

Senbazuru: A Prototype Spreadsheet Database Management System

Zhe Chen, Michael Cafarella, Jun Chen, Daniel Prevo, Junfeng Zhuang
Computer Science Department, University of Michigan

Abstract

Senbazuru is a prototype spreadsheet database management system (SSDBMS). It is able to extract relational information from spreadsheets, which opens up opportunities for integration among spreadsheets and with other relational sources.

We demonstrate that Senbazuru allows users to search for relevant spreadsheets in a large corpus, probabilistically constructs a relational version of the data, and offers several relational operations over the resulting extracted data (including joins to other spreadsheet data).

Our demonstration is available on two clients: a JavaScript website and a touch interface on the iPad.

Interface



Search -- Using a textual search-and-rank interface, the search component allows a user to quickly locate relevant datasets in a huge Web spreadsheet corpus.

Extract -- Spreadsheets often exhibit a implicit hierarchical structure between attributes and values. The extract component automatically infer the implicit structure of spreadsheets based on a probabilistic model.

5				
6	Sex, age, and race	1990 \1	2000	
7				
19	Total smokers \3	25.5	23.2	
20	Male, total	28.4	25.6	
21	18 to 24 years	26.6	28.1	
22	25 to 34 years	31.6	28.9	
23	35 to 44 years	34.5	30.2	
24	45 to 64 years	29.3	26.4	
25	65 years and over	14.6	10.2	
26	White, total	24.0	25.7	
27	18 to 24 years	21.4	30.4	
28	25 to 34 years	31.6	29.7	
29	35 to 44 years	33.5	30.6	
30	45 to 64 years	28.7	25.8	
31	65 years and over	13.7	9.8	
32	Black, total	32.5	26.2	
33	18 to 24 years	21.3	20.9	
34	25 to 34 years	33.8	23.2	
35	35 to 44 years	42.0	30.7	
36	45 to 64 years	36.7	32.2	
37	65 years and over	21.5	14.2	

Repair -- Automatic extraction often emits errors, but our interface allows users to manually repair extraction errors effectively and efficiently.

Spreadsheet

Data Tree

Raw Tuples

Relational Table

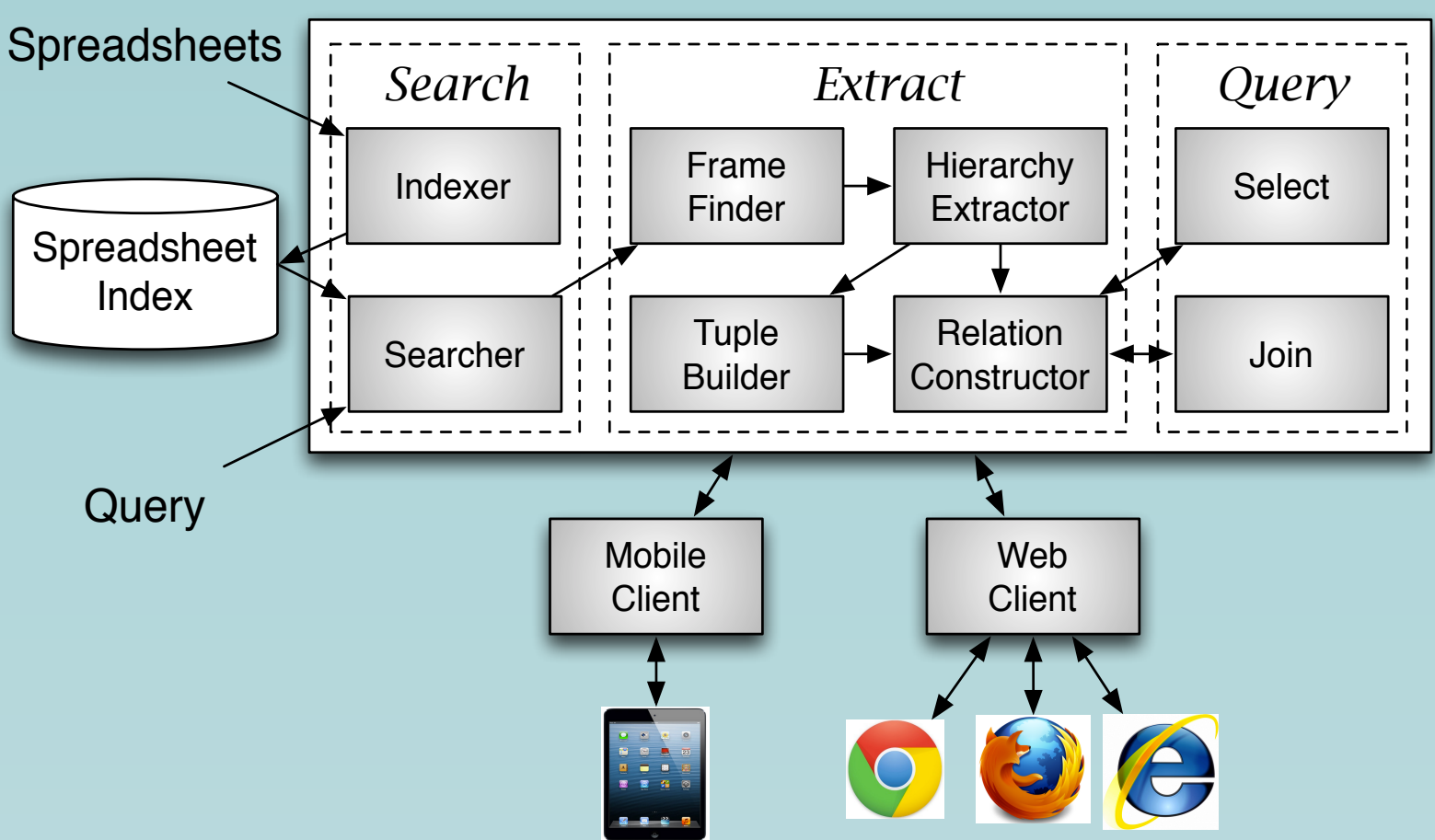
Filter

Join

Female, total,	Black, total,	18 to 24 years,	1965.0,	37.1
Female, total,	Black, total,	18 to 24 years,	1974.0,	35.6
Female, total,	Black, total,	18 to 24 years,	1979.0,	31.8
Female, total,	Black, total,	18 to 24 years,	1983.0,	32.0
Female, total,	Black, total,	18 to 24 years,	1985.0,	23.7
Female, total,	Black, total,	18 to 24 years,	1987.0,	20.4
Female, total,	Black, total,	18 to 24 years,	1988.0,	21.8

Integrate -- Users are also able to integrate two spreadsheets based on the derived relational tables.

Framework



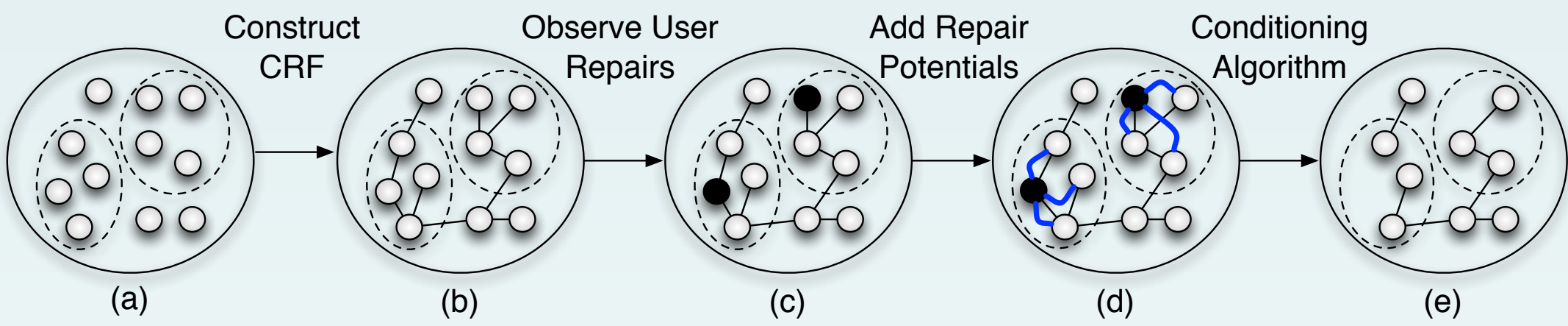
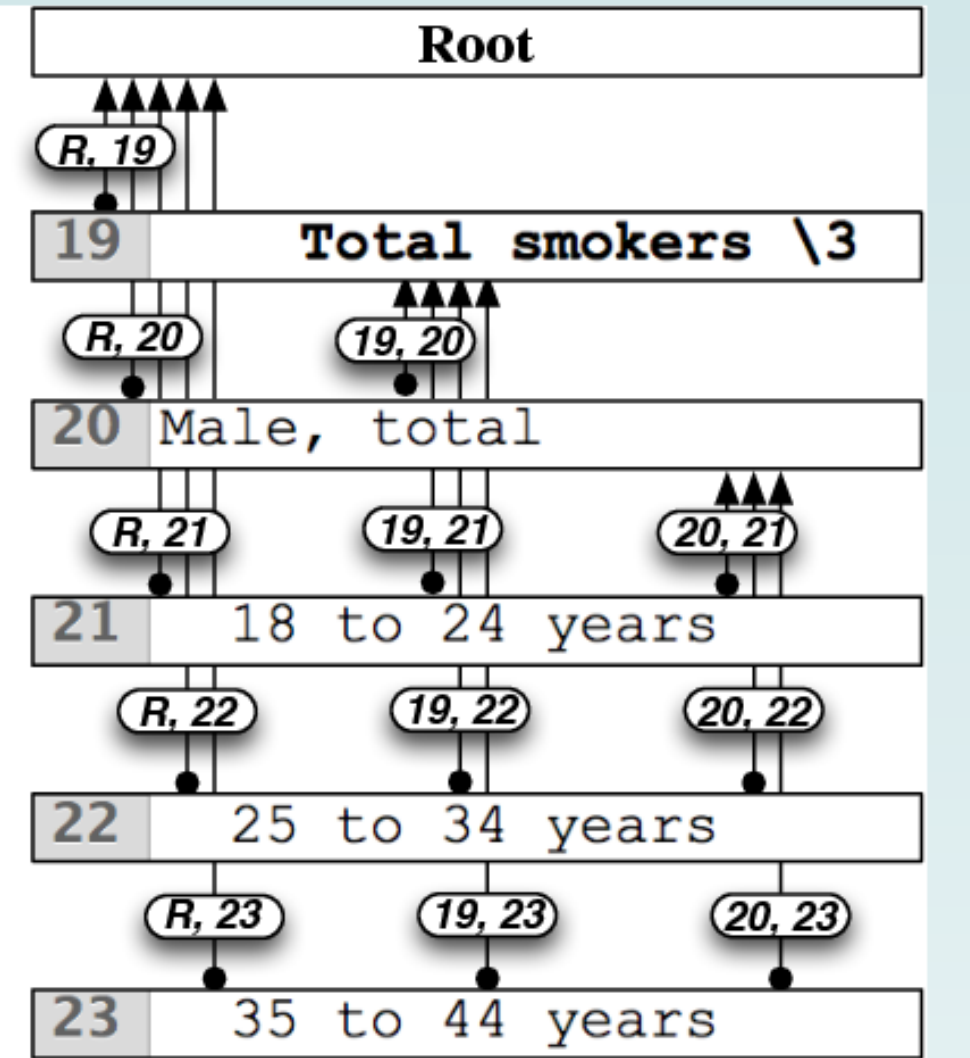
Senbazuru consists of three functional components we view as critical for a useful SSDBMS:

- Search** -- For each spreadsheet in the corpus, the **indexer** extract text from each cell and uses Lucene to index the text. When a query arrives, the **searcher** uses the inverted index to find relevant datasets.
- Extract** -- The extract component is composed of a background extraction pipeline that automatically obtains relational data from spreadsheets, and a repair interface that allows users to manually repair extraction errors.

- Frame Finder** -- It identifies the attribute and value regions in a data frame spreadsheet.

6	Sex, age, and race	1990 \1	2000
7			
19	Total smokers \3	25.5	23.2
20	Male, total	28.4	25.6
21	18 to 24 years	26.6	28.1
22	25 to 34 years	31.6	28.9
23	35 to 44 years	34.5	30.2
24	45 to 64 years	29.3	26.4
25	65 years and over	14.6	10.2
26	White, total	28.0	25.7

- Hierarchy Extractor** -- It recovers the implicit attribute hierarchies for each attribute region.
 - Create random variables.
 - Establish the correlation between variables.
 - Incorporate user repairs.



- Tuple Builder** -- It generates a relational tuple for each value in the value region.

1990	Male	White	45 to 64 years	65 years and over	28.7
------	------	-------	----------------	-------------------	------

- Relation Constructor** -- It assembles the relational tuples into relational tables.
- Query** -- The query component supports basic relational operators, especially **selection** and **join**, which the user can apply to spreadsheet-derived relations.

