

Homework 1
Solutions

Notes:

1. Some sample appropriate answers are given for queries. The rest of the red text is prose explaining tradeoffs, and how design decisions can be made.
2. It is important to *load* the data and look at it before beginning.
3. Optional attributes in SQL queries are written in []. Do not copy paste the brackets as you will get a syntax error.

Part 1: SQL and Relational Algebra

If you wish to check your syntax, you can load this dataset into MySQL by following the directions on CCLE under Homework 1.

Hints: For some of these queries, you will need to use functions on attributes. Check out the list of date and time functions here, as it should be very useful: <https://dev.mysql.com/doc/refman/5.7/en/date-and-time-functions.html>. **You do not need to memorize them!** Note that these are *not* aggregation functions because aggregation functions may take multiple inputs and produce one output per group. These functions take one input and return one output without using any grouping variables. This does not mean that your queries will not use the aggregation functions we discussed in class though.

Bay Area Rapid Transit (BART) is a subway system that stops at various places on the Bay Area Peninsula, City of San Francisco, underwater to Oakland and the East Bay. When a user inserts a ticket, or scans a pass card on a turnstile, BART records that a user entered the subway system. On exit, the passenger does the same thing to exit the station (inserts a ticket or scans a fare card) and BART records that the user's journey is complete. Throughput can be defined as the number of passengers that entered at origin A , and exited at destination B .

In this exercise, we will do various queries on this data to answer several interesting questions.

This dataset consists of BART ridership data from 2017. The schema for the two tables in this database are provided below in case you do not wish to use the data.

-- This table is for your own information to see which station is which.

```
CREATE TABLE station(  
    abbreviation character(5),  
    description varchar(1000),  
    location varchar(23),  
    name varchar(50)  
);
```

```
CREATE TABLE rides2017(  
    origin character(4),  
    destination character(4),  
    throughput int,  
    datetime timestampz  
);
```

Exercises

- (a) Write a query to compute the total throughput (passengers, or number of trips) by *time of day* for the year of 2017. The result should contain only the hour of day, named **hour**, and the number of trips named **trips**.

```
SELECT
    EXTRACT(HOUR FROM datetime) AS hour,
    -- datetime::time is also fine since each row is an hour.
    SUM(throughput) AS trips
FROM hw1.rides2017
GROUP BY hour;
```

If you run the query, you should get the following results.

```
+-----+-----+
| hour   | trips  |
+-----+-----+
| 00:00:00 | 319034 |
| 01:00:00 | 67497  |
| 02:00:00 | 10274  |
| 03:00:00 | 3989   |
| 04:00:00 | 60098  |
| 05:00:00 | 448126 |
| 06:00:00 | 1269004 |
| 07:00:00 | 2945005 |
| 08:00:00 | 4411567 |
| 09:00:00 | 3494007 |
| 10:00:00 | 1900128 |
| 11:00:00 | 1447244 |
| 12:00:00 | 1422823 |
| 13:00:00 | 1474797 |
| 14:00:00 | 1632535 |
| 15:00:00 | 2071627 |
| 16:00:00 | 3062411 |
| 17:00:00 | 4620962 |
| 18:00:00 | 4362644 |
| 19:00:00 | 2472355 |
| 20:00:00 | 1377196 |
| 21:00:00 | 1062459 |
| 22:00:00 | 920195  |
| 23:00:00 | 685018  |
+-----+-----+
```

- (b) Write a query that lists the one pair of station codes that had the largest throughput on the weekdays in 2017.

```
SELECT
    origin,
    destination,
    SUM(throughput) AS total
FROM hw1.rides2017
WHERE EXTRACT(DOW FROM datetime) >= 1 and EXTRACT(DOW FROM datetime) <= 5
GROUP BY origin, destination
ORDER BY total DESC
LIMIT 1;
```

You should get one row, the Balboa to Montgomery trip (San Francisco).

```
+-----+-----+-----+
| origin | destination | total |
+-----+-----+-----+
| BALB   | MONT        | 217690 |
+-----+-----+-----+
```

- (c) Write a query that returns the 5 destinations that saw the highest average throughput on Mondays between 7am and 10am, and rank them from highest to lowest. Return the destinations and their averages.

```
SELECT
    destination,
    AVG(throughput)::DECIMAL(8,4) AS avg_trips
FROM hw1.rides2017
WHERE EXTRACT(DOW FROM datetime) = 1 AND EXTRACT(HOUR FROM datetime) >= 7
    AND EXTRACT(HOUR FROM datetime) < 10
    -- BETWEEN 7 AND 10 would not be accepted on the exam
    -- because 10:59:59am is evaluated to 10, yet it's after 10am. Careful!
GROUP BY destination
ORDER BY avg_trips DESC
LIMIT 5;
```

You should get the following 5 destination stations:

```
+-----+-----+
| destination | avg_trips |
+-----+-----+
| EMBR        | 188.0382  | -- Embarcadero Station (298 Market St, SF)
| MONT        | 175.1023  | -- Montgomery Station (598 Market St, SF)
| CIVC        | 71.5688   | -- SF Civic Center (1150 Market St)
| POWL        | 62.9843   | -- Powell Station (899 Market St, SF)
| 12TH        | 39.9817   | -- Oakland City Center (1245 Broadway)
+-----+-----+
```

If instead you used the **BETWEEN** keyword, your results are as follows. Notice the difference?! This is why we won't be able to accept it on the exam.

destination	avg_trips	
EMBR	158.7725	-- Embarcadero Station (298 Market St, SF)
MONT	150.2950	-- Montgomery Station (598 Market St, SF)
CIVC	60.9626	-- SF Civic Center (1150 Market St)
POWL	57.2327	-- Powell Station (899 Market St, SF)
12TH	33.3756	-- Oakland City Center (1245 Broadway)

- (d) Suppose we take the **result** from part (a) and call it **hourly_ridership**. Given the following query, write the equivalent expression in the relational algebra.

```
SELECT
    hour,
    trips / 100
FROM hourly_ridership
WHERE (hour >= 7 AND hour < 10) OR (hour >= 17 AND hour < 19);
```

$$\Pi_{\text{hour, trips}/100} \left(\sigma_{(\text{hour} \geq 7 \wedge \text{hour} < 10) \vee (\text{hour} \geq 17 \wedge \text{hour} < 19)} (\text{hourly_ridership}) \right)$$

One could also split up the predicate:

$$\Pi_{\text{hour, trips}/100} \left(\sigma_{(\text{hour} \geq 7 \wedge \text{hour} < 10)} (\text{hourly_ridership}) \cup \sigma_{(\text{hour} \geq 17 \wedge \text{hour} < 19)} (\text{hourly_ridership}) \right)$$

It could also be expanded by recalling that

$$\sigma_{p \wedge q}(R) = \sigma_p(\sigma_q(R))$$

- (e) Suppose we want to study how the weather affects how busy particular stations are. In the **Occupancy** relation, we have the name of a station called **Station**, the **DateTime**, and the number of people passing through the station called **Riders**, as attributes. In the **Weather** relation, we have the **Station**, **DateTime**, **Temperature**, **Condition** (for simplicity assume a string, like “cloudy”) attributes. We are *only* interested in comparing occupancy during “sunny” hours and “rainy” hours, and we only care about **Station**, **DateTime**, **Riders** and **Condition**. Write an expression in the relational algebra that represents this context.

$$\Pi_{\text{Station, DateTime, Riders, Condition}} \left(\sigma_{\text{condition} = \text{“cloudy”} \vee \text{condition} = \text{“sunny”}} (\text{Occupancy} \bowtie \text{Weather}) \right)$$

The following is also valid, though more complex.

$$\Pi_{\text{Station, DateTime, Condition}} \left(\left(\text{Occupancy} \bowtie \rho_{\text{Filtered}} \left(\Pi_{\text{Station, DateTime, Condition}} \left(\sigma_{\text{condition} = \text{“sunny”} \vee \text{condition} = \text{“rainy”}} (\text{weather}) \right) \right) \right) \right)$$

Part 2: Schemas and Architecture

Suppose you are working for the data team at *Bird Scooter*, a Santa Monica based startup that aims to revolutionize how people get around on wheels.

How the Bird Scooter service works: a user installs an app on their phone and enters their credit card information. The app shows a map of deployed scooters nearby. The user scans a QR code on the scooter using the app, which activates the scooter for use and begins the clock for per-minute and per-use billing. The user rides the scooter for a distance, for a certain number of minutes. When the user is done with the scooter, he/she leaves it somewhere, and marks the trip as complete in the app.

Each scooter has an identifier, and assume that since Bird is a startup, it only has 10,000 scooters deployed. Each scooter also has a flag that specifies a scooter as online (deployed), offline (not in use for whatever reason), or lost/stolen. Finally, each scooter is assigned to a home location that rarely changes. For example, some scooters may be assigned to UCLA, some may be assigned to Santa Monica, and others to Austin. Occasionally, a Santa Monica scooter may be reassigned to UCLA depending on demand, but we expect that such a change is very rare.

Assume that Bird currently only 500,000 registered users. A user is someone that has installed the app. Each user has an identifier, but not all users have a credit card number on file as many users install the app and then never ride a scooter. There is also an expiration date [present if a credit card number is present] and an email address.

Assume, for simplicity, that the app communicates over the Internet directly to a database server. Being a startup, this is probably how it is done, actually.

Exercises

- (a) Develop a schema for the **scooter** and **user** table based on the requirements and use case described on the previous page. Write the schema as a **CREATE TABLE** statement. Specify a primary key, or composite primary key using the correct syntax. Mark (with a comment) which attributes, if any, are foreign keys (to determine this, you may have to answer the other parts first). Try to minimize storage space.

The scooter Table

We need a unique identifier for each scooter, and this identifier (for now) must accept values up to 10,000. Trivially, this ID is unsigned. The proper data type is **smallint**. We could also make it a **serial** if we wish but it is not strictly necessary if we have our own format for assigning IDs. Of course, if we choose a format that is not a number, we will need to use another data type entirely. This will serve as the **PRIMARY KEY**.

To specify if a scooter is offline, online or lost/stolen, we use an **enum** as the flag. We could also use a **bit(2)**, one bit for online/offline, and one for stolen/lost or not stolen/lost, but this may cause data integrity concerns (if both bits are set, it means online and lost/stolen, which does not make sense). Later on, we can add more bits or more values to the **enum**.

Home location is a bit more vague. We could use an **enum**, but as Bird grows, we will need to add values to this **enum** frequently, and this may be cause for concern if we are not careful with how we add values to it. We could also use some kind of integer (start with **smallint**) as an identifier for a particular home location, and then join on this column with a fact table containing a mapping from the integer to the string with the name of the location, if we need it.

Thus, we can write the following schema:

```
CREATE TYPE scooter_status AS ENUM('online', 'offline', 'lost/stolen');
CREATE TABLE scooter (
    scooter_id smallint [or serial],
    -- serial can be omitted if we have a way of creating IDs.
    status      scooter_status DEFAULT 'offline'
    home        smallint,
    PRIMARY KEY (scooter_id)
    -- no foreign keys expected
);
```

Note that if we have a single column as a primary key, we can just specify it as such after defining the column:

```
CREATE TYPE scooter_status AS ENUM('online', 'offline', 'lost/stolen');
CREATE TABLE scooter (
    scooter_id smallint [or serial] PRIMARY KEY,
    -- serial can be omitted if we have a way of creating IDs.
    status      scooter_status DEFAULT 'offline'
    home        smallint
    -- no foreign keys expected
);
```

In the above schema, I added a `DEFAULT` value for `status`, but this is not truly necessary since it is not specified in the problem. Instead of using a `DEFAULT` value, we can force the user to set the flag by specifying `status` as `NOT NULL`.

While it may seem like a good idea to also specify a scooter's last known location, this is redundant, and would require us to `UPDATE` many rows over and over again each time a scooter is used and parked. `scooter` is a simple lookup table, and it should not have to be updated except to add new scooters, or move scooters from locale to locale.

I did not specify anything about `NULL` values in the home location. Perhaps when we add a scooter, we don't need to assign it to a location, then I can have `NULL`s, otherwise, I would specify it as `NOT NULL`.

Best practices of NOT NULL: It depends on the business case. *However*, it is much easier to start with a column being marked `NOT NULL` and then later removing the restriction than it is to allow `NULL` values and then later place the `NOT NULL` restriction, because what would happen to all of the existing `NULL`s? Something to keep in mind. There is actually a lot of history behind this.

The user Table

Each user has a unique identifier. For simplicity, we can assume this to be a `integer` since there are currently 500,000 users and this gives us plenty of room to grow. We could also specify that it is a `serial`. Of course, we can use a different data type for the IDs entirely, such as a `char` like in the `youtube` examples, or `uuid`. Not all users have a credit card associated with their accounts, so credit card number can be `NULL`. We **should** secure the credit card number by using industry best practices using some kind of hash, but it is not necessary for CS 143 (unless ever explicitly stated). There are federal regulations in terms of how to protect credit card numbers, but we will just use MD5 as an example for this solution. The length of the credit card number depends on which credit cards we accept. By using an integer type, we can avoid worrying about this nuance, or can we? What happens if a credit card number starts with 0? For our purposes, we will ignore this fact, but it is important that we *clearly* understand the formats of our attributes before we assign data types. We also need an expiration date. Since not all users have a credit card number, they also will not have an expiration date. Expiration dates are written as `mm/yy` but we do not have this exact type in PostgreSQL, so we can improvise by using a `date` column type. There is an edge case here though: if a card expires 4/2019, and it is 4/30/18 at home, but you are in a different timezone where it is already 5/1/18, the card may not work. We could research how expiration dates work, or we can just be safe and use a `timestamptz`. Once the local time hits 5/1/18, the card will no longer work. Finally, we have an email address field. You can make up something reasonable here, but RFC 5321 states a max of about 255 characters.

Point: When designing a schema, it is crucial to understand the format of the data that will be inserted into the table, and how it may change over time.

Putting all of this together, *one* good solution is

```
CREATE TABLE user (  
    user_id    integer [or serial] PRIMARY KEY,  
    ccnumber   varchar(11),  
    -- if we want a hash, we could use char(32/64) or BINARY(32/64) on \texttt{uuid}  
    expiration timestamptz,  
    -- our application would just have to ignore the day part, and the time part.  
    -- date may also work, depends on the edge case.  
    email      varchar(255) NOT NULL  
    -- It wasn't clear in the problem, but the email address should be required.  
);
```

- (b) Each interaction between a user and a Bird Scooter is called a *trip*, and we will create a schema for this `trip` table. Assume each trip has a unique identifier. The app will use this database table to determine where an available scooter is located. It will also (somehow) be used to determine the amount to charge the user. Additionally, data scientists would like to be able to use this table to determine daily and hourly trends in when users start and end scooter use and also identify scooter hotspots: areas where people frequently activate scooters and park them (just assume a location is a GPS coordinate: latitude and longitude, which can be represented numerically). Write the `CREATE TABLE` statement for this table. Specify the primary key. Mark (with a comment) which attributes, if any, are foreign keys (to determine this, you may have to answer the other parts first). Try to minimize storage space.

Again, each trip has a unique identifier which we can define as some kind of integer, without more information to the contrary. We do not know how many trips we expect and would have to monitor the situation carefully. We will just use an `integer` with the expectation that we will have to promote. Each trip is also associated with a user, otherwise we cannot bill them! Each trip has a start `timestampz` and an end `timestampz`. The data scientists can use these values to study time trends and finance can use it to bill the customer, thus both of these must be `NOT NULL`. Each trip also has a start location and an end location. Let's assume that this data must be collected and thus must be `NOT NULL`. If we have no GPS signal from the user, we cannot locate scooters for them anyway. Finally, we may also want to know which scooter was used, say, to locate a stolen or missing scooter.

One appropriate solution:

```
CREATE TABLE trip (  
    trip_id          integer [or serial],  
    user_id          integer,  
    -- user_id is a foreign key that references user table.  
    scooter_id       smallint NOT NULL,  
    -- scooter_id is a foreign key that references scooter table.  
    trip_start_time  timestampz NOT NULL,  
    trip_end_time    timestampz NOT NULL,  
    trip_start_loc   point NOT NULL,  
    trip_end_loc     point NOT NULL,  
    PRIMARY KEY(trip_id, user_id), -- trip_id is also appropriate.  
    FOREIGN KEY (user_id) REFERENCES user(user_id),  
    FOREIGN KEY (scooter_id) REFERENCES scooter(scooter_id)  
);
```


- (c) For the `trip` table, there are two ways that we can write data to the table from the app. First, we could insert a row when the user activates the scooter, and then modify the row when the user parks the scooter and ends the session. Second, we can simply cache the ride data on the phone, and at the end of the session transmit this data to the database server to be inserted as a row. What are the advantages and disadvantages of both of these methods? Which would you (an employee at Bird) prefer and why? Is there an even better way?

In the first method, we write a row to the database with a `NULL` end time and location. First off, this violates the schema for the `trip` table. Even if it did not violate it though, this constant modification of rows in the table each time a scooter is activated and parked induces a lot of load on the database system because each row is written once, and then modified once. This table thus becomes more **read and write heavy** than it needs to be.

On the other hand, if instead we cache the identifiers, start time and start location of a trip and then only write a row to the table once we park the scooter, we will have a data integrity nightmare. What happens if the phone loses network connectivity? What happens if the user switches off the phone? We lose the entire trip and the user does not pay. One way around this problem is to simply charge the user the maximum daily rate (like when you lose a parking voucher in a paid lot), but we *still* do not have a record that a ride even started!

Between these two cases, the first method would be preferable to the business; however, there are at least two better ways to handle this situation.

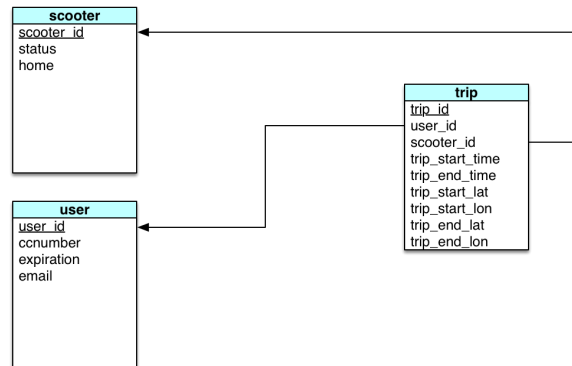
In the first, we split the `trips` table into two: `trip_starts` and `trip_ends` and we write a row once when the scooter is activated and we do not let the user scoot away until the record is written. Then when the scooter is parked, we write a row to the `trip_ends` table. To charge the user at the end of the month, we would then `LEFT JOIN` these two tables together on `trip_id` and compute the amount of time and a charge. The data scientists would also use this joined table. The start table would become **write heavy** and the end table would still be **read and write heavy** (the map of scooters makes this table read heavy), but the situation is not as severe as it was before.

```
CREATE TABLE trip_starts (  
    trip_id      smallint [or serial],  
    user_id      smallint,  
    scooter_id   smallint NOT NULL,  
    start_time   timestampz NOT NULL,  
    start_loc    point NOT NULL,  
    PRIMARY KEY (trip_id, user_id),  
    FOREIGN KEY (user_id) REFERENCES user(user_id),  
    FOREIGN KEY (scooter_id) REFERENCES scooter(scooter_id)  
);
```

The second method uses one table called `trips` still, but adds an additional `action` flag as a column (bit or enum) to mark the row as a `TRIP_START` or `TRIP_END`. The primary key then becomes composite: `trip_id` and `action`. This table is still **read write heavy**.

```
CREATE TYPE action_type AS ENUM('trip_start', 'trip_end');
CREATE TABLE trips (
    trip_id      smallint,
    user_id      smallint,
    scooter_id   smallint NOT NULL,
    action_time  timestampz NOT NULL,
    action_loc   point,
    action       action_type,
    PRIMARY KEY(trip_id, user_id, action),
    FOREIGN KEY (user_id) REFERENCES user(user_id),
    FOREIGN KEY (scooter_id) REFERENCES scooter(scooter_id)
);
```

- (d) Using the tables you just developed in all of Part 1, draw the schema diagram. For an example of a schema diagram, see Figure 2.8 in the text (page 47), or page 23 in the lecture slides for Lecture 2.



Some other things to think about:

- Would you include the number of minutes the trip lasted in the `trip` table? Why or why not?
No, because this is redundant information and can cause data integrity issues if for some reason the difference in time does not match the trip length entered into the row, though we can add this as a constraint.
- What are the advantages and disadvantages of including the `charge` to the user in the `trips` table?
The advantage is that it would allow fast querying of charges, and may not be a bad idea since a user is only charged according to the current fee structure. In other words, we would never be able to recalculate the charge. The disadvantage is redundancy. The charge is a function of the start time, end time and current fee, so it is not necessary to store this in the table.