

Effortless Data Exploration with zenvisage: An Expressive and Interactive Visual Analytics System

Tarique Siddiqui¹ Albert Kim² John Lee¹ Karrie Karahalios¹ Aditya Parameswaran¹

¹University of Illinois (UIUC) ²MIT
{tsiddiq2,lee98,kkarahal,adityagp}@illinois.edu alkim@csail.mit.edu

ABSTRACT

Data visualization is by far the most commonly used mechanism to explore and extract insights from datasets, especially by novice data scientists. And yet, current visual analytics tools are rather limited in their ability to operate on collections of visualizations—by composing, filtering, comparing, and sorting them—to find those that depict desired trends or patterns. The process of visual data exploration remains a tedious process of *trial-and-error*. We propose zenvisage, a visual analytics platform for effortlessly finding desired visual patterns from large datasets. We introduce zenvisage’s general purpose visual exploration language, ZQL (“zee-quel”) for specifying the desired visual patterns, drawing from use-cases in a variety of domains, including biology, mechanical engineering, climate science, and commerce. We formalize the expressiveness of ZQL via a visual exploration algebra—an algebra on collections of visualizations—and demonstrate that ZQL is as expressive as that algebra. zenvisage exposes an interactive front-end that supports the issuing of ZQL queries, and also supports interactions that are “short-cuts” to certain commonly used ZQL queries. To execute these queries, zenvisage uses a novel ZQL graph-based query optimizer that leverages a suite of optimizations tailored to the goal of processing collections of visualizations in certain pre-defined ways. Lastly, a user survey and study demonstrates that data scientists are able to effectively use zenvisage to eliminate error-prone and tedious exploration and directly identify desired visualizations.

1. INTRODUCTION

Interactive visualization tools, such as Tableau [3] and Spotfire [2], have paved the way for the democratization of data exploration and data science. These tools have witnessed an ever-expanding user base—as a concrete example, Tableau’s revenues last year were in the hundreds of millions of US Dollars and is expected to reach tens of billions soon [6]. Using such tools, or even tools like Microsoft Excel, the standard data analysis recipe is as follows: the data scientists load a dataset into the tool, select visualizations to examine, study the results, and then repeat the process until they find ones that match their desired pattern or need. Thus, using this repeated process of manual examination, or *trial-and-error*, data scientists are able to formulate and test hypothesis, and derive insights. The key premise of this work is that to find desired patterns in datasets, *manual examination of each visualization in a collection is simply unsustainable*, especially on large, complex datasets. Even on moderately sized datasets, a data scientist may need to examine as many as tens of thousands of visualizations, all to test a single hypothesis, a severe impediment to data exploration.

To illustrate, we describe the challenges of several collaborator groups who have been hobbled by the ineffectiveness of current data exploration tools:

Case Study 1: Engineering Data Analysis. Battery scientists at Carnegie Mellon University perform visual exploration of datasets of solvent properties to design better batteries. A specific task may involve finding solvents with desired behavior: e.g., those whose solvation energy of Li^+ vs. the boiling point is a *roughly increasing trend*. To do this using current tools, these scientists manually examine the plot of Li^+ solvation energy vs. boiling point for each of the thousands of solvents, to find those that match the desired pattern of a roughly increasing trend.

Case Study 2: Advertising Data Analysis. Advertisers at ad analytics firm Turn, Inc., often examine their portfolio of advertisements to see if their campaigns are performing as expected. For instance, an advertiser may be interested in seeing if there are any keywords that are *behaving unusually* with respect to other keywords in Asia—for example, maybe most keywords have a specific trend for click-through rates (CTR) over time, while a small number of them have a different trend. To do this using the current tools available at Turn, the advertiser needs to manually examine the plots of CTR over time for each keyword (thousands of such plots), and remember what are the typical trends.

Case Study 3: Genomic Data Analysis. Clinical researchers at the NIH-funded genomics center at UIUC and Mayo Clinic are interested in studying data from clinical trials. One such task involves finding pairs of genes that *visually explain the differences* in clinical trial outcomes (positive vs. negative)—visualized via a scatterplot with the x and y axes each referring to a gene, and each outcome depicted as a point in the scatterplot—with the positive outcomes depicted in one color, and the negative ones as another. Current tools require the researchers to generate and manually evaluate tens of thousands of scatter plots of pairs of genes for whether the outcomes can be clearly distinguished in the scatter plot.

Case Study 4: Environmental Data Analysis. Climate scientists at the National Center for Supercomputing Applications at Illinois are interested in studying the nutrient and water property readings on sensors within buoys at various locations in the Great Lakes. Often, they find that a sensor is displaying unusual behavior for a specific property, and want to figure out *what is different* about this sensor relative to others, and if other properties for this sensor are *showing similar behavior*. In either case, the scientists would need to separately examine each property for each sensor (in total 100s of thousands of visualizations) to identify explanations or similarities.

Case Study 5: Server Monitoring Analysis. The server monitoring team at Facebook has noticed a spike in the per-query response time for Image Search in Russia on August 15, after which the response time flattened out. The team would like to identify if there are *other attributes that have a similar behavior* with per-query response time, which may indicate the reason for the spike and subsequent flattening. To do this, the server monitoring team generates

visualizations for different metrics as a function of the date, and assess if any of them has similar behavior to the response time for Image Search. Given that the number of metrics is likely in the thousands, this takes a very long time.

Case Study 6: Mobile App Analysis. The complaints section of the Google mobile platform team have noticed that a certain mobile app has received many complaints. They would like to figure out *what is different* about this app relative to others. To do this, they need to plot various metrics for this app to figure out why it is behaving anomalously. For instance, they may look at network traffic generated by this app over time, or at the distribution of energy consumption across different users. In all of these cases, the team would need to generate several visualizations manually and browse through all of them in the hope of finding what could be the issues with the app.

Thus, in these examples, the recurring theme is the manual examination of a large number of generated visualizations for a specific visual pattern. Indeed, we have found that in these scenarios—data exploration can be a tedious and time-consuming process with current visualization tools.

Key Insight. The goal of this paper is to develop *zenvisage*, a visual analytics system that can *automate the search for desired visual patterns*. Our key insight in developing *zenvisage* is that the data exploration needs in all of these scenarios can be captured within a common set of operations on collections of visualizations. These operations include: *composing* collections of visualizations, *filtering* visualizations, based on some conditions, *comparing* visualizations, and *sorting* them based on some condition. The conditions include similarity or dissimilarity to a specific pattern, “typical” or anomalous behavior, or the ability to provide explanatory or discriminatory power. These operations and conditions form the kernel of a new data exploration language, ZQL (“zee-quel”), that forms the foundation upon which *zenvisage* is built.

Key Challenges. We encountered many challenges in building the *zenvisage* visual analytics platform, a substantial advance over the manually-intensive visualization tools like Tableau and Spotfire; these tools enable the examination of *one visualization at a time*, without the ability to automatically identify relevant visualizations from a collection of visualizations.

First, there were many challenges in developing ZQL, the underlying query language for *zenvisage*. Unlike relational query languages that operate directly on data, ZQL operates on collections of visualizations, which are themselves aggregate queries on data. Thus, in a sense ZQL is a query language that operates on other queries a first class citizen. This leads to a number of challenges that are not addressed in a relational query language context. For example, we had to develop a natural way to users to specify a collection of visualizations to operate on, without having to explicitly list them; even though the criteria on which the visualizations were compared varied widely, we had to develop a small number of general mechanisms that capture all of these criteria; often, the visualizations that we operated on had to be modified in various ways—e.g., we might be interested in visualizing the sales of a product whose profits have been dropping—composing these visualizations from existing ones is not straightforward; and lastly, drilling down into specific visualizations from a collection also required special care. Our ZQL language is a synthesis of desiderata after discussions with data scientists from a variety of domains, and has been under development for the past two years. To further show that ZQL is *complete* under a new visual exploration algebra that we develop, involved additional challenges.

Second, in terms of front-end development, *zenvisage*, being an

interactive analytics tool, needs to support the ability for users to interactively specify ZQL queries—specifically, interactive shortcuts for commonly used ZQL queries, as well as the ability to pose extended ZQL queries for more complex needs. Identifying common interaction “idioms” for these needs took many months.

Third, an important challenge in building *zenvisage* is the backend that supports the execution of ZQL. A single ZQL query can lead to the generation of 10000s of visualizations—executing each one independently as an aggregate query, would take several hours, rendering the tool somewhat useless. (As it turns out, this time would be what an analyst aiming to discover the same pattern would have to spend with present visualization tools, so the naive automation may still help reducing the amount of manual effort.) *zenvisage*’s query optimizer operates as a wrapper over any traditional relational database system. This query optimizer compiles ZQL queries down to a directed acyclic graph of operations on collections of visualizations, followed with the optimizer using a combination of intelligent speculation and combination, to issue queries to the underlying database. We also demonstrate that the underlying problem is NP-HARD. Our query optimizer leads to substantial improvements over the naive schemes adopted within relational database systems for multi-query optimization.

Related Work. There are a number of tools one could use for interactive analysis; here, we briefly describe why those tools are inadequate for the important need of automating the search for desired visual insights. We describe related work in detail in Section 8.

To start, visualization tools like Tableau and Spotfire only generate and provide one visualization at a time, while *zenvisage* analyzes collections of visualizations at a time, and identifies relevant ones from that collection—making it substantially more powerful.

While we do use relational database systems as a computation layer, it is cumbersome to near-impossible to express these user needs in SQL. As an example, finding visualizations of solvents for whom a given property follows a roughly increasing trend is impossible to write within native SQL, and would require custom UDFs—these UDFs would need to be hand-written for every ZQL query. Similarly, finding visualizations of keywords where CTR over time in Asia is behaving unusually with respect to other keywords is challenging to write within SQL. For the small space of queries where it is possible to write the queries within SQL these queries require non-standard constructs, and are both complex and cumbersome to write, even for expert SQL users, and are optimized very poorly (see Section 8). It is also much more natural for end-users to operate directly on visualizations than on data. Indeed, users who have never programmed or written SQL before find it easy to understand and write a subset of ZQL queries, as we will show subsequently.

Statistical, data mining, and machine learning certainly provide functionality beyond *zenvisage* in supporting prediction and statistics; these functionalities are exposed as “one-click” algorithms that can be applied on data. However, no functionality is provided for searching for desired patterns; no querying functionality beyond the one-click algorithms, and no optimization. To use such tools for ZQL, many lines of code and hand-optimization is needed. As such, these tools are beyond the reach of novice data scientists who simply want to explore and visualize their datasets.

Outline. We first describe our query language for *zenvisage*, ZQL (Section 2), and formalize the notion of a *visual exploration algebra*, an analog of relational algebra, describing a core set of capabilities for any language that supports visual data exploration and demonstrate that ZQL is *complete* in that it subsumes these capabilities (Section 3). We describe the graph-based query translator and optimizer for ZQL (Section 4). We then describe our initial

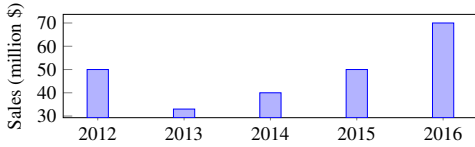


Figure 1: Sales over year visualization for the product chair.

prototype of *zenvisage* (Section 5). We describe our performance experiments (Section 6), and present a user survey and study focused on evaluating the effectiveness and usability of *zenvisage* (Section 7). In the appendix, we present additional details of our query language, along with complete examples, and additional details on user study.

2. QUERY LANGUAGE

zenvisage’s query language, ZQL, provides users with a powerful mechanism to operate on collections of visualizations. In fact, ZQL treats visualizations as a *first-class citizen*, enabling users to operate at a high level on collections of visualizations much like one would operate on relational data with SQL. For example, a user may want to filter out all visualizations where the visualization shows a roughly decreasing trend from a collection, or a user may want to create a collection of visualizations which are most similar to a visualization of interest. Regardless of the query, ZQL provides an intuitive, yet flexible specification mechanism for users to express the desired patterns of interest (in other words, their *exploration needs*) using a small number of ZQL lines. Overall, ZQL provides users the ability to compose collections of visualizations, filter them, and sort and compare them in various ways.

ZQL draws heavy inspiration from the Query by Example (QBE) language [49] and uses a similar table-based specification interface. Although ZQL components are not fundamentally tied to the tabular interface, we found that our end-users feel more at home with it; many of them are non-programmers who are used to spreadsheet tools like Microsoft Excel. Users may either directly write ZQL, or they may use the *zenvisage* front-end, which supports interactions that are transformed internally into ZQL.

We now provide a formal introduction to ZQL in the rest of this section. We introduce many sample queries to make it easy to follow along, and we use a relatable fictitious product sales-based dataset throughout this paper in our query examples—we will reveal attributes of this dataset as we go along.

2.1 Formalization

For describing ZQL, we assume that we are operating on a single relation or a star schema where the attributes are unique (barring key-foreign key joins), allowing ZQL to seamlessly support natural joins. In general, ZQL could be applied to arbitrary collections of relations by letting the user precede an attribute A with the relation name R , e.g., RA . For ease of exposition, we focus on the single relation case.

2.1.1 Overview

The concept of visualizations. We start by defining the notion of a visualization. We use a sample visualization in Figure 1 to guide our discussion. Of course, different visual analysis tasks may require different types of visualizations (instead of bar charts, we may want scatter plots or trend lines), but across all types a visualization is defined by the following five main components: (i) the x-axis attribute, (ii) the y-axis attribute, (iii) the subset of data used, (iv) the type of visualization (e.g., bar chart, scatter plot), and (v) the binning and aggregation functions for the x- and y- axes.

Name	X	Y	Z	Viz
*fl	'year'	'sales'	'product'.'chair'	bar.(y=agg('sum'))

Table 1: Query for the bar chart of sales over year for the product chair.

Name	X	Y	Z	Viz
fl	'year'	'sales'	'product'.	bar.(y=agg('sum'))

Table 2: Query for the bar chart of sales over year for each product.

Visualization collections in ZQL: ZQL has four columns to support the specification of visualizations that the five aforementioned components map into: (i) X , (ii) Y , (iii) Z , and (iv) Viz .

Table 1 gives an example of a valid ZQL query that uses these columns to specify a bar chart visualization of overall sales over the years for the product chair (i.e., the visualization in Figure 1)—ignore the Name column for now. The details for each of these columns are presented subsequently. In short, the x axis (X) is the attribute *year*, the y axis (Y) is the attribute *sales*, and the subset of data (Z) is the product chair, while the type of visualization is a bar chart (*bar*), and the binning and aggregation functions indicate that the y axis is an aggregate (*agg*)—the sum of sales.

In addition to specifying a single visualization, users may often want to retrieve multiple visualizations. ZQL supports this in two ways. Users may use multiple rows, and specify one visualization per row. The user may also specify a *collection* of visualizations in a single row by iterating over a collection of values for one of the X , Y , Z , and Viz columns. Table 2 gives an example of how one may iterate over all products (using the notation $*$ to indicate that the attribute product can take on all values), returning a separate sales bar chart for each product.

High-level structure of ZQL. Starting from these two examples, we can now move onto the general structure of ZQL queries. Overall, each ZQL query consists of multiple rows, where each row operates on collections of visualizations. Each row contains three sets of columns, as depicted in Table 3: (i) the first column corresponds to an identifier for a visualization collection, (ii) the second set of columns defines a visualization collection, while (iii) the last column corresponds to some operation on the visualization collection. All columns can be left empty if needed (in such cases, to save space, for convenience, we do not display these columns in our paper). For example, the last column may be empty if no operation is to be performed, like it was in Table 1 and 2. We have already discussed (ii); now we will briefly discuss (i) and (iii), corresponding to *Name* and *Process* respectively.

Identifiers and operations in ZQL. The *Process* column allows the user to operate on the defined collections of visualizations, applying high-level filtering, sorting, and comparison. The *Name* column provides a way to label and combine specified collections of visualizations, so users may refer to them in the *Process* column. Thus, by repeatedly using the X , Y , Z , and Viz columns to compose visualizations and the *Process* column to process those visualizations, the user is able derive the exact set of visualizations she is looking for. Note that the result of a ZQL query is the *data* used to generate visualizations. The *zenvisage* front-end then uses this data to render the visualizations for the user to peruse.

2.1.2 X , Y , and Z

The X and Y columns specify the attributes used for the x- and y- axes. For example, Table 1 dictates that the returned visual-

Name	X	Y	Z	Viz	Process
Identifier	Visualization Collection			Operation	

Table 3: ZQL query structure.

Name	X	Y	Z	Viz
...	...	{'sales', 'profit'}

Table 4: Query for the sales and profit bar charts for the product chair (missing values are the same as that in Table 1)

Name	X	Y	Z	Viz
...	{'year', 'month'}	{'sales', 'profit'}

Table 5: Query for the sales and profit bar charts over years and months for chairs (missing values are the same as in Table 1).

Name	X	Y	Z	Z2	Viz
...	'location'. 'US'	...

Table 6: Query which returns the overall sales bar chart for the chairs in US (all missing values are the same as that in Table 1).

ization should have ‘year’ for its x-axis and ‘sales’ for its y-axis. As mentioned, the user may also specify a collection of values for the X and Y columns if they wish to refer to a collection of visualizations in one ZQL row. Table 4 refers the collection of both sales-over-years and profit-over-years bar charts for the chair—the missing values in this query (“...”) are the same as Table 1. As we can see, a collection is constructed using {}. If the user wishes to denote all possible values, the shorthand * symbol may be used, as is shown by Table 2. In the case that multiple columns contain collections, a Cartesian product is performed, and visualizations for every combination of values is returned. For example, Table 5 would return the collection of visualizations with specifications: {(X: ‘year’, Y: ‘sales’), (X: ‘year’, Y: ‘profit’), (X: ‘month’, Y: ‘sales’), (X: ‘month’, Y: ‘profit’)}. Additionally, ZQL allows composing multiple attributes in the X and Y columns by supporting Polaris table algebra [4] over the operators: +, x, / (Appendix A).

With the Z column, the user can select which subset of the data they wish to construct their visualizations from. ZQL uses the $\langle attribute \rangle . \langle attribute-value \rangle$ notation to denote the selection of data. Consequently, the query in Table 1 declares that the user wishes to retrieve the sales bar chart only for the chair product. Note that unlike the X and Y columns, both the attribute and the attribute value must be specified for the Z column; otherwise, a proper subset of the data would not be identified. Collections are allowed for both the attribute and the attribute value in the Z column. Table 2 shows an example of using the * shorthand to specify a collection of bar charts, one for each product. A Z column which has a collection over attributes might look like: {'location', 'product'}.* (i.e., a visualization for every product and a visualization for every location). In addition, the Z column allows users to specify predicate constraints using syntax like ‘weight’.[? < 10]; this specifies all items whose weight is less than 10 lbs. To evaluate, the ? is replaced with the attribute and the resulting expression is passed to SQL’s WHERE clause. The predicate constraint syntax has an analogous predicate collection syntax, which creates a collection of the values which satisfy the condition. ‘weight’.[? < 10] specifies that the resulting visualizations must only contain items with less than 10 lbs, ‘weight’.[? < 10] creates a collection of values, one for each item which is less than 10 lbs.

ZQL supports multiple constraints on different attributes through the use of multiple Z columns. In addition to the basic Z column, the user may choose to add Z2, Z3, ... columns depending on how many constraints she requires. Table 6 gives an example of a query which looks at sales plots for chairs only in the US. Note that Z columns are combined using conjunctive semantics.

2.1.3 Viz

The Viz column decides the visualization type, binning, and aggregation functions for the row. Elements in this column have the

Name	X	Y	Viz
*f1	'weight'	'sales'	bin2d.(x=nbin(20), y=nbin(20))

Table 7: Query which returns the heat map of sales vs. weights across all transactions.

Name	X	Y	Z
f1	'year'	'sales'	'product'. 'chair'
f2	'year'	'profit'	'location'. 'US'
*f3 <- f1 + f2			'weight'. [? < 10]

Table 8: Query which returns the sales for chairs or profits for US visualizations for all items less than 10 lbs.

format: $\langle type \rangle . \langle bin+aggr \rangle$. All examples so far have been bar charts with no binning and SUM aggregation for the y-axis, but other variants are supported. The visualization types are derived from the Grammar of Graphics [48] specification language, so all plots from the geometric transformation layer of ggplot [47] (the tool that implements Grammar of Graphics) are supported. For instance, scatter plots are requested with point and heat maps with bin2d. As for binning, binning based on bin width (bin) and number of bins (nbin) are supported for numerical attributes—we may want to use binning, for example, when we are plotting the total number of products whose prices lie within 0-10, 10-20, and so on.

Finally, ZQL supports all the basic SQL aggregation functions such as AVG, COUNT, and MAX. Table 7 is an example of a query which uses a different visualization type, heat map, and creates 20 bins for both x- and y- axes.

Like the earlier columns, the Viz column also allows collections of values. Similar to the Z column, collections may be specified for both the visualization type or the binning and aggregation. If the user wants to view the same data binned at different granularities, she might specify a bar chart with several different bin widths: bar.(x={bin(1), bin(5), bin(10)}, y=agg('sum')). On the other hand, if the user wishes to view the same data in different visualizations, she might write: {bar.(y=agg('sum')), point.()}. The Viz column allows users powerful control over the structure of the rendered visualization. However, there has been work from the visualization community which automatically tries to determine the most appropriate visualization type, binning, and aggregation for a dataset based on the x- and y- axis attributes [25, 33]. Thus, we can frequently leave the Viz column blank and zervisage will use these rules of thumb to automatically decide the appropriate setting for us. With this in mind, we omit the Viz column from the remaining examples with the assumption that zervisage will determine the “best” visualization structure for us.

2.1.4 Name

Together, the values in the X, Y, Z, and Viz columns of each row specify a collection of visualizations. The Name column allows us to label these collections so that they can be referred to be in the Process column. For example, f1 is the label or identifier given to the collection of sales bar charts in Table 2. The * in front of f1 signifies that the the collection is an output collection; that is, ZQL should return this collection of visualizations to the user.

However, not all rows need to have a * associated with their Name identifier. A user may define intermediate collections of visualizations if she wishes to further process them in the Process column before returning the final results. In the case of Table 8, f1 and f2 are examples of intermediate collections.

Also in Table 8, we have an example of how the Name column allows us to perform high-level set-like operations to combine visualization collections directly. For example, f3 <- f1 + f2 assigns f3 to the collection which includes all visualizations in f1 and f2 (similar to set union). This can be useful if the user wishes to combine variations of values without considering the full Cartesian

product. Our example in Table 8, the user is able to combine the sales for chairs plots with the profits for the US plots without also having to consider the sales for the US plots or the profits for chairs plots; she would have had to do so if she had used the specification: (Y: {'sales', 'profit'}, Z: {'product'.chair, 'location'.US}).

An interesting aspect of Table 8 is that the X and Y columns of the third row are devoid of values, and the Z column refer to the seemingly unrelated weight attribute. The values in the X, Y, Z, and Viz columns all help to specify a particular collection of visualizations from a larger collection. When this collection is defined via the Name column, we no longer need to fill in the values for X, Y, Z, or Viz, except to select from the collection—here, ZQL only selects the items which satisfy the constraint, $\text{weight} < 10$.

Other set-like operators include $f1 - f2$ for set minus and $f1 \wedge f2$ for intersection.

2.1.5 Process

The real power of ZQL as a query language comes not from its ability to effortlessly specify collections of visualizations, but rather from its ability to operate on these collections somewhat declaratively. With ZQL’s processing capabilities, users can filter visualizations based on trend, search for similar-looking visualizations, identify representative visualizations, and determine outlier visualizations. Naturally, to operate on collections, ZQL must have a way to iterate over them; however, since different visual analysis tasks might require different forms of traversals over the collections, we expose the iteration interface to the user.

Iterations over collections. Since collections may be composed of varying values from multiple columns, iterating over the collections is not straight-forward. Consider Table 9—the goal is to return profit by year visualizations for the top-10 products whose profit by year visualizations look the most different from the sales by year visualizations. This may indicate a product that deserves special attention. While we will describe this query in detail below, at a high level the first row assembles the visualizations for profit over year for all products ($f1$), the second row assembles the visualizations for sales over year for all products ($f2$), followed by operating (via the Process column) on these two collections by finding the top-10 products who sales over year is most different from profit over year, while the third row displays the profit over year for those top-10 products. A array-based representation of the visualization collections $f1$ and $f2$, would look like the following:

$$f1 = \left\{ \begin{array}{l} \text{X: 'year', Y: 'profit'} \\ \text{Z: 'product.chair'} \\ \text{Z: 'product.table'} \\ \text{Z: 'product.stapler'} \\ \vdots \end{array} \right\} \quad f2 = \left\{ \begin{array}{l} \text{X: 'year', Y: 'sales'} \\ \text{Z: 'product.chair'} \\ \text{Z: 'product.table'} \\ \text{Z: 'product.stapler'} \\ \vdots \end{array} \right\}$$

We would like to iterate over the products—the Z dimension values—of both $f1$ and $f2$ to make our comparisons. Furthermore, we must iterate over the products in the *same order* for both $f1$ and $f2$ to ensure that a product’s profit visualization correctly matches with its sales visualization. Using a single index for this would be complicated and need to take into account the sizes of each of the columns. While there may be other ways to architect this iteration for a single attribute, it is virtually impossible to do when there are multiple attributes that are varying. Instead, ZQL opts for a more powerful *dimension-based* iteration, which assigns each column (or dimension) a separate iterator called an *axis variable*. This dimension-based iteration is a powerful idea that extends to any number of dimensions. As shown in Table 9, axis variables are defined and assigned using the syntax: $\langle \text{variable} \rangle \leftarrow \langle \text{collection} \rangle$; axis variable $v1$ is assigned to the Z dimension of $f1$ and iterates over all

product values. For cases in which multiple collections must traverse over a dimension in the same order, an axis variable must be shared across those collections for that dimension; in Table 9, $f1$ and $f2$ share $v1$ for their Z dimension, since we want to iterate over the products in lockstep.

Operations on collections. With the axis variables defined, the user can then formulate the high-level operations on collections of visualizations as an optimization function which maximizes/minimizes for their desired pattern. Given that $\text{argmax}_x[k = 10] g(x)$ returns the top-10 x values which maximizes the function $g(x)$, and $D(x, y)$ returns the “distance” between x and y , now consider the expression in the Process column for Table 9. Colloquially, the expression says to find the top-10 $v1$ values whose $D(f1, f2)$ values are the largest. The $f1$ and $f2$ in $D(f1, f2)$ refer to the collections of visualizations in the first and second row and are bound to the current value of the iteration for $v1$. In other words, for each product $v1$ in $v1$, retrieve the visualizations $f1[z: v1]$ from collection $f1$ and $f2[z: v1]$ from collection $f2$ and calculate the “distance” between these visualizations; then, retrieve the 10 $v1$ values for which this distance is the largest—these are the products, and assign $v2$ to this collection. Subsequently, we can access this set of products in Z column of the third line of Table 9.

Formal structure. More generally, the basic structure of the Process column is:

$$\begin{aligned} &\langle \text{argopt} \rangle_{\langle \text{axvar} \rangle} [\langle \text{limiter} \rangle] \langle \text{expr} \rangle \quad \text{where} \\ &\langle \text{expr} \rangle \rightarrow (\max | \min | \sum | \prod)_{\langle \text{axvar} \rangle} \langle \text{expr} \rangle \\ &\quad \rightarrow \langle \text{expr} \rangle (+ | - | \times | \div) \langle \text{expr} \rangle \\ &\quad \rightarrow T(\langle \text{nmvar} \rangle) \\ &\quad \rightarrow D(\langle \text{nmvar} \rangle, \langle \text{nmvar} \rangle) \\ &\langle \text{argopt} \rangle \rightarrow (\text{argmax} | \text{argmin} | \text{argany}) \\ &\langle \text{limiter} \rangle \rightarrow (k = \mathbb{N} | t > \mathbb{R} | p = \mathbb{R}) \end{aligned}$$

where $\langle \text{axvar} \rangle$ refers to the axis variables, and $\langle \text{nmvar} \rangle$ refers to collections of visualizations. $\langle \text{argopt} \rangle$ may be one of argmax , argmin , or argany , which returns the values which have the largest, smallest, and any expressions. The $\langle \text{limiter} \rangle$ limits the number of results: $k = \mathbb{N}$ returns only the top- k values; $t > \mathbb{R}$ returns only values who are larger than a threshold value t (may also be smaller, greater than equal, etc.); $p = \mathbb{R}$ returns the top p -percentile values. T and D are two simple *functional primitives* supported by ZQL that can be applied to visualizations to find desired patterns:

- $[T(f) \rightarrow \mathbb{R}]$: T is a function which takes a visualization f and returns a real number measuring some visual property of the trend of f . One such property is “growth”, which returns a positive number if the overall trend is “upwards” and a negative number otherwise; an example implementation might be to measure the slope of a linear fit to the given input visualization f . Other properties could measure the skewness, or the number of peaks, or noisiness of visualizations.
- $[D(f, f') \rightarrow \mathbb{R}]$: D is a function which takes two visualizations f and f' and measures the distance (or dissimilarity) between these visualizations. Examples of distance functions may include pointwise distance functions like Euclidean distance, Earth Mover’s Distance, or the Kullback-Leibler Divergence. The distance D could also be measured using the difference in the number of peaks, or slopes, or some other property.

ZQL supports many different implementations for these two functional primitives, and the user is free to choose any one. If the user does not select one, zenvisage will automatically detect the “best” primitive based on the data characteristics. Furthermore, if ZQL does not have an implementation of the T or D function that the user is looking for, the user may write and use their own function.

Name	X	Y	Z	Process
f1	'year'	'profit'	v1 <- 'product'.*	
f2	'year'	'sales'	v1	v2 <- argmax _{v1} [k = 10]D(f1, f2)
*f3	'year'	'profit'	v2	

Table 9: Query which returns the top 10 profit visualizations for products which are most different from their sales visualizations.

Name	X	Y	Z	Process
f1	'year'	'sales'	v1 <- 'product'.*	v2 <- argmax _{v1} [t < 0]T(f1)
*f2	'year'	'sales'	v2	

Table 10: Query which returns the sales visualizations for all products which have a negative trend.

Concrete examples. With just dimension-based iteration, the optimization structure of the Process column, and the functional primitives T and D , we found that we were able to support the majority of the visual analysis tasks required by our users. Common patterns include filtering based on overall trend (Table 10), searching for the most similar visualization (Table 11), and determining outlier visualizations (Table 12). Table 10 describes a query where in the first row, the variable $v2$ selects all products whose trend is decreasing, and the second row visualizes these product’s sales over year. Table 11 starts with the visualization sales over year for chair in the first row, then in the second row computes the visualizations of sales over year for all products, and in the process column computes the similarity with chair, assigning the top 10 to $v2$, and the third row visualizes the sales over year for these products. Table 12 starts with the visualization collection of sales over year for all products in the first row, followed by another collection of the same in the second row, and in the process column computes the sum of pairwise distances, assigning the 10 products whose visualizations are most distant to others to $v3$, after which they are visualized. Table 13 features a realistic query inspired by one of our case studies. The overall goal of the query is to find the products which have positive sales and profits trends in locations and categories which have overall negative trends; the user may want to look at this set of products to see what makes them so special. Rows 1 and 2 specify the sales and profit visualizations for all locations and categories respectively, and the processes for these rows filter down to the locations and categories which have negative trends. Then rows 3 and 4 specify the sales and profit visualizations for products in these locations and categories, and the processes filter the visualizations down to the ones that have positive trends. Finally, row 5 takes the list of output products from the processes in rows 3 and 4 and takes the intersection of the two returning the sales and profits visualizations for these products.

Pluggable functions. While the general structure of the Process column does cover the majority of the use cases requested by our users, users may want to write their own functions to run in a ZQL query. To support this, ZQL exposes a java-based API for users to write their own functions. In fact, we use this interface to implement the k -means algorithm for ZQL. While the pluggable functions do allow virtually any capabilities to be implemented, it is preferred that users write their queries using the syntax of the Process column; pluggable functions are considered black-boxes and cannot be automatically optimized by the ZQL compiler.

2.2 Discussion of Capabilities and Limitations

Although ZQL can capture a wide range of visual exploration queries, it is not limitless. Here, we give a brief description of what ZQL can do. A more formal quantification can be found in Section 3.

ZQL’s primary goal is to support queries over visualizations—which are themselves aggregate group-by queries on data. Using

these queries, ZQL can compose a collection of visualizations, filter them in various ways, compare them against benchmarks or against each other, and sort the results. The functions T and D , while intuitive, support the ability to perform a range of computations on visualization collections—for example, any filter predicate on a single visualization, checking for a specific visual property, can be captured under T . With the pluggable functions, the ability to perform sophisticated computation on visualization collections is enhanced even further. Then, via the dimension-based iterators, ZQL supports the ability to chain these queries with each other and compose new visualization collections. These simple set of operations offer unprecedented power in being able to sift through visualizations to identify desired trends.

Since ZQL already operates one layer above the data—on the visualizations—it does not support the *creation of new derived data*: that is, ZQL does not support the generation of derived attributes or values not already present in the data. The new data that is generated via ZQL is limited to those from binning and aggregating via the Viz column. This limits ZQL’s ability to perform prediction—since feature engineering is an essential part of prediction; it also limits ZQL’s ability to compose visualizations on combinations of attributes at a time, e.g., $\frac{A1}{A2}$ on the X axis. Among other drawbacks of ZQL: ZQL does not support (i) recursion; (ii) any data modification; (iii) non-foreign-key joins nor arbitrary nesting; (iv) dimensionality reduction or other changes to the attributes; (v) other forms of processing visualization collections not expressible via T , D or the black box; (vi) multiple-dimensional visualizations; (vii) intermediate variable definitions; (viii) merging of visualizations (e.g., by aggregating two visualizations); and (ix) statistical tests.

3. EXPRESSIVENESS

In this section, we formally quantify the expressive power of ZQL. To this end, we formulate an algebra, called the *visual exploration algebra*. Like relational algebra, visual exploration algebra contains a basic set of operators that we believe all visual exploration languages should be able to express. At a high level, the operators of our visual exploration algebra operate on sets of visualizations and are not mired by the data representations of those visualizations, nor the details of how the visualizations are rendered. Instead, the visual exploration algebra is primarily concerned with the different ways in which visualizations can be selected, refined, and compared with each other.

Given a function T that operates on a visualization at a time, and a function D that operates on a pair of visualizations at a time, both returning real-valued numbers, a visual exploration language L is defined to be *visual exploration complete* $VEC_{T,D}(L)$ with respect to T and D if it supports all the operators of the visual exploration algebra. These functions T and D (also defined previously) are “functional primitives” without which the resulting algebra would have been unable to manipulate visualizations in the way

Name	X	Y	Z	Process
f1	'year'	'sales'	'product'. 'chair'	
f2	'year'	'sales'	v1 <- 'product'. (* - 'chair')	v2 <- $\text{argmin}_{v_1} [k=10] D(f1, f2)$
*f3	'year'	'sales'	v2	

Table 11: Query which returns the sales visualizations for the 10 products whose sales visualizations are the most similar to the sales visualization for the chair.

Name	X	Y	Z	Process
f1	'year'	'sales'	v1 <- 'product'. *	
f2	'year'	'sales'	v2 <- 'product'. *	v3 <- $\text{argmax}_{v_1} [k=10] \sum_{v_2} D(f1, f2)$
*f3	'year'	'sales'	v3	

Table 12: Query which returns the sales visualizations for the 10 products whose sales visualizations are the most different from the others.

we need for data exploration. Unlike relational algebra, which does not have any “black-box” functions, visual exploration algebra requires these functions for operating on visualizations effectively. That said, these two functions are flexible and configurable and up to the user to define (or left as system defaults). Next, we formally define the visual exploration algebra operators and prove that ZQL is visual exploration complete.

3.1 Ordered Bag Semantics

In visual exploration algebra, relations have bag semantics. However, since users want to see the most relevant visualizations first, ordering is critical. So, we adapt the operators from relational algebra to preserve ordering information.

Thus, we operate on *ordered bags* (i.e., a bag that has an inherent order). We describe the details of how to operate on ordered bags below. We use the variables R, S to denote ordered bags. We also use the notation $R = [t_1, \dots, t_n]$ to refer to an ordered bag, where t_i are the tuples.

The first operator that we define is an indexing operator, much like indexing in arrays. The notation $R[i]$ refers to the i th tuple within R , and $R[i : j]$ refers to the ordered bag corresponding to the list of tuples from the i th to the j th tuple, both inclusive. In the notation $[i : j]$ if either one of i or j is omitted, then it is assumed to be 1 for i , and n for j , where n is the total number of tuples.

Next, we define a union operator \cup : $R \cup S$ refers to the concatenation of the two ordered bags R and S . If one of R or S is empty, then the result of the union is simply the other relation. We define the union operation first because it will come in handy for subsequent operations.

We define the σ operator like in relational algebra, via a recursive definition:

$$\sigma_\theta(R) = \sigma_\theta([R[1]]) \cup \sigma_\theta(R[2 :])$$

where σ_θ when applied to an ordered bag with a single tuple ($[t]$) behaves exactly like in the relational algebra case, returning the same ordered bag ($[t]$) if the condition is satisfied, and the empty ordered bag ($[]$) if the condition is not satisfied. The π operator for projection is defined similarly to σ in the equation above, with the π operator on an ordered bag with a single tuple simply removing the irrelevant attributes from that tuple, like in relational algebra.

Then, we define the \setminus operator, for ordered bag difference. Here, the set difference operator operates on every tuple in the first ordered bag and removes it if it finds it in the second ordered bag. Thus:

$$R \setminus S = ([R[1]] \setminus S) \cup (R[2 :] \setminus S)$$

where $[t] \setminus S$ is defined like in relational algebra, returning $[t]$ if t is not in S , and $[]$ otherwise. The intersection operator \cap is defined similarly to \cup and \setminus .

Now, we can define the duplicate elimination operator as follows:

$$\delta(R) = [R[1]] \cup (R[2 :] \setminus [R[1]])$$

Thus, the duplication elimination operator preserves ordering, while maintaining the first copy of each tuple at the first position that it was found in the ordered bag.

Lastly, we have the cross product operator, as follows:

$$R \times S = ([R[1]] \times S) \cup (R[2 :] \times S)$$

where further we have

$$[t] \times S = ([t] \times [S[1]]) \cup ([t] \times S[2 :])$$

where $[t] \times [u]$ creates an ordered bag with the result of the cross product as defined in relational algebra.

Given these semantics for ordered bags, we can develop the visual exploration algebra.

3.2 Basic Notation

Assume we are given a k -ary relation \mathcal{R} with attributes (A_1, A_2, \dots, A_k) .

Let \mathcal{X} be the unary relation with attribute X whose values are the names of the attributes in \mathcal{R} that can appear on the x-axis. If the x-axis attributes are not specified by the user for relation \mathcal{R} , the default behavior is to include all attributes in \mathcal{R} : $\{A_1, \dots, A_k\}$. Let \mathcal{Y} be defined similarly with Y for attributes that can appear on the y-axis. Given \mathcal{R} , \mathcal{X} , and \mathcal{Y} , we define \mathcal{V} , the *visual universe*, as follows: $\mathcal{V} = \nu(\mathcal{R}) = \mathcal{X} \times \mathcal{Y} \times \left(\times_{i=1}^k \pi_{A_i}(\mathcal{R}) \cup \{*\} \right)$ where π is the projection operator from relational algebra and $*$ is a special wildcard symbol, used to denote all values of an attribute. Table 14 shows an example of what a sample \mathcal{R} and corresponding \mathcal{X} , \mathcal{Y} , and \mathcal{V} would look like. At a high level, the visual universe specifies all subsets of data that may be of interest, along with the intended attributes to be visualized. Unlike relational algebra, visual exploration algebra mixes schema and data elements, but in a special way in order to operate on a collection of visualizations.

Any subset relation $V \subseteq \mathcal{V}$ is called a *visual group*, and any $k+2$ -tuple from \mathcal{V} is called a *visual source*. The last k portions (or attributes) of a tuple from \mathcal{V} comprise the *data source* of the visual source. Overall, a visual source represents a visualization that can be rendered from a selected data source, and a set of visual sources is a visual group. The X and Y attributes of the visual source determine the x- and y- axes, and the selection on the data source is determined by attributes A_1, \dots, A_k . If an attribute has the wildcard symbol $*$ as its value, no subselection is performed on that attribute for the data source. For example, the third row of Table 14d is a visual source that represents the visualization with year as the x-axis and sales as the y-axis for chair products. Since the value of location is $*$, all locations are considered valid or pertinent for the data source. In relational algebra, the data source for the third row can be written as $\sigma_{\text{product}=\text{chair}}(\mathcal{R})$. The $*$ symbol therefore attempts to emulate the lack of presence of a selection condition on that attribute in the σ operator of the relational algebra. Readers familiar with OLAP will notice the similarity between the use of the symbol $*$ here and the GROUPING SETS functionality in SQL.

Name	X	Y	Z	Z2	Z3	Process
f1	'year'	'sales'	v1 <- 'location'.*			v2 <- argany _{v1} t < 0 T(f1)
f2	'year'	'profit'	v3 <- 'category'.*			v4 <- argany _{v3} t < 0 T(f2)
f3	'year'	'profit'	v5 <- 'product'.*	'location'.[? IN v2]	'category'.[? IN v4]	v6 <- argany _{v5} t > 0 T(f3)
f4	'year'	'sales'	v5	'location'.[? IN v2]	'category'.[? IN v4]	v7 <- argany _{v5} t > 0 T(f4)
*f5	'year'	{'profit', 'sales'}	v6 ~ v7			

Table 13: Query which returns the profit and sales visualizations for products which have positive trends in profit and sales in locations and categories which have overall negative trends.

year	month	product	location	sales	profit
2016	4	chair	US	623,000	314,000
2016	3	chair	US	789,000	410,000
2016	4	table	US	258,000	169,000
2016	4	chair	UK	130,000	63,000
		⋮			

(a) Example \mathcal{R}

X	Y
year	sales
month	profit

(b) \mathcal{X}

X	Y	year	month	product	location	sales	profit
year	sales	*	*	*	*	*	*
year	profit	*	*	*	*	*	*
year	sales	*	*	chair	*	*	*
year	sales	*	*	chair	US	*	*
				⋮			

(c) \mathcal{Y}

(d) \mathcal{V} for \mathcal{R}

Table 14: An example relation \mathcal{R} and its resultant \mathcal{X} , \mathcal{Y} , and \mathcal{V} .

Note that infinitely many visualizations can be produced from a single visual source, due to different granularities of binning, aggregation functions, and types of visualizations that can be constructed, since a visualization generation engine can use a visualization rendering grammar like ggplot [47] that provides that functionality. Our focus in defining the visual exploration algebra is to specify the inputs to a visualization and attributes of interest as opposed to the aesthetic aspects and encodings. Thus, for our discussion, we assume that each visual source maps to a singular visualization. Even if the details of the encoding and aesthetics are not provided, standard rules may be applied for this mapping as alluded earlier [25, 43] in Section 2.1. Furthermore, a visual source does not specify the data representation of the underlying data source; therefore the expressive power of visual exploration algebra is not tied to any specific backend data storage model. The astute reader will have noticed that the format for a visual source looks fairly similar to a collections of visualizations in ZQL; this is no accident. In fact, we will use the visualization collections of ZQL as a proxy to visual sources when proving that ZQL is visual exploration complete.

3.3 Functional Primitives

Earlier, we mentioned that a visual exploration algebra is visual exploration complete with respect to two functional primitives: T and D . Here we define the formal types for these functional primitives with respect to visual exploration algebra.

The function $T : \mathcal{V} \rightarrow \mathbb{R}$ returns a real number given a visual source. This function can be used to assess whether a trend: defined by the visualization corresponding to a specific visual source, is “increasing”, or “decreasing”, or satisfies some other fixed property. Many such T can be defined and used within the visual exploration algebra.

The function $D : \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$ returns a real number given a pair of visual sources. This function can be used to compare pairs of visualizations (corresponding to the visual sources) with respect to each other. The most natural way to define D is via some notion of distance, e.g., Earth Mover’s or Euclidian distance, but once again, the definition can be provided by the user or assumed as a fixed black box.

3.4 Visual Exploration Algebra Operators

Similar to how operators in ordered bag algebra operate on and result in ordered bags, operators in visual exploration algebra operate on and result in visual groups. Many of the symbols for operators in visual exploration algebra are also derived from relational algebra, with some differences. To differentiate, operators in visual exploration algebra are superscripted with a v (e.g., σ^v , τ^v).

The unary operators for visual exploration algebra include (i) σ^v for selection, (ii) τ^v for sorting a visual group based on the trend-estimating function T , (iii) μ^v for limiting the number of visual sources in a visual group, and (iv) δ^v for duplicate visual source removal. The binary operators include (i) \cup^v for union, (ii) \setminus^v for difference, (iii) β^v for replacing the attribute values of the visual sources in one visual group’s with another’s, (iv) ϕ^v to reorder the first visual group based on the visual sources’ distances to the visual sources of another visual group based on metric D , and (v) η^v to reorder the visual sources in a visual group based on their distance to a reference visual source from a singleton visual group based on D . These operators are described below, and listed in Table 15.

3.4.1 Unary Operators.

$\sigma_\theta^v(V)$: σ^v selects a visual group from V based on selection criteria θ , like ordered bag algebra. However, σ^v has a more restricted θ ; while \vee and \wedge may still be used, only the binary comparison operators $=$ and \neq are allowed. As an example, $\sigma_\theta^v(V)$ where $\theta = (X='year' \wedge Y='sales' \wedge year=* \wedge month=* \wedge product \neq * \wedge location='US' \wedge sales=* \wedge profit=*)$ from Table 16 on \mathcal{V} from Table 14 would result in the visual group of time vs. sales visualizations for different products in the US.

In this example, note that the product is specifically set to not equal $*$ so that the resulting visual group will include all products. On the other hand, the location is explicitly set to be equal to US. The other attributes, e.g., sales, profit, year, month are set to equal $*$: this implies that the visual groups are not employing any additional constraints on those attributes. (This may be useful, for example when those attributes are not relevant for the current visualization or set of visualizations.) As mentioned before, visual groups have the semantics of ordered bags. Thus, σ^v operates on one tuple at a time in the order they appear in V , and the result is in the same order the tuples are operated on.

$\tau_{F(T)}^v(V)$: τ^v returns the visual group sorted in an increasing order based on applying $F(T)$ on each visual source in V , where $F(T)$ is a procedure that uses function T . For example, $\tau_{-T}^v(V)$ might return the visualizations in V sorted in decreasing order of estimated slope. This operator is not present in the ordered bag semantics, but may be relevant when we want to reorder the ordered bag using a different criterion. The function F may be any higher-order function with no side effects. For a language to visual exploration complete, the language must be able to support any arbitrary F .

$\mu_k^v(V)$: μ^v returns the first k visual sources of V ordered in the same way they were in V . μ^v is equivalent to the LIMIT statement in SQL. μ^v is often used in conjunction with τ^v to retrieve the top- k visualizations with greatest increasing trends (e.g. $\mu_k^v(\tau_{-T}^v(V))$).

Operator	Name	Derived from Bag Algebra	Meaning	Unary/Binary
σ^v	Selection	Yes	Subselects visual sources	Unary
τ^v	Sort	No	Sorts visual sources in increasing order	Unary
μ^v	Limit	Yes	Returns first k visual sources	Unary
δ^v	Dedup	Yes	Removes duplicate visual sources	Unary
$\cup^v / \setminus^v / \cap^v$	Union/Diff/Int	Yes	Returns the union of/differences between/intersection of two visual groups	Binary
β^v	Swap	No	Returns a visual group in which values of an attribute in one visual group is replaced with values of the same attribute in another visual group	Binary
ϕ^v	Dist	No	Sorts a visual group based on pairwise distance to another visual group	Binary
η^v	Find	No	Sorts a visual group in increasing order based on their distances to a single reference visual source	Binary

Table 15: Visual Exploration Algebra Operators

When instead of a number k , the subscript to μ^v is actually $[a : b]$, then the items of V that are between positions a and b in V are returned. Thus μ^v offers identical functionality to the $[a : b]$ in ordered bag algebra, with the convenient functionality of getting the top k results by just having one number as the subscript. Instead of using μ^v , visual exploration algebra also supports the use of the syntax $V[i]$ to refer to the i th visual source in V , and $V[a : b]$ to refer to the ordered bag of visual sources from positions a to b .

$\delta^v(V)$: δ^v returns the visual sources in V with the duplicates removed, in the order of their first appearance. Thus, δ^v is defined identically to ordered bag algebra.

3.4.2 Binary Operators.

$V \cup^v U \mid V \setminus^v U \mid V \cap^v U$: Returns the union / difference / intersection of V and U . These operations are just like the corresponding operations in ordered bag algebra.

$\beta_A^v(V, U)$: β^v returns a visual group in which values of attribute A in V are replaced with the values of A in U . Formally, assuming A_i is the i th attribute of V and V has n total attributes: $\beta_{A_i}^v(V, U) = \pi_{A_1, \dots, A_{i-1}, A_{i+1}, \dots, A_n}(V) \times \pi_{A_i}(U)$. This can be useful for when the user would like to change an axis: $\beta_X(V, \sigma_{X=\text{year}}^v(V))$ will change the visual sources in V to have year as their x-axis. β^v can also be used to combine multiple dimensions as well. If we assume that V has multiple Y values, we can do $\beta_X^v(V, \sigma_{X \neq *}(V))$ to have the visual sources in V vary over both X and Y . This operator allows us to start with a set of visualizations and then “pivot” to focus on a different attribute, e.g., start with sales over time visualizations and pivot to look at profit. Thus, the operator allows us to transform the space of visual sources.

$\phi_{F(D), A_1, \dots, A_j}^v(V, U)$: ϕ^v sorts the visual sources in V in increasing order based on their distances to the corresponding visual sources in U . More specifically, ϕ^v computes $F(D)(\sigma_{A_1=a_1 \wedge \dots \wedge A_j=a_j}^v(V), \sigma_{A_1=a_1 \wedge \dots \wedge A_j=a_j}^v(U)) \forall a_1, \dots, a_j \in \pi_{A_1, \dots, A_j}(V)$ and returns an increasingly sorted V based on the results. If $\sigma_{A_1=a_1 \wedge \dots \wedge A_j=a_j}^v$ for either V or U ever returns a non-singleton visual group for any tuple (a_1, \dots, a_j) , the result of the operator is undefined.

$\eta_{F(D)}^v(V, U)$: η^v sorts the visual sources in V in increasing order based on their distances to a single reference visual source in singleton visual group U . Thus, $U = [t]$. η^v computes $F(D)(V[i], U[1]) \forall i \in \{1, \dots, |V|\}$, and returns a reordered V based on these values, where $F(D)$ is a procedure that uses D . If U has more than one visual source, the operation is undefined. η^v is useful for queries in which the user would like to find the top- k most similar visualizations to a reference: $\mu_k^v(\eta_D^v(V, U))$, where V is the set of candidates and U contains the reference. Once again, this operator is similar to τ^v , except that it operates on the results of the comparison of

X	Y	year	month	product	location	sales	profit
year	sales	*	*	chair	US	*	*
year	sales	*	*	table	US	*	*
year	sales	*	*	stapler	US	*	*
year	sales	*	*	printer	US	*	*

Table 16: Results of performing unary operators on V from Table 14: $\sigma_\theta^v(V)$ where $\theta = (X='year' \wedge Y='sales' \wedge year=* \wedge month=* \wedge product \neq * \wedge location='US' \wedge sales=* \wedge profit=*)$

individual visual sources to a specific visual source.

3.5 Proof of Visual Exploration Completeness

We now attempt to quantify the expressiveness of ZQL within the context of visual exploration algebra and the two functional primitives T and D . More formally, we prove the following theorem:

THEOREM 3.1. *Given well-defined functional primitives T and D , ZQL is visual exploration complete with respect to T and D : $VEC_{T,D}(\text{ZQL})$ is true.*

Our proof for this theorem involves two major steps:

- Step 1.** We show that a visualization collection in ZQL has as much expressive power as a visual group of visual exploration algebra, and therefore a visualization collection in ZQL serves as an appropriate proxy of a visual group in visual exploration algebra.
- Step 2.** For each operator in visual exploration algebra, we show that there exists a ZQL query which takes in visualization collection semantically equivalent to the visual group operands and produces visualization collection semantically equivalent to the resultant visual group.

LEMMA 3.2. *A visualization collection of ZQL has at least as much expressive power as a visual group in visual exploration algebra.*

PROOF. A visual group V , with n visual sources, is a relation with $k+2$ columns and n rows, where k is the number of attributes in the original relation. We show that for any visual group V , we can come up with a ZQL query q which can produce a visualization collection that represents the same set of visualizations as V .

Name	X	Y	Z1	...	Zk
f1	$\pi_X(V[1])$	$\pi_Y(V[1])$	$E_{1,1}$...	$E_{1,k}$
...
fn	$\pi_X(V[n])$	$\pi_Y(V[n])$	$E_{n,1}$...	$E_{n,k}$
*fn+1 <= f1+...+fn					

Table 17: ZQL query q which produces a visualization collection equal in expressiveness to visual group V .

Query q has the format given by Table 17, where $V[i]$ denotes the i th tuple of relation V and:

$$E_{i,j} = \begin{cases} "" & \text{if } \pi_{A_j}(V[i]) = * \\ A_j.\pi_{A_j}(V[i]) & \text{otherwise} \end{cases}$$

Here, A_j refers to the j th attribute of the original relation. The i th visual source of V is represented with the \mathbf{fi} from q . The X and Y values come directly from the visual source using projection. For the Z_j column, if the A_j attribute of visual source has any value than other than $*$, we must filter the data based on that value, so $E_{i,j} = A_j.\pi_{A_j}(V[i])$. However, if the A_j attribute is equal to $*$, then the corresponding element in \mathbf{fi} is left blank, signaling no filtering based on that attribute.

After, we have defined a visualization collection \mathbf{fi} for each i th visual source in V , we take the sum (or concatenation) across all these visualization collections as defined in Appendix A.3, and the resulting $\mathbf{fn}+1$ becomes equal to the visual group V . \square

LEMMA 3.3. $\sigma_\theta^v(V)$ is expressible in ZQL for all valid constraints θ and visual groups V .

PROOF. We prove this by induction.

The full context-free grammar (CFG) for θ in σ_θ^v can be given by:

$$\theta \rightarrow E \mid E \wedge E \mid E \vee E \mid \varepsilon \quad (1)$$

$$E \rightarrow C \mid (E) \mid E \wedge C \mid E \vee C \quad (2)$$

$$C \rightarrow T_1 = B_1 \mid T_1 \neq B_1 \mid T_2 = B_2 \mid T_2 \neq B_2 \quad (3)$$

$$T_1 \rightarrow X \mid Y \quad (4)$$

$$B_1 \rightarrow A_1 \mid \dots \mid A_k \quad (5)$$

$$T_2 \rightarrow A_1 \mid \dots \mid A_k \quad (6)$$

$$B_2 \rightarrow \text{string} \mid \text{number} \mid * \quad (7)$$

where ε represents an empty string (no selection), and X , Y , and A_1, \dots, A_k refer to the attributes of V .

To begin the proof by induction, we first show that ZQL is capable of expressing the base expressions $\sigma_C^v(V)$: $\sigma_{T_1=B_1}^v(V)$, $\sigma_{T_1 \neq B_1}^v(V)$, $\sigma_{T_2=B_2}^v(V)$, and $\sigma_{T_2 \neq B_2}^v(V)$. The high level idea for each of these proofs is to be come up with a *filtering visual group* U which we take the intersection with to arrive at our desired result: $\exists U, \sigma_C^v(V) = V \cap^v U$.

In the first two expressions, T_1 and B_1 refer to filters on the X and Y attributes of V ; we have the option of either selecting a specific attribute ($T_1 = B_1$) or excluding a specific attribute ($T_1 \neq B_1$). Tables 18 and 19 show ZQL queries which express $\sigma_{T_1=B_1}^v(V)$ for $T_1 \rightarrow X$ and $T_1 \rightarrow Y$ respectively. The ZQL queries do the approximate equivalent of $\sigma_{T_1=B_1}^v(V) = V \cap^v \sigma_{T_1=B_1}^v(V)$.

Name	X	Y	Z1	...	Zk
f1	-	-	-	...	-
f2 <- f1	-	y1 <- -	v1 <- A1.-	...	vk <- Ak.-
f3	B1	y1	v1	...	vk
*f4 <- f1 ^ f3					

Table 18: ZQL query which expresses $\sigma_{X=B_1}^v(V)$.

Name	X	Y	Z1	...	Zk
f1	-	-	-	...	-
f2 <- f1	x1 <- -	-	v1 <- A1.-	...	vk <- Ak.-
f3	x1	B1	v1	...	vk
*f4 <- f1 ^ f3					

Table 19: ZQL query which expresses $\sigma_{Y=B_1}^v(V)$.

We have shown with Lemma 3.2 that a visualization collection is capable of expressing a visual group, so we assume that $\mathbf{f1}$, the visualization collection which represents the operand V , is given

to us for both of these tables. Since we do not know how $\mathbf{f1}$ was derived, we use - for its axis variable columns. The second rows of these tables derive $\mathbf{f2}$ from $\mathbf{f1}$ and bind axis variables to the values of the non-filtered attributes. Here, although the set of visualizations present in $\mathbf{f2}$ is exactly the same as $\mathbf{f1}$, we now have a convenient way to iterate over the non-filtered attributes of $\mathbf{f1}$ (for more information on derived visualization collections, please refer to Appendix A.3). The third row combines the specified attribute B_1 with the non-filtered attributes of $\mathbf{f2}$ to form the *filtering visualization collection* $\mathbf{f3}$, which expresses the filtering visual group U from above. We then take the intersection between $\mathbf{f1}$ and the filtering visualization collection $\mathbf{f3}$ to arrive at our desired visualization collection $\mathbf{f4}$, which represents the resultant visual group $\sigma_{T_1=B_1}^v(V)$. Although, we earlier said that we would come up with $\mathbf{f3} = \sigma_{T_1=B_1}^v(V)$, in truth, we come up with $\mathbf{f3} = B_1 \times \pi_{Y,A_1,\dots,A_k}(V)$ for $T_1 \rightarrow X$ and $\mathbf{f3} = \pi_{X,A_1,\dots,A_k}(V) \times B_1$ for $T_1 \rightarrow Y$ because they are easier to express in ZQL; regardless we still end up with the correct resulting set of visualizations.

Tables 20 and 21 show ZQL queries which express $\sigma_{T_1 \neq B_1}^v(V)$ for $T_1 \rightarrow X$ and $T_1 \rightarrow Y$ respectively. Similar to the queries above, these queries perform the approximate equivalent of $\sigma_{T_1 \neq B_1}^v(V) = V \cap^v \sigma_{T_1 \neq B_1}^v(V)$. We once again assume $\mathbf{f1}$ is a given visualization collection which represents the operand V , and we come up with a filtering visualization collection $\mathbf{f3}$ which mimics the effects of (though is not completely equivalent to) $\sigma_{T_1 \neq B_1}^v(V)$. We then take the intersection between $\mathbf{f1}$ and $\mathbf{f3}$ to arrive at $\mathbf{f4}$ which represents the resulting $\sigma_{T_1 \neq B_1}^v(V)$.

The expressions $\sigma_{T_2=B_2}^v$ and $\sigma_{T_2 \neq B_2}^v$ refer to filters on the A_1, \dots, A_k attributes of V . Specifically, T_2 is some attribute $A_j \in \{A_1, \dots, A_k\}$ and B_2 is the attribute value which is selected or excluded. Here, we have an additional complication to the proof since any attribute A_j can also filter for or exclude $*$. First, we show ZQL is capable of expressing $\sigma_{T_2=B_2'}^v$ and $\sigma_{T_2 \neq B_2'}^v$ for which $B_2' \neq *$; that is B_2' is any attribute value which is not $*$. Tables 22 and 23 show the ZQL queries which express $\sigma_{T_2=B_2'}^v(V)$ and $\sigma_{T_2 \neq B_2'}^v(V)$ respectively. Note the similarity between these queries and the queries for $\sigma_{T_1=B_1}^v(V)$ and $\sigma_{T_1 \neq B_1}^v(V)$.

For $\sigma_{T_2=*}^v(V)$ and $\sigma_{T_2 \neq *}^v(V)$, Tables 24 and 25 show the corresponding queries. In Table 24, we explicitly avoid setting a value for Z_j for $\mathbf{f3}$ to emulate $A_j = *$ for the filtering visualization collection. In Table 25, $\mathbf{f3}$'s Z_j takes on all possible values from A_j , $*$, but that means that a value is set for Z_j (i.e., $T_2 \neq *$).

Now that we have shown how to express the base operations, we next assume ZQL is capable of expressing any arbitrary complex filtering operations $\sigma_{E'}^v$ where E' comes from Line 2 of the CFG. Specifically, we assume that given a visualization collection $\mathbf{f1}$ which expresses V , there exists a filtering visualization collection $\mathbf{f2}$ for which $\sigma_{E'}^v(V) = \mathbf{f1} \wedge \mathbf{f2}$. Given this assumption, we now must take the inductive step, apply Line 2, and prove that $\sigma_{E \rightarrow (E')}^v(V)$, $\sigma_{E \rightarrow E' \wedge C}^v(V)$, and $\sigma_{E \rightarrow E' \vee C}^v(V)$ are all expressible in ZQL for any base constraint C .

$\sigma_{E \rightarrow (E')}^v(V)$: This case is trivial. Given $\mathbf{f1}$ which represents V and $\mathbf{f2}$ which is the filtering visualization collection for E' , we simply the intersect the two to get $\mathbf{f3} \leftarrow \mathbf{f1} \wedge \mathbf{f2}$ which represents $\sigma_{E \rightarrow (E')}^v(V)$.

$\sigma_{E \rightarrow E' \wedge C}^v$: Once again assume we are given $\mathbf{f1}$ which represents V and $\mathbf{f2}$ which is the filtering visualization collection of E' . Based on the base expression proofs above, we know that given any base constraint C , we can find a filtering visualization collection for it; call this filtering visualization collection $\mathbf{f3}$. We can then see that $\mathbf{f2} \wedge \mathbf{f3}$ is the filtering visualization collection of $E \rightarrow E' \wedge C$, and $\mathbf{f4}$

Name	X	Y	Z1	...	Zk
f1	-	-	-	...	-
f2 <- f1	x1 <- -	y1 <- -	v1 <- A1.-	...	vk <- Ak.-
f3	x2 <- x1 - {B1}	y1	v1	...	vk
*f4 <- f1^f3					

Table 20: ZQL query which expresses $\sigma_{X \neq B_1}^v(V)$.

Name	X	Y	Z1	...	Zk
f1	-	-	-	...	-
f2 <- f1	x1 <- -	y1 <- -	v1 <- A1.-	...	vk <- Ak.-
f3	x1	y2 <- y1 - {B1}	v1	...	vk
*f4 <- f1^f3					

Table 21: ZQL query which expresses $\sigma_{Y \neq B_1}^v(V)$.

Name	X	Y	Z1	...	Zj	...	Zk
f1	-	-	-	...	-	...	-
f2 <- f1	x1 <- -	y1 <- -	v1 <- A1.-	...	-	...	vk <- Ak.-
f3	x1	y1	v1	...	B'_2	...	vk
*f4 <- f1^f3							

Table 22: ZQL query which expresses $\sigma_{A_j=B'_2}^v(V)$ when $B'_2 \neq *$.

Name	X	Y	Z1	...	Zj	...	Zk
f1	-	-	-	...	-	...	-
f2 <- f1	x1 <- -	y1 <- -	v1 <- A1.-	...	vj <- Aj.-	...	vk <- Ak.-
f3	x1	y1	v1	...	uj <- vj - {B'_2}	...	vk
*f4 <- f1^f3							

Table 23: ZQL query which expresses $\sigma_{A_j \neq B'_2}^v(V)$ when $B'_2 \neq *$.

Name	X	Y	Z1	...	Zj	...	Zk
f1	-	-	-	...	-	...	-
f2 <- f1	x1 <- -	y1 <- -	v1 <- A1.-	...	-	...	vk <- Ak.-
f3	x1	y1	v1	...	-	...	vk
*f4 <- f1^f3							

Table 24: ZQL query which expresses $\sigma_{A_j=*}^v(V)$.

Name	X	Y	Z1	...	Zj	...	Zk
f1	-	-	-	...	-	...	-
f2 <- f1	x1 <- -	y1 <- -	v1 <- A1.-	...	-	...	vk <- Ak.-
f3	x1	y1	v1	...	vj <- Aj.*	...	vk
*f4 <- f1^f3							

Table 25: ZQL query which expresses $\sigma_{A_j \neq *}^v(V)$.

$\langle - f1 \wedge (f2 \wedge f3) \rangle$ represents $\sigma_{E \rightarrow E \wedge C}^v(V)$.

$\sigma_{E \rightarrow E \vee C}^v$: Once again assume we are given f1 which represents V , f2 which is the filtering visualization collection of E' , and we can find a filtering visualization collection f3 for C . We can then see that $f2 + f3$ is the filtering visualization collection of $E \rightarrow E' \vee C$, and $f4 \langle - f1 \wedge (f2 + f3) \rangle$ represents $\sigma_{E \rightarrow E \vee C}^v(V)$.

With this inductive step, we have shown that for all complex constraints E of the form given by Line 2 of the CFG, we can find a ZQL query which expresses $\sigma_E^v(V)$. Given this, we can finally show that ZQL is capable of expressing $\sigma_\theta^v(V)$ for all θ : $\sigma_{\theta \rightarrow E}^v(V)$, $\sigma_{\theta \rightarrow E \wedge E'}^v(V)$, $\sigma_{\theta \rightarrow E \vee E'}^v(V)$, and $\sigma_{\theta \rightarrow \varepsilon}^v(V)$.

$\sigma_{\theta \rightarrow E}^v(V)$: This case is once again trivial. Assume, we are given f1 which represents V , and f2, which is the filtering visualization collection of E , f3 $\langle - f1 \wedge f2 \rangle$ represents $\sigma_{\theta \rightarrow E}^v(V)$.

$\sigma_{\theta \rightarrow E \wedge E'}^v(V)$: Assume, we are given f1 which represents V , f2, which is the filtering visualization collection of E , and f3, which is the filtering visualization collection of E' . $f2 \wedge f3$ is the filtering visualization collection of $\theta \rightarrow E \wedge E'$, and $f4 \langle - f1 \wedge (f2 \wedge f3) \rangle$ represents $\sigma_{\theta \rightarrow E \wedge E'}^v(V)$.

$\sigma_{\theta \rightarrow E \vee E'}^v(V)$: Assume, we are given f1 which represents V , f2 which is the filtering visualization collection of E , and f3 which is the filtering visualization collection of E' . $f2 + f3$ is the filtering visualization collection of $\theta \rightarrow E \vee E'$, and $f4 \langle - f1 \wedge (f2 + f3) \rangle$ represents $\sigma_{\theta \rightarrow E \vee E'}^v(V)$.

$\sigma_{\theta \rightarrow \varepsilon}^v(V)$: This is the case in which no filtering is done. Therefore, given f1 which represents V , we can simply return f1. \square

PROOF. Assume f1 is the visualization collection which represents V . Query q given by Table 26 produces visualization collection f3 which expresses $\tau_{F(T)}^v(V)$. \square

LEMMA 3.5. $\mu_{[a;b]}^v(V)$ is expressible in ZQL for all valid intervals $a : b$ and visual groups V .

PROOF. Assume f1 is the visualization collection which represents V . Query q given by Table 27 produces visualization collection f2 which expresses $\mu_{[a;b]}^v(V)$. \square

LEMMA 3.6. $\delta^v(V)$ is expressible in ZQL for all valid visual groups V .

PROOF. Assume f1 is the visualization collection which represents V . Query q given by Table 28 produces visualization collection f2 which expresses $\delta^v(V)$. \square

LEMMA 3.7. $V \cup^v U$ is expressible in ZQL for all valid visual groups V and U .

PROOF. Assume f1 is the visualization collection which represents V and f2 represents U . Query q given by Table 29 produces visualization collection f3 which expresses $V \cup^v U$. \square

LEMMA 3.8. $V \setminus^v U$ is expressible in ZQL for all valid visual groups V and U .

PROOF. Assume f1 is the visualization collection which represents V and f2 represents U . Query q given by Table 30 produces visualization collection f3 which expresses $V \setminus^v U$. The proof for \cap^v can be shown similarly. \square

LEMMA 3.9. $\beta_A^v(V, U)$ is expressible in ZQL for all valid attributes A in \mathcal{V} and visual groups V and U .

LEMMA 3.4. $\tau_{F(T)}^v(V)$ is expressible in ZQL for all valid functionals F of T and visual groups V .

Name	X	Y	Z1	...	Zk	Process
f1	-	-	-	...	-	-
f2 <- f1	x1 <- _	y1 <- _	v1 <- A1._	..	vk <- Ak._	x2, y2, u1, ..., uk <- argmin_{x1,y1,v1,...,vk}[k = \infty]F(T)(f2)
*f3	x2	y2	u1	...	uk	

Table 26: ZQL query q which expresses $\tau_{F(T)}^V(V)$.

Name	X	Y	Z1	...	Zk	Process
f1	-	-	-	...	-	
*f2 <- f1[a:b]						

Table 27: ZQL query q which expresses $\mu_{[a:b]}^V(V)$.

Name	X	Y	Z1	...	Zk	Process
f1	-	-	-	...	-	
*f2 <- f1						

Table 28: ZQL query q which expresses $\delta^V(V)$.

Name	X	Y	Z1	...	Zk	Process
f1	-	-	-	...	-	
f2	-	-	-	...	-	
*f3 <- f1+f2						

Table 29: ZQL query q which expresses $V \cup^V U$.

Name	X	Y	Z1	...	Zk	Process
f1	-	-	-	...	-	
f2	-	-	-	...	-	
*f3 <- f1-f2						

Table 30: ZQL query q which expresses $V \setminus^V U$.

Name	X	Y	Z1	...	Zk
f1	-	-	-	...	-
f2	-	-	-	...	-
f3 <- f1		y1 <- _	v1 <- A1._	...	vk <- Ak._
f4 <- f2	x1 <- _				
*f5	x1 ²	y1 ¹	v1 ¹	...	vk ¹

Table 31: ZQL query q which expresses $\beta_A^V(V, U)$ where $A = X$.

Name	X	Y	Z1	...	Zk
f1	-	-	-	...	-
f2	-	-	-	...	-
f3 <- f1	x1 <- _		v1 <- A1._	...	vk <- Ak._
f4 <- f2		y1 <- _			
*f5	x1 ¹	y1 ²	v1 ¹	...	vk ¹

Table 32: ZQL query q which expresses $\beta_A^V(V, U)$ where $A = Y$.

Name	X	Y	Z1	...	Zj-1	Zj	Zj+1	...	Zk
f1	-	-	-	...	-	-	-	...	-
f2	-	-	-	...	-	-	-	...	-
f3 <- f1	x1 <- _	y1 <- _	v1 <- A1._	...	vj-1 <- Aj-1._		vj+1 <- Aj+1._	...	vk <- Ak._
f4 <- f2						uj <- Aj._			
*f5	x1 ¹	y1 ¹	v1 ¹	...	vj-1 ¹	uj ²	vj+1 ¹	...	vk ¹

Table 33: ZQL query q which expresses $\beta_A^V(V, U)$ where $A = A_j$ and A_j is an attribute from \mathcal{R} .

PROOF. Assume f_1 is the visualization collection which represents V and f_2 represents U . There are three cases we must handle depending on the value of A due to the structure of columns in ZQL: (i) $A = X$ (ii) $A = Y$ (iii) $A = A_j$ where A_j is an attribute from the original relation \mathcal{R} . For each of the three cases, we produce a separate query which expresses β^v . For $A = X$, the query given by Table 31 produces f_5 which is equivalent to $\beta_X^v(V, U)$. We use the superscripts in the last row so that cross product conforms to the ordering defined in Section 3.4. For more information about the superscripts, please refer to Appendix A.3. For $A = Y$, the query given by Table 32 produces f_5 which is equivalent to $\beta_Y^v(V, U)$, and for $A = A_j$, the query given by Table 33 produces f_5 which is equivalent to $\beta_{A_j}^v(V, U)$. \square

LEMMA 3.10. $\phi_{F(D), A_1, \dots, A_j}^v(V, U)$ is expressible in ZQL for all valid attributes A_1, \dots, A_j and visual groups V and U .

PROOF. Assume f_1 is the visualization collection which represents V , and f_2 represents U . Without loss of generality, assume the attributes we want to match on (A_1, \dots, A_j) are the first j attributes of \mathcal{R} . Query q given by Table 34 produces visualization collection f_5 which expresses $\phi_{F(D), A_1, \dots, A_j}^v(V, U)$. In the table, we first retrieve the values for (A_1, \dots, A_j) using f_3 and reorder f_2 based on these values to get f_4 . We then compare the visualizations in f_3 and f_4 with respect to (A_1, \dots, A_j) using the distance function $F(D)$ and retrieve the increasingly sorted (A_1, \dots, A_j) values from the *argmin*. We are guaranteed that visualizations in f_3 and f_4 match up perfectly with respect to (A_1, \dots, A_j) since the definition in Section 3.4 allows exactly one visual source to result from any $\sigma_{A_1=a_1 \wedge \dots \wedge A_j=a_j}^v$. Finally, we reorder f_1 according to these values to retrieve f_5 . For more information on the *order* operation, please refer to Appendix A.3. \square

LEMMA 3.11. $\eta_{F(D)}^v(V, U)$ is expressible in ZQL for all valid functionals F of D and visual groups V and singleton visual groups U .

PROOF. Assume f_1 is the visualization collection which represents V and f_2 represents U . Query q given by Table 35 produces visualization collection f_4 which expresses $\eta_{F(D)}^v(V, U)$. \square

Although we have come up with a formalized algebra to measure the expressiveness of ZQL, ZQL is actually more expressive than visual exploration algebra. For example, ZQL allows the user to nest multiple levels of iteration in the Process column as in Table 12. Nevertheless, visual exploration algebra serves as a useful minimum metric for determining the expressiveness of visual exploration languages. Other visual analytics tools like Tableau are capable of expressing the selection operator σ^v in visual exploration algebra, but they are incapable of expressing the other operators which compare and filter visualizations based on functional primitives T and D . General purpose programming languages with analytics libraries such as Python and Scikit-learn [37] are visual exploration complete since they are Turing-complete, but ZQL's declarative syntax strikes a novel balance between simplicity and expressiveness which allows even non-programmers to become data analysts as we see in Section 7.

4. QUERY EXECUTION

In zenvisage, ZQL queries are automatically parsed and executed by the back-end. The ZQL compiler translates ZQL queries into a combination of SQL queries to fetch the visualization collections and processing tasks to operate on them. We present a basic graph-based translation for ZQL and then provide several optimizations to the graph which reduce the overall runtime considerably.

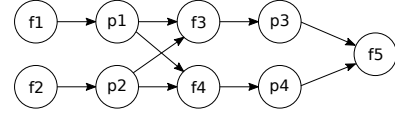


Figure 2: The query plan for the query presented in Table 13.

4.1 Basic Translation

Every valid ZQL query can be transformed into a query plan in the form of a directed acyclic graph (DAG). The DAG contains c -nodes (or *collection nodes*) to represent the collections of visualizations in the ZQL query and p -nodes (or *process nodes*) to represent the optimizations (or *processes*) in the Process column. Directed edges are drawn between nodes that have a dependency relationship. Using this query plan, the ZQL engine can determine at each step which visualization collection to fetch from the database or which process to execute. The full steps to build a query plan for any ZQL query is as follows: (i) Create a c -node or *collection node* for every collection of visualizations (including singleton collections). (ii) Create a p -node or *processor node* for every optimization (or *process*) in the Process column. (iii) For each c -node, if any of its axis variables are derived as a result of a process, connect a directed edge from the corresponding p -node. (iv) For each p -node, connect a directed edge from the c -node of each collection which appears in the process. Following these steps, we can translate our realistic query example in Table 13 to its query plan presented in Figure 2. Here, the c -nodes are annotated with $f\#$, and the p -nodes are annotated with $p\#$ (the i th p -node refers to the process in the i th row of the table). Further, f_1 is a root node with no dependencies since it does not depend on any process, whereas f_5 depends on the results of both p_3 and p_4 and has edges coming from both of them. Once the query plan has been constructed, the ZQL engine can execute it using the simple algorithm presented in Algorithm 1.

ALGORITHM 1. Algorithm to execute ZQL query plan:

1. Search for a node with either no parents or one whose parents have all been marked as done.
2. Run the corresponding task for that node and mark the node as done.
3. Repeat steps 1 and 2 until all nodes have been marked as done.

For c -nodes, the corresponding task is to retrieve the data for visualization collection, while for p -nodes, the corresponding task is to execute the process.

c -node translation: At a high level, for c -nodes, the appropriate SQL group-by queries are issued to the database to compose the data for multiple visualizations at once. Specifically, for the simplest setting where there are no collections specified for X or Y , a SQL query in the form of:

```
SELECT X, A(Y), Z, Z2, ... WHERE C(X, Y, Z, Z2, ...)
GROUP BY X, Z, Z2, ... ORDER BY X, Z, Z2, ...
```

is issued to the database, where X is the X column attribute, Y is the Y column attribute, $A(Y)$ is the aggregation function on Y (specified in the Viz column), $Z, Z2, \dots$ are the attributes/dimensions we are iterating over in the Z columns, while $C(X, Y, Z, Z2, \dots)$ refers to any additional constraints specified in the Z columns. The ORDER BY is inserted to ensure that all rows corresponding to a visualization are grouped together, in order. As an example, the SQL query for the c -node for f_1 in Table 12 would have the form:

```
SELECT year, SUM(sales), product
GROUP BY year, product ORDER BY year, product
```

If a collection is specified for the y -axis, each attribute in the collection is appended to the SELECT clause. If a collection is specified for the x -axis, a separate query must be issued for every X attribute in the collection. The results of the SQL query are then packed into

Name	X	Y	Z1	...	Zj	Process
f1	-	-	-	...	-	u1, ..., uj <- argmin_{v1, ..., vj} [k = ∞] F(D)(f3, f4)
f2	-	-	-	...	-	
f3 <- f1			v1 <- A1, _	...	vj <- Aj, _	
f4 <- f2.order			v1 >	...	vj >	
*f5 <- f1.order			u1 >	...	uj >	

Table 34: ZQL query q which expresses $\phi_{F(D), A_1, \dots, A_j}^v(V, U)$.

Name	X	Y	Z1	...	Zk	Process
f1	-	-	-	...	-	x2, y2, u1, ..., uk <- argmin_{x1, y1, v1, ..., vk} [k = ∞] F(D)(f3, f2)
f2	-	-	-	...	-	
f3 <- f1	x1 <- _	y1 <- _	v1 <- A1, _	...	vk <- Ak, _	
*f4	x2	y2	u1	...	uk	

Table 35: ZQL query q which expresses $\eta_{F(D)}^v(V, U)$.

a m -dimensional array (each dimension in the array corresponding to a dimension in the collection) and labeled with its $f\#$ tag.

p -node translation: At a high level, for p -nodes, depending on the structure of the expression within the process, the appropriate pseudocode is generated to operate on the visualizations. To illustrate, say our process is trying to find the top-10 values for which a trend is maximized/minimized with respect to various dimensions (using T), and the process has the form:

$$\langle argopt \rangle_{v0} [k = k'] \left[\langle op1 \rangle_{v1} \left[\langle op2 \rangle_{v2} \cdots \left[\langle opm \rangle_{vm} T(f1) \right] \right] \right] \quad (8)$$

where $\langle argopt \rangle$ is one of argmax or argmin , and $\langle op \rangle$ refers to one of $(\max | \min | \sum | \prod)$. Given this, the pseudocode which optimizes this process can automatically be generated based on the actual values of $\langle argopt \rangle$, $\langle op \rangle$, and the number of operations. In short, for each $\langle op \rangle$ or dimension traversed over, the ZQL engine generates a new nested for loop. Within each for loop, we iterate over all values of that dimension, evaluate the inner expression, and then eventually apply the overall operation (e.g., \max , \sum). For Equation 8, the generated pseudocode would look like the

one given by Listing 1. Here, f refers to the visualization collection being operated on by the p -node, which the parent c -node should have already retrieved.

```
f = make_ndarray(SQL(...))
tmp0 = make_array(size=len(v0))
for i0 in [1 .. len(v0)]:
  tmp1 = make_array(size=len(v1))
  for i1 in [1 .. len(v1)]:
    tmp2 = make_array(size=len(v2))
    for i2 in [1 .. len(v2)]:
      ...
      tmpm = make_array(size=len(vn))
      for im in [1 .. len(vn)]:
        tmpm[im] = T(f0, i1, i2, ..., im)
        tmpm-1[im-1] = opm(tmpm)
      ...
    tmp1[i1] = op1(tmp2)
  tmp0[i0] = op0(tmp1)
return argopt(tmp0)[k']
```

Listing 1: Pseudocode for a process in the form of Equation 8.

Although this is the translation for one specific type of process, it is easy to see how the code generation would generalize to other patterns.

4.2 Optimizations

We now present several optimizations to the previously introduced basic translator. In preliminary experiments, we found that the SQL queries for the c -nodes took the majority of the runtime for ZQL queries, so we concentrate our efforts on reducing the cost of these c -nodes. However, we do present one p -node-based optimization for process-intensive ZQL queries. We start with the simplest optimization schemes, and add more sophisticated variations later.

4.2.1 Parallelization

One natural way to optimize the graph-based query plan is to take advantage of the multi-query optimization (MQO) [41] present

in databases and issue in parallel the SQL queries for independent c -nodes—the c -nodes for which there is no dependency between them. With MQO, the database can receive multiple SQL queries at the same time and share the scans for those queries, thereby reducing the number of times the data needs to be read from disk.

To integrate this optimization, we make two simple modifications to Algorithm 1. In the first step, instead of searching for a single node whose parents have all been marked done, search for *all* nodes whose parents have been marked as done. Then in step 2, issue the SQL queries for all c -nodes which were found in step 1 in parallel at the same time. For example, the SQL queries for $f1$ and $f2$ could be issued at the same time in Figure 2, and once $p1$ and $p2$ are executed, SQL queries for $f3$ and $f4$ can be issued in parallel.

4.2.2 Speculation

While parallelization gives the ZQL engine a substantial increase in performance, we found that many realistic ZQL queries intrinsically have a high level of interdependence between the nodes in their query plans. To further optimize the performance, we use *speculation*, i.e., the ZQL engine pre-emptively issues SQL queries to retrieve the superset of visualizations for each c -node, considering all possible outcomes for the axis variables. Specifically, by tracing the provenance of each axis variable back to the root, we can determine the superset of all values for each axis variable; then, by considering the cartesian products of these sets, we can determine a superset of the relevant visualization collection for a c -node. After the SQL queries have returned, the ZQL engine proceeds through the graph as before, and once all parent p -nodes for a c -node have been evaluated, the ZQL engine isolates the correct subset of data for that c -node from the pre-fetched data.

For example, in the query in Table 13, $f3$ depends on the results of $p1$ and $p2$ since it has constraints based on $v2$ and $v4$; specifically $v2$ and $v4$ should be locations and categories for which $f1$ and $f2$ have a negative trend. However, we note that $v2$ and $v4$ are derived as a result of $v1$ and $v3$, specified to take on all locations and categories in rows 1 and 2. So, a superset of $f3$, the set of profit over year visualizations for various products for all locations and categories (as opposed to just those that satisfy $p1$ and $p2$), could be retrieved pre-emptively. Later, when the ZQL engine executes $p1$ and $p2$, this superset can be filtered down correctly.

One downside of speculation is that a lot more data must be retrieved from the database, but we found that blocking on the retrieval of data was more expensive in runtime than retrieving extra data. Thus, speculation ends up being a powerful optimization which compounds the positive effects of parallelization.

4.2.3 Query Combination

From extensive modeling of relational databases, we found that the overall runtime of concurrently running issuing SQL queries is heavily dependent on the number of queries being run in parallel. Each additional query constituted a T_q increase in the overall runtime (e.g., for our settings of PostgreSQL, we found $T_q = \sim 900\text{ms}$).

To reduce the total number of running queries, we use *query combination*; that is, given two SQL queries Q_1 and Q_2 , we combine these two queries into a new Q_3 which returns the data for both Q_1 and Q_2 . In general, if we have Q_1 (and Q_2) in the form of:

```
SELECT X1, A(Y1), Z1 WHERE C1(X1, Y1, Z1)
GROUP BY X, Z1 ORDER BY X, Z1
```

we can produce a combined Q_3 which has the form:

```
SELECT X1, A(Y1), Z1, C1, X2, A(Y2), Z2, C2
WHERE C1 or C2
GROUP BY X1, Z1, C1, X2, Z2, C2
ORDER BY X1, Z1, C1, X2, Z2, C2
```

where $C1 = C1(X1, Y1, Z1)$ and $C2$ is defined similarly. From the combined query Q_3 , it is possible to regenerate the data which would have been retrieved using queries Q_1 and Q_2 by aggregating over the non-related groups for each query. For Q_1 , we would select the data for which $C1$ holds, and for each $(X1, Z1)$ pair, we would aggregate over the $X2, Z2$, and $C2$ groups.

While query combining is an effective optimization, there are limitations. We found that the overall runtime also depends on the number of unique group-by values per query, and the number of unique group-by values for a combined query is the product of the number of unique group-by values of the constituent queries. Thus, the number of average group-by values per query grows super-linearly with respect to the number of combinations. However, we found that as long as the combined query had less than M_G unique group-by values, it was more advantageous to combine than not (e.g., for our settings of PostgreSQL, we found $M_G = 100k$).

Formulation. Given the above findings, we can now formulate the problem of deciding which queries to combine as an optimization problem: *Find the best combination of SQL queries that minimizes: $\alpha \times (\text{total number of combined queries}) + \sum_i (\text{number of unique group-by values in combined query } i)$, such that no single combination has more than M_G unique group-by values.*

The cost of adding a thread, α , is generally more than M_G —for instance, in PostgreSQL we found $\alpha > 100k (M_G)$ for different experimental settings. By further assuming that the cost of processing all group by values $< M_G$ is same, we can simplify the problem to finding the minimum number of combined queries such that the maximum number of group by values per combined query is less than M_G . We prove that the solution to this problem is NP-HARD by reduction from the PARTITION PROBLEM.

PROOF. Let g_1, g_2, \dots, g_n be the group by values for the queries Q_1, Q_2, \dots, Q_n we want to combine. We want to find minimum number m of combined queries, such that each combined query G_i has at most M_G maximum group by values. Recall that in the Partition problem, we are given an instance of n numbers a_1, a_2, \dots, a_n , and we are asked to decide if there is a set S such that $\sum_{a_i \in S} a_i = \sum_{a_i \notin S} a_i$. Let $A = \sum a_i$ and consider an instance of Query Combination problem with $g_i = M_G^{2 \times a_i}$. With this setting, it is easy to see that the answer to the Partition instance is YES if and only if the minimum number of combined queries is 2. \square

Wrinkle and Solution. However, a wrinkle to the above formulation is that it assumes no two SQL queries share a group-by attribute. If two queries have a shared group-by attribute, it may be more beneficial to combine those two, since the number of group-by values does not go up on combining them. Overall, we developed the metric *EFGV* or the effective increase in the number of group-by values to determine the utility of combining query Q' to query Q : $EFGV_Q(Q') = \prod_{g \in G(Q')} \#(g)^{[g \notin G(Q)]}$ where $G(Q)$ is the set of group-by values in Q , $\#(g)$ calculates the number of unique group-by values in g , and $[g \notin G(Q)]$ returns 1 if $g \notin G(Q)$ and 0

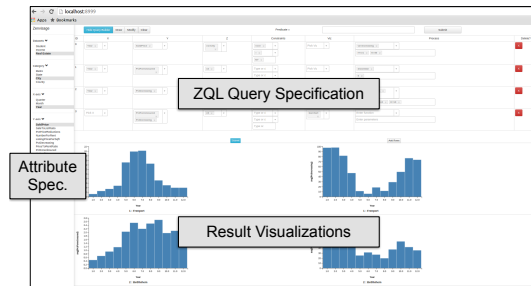


Figure 3: zenvisage basic functionalities

otherwise. In other words, this calculates the product of group-by values of the attributes which are in Q' but not Q . Using the *EFGV* metric, we then apply a variant of agglomerative clustering [10] to decide the best choice of queries to combine. As we show in the experiments section, this technique leads to very good performance.

4.2.4 Cache-Aware Execution

Although the previous optimizations were all I/O-based optimizations for ZQL, there are cases in which optimizing the execution of p -nodes is important as well. In particular, when a process has multiple nested for loops, the cost of the p -node may start to dominate the overall runtime. To address this problem, we adapt techniques developed in high-performance computing—specifically, cache-based optimizations similar to those used in matrix multiplication [19]. With cache-aware execution, the ZQL engine partitions the iterated values in the for loops into blocks of data which fit into the L3 cache. Then, the ZQL engine reorders the order of iteration in the for loops to maximize the time that each block of data remains in the L3 cache. This allows the system to minimize the amount of data the cache needs to eject and thus the amount of data that needs to be copied from main memory to the cache, minimizing the time taken by the p -nodes.

5. zenvisage SYSTEM DESCRIPTION

We now give a brief description of the zenvisage system.

Front-end. The zenvisage front-end is designed as a lightweight web-based client application. It provides a GUI to compose ZQL queries, and displays the resulting visualizations using Vega-lite [25]. A screenshot of zenvisage in action is shown in Figure 3. A list of attributes, divided into qualitative and quantitative, is provided on the left; a table to enter ZQL queries, with auto-completion, is on top, and the resulting visualizations are rendered at the bottom. Users also have the option of hiding the ZQL specification table and instead using a simpler drop-down menu-based interface complemented by a sketching canvas. The sketching canvas allows users to draw their desired trend that can then be used to search for similar trends. The menu-based interface makes it easy for users to perform some of the more common visual exploration queries, such as searching for representative or outlier visualizations. Furthermore, the user may drag-and-drop visualizations from the results onto the sketching canvas, enabling further interaction with the results.

Back-end. The zenvisage front-end issues ZQL queries to the back-end over a REST protocol. The back-end (written in node.js) receives the queries and forwards them to the ZQL engine (written in Java), which is responsible for parsing, compiling, and optimizing the queries as in Section 4. SQL queries issued by the ZQL engine are submitted to one of our back-end databases (which currently include PostgreSQL and Vertica), and the resultant visualization data is returned back to the front-end encoded in JSON.

6. EXPERIMENTAL STUDY

In this section, we evaluate the runtime performance of the ZQL engine. We present the runtimes for executing both synthetic and

realistic ZQL queries and show that we gain speedups of up to $3\times$ with the optimizations from Section 4. We also varied the characteristics of a synthetic ZQL query to observe their impact on our optimizations. Finally, we show that disk I/O was a major bottleneck for the ZQL engine, and if we switched our back-end database to a column-oriented database and cache the dataset in memory, we can achieve interactive run times for datasets as large as 1.5GB.

Setup. All experiments were conducted on a 64-bit Linux server with 8 3.40GHz Intel Xeon E3-1240 4-core processors and 8GB of 1600 MHz DDR3 main memory. We used PostgreSQL with working memory size set to 512 MB and shared buffer size set to 256MB for the majority of the experiments; the last set of experiments demonstrating interactive run times additionally used Vertica Community Edition with a working memory size of 7.5GB.

PostgreSQL Modeling. For modeling the performance on issuing multiple parallel queries with varying number of group by values, we varied the number of parallel queries issued ($\#Q$) from 1 to 100, and the group by values per query ($\#V$) from 10 to 100000, and recorded the response times (T). We observed that the time taken for a batch of queries was practically linearly dependent to both the number of queries as well as the group by values. Fitting a linear equation by performing multiple regression over the observed data, we derived the following cost-model,

$$T(ms) = 908 \times (\#Q) + 1.22 \times \frac{(\#V)}{100} + 1635$$

As per the above model, adding a thread leads to the same rise in response time as increasing the number of group by values by 75000 over the existing threads in the batch. In other words, it is better to merge queries with small number of group by values. Moreover, since there is a fixed cost (1635 ms) associated with every batch of queries, we tried to minimize the number of batches by packing as many queries as possible within the memory constraints.

Optimizations. The four versions of the ZQL engine we use are: (i) NO-OPT: The basic translation from Section 4. (ii) PARALLEL: Concurrent SQL queries for independent nodes from Section 4.2.1. (iii) SPECULATE: Speculates and pre-emptively issues SQL queries from Section 4.2.2. (iv) SMARTFUSE: Query combination with speculation from Section 4.2.3. In our experiments, we consider NO-OPT and the MQO-dependent PARALLEL to be our baselines, while SPECULATE and SMARTFUSE were considered to be completely novel optimizations. For certain experiments later on, we also evaluate the performance of the caching optimizations from Section 4.2.4 on SMARTFUSE.

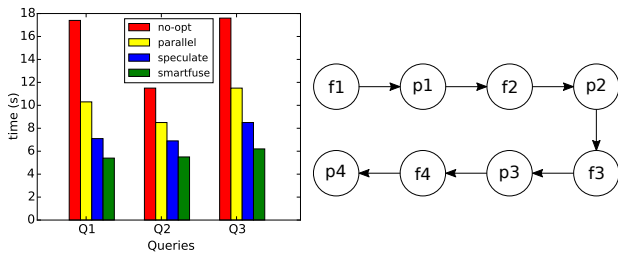


Figure 4: Runtimes for queries on real dataset (left) and single chain synthetic query (right)

6.1 Realistic Queries

For our realistic queries, we used 20M rows of a real 1.5GB airline dataset [1] which contained the details of flights within the USA from 1987-2008, with 11 attributes. On this dataset, we performed 3 realistic ZQL queries inspired by the case studies in our introduction. Descriptions of the queries can be found in Table 36.

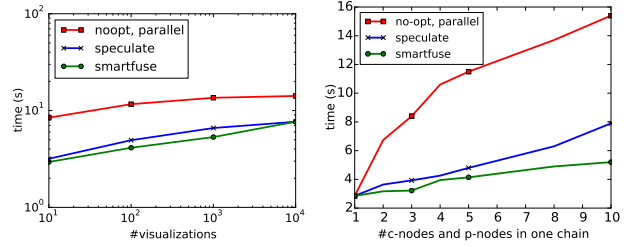


Figure 5: Effect of number of visualizations (left) and length of the chain (right) on the overall runtimes.

Figure 4 (left) depicts the runtime performance of the three realistic ZQL queries, for each of the optimizations. For all queries, each level of optimization provided a substantial speedup in execution time compared to the previous level. *Simply by going from NO-OPT to PARALLEL, we see a 45% reduction in runtime. From PARALLEL to SPECULATE and SPECULATE to SMARTFUSE, we see 15-20% reductions in runtime.* A large reason for why the optimizations were so effective was because ZQL runtimes are heavily dominated by the execution time of the issued SQL queries. In fact, we found that for these three queries, 94-98% of the overall runtime could be contributed to the SQL execution time. We can see from Table 36, SMARTFUSE always managed to lower the number of SQL queries to 1 or 2 after our optimizations, thereby heavily impacting the overall runtime performance of these queries.

6.2 Varying Characteristics of ZQL Queries

We were interested in evaluating the efficacy of our optimizations with respect to four different characteristics of a ZQL query: (i) the number of visualizations explored, (ii) the complexity of the ZQL query, (iii) the level of interconnectivity within the ZQL query, and (iv) the complexity of the processes. To control for all variables except these characteristics, we used a synthetic chain-based ZQL query to conduct these experiments. Every row of the chain-based ZQL query specified a collection of visualizations based on the results of the process from the previous row, and every process was applied on the collection of visualizations from the same row. Therefore, when we created the query plan for this ZQL query, it had the chain-like structure depicted by Figure 4 (right). Using the chain-based ZQL query, we could then (i) vary the number of visualizations explored, (ii) use the length of the chain as a measure of complexity, (iii) introduce additional independent chains to decrease interconnectivity, and (iv) increase the number of loops in a p -node to control the complexity of processes.

To study these characteristics, we used a synthetic dataset with 10M rows and 15 attributes (10 dimensional and 5 measure) with cardinalities of dimensional attributes varying from 10 to 10000. By default, we set the input number of visualizations per chain to be 100, with 10 values for the X attribute, number of c -nodes per chain as 5, the process as T (with a single for loop) with a selectivity of .50, and number of chains as 1.

Impact of number of visualizations. Figure 5 (left) shows the performance of NO-OPT, SPECULATE, and SMARTFUSE on our chain-based ZQL query as we increased the number of visualizations that the query operated on. The number of visualizations was increased by specifying larger collections of Z column values in the first c -node. We chose to omit PARALLEL here since it performs identically to NO-OPT. With the increase in visualizations, the overall response time increased for all versions because the amount of processing per SQL query increased. SMARTFUSE showed better performance than SPECULATE up to 10k visualizations due to reduction in the total number of SQL queries issued. However, at 10k visualization, we reached the threshold of the number of unique

	Query Description	# c -nodes	# p -nodes	# T	# D	# Visualizations	# SQL Queries: NO-OPT	# SQL Queries: SMARTFUSE
1	Plot the related visualizations for airports which have a correlation between arrival delay and traveled distances for flights arriving there.	6	3	670	93,000	18,642	6	1
2	Plot the delays for carriers whose delays have gone up at airports whose average delays have gone down over the years.	5	4	1,000	0	11,608	4	1
3	Plot the delays for the outlier years, outlier airports, and outlier carriers with respect to delays.	12	3	0	94,025	4,358	8	2

Table 36: Realistic queries for the airline dataset with the # of c -nodes, # of p -nodes, # of T functions calculated, # of D functions calculated, # of visualizations explored, # of SQL queries issued with NO-OPT, and # of SQL queries issued with SMARTFUSE per query.

group-by values per combined query (100k for PostgreSQL), so it was less optimal to merge queries. At that point, SMARTFUSE behaved similarly to SPECULATE.

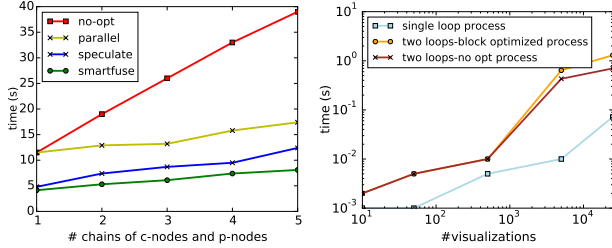


Figure 6: Effect of number of independent chains (left) and the number of loops in a p -node (right) on the overall runtimes.

Impact of the length of the chain. We varied the length of the chain in the query plan (or the number of rows in the ZQL query) to simulate a change in the complexity of the ZQL query and plotted the results in Figure 5 (right). As the number of nodes in the query plan grew, the overall runtimes for the different optimizations also grew. However, while the runtimes for both NO-OPT and SPECULATE grew at least linearly, the runtime for SMARTFUSE grew sublinearly due to its query combining optimization. While the runtime for NO-OPT was much greater than for SPECULATE, since the overall runtime is linearly dependent on the number of SQL queries run in parallel, we see a linear growth for SPECULATE.

Impact of the number of chains. We increased the number of independent chains from 1 to 5 to observe the effect on runtimes of our optimizations; the results are presented in Figure 6 (left). While NO-OPT grew linearly as expected, all PARALLEL, SPECULATE, and SMARTFUSE were close to constant with respect to the number of independent chains. We found that while the overall runtime for concurrent SQL queries did grow linearly with the number of SQL queries issued, they grew much slower compared to issuing those queries sequentially, thus leading to an almost flat line in comparison to NO-OPT.

Impact of process complexity. We increased the complexity of processes by increasing the number of loops in the first p -node from 1 to 2. For the single loop, the p -node filtered based on a positive trend via T , while for the double loop, the p -node found the outlier visualizations. Then, we varied the number of visualizations to see how that affected the overall runtimes. Figure 6 (right) shows the results. For this experiment, we compared regular SMARTFUSE with cache-aware SMARTFUSE to see how much of a cache-aware execution made. We observed that there was not much difference between cache-aware SMARTFUSE and regular SMARTFUSE below 5k visualizations when all data could fit in cache. After 5k visualizations, not all the visualizations could be fit into the cache the same time, and thus the cache-aware execution of the p -node had an improvement of 30-50% as the number of visualizations increased from 5k to 25k. This improvement, while substantial, is only a minor change in the overall runtime.

6.3 Interactivity

The previous figures showed that the overall execution times of ZQL queries took several seconds, even with SMARTFUSE, thus perhaps indicating ZQL is not fit for interactive use with large datasets. However, we found that this was primarily due to the disk-based I/O bottleneck of SQL queries. In Figure 7 (left), we show the SMARTFUSE runtimes of the 3 realistic queries from before on varying size subsets of the airline dataset, with the time that it takes to do a single group-by scan of the dataset. As we can see, the runtimes of the queries and scan time are *virtually the same*, indicating that SMARTFUSE comes very close to the optimal I/O runtime (i.e., a “fundamental limit” for the system).

To further test our hypothesis, we ran our ZQL engine with Vertica with a large working memory size to cache the data in memory to avoid expensive disk I/O. The results, presented in Figure 7 (right), showed that there was a 50 \times speedup in using Vertica over PostgreSQL with these settings. Even with a large dataset of 1.5GB, we were able to achieve sub-second response times for many queries. Furthermore, for the dataset with 120M records (11GB, so only 70% could be cached), we were able to reduce the overall response times from 100s of seconds to less than 10 seconds. Thus, once again zenvisage returns results in a small multiple of the time it takes to execute a single group-by query.

Overall, SMARTFUSE will be interactive on moderate sized datasets on PostgreSQL, or on large datasets that can be cached in memory and operated on using a columnar database—which is standard practice adopted by visual analytics tools [44]. Improving on interactivity is impossible due to fundamental limits to the system; in the future, we plan to explore returning approximate answers using samples, since even reading the entire dataset is prohibitive.

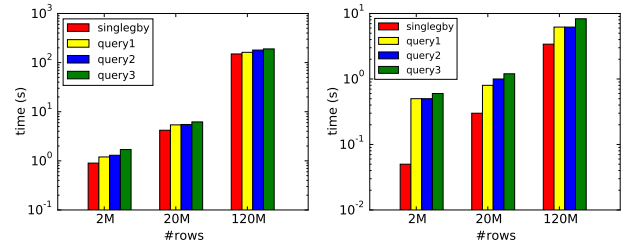


Figure 7: SMARTFUSE on PostgreSQL (left) and Vertica (right)

7. USER STUDY

We conducted a user study to evaluate the utility of zenvisage for data exploration versus two types of systems—first, visualization tools, similar to Tableau, and second, general database and data mining tools, which also support interactive analytics to a certain extent. In preparation for the user study, we conducted interviews with data analysts to identify the typical exploration tasks and tools used in their present workflow. Using these interviews, we identified a set of tasks to be used in the user study for zenvisage. We describe these interviews first, followed by the user study details.

7.1 Analyst Interviews and Task Selection

We hired seven data analysts via Upwork [4], a freelancing platform—we found these analysts by searching for freelancers who

had the keywords *analyst* or *tableau* in their profile. We conducted one hour interviews with them to understand how they perform data exploration tasks. The interviewees had 3—10 years of prior experience, and told about every step of their workflow; from receiving the dataset to presenting the analysis to clients. The rough workflow of all interviewees identified was the following: first, data cleaning is performed; subsequently, the analysts perform data exploration; then, the analysts develop presentations using their findings. We then drilled down onto the data exploration step.

We first asked the analysts what types of tools they use for data exploration. Analysts reported nine different tools—the most popular ones included Excel (5), Tableau (3), and SPSS (2). The rest of the tools were reported by just one analyst: Python, SQL, Alteryx, Microsoft Visio, Microsoft BI, SAS. Perhaps not surprisingly, analysts use both visualization tools (Tableau, Excel, BI), programming languages (Python), statistical tools (SAS, SPSS), and relational databases (SQL) for data exploration.

Then, to identify the common tasks used in data exploration, we used a taxonomy of abstract exploration tasks proposed by Amar et al. [9]. Amar et al. developed their taxonomy through summarizing the analytical questions that arose during the analysis of five different datasets, independent of the capabilities of existing tools or interfaces. The exploration tasks in Amar et al. include: filtering (f), sorting (s), determining range (r), characterizing distribution (d), finding anomalies (a), clustering (c), correlating attributes (co), retrieving value (v), computing derived value (dv), and finding extrema (e). When we asked the data analysts which tasks they use in their workflow, the responses were consistent in that all of them use all of these tasks, except for three exceptions—c, reported by four participants, and e, d, reported by six participants.

Given these insights, we selected a small number of appropriate tasks for our user study encompassing eight of the ten exploration tasks described above: f, s, r, d, a, c, co, v. The other two—dv and e—finding derived values and computing extrema, are important tasks in data analysis, but existing tools (e.g., Excel) already provide adequate capabilities for these tasks, and we did not expect *zenvisage* to provide additional benefits.

7.2 User Study Methodology

The goal of our user study was to evaluate *zenvisage* with other tools, on its ability to effectively support data exploration.

Participants. We recruited 12 graduate students as participants with varying degrees of expertise in data analytics. Table 37 depicts the participants’ experience with different categories of tools.

Tools	Count
Excel, Google spreadsheet, Google Charts	8
Tableau	4
SQL, Databases	6
Matlab,R,Python,Java	8
Data mining tools such as weka, JNP	2
Other tools like D3	2

Table 37: Participants’ prior experience with data analytic tools

Baselines. For the purposes of our study, we explicitly wanted to do a head-to-head qualitative and quantitative comparison with visual analytics tools, and thus we developed a baseline tool to compare *zenvisage* against directly. Further, via qualitative interviews, we compared *zenvisage* versus against other types of tools, such as databases, data mining, and programming tools. Our baseline tool was developed by replicating the visualization selection capabilities of visual analytics tools with a styling scheme identical to *zenvisage* to control for external factors. The tool allowed users to specify the X-axis, Y-axis, dimensions, and filters. The tool would then populate all visualizations meeting the specifications.

Comparison Points. There are no tools that offer the same functionalities as *zenvisage*. Visual analytics tools do not offer the ability to search for specific patterns, or issue complex visual exploration queries; data mining toolkits do not offer the ability to search for visual patterns and are instead tailored for general machine learning and prediction. Since visual analytics tools are closer in spirit and functionality to *zenvisage*, we decided to implement a visual analytics tool as our baseline. Thus, our baseline tool replicated the basic query specification and output visualization capabilities of existing tools such as Tableau. We augmented the baseline tool with the ability to specify an arbitrary number of filters, allowing users to use filters to drill-down on specific visualizations. This baseline visualization tool was implemented with a styling scheme similar to *zenvisage* to control for external factors. As depicted in Figure 8, the baseline allowed users to visualize data by allowing them to specify the x-axis, y-axis, category, and filters. The baseline tool would populate all the visualizations, which fit the user specifications, using an alpha-numeric sort order. In addition to task-based comparisons with this baseline, we also explicitly asked participants to compare *zenvisage* with existing data mining and visual analytics tools that they use in their workflow.

Dataset. We used a housing dataset from Zillow.com [5], consisting of housing sales data for different cities, counties, and states from 2004-15, with over 245K rows, and 15 attributes. We selected this dataset since participants could relate to the dataset and understand the usefulness of the tasks.

Tasks. We designed the user study tasks with the case studies from Section 1 in mind, and translated them into the housing dataset. Further, we ensured that these tasks together evaluate eight of the exploration tasks described above—f, s, r, d, a, c, co, and v. One task used in the user study is as follows: “Find three cities in the state of NY where the Sold Price vs Year trend is very different from the state’s overall trend.” This query required the participants to first retrieve the trend of NY (v) and characterize its distribution (d), then separately filter to retrieve the cities of NY (f), compare the values to find a negative correlation (co), sort the results (s), and report the top three cities on the list.

Study Protocol. The user study was conducted using a within-subjects study design [11], forming three phases. First, participants described their previous experience with data analytics tools. Next, participants performed exploration tasks using *zenvisage* (Tool A) and the baseline tool (Tool B), with the orders randomized to reduce order effects. Participants were provided a 15-minute tutorial-cum-practice session per tool to get familiarized before performing the tasks. Finally, participants completed a survey that both measured their satisfaction levels and preferences, along with open-ended questions on the strengths and weaknesses of *zenvisage* and the baseline, when compared to other analytics tools they may have used. The average study session lasted for 75 minutes on average. Participants were paid ten dollars per hour for their participation. After the study, we reached out to participants with backgrounds in data mining and programming, and asked if they could complete a follow-up interview where they use their favorite analytics tool for performing one of the tasks, via email.

Metrics. Using data that we recorded, we collected the following metrics: completion time, accuracy, and the usability ratings and satisfaction level from the survey results. In addition, we also explicitly asked participants to compare *zenvisage* with tools that they use in their workflow. For comparisons between *zenvisage* and general database and data mining tools via follow-up interviews, we used the number of lines of code to evaluate the differences.

Ground Truth. Two expert data analysts prepared the ground truth

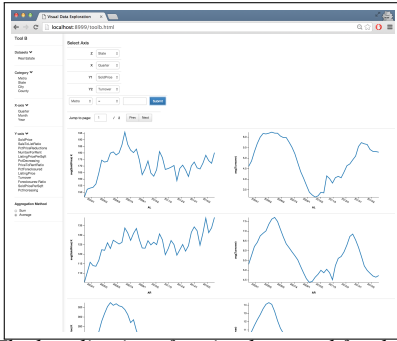


Figure 8: The baseline interface implemented for the user study.

for each the tasks in the form of ranked answers, along with score cut-offs on a 0 to 5 scale (5 highest). Their inter-rater agreement, measured using Kendall’s Tau coefficient, was 0.854. We took the average of the two scores to rate the participants’ answers.

7.3 Key Findings

Three key findings emerged from the study and are described below. We use μ , σ , χ^2 to denote average, standard deviation, and Chi-square test scores, respectively.

Finding 1: *zenvisage enables faster and more accurate exploration than existing visualization tools.* Since all of our tasks involved generating multiple visualizations and comparing them to find desired ones, participants were not only able to complete the tasks *faster*— $\mu=115s$, $\sigma=51.6$ for zenvisage vs. $\mu=172.5s$, $\sigma=50.5$ for the baseline—but also more *accurately*— $\mu=96.3\%$, $\sigma=5.82$ for zenvisage vs. $\mu=69.9\%$, $\sigma=13.3$ for the baseline. A one-way between-subjects ANOVA, followed by a post-hoc Tukey’s test [45], we found that zenvisage had statistically significant faster task completion times compared to the baseline interface, with p value of 0.0069. The baseline requires considerable manual exploration to complete the same task, explaining the high task completion times; in addition, participants frequently compromised by selecting sub-optimal answers before browsing the entire list of results for better ones, explaining the low accuracy. On the other hand, zenvisage is able to automate the task of finding desired visualizations, considerably reducing manual effort. Also of note is the fact that the accuracy with zenvisage is close to 100%—indicating that a short 15 minute tutorial on ZQL was enough to equip users with the knowledge they needed to address the tasks—and that too, within 2 minutes (on average).

When asked about using zenvisage vs. the baseline in their current workflow, 9 of the 12 participants stated that they would use zenvisage in their workflow, whereas only two participants stated that they would use our baseline tool ($\chi^2 = 8.22$, $p < 0.01$). When the participants were asked how, one participant provided a specific scenario: “If I am doing my social science study, and I want to see some specific behavior among users, then I can use tool A [zenvisage] since I can find the trend I am looking for and easily see what users fit into the pattern.” (P7). In response to the survey question “I found the tool to be effective in visualizing the data I want to see”, the participants rated zenvisage higher ($\mu=4.27$, $\sigma=0.452$) than the baseline ($\mu=2.67$, $\sigma=0.890$) on a five-point Likert scale. A participant experienced in Tableau commented: “In Tableau, there is no pattern searching. If I see some pattern in Tableau, such as a decreasing pattern, and I want to see if any other variable is decreasing in that month, I have to go one by one to find this trend. But here I can find this through the query table.” (P10).

Finding 2: *zenvisage complements existing database and data mining systems, and programming languages.* When explicitly asking participants about comparing zenvisage with the tools they

use on a regular basis for data analysis, all participants acknowledged that zenvisage adds value in data exploration not encompassed by their tools. ZQL augmented with inputs from the sketching canvas proved to be extremely effective. For example P8 stated: “you can just [edit] and draw to find out similar patterns. You’ll need to do a lot more through Matlab to do the same thing.” Another experienced participant mentioned the benefits of not needing to know much programming to accomplish certain tasks: “The obvious good thing is that you can do complicated queries, and you don’t have to write SQL queries... I can imagine a non-cs student [doing] this.” (P9). When asked about the specific tools they would use to solve the user study tasks, all participants reported a programming language like Matlab or Python. This is despite half of the participants reporting using a relational database regularly, and a smaller number of participants (2) reporting using a data mining tool regularly. Additionally, multiple participants even with extensive programming experience reported that zenvisage would take less time and fewer lines of code for certain data exploration tasks. (Indeed, we found that all participants were able to complete the user study tasks in under 2 minutes.) In follow-up email interviews, we asked a few participants to respond with code from their favorite data analytics tool for the user study tasks. Two participants responded — one with Matlab code, one with Python code. Both these code snippets were much longer than ZQL: as a concrete example, the participant accomplished the same task with 38 lines of Python code compared to 4 lines of ZQL. While comparing code may not be fair, the roughly order of magnitude difference demonstrates the power of zenvisage over existing systems.

Finding 3: *zenvisage can be improved.* While the participants looked forward to using custom query builder in their own workflow, a few of them were interested in directly exposing the commonly-used trends/patterns such as outliers, through the drag and drop interface. Some were interested in knowing how they could integrate custom functional primitives (we could not cover it in the tutorial due to time constraints). In order to improve the user experience, participants suggested adding instructions and guidance for new users as part of the interface. Participants also commented on the unrefined look and feel of the tool, as well as the lack of a diverse set of usability related features, such as bookmarking and search history, that are offered in existing systems.

8. RELATED WORK

We now discuss related prior work in a number of areas. We begin with analytics tools — visualization tools, statistical packages and programming libraries, and relational databases. Then, we talk about other tools that overlap somewhat with zenvisage.

Visual Analytics Tools. Visualization tools, such as ShowMe, Spotfire, and Tableau [43, 34, 8], along with similar tools from the database community [18, 31, 32, 26] have recently gained in popularity, catering to data scientists who lack programming skills. Using these tools, these scientists can select and view one visualization at a time. However, these tools do not operate on collections of visualizations at a time—and thus they are much less powerful and the optimization challenges are minimal. zenvisage, on the other hand, supports queries over collections of visualizations, returning results not much slower than the time to execute a single query (See Section 6). Since these systems operate one visualization at a time, users are also not able to directly identify desired patterns or needs.

Statistical Packages and Programming Libraries: Statistical tools (e.g., KNIME, RapidMiner, SAS, SPSS) support the easy application of data mining and statistical primitives—including prediction algorithms and statistical tests. While these tools support the selection of a prediction algorithm (e.g., decision trees) to apply, and

```

with ranking as (
with distances as (
with distance_product_year as (
with aggregate_product_year as (
  select product, year, avg(profit) as avg_profit
  from table group by product, year )
  select s.product as source, d.product as destination, s.year,
  power(s.avg_profit - d.avg_profit,2) as distance_year
  from aggregate_product_year s, aggregate_product_year d
  where s.product!=d.product and s.year=d.year )
  select source, destination, sum(distance_year) as distance
  from distance_product_year groupby source, destination )
  select source, destination, distance,
  rank() over (partition by source order by distance asc)
  rank from distances )
  select source, destination, distance
  from ranking where rank < 10;

```

Table 38: Verbose SQL query

the appropriate parameters, they offer no querying capabilities, and as a result do not need extensive optimization. As a result, these tools cannot support user needs like those describe in the examples in the introduction. Similarly, programming libraries such as Weka [22] and Scikit-learn [37] embed machine learning within programs. However, manually translating the user desired patterns into code that uses these libraries will require substantial user effort and hand-optimization. In addition, writing new code and hand-optimization will need to be performed every time the exploration needs change. Additionally, for both statistical tools and programming libraries, there is a need for programming ability and understanding of machine learning and statistics to be useful—something we cannot expect all data scientists to possess.

Relational Databases. Relational databases can certainly support interactive analytics via SQL. In *zenvisage*, we use relational databases as a backend computational component, augmented with an engine that uses SMARTFUSE to optimize accesses to the database, along with efficient processing code. Thus, one can certainly express some ZQL queries by writing multiple SQL queries (via procedural SQL), using complex constructs only found in some databases, such as common table expressions (CTE) and window functions. As we saw in Section 7, these SQL queries are very cumbersome to write, and are not known to most users of databases—during our user study, we found that all participants who had experience with SQL were not aware of these constructs; in fact, they responded that they did not know of any way of issuing ZQL queries in SQL, preferring instead to express these needs in Python. In Table 38, we list the verbose SQL query that computes the following: for each product, find 10 other products that have most similar profit over year trends. The equivalent ZQL query takes two lines. And we were able to write the SQL query only because the function D is Euclidean distance: for other functions, we are unable to come up with appropriate SQL rewritings. On the other hand, for ZQL, it is effortless to change the function by selecting it from a drop-down menu. Beyond being cumbersome to write, the constructs required lead to severe performance penalties on most databases—for instance, PostgreSQL materializes intermediate results when executing queries with CTEs. To illustrate, we took the SQL query in Table 38, and compared its execution with the execution of the equivalent ZQL. As depicted in Figure 9, the time taken by PostgreSQL increases sharply as the number of visualizations increases, taking up to 10X more time as compared to ZQL query executor. This indicates that *zenvisage* is still important even for the restricted cases where we are able to correctly write the queries in SQL.

OLAP Browsing. There has been some work on interactive browsing of data cubes [39, 40]. The work focuses on suggestions for raw aggregates to examine that are informative given past browsing, or those that show a generalization or explanation of a specific cell—an easier problem meriting simpler techniques—not addressing the

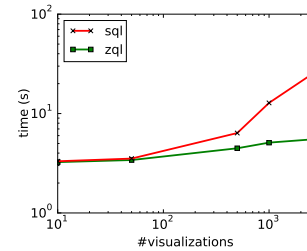


Figure 9: ZQL vs SQL: we want to find top 10 similar products for every product on varying the number of products from 10—5000.

full exploration capabilities provided by ZQL.

Data Mining Languages: There has been some limited work in data mining query languages, all from the early 90s, on association rule mining (DMQL [21], MSQL [24]), or on storing and retrieving models on data (OLE DB [35]), as opposed to a general-purpose visual data exploration language aimed at identifying visual trends.

Visualization Suggestion Tools: There has been some recent work on building systems that suggest visualizations. Voyager [25] recommends visualizations based on aesthetic properties of the visualizations, as opposed to queries. SeeDB [46] recommends visualizations that best display the difference between two sets of data. SeeDB and Voyager can be seen to be special cases of *zenvisage*. The optimization techniques outlined are a substantial generalization of the techniques described in SeeDB; while the techniques in SeeDB are special-cased to one setting (a simple comparison), here, our goal is to support and optimize all ZQL queries.

Multi-Query Optimization: There has been a lot of work on Multi-Query Optimization (MQO), both classic, e.g., [41, 42, 38], and recent work, e.g., [16, 23, 27, 17]. Overall, the approach adopted is to batch queries, decompose into operators, and build “meta”-query plans that process multiple queries at once, with sharing at the level of scans, or at the level of higher level operators (either via simultaneous pipelining or a true global query plan [23]). Unlike these techniques which require significant modifications to the underlying database engine—indeed, some of these systems do not even provide full cost-based optimization and only support hand-tuned plans [16], in this paper, we adopted two syntactic rewriting techniques that operate outside of any relational database as a backend without requiring any modification, and can thus seamlessly leverage improvements to the database. Our third optimization is tailored to the ZQL setting and does not apply more broadly.

Anomaly Discovery: Anomaly detection is a well-studied topic [14, 7, 36]. Our goal in that *zenvisage* is expected to be interactive, especially on large datasets; most work in anomaly detection focuses on accuracy at the cost of latency and is typically a batch operation. In our case, since interactivity is of the essence, and requests can come at any time, the emphasis is on scalable on-the-fly data processing aspects.

Time Series Similarity and Indexing: There has been some work on indexing of time series data, e.g., [30, 20, 13, 29, 12, 15, 28]; for the attributes that are queried frequently, we plan to reuse these techniques for similarity search. For other attributes, indexing and maintaining all trends is impossible, since the number of trends grows exponentially with the number of indexed attributes.

9. CONCLUSION

We propose *zenvisage*, a visual analytics tool for effortlessly identifying desired visual patterns from large datasets. We described the formal syntax of the query language ZQL, motivated by many real-world use-cases, and demonstrated that ZQL is visual exploration algebra-complete. *zenvisage* enables users to effec-

tively and accurately perform visual exploration tasks, as shown by our user study, and complements other tools. In addition, we show that our optimizations for ZQL execution lead to considerable improvements over leveraging the parallelism inherent in databases. Our work is a promising first step towards substantially simplifying and improving the process of interactive data exploration for novice and expert analysts alike.

10. REFERENCES

- [1] Airline dataset (<http://stat-computing.org/dataexpo/2009/the-data.html>). [Online; accessed 30-Oct-2015].
- [2] Spotfire, <http://spotfire.com>. [Online; accessed 17-Aug-2015].
- [3] Tableau public (www.tableaupublic.com/). [Online; accessed 3-March-2014].
- [4] Upwork (<https://www.upwork.com/>). [Online; accessed 3-August-2016].
- [5] Zillow real estate data (<http://www.zillow.com/research/data/>). [Online; accessed 1-Feb-2016].
- [6] Tableau q2 earnings: Impressive growth in customer base and revenues. <http://www.forbes.com/sites/greatspeculations/2015/07/31/tableau-q2-earnings-impressive-growth-in-customer-base-and-revenues>.
- [7] M. Agyemang, K. Barker, and R. Alhajj. A comprehensive survey of numeric and symbolic outlier mining techniques. *Intell. Data Anal.*, 10(6):521–538, Dec. 2006.
- [8] C. Ahlberg. Spotfire: An information exploration environment. *SIGMOD Rec.*, 25(4):25–29, Dec. 1996.
- [9] R. Amar, J. Eagan, and J. Stasko. Low-level components of analytic activity in information visualization. In *INFOVIS.*, pages 111–117. IEEE, 2005.
- [10] M. R. Anderberg. *Cluster analysis for applications: probability and mathematical statistics: a series of monographs and textbooks*, volume 19. Academic press, 2014.
- [11] K. S. Bordens and B. B. Abbott. *Research design and methods: A process approach*. McGraw-Hill, 2002.
- [12] K. Chakrabarti, E. Keogh, S. Mehrotra, and M. Pazzani. Locally adaptive dimensionality reduction for indexing large time series databases. *ACM Trans. Database Syst.*, 27(2):188–228, June 2002.
- [13] K.-P. Chan and A.-C. Fu. Efficient time series matching by wavelets. In *Data Engineering, 1999. Proceedings., 15th International Conference on*, pages 126–133, Mar 1999.
- [14] V. Chandola, A. Banerjee, and V. Kumar. Anomaly detection: A survey. *ACM Comput. Surv.*, 41(3):15:1–15:58, July 2009.
- [15] C. Faloutsos, M. Ranganathan, and Y. Manolopoulos. Fast subsequence matching in time-series databases. *SIGMOD Rec.*, 23(2):419–429, May 1994.
- [16] G. Giannikis et al. Workload optimization using shareddb. In *SIGMOD*, pages 1045–1048. ACM, 2013.
- [17] G. Giannikis et al. Shared workload optimization. *Proceedings of the VLDB Endowment*, 7(6):429–440, 2014.
- [18] H. Gonzalez et al. Google fusion tables: web-centered data management and collaboration. In *SIGMOD Conference*, pages 1061–1066, 2010.
- [19] K. Goto and R. A. Geijn. Anatomy of high-performance matrix multiplication. *ACM Transactions on Mathematical Software (TOMS)*, 34(3):12, 2008.
- [20] D. Gunopulos and G. Das. Time series similarity measures and time series indexing (abstract only). *SIGMOD Rec.*, 30(2):624–, May 2001.
- [21] J. Han et al. Dmql: A data mining query language for relational databases. In *Proc. 1996 SIGMOD*, volume 96, pages 27–34, 1996.
- [22] G. Holmes, A. Donkin, and I. H. Witten. Weka: A machine learning workbench. In *Conf. on Intelligent Information Systems '94*, pages 357–361. IEEE, 1994.
- [23] I. Psaroudakis et al. Sharing data and work across concurrent analytical queries. *VLDB*, 6(9):637–648, 2013.
- [24] T. Imielinski and A. Virmani. A query language for database mining. *Data Mining and Knowledge Discovery*, 3(4):373–408, 2000.
- [25] K. Wongsuphasawat et al. Voyager: Exploratory analysis via faceted browsing of visualization recommendations. *IEEE TVCG*, 2015.
- [26] S. Kandel et al. Profiler: integrated statistical analysis and visualization for data quality assessment. In *AVI*, pages 547–554, 2012.
- [27] A. Kementsietsidis et al. Scalable multi-query optimization for exploratory queries over federated scientific databases. *PVLDB*, 1(1):16–27, 2008.
- [28] E. Keogh. A decade of progress in indexing and mining large time series databases. In *Proceedings of the 32Nd International Conference on Very Large Data Bases, VLDB '06*, pages 1268–1268. VLDB Endowment, 2006.
- [29] E. Keogh, K. Chakrabarti, M. Pazzani, and S. Mehrotra. Dimensionality reduction for fast similarity search in large time series databases. *Knowledge and Information Systems*, 3(3):263–286, 2001.
- [30] E. Keogh, K. Chakrabarti, M. Pazzani, and S. Mehrotra. Locally adaptive dimensionality reduction for indexing large time series databases. *SIGMOD Rec.*, 30(2):151–162, May 2001.
- [31] A. Key, B. Howe, D. Perry, and C. Aragon. Vizdeck: Self-organizing dashboards for visual analytics. *SIGMOD '12*, pages 681–684, 2012.
- [32] M. Livny et al. Devise: Integrated querying and visualization of large datasets. In *SIGMOD Conference*, pages 301–312, 1997.
- [33] J. Mackinlay. Automating the design of graphical presentations of relational information. *ACM Trans. Graph.*, 5(2):110–141, Apr. 1986.
- [34] J. D. Mackinlay et al. Show me: Automatic presentation for visual analysis. *IEEE Trans. Vis. Comput. Graph.*, 13(6):1137–1144, 2007.
- [35] A. Netz et al. Integrating data mining with sql databases: Ole db for data mining. In *ICDE'01*, pages 379–387. IEEE, 2001.
- [36] A. Patcha and J.-M. Park. An overview of anomaly detection techniques: Existing solutions and latest technological trends. *Comput. Netw.*, 51(12):3448–3470, Aug. 2007.
- [37] Pedregosa et al. Scikit-learn: Machine learning in python. *The Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [38] P. Roy et al. Efficient and extensible algorithms for multi query optimization. In *ACM SIGMOD Record*, volume 29, pages 249–260. ACM, 2000.
- [39] S. Sarawagi. Explaining differences in multidimensional aggregates. In *VLDB*, pages 42–53, 1999.

- [40] G. Sathe and S. Sarawagi. Intelligent rollups in multidimensional olap data. In *VLDB*, pages 531–540, 2001.
- [41] T. K. Sellis. Multiple-query optimization. *ACM TODS*, 13(1):23–52, 1988.
- [42] K. Shim et al. Improvements on a heuristic algorithm for multiple-query optimization. *Data & Knowledge Engineering*, 12(2):197–222, 1994.
- [43] C. Stolte et al. Polaris: a system for query, analysis, and visualization of multidimensional databases. *Commun. ACM*, 51(11):75–84, 2008.
- [44] P. Terlecki et al. On improving user response times in tableau. In *SIGMOD*, pages 1695–1706. ACM, 2015.
- [45] J. W. Tukey. Comparing individual means in the analysis of variance. *Biometrics*, pages 99–114, 1949.
- [46] M. Vartak et al. Seedb: Efficient data-driven visualization recommendations to support visual analytics. *VLDB*, 8(13), Sept. 2015.
- [47] H. Wickham. ggplot: An implementation of the grammar of graphics. *R package version 0.4. 0*, 2006.
- [48] L. Wilkinson. *The grammar of graphics*. Springer Science & Business Media, 2006.
- [49] M. M. Zloof. Query-by-example: A data base language. *IBM Systems Journal*, 16(4):324–343, 1977.

APPENDIX

Here, we provide additional details on X, Y, Z, and Process columns (Appendix A), four real-world complex examples (Appendix B), and provide a participant’s python code implementation (Appendix C) for a ZQL task to support the finding 2 in the user-study (Section 7).

A. QUERY LANGUAGE FORMALIZATION: ADDITIONAL DETAILS

In this section, we present some additional details on our formalization that we did not cover in the main body of the paper.

A.1 Additional Details on X and Y Columns

In addition to using a single attribute for an X or Y column, ZQL also allows the use of the Polaris table algebra [3] in the X and Y columns to to arbitrarily compose multiple attributes into a single attribute; all three operators are supported: +, ×, /. Table 39 shows an example of using the + operator to visualize both profits and sales on a single y-axis. Note that this is different from the example given in Table 4, which generates two visualizations, as opposed to a single visualization. An example using both table algebra and sets is given in Table 40, which uses the × operator to return the set of visualizations which measures the sales for the Cartesian product of ‘product’ and one of ‘county’, ‘state’, and ‘country’.

Name	X	Y	Z
*f1	‘product’	‘profit’ + ‘sales’	location.‘US’

Table 39: A ZQL query for a visualization which depicts both profits and sales on the y-axis for products in the US.

Name	X	Y
*f1	‘product’ × {‘county’, ‘state’, ‘country’}	‘sales’

Table 40: A ZQL query for the set of visualizations which measures the sales for one of (‘product’, ‘county’), (‘product’, ‘state’), and (‘product’, ‘country’).

A.2 Additional Details on the Z Column

ZQL also allows the iteration over attributes in the Z column as shown in Table 41. The result of this query is the set of all sales

over time visualizations for every possible slice in every dimension except ‘time’ and ‘sales’. Since both attribute and attribute value can vary in this case, we need separate variables for each component, and the full attribute name, value pair (z1.v1) must be specified. Note that the resulting set of visualizations comes from the Cartesian product of possible attribute and attribute value pairs. The first * symbol refers to all possible attributes, while the second * symbol refers to all possible attribute values given an attribute. If the user wishes to specify specific subsets of attribute values for attributes, she must name them individually.

Name	X	Y	Z
f1	‘year’	‘sales’	z1.v1 <- (\ {‘year’, ‘sales’}).*

Table 41: A ZQL query which returns the set of sales over year visualizations for each attribute that is not time or sales.

A.3 Additional Details on the Name Column

For deriving visualization collections based on other collections in the Name column, in addition to + operation, ZQL supports following operations. (i) f3 <- f1-f2: where f3 refers to the list of visualizations in f1 with the exception of the visualizations which appear in f2, (ii) f2 <- [f1[i]]: where f2 refers to the ith visualization in f1, (iii) f2 <- f1[i:j]: where f2 refers to the list of visualizations starting from ith visualization to the jth visualization in f1, (iv) f2 <- f1.uni: where f2 refers to the set of visualizations derived from f1 by removing duplicate visualizations (only the first appearance of each visualization is kept), and (v) f3 IN f1~f2: where f3 refers to the list of visualizations in f1 which also appear in f2; an “intersection” between lists of visualizations in some sense. These operations are useful if the user wants to throw away some visualizations, or create a new larger set of visualizations from smaller sets of visualizations.

After a visualization collection has been derived using the Name column, the user may also define axis variables in the X, Y, and Z columns using the special _ symbol to bind to the derived collection. For example in Table 43, v2 is defined to be the iterator which iterates over the set of product values which appear in derived collection f3; in this case, v2 iterates over all possible products. y1 is defined to be the iterator over all the values in the Y column of f3. Although in the case of Table 43, the only value y1 takes on is ‘sales’, y1 and v2 are considered to be declared together, so the iterations for y1, v2 will look like: [(‘sales’, ‘chair’), (‘sales’, ‘table’), ...]. Also in this case, the variable y1 is not used, however, there may be other cases where it may be useful to iterate over multiple axis variables. The defined axis variables can then be used to create other visualization collections or within the Process column as shown in the 4th row of Table 43.

Finally, visualization collections may also be ordered based on the values of axis variables: f2 <- f1.order. Here, f1 is ordered based on the axis variables which appear together with the -> symbol. Table 42 shows an example of such an operator in use.

A.4 Additional Details on the Process Column

Although visualization collections typically outnumber processes, there may occur cases in which the user would like to specify multiple processes in one line. To accomplish this, the user simply delimits each process with a comma and surrounds each declaration of variables with parentheses. Table 44 gives an example of this.

B. ADDITIONAL COMPLETE EXAMPLES

To demonstrate the full expressive power of ZQL, we present four realistic, complex example queries. We show that even with

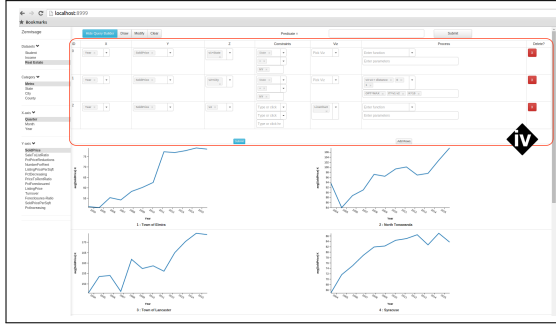


Figure 10: Custom query builder with bar charts



Figure 11: Custom query builder with scatterplots

Name	X	Y	Z	Process
f1	'year'	'sales'	v1 <- 'product'.*	u1 <- argmin _{v1} [k = ∞]T(f1)
*f2=f1.order			u1 ->	

Table 42: A ZQL query which reorders the set of sales over years visualizations for different products based on increasing overall trend.

Name	X	Y	Z	Process
f1	'year'	'sales'	v1 <- 'product'.(* - 'stapler')	v3 <- argmax _{v2} [k = 10]D(f3, f4)
f2	'year'	'sales'	'stapler'	
f3=f1+f2		y1 <- _	v2 <- 'product'._	
f4	'year'	'profit'	v2	
*f5	'year'	'sales'	v3	

Table 43: A ZQL query which returns the sales over years visualizations for the top 10 products which have the most different sales over years visualizations and profit over years visualizations.

Name	X	Y	Z	Process
f1	-	-	-	(v2 <- argmax _{v1} [k = 1]D(f1, f2)), (v3 <- argmin _{v1} [k = 1]D(f1, f2))
f2	'year'	'sales'	v1 <- 'product'.*	
*f3	'year'	'sales'	v2	
*f4	'year'	'sales'	v3	

Table 44: A ZQL query which returns the sales over years visualizations for the product that looks most similar to the user-drawn input and most dissimilar to the user-drawn input.

complicated scenarios, the user is able to capture the insights she wants with a few meaningful lines of ZQL.

Query 1. The stapler has been one of the most profitable products in the last years for GlobalMart. The Vice President is interested in learning about other products which have had similar profit trends. She wishes to see some representative sales over the years visualizations for these products.

Table 45 shows what the query that the Vice President would write for this scenario. She first filters down to the top 100 products which have the most similar to profit over year visualizations to that of the stapler's using the *argmin* in the second row. Then, from the resulting set of products, v2, she picks the 10 most representative set of sales over visualizations using *R*, and displays those visualizations in the next line with f4. Although the Vice President does not specify the exact distance metric for *D* or specify the exact algorithm for *R*, she knows *zenvisage* will select the most reasonable default based on the data.

Query 2. The Vice President, to her surprise, sees that there a few products whose sales has gone up over the last year, yet their profit has declined. She also notices some product's sales have gone down, yet their profit has increased. To investigate, the Vice President would like to know about the top 10 products who have the most discrepancy in their sales and profit trends, and she would like to visualize those trends.

This scenario can be addressed with the query in Table 46. The Vice President names the set of visualizations for profit over month f1 and the sales over month visualizations f2. She then compares the visualizations in the two set using the *argmax* and retrieves the top 10 products whose visualizations are the most different. For these visualizations, she plots both the sales and profit over months; y1 <- {'sales', 'profit'} is a shortcut to avoid having to separates rows for sales and profit. Note that the Vice President was careful to constrain ZQL to only look at the data from 2015.

Query 3. The Vice President would like to know more about the

differences between a product whose sales numbers do not change over the year and a product that has the largest growth in the number of sales. To address this question, she writes the query in Table 47. The first *R* function call returns the one product whose sales over year visualization is most representative for all products; in other words, v2 is set to the product that has the most average number of sales. The task in the second row selects the product v3 which has the greatest upward trending slope *T* for sales. Finally, the Vice President tries to find the y-axes which distinguish the two products the most, and visualizes them. Although we know v2 and v3 only contain one value, they are still sets, so *argmax* must iterate over them and output corresponding values v4 and v5.

Query 4: Finally, the Vice President wants to see a pair of dimensions whose correlation pattern (depicted as a scatterplot) is the most unusual, compared to correlation patterns of other pairs of attributes. To address this question, she writes the query in Table 48. She keeps the Z column empty as she does not want to slice the data. Both X and Y refer to a set M consisting of all the attributes in the dataset she wishes to explore. The task in the second row selects the X and Y attributes whose sum of distances from other visualizations (generated by considering all pairs of attributes) is the maximum.

C. USER STUDY: ADDITIONAL DETAILS ON FINDING 2

Participant P6's python code implementation for a task where we want to find a pair of X and Y axes where two states 'CA' and 'NY' differ the most:

```
[language=python]
import pandas
import numpy as np
def generate_maps(date_list, d, Y, Z):
    d = d[d['State']==Z][np.append(date_list, Y)]
    maps = {}
```

Name	X	Y	Z	Viz	Process
f1	'year'	'profit'	'product'. 'stapler'	bar.(y=agg('sum'))	v2 <- argmin _{v1} [k = 100]D(f1, f2) v3 <- R(10, v2, f3)
f2	'year'	'profit'	v1 <- 'product'.(* \ {'stapler'})	bar.(y=agg('sum'))	
f3	'year'	'sales'	v2	bar.(y=agg('sum'))	
*f4	'year'	'sales'	v3	bar.(y=agg('sum'))	

Table 45: The ZQL query which returns 10 most representative sales over year visualizations for products which have similar profit over year visualizations to that of the stapler's.

Name	X	Y	Z	Z2	Viz	Process
f1	'month'	'profit'	v1 <- 'product'.*	'year'.2015	bar.(y=agg('sum'))	v2 <- argmax _{v1} [k = 10]D(f1, f2)
f2	'month'	'sales'	v1	'year'.2015	bar.(y=agg('sum'))	
*f3	'month'	y1 <- {'sales', 'profit'}	v2	'year'.2015	bar.(y=agg('sum'))	

Table 46: The ZQL query which returns the sales over month and profit over month visualizations for 2015 for the top 10 products which have the biggest discrepancies in their sales and profit trends.

```

for id, item in d.iterrows():
    date = ""
    for k in date_list:
        date += str(item[k])
    if date not in maps:
        maps[date] = []
        maps[date].append(item[Y])
    maps = dict([(k, np.mean(v)) for k, v in maps.items()])
return maps
def filter(d, X, Y, Z):
    """
    X : Month, Year, Quater
    Y : SoldPrice, ListingPrice, Turnover_rate
    Z : State Name such as CA
    """
    maps = {}
    if X == 'Year':
        date_list = ['Year']
    elif X == 'Quater':
        date_list = ['Year', "Quater"]
    elif X == 'Month':
        date_list = ['Year', "Quater", "Month"]
    return generate_maps(date_list, d, Y, Z)
def mapping(map1, map2):
    """ calculate distance"""
    t = 0.0
    for k, v in map1.items():
        t += (map2[k] - v) * (map2[k] - v)
    return t
if __name__ == "__main__":
    import matplotlib.pyplot as plt
    import numpy.linalg as LA
    d = pandas.read_csv("./tarique_data")
    XSet = ["Year", "Quater", "Month"]
    YSet = ["SoldPrice", "ListingPrice", "Turnover_rate"]
    result = [(X, Y, mapping(filter(d, X, Y, 'CA'),
    filter(d, X, Y, 'NY')))) for X in XSet for Y in YSet]
    best_x, best_y, difference = sorted(result,
    cmp=lambda x, y: -cmp(x[2], y[2]))[0]
    CA, NY = filter(d, best_x, best_y, 'CA'),
    filter(d, best_x, best_y, "NY")
    xset = CA.keys()
    xset.sort()
    y_CA, y_NY = [CA[x] for x in xset],
    [NY[x] for x in xset]
    plt.plot(range(len(xset)), y_CA, label='CA')
    plt.plot(range(len(xset)), y_NY, label='NY')
    plt.legend()
    plt.show()

```


Name	X	Y	Z	Viz	Process
f1	'year'	'sales'	v1 <- 'product'.*	bar.(y=agg('sum'))	v2 <- R(1,v1,f1)
f2	'year'	y1 <- M	v2	bar.(y=agg('sum'))	v3 <- argmax _{x1} [k = 1]T(f1)
f3	'year'	y1	v3	bar.(y=agg('sum'))	y2,v4,v5 <- argmax _{y1,v2,v3} [k = 10]D(f2,f3)
*f4	'year'	y2	v6 <- (v4 v5)	bar.(y=agg('sum'))	

Table 47: The ZQL query which returns varying y-axes visualizations where the following two products differ the most: one whose sales numbers do not change over the year and another which has the largest growth in the number of sales.

Name	X	Y	Z	Viz	Process
f1	x1 <- M	y1 <- M			x3,y3 <- argmax _{x1,y1} [k = 1]sum _{x2,y2} D(f1,f2)
f2	x2 <- M	y2 <- M			
*f3	x3	y3		scatterplot	

Table 48: The ZQL query which returns scatter plot visualization between a pair of attributes whose pattern is most unusual, i.e very different from the patterns made by any other pair of attributes in M.