - Tanish Gupta, Rishi Jain, Daksh Khandelwal

20th March 2023

## 1 Introduction

Our project is about building a clone of a subset of StackExchange. For Milestone 1 we have built the basic setup of our database. We defined the underlying schema, key constraints and the associated ER-Diagrams. We computed the functional dependecies, and normalized the relations in our database to BCNF. The major chunk of the milestone was to study the dataset, and derive meanings and logical relationships between the tables.

The dataset has been taken from here, which is the official dump of stackoverflow. The XML files were converted to CSV, thanks to this repository, so that they could be easily read and imported into the Postgresql database. Functional dependencies and BCNF was studied from lecture slides and db-book.

Now we are set up with the basic structure and ready to work for the next stage. In Milestone-2, we aim to write the SQL queries that will drive our application and optimize them to obtain results faster.

## 2 Functional Dependencies

For the Badges relation, we have the following basic set of non-trivial FDs:

- ID → UserId
- ID → Class
- ID → Name
- ID → TagBased
- ID → Date

For the Users relation, we have the following basic set of non-trivial FDs:

- ID → Reputation
- ID → CreationDate
- ID → DisplayName
- ID → LastAccessDate
- ID → WebsiteUrl
- ID → Location
- ID → AboutMe
- ID → Views
- ID → UpVotes
- ID → DownVotes
- ID → AccountId

For the Comments relation, we have the following basic set of non-trivial FDs:

- ID → PostId
- ID → Score
- ID → Text
- ID → CreationDate
- ID → UserId
- ID → ContentLicense
- ID → UserDisplayName

For the Post History relation, we have the following basic set of non-trivial FDs:

- ID → PostHistoryTypeId
- ID → PostId
- ID → RevisionGUID
- ID → CreationDate
- ID → UserId
- ID → Text
- ID → ContentLicense
- ID → Comment
- ID → UserDisplayName

For the Posts relation, we have the following basic set of non-trivial FDs:

- ID → PostTypeId
- ID → AcceptedAnswerId
- ID → CreationDate
- ID → Score
- ID → ViewCount
- ID → Body
- ID → OwnerUserId
- ID → LastEditorUserId
- ID → LastEditDate
- ID → LastActivityDate
- ID → Title
- ID → Tags
- ID → AnswerCount
- ID → CommentCount
- ID → ContentLicense
- ID → FavoriteCount

- ID → ParentId

- ID → OwnerDisplayName

- ID → ClosedDate

- ID → LastEditorDisplayName

- ID → CommunityOwnedDate

For the Votes relation, we have the following basic set of non-trivial FDs:

- ID → PostId

- ID → VoteTypeId

- ID → CreationDate

- ID → UserId

- ID → BountyAmount

For the Tags relation, we have the following basic set of non-trivial FDs:

- ID → TagName

- ID → Count

- ID → ExcerptPostId

- ID → WikiPostId

For the Post Links relation, we have the following basic set of non-trivial FDs:

- ID → CreationDate

- ID → PostId

- ID → RelatedPostId

- ID → LinkTypeId

# 3   FD Preserving Normalizations

**Boyce Codd Normal Form (BCNF):** A relation schema R is in BCNF with respect to a set F of functional dependencies if for all functional dependencies in F$^+$ (closure of FDs) of the form $\alpha \to \beta$, where $\alpha \subseteq R$ and $\beta \subseteq R$, atleast one of the following holds:

- $\alpha \to \beta$ is trivial i.e. $\beta \subseteq \alpha$

- $\alpha$ is a superkey for R

We claim that the submitted dataset is in BCNF form. Let's prove the claim for one of the relations say `Badges`.

## 3.1   Badges Table

For the Badges relation, we have the following basic set of non-trivial FDs:

- ID → UserId

- ID → Class

- ID → Name

- ID → TagBased

- ID → Date

There are other trivial dependencies that hold here.

Note that Id is a superkey here, since Id can functionally determine all other attributes of the relation. So, the above FDs satisfy the conditions for BCNF.

Now consider any non trivial dependency $\alpha \to \beta$. Here, $\beta \subseteq \alpha$. The closure may include new FDs only by one of the following rules:

- Reflexive Rule: All FDs formed using reflexive rule satisfy the first condition of BCNF.

- Augmentation Rule: If the original FD is a trivial FD to which this rule is applied, then the new FD is also trivial (i.e. if $\alpha \to \beta$ is trivial, then so is $\gamma\alpha \to \gamma\beta$ since $\gamma\beta \subseteq \gamma\alpha$ follows from $\beta \subseteq \alpha$). In the other case, if original FD is not a trivial FD, then it must contain ID on its LHS i.e. ID $\subseteq \alpha$ (This can be proven by induction on the length of derivation of FD, since we have only 3 rules for FD closure). Thus, $\alpha$ is a superkey. So, both cases satisfy BCNF conditions.

- Transitivity Rule: Again, we can proceed on case by case basis. If both the original FDs are trivial, then the transitive FD is also trivial. If either of them is non-trivial, then that particular FD should have ID in LHS (Again, this can be proven by induction), and hence, $\alpha$ will be a superkey. Thus, BCNF conditions are satisfied.

Therefore, every FD satisfies one of the BCNF condition, and hence this relation is in BCNF.

Another way to show that all the FDs satsify the BCNF conditions is as follows:

We can compute the closure of FDs using the closure for each subset of attributes i.e. we can use the following algorithm to compute the $\alpha^+$ for any $\alpha \subseteq R$

$result := \alpha;$
**repeat**
  **for** each functional dependency $\beta \to \gamma$ in F **do**
    **if** $\beta \subseteq$ result **then** result := result $\cup \gamma$;
**until** (result does not change)

Now, there are 2 cases for $\alpha$:

- If ID $\subseteq \alpha$, then it easily follows, since all FDs formed using these subsets will have ID on LHS, and since ID is a superkey, the 2nd condition of BCNF is satsified.

- For the other case, if ID $\not\subseteq \alpha$, then, since all the non-trivial functional dependencies have ID on their LHS (can be proven by induction for the ones not listed above), the if condition in the for loop will not be true for any of the non-trivial dependencies. The only possible unions to the result can be through trivial FDs (which do not contain ID on LHS since otherwise $\beta \not\subseteq$ result), and hence any new FDs will also be trivial. Thus, 1st condition of BCNF will be satisfied.

Hence, the relation is in BCNF.

Since in all the relations, the non-trivial FDs are of similar form (ID - primary key on LHS), it follows that all the relations are in BCNF. Similar arguments may follow for their proof as well.

# 4 ER Diagram

A few notations in order to understand the ER diagram attached:

- Entity sets are represented by rectangles and Relationship sets are represented with diamonds.

- PK stands for Primary Key and FK stands for Foreign key.

- We indicate cardinality constraints on a relationship by drawing either a directed line ($\to$) or an undirected line (—) between the relationship set and the entity set in question.

- One-to-one: We draw a directed line from the relationship set to both entity sets.
- One-to-many: We draw a directed line from the relationship set to the "one" side of the relationship.
- Many-to-one: We draw a directed line from the relationship set to the "one" side of the relationship.
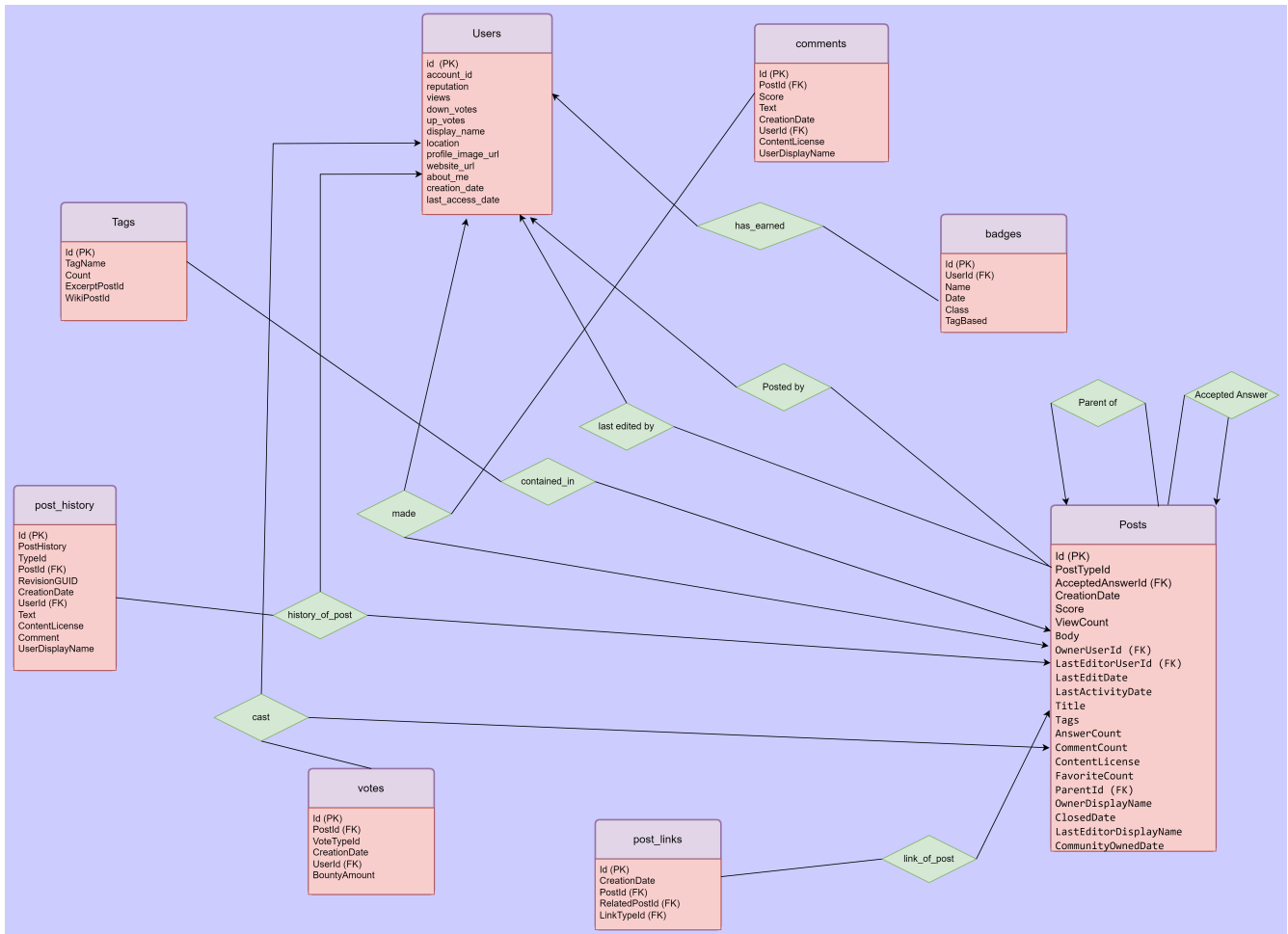- Many-to-many: We draw an undirected line from the relationship set to both entity sets.



Figure 1: ER Diagram

w

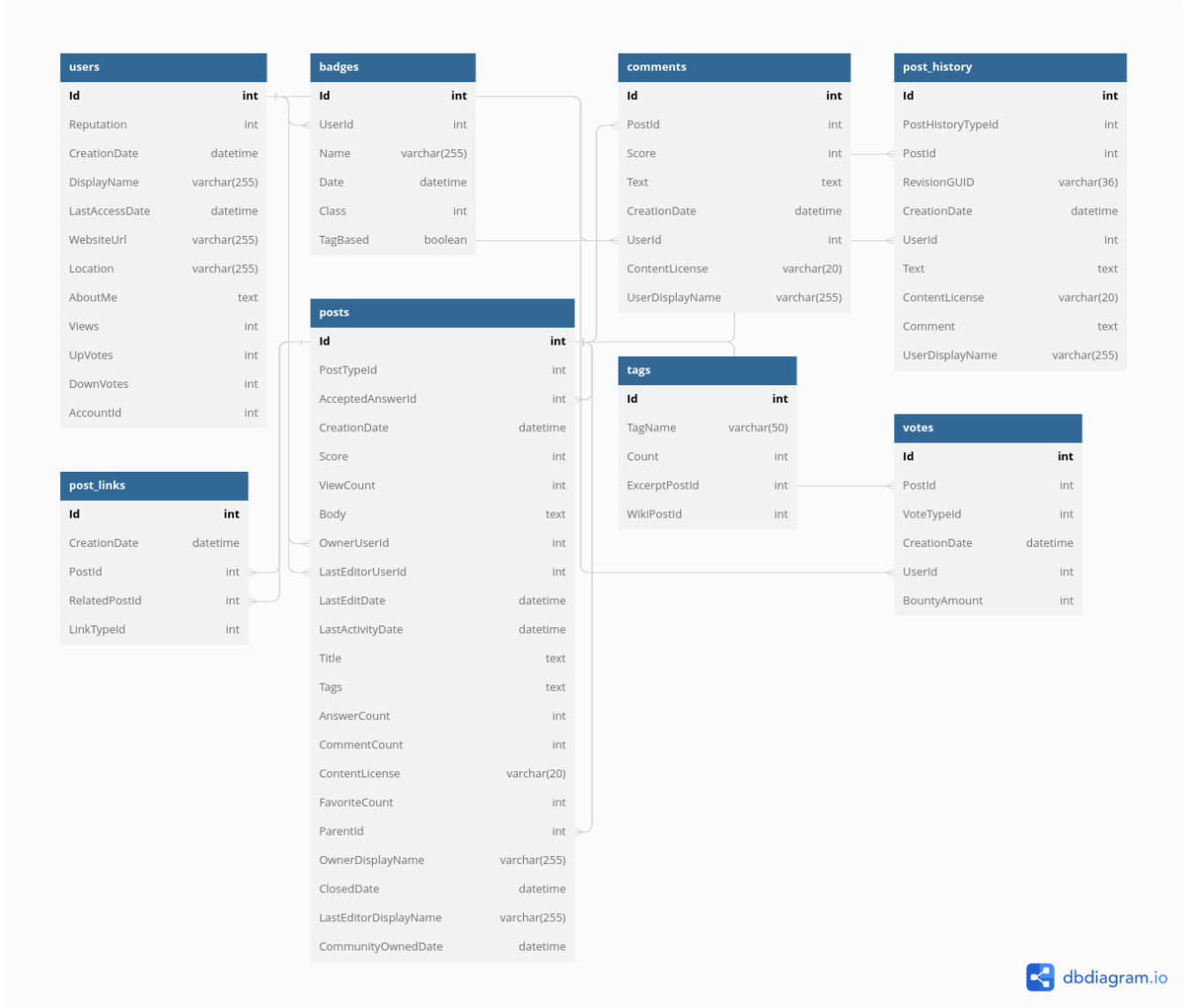# 5 Relational Schema

The Relational Schema is as follows :-



Figure 2: Relational Schema

# 6 Code snippet

The following code was used to generate the relational schema :-

```
1  BEGIN TRANSACTION;
2
3  DROP TABLE IF EXISTS users;
4  CREATE TABLE users (
5    Id INT NOT NULL PRIMARY KEY,
6    AccountId INT ,
7    Reputation INT NOT NULL,
8    Views INT NOT NULL,
9    DownVotes INT NOT NULL,
10   UpVotes INT NOT NULL,
11   DisplayName VARCHAR(255) NOT NULL,
12   Location VARCHAR(255) ,
13   WebsiteUrl VARCHAR(255) ,
```

```sql
14    AboutMe TEXT,
15    CreationDate TIMESTAMP NOT NULL,
16    LastAccessDate TIMESTAMP NOT NULL
17  );
18
19  DROP TABLE IF EXISTS badges;
20  CREATE TABLE badges (
21    Id INT NOT NULL PRIMARY KEY,
22    UserId INT NOT NULL REFERENCES users(Id),
23    Class INT NOT NULL,
24    Name VARCHAR(255) NOT NULL,
25    TagBased BOOLEAN NOT NULL,
26    Date TIMESTAMP NOT NULL
27  );
28
29
30  DROP TABLE IF EXISTS comments;
31  CREATE TABLE comments (
32    Id INT NOT NULL PRIMARY KEY,
33    PostId INT NOT NULL REFERENCES posts(Id),
34    UserId INT REFERENCES users(Id),
35    --If userid doesn't exist, then the user commented as a guest.
36    --She needs name and email, but email is not included in public data for privacy reasons.
37    Score INT NOT NULL,
38    ContentLicense VARCHAR(20),
39    UserDisplayName VARCHAR(255)  ,
40    Text TEXT NOT NULL,
41    CreationDate TIMESTAMP NOT NULL
42  );
43
44  DROP TABLE IF EXISTS post_history;
45  CREATE TABLE post_history (
46    Id INT NOT NULL PRIMARY KEY,
47    PostId INT NOT NULL REFERENCES posts(Id),
48    UserId INT REFERENCES users(Id),
49    PostHistoryTypeId INT NOT NULL,
50    UserDisplayName VARCHAR(255)  ,
51    ContentLicense VARCHAR(20),
52    RevisionGUID VARCHAR(36) NOT NULL,
53    Text TEXT ,
54    Comment TEXT,
55    CreationDate TIMESTAMP NOT NULL
56
57  );
58
59  DROP TABLE IF EXISTS post_links;
60  CREATE TABLE post_links (
61    Id INT NOT NULL PRIMARY KEY,
62    RelatedPostId INT NOT NULL REFERENCES posts(Id),
63    PostId INT NOT NULL REFERENCES posts(Id),
64    LinkTypeId INT NOT NULL,
65    CreationDate TIMESTAMP NOT NULL
66
67  );
68
69  DROP TABLE IF EXISTS posts;
70  CREATE TABLE posts (
71    Id INTEGER NOT NULL PRIMARY KEY,
72    OwnerUserId INTEGER REFERENCES users(Id),
73    LastEditorUserId INTEGER REFERENCES users(Id),
74    PostTypeId INTEGER NOT NULL,
75    AcceptedAnswerId INTEGER REFERENCES posts(Id),
76    Score INTEGER NOT NULL,
77    ParentId INTEGER REFERENCES posts(Id),
78    ViewCount INTEGER,
79    AnswerCount INTEGER,
80    CommentCount INTEGER,
81    OwnerDisplayName TEXT,
82    LastEditorDisplayName TEXT,
83    Title TEXT,
84    Tags TEXT,
85    ContentLicense TEXT,
86    Body TEXT,
```

```
 87    FavoriteCount INTEGER,
 88    CreationDate TIMESTAMP NOT NULL,
 89    CommunityOwnedDate TIMESTAMP,
 90    ClosedDate TIMESTAMP,
 91    LastEditDate TIMESTAMP,
 92    LastActivityDate TIMESTAMP NOT NULL
 93  ) ;
 94
 95  DROP TABLE IF EXISTS tags ;
 96  CREATE TABLE tags (
 97      Id INTEGER NOT NULL PRIMARY KEY,
 98      ExcerptPostId INTEGER ,
 99      WikiPostId INTEGER ,
100      TagName VARCHAR(50) NOT NULL,
101      Count INTEGER NOT NULL
102  ) ;
103
104  DROP TABLE IF EXISTS votes ;
105  CREATE TABLE votes (
106      Id INT NOT NULL PRIMARY KEY,
107      UserId INT   REFERENCES users (Id),
108      PostId INT NOT NULL REFERENCES posts (Id),
109      VoteTypeId INT NOT NULL,
110      BountyAmount INT  ,
111      CreationDate TIMESTAMP NOT NULL
112  ) ;
113
114  END TRANSACTION;
```