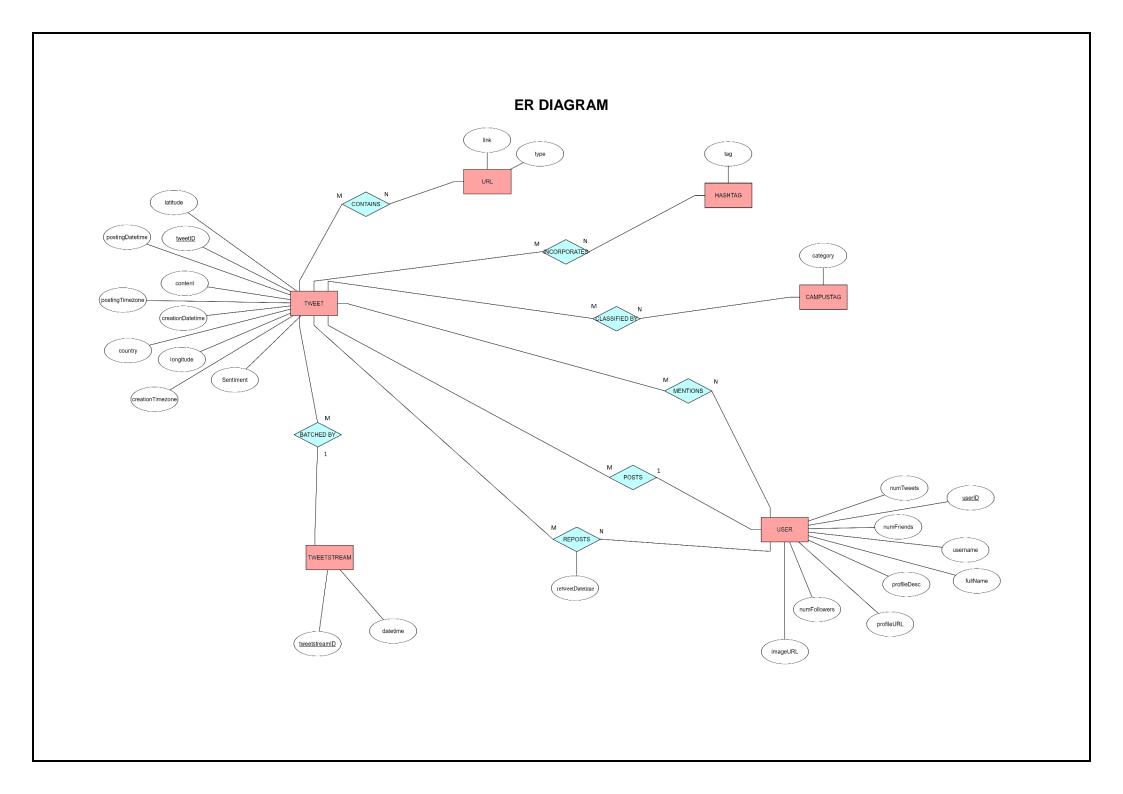# ASSIGNMENT 1

Kaustubh Jagtap – A0168820B

NATIONAL UNIVERSITY OF SINGAPORE

# ER DIAGRAM

# RELATIONS

TWEETSTREAM (<u>tweetstreamID</u>, date_time)

USER (<u>userID</u>, username, fullName, numFriends, numFollowers, profileDesc, profileURL, imageURL, numTweets)

TWEET (<u>tweetID</u>, content, creationDatetime, creationTimezone, postingDatetime, postingTimezone, country, longitude, latitude, sentiment, **tweetStreamID**, **userID**)

MENTION (**<u>userID</u>**, **<u>tweetID</u>**)

REPOST (**<u>tweetID</u>**, **<u>userID</u>**, retweetDatetime)

CAMPUSTAG (**<u>tweetID</u>**, <u>category</u>)

URL (**<u>tweetID</u>**, <u>link</u>, type)

HASHTAG (**<u>tweetID</u>**, <u>tag</u>)

**BOLD – Foreign Key**
<u>Underline</u> – Primary Key

# DATA TYPES

| Relation | Attribute | Comment | Data Type |
|---|---|---|---|
| TWEETSTREAM | tweetstreamID | Primary | BIGINT |
| | date_time | | DATETIME |
| TWEET | tweetID | Primary | BIGINT |
| | content | | TEXT |
| | creationDatetime | | DATETIME |
| | creationTimezone | | TINYINT |
| | postingDatetime | | DATETIME |
| | postingTimezone | | TINYINT |
| | country | | VARCHAR (74) |
| | longitude | | DECIMAL (9,6) |
| | latitude | | DECIMAL (9,6) |
| | sentiment | | TINYINT |
| | tweetStreamID | Foreign | BIGINT |
| | userID | Foreign | BIGINT |
| USER | userID | Primary | BIGINT |
| | username | | VARCHAR (50) |
| | fullName | | VARCHAR (100) |
| | profileURL | | VARCHAR (100) |
| | profileDesc | | TEXT |
| | numFollowers | | INT |
| | numFriends | | INT |
| | imageURL | | VARCHAR (100) |
| | numTweets | | INT |
| MENTION | userID | Primary, Foreign | BIGINT |
| | tweetID | Primary, Foreign | BIGINT |
| REPOST | tweetID | Primary, Foreign | BIGINT |
| | userID | Primary, Foreign | BIGINT |
| | retweetDateTime | | DATETIME |
| CAMPUSTAG | tweetID | Primary, Foreign | BIGINT |
| | category | Primary | CHAR (15) |
| URL | tweetID | Primary, Foreign | BIGINT |
| | link | Primary | VARCHAR (100) |
| | type | | TINYINT |
| HASHTAG | tweetID | Primary, Foreign | BIGINT |
| | tag | Primary | VARCHAR (50) |

## Assumptions and Clarifications

- TWEET.country will be stored as VARCHAR (74) since the longest country name is 74 characters.

- TWEET.content of the tweet has to be stored as TEXT since it can go up to 280 characters.

- URL.type is for now just 0 and 1, to signify whether it is external or twitter media. This leaves space to introduce a higher level of classification (say within external), and this can be done by using more integers to signify the types.

- Another table to store followers and friends between users could have been created, but I decided against this since this would make the database unnecessarily large and clunky. (Imagine having to store 500 followers for every person within Singapore, we would have to store the entire user profile of everyone in that network, even globally, and this would have a chain effect). For our purpose of sentiment analysis, simply knowing the number of friends and followers is enough to get a gauge of the person's influence.

- A separate table has to be created for REPOST because we need a way for the original tweet to be referenced. Also, no additional content or data is created for a retweet, other than the date and time of the new tweet.

**SQL CODE**

```sql
1    -- Name: Kaustubh Jagtap
2    -- ID: A0168820B
3
4    CREATE DATABASE IF NOT EXISTS assignment1;
5    USE assignment1;
6
7    CREATE TABLE tweetstream (tweetstreamID BIGINT, date_time DATETIME, PRIMARY KEY (tweetstreamID));
8
9    CREATE TABLE user (userID BIGINT, username VARCHAR (50), fullName VARCHAR (100), profileURL VARCHAR (100), profileDesc
10                 TEXT, numFollowers INT, numFriends INT, imageURL VARCHAR (100), numTweets INT, PRIMARY KEY (userID));
11
12   CREATE TABLE tweet (tweetID BIGINT, content TEXT, creationDatetime DATETIME, creationTimezone TINYINT, postingDatetime
13                 DATETIME, postingTimezone TINYINT, country VARCHAR(74), longitude DECIMAL(9,6), latitude DECIMAL(9,6),
14                 sentiment TINYINT, tweetstreamID BIGINT, userID BIGINT, PRIMARY KEY (tweetID), FOREIGN KEY
15                 (tweetstreamID) REFERENCES tweetstream (tweetstreamID), FOREIGN KEY (userID) REFERENCES user (userID));
16
17   CREATE TABLE mention (userID BIGINT, tweetID BIGINT, PRIMARY KEY (userID, tweetID), FOREIGN KEY (userID) REFERENCES user
18                 (userID), FOREIGN KEY (tweetID) REFERENCES tweet (tweetID));
19
20   CREATE TABLE repost (userID BIGINT, tweetID BIGINT, retweetDatetime DATETIME, PRIMARY KEY (userID, tweetID), FOREIGN KEY
21                 (userID) REFERENCES user (userID), FOREIGN KEY (tweetID) REFERENCES tweet (tweetID));
22
23   CREATE TABLE campustag (tweetID BIGINT, category CHAR(15), PRIMARY KEY (tweetID, category), FOREIGN KEY (tweetID)
24                 REFERENCES tweet (tweetID));
25
26   CREATE TABLE url (tweetID BIGINT, link VARCHAR(100), type TINYINT, PRIMARY KEY (tweetID, link), FOREIGN KEY (tweetID)
27                 REFERENCES tweet (tweetID));
28
29   CREATE TABLE hashtag (tweetID BIGINT, tag VARCHAR (50), PRIMARY KEY (tweetID, tag), FOREIGN KEY (tweetID) REFERENCES tweet
30                 (tweetID));
31
32   -- to run: mysql> source [path-to-file]\[filename.sql]
```

**Tables in Database**

```
mysql> desc tweetstream;
+----------------+-------------+------+-----+---------+-------+
| Field          | Type        | Null | Key | Default | Extra |
+----------------+-------------+------+-----+---------+-------+
| tweetstreamID  | bigint(20)  | NO   | PRI | NULL    |       |
| date_time      | datetime    | YES  |     | NULL    |       |
+----------------+-------------+------+-----+---------+-------+
2 rows in set (0.00 sec)
```

```
mysql> desc user;
+--------------+--------------+------+-----+---------+-------+
| Field        | Type         | Null | Key | Default | Extra |
+--------------+--------------+------+-----+---------+-------+
| userID       | bigint(20)   | NO   | PRI | NULL    |       |
| username     | varchar(50)  | YES  |     | NULL    |       |
| fullName     | varchar(100) | YES  |     | NULL    |       |
| profileURL   | varchar(100) | YES  |     | NULL    |       |
| profileDesc  | text         | YES  |     | NULL    |       |
| numFollowers | int(11)      | YES  |     | NULL    |       |
| numFriends   | int(11)      | YES  |     | NULL    |       |
| imageURL     | varchar(100) | YES  |     | NULL    |       |
| numTweets    | int(11)      | YES  |     | NULL    |       |
+--------------+--------------+------+-----+---------+-------+
9 rows in set (0.00 sec)
```

```
mysql> desc tweet;
+------------------+--------------+------+-----+---------+-------+
| Field            | Type         | Null | Key | Default | Extra |
+------------------+--------------+------+-----+---------+-------+
| tweetID          | bigint(20)   | NO   | PRI | NULL    |       |
| content          | text         | YES  |     | NULL    |       |
| creationDatetime | datetime     | YES  |     | NULL    |       |
| creationTimezone | tinyint(4)   | YES  |     | NULL    |       |
| postingDatetime  | datetime     | YES  |     | NULL    |       |
| postingTimezone  | tinyint(4)   | YES  |     | NULL    |       |
| country          | varchar(74)  | YES  |     | NULL    |       |
| longitude        | decimal(9,6) | YES  |     | NULL    |       |
| latitude         | decimal(9,6) | YES  |     | NULL    |       |
| sentiment        | tinyint(4)   | YES  |     | NULL    |       |
| tweetstreamID    | bigint(20)   | YES  | MUL | NULL    |       |
| userID           | bigint(20)   | YES  | MUL | NULL    |       |
+------------------+--------------+------+-----+---------+-------+
12 rows in set (0.00 sec)
```

```
mysql> desc mention;
+---------+------------+------+-----+---------+-------+
| Field   | Type       | Null | Key | Default | Extra |
+---------+------------+------+-----+---------+-------+
| userID  | bigint(20) | NO   | PRI | NULL    |       |
| tweetID | bigint(20) | NO   | PRI | NULL    |       |
+---------+------------+------+-----+---------+-------+
2 rows in set (0.00 sec)

mysql> desc repost;
+------------------+------------+------+-----+---------+-------+
| Field            | Type       | Null | Key | Default | Extra |
+------------------+------------+------+-----+---------+-------+
| userID           | bigint(20) | NO   | PRI | NULL    |       |
| tweetID          | bigint(20) | NO   | PRI | NULL    |       |
| retweetDatetime  | datetime   | YES  |     | NULL    |       |
+------------------+------------+------+-----+---------+-------+
3 rows in set (0.02 sec)

mysql> desc campustag;
+----------+------------+------+-----+---------+-------+
| Field    | Type       | Null | Key | Default | Extra |
+----------+------------+------+-----+---------+-------+
| tweetID  | bigint(20) | NO   | PRI | NULL    |       |
| category | char(15)   | NO   | PRI | NULL    |       |
+----------+------------+------+-----+---------+-------+
2 rows in set (0.00 sec)

mysql> desc url;
+---------+--------------+------+-----+---------+-------+
| Field   | Type         | Null | Key | Default | Extra |
+---------+--------------+------+-----+---------+-------+
| tweetID | bigint(20)   | NO   | PRI | NULL    |       |
| link    | varchar(100) | NO   | PRI | NULL    |       |
| type    | tinyint(4)   | YES  |     | NULL    |       |
+---------+--------------+------+-----+---------+-------+
3 rows in set (0.00 sec)

mysql> desc hashtag;
+---------+-------------+------+-----+---------+-------+
| Field   | Type        | Null | Key | Default | Extra |
+---------+-------------+------+-----+---------+-------+
| tweetID | bigint(20)  | NO   | PRI | NULL    |       |
| tag     | varchar(50) | NO   | PRI | NULL    |       |
+---------+-------------+------+-----+---------+-------+
2 rows in set (0.00 sec)
```

# QUERIES

A. <u>For every country available in our data, how many tweets originated from there?</u>

SELECT country, count(*) as numTweets
FROM tweet GROUP BY country;

B. <u>What are the 10 most common user hashtags used in our collected tweets?</u>

SELECT tag, count(*) as numTags
FROM hashtag GROUP BY tag
ORDER BY numTags
DESC LIMIT 10;

C. <u>Who is the most influential twitter user in our stream?</u> **(Use number of followers as proxy for influence)**

SELECT * FROM user
WHERE numFollowers =
(SELECT MAX(numFollowers) FROM user);

D. <u>How many tweets are there for each sentiment?</u>

SELECT sentiment, count(*) AS numTweets
FROM tweet GROUP BY sentiment;

E. (i) <u>Which faculty has the most tweets with a positive sentiment?</u>

SELECT category AS campus, count(*) AS numTweets
FROM campustag INNER JOIN tweet
ON tweet.tweetID = campustag.tweetID
WHERE sentimentID = 1
GROUP BY campus
ORDER BY numTweets
DESC LIMIT 1;

(ii) <u>Which faculty has the least tweets with a positive sentiment?</u>

SELECT category AS campus, count(*) AS numTweets
FROM campustag INNER JOIN tweet
ON tweet.tweetID = campustag.tweetID
WHERE sentimentID = 1
GROUP BY campus
ORDER BY numTweets
ASC LIMIT 1;

(iii) <u>Which faculty has the most tweets with a negative sentiment?</u> – **extension of question**

```
SELECT category AS campus, count(*) AS numTweets
FROM campustag INNER JOIN tweet
ON tweet.tweetID = campustag.tweetID
WHERE sentimentID = -1
GROUP BY campus
ORDER BY numTweets
DESC LIMIT 1;
```

F.  (i) <u>What were the most common open-day related hashtags (other than #NUSOpenDay18)?</u>

```
CREATE TABLE temp AS
(SELECT tweet.tweetID FROM tweet INNER JOIN HashTag
ON Tweet.tweetID = HashTag.tweetID
WHERE cast('2018-03-09' AS DATE) <= postingDate <= cast('2018-03-11' AS DATE)
AND tag = '#NUSOpenDay18'
GROUP BY tweetID);


SELECT tag, count(*) AS numTags FROM hashtag
WHERE tweetID IN (SELECT tweeID from temp)
AND tag != '#NUSOpenDay18'
GROUP BY tag ORDER BY numTags
DESC LIMIT 3;
```

(ii) <u>What were the top 3 retweeted tweets and who were their users?</u>

```
SELECT username, tweetID, count(*) AS numTweets
FROM repost INNER JOIN user on repost.userID = user.userID
WHERE tweetID in (SELECT tweetID from temp)
GROUP BY username, tweetID
ORDER BY numTweets
DESC LIMIT 3;
```

(iii) <u>What was the breakdown of sentiment across these tweets?</u>

```
SELECT sentiment, count(*) from tweet
WHERE tweetID IN (SELECT tweetID FROM temp)
GROUP BY sentiment;

DROP TABLE temp;
```