

## **Flaws in Bad giant table:**

### **Overview:**

The giant table is the collection of tweets that has been tweeted in twitter. It has 27 attributes to describe the tweets and tweet\_id as the **primary key**. They attributes are,

1. created\_at
2. text
3. tweet\_id
4. in\_reply\_to\_screen\_name
5. in\_reply\_to\_status\_id
6. in\_reply\_to\_user\_id
7. retweet\_count
8. tweet\_source
9. retweet\_of\_tweet\_id
10. hashtag1
11. hashtag2
12. hashtag3
13. hashtag4
14. hashtag5
15. hashtag6
16. user\_id
17. user\_name
18. user\_screen\_name
19. user\_location
20. user\_utc\_offset
21. user\_time\_zone
22. user\_followers\_count
23. user\_friends\_count
24. user\_lang character
25. user\_description
26. user\_status\_count
27. user\_created\_at

## **Flaws in Giant Table:**

### **1.User and tweet dependency**

tweet\_id -> { user\_name, user\_screen\_name, user\_location, user\_utc\_offset, user\_time\_zone, user\_followers\_count, user\_friends\_count, user\_lang, user\_description, user\_status\_count, user\_created\_at }

{ user\_name, user\_screen\_name, user\_location, user\_utc\_offset, user\_time\_zone, user\_followers\_count, user\_friends\_count, user\_lang, user\_description, user\_status\_count, user\_created\_at } -> user\_id

So,

tweet\_id -> user\_id

### **Problem and Explanation:**

Tweet ID will give the entire details of User along with the user\_id.

And the user\_id can give the entire detail of Users. So, the User and Tweet has the Reflexivity Functional Dependency.

### **2.In Reply section with additional user information**

tweet\_id -> { in\_reply\_to\_screen\_name, in\_reply\_to\_status\_id, in\_reply\_to\_user\_id }

in\_reply\_to\_user\_id -> { in\_reply\_to\_screen\_name }

So,

tweet\_id -> { in\_reply\_to\_status\_id, in\_reply\_to\_user\_id }

in\_reply\_to\_user\_id -> { in\_reply\_to\_screen\_name }

### **Problem and Explanation:**

user\_id can give the user screen name which is additionally add in the tweet column. That can be removed. Using Augmentation Rule of Functional dependency.

### **3. UTC and TimeZone dependency:**

user\_id -> { user\_utc\_offset, user\_time\_zone }

user\_utc\_offset -> {user\_time\_zone }

So,

user\_id -> {user\_utc\_offset}

user\_utc\_offset -> {user\_time\_zone }

### **Problem and Explanation:**

user\_utc\_offset can give the user\_time\_zone which is additionally add in the tweet column. That can be removed using Augmentation Rule of Functional dependency. And have it in the separate table.

### **4. Hash Tag:**

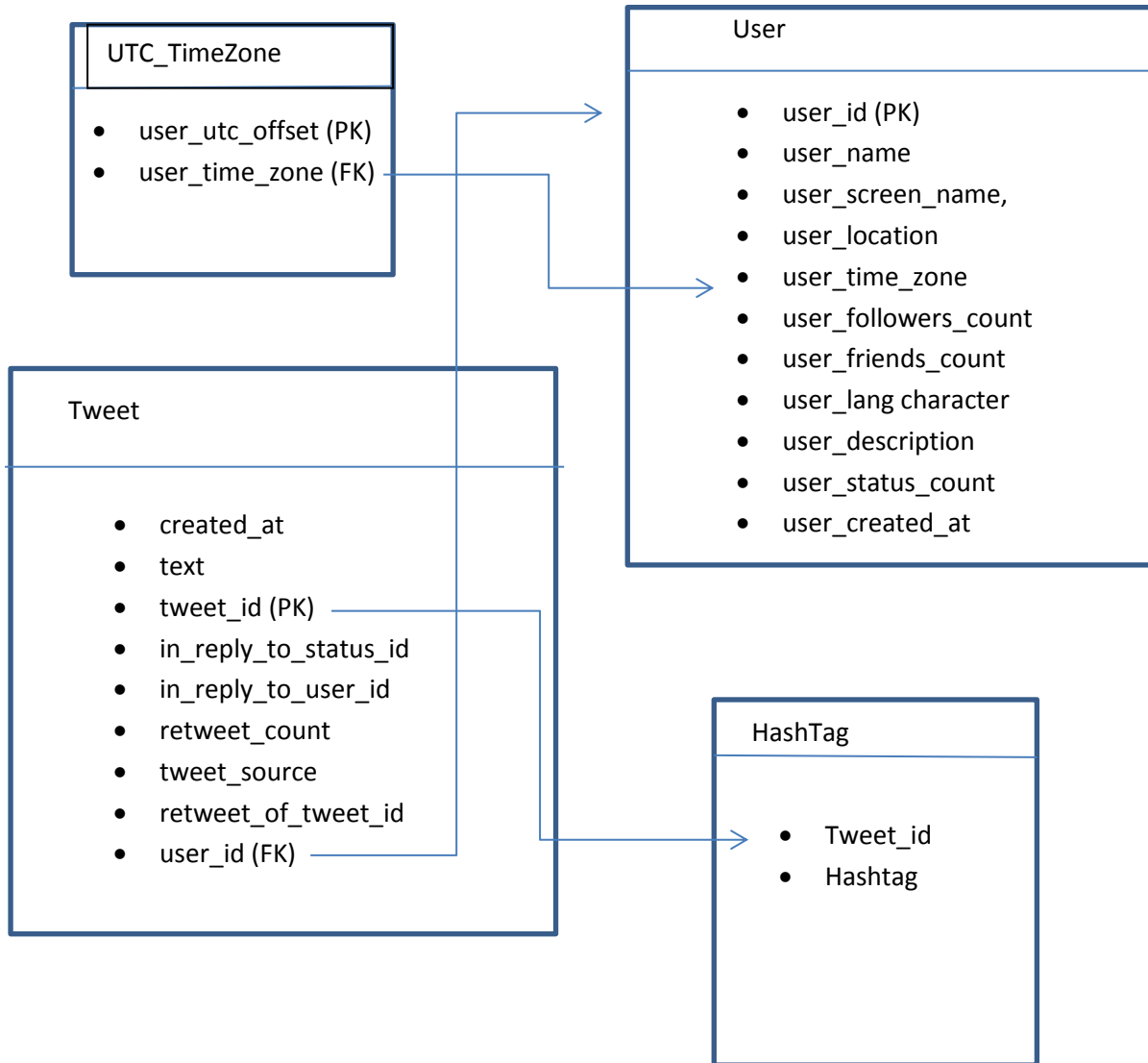
tweet\_id -> {hashtag1, hashtag2, hashtag3, hashtag4, hashtag5, hashtag6}

### **Problem and Explanation:**

Here, the hashtag count is limited by only 6. We can't add any further hashtag if needed. And NULL is stored in many cases which may reduce the efficiency of the database.

## New Database Design:

### Class Diagram:



## **Problems Resolved:**

Created tables without losing the functional dependencies and improved efficiency, which is the way to implement the 3NF.

### **1. Tweet and User dependency:**

As we have created table separately for tweet and user. Now the search can be done efficiently

### **2. In Reply section with additional user information**

Removed the column in\_reply\_to\_screen\_name from the tweet table.

### **3. UTC and TimeZone dependency:**

Separated the table for timezone and UTC column.

### **4. Hashtag**

hashtag table is created separately where we can add any number of hash tags to a tweet and efficiently identify the hashtags using tweet\_id.

## **Solutions to Queries:**

### **Tweets, users and languages**

#### **Question 1**

"110574"

#### **Explanation:**

Total number of tweets combines the normal tweets, replies, retweets

## **Question 2**

<b>User language</b>	<b>Number of tweets</b>
"ar"	"1405"
"ca"	"7"
"cs"	"2"
"de"	"75"
"el"	"2"
"en"	"74077"
"es"	"17910"
"eu"	"1"
"fi"	"1"
"fil"	"9"
"fr"	"231"
"hu"	"1"
"id"	"1423"
"it"	"26"
"ja"	"9861"
"ko"	"691"
"msa"	"14"
"nl"	"61"
"no"	"1"
"pl"	"4"
"pt"	"3977"
"ru"	"396"
"sv"	"2"
"th"	"233"
"tr"	"123"
"ur"	"1"
"zh-cn"	"26"
"zh-tw"	"14"

## **Explanation:**

Each language with number of tweets in each language.

### Question 3

**User language      fraction of tweets      fraction of users**

"ar"	"0.01270642284804746143"	"0.01353734090011273088"
"ca"	"0.000063306021306998028470"	"0.000067445826548604353146"
"cs"	"0.000018087434659142293848"	"0.000019270236156744100899"
"de"	"0.00067827879971783602"	"0.00072263385587790378"
"el"	"0.000018087434659142293848"	"0.000019270236156744100899"
"en"	"0.66993144862264185071"	"0.71374064189156638115"
"es"	"0.16197297737261924141"	"0.17256496478364342355"
"eu"	"0.000009043717329571146924"	"0.000009635118078372050449"
"fi"	"0.000009043717329571146924"	"0.000009635118078372050449"
"fil"	"0.000081393455966140322318"	"0.000086716062705348454045"
"fr"	"0.00208909870313093494"	"0.00222571227610394365"
"hu"	"0.000009043717329571146924"	"0.000009635118078372050449"
"id"	"0.01286920975997974207"	"0.01371077302552342779"
"it"	"0.00023513665056884982"	"0.00025051307003767331"
"ja"	"0.08918009658690107982"	"0.09501189937082678948"
"ko"	"0.00624920867473366252"	"0.00665786659215508686"
"msa"	"0.00012661204261399606"	"0.00013489165309720871"
"nl"	"0.00055166675710383996"	"0.00058774220278069508"
"no"	"0.000009043717329571146924"	"0.000009635118078372050449"
"pl"	"0.000036174869318284587697"	"0.000038540472313488201798"
"pt"	"0.03596686381970445132"	"0.03831886459768564464"
"ru"	"0.00358131206251017418"	"0.00381550675903533198"
"sv"	"0.000018087434659142293848"	"0.000019270236156744100899"
"th"	"0.00210718613779007723"	"0.00224498251226068775"
"tr"	"0.00111237723153725107"	"0.00118511952363976221"
"ur"	"0.000009043717329571146924"	"0.000009635118078372050449"
"zh-cn"	"0.00023513665056884982"	"0.00025051307003767331"
"zh-tw"	"0.00012661204261399606"	"0.00013489165309720871"

**Explanation:**

Each language with fraction of tweets in each language and fraction of users in each language.

fraction of tweets in each language = No of tweet in a language / Total tweets

fraction of users in each language = No of users in a language / Total users

**Retweeting habits****Question 1**

"0.33079204876372384105"

**Explanation:**

Fraction of tweets are retweeted.

Fraction of tweets Retweeted = No of tweet Retweeted / Total Tweets

**Question 2**

"0.00109677012472856706"

**Explanation:**

Average value for retweet per tweet.

**Question 3**

"0.66920795123627615895"

**Explanation:**

Fraction of tweets are not retweeted.

Fraction of tweets not Retweeted = No of tweet not Retweeted / Total Tweets



#### **Question 4**

"0.93619657423987555845"

#### **Explanation:**

Fraction of tweets less than the average value of retweet.

Fraction of tweets not Retweeted = No of tweet less than avg value Retweeted  
/ Total Tweets

#### **Hashtags**

#### **Question 1**

"10158 "

#### **Explanation:**

No of distinct tags

#### **Question 2**

#### **Tags**

#### **No of item used**

"ReasonsIFailAtBeingAGirl"	"467"
"RED"	"240"
"oomf"	"190"
"HonestyHour"	"172"
"TeamFollowBack"	"139"
"EresGuapaSi"	"130"
"10PeopleYouTrulyLove"	"126"
"TweetLikeAGirl"	"98"
"ImSingleBecause"	"97"
"WeAllGotThatOneFriend"	"96"

#### **Explanation:**

Top 10 hash tags frequently used.

### **Question 3**

#### **Explanation:**

No of distinct tags

### **Replies**

#### **Question 1**

"25776"

#### **Explanation:**

query for finding No of replies

#### **Question 2**

"0.4816997351146120281765611460"

#### **Explanation:**

probability of two user having same language

= sum ((No of users in each language C 2) / (Total No of users C 2))

ie. 2 represent user1 and user 2 as mentioned in question.

nCr formula is used

example for nCr:

$$5C2 = (5 * 4) / (2 * 1)$$

### **Question 3**

"0.4816997351146120281765611460"

### **Explanation:**

**As per my understanding question 2 and question 3 will produce same result**

probability of two user having same language

= sum ((No of users in each language C 2) / (Total No of users C 2))

ie. 2 represent user1 and user 2 as mentioned in question.

nCr formula is used

example for nCr:

$$5C2 = (5 * 4) / (2 * 1)$$