

CptS - 451 Introduction to Database Systems Spring 2021

Project Description

In your semester long CptS 451 course project you would develop a data search application for “Yelp.com’s business review data”. The emphasis would be on the database infrastructure of the application.

Learner Objectives:

At the conclusion of this assignment you will gain experience in:

- ✓ Database modeling and design
- ✓ Populating the database with large datasets
- ✓ Querying large databases
- ✓ Optimizing query performance through indexes
- ✓ JSON parsing
- ✓ Database Application Development

Overview:

In 2013, Yelp.com has announced the “Yelp Dataset Challenge” and invited students to use this data in an innovative way and break ground in research. In your project you would query this dataset to extract useful information for local businesses and individual users.

The Yelp data is available in JSON format. The original Yelp dataset includes and **1.32M** tips by **1.97M** users for **209.39K** businesses from United States, Canada, UK, and Germany. (<https://www.yelp.com/dataset>) In your project you will use a smaller dataset that your instructor created. This simplified dataset includes only **19,983** businesses, **189,298** users, **287,288** tips written for those businesses, and **3,786,310** check-ins to those businesses.

You will be given sample (Python) code to parse some of the Yelp JSON files (available on Canvas). The Yelp JSON files that you will use in this project are available at the instructor’s website at:

https://eecs.wsu.edu/~arslanay/CptS451/project/yelp_dataset/yelp_CptS451_2020.zip

(Note: Please make sure to use the dataset available on the above link, not the one from the Yelp.com website)

See Appendix-B for an overview of the Yelp Academic Dataset.

Requirements:

You will develop a target application which runs queries on the Yelp data and extracts useful information. The primary users for this application will be potential customers seeking for businesses.

Using this application the users can gather information about:

- the businesses in a particular state, city, and/or zipcode,
- the businesses that belong to certain categories,
- detailed information about a business,
- ratings and popularity of businesses

You may design your application either as a standalone or a web-based application.

A detailed description of the application and example screenshots are available in **Appendix-A**. In evaluating your work instructor will primarily focus on how you design your database and how efficiently you can search

the database. However, your GUI should provide the basic functionality for easy search of the business. Creativity is encouraged! Additional functionality will be considered for extra credit.

Submission Instructions:

You will submit the deliverables for milestones on **Canvas** (canvas.wsu.edu). For each milestone you will create a .zip files that contains all deliverables for that milestone, name the .zip files as `<yourteamname>_milestoneX.zip`, and submit it to the corresponding milestone dropbox on Canvas. Specific submission details for each milestone will be provided in milestone descriptions.

Below is a summary of the milestone tasks. We will post detailed descriptions of the milestones when they are assigned.

Project Milestones:

I. Milestone-0: (no submission required)

Download and install PostgreSQL Database Server. You may download the latest version from the link <https://www.postgresql.org/> (Download and install the latest version of PostgreSQL Core Distribution.)

II. Milestone-1:

1) Parse JSON Data:

Download the Yelp dataset from:

https://eecs.wsu.edu/~arslanay/CptS451/project/yelp_dataset/yelp_CptS451_2020.zip. Look at each JSON file and understand what information the JSON objects provide. Pay attention to the data items in JSON objects that you will need for your application.

Download the sample program from Canvas (*Project/ Sample JSON Parsing Code*). The sample code:

- reads JSON objects form a file and extracts certain key and value pairs from JSON objects, and
- writes the extracted data into a text file.

Please note that the sample code includes examples of extracting simple key, value pairs from business JSON objects. In a JSON object the key value can be another JSON object (for example: categories and attributes in business data), therefore you need to recursively parse those objects until you extract all data stored in JSON objects. You will write the code for parsing business, user, tip, and check-in JSON objects.

- 2) i) Design a database schema that models the database for the described application scenario in Appendix-A and provide the ER diagram for your database design. Your schema should be precise but complete. It should be designed in such a way that all queries/data retrievals on/from the database run efficiently and effectively. In Milestone2 you will revise your ER model.

ii) Translate your ER model into relations and produce DDL SQL statements for creating the corresponding tables in a relational DBMS. Note the constraints, including key constraints, referential integrity constraints, not NULL constraints, etc. needed for the relational schema to capture and enforce the semantics of your ER design.

- 2) Build a very simple database application (either web or standalone) which runs simple queries on the given simple database. The goal of this exercise is to get you started in database programming early on.

The instructor will provide a video which explains how to establish connectivity with PostgreSQL in C# using Npgsql. Instructor will provide the queries you need to run on your table (see Milestone 1 specification).

Milestone-1 Deliverables:

1. (25%) Source code for parsing all JSON data. Only submit your source code, not the data files.
2. (40%) The E-R diagram and relations (CREATE TABLE statements) for your database design. To create your ER diagram, I suggest you use draw.io tool (<https://www.draw.io/>). You may also use your favorite drawing tool (e.g., Visio, Word, PowerPoint). Should be submitted in .pdf format. Name the diagram "<your-team-name>_ER_v1.pdf" and the SQL statements "<your-team-name>_schema.sql".
3. (35%) Source code for your application. Only submit your source code, not the data files.

III. Milestone-2:

- 1) Revise your database schema (ER model and relations).
- 2) Populate your database with the Yelp data. Generate INSERT statements for your tables and run those to insert data into your DB. You will also write and additional scripts to update the information stored in your database.

Write triggers to ensure the validity and consistency of the information stored in your database. Details will be available in Milestone2 specification.

- 3) Build the alpha-prototype of your application.

Milestone-2 Deliverables:

(Weights of the deliverables are TBA)

1. The revised E-R diagram. **Should be submitted in .pdf format.** Name this file "<your-team-name>_ER_v2.pdf"
2. SQL script file containing all SQL statements (i.e., CEATE TABLE statements, UPDATE statements, and TRIGGERS) . Name this file "<your-team-name>_SQL.sql"
3. Source code for parsing/inserting Yelp data into the database.
4. Alpha version of your Yelp application.

Check the Milestone2 specification for additional deliverables.

You will demonstrate your Milestone2 to the instructor and the TA.

IV. Milestone-3:

In this milestone you will complete your Yelp application A detailed description of the application requirements is provided in Appendix-A.

Milestone-3 Deliverables:

The source code of your application. **Please only upload your source code, not your DB files.**

You will demonstrate your final project to the instructor and the TA. The demonstration schedule will be announced in mid-April.

References:

1. Yelp Dataset Challenge, <https://www.yelp.com/dataset>
2. Samples for users of the Yelp Academic Database, <https://github.com/Yelp/dataset-examples>
3. Yelp Challenge, University of Washington Student Paper 1
<http://courses.cs.washington.edu/courses/cse544/13sp/final-projects/p08-fants.pdf>
4. Yelp Challenge, University of Washington Student Paper 2,
<http://courses.cs.washington.edu/courses/cse544/13sp/final-projects/p10-michelmj.pdf>

Appendix-A

Application Specification

The primary users for this application will be potential customers seeking for businesses. Using this application the users can gather information about:

- the businesses in a particular state, city, and/or zipcode,
- the businesses that belong to certain categories,
- detailed information about businesses,
- ratings and popularity of businesses,
- the businesses that their friends visited and reviewed, etc.

You may design your application either as a standalone or a web-based application. Below you will find screenshots to help you visualize the required functionality.

The application will have 2 main windows:

A. User Information:

Users can view their own information, the list of their own friends, the latest tips that each friend has provided, etc.

Use Case:

1. The user enters their name and chooses their own user id (among the users who has the same name). The system displays the following for the selected user:
 - user's profile information (including, their *name*, *average stars*, *the date he/she joined yelp*, *number of fans*, *average stars*, *count of votes*, *total tip count*, *total tip likes*, and *user's location (lat/long coordinates)*).
 - user's friends (name and star rating of each friend and the date he/she yelps since)
 - the latest tip that each of those friends posted. Note that this is different than the list of most recent reviews by friends. You need to return the latest tip that each friend has posted. (See Figure-1 for the properties you should include for each friend tip.)

Note: The user's location (lat/long) , total tip count, and total tip likes are not part of the JSON data, but are calculated/collected and stored in the databases as the application is used. The latitude and longitude coordinates are entered by the user in the "User" tab. "Tip count" is the count of tips user provided for businesses and the "total tip like" is the sum of all like counts for user's tips.

User Information

Set Current User: Lee Ann

User Information:

Name: Lee Ann

Stars: 3.54 Fans: 10

Yelping Since: 11/12/2011

Votes: funny: 101 cool: 135 useful: 307

Tip Count: 323

Total Tip Likes: 24475

Location: lat: 33.50730133056 long: -112.0380020141

Friends

Name	Total Likes	Avg Stars	Yelping Since
Kevin	601	3.96	3/21/2012 11:28:56 PM
gilbert	935	1.93	2/28/2011 6:45:22 PM
Rich	0	3.31	5/9/2012 11:32:28 PM
Shellie	7010	4.06	5/6/2012 2:07:27 AM
Taylor	12711	3.84	6/14/2011 11:56:37 PM
Georgie	9701	3.68	7/2/2010 4:21:38 PM
April	21995	3.79	8/17/2011 11:38:09 PM
Liz	13121	3.82	10/19/2011 1:47:26 PM
Kate	7141	4.08	7/25/2012 3:47:08 PM
Christina	874	4.06	5/31/2013 8:43:16 AM
C.	96	4.04	7/18/2012 12:45:24 AM
Christine	543	4.09	6/13/2010 5:58:41 PM
Erica	68	3.39	6/12/2007 7:54:55 PM
Jennifer	2443	4.12	4/9/2011 11:10:01 PM
Christina	0	3.79	12/7/2006 3:25:18 PM
Chris	4500	3.81	7/18/2008 9:17:35 PM
Had	6677	3.98	10/6/2010 7:09:22 PM
Matthew	4759	3.79	11/12/2011 2:18:35 AM
Matthew	0	3.8	2/15/2008 10:25:12 PM

Latest tips of my friends?

User Name	Business	City	Text	Date
Chris	Tamales Guadalajara	Tolleson	Good find, good food, great portion size and bet reasonable prices	11/12/2016 8:11
Sarah	Kokobelli Bagel Cafe	Mesa	Repeatedly out of certain bagels very early in the day so love this place but it gets annoying	1/20/2018 5:30
Thomas	Hula's Modern Tiki	Phoenix	Great decor and spotless everywhere! Very impressed	1/29/2012 11:02
April	Thallicious	Chandler	Whatever you do, don't order, sit down, and then browse through Yelp pics of this place. I almost walked the fuck it after seeing TWO roach pictures. Wtf. I very cautiously ate my food but I was cringing the entire time. Literally ate and ran out of there!	11/9/2017 7:30
Liz	Elevate Coffee	Phoenix	Great live band tonight. Saturday nights from 7:30-9:30 enjoy some local talent for free.	4/19/2015 5:09
Matthew	Koreatown	Mesa	Always check the specials (chalkboard at the entrance, or ask your server). Koreatown has recently added Korean Fried Chicken and its the best you'll find. Original, Garlic and Spicy Or you can order half n half.	12/4/2017 12:4
Michele	On The Border Mexican Grill & Cantina	Mesa	Great deal \$8.99 all u can eat enchiladas n yummy too!	3/20/2013 2:00
Christine	Macayo's Mexican Table	Mesa	Newly remodeled and lookin good!	1/4/2016 11:46
Erica	Postino Annex	Tempe	Only 2/3 full and a 40 minute wait? Need to work out the staffing situation...	10/17/2014 8:35
Lindsey	Brat Haus	Scottsdale	This is by far one of my favorite patios in Old Town! Great food and drinks too!	1/18/2015 1:43
Christie	Schlotzsky's	Mesa	They get very busy about 12-12:30!	11/18/2013 7:3
Taylor	Oscar's Pier 83	Glendale	Closed during renovation. Reopens Nov. 1 2018. Looking forward to great seafood then!!!	10/21/2018 10:
Kevin	Cornerstone Christian Fellowship	Chandler	July 28th Michael Irving will be here!	6/23/2013 7:01
Shellie	Buddyz A Chicago Pizzeria	Gilbert	Once we tried Buddyz, we haven't been to another pizza spot! It's THAT good!	9/22/2017 10:51
gilbert	China City Super Buffet	Mesa	I hate this place now, really has gone down in quality	7/19/2017 4:51
Scott	Giordanos - Arrowhead	Peoria	It only took an hour to get our stuffed pizza. According to our server the wait times were over 2 hours on Friday.	12/25/2016 1:0
Kari	The Madison Improvement Club	Phoenix	Spin yoga meditate	1/26/2013 4:46
Jennifer	Kneaders Bakery & Cafe	Surprise	This location just opened on Thursday, October 2nd, so it's a bit crowded right now at lunch! You can bypass the regular line if you are just buying bread or a pastry or even holiday gifts at the first register by the door.	10/7/2014 5:50
Marshall	Angry Crab Peoria	Peoria	Great place! I'm happy I have one near where I live! Service, eh...	11/21/2015 11:

Figure 1 – User information window

B. Business Search:

Users can search for businesses which are within a certain state, city, and zip and which belong to the selected categories and/or attributes. The application allows users to retrieve and display various information about a selected business.

Use Cases:

1. User selects a state, city, and zipcode. When search button is pressed, the businesses in that state/city/zipcode are retrieved (see Figure-2). The following information is displayed for each business in the search result.

- Business name
- Address, city, state
- Distance to user's location.
- Business rating (stars)
- Number of tips provided for the business
- Number of check-ins to the business

Notes:

1. To get the number of tips and check-ins, you should
 - (i) query the tips table to calculate the number of tips;
 - (ii) query the check-in table to calculate the number of check-ins for each business; and
 - (iii) update those attribute values in the business table.
2. Distance: You should calculate the earth distance between the lat/long coordinated of the user and each business using the Haversine Formula. **The distance calculation should be implemented in SQL (as a function).**

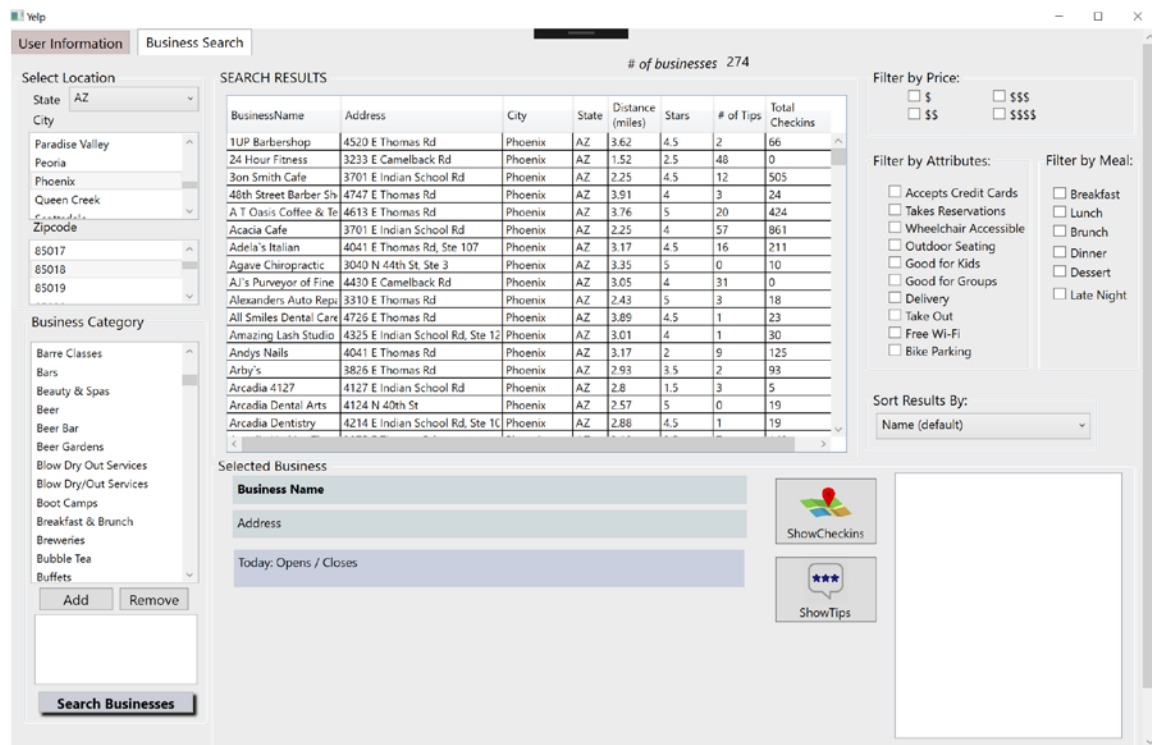


Figure 2 – Searching for the businesses in a Phoenix, AZ, 85018

3. The user might refine the results by specifying one or more business categories. The search will return the businesses which belong **to the ALL categories** specified by the user (i.e., AND condition) Note that, the more categories are selected the more restrictive the search will be. (see Figures-3 and -4)

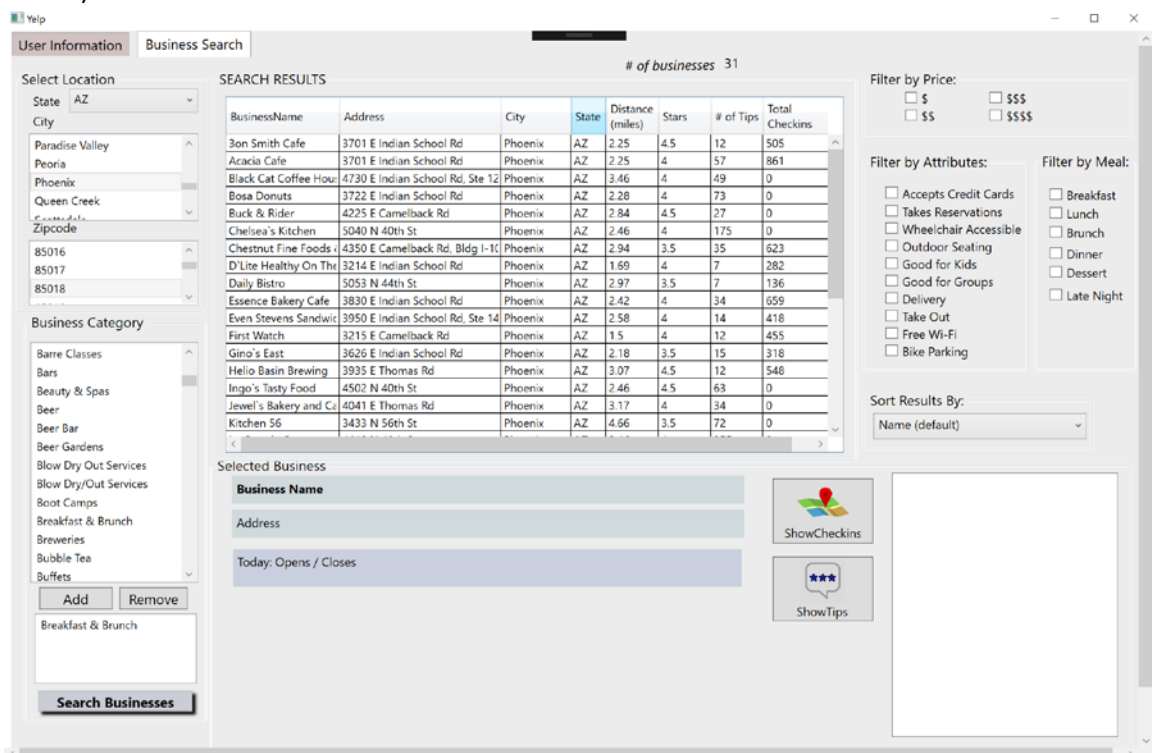
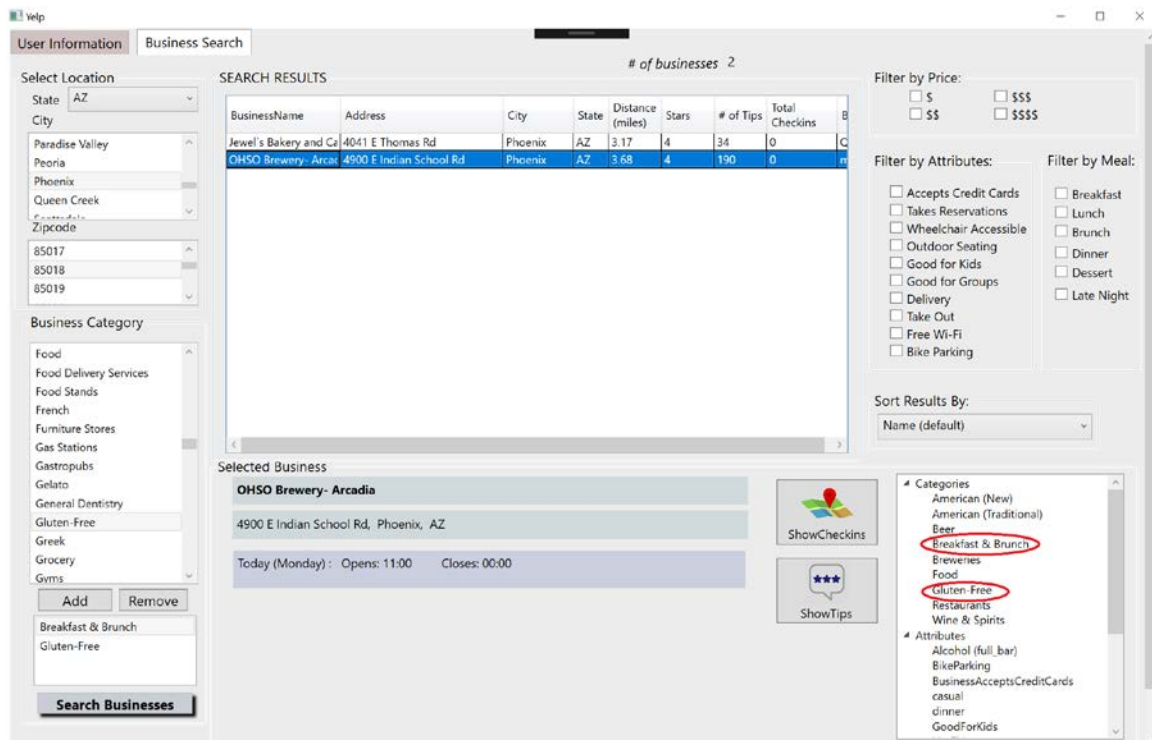
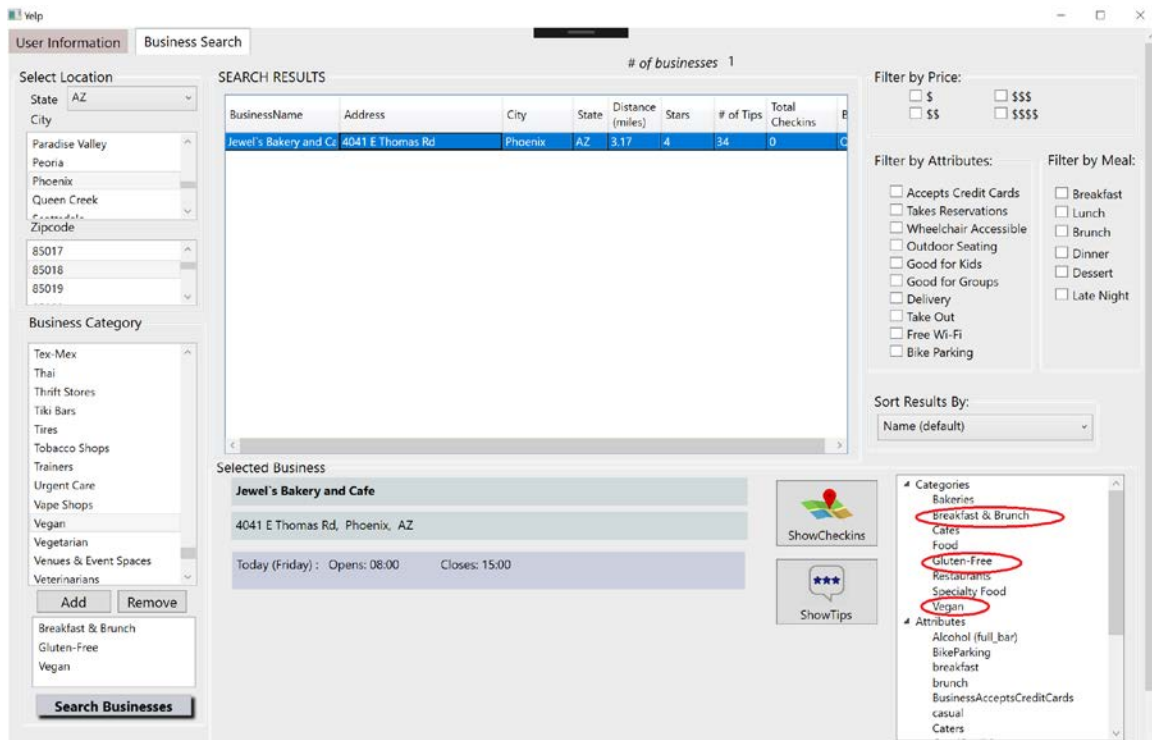


Figure 3 – Searching for the businesses with category 'Breakfast & Brunch' in Phoenix, AZ, 85018. The search result includes 31 businesses.



(a) Searching for the businesses with categories 'Breakfast & Brunch' and 'Gluten-Free' in Phoenix, AZ, 85018. The search result includes 2 businesses.



(b) Searching for the businesses with categories 'Breakfast & Brunch', 'Gluten-Free', and 'Vegan' in Phoenix, AZ, 85018. The search result includes a single business.

Figure 4 – Searching for the businesses that match multiple categories.

4. When the user selects a business in the search results, the following are displayed for the selected business:
 - i. the name, address of the business
 - ii. the open/close times of the business for the current day of the week;
 - iii. the categories and attributes of the business
 (see Figures 4, 5, and 6 for examples)
5. The user may also refine results by specifying various attributes including:
 - i. Price range (1 to 4) – (see RestaurantsPriceRange2 attribute)
 - ii. Accepts credit cards (see BusinessAcceptsCreditCards attribute)
 - iii. Takes reservations (see RestaurantsReservations attribute)
 - iv. Wheelchair accessible (see WheelchairAccessible attribute)
 - v. Outdoor seating (see OutdoorSeating attribute)
 - vi. Good for kids (see GoodForKids attribute)
 - vii. Good for groups (see RestaurantsGoodForGroups attribute)
 - viii. Delivery (see RestaurantsDelivery attribute)
 - ix. Take out (see RestaurantsTakeOut attribute)
 - x. Wifi (free Wifi only) – (see WiFi attribute)
 - xi. Bike parking (see BikeParking attribute)
 - xii. Meals (breakfast, brunch, lunch, dinner, desert, and latenight attributes)

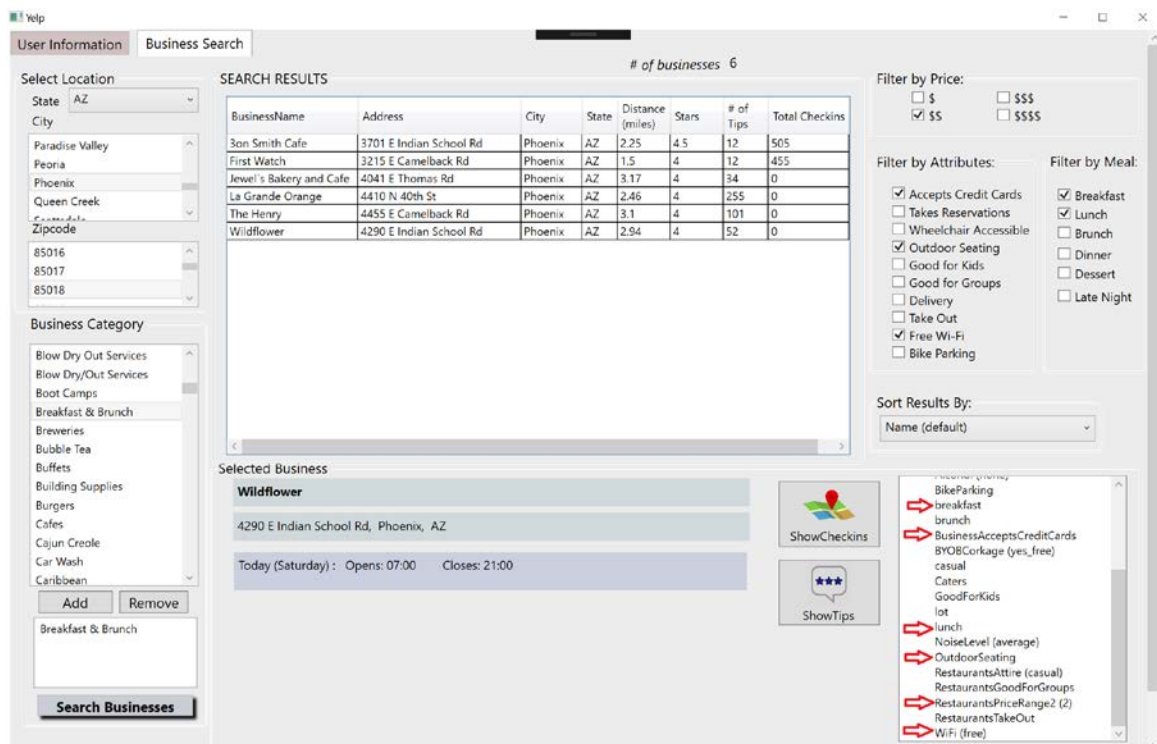


Figure 5 – Searching for the businesses with category ‘Breakfast & Brunch’ in Phoenix, AZ, 85018 which has ‘price range 2’, accepts credit cards, has outdoor seating, has free Wi-Fi, and serves breakfast and lunch.

The details of the selected business (*Wildflower*) are displayed.

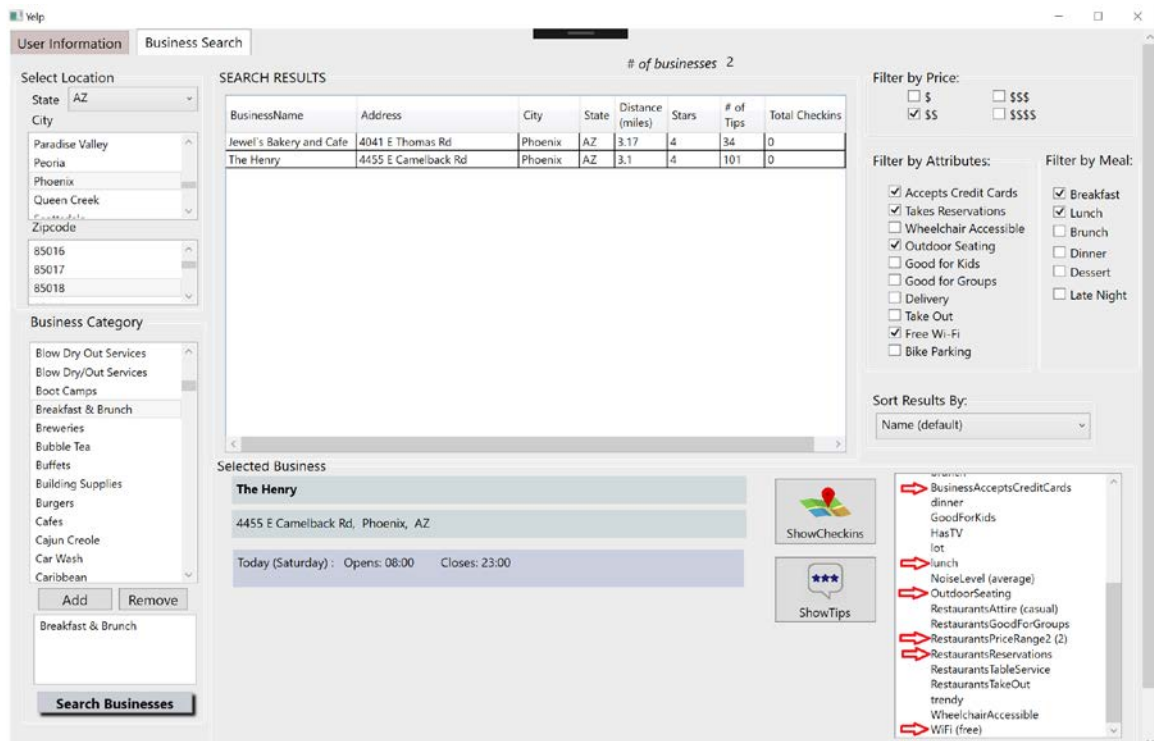


Figure 6 – Searching for the businesses with category ‘Breakfast & Brunch’ in Phoenix, AZ, 85018 which has ‘price range 2’, accepts credit cards, takes reservations, has outdoor seating, has free Wi-Fi, and serves breakfast and lunch. Compare this to Figure-5: adding an additional search attribute (takes reservations) narrows down the results.

6. The user may sort the results based on the following attribute values.
 - a. Business name (default sort order)
 - b. Highest rating (stars)
 - c. Most number of tips
 - d. Most check-ins
 - e. Nearest

In Figure-7 the business search results are sorted by number of check-ins (in descending order). All sorting should be done in the SQL query (you can’t use the sorting features of the data-grid.)

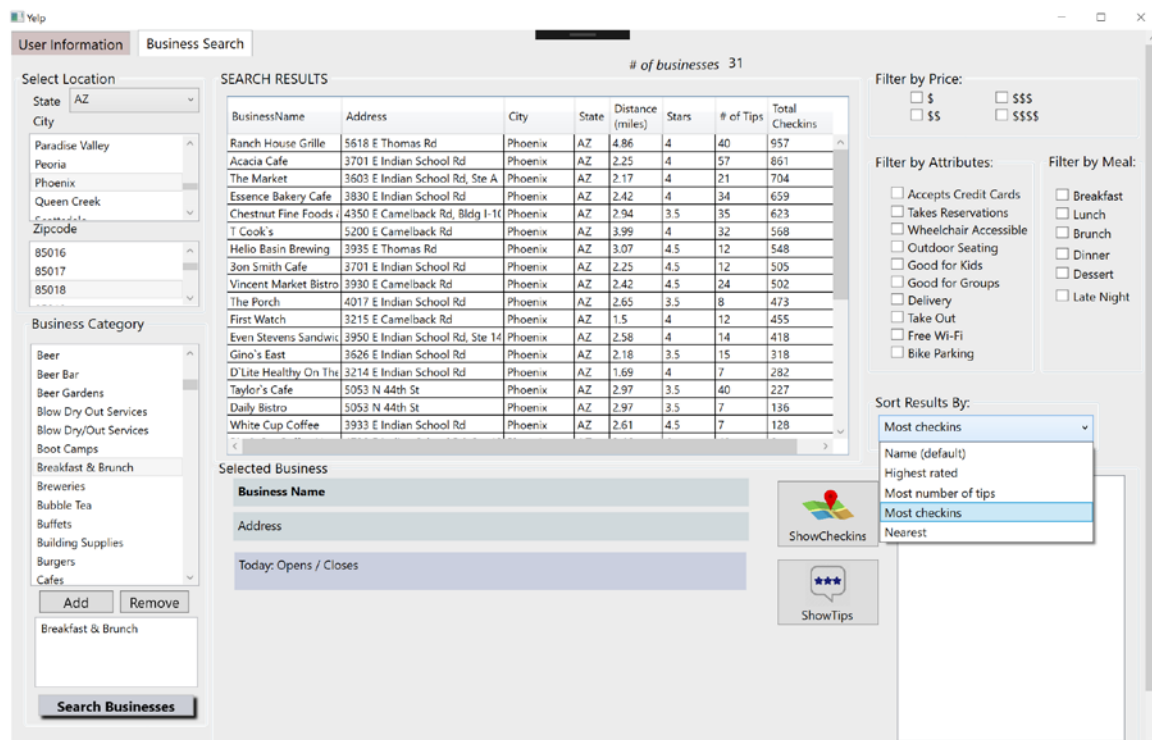


Figure 7 – Business search results can be sorted by various values. In this figure, results are sorted by the number of check-ins (in descending order).

7. The user may select a certain business in the search results (by simply clicking on a business) and
 - (a) display the tips provided for the business;
 - (b) display the tips that users friends have provided;
 - (c) add a new tip for the selected business. (Note: We assume that the tip is posted by the user selected in the 'User Information' tab. The timestamp of the new tip should be the time that the tip is added.) (See Figure-8)
 - (d) like a tip (in Figure-8 the user liked their own tip and the like count was incremented to 1)

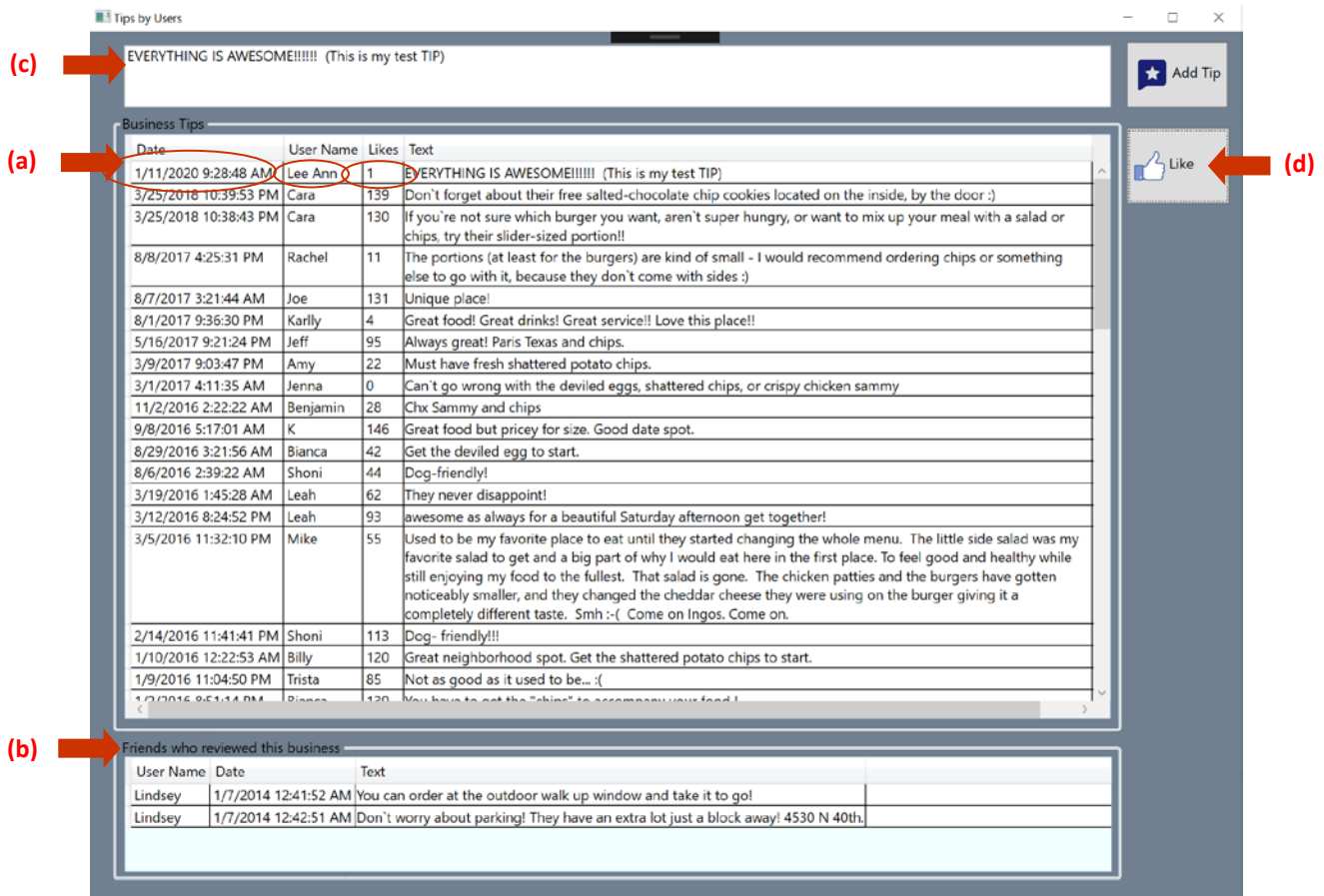


Figure 8 – The tips for the business ‘Ingo’s Tasty Foods’ in Phoenix, AZ, 85018. The most recent tips appear at the top. For each tip, you should display the name of the user who provided the tip, the date tip is provided, the likes for the tip, and the tip text.

8. The user may select a certain business in the search results (by simply clicking on a business) and
 - (a) display the total number of check-ins for each month as a chart;
 - (b) check-in to the selected business.

Please see Figure-10 for an example

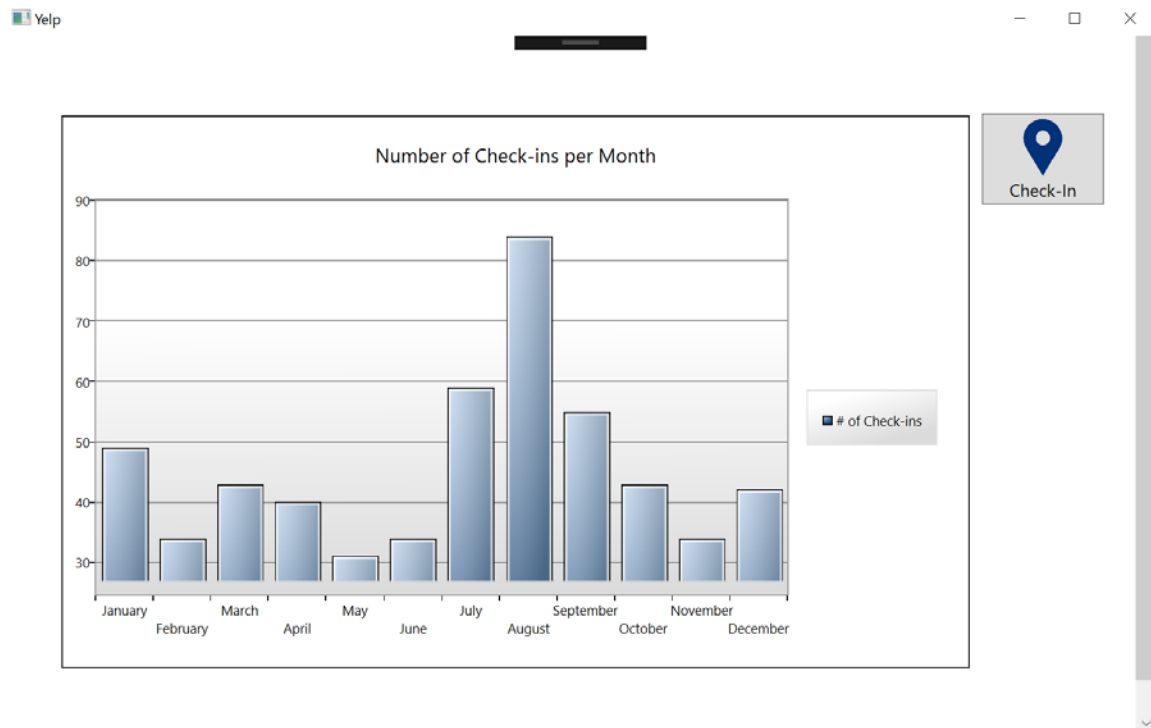
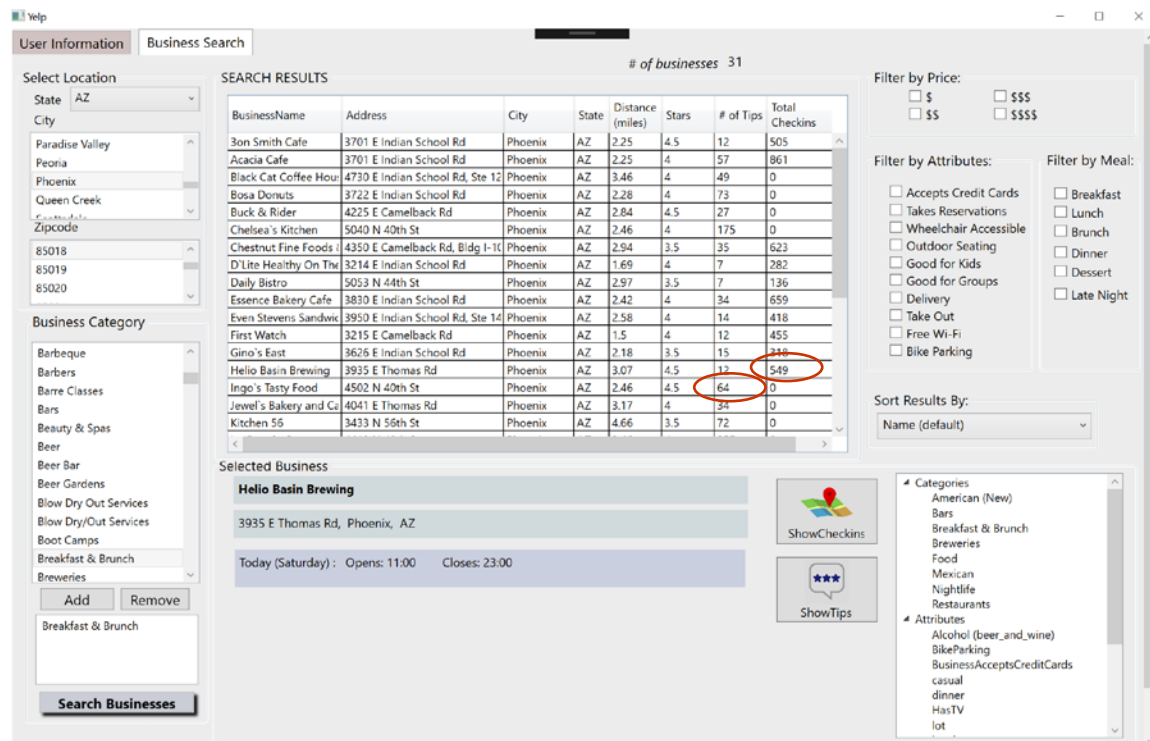


Figure 8 – The check-in chart for ‘Helio Basin Brewing’ in Phoenix, AZ, 85018. The graph shows the number of check-ins per month. When the user checks-in, the graph will be updated (the count for current month will be increased by one).

9. After a check-in to a business, the check-in count of the business should be updated (by a trigger). Similarly, when a new tip is provided for a business, the total tip count of the business should be incremented. (See Figure-9)
You will write the triggers that increment the check-in and tip count in milestone-2.



SEARCH RESULTS # of businesses 31

BusinessName	Address	City	State	Distance (miles)	Stars	# of Tips	Total Checkins
Bon Smith Cafe	3701 E Indian School Rd	Phoenix	AZ	2.25	4.5	12	505
Acacia Cafe	3701 E Indian School Rd	Phoenix	AZ	2.25	4	57	861
Black Cat Coffee House	4730 E Indian School Rd, Ste 12	Phoenix	AZ	3.46	4	49	0
Bosa Donuts	3722 E Indian School Rd	Phoenix	AZ	2.28	4	73	0
Buck & Rider	4225 E Camelback Rd	Phoenix	AZ	2.84	4.5	27	0
Chelsea's Kitchen	5040 N 40th St	Phoenix	AZ	2.46	4	175	0
Chestnut Fine Foods	4350 E Camelback Rd, Bldg 1-10	Phoenix	AZ	2.94	3.5	35	623
D Lite Healthy On The Go	3214 E Indian School Rd	Phoenix	AZ	1.69	4	7	282
Daily Bistro	5053 N 44th St	Phoenix	AZ	2.97	3.5	7	136
Essence Bakery Cafe	3830 E Indian School Rd	Phoenix	AZ	2.42	4	34	659
Even Stevens Sandwiches	3950 E Indian School Rd, Ste 14	Phoenix	AZ	2.58	4	14	418
First Watch	3215 E Camelback Rd	Phoenix	AZ	1.5	4	12	455
Gino's East	3626 E Indian School Rd	Phoenix	AZ	2.18	3.5	15	248
Helio Basin Brewing	3935 E Thomas Rd	Phoenix	AZ	3.07	4.5	12	549
Ingo's Tasty Food	4502 N 40th St	Phoenix	AZ	2.46	4.5	64	0
Jewell's Bakery and Cafe	4041 E Thomas Rd	Phoenix	AZ	3.17	4	34	0
Kitchen 55	3433 N 56th St	Phoenix	AZ	4.66	3.5	72	0

Selected Business

Helio Basin Brewing

3935 E Thomas Rd, Phoenix, AZ

Today (Saturday) : Opens: 11:00 Closes: 23:00

Filter by Price:

☐ \$ ☐ \$\$\$

☐ \$\$ ☐ \$\$\$\$

Filter by Attributes:

☐ Accepts Credit Cards ☐ Breakfast

☐ Takes Reservations ☐ Lunch

☐ Wheelchair Accessible ☐ Brunch

☐ Outdoor Seating ☐ Dinner

☐ Good for Kids ☐ Dessert

☐ Good for Groups ☐ Late Night

☐ Delivery

☐ Take Out

☐ Free Wi-Fi

☐ Bike Parking

Sort Results By:

Name (default)

Categories

- American (New)
- Bars
- Breakfast & Brunch
- Breweries
- Food
- Mexican
- Nightlife
- Restaurants

Attributes

- Alcohol (beer_and_wine)
- BikeParking
- BusinessAcceptsCreditCards
- casual
- dinner
- HasTV
- lot

Figure 8 – After check-in, the check-in and tip counts are automatically updated. Compare the circled counts to those in Figure-3. (Check-in count increased from 63 to 64, and the tip count increased from 548 to 549).

Please note that all data displayed on the GUI should be kept in the database and should be retrieved from it when needed. You are not allowed to create internal data structures to store data.

You may design your application either as a standalone or a web-based application.

Appendix-B

Yelp's Academic Dataset

Yelp has made available a dataset which contains user reviews for **192.6K** businesses from United States, Canada, UK, and Germany. The purpose was to provide a real-world data set to promote research in various areas of research. The dataset includes 6 types of data objects: *business*, *review*, *user*, *tip*, *check-in*, and *photos*. Every object contains a 'type' field, which tells whether it is a *business*, a *user*, or a *review*. *Business* objects contain basic information about local businesses. *Review* objects contain the details of the reviews by users for the businesses. *Review*'s `user_id` associates the reviews with the *user* objects. Similarly, *review*'s `business_id` associates each review with the *businesses*.

Detailed description of the data objects is available at:

<https://www.yelp.com/dataset/documentation/main>

In your project, you will only parse *business*, *user*, *tip*, and *check-in* objects.

Usage of this dataset is governed by the Academic Dataset Terms of Use.