

# CS 513 – FINAL PROJECT – PHASE II

TEAM 169 - CHIRANJEEVI (CHIRU) KOILOTH, CHUN CHIEH (WILLY) CHANG, AARON BARRIE

## TEAM COMPOSITION

1. Chiranjeevi (Chiru) Koiloth (koiloth2@illinois.edu)
2. Chun Chieh (Willy) Chang (cc132@illinois.edu)
3. Aaron Barrie (abarrie2@illinois.edu)

## USE CASE DESCRIPTION

The following section reiterates the Main Use Case described in the team's Phase I submission. As the U1 definition is built on top of the U0 definition, both are included below.

### U0 – ZERO CLEANING

The team wishes to build a web-based application to deliver the contents of this dataset for an embedded system where the image references indicated in the MenuPage entity of this dataset are unavailable for distribution. Each menu can be identified from a navigable list that provides search and filtering based on:

- Menu Name
- Location
- Special Occasion
- Date of Publication
- Keyword Search for Notes

Clicking on a menu will open a navigable document that has been reconstructed based on the information from the Menu Page, Menu Item, and Dish entities.

The Menu Page entity will provide the Menu Page ID and page dimensions, which will be scaled to fit the current browser window.

All Menu Items associated with a given Menu Page will be plotted using a textbox with an origin point located at the relative X/Y described by xpos/ypos. This textbox will include:

- Dish Name (from the Dish entity)
- Dish Price (from the MenuItem entity)
- Dish Alternate Price (from the MenuItem entity, when available)
- Dish Prices will be decorated with a Currency Symbol from the Menu entity

Additionally, this application will support the ability to filter the list of Menu's based on a search for Dish names. If a menu is selected when this filter is used, then the application will automatically navigate the user to the page on which that Dish is found, with a highlight placed around the location of the Menu Item textbox on that page.

Any data points presented in the original data set will be done as-is. Any data points that are inaccessible due to integrity constraint violations will result in a "Missing Page" notice to the end user in the event of a missing page

reference, or “Missing Data” notice to the end user in the event of a missing dish reference. When currency markers are unavailable, the prices will be listed as-is without currency designation.

---

## U1 – MAIN USE CASE

The team wishes to extend the functionality of the application described in U0 by adding the ability to offer additional Organization- and Dish-based analysis.

To offer Organization-based analysis, the data points for Name/Sponsor on the Menu entity need to be appropriately clustered to accurately present a comprehensive view of how the menu characteristics for a single organization change over time.

To offer Dish-based analysis, the data points for the Dish Names from Dish.csv and prices from MenuItem.csv need to be cleaned to provide an accurate view for how the popularity of specific dishes changed over time, and how price trends for dishes change over time.

## 1. DESCRIPTION OF DATA CLEANING PERFORMED

The general approach to the data cleanup process ended up being a three-step process when viewed from a high-level.

1. Load each source file in SQLite and perform baseline IC data consistency checks.
  - a. Inputs for this step are the four raw data files acquired from the course website.
  - b. Outputs for this step are IC exception metrics across all four input files.
2. Load each source file in OpenRefine and perform internal file cleanup.
  - a. Inputs for this step are the four raw data files acquired from the course website.
  - b. This work involves performing cleanup operations that do not require any awareness of relational consistency between each of the input files.
  - c. Outputs for this step are four cleaned data files corresponding to the four raw data files provided at input.
3. Load each OpenRefine cleaned file in SQLite and perform relational data cleanup.
  - a. Inputs for this step are the four cleaned data files produced as an output of step 1.
  - b. This work involves performing cleanup operations that are primarily focused on relational consistency between each of the four input data sets.
  - c. Outputs for this step are four final-state cleaned data files corresponding to the four input files.
  - d. Additional outputs for this step are IC exception metrics across all four output files.

The rationale for this approach is straightforward: OpenRefine excels at operations that only require awareness of operations within the single data file being inspected. SQLite excels at operations that require awareness of relationships between each of the input files in the source data set.

## OPENREFINE ACTIONS

---

### MENU

1. Sponsor, event, venue, place, occasion, and location
  1. Using text transformations in Edit cells
    1. Trim leading and trailing whitespaces.

**Text transform on 14 cells in column sponsor: value.trim()**  
Undo

2. Collapse consecutive whitespaces.

**Text transform on 127 cells in column sponsor: value.replace(/[p{Zs}\s]+/, ' ') Undo**

3. Convert the whole column to uppercase.

**Text transform on 8,535 cells in column sponsor: value.toUpperCase() Undo**

4. Remove the special characters '~!@#%&\*()[]\_+V-=?:,;<>' using General Refine Expression Language(GREL)

**Custom text transform on column sponsor**

Expression Language General Refine Expression Language (GREL)

`value.replace(/[~!@#%&*()[]_+V-=?:,;<>']/, "")` No syntax error.

**Preview** History Starred Help

row	value	value.replace(/[~!@#%&*()[]_+V-=?:,;<>']/, "...
1.	HOTEL EASTMAN	HOTEL EASTMAN
2.	REPUBLICAN HOUSE	REPUBLICAN HOUSE
3.	NORDDEUTSCHER LLOYD BREMEN	NORDDEUTSCHER LLOYD BREMEN
4.	NORDDEUTSCHER LLOYD BREMEN	NORDDEUTSCHER LLOYD BREMEN
5.	NORDDEUTSCHER LLOYD BREMEN	NORDDEUTSCHER LLOYD BREMEN
6.	CANADIAN PACIFIC RAILWAY COMPANY	CANADIAN PACIFIC RAILWAY COMPANY

On error ☒ keep original ☐ Re-transform up to  times until no change  
☐ set to blank  
☐ store error

OK Cancel

**Text transform on 3,564 cells in column sponsor: grel:value.replace(/[~!@#%&\*()[]\_+V-=?:,;<>']/, "") Undo**

2. Create a text facet and implement the cluster operation
  1. Using the key-collision method with different keying functions
    1. Fingerprint function

**Cluster and edit column "sponsor"**

Find groups of different cell values that might be other representations of the same thing. For example, "New York" and "new york" likely refer to the same concept and just differ by capitalization, and "Gödel" and "Godel" probably refer to the same person. [Find out more...](#)

Method Key collision Keying function Fingerprint 154 clusters found

Cluster size	Row count	Values in cluster
4	24	<ul style="list-style-type: none"> <li>RED STAR LINE ANTWERPEN NY (11 rows)</li> <li>RED STAR LINE ANTWERPEN NY (7 rows)</li> <li>RED STAR LINE ANTWERPEN NY (5 rows)</li> <li>RED STAR LINE ANTWERPEN NY</li> </ul>
3	29	<ul style="list-style-type: none"> <li>GRAMERCY PARK HOTEL (19 rows)</li> <li>HOTEL GRAMERCY PARK (9 rows)</li> <li>GRAMERCY PARK HOTEL HOTEL GRAMERCY PARK</li> </ul>
3	25	<ul style="list-style-type: none"> <li>COMPAGNIE GÉNÉRALE TRANSATLANTIQUE (17 rows)</li> <li>COMPAGNIE GÉNÉRALE TRANSATLANTIQUE (5 rows)</li> <li>COMPAGNIE GÉNÉRALE TRANSATLANTIQUE (3 rows)</li> </ul>
3	667	<ul style="list-style-type: none"> <li>NORDDEUTSCHER LLOYD BREMEN (635 rows)</li> <li>NORDDEUTSCHER LLOYD BREMEN (31 rows)</li> <li>BREMEN NORDDEUTSCHER LLOYD</li> </ul>
3	24	<ul style="list-style-type: none"> <li>HOTEL KNICKERBOCKER (22 rows)</li> <li>HOTEL KNICKERBOCKER</li> <li>KNICKERBOCKER HOTEL</li> </ul>

**# Choices in cluster** 2 — 4

**# Rows in cluster** 0 — 670

**Average length of choices** 3 — 93

**Length variance of choices** 0 — 9.43

Select all Deselect all Export clusters Merge selected & re-cluster Merge selected & Close Close

**Mass edit 2,933 cells in column sponsor Undo**

2. n-Gram fingerprint function with n-Gram size 2

**Cluster and edit column "sponsor"**

Find groups of different cell values that might be other representations of the same thing. For example, "New York" and "new york" likely refer to the same concept and just differ by capitalization, and "Gödel" and "Gode!" probably refer to the same person. [Find out more...](#)

Method: **Key collision**      Keying function: **n-Gram fingerprint**      n-Gram size: **2**      75 clusters found

Cluster size	Row count	Values in cluster	Merge?	New cell value
3	4	<ul style="list-style-type: none"> <li>RED STAR LINE S.S.FRIESLAND (2 rows)</li> <li>RED STAR LINE SS FRIESLAND</li> <li>RED STAR LINESS FRIESLAND</li> </ul>	<input checked="" type="checkbox"/>	RED STAR LIN
3	4	<ul style="list-style-type: none"> <li>U.S.S. RALEIGH (2 rows)</li> <li>U.S.S.RALEIGH</li> <li>U.S.S.S.RALEIGH</li> </ul>	<input checked="" type="checkbox"/>	U.S.S. RALEIGH
3	8	<ul style="list-style-type: none"> <li>HOFBRAU HAUS (5 rows)</li> <li>HOFBRAUHAUS (2 rows)</li> <li>HOF BRAU HAUS</li> </ul>	<input checked="" type="checkbox"/>	HOFBRAU HA
3	10	<ul style="list-style-type: none"> <li>NIPPON YUSEN KAISHA S.S.KOBE MARU (5 rows)</li> <li>NIPPON YUSEN KAISHA S.S. KOBE MARU (4 rows)</li> <li>NIPPON YUSEN KAISHA S.S. KOBE MARU</li> </ul>	<input checked="" type="checkbox"/>	NIPPON YUSE
2	58	<ul style="list-style-type: none"> <li>(57 rows)</li> <li>L</li> </ul>	<input checked="" type="checkbox"/>	
2	2	<ul style="list-style-type: none"> <li>SOCIETA LA PIEMONTESE</li> <li>SOCIETALA PIEMONTESE</li> </ul>	<input checked="" type="checkbox"/>	SOCIETA LA P

Select all   Deselect all   Export clusters   Merge selected & re-cluster   Merge selected & Close   Close

# Choices in cluster: 2 — 3  
# Rows in cluster: 0 — 790  
Average length of choices: 0 — 49  
Length variance of choices: 0 — 2

**Mass edit 1,787 cells in column sponsor** **Undo**

### 3. Metaphone3 function

**Cluster and edit column "sponsor"**

Find groups of different cell values that might be other representations of the same thing. For example, "New York" and "new york" likely refer to the same concept and just differ by capitalization, and "Gödel" and "Gode!" probably refer to the same person. [Find out more...](#)

Method: **Key collision**      Keying function: **Metaphone3**      389 clusters found

Cluster size	Row count	Values in cluster
12	729	<ul style="list-style-type: none"> <li>NORDEUTSCHER LLOYD BREMEN (667 rows)</li> <li>NORDEUTSCHER LLOYD BREMEN (29 rows)</li> <li>NORDEUTSCHER LLOYD BREMEN (21 rows)</li> <li>NORDEUTSCHER LLOYD BREMEN S.S.BARBAROSSA (3 rows)</li> <li>NORDEUTSCHER LLOYD BREMEN LINE (2 rows)</li> <li>NORDEUTSCHER LLOYD PRINZ FRIEDRICH WILHELM</li> <li>NORDEUTSCHER LLOYD BREMAN</li> <li>NORDEUTSCHER LLOYD BREMEN BARBAROSSA</li> <li>NORDEUTSCHER LLOYD BREMENAMERIKA</li> <li>NORDEUTSCHER LLOYD BREMENAMERIKA</li> <li>NORDEUTSCHER LLOYD BREMEN</li> <li>NORDEUTSCHER LLOYD BREMEN ON BOARD S.S. GEORGE WASHINGTON.</li> </ul>
12	246	<ul style="list-style-type: none"> <li>OCEANIC STEAMSHIP COMPANY (114 rows)</li> <li>OCEANIC STEAMSHIP CO. (49 rows)</li> <li>OCEANIC STEAMSHIP CO. VENTURA (33 rows)</li> <li>OCEANIC STEAMSHIP CO. SONOMA (11 rows)</li> <li>OCEANIC STEAMSHIP CO. S.S.ZEALANDIA (9 rows)</li> <li>OCEANIC STEAMSHIP COMPANY SIERRA (7 rows)</li> <li>OCEANIC STEAMSHIP COMPANY SONOMA (7 rows)</li> <li>OCEANIC STEAMSHIP (5 rows)</li> <li>OCEANIC STEAMSHIP CO. SIERRA (5 rows)</li> <li>OCEANIC STEAMSHIP COMPANY S.S.ZEALANDER (3 rows)</li> <li>OCEANIC STEAMSHIP COMPANY VENTURA (2 rows)</li> <li>OCEANIC STEAMSHIP COMPANY S.S.ZEALANDIA</li> </ul>

Select all   Deselect all   Export clusters   Merge selected & re-cluster   Merge selected & Close   Close

# Choices in cluster: 2 — 12  
# Rows in cluster: 0 — 810  
Average length of choices: 3 — 88  
Length variance of choices: 0 — 25

**Mass edit 5,118 cells in column sponsor**

### 4. Beider-Morse functions

**Cluster and edit column "sponsor"**

Find groups of different cell values that might be other representations of the same thing. For example, "New York" and "new york" likely refer to the same concept and just differ by capitalization, and "Gödel" and "Gode!" probably refer to the same person. [Find out more...](#)

Method: **Key collision**      Keying function: **Beider-Morse**      6 clusters found

Cluster size	Row count	Values in cluster	Merge?	New cell value
2	5	<ul style="list-style-type: none"> <li>RATHSKELLER (4 rows)</li> <li>RATSKELLER</li> </ul>	<input type="checkbox"/>	RATHSKELLER
2	6	<ul style="list-style-type: none"> <li>DENNETTS (5 rows)</li> <li>DINNER</li> </ul>	<input type="checkbox"/>	DENNETTS
2	219	<ul style="list-style-type: none"> <li>NIPPON YUSEN KAISHA (215 rows)</li> <li>NIPPON YUSEN KAISA (4 rows)</li> </ul>	<input type="checkbox"/>	NIPPON YUSEN KAISHA
2	5	<ul style="list-style-type: none"> <li>PROMENADE CAFÉ (3 rows)</li> <li>PROMENADE CAFÉS (2 rows)</li> </ul>	<input type="checkbox"/>	PROMENADE CAFÉ
2	739	<ul style="list-style-type: none"> <li>NORDEUTSCHER LLOYD BREMEN (729 rows)</li> <li>NORDEUTSCHERR LLOYD BREMEN (10 rows)</li> </ul>	<input type="checkbox"/>	NORDEUTSCHER LLOYD
2	3	<ul style="list-style-type: none"> <li>THE NEW HAVEN RAILROAD (2 rows)</li> <li>THE NEW HAVEN R. R. RAILROAD</li> </ul>	<input type="checkbox"/>	THE NEW HAVEN RAILROA

Select all   Deselect all   Export clusters   Merge selected & re-cluster   Merge selected & Close   Close

# Rows in cluster: 0 — 740  
Average length of choices: 7 — 27  
Length variance of choices: 0.5 — 3

**Mass edit 977 cells in column sponsor** **Undo**

2. Using Nearest Neighbor method with different distance functions
  1. Levenshtein distance function with radius 1.0 and block chars 6

**Cluster and edit column "sponsor"**

Find groups of different cell values that might be other representations of the same thing. For example, "New York" and "new york" likely refer to the same concept and just differ by capitalization, and "Gödel" and "Godel" probably refer to the same person. [Find out more...](#)

Method: **Nearest neighbor** Distance function: **Levenshtein** Radius: **1.0** Block chars: **6** **44 clusters found**

Cluster size	Row count	Values in cluster	Merge?	New name
3	744	<ul style="list-style-type: none"> <li>NORDEUTSCHER LLOYD BREMEN (739 rows)</li> <li>NORDEUTSCHER LLOYD BREMEN (4 rows)</li> <li>NORDEUTSCHER LLOYDS BREMEN</li> </ul>	<input checked="" type="checkbox"/>	NORD
2	86	<ul style="list-style-type: none"> <li>THE WALDORF ASTORIA (85 rows)</li> <li>THE WALFORF ASTORIA</li> </ul>	<input checked="" type="checkbox"/>	THE W
2	3	<ul style="list-style-type: none"> <li>EL TOVAR HOTEL (2 rows)</li> <li>EL COVAR HOTEL</li> </ul>	<input checked="" type="checkbox"/>	EL TO
2	18	<ul style="list-style-type: none"> <li>GRAND HOTEL (17 rows)</li> <li>GRAD HOTEL</li> </ul>	<input checked="" type="checkbox"/>	GRAN
2	4	<ul style="list-style-type: none"> <li>WILLARDS HOTEL (3 rows)</li> <li>WILLARD HOTEL</li> </ul>	<input checked="" type="checkbox"/>	WILLA
2	3	<ul style="list-style-type: none"> <li>CLOVER CLUB OF BOSTON (2 rows)</li> <li>COVER CLUB OF BOSTON</li> </ul>	<input checked="" type="checkbox"/>	CLOV
2	14	<ul style="list-style-type: none"> <li>HOTEL ANDAMERICA LINE (13 rows)</li> </ul>	<input checked="" type="checkbox"/>	HOTEL

Select all Deselect all Export clusters Merge selected & re-cluster Merge selected & Close Close

**Mass edit 1,588 cells in column sponsor Undo**

2. PPM distance function with radius 1.0 and block chars 6

**Cluster and edit column "sponsor"**

Find groups of different cell values that might be other representations of the same thing. For example, "New York" and "new york" likely refer to the same concept and just differ by capitalization, and "Gödel" and "Godel" probably refer to the same person. [Find out more...](#)

Method: **Nearest neighbor** Distance function: **PPM** Radius: **1.0** Block chars: **6** **163 clusters found**

Cluster size	Row count	Values in cluster	Merge?	New name
3	9	<ul style="list-style-type: none"> <li>BROWN PALACE HOTEL (6 rows)</li> <li>THE BROWN PALACE HOTEL (2 rows)</li> <li>THE H.C. BROWN PALACE HOTEL</li> </ul>	<input checked="" type="checkbox"/>	BR
3	5	<ul style="list-style-type: none"> <li>YORK HOTEL (3 rows)</li> <li>DEER PARK HOTEL</li> <li>PARK HOTEL</li> </ul>	<input type="checkbox"/>	YC
3	23	<ul style="list-style-type: none"> <li>GRAND HOTEL (18 rows)</li> <li>GRANDON HOTEL (4 rows)</li> <li>STRAND HOTEL</li> </ul>	<input checked="" type="checkbox"/>	GR
3	6	<ul style="list-style-type: none"> <li>MICHIGAN SOCIETY OF THE SONS OF THE AMERICAN REVOLUTION (4 rows)</li> <li>MASS SOCIETY OF THE SONS OF THE AMERICAN REVOLUTION</li> <li>OREGON SOCIETY OF THE SONS OF THE AMERICAN REVOLUTION</li> </ul>	<input checked="" type="checkbox"/>	MIC
3	4	<ul style="list-style-type: none"> <li>THE NEW WILLARD HOTEL (2 rows)</li> <li>NEW WILLARD HOTEL</li> <li>THE WILLARD HOTEL</li> </ul>	<input checked="" type="checkbox"/>	TH

Select all Deselect all Export clusters Merge selected & re-cluster Merge selected & Close Close

**Mass edit 3,280 cells in column sponsor Undo**

3. Physical\_description

1. Split the Physical\_description column using the semicolon ';':

**Split column physical\_description 1 into several columns**

**How to split column**

☒ by separator  
 Separator:  ☐ regular expression  
 Split into  columns at most (leave blank for no limit)

☐ by field lengths  
  
 List of integers separated by commas, e.g., 5, 7, 15

**After Splitting**

☒ Guess cell type  
☒ Remove this column

OK Cancel

4. Date

1. Using Common transforms data to date

**Text transform on 16,959 cells in column date: value.toDate()**  
[Undo](#)

2. Change the date less than 1851 and more than 2015 to Null
3. Using GREL to format data to "YYYY-MM-DD"

**Custom text transform on column date**

Expression Language General Refine Expression Language (GREL)

`value.toString("yyyy-MM-dd")` No syntax error.

**Preview** History Starred Help

row	value	value.toString("yyyy-MM-dd")
1.	1900-04-15T00:00:00Z	1900-04-15
2.	1900-04-15T00:00:00Z	1900-04-15
3.	1900-04-16T00:00:00Z	1900-04-16
4.	1900-04-16T00:00:00Z	1900-04-16
5.	1900-04-16T00:00:00Z	1900-04-16
6.	1900-04-16T00:00:00Z	1900-04-16

On error ☒ keep original ☐ set to blank ☐ store error ☐ Re-transform up to  times until no change

OK Cancel

**Text transform on 16,954 cells in column date:**  
**grel:value.toString("yyyy-MM-dd")** [Undo](#)

---

## MENUPAGE

1. We only uppercased the UUID to make sure the format was consistent.

**Text transform on 66,936 cells in column uuid:**  
**value.toUpperCase()** [Undo](#)

---

## MENUITEM

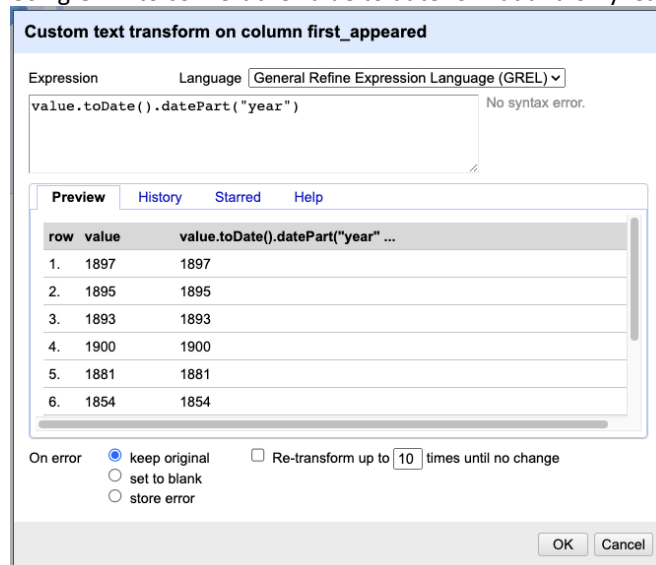
No OpenRefine operations were performed on this dataset. Cleanup operations were attempted on the date fields but no changes resulted as no date format inconsistencies were found.

---

## DISH

1. To resolve the MenuItem -> Dish relationships after the clustering took place. We have to retain the original name column and create a new column called clean\_name (Duplicated from the name column)
  1. Using text transformations in Edit cells
    1. Trim leading and trailing whitespace
    2. Collapse consecutive whitespace
    3. Remove the special characters '~!@#%&\*()[]\_|+V-=?,:;"<>' using General Refine Expression Language(GREL)
    4. Transform the text to lowercase

2. Create a text facet and implement the cluster operation using key-collision methods with different keying functions.
  1. Fingerprint function
  2. n-Gram fingerprint function with n-Gram size 2
  3. Metaphone3 function
2. first\_appeared, last\_appeared
  1. Using GREL to convert the value to date format and only leave the year part.



**Text transform on 367,905 cells in column first\_appeared:**  
**grel:value.toDate().datePart("year")** [Undo](#)

2. Change the illegal year which is less than 1851 or greater than 2928 to 0.

## SQLITE ACTIONS

### BASELINE MEASUREMENTS

The first steps performed in SQLite were to use the raw (uncleaned) data files and load them into the Initial SQLite Schema (see additional details below in the SQLite Schemas section). A set of count queries were run validating each of set of predefined IC check queries to determine the number of records in violation of each IC.

The baseline metrics are represented in the section titled "Document Data Quality Changes" as well as included in the supplementary file named: CS513-PhaseII-QualityMeasures.xlsx

### INTERMEDIATE CLEANUP

The data cleanup steps performed in OpenRefine were valuable for performing cleanup that was sensible within the boundaries of each input file. All relational cleanup work was performed in SQLite. The following cleanup operations were performed in SQLite.

1. Remove Duplicate (Clustered) Dish Records While Preserving Statistics

The OpenRefine cleanup operation saved a new name\_clean property that represented the output of the clustering operation on the Dish name property. It is important that we preserve only a single row representing the cluster. However, in order to achieve the desired output, we need to update all MenuItem references to Dish.id to reference the updated Dish.clustered\_id.

2. Recalculate Menu Dish Counts
3. Recalculate Dish Aggregate Data
  - a. Rather than rely on supplied data (which was often incorrect), this process recalculates the following Dish record attributes to represent accurate data based on MenuItem, MenuPage, and Menu datasets:
    - i. Menus\_appeared
    - ii. Times\_appeared
    - iii. First\_appeared
    - iv. Last\_appeared
    - v. Lowest\_price
    - vi. Highest\_price
4. Remove Unreferenced Records
  - a. Remove Menu Pages that reference a non-existent menu\_id.
  - b. Remove Menu Items that reference a non-existent menu\_page\_id.
  - c. Remove Menu Items that reference a non-existent dish\_id.
  - d. Remove Dishes that are not referenced by any Menu Items.
5. Extract categorized properties from physical\_description and push to the following attributes on each Menu record:
  - a. Misc Properties (extracted\_misc)
  - b. Damage Descriptions (extracted\_damage)
  - c. Dimension Descriptions (extracted\_dimensions)
  - d. Document Format Descriptions (extracted\_format)
6. Remove unneeded attributes:
  - a. Attributes that were always empty and do not contribute to any Use Case operations should be removed.
    - i. Menu.keywords
    - ii. Menu.language
    - iii. Menu.location\_type
    - iv. Dish.description

---

## FINAL MEASUREMENTS

The final steps performed in SQLite were to run the same predefined IC check queries to determine the number of records in violation of each IC in order to quantify the data quality improvements achieved.

The cleansed metrics are represented in the section titled “Document Data Quality Changes” as well as included in the supplementary file named: CS513-PhaseII-QualityMeasures.xlsx

---

## DETAILS STEPS

1. Add clustered\_id to Dish, seed with existing id value;
2. Populated DishIntermediate with starting data from OpenRefine Dish export, only generate one row per unique name\_clean, and assign the lowest (first occurring) Dish.id value as DishIntermediate.clustered\_id.



3. Populate MenuIntermediate with starting data from OpenRefine Menu export.
4. Override Dish.name with the cleaned version from name\_clean (**423397 records affected**).
5. Update all MenuItem records to reference the clustered\_id that corresponds to the original Dish.id value that was stored in MenuItem.dish\_id (**236813 records effected**).
6. Delete all Dish records that were “clustered” into another record (i.e. remove all rows where Dish.id is not used as any instance of Dish.clustered\_id) (**236813 records effected**).
7. Calculate DishIntermediate.menus\_appeared based on relationship from Dish through MenuItem through MenuPage to Menu.
8. Calculate DishIntermediate.times\_appeared based on relationship from Dish through to MenuItem.
9. Calculate DishIntermediate.first\_appeared by extracting the year from the earliest date for a Dish as it relates to Menu through MenuItem and MenuPage (using Menu.date\_clean).
10. Calculate DishIntermediate.last\_appeared by extracting the year from the latest date for a Dish as it relates to Menu through MenuItem and MenuPage (using Menu.date\_clean).
11. Calculate DishIntermediate.lowest\_price by extracting the minimum value from price and higher\_price for a Dish as it relates to MenuItem.
12. Calculate DishIntermediate.highest\_price by extracting the maximum value from price and higher\_price for a Dish as it relates to MenuItem.
13. Push recalculated menus\_appeared from DishIntermediate back on to corresponding Dish records. (**57582 records effected**)
14. Push recalculated times\_appeared from DishIntermediate back on to corresponding Dish records. (**58224 records effected**)
15. Push recalculated first\_appeared from DishIntermediate back on to corresponding Dish records. (**29035 records effected**)
16. Push recalculated last\_appeared from DishIntermediate back on to corresponding Dish records. (**32186 records effected**)
17. Push recalculated lowest\_price from DishIntermediate back on to corresponding Dish records. (**70989 records effected**)
18. Push recalculated highest\_price from DishIntermediate back on to corresponding Dish records. (**71774 records effected**)
19. Remove MenuPage records that do not reference an existing menu.id value. (**5803 records effected**)
20. Remove MenuItem records that do not reference an existing menupage.id value. (**5373 records effected**)
21. Remove MenuItem records that do not reference an existing dish.id value. (**244 records effected**)
22. Remove Dish records that are not referenced by any MenuItem.dish\_id values. (**2285 records effected**)
23. Calculate MenuIntermediate.dish\_count based on the number of MenuItem records that reference each menu through MenuPage.
24. Push recalculated dish\_count from MenuIntermediate back on to corresponding Menu records. (**4775 records effected**)
25. Trim leading/trailing whitespace in OpenRefine extracted physical\_description1, physical\_description2, physical\_description3, physical\_description4, physical\_description5, physical\_description6, physical\_description7 fields.
26. Populate Menu attributes: extracted\_misc, extracted\_dimensions, extracted\_damage, extracted\_format based on physical\_description# flags manually categorized in an Excel spreadsheet (see: **property-cleanup.xlsx** submitted with this assignment).
27. Drop physical\_description1 through physical\_description7 fields from Menu.
28. Drop DishIntermediate table;
29. Drop MenuIntermediate table;

30. Drop fields that are completely empty:

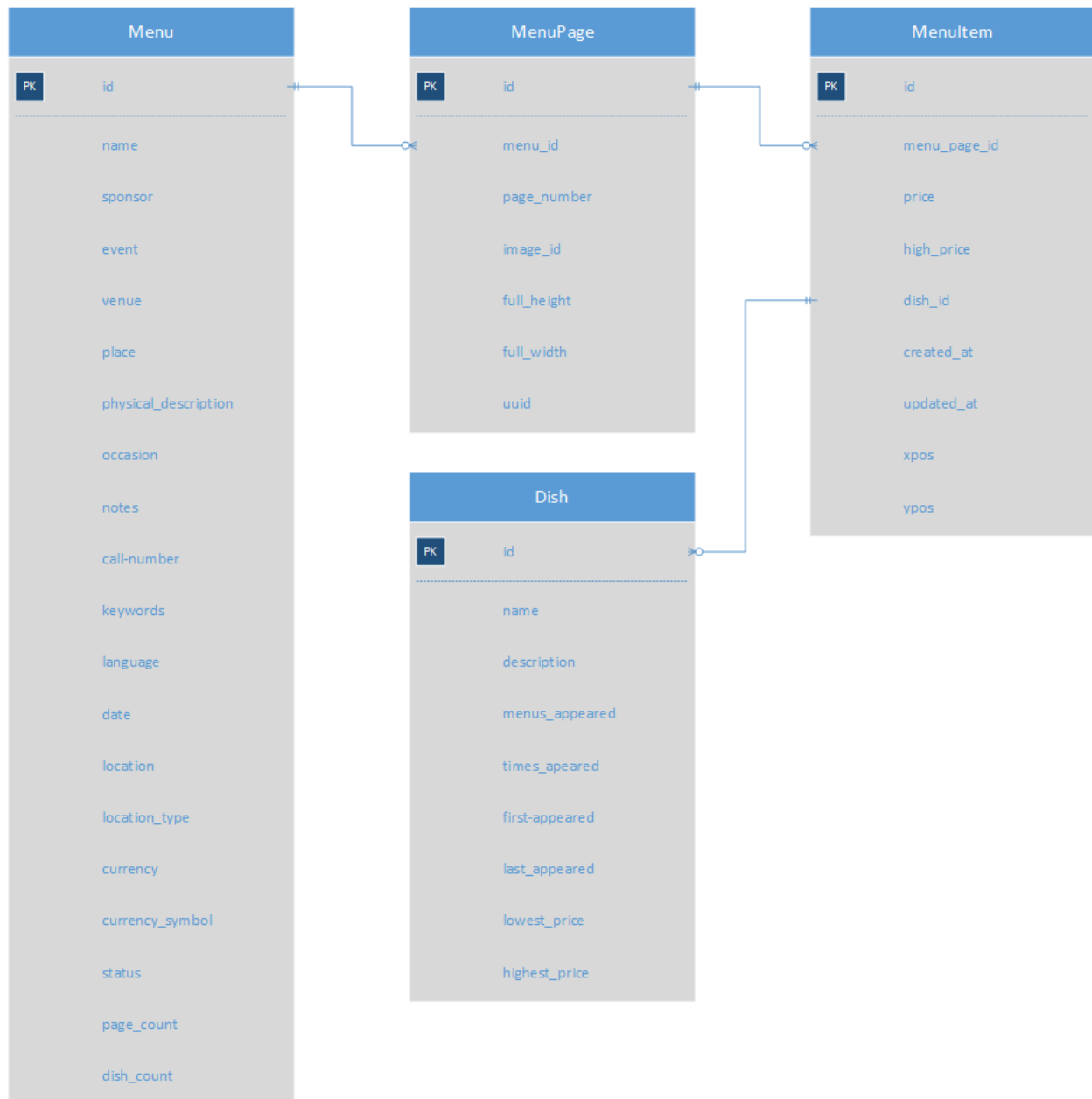
- a. Menu:
  - i. keywords
  - ii. language
  - iii. location\_type
- b. Dish
  - i. Description

## SQLITE SCHEMA ITERATIONS

A set of three SQLite data model iterations were used during the cleanup activity for this project.

1. Initial Schema
  - a. This schema represents the schema of the original data set and is identical to the schema included in our Phase I report.
2. Intermediate Schema
  - a. This schema represents the schema of the data as it was emitted from OpenRefine after the OpenRefine-based cleanup operations were completed.
  - b. For any notably transformative cleanup operations in OpenRefine (e.g., Clustering, or Flag-Expansion for the Physical Description attributes), the original attribute was preserved, and the new/cleansed attribute was added. This was necessary in order to preserve MenuItem to Dish relationships when deduplicating.
  - c. This schema also contains intermediate tables (MenuItemIntermediate and DishIntermediate) whose sole purpose is to support performant relational statistic recalculations for aggregate values on MenuItem and Dish tables.
3. Final Schema
  - a. This schema represents the final state of the cleansed data set.
  - b. This includes the elimination of both unwanted attributes (e.g., attributes with no valid data) as well as intermediate attributes (e.g., attributes we created to support FK-preservation during deduplication).

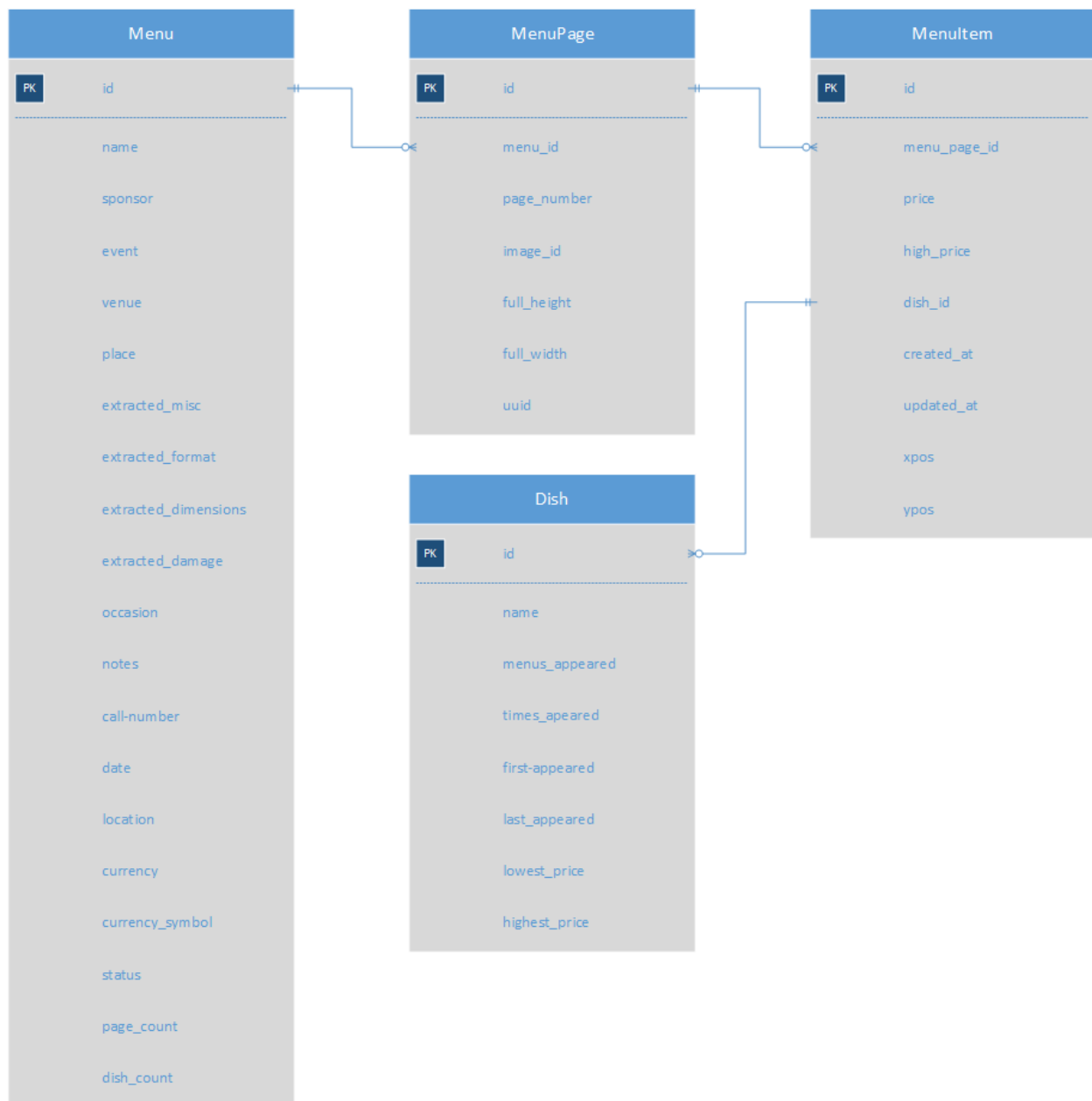
## INITIAL SCHEMA



## INTERMEDIATE SCHEMA



## FINAL SCHEMA



## SQLITE CONSTRAINTS

## MENU TABLE

## MISSING ID

```
SELECT id FROM Menu WHERE id is null or id = ''
```

## MISSING NAME

```
SELECT id FROM Menu WHERE name is null or name = ''
```

---

#### MISSING SPONSOR

SELECT id FROM Menu WHERE sponsor is null or sponsor = ''

---

#### MISSING EVENT

SELECT id FROM Menu WHERE event is null or event = ''

---

#### MISSING VENUE

SELECT id FROM Menu WHERE venue is null or venue = ''

---

#### MISSING PHYSICAL DESCRIPTION

SELECT id FROM Menu WHERE physical\_description is null or physical\_description = ''

---

#### MISSING PLACE

SELECT id FROM Menu WHERE place is null or place = ''

---

#### MISSING OCCASION

SELECT id FROM Menu WHERE occasion is null or occasion = ''

---

#### MISSING NOTES

SELECT id FROM Menu WHERE notes is null or notes = ''

---

#### MISSING CALL NUMBER

SELECT id FROM Menu WHERE call\_number is null or call\_number = ''

---

#### MISSING KEYWORDS

SELECT id FROM Menu WHERE keywords is null or keywords = ''

---

#### MISSING LANGUAGE

SELECT id FROM Menu WHERE language is null or language = ''

---

#### MISSING DATE

SELECT id FROM Menu WHERE date is null or date = ''

---

#### MISSING LOCATION

SELECT id FROM Menu WHERE location is null or location = ''

---

#### MISSING LOCATION TYPE

SELECT id FROM Menu WHERE location\_type is null or location\_type = ''

---

#### MISSING CURRENCY

SELECT id FROM Menu WHERE currency is null or currency = ''

---

#### MISSING CURRENCY SYMBOL

SELECT id FROM Menu WHERE currency\_symbol is null or currency\_symbol = ''

---

#### MISSING STATUS

SELECT id FROM Menu WHERE status is null or status = ''

---

#### MISSING PAGE COUNT

SELECT id FROM Menu WHERE page\_count is null or page\_count = ''

---

#### MISSING DISH COUNT

SELECT id FROM Menu WHERE dish\_count is null or dish\_count = ''

---

#### NON-UNIQUE ID

SELECT COUNT(\*), id FROM Menu GROUP BY id HAVING COUNT(\*) >= 2

---

#### NO PAGES

SELECT id FROM Menu WHERE page\_count = 0 or page\_count IS NULL

---

#### NO DISHES

SELECT id FROM Menu WHERE dish\_count = 0 or dish\_count IS NULL

---

#### UNREALISTICALLY EARLY DATE

SELECT id FROM Menu WHERE datetime(date) < datetime('1850-01-01')

---

#### UNREALISTICALLY LATE DATE

SELECT id FROM Menu WHERE datetime(date) > datetime('2023-08-01')

---

#### MENU NOT REFERENCED BY ANY MENU PAGES

SELECT id FROM Menu WHERE id NOT IN (select menu\_id from MenuPage)

---

#### DISH TABLE

---

##### MISSING ID

SELECT id FROM Dish WHERE id is null or id = ''

---

##### MISSING NAME

SELECT id FROM Dish WHERE name is null or name = ''

---

##### MISSING DESCRIPTION

SELECT id FROM Dish WHERE description is null or description = ''

---

##### MISSING MENUS APPEARED

SELECT id FROM Dish WHERE menus\_appeared is null or menus\_appeared = ''

---

#### MISSING TIMES APPEARED

SELECT id FROM Dish WHERE times\_appeared is null or times\_appeared = ''

---

#### MISSING FIRST APPEARED

SELECT id FROM Dish WHERE first\_appeared is null or first\_appeared = ''

---

#### MISSING LAST APPEARED

SELECT id FROM Dish WHERE last\_appeared is null or last\_appeared = ''

---

#### MISSING LOWEST PRICE

SELECT id FROM Dish WHERE lowest\_price is null or lowest\_price = ''

---

#### MISSING HIGHEST PRICE

SELECT id FROM Dish WHERE highest\_price is null or highest\_price = ''

---

#### NON-UNIQUE ID

SELECT COUNT(\*), id FROM Dish GROUP BY id HAVING COUNT(\*) >= 2)

---

#### ZERO MENUS APPEARED

SELECT id FROM Dish WHERE menus\_appeared = 0 or menus\_appeared IS NULL)

---

#### ZERO TIMES APPEARED

SELECT id FROM Dish WHERE times\_appeared = 0 or times\_appeared IS NULL)

---

#### ZERO FIRST YEAR APPEARED

SELECT id FROM Dish WHERE first\_appeared = 0 or first\_appeared IS NULL)

---

#### ZERO LAST YEAR APPEARED

SELECT id FROM Dish WHERE last\_appeared = 0 or last\_appeared IS NULL)

---

#### NO LOW PRICE

SELECT id FROM Dish WHERE lowest\_price = 0 or lowest\_price IS NULL)

---

#### NO HIGH PRICE

SELECT id FROM Dish WHERE highest\_price = 0 or highest\_price IS NULL)

---

#### DOES NOT APPEAR IN ANY MENUITEM RECORDS

SELECT id FROM Dish WHERE id NOT IN (select dish\_id from MenuItem))

---

#### HIGH PRICE IS LESS THAN LOW PRICE

---



```
SELECT id FROM Dish WHERE highest_price < lowest_price)
```

---

#### UNREALISTICALLY HIGH LOW-PRICE

```
SELECT id FROM Dish WHERE lowest_price > 500)
```

---

#### UNREALISTICALLY HIGH HIGH-PRICE

```
SELECT id FROM Dish WHERE highest_price > 500)
```

---

#### NEGATIVE LOW PRICE

```
SELECT id FROM Dish WHERE lowest_price < 0)
```

---

#### NEGATIVE HIGH PRICE

```
SELECT id FROM Dish WHERE highest_price < 0)
```

---

#### TIMES APPEARED NOT IN SYNC WITH RELATIONSHIPS

```
SELECT id FROM Dish WHERE id NOT IN (select dish_id from MenuItem) and times_appeared > 0);
```

---

#### MENUPAGE TABLE

---

##### MISSING ID

```
SELECT id FROM MenuPage WHERE id is null or id = ''
```

---

##### MISSING MENU ID

```
SELECT id FROM MenuPage WHERE menu_id is null or menu_id = ''
```

---

##### MISSING PAGE NUMBER

```
SELECT id FROM MenuPage WHERE page_number is null or page_number = ''
```

---

##### MISSING IMAGE ID

```
SELECT id FROM MenuPage WHERE image_id is null or image_id = ''
```

---

##### MISSING FULL HEIGHT

```
SELECT id FROM MenuPage WHERE full_height is null or full_height = ''
```

---

##### MISSING FULL WIDTH

```
SELECT id FROM MenuPage WHERE full_width is null or full_width = '') UNION
```

---

##### MISSING UUID

```
SELECT id FROM MenuPage WHERE uuid is null or uuid = ''
```

---

##### NON-UNIQUE ID

```
SELECT COUNT(*), id FROM MenuPage GROUP BY id HAVING COUNT(*) >= 2)
```

---

#### ZERO MENU ID

SELECT id FROM MenuPage WHERE menu\_id = 0 or menu\_id IS NULL

---

#### INVALID MENU ID (FK VIOLATION)

SELECT id FROM MenuPage WHERE menu\_id NOT IN (SELECT ID FROM Menu)

---

#### ZERO OR NEGATIVE PAGE NUMBER

SELECT id FROM MenuPage WHERE page\_number <= 0 or page\_number IS NULL

---

#### PAGE NUMBER OUT OF MENU PAGE RANGE

SELECT mp.id FROM MenuPage mp LEFT JOIN Menu m ON mp.menu\_id = m.id WHERE mp.page\_number > m.page\_count)

---

#### MORE PAGES FOR THIS MENU IN MENUPAGE THAN REPRESENTED IN MENU.PAGE\_COUNT

SELECT mp.id FROM MenuPage mp LEFT JOIN Menu m ON mp.menu\_id = m.id WHERE mp.page\_number > m.page\_count)

---

#### NOT REFERENCED BY ANY MENUITEM

SELECT id FROM MenuPage WHERE id NOT IN (select menu\_page\_id from MenuItem)

---

#### MORE PAGES IN MENU THAN REPRESENTED IN MENUPAGE

SELECT m.id FROM Menu m LEFT JOIN (SELECT COUNT(\*) as page\_count, menu\_id FROM MenuPage GROUP BY menu\_id) mp on m.id = mp.menu\_id WHERE mp.page\_count < m.page\_count)

---

#### FEWER PAGES IN MENU THAN REPRESENTED IN MENUPAGE

SELECT m.id FROM Menu m LEFT JOIN (SELECT COUNT(\*) as page\_count, menu\_id FROM MenuPage GROUP BY menu\_id) mp on m.id = mp.menu\_id WHERE mp.page\_count > m.page\_count);

---

### MENUITEM TABLE

---

#### MISSING ID

SELECT id FROM MenuItem WHERE id is null or id = ''

---

#### MISSING MENU PAGE ID

SELECT id FROM MenuItem WHERE menu\_page\_id is null or menu\_page\_id = ''

---

#### MISSING PRICE

SELECT id FROM MenuItem WHERE price is null or price = ''

---

#### MISSING HIGH PRICE

SELECT id FROM MenuItem WHERE high\_price is null or high\_price = ''

---

#### MISSING DISH ID

```
SELECT id FROM MenuItem WHERE dish_id is null or dish_id = ''
```

---

#### MISSING CREATED AT

```
SELECT id FROM MenuItem WHERE created_at is null or created_at = ''
```

---

#### MISSING UPDATED AT

```
SELECT id FROM MenuItem WHERE updated_at is null or updated_at = ''
```

---

#### MISSING X POSITINO

```
SELECT id FROM MenuItem WHERE xpos is null or xpos = ''
```

---

#### MISSING Y POSITION

```
SELECT id FROM MenuItem WHERE ypos is null or ypos = ''
```

---

#### NON-UNIQUE ID

```
SELECT COUNT(*), id FROM MenuItem GROUP BY id HAVING COUNT(*) >= 2
```

---

#### ZERO MENU PAGE ID

```
SELECT id FROM MenuItem WHERE menu_page_id = 0 or menu_page_id IS NULL
```

---

#### ZERO DISH ID

```
SELECT id FROM MenuItem WHERE dish_id = 0 or dish_id IS NULL
```

---

#### INVALID MENU PAGE ID (FK CONSTRAINT)

```
SELECT id FROM MenuItem WHERE menu_page_id NOT IN (SELECT ID FROM MenuPage)
```

---

#### INVALID DISH ID (FK CONSTRAINT)

```
SELECT id FROM MenuItem WHERE dish_id NOT IN (SELECT ID FROM Dish)
```

---

#### UPDATED TIMESTAMP EARLIER THAN CREATION TIMESTAMP

```
SELECT id FROM MenuItem WHERE datetime(updated_at) < datetime(created_at)
```

---

#### PRICE ON ITEM IS LOWER THAN DISH LOWEST PRICE

```
SELECT mi.id from MenuItem mi INNER JOIN Dish d on mi.dish_id = d.id WHERE mi.price < d.lowest_price and  
mi.price != '' and mi.price is not null)
```

---

#### PRICE ON ITEM IS HIGHER THAN DISH HIGHEST PRICE

```
SELECT mi.id from MenuItem mi INNER JOIN Dish d on mi.dish_id = d.id WHERE mi.price > d.highest_price and  
mi.price != '' and mi.price is not null)
```

---

#### HIGH PRICE ON ITEM IS LOWER THAN DISH LOWEST PRICE

SELECT mi.id from MenuItem mi INNER JOIN Dish d on mi.dish\_id = d.id WHERE mi.high\_price is not null and mi.high\_price != " and mi.high\_price < d.lowest\_price)

#### HIGH PRICE ON ITEM IS HIGHER THAN DISH HIGHEST PRICE

SELECT mi.id from MenuItem mi INNER JOIN Dish d on mi.dish\_id = d.id WHERE mi.high\_price is not null and mi.high\_price != " and mi.high\_price > d.highest\_price);

## 2. DOCUMENT DATA QUALITY CHANGES

All IC checks were chained together in a series of queries to be executed in the SQLite environment enabling us to compare the values produced by these checks prior to the data cleanup being executed and after the data cleanup being executed. Please review the table below for a summary of outcomes. Each constraint will correspond to a constraint defined in the previous section.

Table	Measure	Total Rows	Raw Data	%	Clean Total Rows	Clean Data	%	Delta	% Delta
Menu	Total Rows	17545	17545	100.00%	17545	17545	100.00%	0	
MenuPage	Total Rows	66937	66937	100.00%	61134	61134	100.00%	5803	
Dish	Total Rows	423397	423397	100.00%	184299	184299	100.00%	239098	
MenuItem	Total Rows	1332726	1332726	100.00%	1327109	1327109	100.00%	5617	
Menu	Missing Attribute: Dish Count	17545	0	0.00%	17545	0	0.00%	0	0.00%
Menu	Missing Attribute: ID	17545	0	0.00%	17545	0	0.00%	0	0.00%
Menu	Missing Attribute: Location	17545	0	0.00%	17545	48	0.27%	-48	-0.27%
Menu	Missing Attribute: Page Count	17545	0	0.00%	17545	0	0.00%	0	0.00%
Menu	Missing Attribute: Status	17545	0	0.00%	17545	0	0.00%	0	0.00%
Menu	Missing Attribute: Date	17545	586	3.34%	17545	586	3.34%	0	0.00%
Menu	Missing Attribute: Sponsor	17545	1561	8.90%	17545	1561	8.90%	0	0.00%
Menu	Missing Attribute: Call Number	17545	1562	8.90%	17545	1562	8.90%	0	0.00%
Menu	Missing Attribute: Physical Description	17545	2777	15.83%	17545	0	0.00%	2777	15.83%
Menu	Missing Attribute: Notes	17545	6932	39.51%	17545	6933	39.52%	-1	-0.01%
Menu	Missing Attribute: Event	17545	9391	53.53%	17545	9409	53.63%	-18	-0.10%
Menu	Missing Attribute: Venue	17545	9414	53.66%	17545	9435	53.78%	-21	-0.12%
Menu	Missing Attribute: Place	17545	9422	53.70%	17545	9507	54.19%	-85	-0.48%
Menu	Missing Attribute: Currency	17545	11089	63.20%	17545	11089	63.20%	0	0.00%
Menu	Missing Attribute: Currency Symbol	17545	11089	63.20%	17545	11089	63.20%	0	0.00%

<b>Menu</b>	Missing Attribute: Occasion	17545	13742	78.32%	17545	13793	78.61%	-51	-0.29%
<b>Menu</b>	Missing Attribute: Name	17545	14348	81.78%	17545	14348	81.78%	0	0.00%
<b>Menu</b>	Missing Attribute: Keywords	17545	17545	100.00%	17545	17545	100.00%	0	0.00%
<b>Menu</b>	Missing Attribute: Language	17545	17545	100.00%	17545	17545	100.00%	0	0.00%
<b>Menu</b>	Missing Attribute: Location Type	17545	17545	100.00%	17545	17545	100.00%	0	0.00%
<b>Menu</b>	Non-Unique ID	17545	0	0.00%	17545	0	0.00%	0	0.00%
<b>Menu</b>	NULL ID	17545	0	0.00%	17545	0	0.00%	0	0.00%
<b>Menu</b>	Zero MenuPage Records related to this Menu	17545	0	0.00%	17545	0	0.00%	0	0.00%
<b>Menu</b>	Zero Pages in Page Count Attribute	17545	0	0.00%	17545	0	0.00%	0	0.00%
<b>Menu</b>	Zero Dishes	17545	32	0.18%	17545	28	0.16%	4	0.02%
<b>Menu</b>	Very Early Date	17545	4	0.02%	17545	4	0.02%	0	0.00%
<b>Menu</b>	Very Late Date	17545	1	0.01%	17545	1	0.01%	0	0.00%
<b>Dish</b>	Missing Attribute: First Appeared	423397	0	0.00%	184299	0	0.00%	0	0.00%
<b>Dish</b>	Missing Attribute: Id	423397	0	0.00%	184299	0	0.00%	0	0.00%
<b>Dish</b>	Missing Attribute: Last Appeared	423397	0	0.00%	184299	0	0.00%	0	0.00%
<b>Dish</b>	Missing Attribute: Menus Appeared	423397	0	0.00%	184299	0	0.00%	0	0.00%
<b>Dish</b>	Missing Attribute: Name	423397	0	0.00%	184299	0	0.00%	0	0.00%
<b>Dish</b>	Missing Attribute: Times Appeared	423397	0	0.00%	184299	0	0.00%	0	0.00%
<b>Dish</b>	Missing Attribute: Highest Price	423397	29100	6.87%	184299	9678	5.25%	19422	1.62%
<b>Dish</b>	Missing Attribute: Lowest Price	423397	29100	6.87%	184299	9678	5.25%	19422	1.62%
<b>Dish</b>	Missing Attribute: Description	423397	423397	100.00%	184299	184299	100.00%	239098	0.00%
<b>Dish</b>	Non-Unique ID	423397	0	0.00%	184299	0	0.00%	0	0.00%
<b>Dish</b>	NULL ID	423397	0	0.00%	184299	0	0.00%	0	0.00%
<b>Dish</b>	Relationally Appeared in Zero Menus	423397	9262	2.19%	184299	0	0.00%	9262	2.19%
<b>Dish</b>	Negative High Price	423397	0	0.00%	184299	0	0.00%	0	0.00%
<b>Dish</b>	Negative Low Price	423397	0	0.00%	184299	0	0.00%	0	0.00%
<b>Dish</b>	Dish not referenced but MenuItem, but	423397	9	0.00%	184299	0	0.00%	9	0.00%

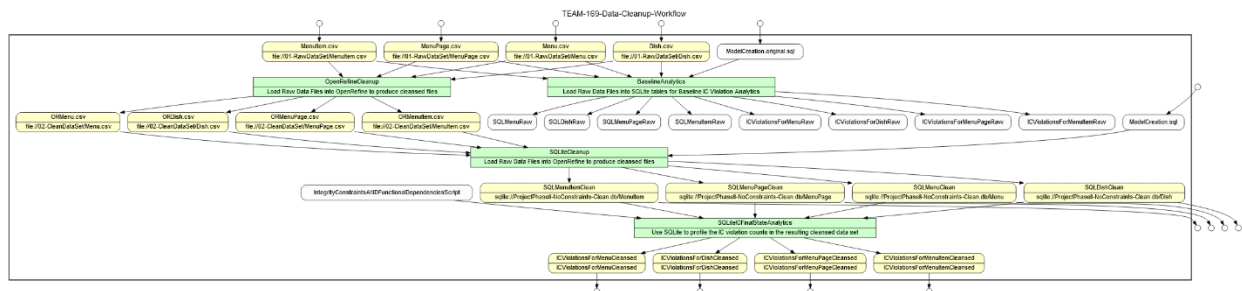
Times Appeared > 0									
Dish	Zero Menu Appeared	423397	2271	0.54%	184299	0	0.00%	2271	0.54%
Dish	Zero Times Appeared	423397	9248	2.18%	184299	0	0.00%	9248	2.18%
Dish	Zero First Appeared	423397	55284	13.06%	184299	18310	9.93%	36974	3.12%
Dish	Zero Last Appeared	423397	55287	13.06%	184299	18317	9.94%	36970	3.12%
Dish	Zero Highest Price	423397	218014	51.49%	184299	48183	26.14%	169831	25.35%
Dish	Zero Lowest Price	423397	222566	52.57%	184299	48373	26.25%	174193	26.32%
MenuPage	Missing Attribute: Id	66937	0	0.00%	61134	0	0.00%	0	0.00%
MenuPage	Missing Attribute: Image Id	66937	0	0.00%	61134	0	0.00%	0	0.00%
MenuPage	Missing Attribute: Menu Id	66937	0	0.00%	61134	0	0.00%	0	0.00%
MenuPage	Missing Attribute: Uuid	66937	0	0.00%	61134	0	0.00%	0	0.00%
MenuPage	Missing Attribute: Full Height	66937	329	0.49%	61134	0	0.00%	329	0.49%
MenuPage	Missing Attribute: Full Width	66937	329	0.49%	61134	0	0.00%	329	0.49%
MenuPage	Missing Attribute: Page Number	66937	1202	1.80%	61134	945	1.55%	257	0.25%
MenuPage	Non-Unique ID	66937	0	0.00%	61134	0	0.00%	0	0.00%
MenuPage	NULL ID	66937	0	0.00%	61134	0	0.00%	0	0.00%
MenuPage	Relational Page Count Mismatch - Actual Count Higher	66937	0	0.00%	61134	0	0.00%	0	0.00%
MenuPage	Relational Page Count Mismatch - Actual Count Lower	66937	0	0.00%	61134	0	0.00%	0	0.00%
MenuPage	Zero Menu ID	66937	0	0.00%	61134	0	0.00%	0	0.00%
MenuPage	Page Number Out of Range	66937	951	1.42%	61134	951	1.56%	0	-0.13%
MenuPage	Relational - Page Count Aggregate	66937	951	1.42%	61134	951	1.56%	0	-0.13%
MenuPage	Invalid Menu ID	66937	5803	8.67%	61134	0	0.00%	5803	8.67%
MenuPage	Page Not Referenced by Menu Item	66937	40347	60.28%	61134	34656	56.69%	5691	3.59%
MenuItem	Missing Attribute: Created At	1332726	0	0.00%	1327109	0	0.00%	0	0.00%
MenuItem	Missing Attribute: Id	1332726	0	0.00%	1327109	0	0.00%	0	0.00%
MenuItem	Missing Attribute: Menu Page Id	1332726	0	0.00%	1327109	0	0.00%	0	0.00%

<b>MenuItem</b>	Missing Attribute: Updated At	1332726	0	0.00%	1327109	0	0.00%	0	0.00%
<b>MenuItem</b>	Missing Attribute: Xpos	1332726	0	0.00%	1327109	0	0.00%	0	0.00%
<b>MenuItem</b>	Missing Attribute: Ypos	1332726	0	0.00%	1327109	0	0.00%	0	0.00%
<b>MenuItem</b>	Missing Attribute: Dish Id	1332726	241	0.02%	1327109	0	0.00%	241	0.02%
<b>MenuItem</b>	Missing Attribute: Price	1332726	445916	33.46%	1327109	444424	33.49%	1492	-0.03%
<b>MenuItem</b>	Missing Attribute: High Price	1332726	1240821	93.10%	1327109	1235450	93.09%	5371	0.01%
<b>MenuItem</b>	Invalid Menu Page ID	1332726	0	0.00%	1327109	0	0.00%	0	0.00%
<b>MenuItem</b>	Non-Unique ID	1332726	0	0.00%	1327109	0	0.00%	0	0.00%
<b>MenuItem</b>	NULL ID	1332726	0	0.00%	1327109	0	0.00%	0	0.00%
<b>MenuItem</b>	Zero Dish ID	1332726	0	0.00%	1327109	0	0.00%	0	0.00%
<b>MenuItem</b>	Zero Menu Page ID	1332726	0	0.00%	1327109	0	0.00%	0	0.00%
<b>MenuItem</b>	Higher Price Lower than Low Price on Dish (relational)	1332726	58	0.00%	1327109	0	0.00%	58	0.00%
<b>MenuItem</b>	Invalid Dish ID	1332726	244	0.02%	1327109	0	0.00%	244	0.02%
<b>MenuItem</b>	Price Lower Than Lowest Price on Dish (relational)	1332726	1172	0.09%	1327109	0	0.00%	1172	0.09%
<b>MenuItem</b>	High Price Higher than Highest Price on Dish (relational)	1332726	13591	1.02%	1327109	0	0.00%	13591	1.02%
<b>MenuItem</b>	Price Higher than Highest Price on Dish (relational)	1332726	149616	11.23%	1327109	0	0.00%	149616	11.23%

### 3. WORKFLOW MODEL

All workflow images have a corresponding higher-resolution PNG or JPG file included in the project submission for situations where the readability is difficult due to image scaling. See Supplementary Materials section for more details.

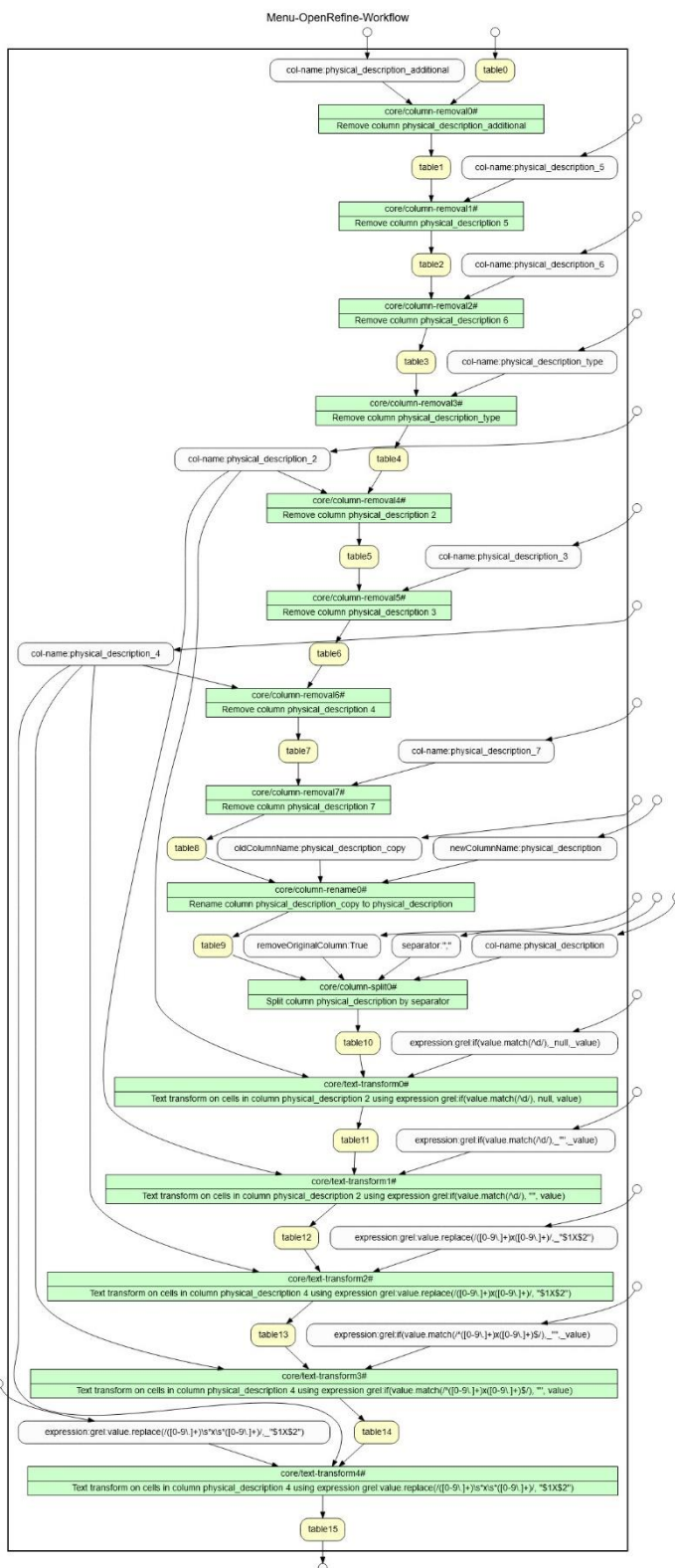
#### OUTER WORKFLOW



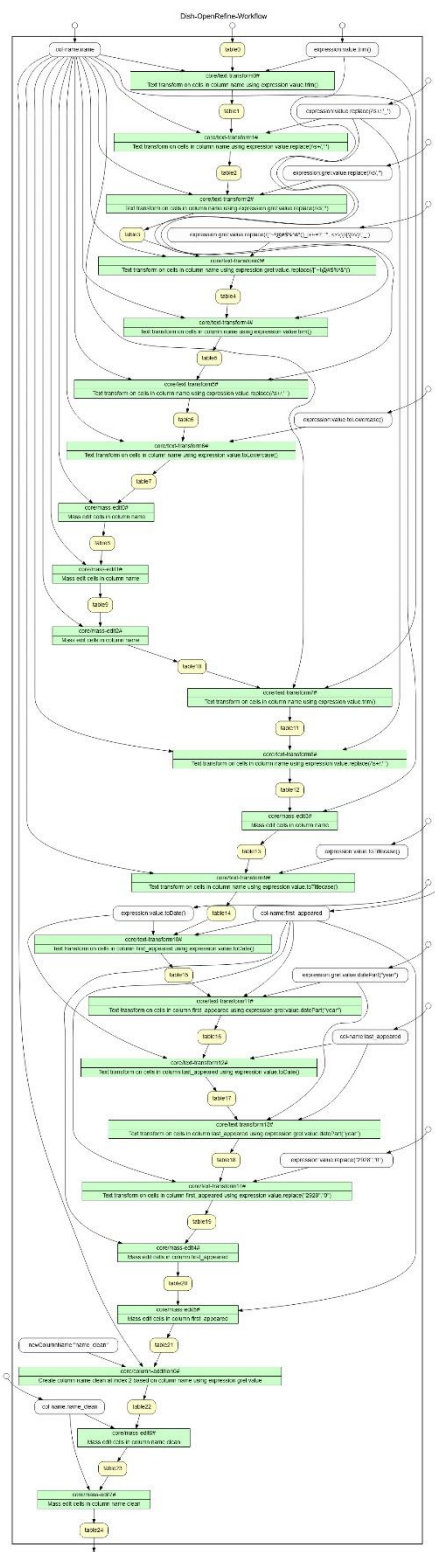
#### INNER WORKFLOW

## OPENREFINE

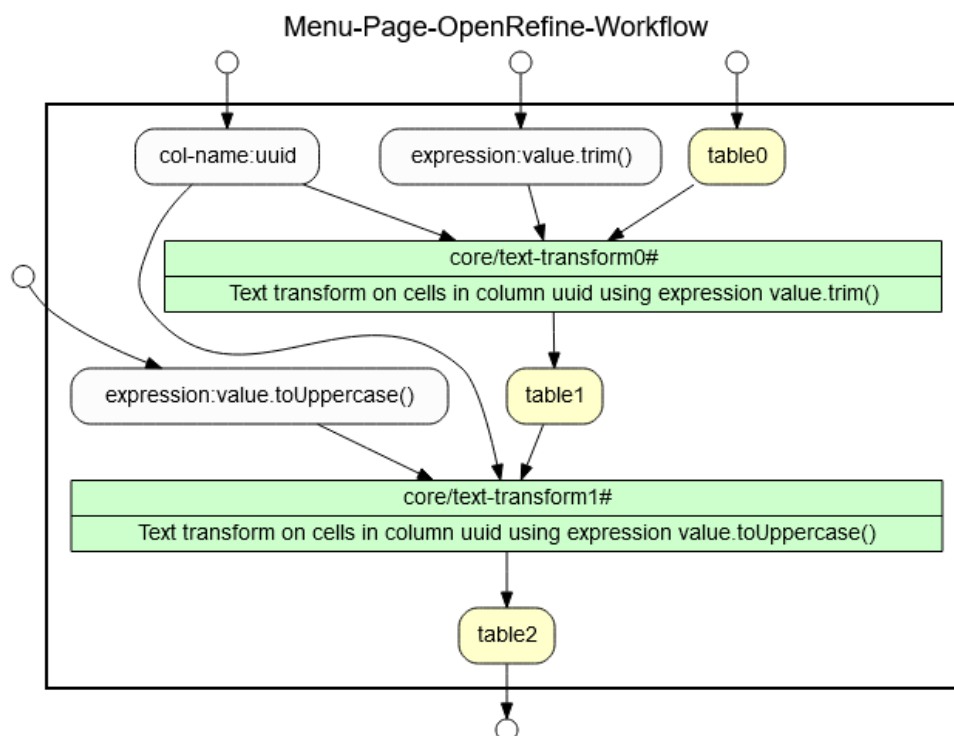
## MENU



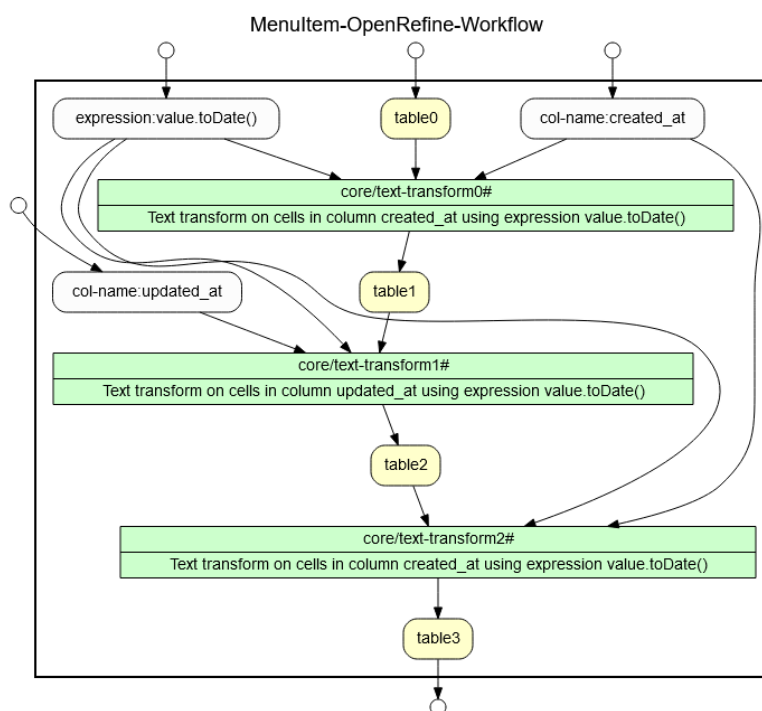




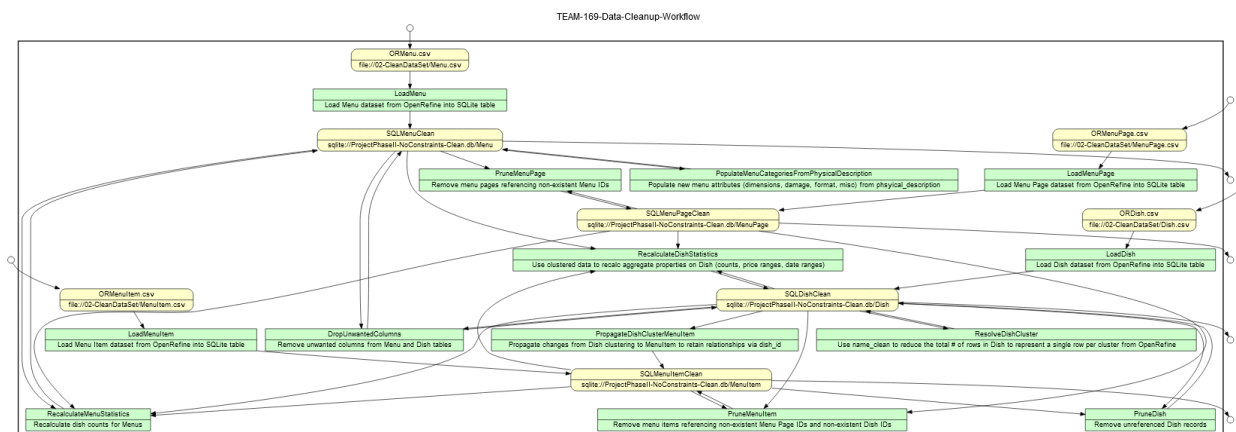
## MENUPAGE



## MENUITEM



## SQLITE



## 4. CONCLUSIONS &amp; SUMMARY

The processes described in this report have demonstrated a high value in improving the data quality, particularly with aggregate data calculation between data entities, in service of our U1 Main Use Case.

Both the OpenRefine and SQLite processes implemented in service of this cleanup project were valuable in achieving the goals we set out to accomplish.

Some moderately cumbersome SQLite conventions were used (i.e., frequent use of UPSERT operations) to achieve high-speed execution of our SQLite based operations. While these SQL expressions were not as concise as ones that would have been described using a SQL platform that supported UPDATE operations with JOIN conditions, the UPSERT conventions used did enable us to implement a SQL-based cleanup process that was able to execute in a very high-speed fashion.

## CHALLENGES

The key challenge experienced in this activity was the late realization that the low and high price aggregate representations produce a moderately inaccurate price range given the inconsistency between currency and currency units across different representations of the same dish.

The team agreed not to pursue any currency normalization (i.e., rebase the aggregate price calculations using a normalized currency representation based on historical foreign exchange rates for the involved currencies). The team made this decision based on the belief that this expectation was not set out by our team in our original Project Plan. Were this activity to be undertaken again, the team would consider either currency normalization or representing low/high price ranges in all available currencies for any Dish that is priced in any Menu using that currency.

## TEAM MEMBER CONTRIBUTIONS

The group regularly engaged in working sessions where the full team was present on a call while work was being conducted. The table below summarizes the areas where each members contributions were central.

TEAM MEMBER	PRIMARY AREAS OF CONTRIBUTION
AARON	SQLite Operations, Yes Workflow Definition, Report preparation
CHIRU	OpenRefine Cleanup Process, Yes Workflow Definition
WILLY	OpenRefine Cleanup Process, Yes Workflow Definition, OpenRefine process documentation

## 5. SUPPLEMENTARY MATERIALS

The following materials were submitted with the report:

- Property-cleanup.xlsx
  - o Spreadsheet representing the categorization work performed on flag values extracted from physical\_description in the Menu data set.
- CS513-PhaseII-QualityMeasures.xlsx
  - o Spreadsheet version of the quantified data quality changes included in section 2.
- OpenRefine-Project.zip
  - o OpenRefine project files representing the cleanup activities and provenance for the OpenRefine based cleanup operations described in the “OpenRefine Actions” section.
- OpenRefine-CleansedOutput.zip
  - o ZIP containing 4x CSV files representing the intermediate state cleaned data for each of the four tables produced as a result of OpenRefine cleanup operations..
- SQLite-CleansedOutput.zip
  - o ZIP containing 4x CSV files representing the final state cleaned data for each of the four tables represented in our initial and final schema.
- ModelCreation.original.sql
  - o Script file used to load and operate the baseline data quality assessment on the raw input files in SQLite.
- ModelCreation.sql
  - o Script file used to load the OpenRefine-cleansed data files, operate the SQLite data cleanup, and operate the finish-line data quality assessment.
- YesWorkflow-Outer.png
  - o Detailed PNG for workflow represented in Workflow Model – Outer Workflow section.
- YesWorkflow-Inner-OpenRefine-Menu.jpg
- YesWorkflow-Inner-OpenRefine-Dish.jpg
- YesWorkflow-Inner-OpenRefine-MenuItem.png
- YesWorkflow-Inner-OpenRefine-MenuPage.png
  - o Detailed PNG/JPG files for workflow represented in Workflow Model – Inner Workflow – OpenRefine section.
- YesWorkflow-Inner-SQLite.png
  - o Detailed PNG for workflow represented in Workflow Model – Inner Workflow – SQLite section.
- Outer-workflow.yw
  - o Outer YesWorkflow workflow file.
- MenuItem-OR.yw
- Menu-OR.yw

- MenuPage-OR.yw
- Dish-OR.yw
- Inner-workflow-sqlite.yw
  - o Inner YesWorkflow workflow definition files.