

Doctoral Dissertation

Curran Kelleher

4/29/2014

Abstract

There is immense potential value in data that is not being realized. While many data sets are available, it is difficult to realize their full value because they are made available using many different formats, protocols and vocabularies. The heterogeneity of formats, protocols and vocabularies makes it difficult to combine data sets together and hinders the development of reusable data visualization software. While it is straightforward to produce static visualizations of just about any data set by customizing existing examples, there is a lack of generalized visualization software that supports the creation of interactive visualizations and visualization dashboards with multiple linked views. The contribution of this dissertation is a collection of data structures and algorithms supporting integration and interactive visualization of many data sets using interactive visualization dashboards with multiple linked views. A proof of concept implementation demonstrates support for several public data sets and well known visualization techniques.

Contents

1	Introduction	3
1.1	Related Work	8
1.1.1	Data Representation	8
1.1.2	Data Integration	9
1.1.3	Data Visualization	10
1.1.4	Web Graphics Technology	12
1.2	Pseudocode Conventions	15
2	Functional Reactive Models	21
2.1	The Model Data Structure	22
2.2	Functional Reactive Change Propagation	25
2.3	Functional Reactive Visualizations	28
3	The Universal Data Cube	30
4	Data Sets	31
4.1	United Nations Population Estimates	31
4.2	United Nations Millenium Development Goals	31
4.3	United States Census Population Estimates	31
4.4	US Central Intelligence Agency World Factbook	32
4.5	US Centers for Disease Control Causes of Death	32
4.6	W3Schools Browser Market Share	33
4.7	Natural Earth	33
5	Visualizations	34
5.1	Data Cubes and Visualization Theory	34
5.2	Bar Chart	35
5.3	Scatter Plot	35

5.4	Line Chart	35
5.5	Choropleth Map	35
5.6	Stacked Area Plot	35
5.7	TreeMap	35
5.8	Node Link Tree	35
5.9	Radial Tree	35
5.10	Icicle Plot	35
5.11	TreeMap	35
5.12	Parallel Coordinates	35
6	Visualization Dashboard Infrastructure	36
6.1	Dashboard Layout using Nested Boxes	36
6.2	Multiple Linked Views	36
6.3	Dynamic Dashboard Configuration	36
7	UDC Visualization Dashboards	37
8	Collaboration and History Navigation	38

Chapter 1

Introduction

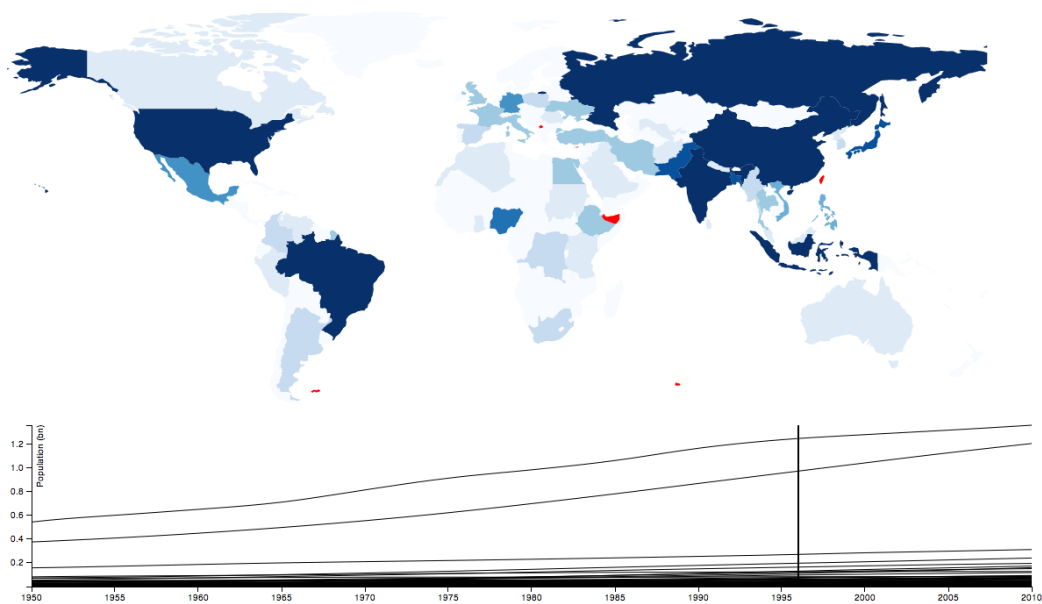
The contributions of this dissertation are novel data structures and algorithms for integration and interactive visualization of many data sets from multiple sources, based on the data cube concept. The proposed data representation framework will allow data sets to be combined together and visualized using interactive visualization dashboards, giving users the sense that the data exists within a single unified structure. The framework is designed to be able to represent and integrate an arbitrary number of data sets created independently of one another, and expose the integrated structure to reusable visualizations. The reusable visualizations can be combined together in dashboard layouts with multiple linked views using existing interaction techniques such as brushing and linking. The proposed data representation and visualization framework is fundamentally new, and will allow heterogeneous data sets to be explored in a unified way that was never before possible.

The overall goal of this work is to build digital telescope into the universe of phenomena on Earth via publicly available data. For example, consider public data sources such as the United Nations, the US Census, the US Bureau of Labor Statistics, or the US Centers for Disease Control. These organizations and hundreds of others around the world provide publicly available data about various topics including population statistics, public health, distribution of wealth, quality of life, economics, the environment, and many others. By unifying these data sources and providing users with tools to explore and present the data visually, a deeper understanding of the world can be gleaned through the lens of public data. The focus of this dissertation is on applications involving public data, however the techniques introduced can be applied to any data sets that can be conceptually modeled as data cubes, regardless of whether they are public or private.

Consider the data from the US Census that covers population statistics for US States from 1950 to 2010. Consider also population statistics from the United Nations covering World Countries from 1970 to 2012. These two data sets may use different identifiers for years and geographic regions, but they cover an overlapping conceptual data space of time, geography and population. From these two data sets it is possible to create a visualization dashboard with a map of the world showing population as color and a corresponding line graph showing population for each region as lines. If the user views the whole world, the UN population data is shown for each country. If the user zooms into the US, US Census data is shown for each state. If

the user selects a point of time in the line graph, the data shown on the map is from that point in time. If the user pans and zooms on the map, the lines in the line graph update to only show the regions visible on the map. This is one example of an interactive visualization dashboard with multiple linked views (the timeline and map views) operating over multiple data sets integrated from different sources (the United Nations population data and the US Census population data).

Figure 1.1: An example visualization dashboard with multiple linked views. Selecting a time on the line chart slices the data displayed on the choropleth map. Zooming and panning on the choropleth map filter the country profiles shown in the line chart. The vertical black line on the line chart indicates the selected year.



TODO Improve this visualization - add color legend, add names to lines

Data cubes, also known as OLAP (OnLine Analytical Processing) cubes, can represent data that contains measures aggregated (typically using sum or average) along categorical hierarchies. The data cube concept emerged from the field of data warehousing as a way to summarize transactional data, allowing analysts to get a bird’s eye view of company activities. The term OLAP stands in contrast to the term OLTP (OnLine Transaction Processing), which is the part of the data warehouse system that ingests and stores data at the level of individual transactions or events. After the ETL (Extract, Transform and Load) phase of the data warehouse flow, the data is analyzed by computing a data cube from the transactional data.

The data cube concept and structure can be used to model existing data sets as well. Publicly available data sets (often termed “statistical data”) may be considered as pre-computed data cubes if they contain aggregated measures (also called “indicators”, “metrics” or “statistics”) across time, geographic space or other dimensions such as gender, age range, ethnicity or industry sector. Any categorization scheme containing distinct entities, organized as an unsorted collection, an ordered collection, or a hierarchy can be modeled as a dimension. Any numeric value that represents an aggregated statistical summary using sum, average, or other aggregation operator can be modeled as a measure.

With this approach, it is possible to model many data sets together using shared dimensions and measures. This will allow integration of many data sets together in a single unified structure. Existing data cube technologies

assume that data cubes will be computed from a relational source, and are not designed to handle integration of pre-computed data cubes that may use inconsistent identifiers for common dimensions and inconsistent scaling factors for common measures. Therefore the application of the data cube concept to integration and visualization of many pre-computed data cubes, while theoretically plausible, requires the development of novel data structures and algorithms that extend the data cube model to handle integration of pre-computed data cubes that may use inconsistent identifiers for common dimensions and inconsistent scaling factors for common measures.

The data cube structure lends itself particularly well to visualization. Long standing perception-based data visualization theory presented by Bertin [6] and Mackinlay [60] identify effective ways to visually encode data based on data fields that are nominal, ordinal or quantitative. Visualization techniques have been explored for hierarchical (tree-based) data as well [34]. Data cubes can contain data of these types. Therefore existing visualization theory can be applied to the data cube model to determine which visualizations are appropriate for representing which data, depending on its data cube structure. This topic is discussed in section ??.

1.1 Related Work

Previous work relating to this dissertation falls into four major categories:

- data representation - structures, models and formats for data
- data integration - merging data from many sources
- data visualization - transforming data into interactive graphics
- Web graphics technology - HTML5 graphics APIs and libraries

1.1.1 Data Representation

Relational database systems provide a mature data management solution and are widely adopted [75]. The relational model has well understood theoretical underpinnings such as relational algebra [15]. Data warehouse systems are typically built on the relational model and are augmented by multi-scale aggregated data structures called data cubes, also known as OLAP (OnLine Analytical Processing) cubes [35, 16]. Data cubes contain summaries of the collection of facts stored in a relational database [11]. For example, a data cube may contain how much profit was made from month to month subdivided by product category, while the relational database may contain the information associated with each individual transaction.

Because data cubes provide an abstraction that handles aggregation, they are a widely used method of data abstraction for supporting visualization

and analysis tasks [86]. Kimball pioneered the area of “Dimensional Modeling”, which concerns constructing data warehouse schemas amenable to data cube construction and analysis [50]. Data cubes have been implemented in a variety of different systems, so effort has been made to discover unified conceptual or mathematical models that can characterize many implementations [23, 92, 91, 57, 2, 37, 7]. The data cube structure has also been used to model user Web browsing sessions to support data mining algorithms for Web prefetching [101].

1.1.2 Data Integration

The field of data integration offers many techniques for combining data from multiple sources based on the relational model [26] as well as from a theoretical perspective [56, 38, 103]. Schema matching is the area of data integration that concerns semantic matching between the attributes of data tables from different sources [74, 31]. Schema matching may be performed manually, however it must be automated in order to scale to hundreds or thousands of different sources. Numerous approaches for automated schema matching have been proposed [82, 24, 49, 68, 61, 25]. Schema matching approaches aimed specifically at Web and Ontology based data integration have also been proposed [40, 69, 27, 62, 47, 70, 90, 94, 71, 30]. Data matching (also known as record linkage) is the area of data integration focusing on resolving different identifiers to the same real-world entity [99, 100, 52, 3, 36, 98]. Record linkage has been applied extensively to public data [46, 45, 43]. Sev-

eral tools have been introduced that aid users in data integration tasks via a graphical user interface [14, 48, 28]. Techniques from both of these areas must be applied in order to integrate data sets from multiple sources and utilize our proposed unified data model.

1.1.3 Data Visualization

The field of information visualization offers several compelling theoretical approaches for visualizing data. Arguably the first significant work concerning data visualization was William Playfair’s “Commercial and Political Atlas”, published in 1786 [73]. In this work, Playfair introduced the Bar Chart, Pie Chart, and Line Graph. The first attempt at a systematic formalization of data visualization was Jacques Bertin’s “Semiology of Graphics” [6]. In this work, Bertin relates data types to visual marks and channels in a coherent system that takes visual perception into account. Bertin’s work has influenced many future theoretical underpinnings of visualization, including Leland Wilkinson’s “Grammar of Graphics” [97] and Jock Mackinlay’s APT (A Presentation Tool) system [60], which led to the development of the commercial visualization package Tableau [39].

As a more concrete manifestation of visualization theory, much effort has been placed on generating taxonomies of visualization techniques. Chi et al. introduced a concrete taxonomy of visualizations [12] based on the Data State Reference Model [13]. Shneiderman introduced a more general taxonomy based on tasks and data types [81]. Card et al. made steps toward

characterizing the entire design space of data visualizations based on Bertin’s theory [10]. Tufte explored numerous visualization techniques for quantitative information in general, many of which can be applied to visualization of data cubes [89].

Much work has been done regarding visualization of data cubes. Stolte et al. introduced a formalism for defining multi-scale visualizations of data cubes throughout their work on the Polaris system [86, 85, 84]. In this work the authors introduce theoretical underpinnings of a visualization system capable of navigating hierarchical data cubes with a combination of data abstraction and visual abstraction. One fundamental concept in this system is the “zoom graph”, a lattice of data cubes that supports arbitrary interactive zoom paths through the multidimensional data cube hierarchies.

Mansmann coined the term “Visual OLAP”, framed it as a fundamentally new paradigm for exploring multidimensional aggregates [66], explored applications of hierarchical visualization techniques to OLAP cubes [65] and extended Visual OLAP to support irregular hierarchies [64]. Cuzzocrea et al. surveyed the area of data cube visualization in depth [19] and have made several contributions regarding semantics-aware OLAP visualization [21] and a hierarchy driven compression technique for OLAP visualization [20]. Scotch et al. developed and evaluated SOVAT, a Spatial OLAP visualization and analysis tool applied to community health assessments [79, 78]. Lee and Ong introduced a visualisation technique for knowledge discovery in OLAP combining elements of bar charts and parallel coordinates [54]. Maniatis et

al. explored how OLAP cubes can be visualized using TableLens and other techniques [63]. Data cubes have also been utilized as the foundational data structure for several “Big Data” visualization systems [58, 59].

Interactions within data visualization environments have been well studied. Becker et al. investigated brushing in scatter plots [5]. Shneiderman et al. explored dynamic queries in general and how these operations fit into a larger context of visual information seeking [80]. Ward introduced a visualization system based on multiple linked views with direct manipulation techniques including brushing and linking [96]. Anselin discussed how interactive visualization systems with linked views can be applied to Geographic Information Systems [4]. Yi et al. conducted a thorough survey of existing taxonomies for visualization and interactions and developed a set of generalized classes of interactions for visualization [102]. Techapichetvanich et al. explored how visualization interactions pertain to OLAP cubes in particular [87]. Sifer et al. introduced a visual interface utilizing coordinated dimension hierarchies for OLAP cubes [83]. Tegarden formulates some requirements for information visualization relevant for business applications, and highlights some unconventional interactive visualizations with potential application to data cube visualization [88].

1.1.4 Web Graphics Technology

The World Wide Web has evolved to become a full fledged application development platform [72]. HTML5 is the latest set of standards and APIs

(Application Programming Interfaces) from the World Wide Web Consortium that define the capabilities of modern Web browsers [41]. HTML5 contains three graphics technologies that can support interactive Web-based visualizations: Canvas, SVG (Scalable Vector Graphics), and WebGL.

HTML5 Canvas provides a 2D immediate mode graphics API [33]. When using the Canvas API, developers must work with a stateful graphics context by issuing commands to manipulate the raster image of a Canvas element within the HTML page. This approach requires developers to manage rendering logic at a low level and manage data structures that correspond to graphical representations. The Canvas API has seen wide adoption for HTML5-based games, however for visualization applications the higher level SVG API has seen wider adoption.

SVG (Scalable Vector Graphics) provides a 2D retained mode Graphics API [22]. SVG uses the HTML DOM (Document Object Model) to represent the definition of persistent graphical elements. When using SVG, developers need only be concerned with updating the DOM. The SVG engine within the browser is responsible for updating the display to correspond with the SVG DOM. In this way, SVG is a higher level API than Canvas. This makes SVG a preferred platform for developing visualizations. However, SVG is less optimizable than Canvas, because developers do not have access to the rendering logic. SVG has performance limitations relating to DOM manipulation overhead.

WebGL provides a 3D graphics API that is essentially a JavaScript in-

terface to OpenGL ES [67]. OpenGL ES is a subset of OpenGL designed for use in embedded systems and mobile devices. Developers using WebGL must use programming techniques inherited from OpenGL such as buffer management, vertex management, shader definition, 3D projection, and lighting techniques. WebGL enables developers to take advantage of the GPU (Graphics Processing Unit) for massively parallel computation using shaders. WebGL supports high performance 2D and 3D graphics, but is much more complicated to use than Canvas or SVG. WebGL has been applied to interactive visualization of volumetric data [17].

Many high level libraries have been built for supporting use of Canvas, SVG, and WebGL. Three.js is a 3D scene graph library that includes rendering engines for all three graphics technologies [9]. Highcharts is a high level visualization library that provides pre-packaged chart types that can be customized to a limited extent. Leaflet is a library for creating tile-based geographic maps with zooming and panning. hBrowse is a generic framework introduced for Web-based hierarchy visualization [51]. Processing.js is a JavaScript port of the graphics language Processing using HTML5 Canvas. Many more libraries for Web-based graphics and visualization exist, but none have come close to the widespread adoption of D3.js.

D3.js is a flexible and powerful visualization library that uses SVG and has a strong community of users [8]. D3 at its core is a DOM manipulation library with heavy use of functional programming. D3 allows concise declarative statements to define the core logic of visualizations. D3 provides

additional APIs for performing common visualization tasks such as defining and using scales, generating labeled axes, and computing layouts from graphs and trees. D3 is at the center of a vibrant developer ecosystem and has seen wide adoption in industry. There are plentiful examples of D3.js usage for creating visualizations. Many supporting libraries for D3 have been created including Chart.js for composing visualization elements, Crossfilter.js for interactive multidimensional filtering, and DC.js for multiple linked views. Several reusable chart libraries have appeared including NVD3, reD3, and dimple.

Several projects have focused explicitly on visualization of public data on the Web. ManyEyes was an experiment in scaling the audience for visualizations by empowering users to create visualizations of their own data [93]. ManyEyes provided a fixed set of pre-packaged visualization tools and allowed users to visualize their own data tables using the provided visualizations. GapMinder is a project aimed at exposing public data (primarily the United Nations Millenium Development Goals Indicators) using visualization [76]. GapMinder includes an animated scatter plot with an interactive time slider, a line chart showing statistics over time, and a world map (see figure ??). The Google Public Data Explorer provides a visual interface to selected public data sets similar to GapMinder, however it does not make the data available to users in a machine-readable format.

1.2 Pseudocode Conventions

Throughout this document, pseudocode is used to express data structures and algorithms. Our pseudocode is similar to that found in the book “Introduction to Algorithms” [18], but differs significantly in that it uses a functional style. Primitive types in our pseudocode include numbers, strings, booleans, arrays, objects and functions. The following examples demonstrate the features of our pseudocode language.

```
1   $x = 5$ 
```

Line numbers appear to the left of each line of pseudocode. Variables can have any name comprised of characters without spaces, and can be assigned a value with the `=` symbol. Variables need not be explicitly declared. The scope of a variable is determined by where it is first assigned. Our pseudocode uses block scope, meaning that every indentation level introduces a new nested scope. On line 1 of the above pseudocode example, the variable x is defined and assigned the value of 5, a numeric literal.

The following pseudocode demonstrates numbers, strings and booleans.

```
1  myNumber = 5
2  myString = 'test'
3  myBoolean = TRUE
4  myOtherBoolean = FALSE
```

All numbers are treated as double precision floating point. Numeric literals in pseudocode become numbers (see line 1). String literals are denoted by single quotes and a monospace font (see line 2). Booleans can be either true or false. True and false are builtin constant boolean values denoted by all capitalized words (see lines 3 and 4). Camel case names starting with a lower case letter are used for most variables in our pseudocode.

```
1  add =  $\lambda(a, b)$ 
2      return a + b
3  result = add(4, 6) // result is assigned the value 10
4  triple =  $\lambda(x)$  return x * 3
5  triple(3) // evaluates to 9
```

The above pseudocode demonstrates how a function is defined and invoked, and also introduces comments. This example defines a function called *add* that adds two numbers together. The λ symbol defines a new anonymous function. Variables can be assigned functions as values using `=`. The comma separated names in parentheses directly following the λ are the arguments to the function. The pseudocode on lines following the λ that is indented one level constitutes the function body (also called the function closure). The function arguments are only visible inside the function closure.

Functions can be invoked using parentheses. The argument values are passed to the function in a comma separated list within parentheses. On line 3, the *add* function is invoked, passing the value 4 as argument *a* and

6 as argument *b*. The value returned by the function is assigned to the variable *result*. The function invocation causes the function body to execute, which adds the two numbers together and returns the resulting number using the “return” keyword on line 2. Lines 4 and 5 demonstrate that a simple anonymous function can be defined in a single line. Text following the `//` symbol is a comment, and is not executed.

```
1  myArray = []
2  myArray.push(5) // myArray now contains [5]
3  myArray.push(7) // myArray now contains [5, 7]
4  myArray.push(9) // myArray now contains [5, 7, 9]
5  myArray[0] // evaluates to 5
6  myArray[2] // evaluates to 9
7  myArray[1] = 3 // myArray now contains [5, 3, 9]
8  myBooleanArray = [TRUE, FALSE, TRUE, TRUE]
9  myStringArray = ['foo', 'bar']
10 numberOfBooleans = myBooleanArray.length // evaluates to 4
11 numberOfStrings = myStringArray.length // evaluates to 2
12 for str ∈ myStringArray
13     log(str) // also prints 'foo', then prints 'bar'
14 myArray.map(triple) // evaluates to [15, 21, 27]
```

The above pseudocode demonstrates arrays. Arrays are ordered lists of elements. Arrays can contain elements of any type. Array literals are denoted

by square brackets and can be empty (as in line 1) or populated (as in lines 8 and 9). Arrays have a built-in function attached to them called *push*, which appends a new element to the end of the array. Lines 2-4 demonstrate how *push* can be used to append items to an array. The dot notation seen on lines 2-4 is used on arrays only to access built-in functions and properties.

Square brackets denote access of array elements by index when placed directly after the variable name of the array. Array indices start at zero. Lines 5 and 6 demonstrate how square bracket notation can be used to access values in an array based on their index. Line 7 demonstrates that square bracket notation can also be used to assign to values in an array. Lines 10 and 11 demonstrate the built-in property *length*, the number of elements in the array.

TODO explain lines 12 - 15

```

1  myObject = {}
2  myObject.first = 'John' // myObject now contains {first : 'John'}
3  myObject['last'] = 'Doe' // now {first : 'John', last : 'Doe'}
4  myOtherObject = {first : 'Jane', last : 'Doe'}
5  box =
6      x : 50
7      y : 60
8      width : 100
9      height : 150
10 properties = box.keys // evaluates to [x,y,width,height]
11 values = properties.map( $\lambda$ (property) return box[property])
12 // values is assigned [50, 60, 100, 150]

```

TODO explain above code

```

1  run( $\lambda$ () log('b'))
2  log('a')
3  // Prints a, then b

```

Our pseudocode assumes a single threaded execution environment with a built-in event loop, which may be implemented using the reactor pattern [77]. The event loop can be used to queue functions to be executed in the future. In our pseudocode, the *run* built in function provides access to the event loop. Calling *run* and passing a function queues that function to be invoked in the future, after the current codepath terminates and all previously

queued functions finish executing. In the above pseudocode, line

```
TODO discuss NIL TODO Discuss apply(fn, args) TODO Discuss  
!bool
```

Chapter 2

Functional Reactive Models

Functional reactive programming allows developers to declaratively specify data dependency graphs [95]. Functional reactive programming has been applied to interactive graphics [29] and robotics [44]. The Model View Controller paradigm is an approach to cleanly separate application operations into three classes. The *model* contains the data structures representing the application state. The *view* handles graphical presentation of the model to the user. The *controller* translates user interactions (such as mouse clicks, key presses, or multi-touch gestures) into changes in the model [53]. The Model View Controller paradigm can be combined with functional reactive programming to enable straightforward creation of reactive systems based on data flow graphs.

2.1 The Model Data Structure

The Model in the Model View Controller (MVC) paradigm is responsible for:

- managing the state of the application,
- allowing the Controller to change the state of the application, and
- notifying the view when the state of the application changes.

One simple and widely used method for structuring a Model is as a set of key-value pairs [55]. This kind of model can fulfill the all of the responsibilities of a Model with three methods:

- *set(key, value)* Set the value for a given key.
- *get(key)* Get the value for a given key.
- *on(key, callback)* Add a change listener for a given key. Here, *callback* is a function that will be invoked synchronously when the value for the given key is changed.

This approach is used by popular JavaScript frameworks such as Backbone and Knockout. Using our pseudocode conventions, we will discuss one way a key-value Model can be implemented. We will first discuss a simple version with only *set* and *get*, then discuss a more complex version that also includes *on*.

```

1  SimplestModel =  $\lambda()$ 
2      values = { }
3      return
4          set :  $\lambda(key, value)$  values[key] = value
5          get :  $\lambda(key)$  return values[key]

```

The above pseudocode implements a key-value model that has only *set* and *get* methods. Line 1 defines the constructor function, *SimplestModel*, which will return a new object that has *set* and *get* methods. Line 2 defines a private variable called *values* that will contain the key-value mapping. Lines 3 - 5 define the *set* and *get* methods, which store and retrieve values from the internal *values* object. Here's an example of how *SimplestModel* might be used.

```

1  mySimplestModel = SimplestModel()
2  mySimplestModel.set('x', 5)
3  mySimplestModel.get('x') // Evaluates to 5

```

Here is a version of the model that implements the *on* method as well:

```

1  SimpleModel =  $\lambda()$ 
2      values = { }
3      callbacks = { }
4      return
5          on :  $\lambda(\textit{key}, \textit{callback})$ 
6              if callbacks[key] == NIL
7                  callbacks[key] = []
8                  callbacks[key].push(callback)
9          set :  $\lambda(\textit{key}, \textit{value})$ 
10             values[key] = value
11             if callbacks[key]  $\neq$  NIL
12                 for callback  $\in$  callbacks[key]
13                     callback()
14             get :  $\lambda(\textit{key})$  return values[key]

```

The above version includes an additional private variable, *callbacks*, which is an object whose keys are property names and whose values are arrays of callback functions. The *on* method defined starting at line 5 adds the given callback to the list of callbacks for the given key (and creates the list if it does not yet exist). The *set* method has been modified to invoke the callback functions associated with the given key when the value for that key is changed. Here is an example of how the *on* method can be used.

```

1  mySimpleModel = SimpleModel()
2  mySimpleModel.on('x', λ()
3      log(mySimpleModel.get('x'))
4  )
5  mySimpleModel.set('x', 5) // Causes line 3 to log 5
6  mySimpleModel.set('x', 6) // Causes line 3 to log 6

```

2.2 Functional Reactive Change Propagation

For complex applications such as interactive visualizations, managing propagation of changes can quickly become complex. For this reason, modular visualization environments based on data flow have become popular [1]. A data flow graph defines a directed acyclic graph of data dependencies. The data flow model is amenable to construction of visual programming languages [42]. While many systems consider data flow as a means to construct data transformation pipelines, the concept also applies to building reactive systems that manage change propagation throughout an application or subsystem in response to user interactions or other events [29].

To provide a solid foundation for dynamic visualization systems, the Model should be able function in the context of data dependency graphs. Developers should be able to declaratively specify data dependencies, and change propagation should be automatically managed. The *when* operator from functional reactive programming propagates changes from one or more reactive functions (such as is found in the JavaScript libraries Bacon.js and

RXJS).

Our Model implementation can be extended with a *when* operator that enables construction of data dependency graphs. This operator will become a foundation for building dynamic interactive visualizations. Since *when* is superior to *on* in that it handles change propagation intelligently, in this final version *on* is not exposed in the public Model API.

```
1  Model = λ()  
2      simpleModel = SimpleModel()  
3      return  
4          set : simpleModel.set  
5          get : simpleModel.get  
6          when : λ(dependencies, fn)  
7              callFn = debounce(λ()  
8                  args = dependencies.map(simpleModel.get)  
9                  if allAreDefined(args)  
10                     apply(fn, args)  
11                 )  
12              callFn()  
13              for key ∈ dependencies  
14                  simpleModel.on(key, callFn)
```

TODO describe Model()

```

1  debounce =  $\lambda(fn)$ 
2      queued = FALSE
3      return  $\lambda()$ 
4          if queued == FALSE
5              queued = TRUE
6              run( $\lambda()$ )
7              queued = FALSE
8              fn()
9          )

```

TODO describe *debounce*()

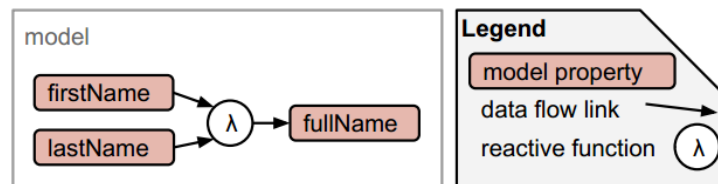
```

1  allAreDefined =  $\lambda(arr)$ 
2      for item  $\in$  arr
3          if item is undefined
4              return FALSE
5      return TRUE

```

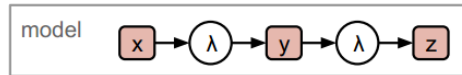
TODO describe *allAreDefined*()

Figure 2.1: A simple data dependency graph using functional reactive models. Here, *fullName* is recomputed whenever *firstName* or *lastName* change.



TODO add pseudocode for each figure

Figure 2.2: A data dependency graph with two hops. When x changes, the change propagates to y then to z .



2.3 Functional Reactive Visualizations

Figure 2.3: The data dependency graph for a dynamic bar chart.

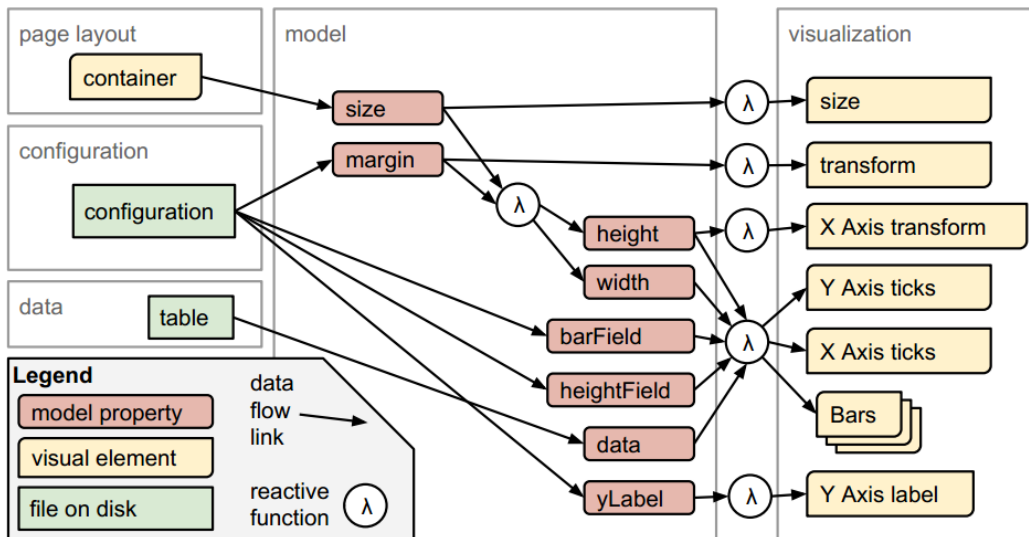


Figure 2.4: The appearance of a dynamic bar chart.

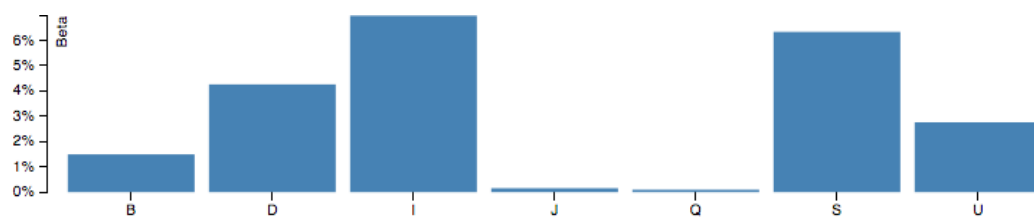
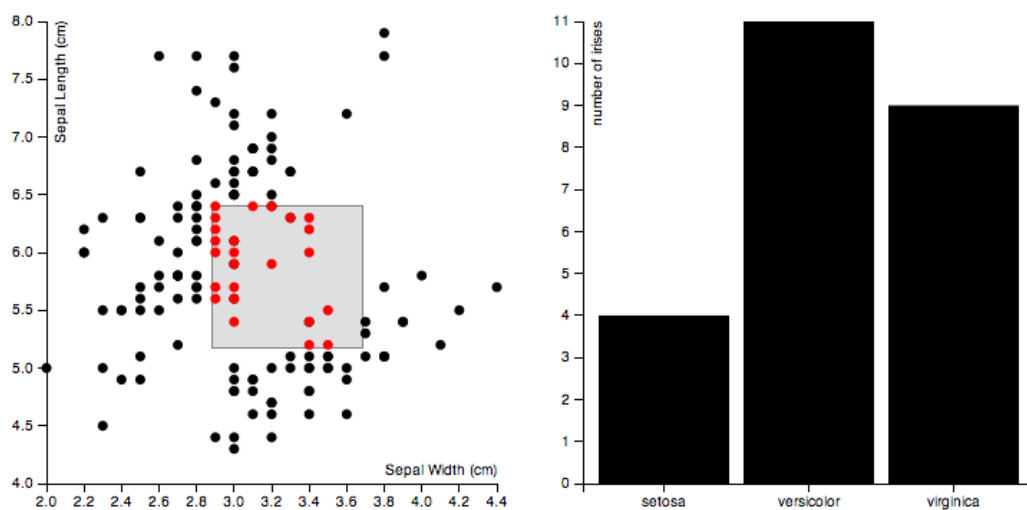


Figure 2.5: An example of multiple linked views using functional reactive models. Brushing in the scatter plot causes the selected data to be aggregated and plotted in the bar chart. Data shown is Fisher's Iris data [32].



Chapter 3

The Universal Data Cube

Chapter 4

Data Sets

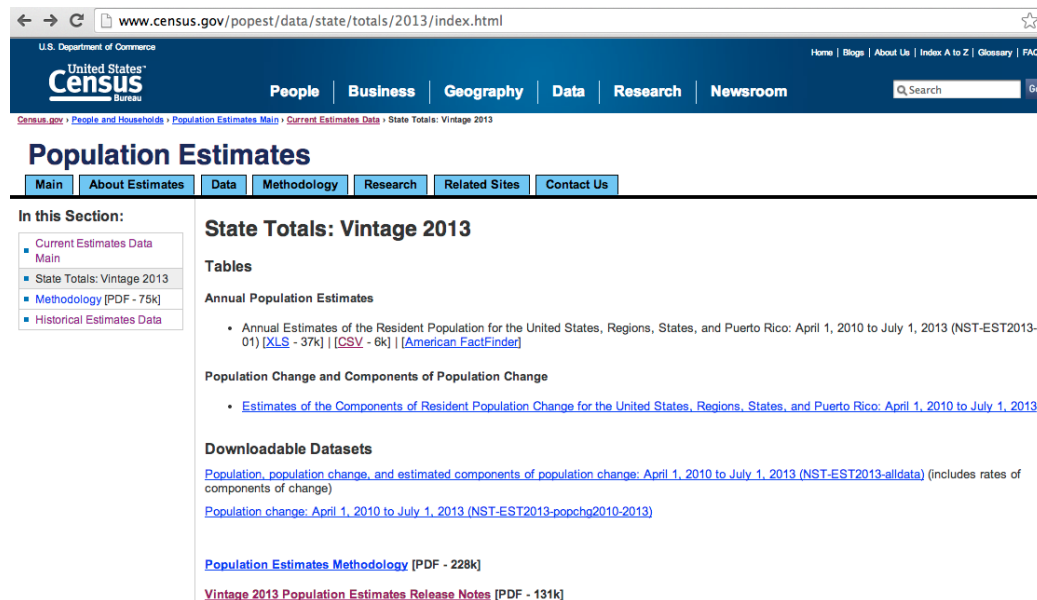
4.1 United Nations Population Estimates

4.2 United Nations Millenium Development Goals

4.3 United States Census Population Estimates

TODO import data from <http://www.census.gov/popest/data/state/totals/2013/index.html>

Figure 4.1: A screenshot of the US Census Population by State data set. This data set is made available in Excel and CSV formats.



4.4 US Central Intelligence Agency World Factbook

TODO import subsets of the data

4.5 US Centers for Disease Control Causes of Death

TODO generalize stacked area and tree vis

4.6 W3Schools Browser Market Share

TODO import this data completely

4.7 Natural Earth

Figure 4.2: A screenshot of the Natural Earth geographic data Web site.



TODO discuss data transformation process

Chapter 5

Visualizations

5.1 Data Cubes and Visualization Theory

Data cubes can contain the following kinds of data:

- *nominal* - unordered collections of categories
- *ordinal* - ordered collections of categories
- *hierarchical* - collections of categories organized as trees
- *quantitative* - continuously varying numeric values

Data cube dimensions can be nominal, ordinal or hierarchical. Data cube measures are always quantitative. This mapping relates data cubes to data types that have been well studied in the literature on visualization theory [6, 60, 34]. **TODO** include pseudocode for each visualization

- 5.2 Bar Chart
- 5.3 Scatter Plot
- 5.4 Line Chart
- 5.5 Choropleth Map
- 5.6 Stacked Area Plot
- 5.7 TreeMap
- 5.8 Node Link Tree
- 5.9 Radial Tree
- 5.10 Icicle Plot
- 5.11 TreeMap
- 5.12 Parallel Coordinates

Chapter 6

Visualization Dashboard Infrastructure

6.1 Dashboard Layout using Nested Boxes

TODO Include nested box layout pseudocode

6.2 Multiple Linked Views

TODO Include generic linking pseudocode using "when"

6.3 Dynamic Dashboard Configuration

TODO Include pseudocode for computing configuration diffs

Chapter 7

UDC Visualization Dashboards

Chapter 8

Collaboration and History Navigation

Bibliography

- [1] Greg Abram and Lloyd Treinish. An extended data-flow architecture for data analysis and visualization. In *Proceedings of the 6th conference on Visualization'95*, page 263. IEEE Computer Society, 1995.
- [2] Rakesh Agrawal, Ashish Gupta, and Sunita Sarawagi. Modeling multidimensional databases. In *Data Engineering, 1997. Proceedings. 13th International Conference on*, pages 232–243. IEEE, 1997.
- [3] Akiko Aizawa and Keizo Oyama. A fast linkage detection scheme for multi-source information integration. In *Web Information Retrieval and Integration, 2005. WIRI'05. Proceedings. International Workshop on Challenges in*, pages 30–39. IEEE, 2005.
- [4] Luc Anselin. Interactive techniques and exploratory spatial data analysis. *Geographical Information Systems: principles, techniques, management and applications*, 1:251–264, 1999.
- [5] Richard A Becker and William S Cleveland. Brushing scatterplots. *Technometrics*, 29(2):127–142, 1987.
- [6] Jacques Bertin. *Semiology of graphics: diagrams, networks, maps*. 1983.
- [7] Markus Blaschka, Carsten Sapia, Gabriele Hofling, and Barbara Dinter. Finding your way through multidimensional data models. In *Database and Expert Systems Applications, 1998. Proceedings. Ninth International Workshop on*, pages 198–203. IEEE, 1998.
- [8] Michael Bostock, Vadim Ogievetsky, and Jeffrey Heer. D3: Data-driven documents. *IEEE Trans. Visualization & Comp. Graphics (Proc. InfoVis)*, 2011.

- [9] Ricardo Cabello. Three. js. URL: <https://github.com/mrdoob/three.js>, 2013.
- [10] Stuart K Card and Jock Mackinlay. The structure of the information visualization design space. In *Information Visualization, 1997. Proceedings., IEEE Symposium on*, pages 92–99. IEEE, 1997.
- [11] Surajit Chaudhuri and Umeshwar Dayal. An overview of data warehousing and olap technology. *ACM Sigmod record*, 26(1):65–74, 1997.
- [12] Ed Huai-hsin Chi. A taxonomy of visualization techniques using the data state reference model. In *Information Visualization, 2000. InfoVis 2000. IEEE Symposium on*, pages 69–75. IEEE, 2000.
- [13] Ed Huai-hsin Chi and John T Riedl. An operator interaction framework for visualization systems. In *Information Visualization, 1998. Proceedings. IEEE Symposium on*, pages 63–70. IEEE, 1998.
- [14] Peter Christen. Febrl: an open source data cleaning, deduplication and record linkage system with a graphical user interface. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1065–1068. ACM, 2008.
- [15] James Clifford and Abdullah Uz Tansel. On an algebra for historical relational databases: two views. In *ACM SIGMOD Record*, volume 14, pages 247–265. ACM, 1985.
- [16] Edgar F Codd, Sharon B Codd, and Clynch T Salley. Providing olap (on-line analytical processing) to user-analysts: An it mandate. *Codd and Date*, 32, 1993.
- [17] John Congote, Alvaro Segura, Luis Kabongo, Aitor Moreno, Jorge Posada, and Oscar Ruiz. Interactive visualization of volumetric data with webgl in real-time. In *Proceedings of the 16th International Conference on 3D Web Technology*, pages 137–146. ACM, 2011.
- [18] Thomas H Cormen, Charles E Leiserson, Ronald L Rivest, Clifford Stein, et al. *Introduction to algorithms*. MIT press, 2009.
- [19] Alfredo Cuzzocrea and Svetlana Mansmann. Olap visualization: models, issues, and techniques. *Encyclopedia of Data Warehousing and Mining*, pages 1439–1446, 2009.

- [20] Alfredo Cuzzocrea, Domenico Saccà, and Paolo Serafino. A hierarchy-driven compression technique for advanced olap visualization of multi-dimensional data cubes. In *Data Warehousing and Knowledge Discovery*, pages 106–119. Springer, 2006.
- [21] Alfredo Cuzzocrea, Domenico Sacca, and Paolo Serafino. Semantics-aware advanced olap visualization of multidimensional data cubes. *International Journal of Data Warehousing and Mining (IJDWM)*, 3(4):1–30, 2007.
- [22] Erik Dahlström, Patrick Dengler, Anthony Grasso, Chris Lilley, Cameron McCormack, Doug Schepers, Jonathan Watt, Jon Ferraiolo, Fujisawa Jun, and Dean Jackson. Scalable vector graphics (svg) 1.1 (second edition). *W3C Recommendation*, 2011.
- [23] Anindya Datta and Helen Thomas. The cube data model: a conceptual model and algebra for on-line analytical processing in data warehouses. *Decision Support Systems*, 27(3):289–301, 1999.
- [24] AnHai Doan, Pedro Domingos, and Alon Y Halevy. Reconciling schemas of disparate data sources: A machine-learning approach. In *ACM Sigmod Record*, volume 30, pages 509–520. ACM, 2001.
- [25] AnHai Doan, Pedro Domingos, and Alon Y Levy. Learning source description for data integration. In *WebDB (Informal Proceedings)*, pages 81–86, 2000.
- [26] AnHai Doan, Alon Halevy, and Zachary Ives. *Principles of data integration*. Elsevier, 2012.
- [27] AnHai Doan and Alon Y Halevy. Semantic integration research in the database community: A brief survey. *AI magazine*, 26(1):83, 2005.
- [28] Mohamed G Elfeky, Vassilios S Verykios, and Ahmed K Elmagarmid. Tailor: A record linkage toolbox. In *Data Engineering, 2002. Proceedings. 18th International Conference on*, pages 17–28. IEEE, 2002.
- [29] Conal Elliott and Paul Hudak. Functional reactive animation. In *ACM SIGPLAN Notices*, volume 32, pages 263–273. ACM, 1997.

- [30] Jérôme Euzenat, Jérôme Euzenat, Pavel Shvaiko, et al. *Ontology matching*. Springer, 2007.
- [31] Ronald Fagin, Phokion G Kolaitis, Renée J Miller, and Lucian Popa. Data exchange: Semantics and query answering. In *Database Theory—ICDT 2003*, pages 207–224. Springer, 2003.
- [32] Ronald A Fisher. The use of multiple measurements in taxonomic problems. *Annals of eugenics*, 7(2):179–188, 1936.
- [33] Steve Fulton and Jeff Fulton. *HTML5 Canvas*. O’Reilly Media, 2013.
- [34] Martin Graham and Jessie Kennedy. A survey of multiple tree visualisation. *Information Visualization*, 9(4):235–252, 2010.
- [35] Jim Gray, Surajit Chaudhuri, Adam Bosworth, Andrew Layman, Don Reichart, Murali Venkatrao, Frank Pellow, and Hamid Pirahesh. Data cube: A relational aggregation operator generalizing group-by, cross-tab, and sub-totals. *Data Mining and Knowledge Discovery*, 1(1):29–53, 1997.
- [36] Lifang Gu, Rohan Baxter, Deanne Vickers, and Chris Rainsford. Record linkage: Current practice and future directions. *CSIRO Mathematical and Information Sciences Technical Report*, 3:83, 2003.
- [37] Marc Gyssens and Laks VS Lakshmanan. A foundation for multi-dimensional databases. In *VLDB*, volume 97, pages 106–115, 1997.
- [38] Alon Halevy, Anand Rajaraman, and Joann Ordille. Data integration: the teenage years. In *Proceedings of the 32nd international conference on Very large data bases*, pages 9–16. VLDB Endowment, 2006.
- [39] Pat Hanrahan, Chris Stolte, and Jock Mackinlay. visual analysis for everyone. *Tableau White paper*, 4, 2007.
- [40] Bin He and Kevin Chen-Chuan Chang. Statistical schema matching across web query interfaces. In *Proceedings of the 2003 ACM SIGMOD international conference on Management of data*, pages 217–228. ACM, 2003.

- [41] Ian Hickson, Robin Berjon, Steve Faulkner, Travis Leithead, Erika Doyle Navara, Edward O'Connor, and Silvia Pfeiffer. *Html5: A vocabulary and associated apis for html and xhtml. W3C Working Draft edition*, 2013.
- [42] Daniel D Hils. Visual languages and computing survey: Data flow visual programming languages. *Journal of Visual Languages & Computing*, 3(1):69–101, 1992.
- [43] CD'Arcy J Holman, A John Bass, Ian L Rouse, and Michael ST Hobbs. Population-based linkage of health records in western australia: development of a health services research linked database. *Australian and New Zealand journal of public health*, 23(5):453–459, 1999.
- [44] Paul Hudak, Antony Courtney, Henrik Nilsson, and John Peterson. Arrows, robots, and functional reactive programming. In *Advanced Functional Programming*, pages 159–187. Springer, 2003.
- [45] Matthew A Jaro. Advances in record-linkage methodology as applied to matching the 1985 census of tampa, florida. *Journal of the American Statistical Association*, 84(406):414–420, 1989.
- [46] Matthew A Jaro. Probabilistic linkage of large public health data files. *Statistics in medicine*, 14(5-7):491–498, 1995.
- [47] Yannis Kalfoglou and Marco Schorlemmer. Ontology mapping: the state of the art. *The knowledge engineering review*, 18(1):1–31, 2003.
- [48] Sean Kandel, Andreas Paepcke, Joseph Hellerstein, and Jeffrey Heer. Wrangler: Interactive visual specification of data transformation scripts. In *PART 5——Proceedings of the 2011 annual conference on Human factors in computing systems*, pages 3363–3372. ACM, 2011.
- [49] Jaewoo Kang and Jeffrey F Naughton. On schema matching with opaque column names and data values. In *SIGMOD Conference*, pages 205–216, 2003.
- [50] Ralph Kimball. *The data warehouse lifecycle toolkit: expert methods for designing, developing, and deploying data warehouses*. Wiley. com, 1998.

- [51] Lukasz Kokoszkiewicz, Julia ANDREEVA, Ivan Antoniev DZHUNOV, Edward KARAVAKIS, Massimo LAMANNA, Jakub MOSCICKI, and Laura SARGSYAN. hbrowse-generic framework for hierarchical data visualization. In *Proceedings of the Computing in High Energy and Nuclear Physics 2012 International Conference*, 2012.
- [52] Nick Koudas, Sunita Sarawagi, and Divesh Srivastava. Record linkage: similarity measures and algorithms. In *Proceedings of the 2006 ACM SIGMOD international conference on Management of data*, pages 802–803. ACM, 2006.
- [53] Glenn E Krasner, Stephen T Pope, et al. A description of the model-view-controller user interface paradigm in the smalltalk-80 system. *Journal of object oriented programming*, 1(3):26–49, 1988.
- [54] Hing-Yan Lee and Hwee-Leng Ong. A new visualisation technique for knowledge discovery in olap. In *KDOOD/TDOOD*, pages 23–25. Citeseer, 1995.
- [55] Avraham Leff and James T Rayfield. Web-application development using the model/view/controller design pattern. In *Enterprise Distributed Object Computing Conference, 2001. EDOC’01. Proceedings. Fifth IEEE International*, pages 118–127. IEEE, 2001.
- [56] Maurizio Lenzerini. Data integration: A theoretical perspective. In *Proceedings of the twenty-first ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 233–246. ACM, 2002.
- [57] Chang Li and X Sean Wang. A data model for supporting on-line analytical processing. In *Proceedings of the fifth international conference on Information and knowledge management*, pages 81–88. ACM, 1996.
- [58] Lauro Lins, James T Klosowski, and Carlos Scheidegger. Nanocubes for real-time exploration of spatiotemporal datasets. *Visualization and Computer Graphics, IEEE Transactions on*, 19(12):2456–2465, 2013.
- [59] Zhicheng Liu, Biye Jiang, and Jeffrey Heer. immens: Real-time visual querying of big data. *Eurovis*, 2013.

- [60] Jock Mackinlay. Automating the design of graphical presentations of relational information. *ACM Transactions on Graphics (TOG)*, 5(2):110–141, 1986.
- [61] Jayant Madhavan, Philip A Bernstein, and Erhard Rahm. Generic schema matching with cupid. In *VLDB*, volume 1, pages 49–58, 2001.
- [62] Jayant Madhavan, S Jeffery, Shirley Cohen, X Dong, David Ko, Cong Yu, and Alon Halevy. Web-scale data integration: You can only afford to pay as you go. In *Proceedings of CIDR*, pages 342–350, 2007.
- [63] Andreas S Maniatis, Panos Vassiliadis, Spiros Skiadopoulos, and Yannis Vassiliou. Advanced visualization for olap. In *Proceedings of the 6th ACM international workshop on Data warehousing and OLAP*, pages 9–16. ACM, 2003.
- [64] Svetlana Mansmann and Marc H Scholl. *Extending visual OLAP for handling irregular dimensional hierarchies*. Springer, 2006.
- [65] Svetlana Mansmann and Marc H Scholl. Exploring olap aggregates with hierarchical visualization techniques. In *Proceedings of the 2007 ACM symposium on Applied computing*, pages 1067–1073. ACM, 2007.
- [66] Svetlana Mansmann and Marc H Scholl. Visual olap: A new paradigm for exploring multidimensional aggregates. In *Proc. of IADIS Int’l Conf. on Computer Graphics and Visualization (CGV)*, pages 59–66, 2008.
- [67] Kouichi Matsuda and Rodger Lea. *WebGL programming guide: interactive 3D graphics programming with WebGL*. Pearson Education, 2013.
- [68] Tova Milo and Sagit Zohar. Using schema matching to simplify heterogeneous data translation. In *VLDB*, volume 98, pages 24–27. Citeseer, 1998.
- [69] Natalya F Noy. Semantic integration: a survey of ontology-based approaches. *ACM Sigmod Record*, 33(4):65–70, 2004.
- [70] Natalya F Noy. Ontology mapping. In *Handbook on ontologies*, pages 573–590. Springer, 2009.

- [71] Natalya F Noy and Mark A Musen. The prompt suite: interactive tools for ontology merging and mapping. *International Journal of Human-Computer Studies*, 59(6):983–1024, 2003.
- [72] Mark Pilgrim. *HTML5: up and running*. ” O’Reilly Media, Inc.”, 2010.
- [73] William Playfair. Commercial and political atlas: Representing, by copper-plate charts, the progress of the commerce, revenues, expenditure, and debts of england, during the whole of the eighteenth century. *London: Corry*, 1786.
- [74] Erhard Rahm and Philip A Bernstein. A survey of approaches to automatic schema matching. *the VLDB Journal*, 10(4):334–350, 2001.
- [75] Raghu Ramakrishnan and Johannes Gehrke. *Database management systems*. Osborne/McGraw-Hill, 2000.
- [76] Hans Rosling, Rönnlund A Rosling, and Ola Rosling. New software brings statistics beyond the eye. *Statistics, Knowledge and Policy: Key Indicators to Inform Decision Making. Paris, France: OECD Publishing*, pages 522–530, 2005.
- [77] Douglas C Schmidt. Reactor: An object behavioral pattern for concurrent event demultiplexing and dispatching. 1995.
- [78] Mathew Scotch, Bambang Parmanto, and Valerie Monaco. Usability evaluation of the spatial olap visualization and analysis tool (sovat). *Journal of Usability Studies*, 2(2):76–95, 2007.
- [79] Matthew Scotch and Bambang Parmanto. Sovat: Spatial olap visualization and analysis tool. In *System Sciences, 2005. HICSS’05. Proceedings of the 38th Annual Hawaii International Conference on*, pages 142b–142b. IEEE, 2005.
- [80] Ben Shneiderman. Dynamic queries for visual information seeking. *Software, IEEE*, 11(6):70–77, 1994.
- [81] Ben Shneiderman. The eyes have it: A task by data type taxonomy for information visualizations. In *Visual Languages, 1996. Proceedings., IEEE Symposium on*, pages 336–343. IEEE, 1996.

- [82] Pavel Shvaiko and Jérôme Euzenat. A survey of schema-based matching approaches. In *Journal on Data Semantics IV*, pages 146–171. Springer, 2005.
- [83] Mark Sifer. A visual interface technique for exploring olap data with coordinated dimension hierarchies. In *Proceedings of the twelfth international conference on Information and knowledge management*, pages 532–535. ACM, 2003.
- [84] Chris Stolte, Diane Tang, and Pat Hanrahan. Polaris: A system for query, analysis, and visualization of multidimensional relational databases. *Visualization and Computer Graphics, IEEE Transactions on*, 8(1):52–65, 2002.
- [85] Chris Stolte, Diane Tang, and Pat Hanrahan. Query, analysis, and visualization of hierarchically structured data using polaris. In *KDD*, volume 2, pages 112–122. Citeseer, 2002.
- [86] Chris Stolte, Diane Tang, and Pat Hanrahan. Multiscale visualization using data cubes. *Visualization and Computer Graphics, IEEE Transactions on*, 9(2):176–187, 2003.
- [87] Kesaraporn Techapichetvanich and Amitava Datta. Interactive visualization for olap. In *Computational Science and Its Applications–ICCSA 2005*, pages 206–214. Springer, 2005.
- [88] David P Tegarden. Business information visualization. *Communications of the AIS*, 1(1es):4, 1999.
- [89] Edward R Tufte and PR Graves-Morris. *The visual display of quantitative information*, volume 2. Graphics press Cheshire, CT, 1983.
- [90] Michael Uschold and Michael Gruninger. Ontologies and semantics for seamless connectivity. *ACM SIGMod Record*, 33(4):58–64, 2004.
- [91] Panos Vassiliadis. Modeling multidimensional databases, cubes and cube operations. In *Scientific and Statistical Database Management, 1998. Proceedings. Tenth International Conference on*, pages 53–62. IEEE, 1998.

- [92] Panos Vassiliadis and Timos Sellis. A survey of logical models for olap databases. *ACM Sigmod Record*, 28(4):64–69, 1999.
- [93] Fernanda B Viegas, Martin Wattenberg, Frank Van Ham, Jesse Kriss, and Matt McKeon. Manyeyes: a site for visualization at internet scale. *Visualization and Computer Graphics, IEEE Transactions on*, 13(6):1121–1128, 2007.
- [94] Holger Wache, Thomas Voegelé, Ubbo Visser, Heiner Stuckenschmidt, Gerhard Schuster, Holger Neumann, and Sebastian Hübner. Ontology-based integration of information-a survey of existing approaches. In *IJCAI-01 workshop: ontologies and information sharing*, volume 2001, pages 108–117. Citeseer, 2001.
- [95] Zhanyong Wan and Paul Hudak. Functional reactive programming from first principles. In *ACM SIGPLAN Notices*, volume 35, pages 242–252. ACM, 2000.
- [96] Matthew O Ward. Xmdvtool: Integrating multiple methods for visualizing multivariate data. In *Proceedings of the Conference on Visualization’94*, pages 326–333. IEEE Computer Society Press, 1994.
- [97] Leland Wilkinson. *The grammar of graphics*. Springer, 2005.
- [98] William E Winkler. Matching and record linkage. *Business survey methods*, 1:355–384, 1995.
- [99] William E Winkler. The state of record linkage and current research problems. In *Statistical Research Division, US Census Bureau*. Citeseer, 1999.
- [100] William E Winkler. Overview of record linkage and current research directions. In *Bureau of the Census*. Citeseer, 2006.
- [101] Qiang Yang, Joshua Zhexue Huang, and Michael Ng. A data cube model for prediction-based web prefetching. *Journal of Intelligent Information Systems*, 20(1):11–30, 2003.
- [102] Ji Soo Yi, Youn ah Kang, John T Stasko, and Julie A Jacko. Toward a deeper understanding of the role of interaction in information visualization. *Visualization and Computer Graphics, IEEE Transactions on*, 13(6):1224–1231, 2007.

- [103] Patrick Ziegler and Klaus R Dittrich. Three decades of data integration-all problems solved? In *IFIP congress topical sessions*, pages 3–12. Springer, 2004.