# CS5855 Final Coursework

- Ramachandran Murugesh

# Design Flaws of the twitter table (bad giant table):

As we can see from the structure of the bad\_giant\_table, the table is a result of merging more than 2 table for the sole purpose of data visibility: it helps with analysing the data without having much trouble to analyse more than 1 table to draw a result. This kind of data is dirty since it has a lot of redundant data, a single primary key and empty values in a few columns of the table. These types of data are purely used for visualization, i.e. to get a broader view of how the data is, i.e. the resultant of joining more than 2 tables. Tables like this (bad\_giant\_table) can't be used for query processing, since these are poorly normalized and unoptimized. For a small dataset, these might work perfectly, but when it comes to dealing with a large dataset, it has a negative impact on the performance of the database.

#### - Having a clear idea of the purpose of the data:

The purpose of switching from file-system to a database system is to efficiently store and retrieve data. To do this effectively, a database designer must know what the data is going to represent, how the data is going to be collected and the volume of the data we are going to dealing with. Because the volume of the data and purpose of the data is going to directly affect the levels of normalization, database design and in general the implementation of the entire database system. None of these pre-requisites have been kept in mind while creating the bad\_giant\_table.

#### Poor normalization:

Every table in a database should be normalized to the third normal form 3NF, since this will have a direct impact on the basic *CRUD* operations – *Create*, *Read*, *Update* and *Delete* operations in a database. Although the bad\_giant\_table is in 1NF (all the column values are atomic, no columns in the bad\_giant\_table has a column which holds more than one value), the table has a lot of redundant data which have to normalized. The following are the flaws in the bad\_giant\_table:

i) The table is not in 2NF, this can be explained with an example. Consider the following record-

Tweet id -260213890372211000, user id -502387030

This user has only tweeted once which corresponds the tweet\_id above. When we delete this record, we not only lose the details of the tweet, we lose the information pertaining to the particular user.

- ii) The column user\_id column is redundant; a particular user can have 'n' number of tweets corresponding to him/her.
- iii) The columns in\_reply\_to\_screen\_name and in\_reply\_to\_user\_id is directly dependent on the column in\_reply\_to\_status\_id which is already mapped to a tweet. Hence, the screen\_name and user\_id can be retrieved that way.
- iv) Some of the Hashtag columns 1-6, have all null values from 1 through 6. If there aren't any hashtags in a tweet, there is no reason to store null values in it.
- v) Functional dependency to represent the bad\_giant\_table:

Since tweet id is unique, it can independently identify each column of the table.

Tweet\_id -> created\_at

Tweet\_id -> text varchar

Tweet\_id -> in\_reply\_to\_screen\_name,

Tweet\_id -> in\_reply\_to\_status\_id,

```
Tweet_id -> in_reply_to_user_id,
Tweet id -> retweet count.
Tweet id -> tweet source,
Tweet id -> retweet of tweet id,
Tweet id -> hashtag1,
Tweet id -> hashtag2,
Tweet id -> hashtag3,
Tweet_id -> hashtag4,
Tweet id -> hashtag5,
Tweet id -> hashtag6,
Tweet id -> user id,
Tweet_id -> user_name,
Tweet_id -> user_screen_name,
Tweet id -> user location,
Tweet id -> user utc offset,
Tweet id -> user time zone,
Tweet id -> user followers count,
Tweet id -> user friends count,
Tweet id -> user lang,
Tweet id -> user description,
Tweet_id -> user_status_count,
Tweet id -> user created at
```

## - Bad Referential Integrity Constraints:

This table (bad\_giant\_table) has only one Primary Key constraint. If very few constraints or no constraints are implemented from the database design stage, the data processing completely relies on business logic.

## Poor Indexing:

As the number of records in the table increase, the time taken to process each query increases exponentially. Having as many columns as we have in the bad\_giant\_table will result in poor retrieving times on *SELECT* statements. The next thing on our mind is to index every column in the table, but this will have an adverse effect on operations like *CREATE*, *UPDATE* and *DELETE* statements.

#### - Failing to use Database Engine Features:

The database engine offers features such as creation of views, indexing, aggregator functions, transactions, procedures and triggers. In the bad\_giant\_table, we have a column by the name user\_status\_count which holds the total number of tweets of a particular user(including retweets). This result can be easily obtained by using the aggregator function *SUM*.

These are the design flaws that directly affect the bad\_giant\_table.

# Improved and Normalized design for Twitter Table (bad giant table):

#### Normalization of the table:

Since the table has no columns that have a group of values, i.e. the table satisfies the property of atomicity. This means that the table is in 1NF.

Since the table is not in 2NF as discussed from the flaws above, splitting the tables into smaller independent tables is a solution.

bad\_giant\_table(created\_at, text, tweet\_id PRIMARY KEY, in\_reply\_to\_screen\_name, in\_reply\_to\_status\_id, in\_reply\_to\_user\_id, retweet\_count, tweet\_source, retweet\_of\_tweet\_id, hashtag1, hashtag2, hashtag3, hashtag4, hashtag5, hashtag6, user\_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)

i) Splitting the user columns from the bad\_giant\_table:

USER\_TABLE(user\_id, user\_name, user\_screen\_name)
User\_id -> user\_name, user\_screen\_name (Each user\_id will have a user\_name and
user\_screen\_name, combination of user\_name and user\_screen\_name is going to yield
the user id.)

ii) Splitting the region and language settings from the bad giant table:

LOCALE\_SETTINGS (user\_id, user\_lang, user\_location, user\_utc\_offset, user\_time\_zone) user\_id -> user\_lang, user\_location, user\_utc\_offset, user\_time\_zone (Each user\_id will have a user\_lang, user\_location, user\_utc\_offset, user\_time\_zone)
In this table user\_id is a foreign key referencing the user\_id in USER\_TABLE. User\_id is also added with a UNIQUE constraint to avoid duplicate values in this table. Because every user is going to have only one user lang, user location, user utc offset and user time zone.

iii) Splitting the profile of the user from the bad giant table:

USER\_PROFILE(user\_id, user\_followers\_count, user\_friends\_count, user\_description,
user\_status\_count, user\_created\_at)
user\_id -> user\_followers\_count, user\_friends\_count, user\_description,
user\_status\_count, user\_created\_at
In this table user\_id is a foreign key referencing the user\_id in USER\_TABLE. User\_id is also
added with a UNIQUE constraint to avoid duplicate values in this table.

iv) Splitting the reply related fields from the bad\_giant\_table:

REPLY(in\_reply\_to\_user\_id, in\_reply\_to\_screen\_name, in\_reply\_to\_status\_id)
in\_reply\_to\_user\_id -> in\_reply\_to\_screen\_name, in\_reply\_to\_status\_id
(in\_reply\_to\_user\_id can uniquely identify in\_reply\_to\_screen\_name, in\_reply\_to\_status\_id)

Since multiple people can reply to a person's tweet, the columns in\_reply\_to\_screen\_name and in\_reply\_to\_user\_id can be stored once in the *REPLY* table and retrieved easily.

v) Splitting the hashtag fields from the bad giant table:

HASHTAG (tweet\_id, hashtag1, hashtag2, hashtag3, hashtag4, hashtag5, hashtag6)
Tweet\_id -> hashtag1, hashtag2, hashtag3, hashtag4, hashtag5, hashtag6
Each tweet\_id corresponds to hashtags from 1 to 6. The purpose of this table to limit the number of empty/null values in hashtag columns 1-6. If there exists no hashtags for a particular tweet id, there won't exist a record in hashtag table.

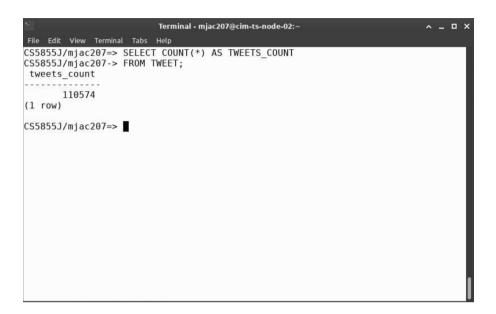
vi) TWEET(tweet\_id, text, user\_id, retweet\_of\_tweet\_id, retweet\_count, in\_reply\_to\_status\_id, tweet\_source, created\_at)
Tweet\_id -> text, user\_id, retweet\_of\_tweet\_id, retweet\_count, in\_reply\_to\_status\_id, tweet\_source, created\_at
Each tweet\_id can uniquely identify the columns above. In this table in\_reply\_to\_status\_id is a foreign key to USER TABLE.

**Note:** The tables *LOCALE\_SETTINGS* and *USER\_PROFILE* are created to reduce the number of columns for each user in the *USER\_TABLE*. The redundancy with user\_lang could not be reduced due to the inconsistency in the data. If each user\_utc\_offset and user\_time\_zone can uniquely identify a user\_lang, a separate table could have been drawn from this. Apart from these, the table has been normalized to the level 3NF.

# **Query Results:**

# Tweets, users and languages:

1) Total Number of tweets in total: (Since each tweet\_id is unique in TWEET table)



2) For every language the number of tweets in that language: Joining *TWEET* table with *LOCALE\_SETTINGS ON* user\_id and applying group by function to calculate the count.

File Edit View	Terminal Tabs Help
user_lang	Manager See Cont.
	1405
ar	1405
ca	7
cs	2 75
de	
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
(28 rows)	14

## Retweeting habits:

1) Fraction of tweets that are retweets. (Calculating the count of retweet\_of\_tweet\_id and dividing it by the total number of tweets present in the table)



2) Computing average number of retweets per tweet. (Assumption: Calculating the average of column – retweet\_count where retweet\_of\_tweet\_id is NOT NULL)



3) Computing fraction of tweets that are never retweeted. (Since we already calculated the fraction of tweets that are retweets, we subtract 1 from it to get the fraction of tweets that are never retweeted.

4) Computing fraction of tweets which are retweeted fewer times than the average number of retweets. (Assumption: Calculating the number of retweets which have a retweet\_count that is less than the average number of retweets – already calculated above. This is divided by the total number of tweets. The distribution of retweet\_count is not the best way to derive the average number of retweets per tweet. Since popular users have a very high value of retweet count which will affect the distribution of the retweet count.

## Hashtags:

1) Computing the number of distinct hashtags in the tweets.

```
Terminal - mjac207@cim-ts-node-02:~
                                                                           ^ _ D X
File Edit View Terminal Tabs Help
CS5855J/mjac207=> SELECT count(*) AS COUNT DISTINCT HASTAGS FROM
CS5855J/mjac207-> ((select hashtag1 AS distinct hashtag FROM HASHTAG WHERE hasht
ag1 IS NOT NULL)
CS5855J/miac207(> UNION
CS5855J/mjac207(> (select hashtag2 AS distinct hashtag FROM HASHTAG WHERE hashta
a2 IS NOT NULL)
CS5855J/miac207(> UNION
CS5855J/mjac207(> (select hashtag3 AS distinct hashtag FROM HASHTAG WHERE hashta
g3 IS NOT NULL)
CS5855J/mjac207(> UNION
CS5855J/mjac207(> (SELECT hashtag4 AS distinct_hashtag FROM HASHTAG WHERE hashta
q4 IS NOT NULL)
CS5855J/mjac207(> UNION
CS5855J/mjac207(> (SELECT hashtag5 AS distinct_hashtag FROM HASHTAG WHERE hashta
g5 IS NOT NULL)
CS5855J/mjac207(> UNION
CS5855J/mjac207(> (SELECT hashtag6 AS distinct hashtag FROM HASHTAG WHERE hashta
g6 IS NOT NULL))count distinct;
count distinct hastags
                  10158
(1 row)
CS5855J/mjac207=>
```

2) Computing the top ten most popular hashtags in the table.

```
Terminal - mjac207@cim-ts-node-02:~
                                                                            ^ _ D X
File Edit View Terminal Tabs Help
CS5855J/mjac207(> UNION ALL
CS5855J/mjac207(> (SELECT hashtag5 AS popular_hashtag FROM HASHTAG WHERE hashtag
5 IS NOT NULL)
CS5855J/miac207(> UNION ALL
CS5855J/mjac207(> (SELECT hashtag6 AS popular hashtag FROM HASHTAG WHERE hashtag
6 IS NOT NULL))hashtag count popular
CS5855J/mjac207-> GROUP BY popular hashtag
CS5855J/mjac207-> ORDER BY hashtag count DESC
CS5855J/mjac207-> LIMIT 10;
     popular hashtag
                            hashtag count
ReasonsIFailAtBeingAGirl
                                       467
 RED
                                       240
                                       190
 oomf
 HonestvHour
                                       172
 TeamFollowBack
                                       139
 EresGuapaSi
                                       130
10PeopleYouTrulyLove
                                       126
TweetLikeAGirl
                                        98
ImSingleBecause
                                        97
 WeAllGotThatOneFriend
                                        96
(10 rows)
CS5855J/mjac207=>
```

## Replies:

1) Computing the tweets that are neither replies nor replied to. (Assumption: the tweets that have a value in the column in\_reply\_to\_status\_id are the ones that have replied to a tweet. The value in the column in\_reply\_to\_status\_id holds the value of the id of the original tweet. We add these two counts and subtract it from the total number of tweets in the table. This will result in a set of tweets that are neither replies nor replied to.



2) Computing the probability of two arbitrary users having the same language setting. (Assumption: Since this is a completely random probability of two users having the same language setting. There are a total of 28 distinct languages. The probability of two users having the same language setting is 1/28 \* 1/28.

