

5. INTEGRATING DATA

As demonstrated in the last chapter, it is feasible today to let casual users with no programming skills extract data from the text-centric Web for reuse. Such capability is already useful for any single web site that does not offer the features that a casual user needs. More empowering to casual users is the ability to combine data from several sites to get values that no single site can offer by itself.

Like conventional data extraction tools, conventional data integration tools have also been built for a different audience than casual users. For a typical conventional data integration task, such as merging huge databases of two institutional libraries together to serve the combined data through a new web site, several experts and programmers are involved to handle different aspects of the task, including aligning the schemas, cleaning up the data, and then building the web site. Each tool employed is specialized for one stage of the process and designed to handle large, complex data sets. In contrast, a casual user might just want to merge two lists of a few dozens of addresses from two web sites together to plot them on a common map. There is much less data and the data is much simpler. Power tools used by experts are too advanced for casual users.

Compared to experts, casual users lack both data modeling skills and programming skills. However, for small, simple data sets, neither skill set may be crucial.

- First, when the data set is small, schemas—useful for efficiently reasoning about a massive quantity of data in the abstract—are not as useful and can introduce overhead cognitive costs. Instead of having to learn about theoretical concepts like schemas, casual users can manipulate data instances directly. Visual semantics are often enough: if the data “looks right” in the user interface, it most probably has been edited correctly and there is no need to verify the data model.

- Second, instead of using programming to process data, casual users can just use direct manipulation techniques. For example, two fields can be aligned by dragging and dropping one onto the other.

These ideas have been built into Potluck, a tool that lets casual users—non-programmers—integrate data all by themselves. This chapter will next describe a hypothetical usage scenario for Potluck. Using various challenges in that scenario as motivations, the user interface of Potluck will be explained. Then Potluck’s implementation will be detailed. Finally, an evaluation of Potluck is reported.

5.1 Scenario

Before describing the user interface of Potluck, this section motivate it with a scenario that illustrates various idiosyncrasies of data integration. Let us be optimistic that within a decade, the Semantic Web will be prevalent and RDF data will be everywhere. Even in this future world, users will *still* face problems integrating data from different sources and tools such as Potluck are still needed.

In 2017, a historian named Henry is documenting the first cases of a rare genetic disease called GD726. These first cases occurred in the Valentine family in the 1820s. He wants to include in his final report a genealogical tree of the Valentine family, annotated with the disease’s infliction, as well as a comprehensive table of the Valentines’ data in an appendix.

Like most historians, Henry is not a programmer but he is experienced in collecting and managing data in his professional work. The proliferation of RDF means that Henry does not need programming skills to scrape HTML himself: all the information needed for his research has been converted into RDF by various independent organizations and individuals, both professionals and enthusiasts. Henry thinks it would be trivial to simply pool the RDF together and call it done.

Henry tracks down various birth certificate issuing offices and death certificate issuing offices where the Valentines lived for their RDF data. He notes that some offices use `dc:date` in their data to mean “birth date,” some to mean “death date,” and some “certificate issuing date.” It would be disastrous to consider all the `dc:dates` the same even if the same predicate URI is used.

Henry also tracks down hospital records, which contain `hospital:tod` (short for “time of death”). Hence, `hospital:tod` is equivalent to some of the `dc:dates`. It would be hard to match `hospital:tod` with `dc:date` based on string analysis alone, yet match for some of the cases only.

The records all have geographical location names, but these names are not fully qualified. Those responsible for digitizing them thought that since all locations were within their country, there was no need to include the country name. As a consequence, Henry needs to append the country name to the many location names in order to map them.

People’s names are encoded in two different forms: “first-name last-name” in some data sets and “last-name, first-name” in others. Nick names are also present (e.g., “Bill” instead of “William”, and “Vicky” instead of “Victoria”).

The hospital records also pose problems. While most of their admittance dates are in ISO 8601 format, a few are of the kind “Easter Day 1824.” Such sloppiness has been observed in industrial and institutional databases, and should be expected on the Semantic Web.

Despite all these problems, there is one good thing about the data: Henry can reliably get the mother and father of each Valentine through the `gen:mother` and `gen:father` predicates, which seem to be very widely adopted. This helps Henry construct a genealogical tree visualization.

However, as males and females both have equal chance of passing on GD726, Henry wants to treat `gen:mother` and `gen:father` the same while tracing the disease through the family. Unfortunately, adding an `owl:sameAs` equivalence between those two predicates will break his genealogical tree.

While all parties involved in this scenario acted logically and responsibly, Henry still ends up with a mess of RDF. To fix up the data, Henry must be able to:

- Merge `dc:dates` into *several* groups (the birth dates and the death dates) even though they all use the same predicate URI. This requires distinguishing the fields by their origins rather than just by their URIs.
- Merge `gen:mother` and `gen:father` together in some situations while keeping them separate in other situations. This precludes the simple approach of adding `owl:sameAs` statements in the data model to implement equivalences.
- Edit the data efficiently to unify its syntax.
- Fix up the data iteratively as he learns more and more about the data.

These are the tasks that must be supported by such a tool as Potluck in order for a casual user such as Henry to be able to integrate data all by himself.

5.2 User Interface

This section describes Potluck’s user interface, showing how it addresses the problems in the scenario above. The reader is encouraged to view this screencast to understand Potluck’s interactivity:

<http://people.csail.mit.edu/dfhuynh/research/media/iswc2007/>.

Figure 5.1 shows the starting screen of Potluck where the user can paste in several URLs and click **Mix Data**. This results in Figure 5.2, which lists data records from the original web pages. The records are interleaved by origins—the pages from which they have been extracted—to ensure that some records of each data set are always visible.

Fields are rendered as *field tags*: `label`, `position`, and `title`. Field tags are color-coded to indicate their origins: blue from one source and pink from another in

5. INTEGRATING DATA

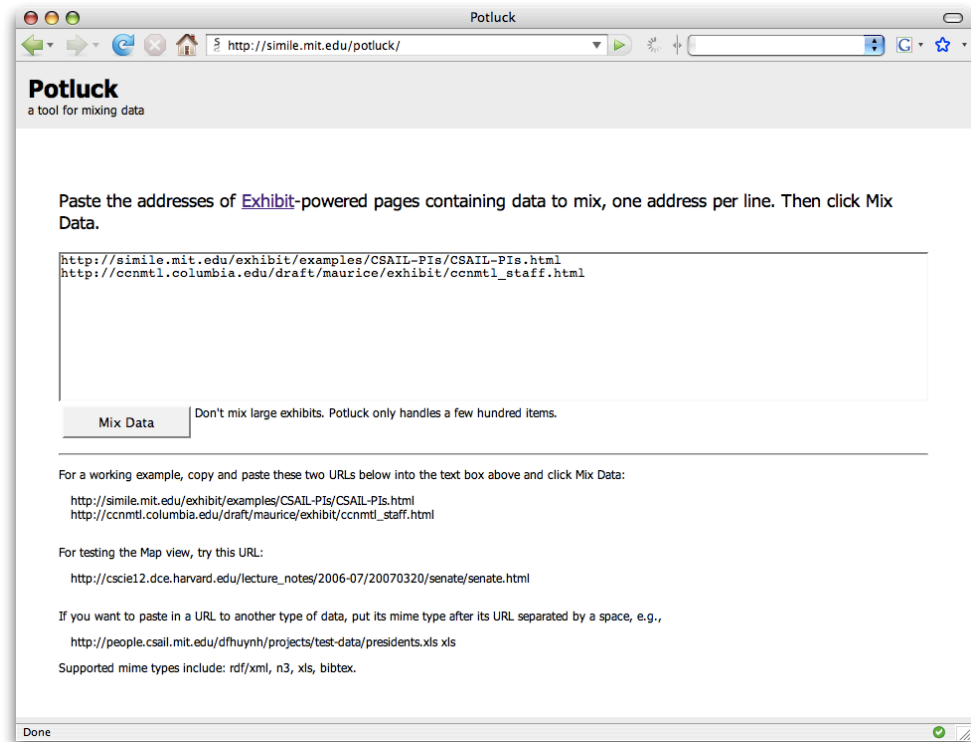


Figure 5.1. The starting screen of Potluck. Clicking **Mix Data** yields the mixed data in a screen like Figure 2.

Figure 5.2. Three core fields, **label**, **type**, and **origin**, are automatically assigned to all records and their tags are colored gray. Fields from different origins having the same name are considered different. For example, while **phone** means office phone, **phone** might mean secretary's phone. Or more dangerously, **dc:date** in the scenario (in section 2) has several distinct meanings. These semantic differences, subtle or significant, might or might not be important to one particular user at one particular moment in time. Keeping the fields apart rather than automatically merging them together allows the user to make the decision whether or not to merge.

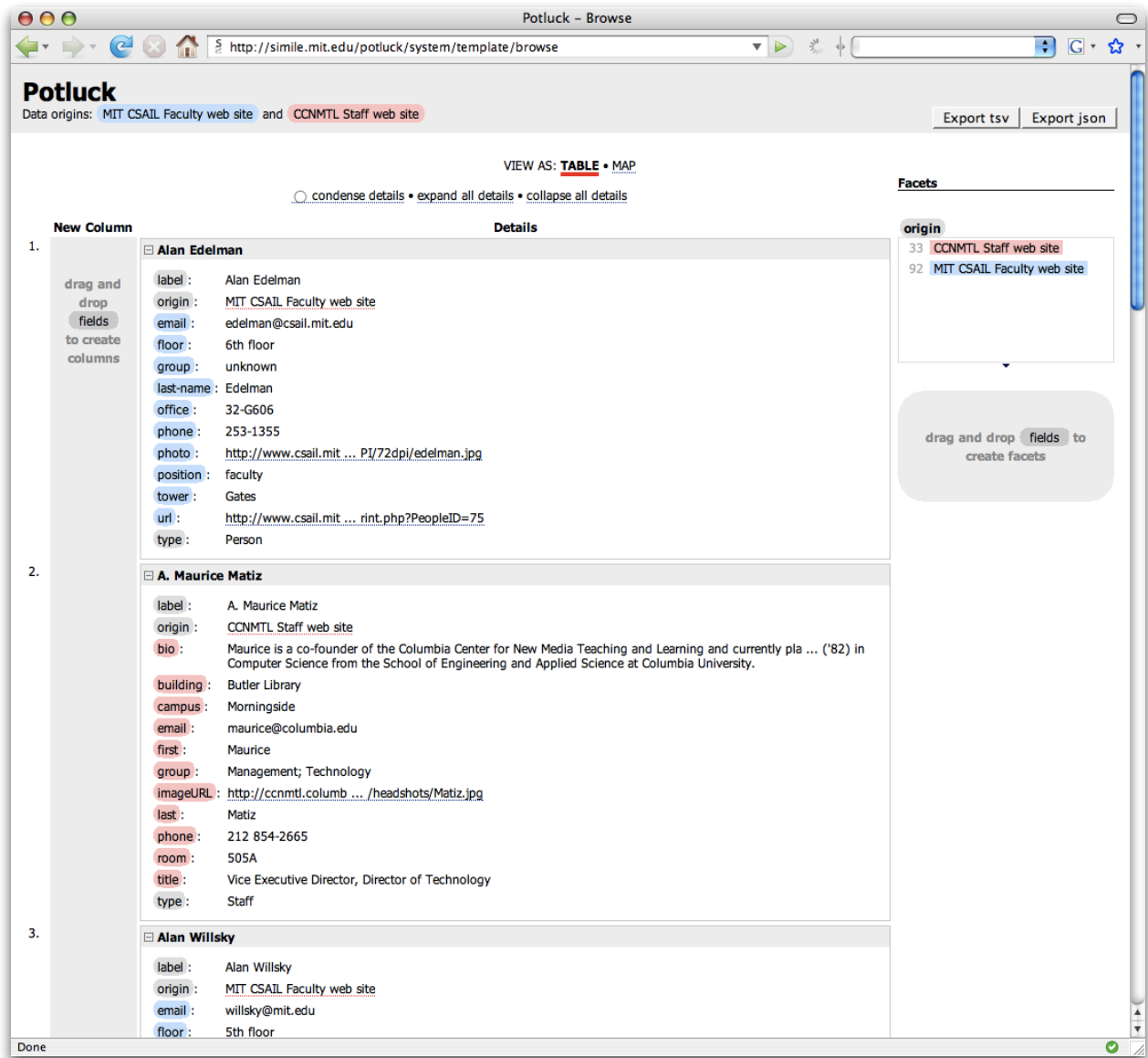

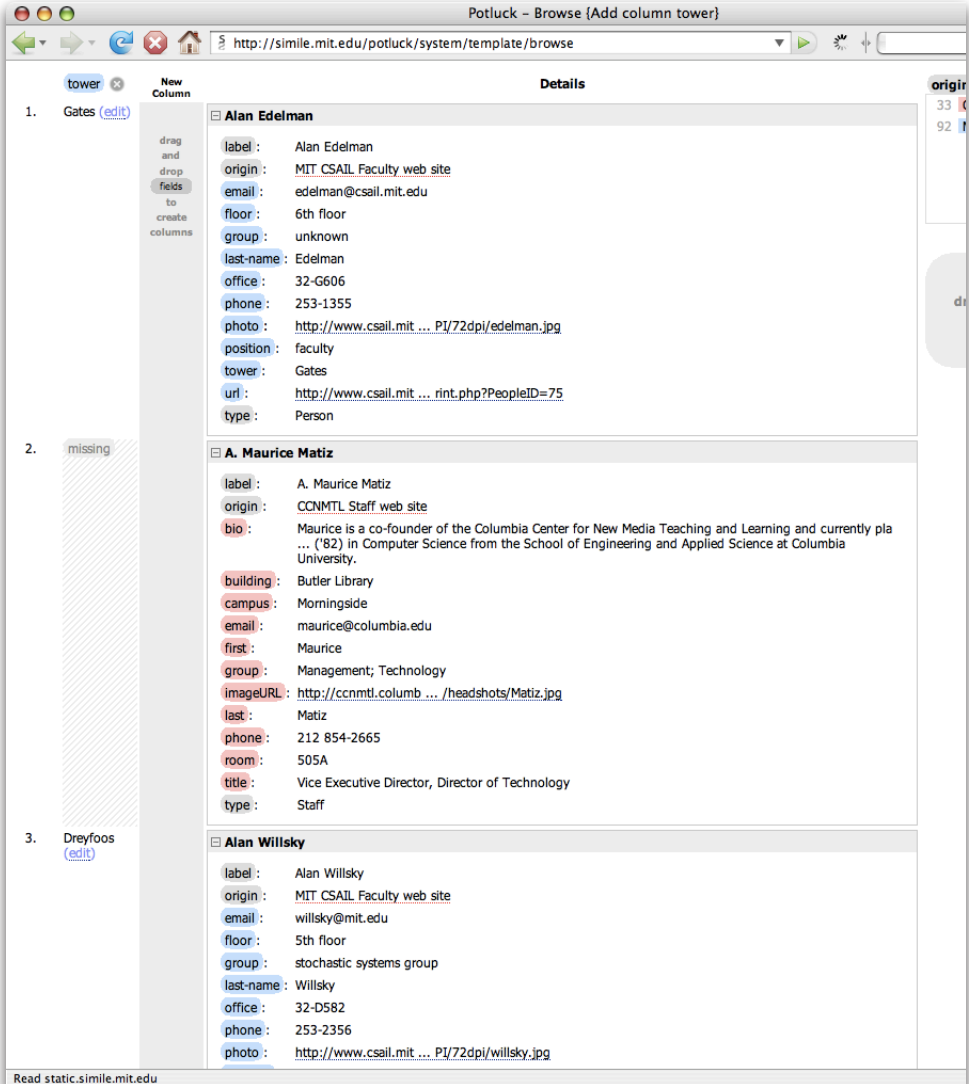


Figure 5.2. Potluck’s user interface shows data that has just been mixed together but not yet processed by the user. Fields are rendered as draggable “field tags,” color-coded to indicate their origins. There are two drop target areas for creating columns and facets.

5.2.1 Creating columns and facets

A field tag can be dragged and dropped onto the gray column to the left (Figure 5.2) to create a new column listing that field, or onto the gray box to the right to create a facet for filtering by that field. Figure 5.3 shows a newly created column. A column or facet can be moved by dragging its field tag and dropping the tag between other columns or facets. Deleting a column or facet (by clicking its ) removes the column or facet from the display but does not delete the corresponding field's data.



The screenshot shows a web browser window titled "Potluck - Browse {Add column tower}" with the URL "http://simile.mit.edu/potluck/system/template/browse". The interface displays a table with three rows. The first row is for "Gates" (edit), the second is "missing", and the third is "Dreyfoos" (edit). A "New Column" panel on the left shows a "tower" tag being dragged. The "Details" panel on the right shows the following data for each record:

Record	label	origin	email	floor	group	last-name	office	phone	photo	position	tower	uri	type
1. Gates	Alan Edelman	MIT CSAIL Faculty web site	edelman@csail.mit.edu	6th floor	unknown	Edelman	32-G606	253-1355	http://www.csail.mit.edu/PI/72dpi/edelman.jpg	faculty	Gates	http://www.csail.mit.edu/rint.php?PeopleID=75	Person
2. missing	A. Maurice Matiz	CCNMTL Staff web site	maurice@columbia.edu		Management; Technology	Matiz		212 854-2665	http://ccnmtl.columbia.edu/headshots/Matiz.jpg				Staff
3. Dreyfoos	Alan Willisky	MIT CSAIL Faculty web site	willisky@mit.edu	5th floor	stochastic systems group	Willisky	32-D582	253-2356	http://www.csail.mit.edu/PI/72dpi/willisky.jpg				

Figure 5.3. Potluck renders a new column to the left when `tower` is dropped into the `New Column` drop target. Since the second record is not from the same origin as the dropped field, its cell in that column shows `missing`.

5.2.2 Merging fields

A field tag can be dropped onto an existing column or facet in order to make that column or facet contain data for both the original field and the newly dropped field. Such an operation creates a *merged field*, whose field tag is rendered as a visual juxtaposition of the original tags, taking on a pill-shaped form `position title`. Figure 5.4 shows several columns and facets of merged fields. Merged field tags can be dragged and dropped just like elemental field tags can in order to create new columns and facets, or to merge into other existing columns and facets.

Creating a merged field does not disturb the elemental fields. Thus, in the scenario, it would be easy to have `gen:mother` and `gen:father` merged together for one purpose while keeping them separate for another purpose, all at the same time.

The screenshot shows the Potluck web interface. At the top, it says "Potluck - Browse [select 'Gates']" and the URL is "http://simile.mit.edu/potluck/system/template/browse". Below the header, it lists data origins: "MIT CSAIL Faculty web site" and "CCNMTL Staff web site". There are buttons for "Export tsv" and "Export json". A notification says "You've just done: select 'Gates'" with an "Undo" button. Below that, it says "VIEW AS: TABLE • MAP".

The main content area shows a table with columns: "tower building", "photo imageURL", "email email", and "phone phone". The records are:

Record	tower building	photo imageURL	email email	phone phone
1. Gates (edit)			edelman@csail.mit.edu (edit)	253-1355 (edit)
2. Lewisohn Hall (edit)			acox@ccnmtl.columbia.edu (edit)	212 854-1851 (edit)
3. Gates (edit)			meyer@csail.mit.edu (edit)	253-6024 (edit)
4. Lewisohn Hall (edit)			fmoretti@columbia.edu (edit)	212 854-1692 (edit)
5. Gates (edit)			agarwal@csail.mit.edu (edit)	253-1448 (edit)

On the right side, there are "Details" and "Facets" panels. The "Details" panel shows the details for the selected record, "Alan Edelman". The "Facets" panel shows a list of facets: "origin", "tower building", and "group group". The "tower building" facet is selected and shows a list of buildings with checkboxes: "3 Armory", "20 Butler Library", "22 Common/other", "11 Dreyfoos", "59 Gates", and "10 Lewisohn Hall". The "group group" facet shows a list of groups with checkboxes: "9 algorithms", "3 complexity theory group", "4 computation structures group", "1 computational biology group", "1 computational cognitive science group", "1 computational genomics group", "3 computer architecture group", "4 cryptography and information security group", "5 database group", "2 decentralized information group", "7 Education", "1 haystack", and "1 learning rich, testable models".

Figure 5.4. A screen shot of Potluck showing several columns and facets of merged fields. The records' details have been collapsed to make space for the columns.

Furthermore, the merging operation is *not* transitive, so that, say, merging fields **mother** and **father** together (to mean **parent**) and then **mother** and **grandmother** together (to mean **female ancestor**) does *not* force all three fields to be merged into **mother/father/grandmother**.

5.2.3 Simultaneous editing

The **edit** link next to each field value opens up the Simultaneous Editing dialog box where the values of that field can be edited *en masse* (Figure 5.5). The concept of simultaneous editing originated from LAPIS [58], a text editor that displays several keyboard cursors simultaneously on a text document, generalizes the user’s editing actions at one cursor, and applies them to the text at the rest of the cursors. Based on the user’s mouse clicks, LAPIS guesses how to divide the text document into records (often into lines or paragraphs) and where the cursors should be placed within those records (e.g., after the second word of the third sentence in each paragraph). Whereas LAPIS has to guess what a record is for the purpose of simultaneous editing, Potluck already has the field values conveniently separate. Potluck groups field values into columns by structural similarity, e.g., the phone numbers in the second column all have area code 212. These columns serve to visually separate out values of different forms, call out outliers (such as “Easter Day 1824” in the scenario), and let the user edit different forms differently. The user can click on any field value to give it keyboard focus, and editing changes made to it are applied to other values in the same column in a similar fashion. The multiple cursors in Figure 5.5 give visual feedback of the simultaneous editing operations in progress.

If a value appears in several records it is shown in only one entry in the dialog box. In the scenario, if the nickname “Bill” appears in three records, the user can

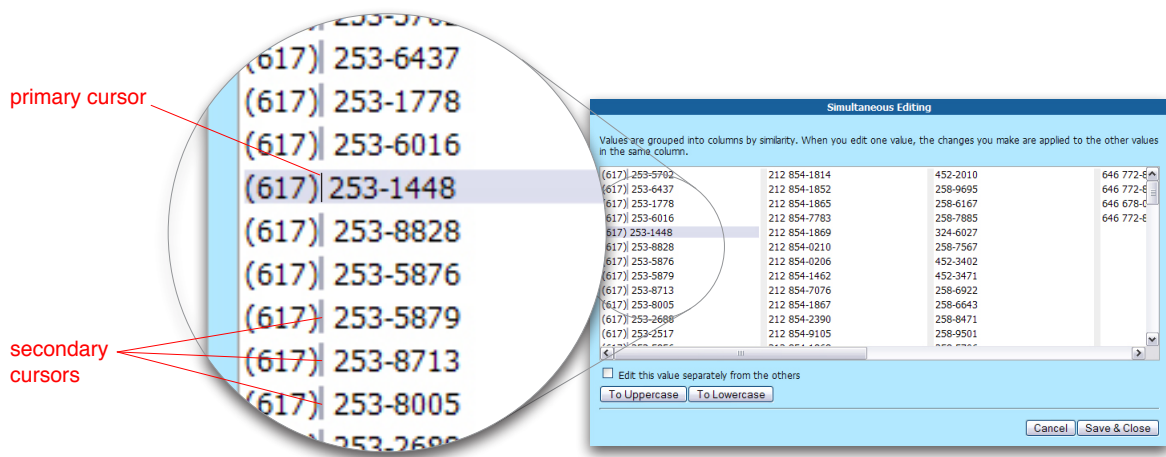


Figure 5.5. Potluck’s Simultaneous Editing dialog box lets the user change several similar values simultaneously by editing any one of them. Multiple keyboard cursors are shown and any editing change to the focused value is immediately reflected in the other values.

click on its single entry in the dialog box, set the checkbox **Edit this value separately from the others**, and change it to “William” to correct all three records.

Simultaneous editing is useful for correcting inconsistencies between data sets that occur many times, such as prefixing area codes to phone numbers and wrapping existing area codes in parentheses. It is also useful for reformatting a field, such as changing “first-name last-name” into “last-name, first-name”, and for making a new field out of an existing field, such as extracting building numbers (32) from within office numbers (32-582).

5.2.4 Faceted browsing

Faceted browsing [76] is a browsing paradigm in which a set of records can be filtered progressively along several dimensions in any arbitrary order. For example, a set of recipes can be filtered by picking an ingredient first, a cooking method second, and a cuisine finally, or by picking a cuisine first, then an ingredient, and a cooking method finally depending on which order suits the user best. Because the data Potluck handles is often multidimensional, faceted browsing is useful in Potluck as it is designed for exploring multidimensional data in flexible, user-controllable ways. Exploration is needed for identifying and selecting out just the subset of data that is useful as well as for isolating on records that need cleaning up. All faceted browsers so far work on single data sets. Potluck extends faceted browsing for the data integration task in which data arrives from many sources.

If within a facet there are records for which the corresponding field is missing, the facet explicitly shows a choice for filtering to those records (Figure 5.6). This visual element, not present in conventional faceted browsing interfaces, also serves

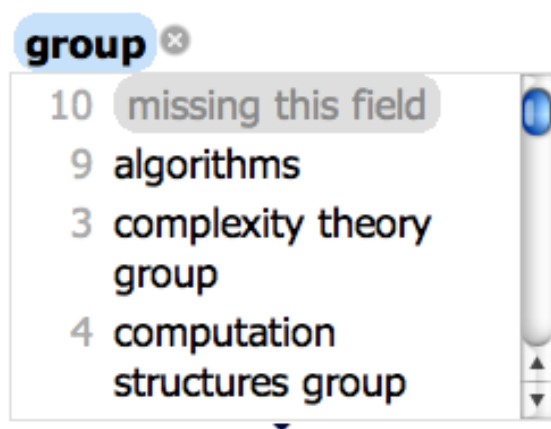


Figure 5.6. If inside a facet there are records for which the corresponding field is missing, the facet shows `missing this field` as a choice so that the user can get to those records.

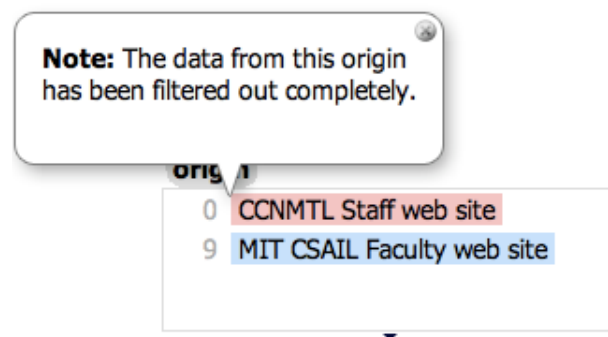


Figure 5.7. The origin facet does not remove choices for which there are no records. Moreover, it pops up messages to call the user’s attention to those filtered out origins.

to remind the user that, if that field is an elemental field instead of a merged field, the field is not present for records in other data sets.

While working with multiple data sets at the same time, it can be easy to forget that an elemental field from one data set does not exist in the others. Whenever a facet choice causes all records from an origin to be filtered out completely, that origin remains in the origin facet and a message is popped up drawing the user's attention to it (Figure 5.7).

5.2.5 Visualizations

Potluck currently provides two visualizations: a tabular view and a map view. Figure 5.8 shows the map view in which any field containing street addresses or latitude/longitude pairs can be dropped onto the map view to plot the records. The map markers can also be color-coded using drag and drop. Faceted browsing is supported concurrently so that the user can construct a map while browsing through the data at the same time.

5.2.6 Miscellany

Potluck provides drop down menus on left clicks as alternatives to drag and drop in order to increase the likelihood that the user succeeds at finding some way to accomplish a task. The browser's Back and Forward buttons can be used to redo and undo user actions. Like contemporary highly interactive web interfaces, Potluck also shows the most recently done or undone action and provides a link to undo or redo it.

5.3 Implementation

Potluck consists of two components: a server-side component implemented as a Java servlet, responsible for retrieving the data within the Exhibit-embedding web pages to mix; and a client-side component implemented in Javascript on top of the Exhibit API, responsible for all the user interface interactivity.

Merged fields are implemented as query unions: when the values of a merged field are requested, the values of each elemental field in that merged field are returned in a single result set. No equivalence is added into the data model so that merging operations will not be transitive and so that the original elemental fields can still be used in isolation even after they have been merged.

Simultaneous editing is implemented in Javascript. Each field value is parsed into a sequence of features. Features are runs of digits, of letters, or of white spaces, or individual punctuation marks and symbols. For example, "733-3647" is broken down into three features: the run of digits "733", the symbol "-", and the run

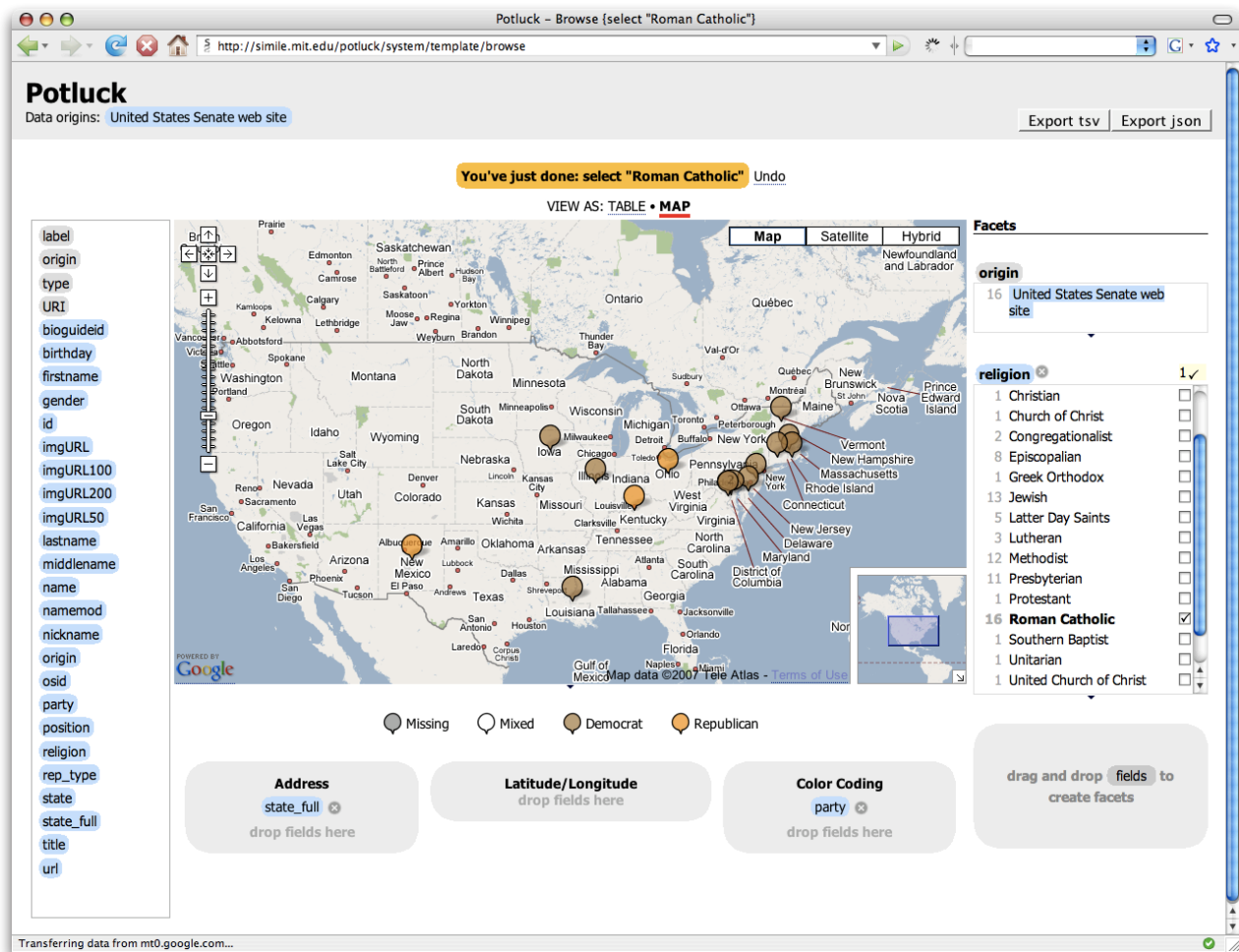


Figure 5.8. Potluck’s map view allows plotting and color-coding records by dropping field tags into drop target areas. Faceted browsing is also offered during map construction.

of digits “3647”. Field values are then clustered into columns by greedily aligning these sequences of features.

As the user moves the keyboard cursor, makes selections, and edits the text of one value, the cursor positions are generalized to be relative to the features of the field value being edited (e.g., “second character from the beginning of the third last feature”), and then those generalized cursor positions are turned into absolute cursor positions of each of the other field values in the same cluster and used to apply the edit. Secondary cursors are rendered using colored `` elements.

As the clipboard Cut and Paste operations cannot be reliably detected, cut-and-paste must be supported in simultaneous editing using a trick. When some text is inserted, if that same piece of text has previously been deleted in one edit action, it is assumed that what has taken place is a cut-and-paste operation. Note that this trick works only for cut-and-paste, not copy-and-paste.

5.4 Evaluation

A user study on Potluck has been conducted to ascertain whether people could learn how to use Potluck as well as to discover usability problems. Another purpose was to observe how people use Potluck in an open-ended task using their own judgement about which fields to merge and edit, and how to display them, so as to determine if casual users could actually perform data integration themselves.

5.4.1 Design and Procedure

This study consists of two tasks: a structured task during which the subjects performed simple steps to familiarize themselves with Potluck, and an unstructured task during which the subjects performed an open ended task based on the skills they had just acquired.

In Task #1, subjects browsed two web pages containing information about 92 people in a lab and 33 people in another lab, and answered questions about these people in ways that required the pages' faceted browsing features (e.g., "how many people are in the Gates tower?"). This warm-up exercise let the subjects learn about the data and about faceted browsing. Then the subjects were asked to use Potluck to mix the data in those web pages and to achieve the following goals (quoted almost verbatim from the study's instructions):

- create a column listing the buildings where people work and make sure the column is filled in with information for people from both labs;
- create a column listing people's phone numbers and edit them to have the form (xxx) xxx-xxxx, using 617 for phone numbers without area code;
- create a column listing people's job titles;
- create a facet of people's job titles, use it to filter for people in directing positions (directors and co-directors), and determine how many such people there are in each lab; and
- create a column of people's last names and sort it in ascending order.

These instructions were *not* worded in low-level details (e.g., click this button) so to allow the subjects the opportunities to learn how to use Potluck's user interface by themselves and to allow us the chance to discover usability problems.

In Task #2, the subjects were asked to use Potluck to mix data from two Exhibit-powered web pages of 40 + 55 publications and then mock up a single web page where hypothetical visitors could conveniently sort and filter through all of those publications as if the data came from a single source. The subjects were left to their own discretion to decide which columns and facets to create, although some examples were given in case the subjects were not familiar with the domain.

5.4.2 Participants and Apparatus

Six subjects (2 male, 4 female) from a university community were recruited by sending an e-mail message to a mailing list and posting paper ads around our college campus. Four were younger than 30, and two older than 30. They were two students (mechanical engineering and computer science), two researchers (applied math and brain and cognitive science), a lawyer, and an applied math consultant.

Five subjects (1 male, 4 female) were also recruited from a local campus' libraries, who worked with data in their daily job. Two were in their 20s, one 30s, and two 40s. There was a desire to observe if librarians, who have more experience working with data, would use Potluck differently.

There were a total of 11 subjects, referred to as G1 to G6 from the general university population and L1 to L5 from the libraries. All browsed the Web at least a few times a day and used Firefox as one of their primary browsers.

Subjects received \$10 each for participating in a 30 – 45 minute study session. All sessions were conducted by one investigator on a single computer (Pentium 4 2.53GHz, 1.00GB) with an 18" LCD flat panel at 1600×1200 resolution in 32-bit color and a Dell two-button mouse with wheel, running Microsoft Windows XP. The study facilitator observed the subjects and took written notes.

5.4.3 Results

All subjects were able to learn Potluck's user interface with little guidance and to complete the user study's tasks within 45 minutes. We now report the results in more details and point out usability issues to address in the future.

5.4.3.1 Columns

Nine subjects out of 11 used only drag and drop to create columns. This indicates that the relevant visual cues were sufficiently strong. One of the other two subjects, G5, used the Create Column menu command at first but adopted drag and drop later. L1 used only the menu command.

G5 and L5 had difficulty understanding that dragging a field tag to create a column automatically filled up the whole column with data wherever the field was available. They continued to drag the same field tag out again and again for each row, paying no attention to the data already shown in the column. The drag feedback can be improved to better indicate that the whole field is being dragged, such as showing ghosted images of several field values near the mouse pointer.

All except one subject merged columns using drag and drop; G2 used the corresponding menu command. G3 and G4 expected the phone fields from both sources in Task #1 to be merged automatically. Potluck can be made to suggest such merging if the field names match precisely.

Most subjects merged **position** and **title** together into one column, but one subject also included **group** to more fully qualify **position**. This was because most **title** values were more specific than most **position** values (e.g., "Codirector of Marketing" vs. "professor"). This operation was actually not what the subject in-

tended (as he verbalized): the operation performed a set union of two fields instead of a string concatenation. But as Potluck rendered the **group** value after the **position** value for each record (e.g., “professor, computer architecture”), the visual outcome looked right and the subject was contented. However, sorting on this merged field would produce random orders and a facet created out of this merged field would list the **group** and **position** values separately, not paired together. Potluck should support string concatenation and suggest it as an alternative to merging whenever the two fields involved come from the same source. Note that in the scenario in section 2, concatenation is probably not the desired choice when the **gen:mother** field is dropped onto the **gen:father** field even though both come from the same source. This is why heuristics should only be used to make suggestions, not to automate.

5.4.3.2 Facets

All subjects used drag and drop to create facets. Two subjects initially created facets using the corresponding menu command, but they discovered the drag and drop alternative and did not revert to the menu. Merging facets was done solely using drag and drop. Note that the field tags on facets do not offer any menu (an oversight in our implementation); only field tags in the details column and in the column headers support menus.

Some subjects tended to drag already merged field tags from columns to create facets while the others dragged elemental field tags from the Details column to create merged facets. The latter behavior forced the user to re-merge fields she has already merged; this is both inefficient and error-prone as some subjects did forget to re-merge fields. Potluck should have automatically suggested or defaulted to the merged field whenever an elemental field that has been merged is used.

G4 did not initially merge facets in Task #1 to filter for people in directing positions. Instead, he created two facets, **position** and **title**, from the two sources separately and used missing this field to achieve the goal. In either facet, he selected directing positions as well as missing this field so that records in the other source were not excluded. This required on his part deeper understanding of how faceted browsing worked. When asked to achieve the goal without using missing this field, he discovered that he could merge facets.

5.4.3.3 Simultaneous editing

All subjects were able to edit several phone numbers using the simultaneous editing feature. G1 anticipated this feature even before clicking **edit**, asking out loud, “can I edit them all together?” She later used the feature to delete first names from people’s full names to get a field of last names. This action properly utilized the simultaneous editing feature’s power but destroyed data (the first names). Potluck can be made to alert the user of this loss and offer a convenient way to apply the edit on a copy of the original field instead.

G4 tried to move the leading “A” from publication titles to the end (e.g., “Tale of Two Cities, A”) using simultaneous editing (a reasonable goal) but the facilitator explained that the feature did not support that case. L2 and G6 tried to swap first

names and last names so that publications could be sorted by their authors' last names. L2 selected a last name in the simultaneous editing dialog box and dragged it to the front of the corresponding first name; unfortunately, a bug prevented this from working. G6 used keyboard shortcuts for cut-and-paste and succeeded. These subjects' actions indicated some intuitiveness in using cut-and-paste and drag-and-drop for simultaneous editing.

G3 expressed that she did not want to see all phone numbers in the simultaneous editing dialog box but only their templates. G5 and L3 edited only the first group of phone numbers, and L4 edited only the first and third groups, neglecting the groups that were not scrolled into view. To avoid such oversight, which pieces of data an edit does and does not affect must be made apparent.

5.4.3.4 *Librarians vs. general subjects*

Among the five librarians, four were catalogers (who characterize physical artifacts such as books and enter their metadata into databases), and one was a programmer responsible for integrating large data sets. While the catalogers showed no significant difference with the general subjects in their use of Potluck, the programmer, L1, was clearly an outlier: he created 10 columns and 7 facets in total. He was very excited about the user interface of Potluck and described his data integration work, consisting of manual data entry and Perl scripting, to be tedious and painful.

G6, who also needed programming skills to deal with some data for his work, expressed equal enthusiasm for Potluck. He used simultaneous editing to swap first name and last name. Thus, while there was no noticeable difference between the subjects from the general population and the librarians, who purportedly work with data and metadata on a daily basis, there was a difference between programmers and non-programmers in how much they appreciated Potluck. Programmers, who have encountered difficulties in dealing with data even with their programming skills, appreciated Potluck more. Non-programmers accomplished the tasks in the study equally well, but were not equally excited perhaps because there was not enough reusable data on the Web for them to feel the need to integrate data themselves. However, when there will be more reusable data in the future, Potluck will level the playing field for non-programmers, making them as effective as programmers for the task of integrating data.

5.5 Summary

This chapter presented several techniques embodied in a tool called Potluck for letting for casual users—those without programming skills and data modeling expertise—integrate data by themselves and obtain usefulness from the integrated data. Potluck is novel in its use of drag and drop for merging fields, its integration and extension of the faceted browsing paradigm for focusing on subsets of data to align, and its application of the simultaneous editing technique for cleaning up data syntactically. Potluck lets the user construct rich visualizations of data in-place as

5. INTEGRATING DATA

the user aligns and cleans up the data. This iterative process of integrating the data while constructing useful visualizations is desirable when the user is unfamiliar with the data at the beginning—a common case—and wishes to get immediate value out of the data without having to spend the overhead of completely and perfectly integrating the data first. A user study on Potluck indicated that it was usable and learnable, and even solicited excitement from programmers who had experienced great difficulties in integrating data even with their programming skills.