

# From Sketching to Natural Language: Expressive Visual Querying for Accelerating Insight

Tarique Siddiqui  
Microsoft Research  
tasidd@microsoft.com

Paul Luh  
University of Illinois (UIUC)  
luh2@illinois.edu

Zesheng Wang  
University of Illinois (UIUC)  
zwang180@illinois.edu

Karrie Karahalios  
University of Illinois (UIUC)  
kkarahal@illinois.edu

Aditya G. Parameswaran  
UC Berkeley  
adityagp@berkeley.edu

## ABSTRACT

Data visualization is the primary means by which data analysts explore patterns, trends, and insights in their data. Unfortunately, existing visual analytics tools offer limited *expressiveness* and *scalability* when it comes to searching for visualizations over large datasets, making visual data exploration labor-intensive and time-consuming. We first discuss our prior work on Zenvisage that helps accelerate exploratory data analysis via an interactive interface and an expressive visualization query language, but offers limited *flexibility* when the pattern of interest is under-specified and approximate. Motivated from our findings from Zenvisage, we develop ShapeSearch, an efficient and flexible pattern-searching tool that enables the search for desired patterns via multiple mechanisms: sketch, natural-language, and visual regular expressions. ShapeSearch leverages a novel *shape querying algebra* that can express a large class of shape queries and supports query-aware and perceptually-aware optimizations to execute shape queries within interactive response times. To further improve the usability and performance of both Zenvisage and ShapeSearch, we discuss a number of open research problems.

## 1. INTRODUCTION

With the pressing need to derive value from data, domain experts, spanning virtually all sectors of society, spend considerable time exploring data to identify patterns and trends. These domain experts often have limited understanding of programming and hence rely heavily on visual analytic tools such as Excel and Tableau to understand their data. The state of the art for domain experts is to load their data into a visualization tool, and repeatedly generate visualizations until the desired patterns or insights are identified. Unfortunately, this repeated process of manual examination to scour for desired insights becomes painful, tedious, and time-consuming as the size and complexity of datasets increase. Even on moderately sized datasets, a domain 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.

We characterize this problem of *visualization search* using examples from genomic data analysis.

**Motivating Example.** *Genomic researchers often study genes during clinical trials, e.g., how genes affect clinical trial outcomes,*

©ACM 2021. This is a minor revision of the paper entitled “ShapeSearch: A Flexible and Efficient System for Shape-based Exploration of Trendlines”, published in SIGMOD’20, 978-1-4503-6735-6/20/06, June 14–19, 2020, Portland, OR, USA. DOI: <https://doi.org/10.1145/3318464.3389722>. Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.



Figure 1: Two Facets of the Visualization Search Problem

*how the behavior of genes get affected on specific medications, etc. As an example, given a dataset consisting of clinical trial outcomes (positive vs. negative), researchers often want to find genes that can visually explain the differences in these outcomes. To do so, current tools require the researchers to manually generate tens of thousands of scatter plots—with the x- and y- axes each referring to a gene, and each outcome depicted as a point in the scatterplot—to determine whether the outcomes can be clearly distinguished in the scatter plot.*

*Similarly, researchers study changes in gene expressions while investigating the impact of drugs on disease treatment. For doing so, they often explore trend line visualizations, one corresponding to each gene, with the x-axis as days, and the y-axis as the expression values. For example, when influenced by an external factor, a gene can get induced (up-regulated), or repressed (down-regulated), or can have both induced or repressed pattern within a certain time window. Based on their domain understanding, researchers first hypothesize the expected change in expression that an affected gene should depict. They, then, generate thousands of visualizations, one for each gene, and manually inspect them for the hypothesized patterns.*

In both of the above scenarios, the common theme is the manual examination of a large number of generated visualizations for a specific visual pattern. As depicted in Figure 1, there are two facets to this visualization search problem. First, it is challenging for users to specify the search space of visualizations they are interested in, which forces them to manually generate a large collection of visualizations. The space of visualizations is determined by the number of possible attributes for X and Y axes, aggregation functions and possible subsets of data (denoted by symbol Z in Figure 1a). This space grows exponentially as the size and number of attributes in the data increases. The second facet deals with *visualization matching*. Given a specific pattern of interest, users are typically interested in a subset of visualizations that closely match this pattern. Unfortunately, existing visualization tools are not *expressive* enough to capture either of the two facets.

Our first attempt to address these challenges resulted in a visual data exploration system, Zenvisage [11, 23, 24]. Zenvisage takes as input a high level specification of what the user wants and automatically identifies the relevant visualizations. It supports an interactive interface that allows users to quickly search for simple patterns

via sketching. For expressing more complex search enumeration and visualization matching, Zenvisage supports ZQL—an expressive visualization exploration language that lets users operate over a collection of visualizations using a core set of primitives (e.g., comparison, filtering, sorting) based on visual patterns.

While Zenvisage is an useful first step in solving the visualization search problem, it only supports standard distance measures such as Euclidean distance for matching visualizations, thereby lacking *flexibility* in terms of how a visualization is matched. For instance, it is unable to support search when the desired shape is under-specified or approximate, e.g., finding products whose sales is decreasing over some 3 month window, without specifying when, or those whose sales has many increasing and decreasing portions, without specifying when these portions occur, their magnitude or their width.

Thus, to support such flexible querying mechanisms, we developed ShapeSearch [25, 26], a pattern querying system that supports multiple mechanisms for helping users express and search for desired visual patterns. ShapeSearch incorporates an expressive shape query algebra consisting of shape-based primitives and operators, for expressing a large variety of patterns in trendlines. We developed this algebra after discussions with domain experts, including those from astronomy and genomics, as well as studying a large corpus of pattern queries collected via Mechanical Turk.

ShapeSearch supports multiple specification mechanisms that are internally translated to a shape query algebra representation: ShapeSearch supports a *natural language interface*, coupled with a sophisticated parser and translator for translating them into the algebra. ShapeSearch also supports a *sketching interface* for simpler patterns, and returns visualizations that precisely match the drawn trends. To support more complex needs, the system provides a *visual regular expression* language for issuing queries that cannot be easily expressed via natural language or sketching. The three interfaces can be used simultaneously and interchangeably, as user needs and pattern complexities evolve.

Finally, for ensuring interactive response times on ad-hoc queries, ShapeSearch leverages a pattern-matching engine that relies on minimal pre-processing or indexing. Directly generating and processing a large collection of visualizations, where each visualization has thousands of values, can lead to a long response time. Instead, ShapeSearch uses *perceptually-aware* pattern scoring mechanisms and *query-aware* optimizations—that help prune a large number of visualizations and/or parts of visualizations, for effective and efficient pattern matching.

**Outline.** The rest of our paper is organized as follows. We first discuss our experiences from our prior work on Zenvisage that motivated us to develop ShapeSearch, describing a simple interactive interface and ZQL (Section 2). We then give an overview of ShapeSearch, discussing how it addresses the limitations of Zenvisage (Section 3). Next, we dive into the details of shape algebra that makes the core of ShapeSearch (Section 4). We then describe efficient algorithms for executing shape queries (Section 5). We discuss how we support natural language queries in ShapeSearch (Section 6). In the end, we discuss future directions to further improve the usability and performance of both Zenvisage and ShapeSearch (Section 7).

## 2. EXPERIENCES FROM ZENVISAGE

Zenvisage is a visual analytics system that supports an interactive interface for searching for visualization with simple patterns, along with an expressive query language for more complex queries. We briefly discuss each of these modes and then describe the findings from our user evaluation.

### 2.1 Interactive Search Interface

Figure 2 shows the interactive search interface of Zenvisage loaded with a real estate dataset.

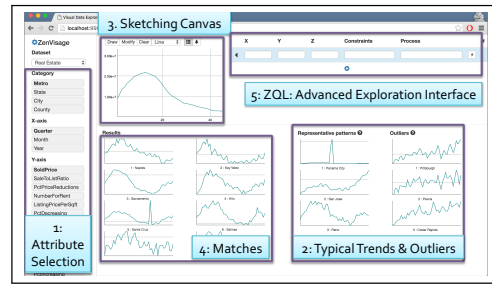


Figure 2: Zenvisage interactive visual query interface

**Attribute Selection.** The first step is attribute selection (Box 1). Here the user can specify the desired X axis attribute, and the desired Y axis attribute for the visualization(s) that the user is interested in. In this case, the user has specified the X axis as quarters (in other words, time), and the Y axis as the real-estate sold price. Additionally, the user specifies the category: this is a variable indexing the space of candidate visualizations the user is operating over. Here, the selected category is “metro”—indicating a metro area or township. We depicted the category as “Z” in Figure 1a.

**Summarization of Typical and Outlier Trends.** As soon as the user selects the X, Y and category, immediately, Zenvisage populates Box 2 with typical or representative trends across categories, and outliers. In this case, there are three typical trends that were found across different metros (i.e., categories): one corresponding to a spike in the middle (Panama City), one to a gradual increasing trend (San Jose), and one to a trend that increased and then decreased (Reno)—most of the other trends were found to be similar to one of these three. The outlier visualizations (Pittsburgh, Peoria, Cedar Rapids) have a large number of seemingly random spikes.

**Drawing or Drag-and-Drop Canvas.** Then, in Box 3, the editable canvas, the user can either draw a shape that they are looking for, or alternatively drag and drop one of the displayed visualizations into the canvas. In this manner, the user indicates that they would like to see a similarity search starting from the shape or pattern that they have drawn or dragged onto the canvas. The user is also free to edit the drawn pattern. In this figure, the user has drawn a trend which is gradually increasing up, then gradually decreasing after that.

**Similarity Search Results.** As soon as the user completes an interaction in Box 3, Box 4 is populated with results corresponding to visualizations (on varying the category) that are most similar to the trend in Box 3, ordered by similarity. The system allows users to choose between three different similarity metrics. Currently, the three metrics Zenvisage provides are Euclidean Distance, DTW, and Segmentation [24].

Overall, this interactive search interface satisfies simple pattern search needs via sketching and drag-and drop, and provides context via representative and outlier patterns. However, it offers limited expressiveness when it comes to more complex data exploration needs. For instance, it is difficult to search for visualizations across a wide range of X and Y attributes (recall that before sketching, we need to set the X and Y axis to specific attributes), or compare two visualizations without using the drawing canvas (e.g., finding 2 products that similar revenue and profit trends over years). Furthermore, one cannot specify multi-step queries involving search for multiple patterns simultaneously, e.g., find products with increasing sales trend in Europe but decreasing sales trend in the US. For supporting these more complex needs, we introduced a second mode, called ZQL, short for Zenvisage Query Language, that users can specify in Box 5 in Figure 4.

### 2.2 ZQL: A Visualization Querying Language

ZQL is a high level language that automates the manual visual data exploration process by allowing users to specify their desired

Name	X	Y	Z	Process
f1	year	soldprice	z1 <- 'state'. '*'	z2 <- $\text{argmax}_{z1} [k = 1] D(f1, f2)$
f2	year	soldpricepersqft	z1	
*f3	year	{soldprice, soldpricepersqft}	z2	

Table 1: A ZQL query retrieving visualizations for a state where the soldprice over year trends are most dissimilar to the soldpricepersqft trend.

Name	X	Y	Z	Process
f1	x1 <- *	y1 <- *	'state'. 'CA'	x2, y2 <- $\text{argmax}_{x1, y1} [k = 1] D(f1, f2)$
f2	x1	y1	'state'. 'NY'	
*f3	x2	y2	'state'. {'NY', 'CA'}	

Table 2: A ZQL query retrieving two different visualizations (among different combinations of x and y) for states of CA and NY that are the most dissimilar.

visualization objective in a few lines. Instead of providing the low-level data retrieval and manipulation operations, users operate at the level of *sets of visualizations, and compare, sort, filter, and transform* visualizations as well as attributes—eventually visualized on either the X or Y axis, or used to sub-select the set of data that is visualized.

We describe the capabilities of ShapeQuery via two examples (depicted via Table 1 and Table 2). Consider the first example where we want to find the states where the soldprice trend is most similar to the soldpricepersqft (i.e., sold price per square foot) trend. Table 1 depicts a 3-line ZQL query for this task. We first compose two collections of visualizations. The first row composes the first collection with X = year, Y = avg(soldprice), and Z = state. \*, consisting of one visualization for each possible state. The Z column corresponds to the Category header in the previous section, indicating the space of visualizations over which the user is operating—in this case, the Z column is fairly simple, there is a single visualization, corresponding to each state. Similarly, the second row composes the second collection with X and Z column stay similar and Y is set to avg(soldpricepersqft).

Once we have composed the two visualization collections (referred via f1 and f2), the Process column is used to compare, sort, and filter the visualizations between the collections. In this example, we iterate the visualizations for each state (notice the variable z1) in f1 and f2 and compare them using a functional primitive D, computing distance, via D(f1, f2). Then, argmin is a sort-filter primitive that sorts the states based on distance scores and selects the top 1 state with minimum scores. Finally, in row 3, we output the overall sales over year visualizations for the selected products as bar-charts. The \* in \*f3 indicates that these visualizations are to be output to the user.

As another example, say we are interested in finding a pair of X and Y axes where the visualizations for two specific states 'NY' and 'CA' differ the most. For doing this, we write a ZQL query depicted in Table 2. In the first line, we fetch all visualizations for the states 'NY' that can be formed by having different combinations of X and Y axes. Similarly in the second row, we retrieve all possible visualizations for the product 'stapler'. In the process column, we iterate over the possible pairs of X and Y axes values, compare the corresponding visualizations in f1 and f2 and finally select the pair of X and Y axis values where the two products differ the most. In the last two rows, we output these visualizations.

Overall, ZQL can capture a wide range of visual exploration queries, including drill-downs and filtering based on specific patterns. We formally describe the expressive power of ZQL using a visual exploration algebra in [23].

## 2.3 Takeaways from User Evaluation

To understand the utility of Zenvisage, we conducted user studies with both novice and experienced data analysts [23], as well as case-studies with collaborating researchers from domains such as genomics, astronomy, and battery science [11].

Our findings show that Zenvisage enables faster and more accurate exploration compared to existing visualization tools such as Tableau, which require considerable manual exploration for find-

ing visualizations with specific patterns. Users who had worked with MATLAB, Python, and R said that ZQL can lead to faster initial exploration of data without requiring to write a lot of code. Those having experience with SQL found ZQL a lot less complicated, less verbose and faster when it comes to comparing subsets of data [23]. Similarly, our collaborating researchers have used Zenvisage for various findings, including the fact that a dip in a light curve was caused by malfunctioning equipment (for astronomy), the fact that a relationship between two specific physical properties of electrolytes was independent of a third one (for battery science), and for reproducing of characteristic gene expression profiles from a recent paper (for genetics) [11].

While Zenvisage offers a promising first step to the problem of painful manual exploration of visualizations, the underlying challenge of visualization search is far from solved. We discovered two main challenges. One pertains to the usability of ZQL. In order to leverage ZQL, domain experts need to learn and switch to a new querying language, a major hindrance to its broader adoption. Domain experts with prior experience with computational notebooks often expressed a need for transitioning between writing code and using ZQL abstractions. Additionally, instead of writing their queries in one step, users often intended to construct them in an incremental manner using prior queries as context. In Section 7, we discuss these issues and potential solutions in more detail, highlighting another system LUX [1] from Lee et al. that partially addresses these issues.

The second challenge with Zenvisage deals with how visualizations are matched. For the rest of the paper, we focus on this challenge and present a new system ShapeSearch to address it.

### 2.3.1 The Problem of Flexible Shape Matching

The sketch-based interface in Zenvisage as well as other similar visualization searching tools [7, 13, 27] offer limited flexibility in terms of how a visualization is matched. For instance, visualization search often involves pattern matching where the desired pattern of interest is under-specified and approximate, e.g., finding stocks whose prices are decreasing for some time, followed by a sharp rise, with the position and intensity of movements being left unspecified, or when the desired shape is complex, e.g., finding gene expression profiles where there is an unspecified number of peaks and valleys followed by a flattening out. We highlight the key characteristics of such pattern matching tasks below.

**Fuzzy Matching.** Domain experts (i) typically search for patterns that are *approximate*, and are often not interested in the specific details or local fluctuations as much as the overall shape, and (ii) they often *do not* specify or even know the exact location of the occurrence of patterns. For example, biologists routinely look for structural changes in gene expression, e.g., rising and falling at different times (Figure 3a), characterizing internal biological processes such as the cell cycle or circadian rhythms, or external perturbation, such as the influence of a drug or presence of a disease.

**Searching Multiple Simple Patterns.** We notice that domain experts often describe complex patterns using a *combination of multiple simple ones*. Each individual pattern is typically described

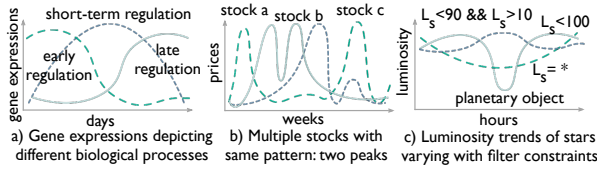


Figure 3: Shapes characterizing real world phenomena

using words such as "increasing", "stable", "falling", that are easy to state in natural language but hard to specify using existing query languages. Moreover, pattern matching tasks often go beyond finding a sequence of patterns, requiring arbitrary combinations, e.g., disjunction, conjunction or quantification, with varying location or width constraints. Examples include finding stocks with at least 2 peaks within a span of 6 months, e.g., the so-called "double/triple top" patterns that indicate future downtrends [2], or finding cities where the temperature rises from November to January and falls during May to July such as Sydney.

**Ad-hoc and Interactive Querying.** Pattern-based queries are often defined *on-the-fly* during analysis, based on other patterns observed. For instance, biologists often search for a pattern in a group of genes similar to a pattern recently discovered in another group [11]. Similarly, astronomers monitor the shape of the luminosity trends of stars over time to search for and characterize new planetary objects (Figure 3c). For example, a dip in brightness often indicates a planetary object passing between the star and the telescope.

To address these issues, we developed ShapeSearch, described next.

### 3. OVERVIEW OF SHAPESEARCH

ShapeSearch provides powerful yet flexible mechanisms for users to search for trendline visualizations with a desired shape. In this section, we first present an overview of ShapeSearch along with user experience.

ShapeSearch supports an interactive interface for composing shape queries. Figure 4 depicts this interface, with an example query on genomics data discussed in the introduction. Here, the user is interested in searching for genes that get suppressed due to the influence of a drug, depicted by a specific shape in their gene expression—first rising, then going down, and finally rising again—with three patterns: up, down, and up, in a sequence. To search for this shape, the user first loads the dataset [6] via form-based options on the left (Figure 4 Box 1), and then selects the space of visualizations to explore by setting the *x* axis as time, the *y* axis as expression values, and the category as gene. Each value of the category attribute results in a candidate visualization with the given *x* and *y* axis. Thus, the category attribute defines the space of visualizations over which we match the shape. ShapeSearch supports three mechanisms for shape specification—natural language, regular expressions (regex for short), and sketching on a canvas:

**Sketching on Canvas.** By drawing the desired shape as a sketch on the canvas (Figure 4 Box 2a), the user can search for visualizations that are *precisely* similar (using a distance measure such as Euclidean distance or Dynamic Time Warping [18]). As soon as the user finishes sketching, ShapeSearch outputs visualizations that are similar to the drawn sketch in the results panel (Figure 4 Box 4).

**Natural Language (NL).** For searching for visualizations that approximately match patterns, users can use natural language. For instance, as in Figure 4 Box 2b, the desired shape in the aforementioned genomics example can be expressed as “*show me genes that are rising, then going down, and then increasing*”. Similarly, scientists analyzing cosmological data can easily search for supernovae (bright stellar explosions) using “*find objects with a sharp peak in*

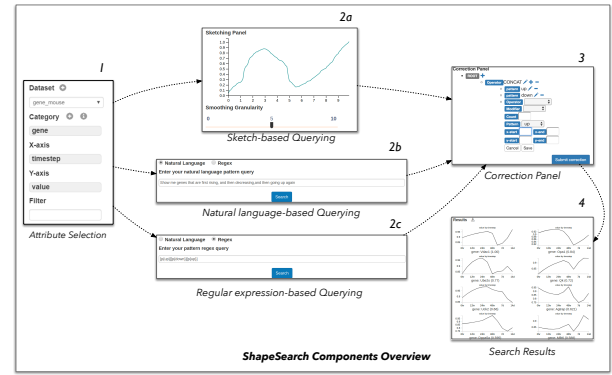


Figure 4: ShapeSearch Interface, consisting of six components. 1) Data upload, attribute selection, and applying filter constraints 2) Query specification: 2a) Sketching canvas 2b) Natural language query interface, and 2c) Regular expression interface, 3) Correction panel, and 4) Results panel

*luminosity*”. We describe in Siddiqui et al. [26] how ShapeSearch translates natural language queries to a structured internal representation.

**Regular Expression (regex).** For queries that involve complex combinations of patterns that are difficult to express using natural language or sketch, the user can issue a regular expression-like query that directly maps to the structured internal representation, consisting of ShapeSearch primitives and operations, described in detail in Section 4.

During exploration, users can choose specification mechanisms interchangeably based on the complexity of the query. For both NL as well as regex, ShapeSearch additionally supports an auto-complete functionality to guide users towards their target query. We use the term *user query* to refer to the submitted query using any of the specification mechanisms.

The ShapeSearch back-end parses and translates the user query into a ShapeQuery, a structured internal representation of the query consisting of operators and primitives supported in our algebra (Section 4). The back-end supports an ambiguity resolver that uses a set of rules for automatically resolving syntactic and semantic ambiguities, as well as forwards the parsed query to the user for further corrections and validation (Figure 4 Box 3). The validated query is finally optimized and executed by the execution engine (Section 4.3), and the top visualizations that best match the ShapeQuery are presented to the user in the results panel (Figure 4 Box 4). Next, we discuss a ShapeQuery algebra that makes the core of ShapeSearch.

### 4. SHAPE ALGEBRA

ShapeQueries help express a large variety of patterns over trendlines with a minimal set of primitives and operators. A ShapeQuery represents a *shape* as a combination of multiple *simple patterns*. A simple pattern can either be precise with specific location constraints, e.g., matching  $y = x$  between  $x = 2$  to  $x = 6$ , or fuzzy, e.g., roughly increasing, where the notion of the pattern is approximate and its location unspecified. Each simple pattern along with its precise or imprecise constraints is called a ShapeSegment. Complex shapes, e.g., rising and then falling, are formed by combining multiple ShapeSegments using one or more *operators*. One can search for multiple patterns in a sequence (concat,  $\otimes$ ) or matching the same sub-region of the trendline (and,  $\odot$ ), or one of many patterns matching a sub-region (or,  $\oplus$ ), described later.

As an example, “rising from  $x=2$  to  $x=5$  and then falling” can be translated into a ShapeQuery  $[x.s=2, x.e=5, p=up] \otimes [p=down]$  consisting of two ShapeSegments separated by a  $\otimes$  operator. The first ShapeSegment captures “rising from  $x = 2$  to  $x = 5$ ”; the second expresses a “falling” pattern. Since the second must “follow”



Table 3: Primitives and Operators in ShapeQuery

Symbol	Name	Type
$x.s$	START X VALUE	Location Sub-Primitive
$y.s$	START Y VALUE	Location Sub-Primitive
$x.e$	END X VALUE	Location Sub-Primitive
$y.e$	END Y VALUE	Location Sub-Primitive
$v$	SKETCH	Location Sub-Primitive
$.$	ITERATOR	Location Sub-Primitive
$p$	PATTERN	Primitive
$s$	POSITION	Pattern Sub-Primitive
$m$	MODIFIER	Primitive
$>$	MORE	Modifier value
$>2$	ATLEAST 2X	Modifier value
$=$	SIMILAR	Modifier value
$\otimes$	CONCAT	Operator
$\odot$	AND	Operator
$\oplus$	OR	Operator
$!$	OPPOSITE	Operator

the first, the two ShapeSegments are combined using the CONCAT operator, denoted by  $\otimes$ . We now describe the shape primitives and operators that constitute the ShapeQuery algebra. Table 3 lists these primitives and operators.

#### 4.1 Shape Primitives and Operators

A ShapeSegment is described using two high level primitives: LOCATION and PATTERN. The LOCATION values can be skipped in order to match the PATTERN anywhere in the trendline. Similarly, users can input the exact trendline to match, or the endpoints of the ShapeSegments to match without specifying the PATTERN.

**Specifying LOCATION.** LOCATION defines the endpoints of the sub-region of the trendline between which a pattern is matched: starting X/Y coordinate ( $x.s/y.s$ ), ending X/Y coordinate ( $x.e/y.e$ ). For example,  $[x.s=2, x.e=10, y.s=10, y.e=100]$  is a simple ShapeQuery to find trendlines whose trend between  $x=2$  to  $x=10$  is similar to the line segment from (2, 10) to (10, 100). Users can also draw a sketch to find trendlines similar to the sketch, a functionality supported in other tools alluded to in the introduction [7, 14, 23]. ShapeSearch translates the pixel values of the user-drawn sketch to the domain values of the X and Y attributes, and adds the transformed vector of ( $x, y$ ) values as a vector  $v$  in the ShapeQuery. As an example, the ShapeQuery  $[v=(2:10, 3:14, \dots, 10:100)]$  finds trendlines that have precisely similar values to  $v$  using some distance measure, e.g., Euclidean distance, or dynamic time warping [18].

**Specifying PATTERN.** PATTERN defines a trend or a semantic feature in a sub-region of the trendline. A number of basic semantic patterns, commonly used for characterizing trendlines, are supported, such as *up*, *down*, *flat*, or the slope ( $\theta$ ) in degrees. For example  $[p=up]$  finds trendlines that are increasing,  $[p=45]$  finds trendlines that are increasing with a slope of about  $45^\circ$ , and  $[x.s=2, x.e=10, p=up]$  finds trendlines that are increasing from  $x=2$  to  $x=10$ . Finally, one can use  $p=*$  to match any pattern and  $p=empty$  to ensure that there are no points over the sub-region.

**Combining PATTERNS.** ShapeQuery supports three operators to combine ShapeSegments:

- CONCAT ( $\otimes$ ) specifies a sequence of two or more ShapeSegments. For example, using  $[p=up] \otimes [p=down]$  one can search for genes that are first rising, and then falling. Note that  $\otimes$  is one of the most frequently used operations, and we sometimes omit  $\otimes$  between ShapeSegments, e.g.,  $[p=up][p=down]$ , to make it succinct to describe.
- AND ( $\odot$ ) simultaneously matches multiple patterns in the same sub-region of the trendline. Unlike CONCAT, all of the patterns must be present in the same sub-region. For example, one can look for genes whose expression values rise twice but do not fall more than once within the same sub-region.
- OR ( $\oplus$ ) searches for one among many patterns in the same sub-region of the trendline, picking the one that matches the sub-region best. For example, one can search for genes whose expressions are either *up*- or *down*-regulated.

Table 4: Pattern Scores

P	Score
<i>up</i>	$\frac{2 \cdot \tan^{-1}(\text{slope})}{\pi}$
<i>down</i>	$-\frac{2 \cdot \tan^{-1}(\text{slope})}{\pi}$
<i>flat</i>	$(1.0 - \ \frac{4 \cdot \tan^{-1}(\text{slope})}{\pi}\ )$
$\theta = x$	$(1.0 - \ \frac{2 \cdot \tan^{-1}(\text{slope}-x)}{(\pi -  \tan^{-1}(x) )}\ )$
$*$	1
<i>empty</i>	-1
$v$	$L_2$ norm (configurable)

Table 5: Operator Scores

O	Score
$\otimes$	$\sum_{i=1}^k \text{score}_i / k$
$\odot$	$\min(\text{score}_1, \dots, \text{score}_k)$
$\oplus$	$\max(\text{score}_1, \dots, \text{score}_k)$

**Comparing Patterns.** In some cases, one may want to compare the pattern in a ShapeSegment with the preceding or succeeding ShapeSegments. To support such use cases, ShapeSearch (i) allows a ShapeSegment to refer to the previous or the next ShapeSegment using  $\$+$  or  $\$-$  respectively, and (ii) compare patterns between the current and referred ShapeSegment using operations  $>$ ,  $<$ , or  $=$ . For example, astronomers can issue a ShapeQuery  $[p=up] \otimes [p < \$-p]$  with  $x=\text{time}$  and  $y=\text{luminosity}$  (brightness) to search for celestial objects that were initially moving rapidly towards earth, but after some point either slowed down or started moving away.  $[p < \$-p]$  ensures that the slope of brightness over time is less than that in the previous sub-region  $[p=up]$ .

**Expressing Complex Patterns.** The aforementioned basic primitives and operators are powerful enough to express more complex ShapeSearch use-cases. We discuss three such complex patterns below, along with shortcuts for their easy specification.

1. *Searching shapes of specific width.* In some cases, users want to find specific shapes irrespective of their start location, e.g., searching for cities with the steepest rise in temperature over a width of 3 months. To express such queries, ShapeSearch supports the ITERATOR ( $.$ ), e.g.,  $[x.s=., x.e=x.s+3, p=up]$  that iterates over all points in the trendline, setting each point as the start  $x$  position, with the  $x$  end position set to 3 units ahead. Internally, for a trendline of length  $n$ , this query can be rewritten as an OR operation over  $(n-3+1)$  ShapeSegments, where, for the  $i$ th ShapeSegment,  $x.s=i$  and  $x.e=i+3$ .

2. *Quantifiers.* One can search for trendlines where a pattern occurs a specific number of times using quantifiers, denoted by  $q$ . For example,  $[p=up, q=\{1, 2\}]$  can be used to search for trendlines where there is an increasing pattern at least once and at most twice. Quantifiers can be internally rewritten using an OR of one or more CONCAT operations. For example, the above query is rewritten as  $(([p=*] \otimes [p=up] \otimes [p=*) \oplus ([p=*] \otimes [p=up] \otimes [p=*) \otimes [p=up] \otimes [p=*)]$ .

3. *Nesting.* A combination of patterns can be constrained to be within a specific sub-region by specifying them as a value of the PATTERN primitive. For example, to search for stocks that increased anytime between February to October, we can use nesting as follows:  $[x.s=2, x.e=10, p=([p=*] [p=up] [p=*)]$ . This can be rewritten using CONCAT operations as follows:  $[x.s=2, p=*) \otimes [p=up] \otimes [p=*) \otimes [x.s=10, p=*)]$ .

#### 4.2 Scoring

A ShapeQuery  $Q$  operates on one trendline,  $V_i$ , at a time, and returns a real number, called *score*, between  $-1$  to  $+1$ . The ShapeQuery  $Q$  operates on  $V_i$  with the help of ShapeSegments ( $S_1, S_2, \dots, S_n$ ) and operators ( $O_1, O_2, \dots, O_m$ ). Each ShapeSegment  $S_i$  operates on  $V_i^{p,q}$ , a sub-region of  $V_i$  starting at  $p = x.s$  and ending at  $q = x.e$  and returns a  $\text{score}_i \in [-1, 1]$  using scoring functions we describe subsequently. One or more ShapeSegments are combined using operators such as  $\otimes, \odot, \oplus$ . Formally, an operator  $O_i$  takes as input the scores  $\text{score}_1, \text{score}_2, \dots, \text{score}_n$  from its  $n$  input ShapeSegments and outputs another  $\text{score}_i$  using scoring functions that capture the behavior of the operators.

For both efficiency and effectiveness, ShapeSearch approximates each sub-region with a line, using the slope to quantify how closely

it captures any given ShapeSegment. As depicted in Table 4, ShapeSearch uses different scoring functions for each pattern primitive that transforms the slope to a value in  $[-1, 1]$  using a  $\tan^{-1}$  function. For example, for an *up* pattern, the function returns a score between  $[0, 1]$  for a trendline with a slope from  $0^\circ$  to  $90^\circ$ , a score of  $[-1, 0]$  for a slope of less than  $0^\circ$  (opposite of *up*).

For execution, ShapeSearch takes the entire trendline, the Abstract Tree Representation (AST) of ShapeQuery, and the list of scoring functions *ScrFunc* as in Tables 4 and 5 as inputs. If the root node of the ShapeQuery tree is a ShapeSegment, ShapeSearch directly computes the score of ShapeSegment on the specified part of the trendline. If the root node is  $\odot$  or  $\oplus$ , ShapeSearch invokes each of the operands (i.e., child sub-trees) to compute their scores on the sub-region independently, combining the scores as per operator's functions. However, if the root node is a CONCAT with  $k$  operands, i.e., child sub-trees, ShapeSearch segments  $L$  into all possible  $k$  sub-regions:  $L_1, L_2, \dots, L_k$  and then, for each segmentation, invokes the  $i$ th operand on  $i$ th segment. Finally, the maximum score across all segmentations is output.

### 4.3 Executing Fuzzy ShapeQueries

A common subclass of ShapeQueries are *fuzzy* ShapeQueries, consisting of at least one ShapeSegment with missing or multiple possible values for  $x.s$  or  $x.e$ . Thus, for fuzzy ShapeQueries, we try all possible values of  $p$  and  $q$ , selecting the sub-region that leads to the best score. This becomes prohibitively expensive as the number of points in the trendline increases. For a CONCAT with  $k$  operands, the exhaustive approach creates  $n^{(k-1)}$  segmentations, where  $n$  is the number of points in the trendline.

**The Dynamic Programming Algorithm.** We can show [26] that for the CONCAT operation, the scoring of the  $j$ th operand on  $j$ th sub-region does not depend on the scoring of the first  $j-1$  operands on the first  $j-1$  sub-regions. We use this idea to develop a faster dynamic programming algorithm (DP) for scoring CONCAT operations over ShapeSegments. Formally, let  $OPT(1, t, (1 : j-1))$  be the best score corresponding to the optimal segmentation over the sub-region between  $x = 1$  to  $x = t$  for the first  $j-1$  operands, and  $SC(t+1, i, j)$  be the score of the  $j$ th operand over the sub-region between  $x = t+1$  and  $x = i$ . Then, the optimal segmentation  $OPT(1, i, (1 : j))$  for the first  $j$  operands over  $x = 1$  and  $x = i$  can be computed using the following recursion:

$$OPT(1, i, (1 : j)) = \underset{t}{MAX} \{ \frac{(j-1) \times OPT(1, t, (1 : j-1)) + SC(t+1, i, j)}{j} \}$$

Unfortunately, even though the DP algorithm is orders of magnitude faster than the exhaustive approach, for trendlines with large number of points, even a ShapeQuery with a single CONCAT operation can be slow, because of its quadratic runtime. We, next, discuss optimizations to further decrease the runtime of CONCAT operation on ShapeSegments.

#### 4.3.1 A Pattern-Aware Bottom-up Approach

While the DP-based optimal approach scores all possible sub-regions for each operand in the CONCAT operation, a more efficient approach could be to select end points to be those where the slope (or pattern) changes drastically. We first illustrate our intuition, and then describe an algorithm that performs segmentation in a pattern-aware manner.

**Intuition.** As depicted in Figure 5, consider two sub-regions  $A$  on the left and  $B$  on the right for the trendline  $L$ . Say the trendline in sub-region  $A$  is inverted V-shaped, i.e., increasing until a point  $P$  and then decreasing. Now, for all possible segmentations where  $[p=up]$ 's sub-region lies completely in  $A$ , there are the following possibilities for  $x.e$  of  $[p=up]$ : 1)  $[p=up]$ 's  $x.e$  point is before  $P$ . 2)  $[p=up]$ 's  $x.e$  point is after  $P$ . 3)  $[p=up]$ 's  $x.e$  point is at  $P$ .

Since  $[p=down]$  follows  $[p=up]$ , we can see that option 1 that sets  $[p=up]$ 's  $x.e < P$  is less likely to be optimal as that will lead to scoring of a part of  $[p=down]$  on an increasing trend. Similarly,

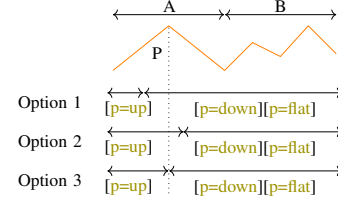


Figure 5: Pattern-aware selection of LOPs

$x.e > P$  is less optimal as that will lead to scoring of a part of  $[p=up]$  on a decreasing trend. Thus, if we have to (greedily) select one point in sub-region  $A$  for  $[p=up]$ 's  $x.e$ ,  $P$  is likely a better choice. We call such a point as *locally optimal point* (LOP).

**A Bottom-up Algorithm.** Based on the above intuition, we develop a much faster algorithm that uses the following assumption to reduce the number of segmentations.

**Assumption 4.1** (Closure). *If a point is not locally optimal for any of the sub-expressions in the CONCAT operation (i.e., a CONCAT on a sub-sequence of the operands), it cannot be  $x.s$  or  $x.e$  of a ShapeSegment in the optimal segmentation.*

That is, local optimality leads to global optimality. Because of this assumption, our proposed algorithm is approximate. However, our empirical results show that despite this assumption, the accuracy of the algorithm is very close to that of DP, while taking orders of magnitude less time.

At a high level, the algorithm starts by dividing the trendline into smaller contiguous sub-regions. Next, it selects locally optimal points (LOPs) over small sub-regions, followed by a bottom-up merging step that uses LOPs over small sub-regions to find LOPs over larger sub-regions.

**Selection of LOPs.** We define a point  $P$  to be a LOP in a sub-region  $A$  for the sub-expression  $S_i$  if it is either the  $x.e$  of the first ShapeSegment or  $x.s$  of the last ShapeSegment of  $S_i$ . For instance, in the above example, it is easy to see that a LOP  $P$  in sub-region  $A$  is the  $x.e$  value of  $[p=up]$  in the optimal segmentation of  $[p=up] \otimes [p=down]$  in  $A$ . Since a CONCAT operation with  $k$  operands can have  $(k^2)$  sub-sequences, there can be a maximum of  $2.k^2$  LOPs in  $A$ .

**Merging.** Next, we incrementally merge nodes in a bottom-up fashion to select LOPs over larger sub-regions. For example, in Figure 6, node 4 depicts the sub-sequences formed by combining sub-sequences from nodes 1 and 2, and node 5 depicts the sub-sequences formed by combining sub-sequences from nodes 3 and 4. When multiple sub-sequences in the children nodes generate the same sub-sequence in the parent node, we select the one with maximum score after concatenation (i.e., the one with the most optimal segmentation), thereby pruning out LOPs corresponding to non-selected sub-sequences. For example, at node 5,  $a \otimes b$  can be computed from 1)  $a$  from node 3 and  $b$  from node 4, 2)  $a \otimes b$  from node 3 and  $b$  from node 4, and 3)  $a$  from node 3 and  $a \otimes b$  from node 4. Among these 3 concatenations, we pick the one that gives the maximum score. This merging process is repeated at each intermediate node. Finally, at the root node, we select the points that result in the maximum score for the entire sequence of operands. More details along with the pseudo-code can be found in [4].

Given the closure assumption, we prove in [4] that the merging process leads to optimal segmentation and that the bottom-up algorithm with  $k$  CONCAT operands is optimal with a time complexity of  $O(nk^4)$ , i.e., linear in the number of points in the trendlines.

## 5. NATURAL LANGUAGE TRANSLATION

So far, we haven't described how natural language queries are parsed into ShapeQueries. We provide a brief overview of the three key steps involved in parsing, and refer readers to our extended report [4] for additional details. We use the following natu-

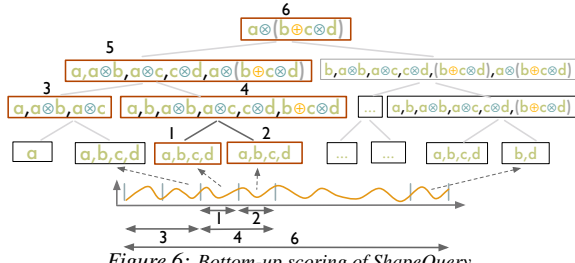


Figure 6: Bottom-up scoring of ShapeQuery

ral language query collected from MTurk for illustration: “show me the trendlines that are increasing from 2 to 5 and then decreasing”.

**Step 1. Primitives and Operators Recognition.** Given a natural language query, the first step is to map words to their corresponding shape primitives and operators. For example, the above query is tagged as “show (noise) me (noise) the (noise) trendlines (noise) that (noise) are (noise) increasing (p) from (noise) 2 (x.s) to (noise) 5 (x.e) and then (⊗) decreasing (p)”. In order to do so, we learn a linear-chain conditional-random field model (CRF) [10] and train it on the same 250 natural language queries we collected via Mechanical Turk (described in [4]) for understanding query characteristics. For each word, we use its part-of-speech (POS) tags along with word-level context as features.

**Step 2. Identifying Pattern Value.** For each of the words predicted of type p, e.g., increasing and decreasing in the above query, we additionally map them to the corresponding semantic pattern supported in ShapeSearch, e.g., “increasing” is mapped to p=up. For this mapping, ShapeSearch computes the similarity between the specified word and synonyms of the supported patterns, first using edit distance and then using wordnet [19]. The semantic pattern with the highest similarity between any of its synonyms and the specified word is selected.

**Step 3. ShapeQuery Generation and Ambiguity Resolution.** Next, we group primitives and operators into a ShapeQuery, first grouping all primitives between two operators into a single ShapeSegment. For the example query, the primitives are grouped as follows: [increasing (p=up), 2 (x.s), 5 (x.e)] and then (⊗) [decreasing (p=down)]. In some cases, this leads to incorrect grouping of primitives, e.g., two patterns in the same ShapeSegment. There could further be semantic ambiguity because of wrong entity tagging, e.g., decreasing (p=up) from 5 (y.s) to 10 (y.e) where x.s and x.e values are wrongly tagged as y.s and y.e respectively. ShapeSearch uses rule-based transformations that try to reorder and change the types of entities to get a correct and meaningful ShapeQuery [26].

The parsed ShapeQuery is sent to the front-end (Box 4 in Figure 4) for users to edit or further refine it if needed. The validated query is then executed to generate the matching trendlines.

## 6. FUTURE DIRECTIONS

We now discuss open research directions for improving the usability and performance of both Zenvisage and ShapeSearch.

### 6.1 Search Enumeration + Shape Matching

In ShapeSearch, users currently need to specify the X and Y attribute before issuing ShapeQueries. However, in certain scenarios, users may not know the X and Y attributes in advance or may want to search for the same shape over different combinations of attributes. Additionally, users may want to issue a multi-step query involving multiple shapes at the same time, finding states with decreasing listing prices trends but increasing soldprice trends of houses. To support such complex data exploration needs, we envision integrating ZQL with ShapeQuery. One simple option is support ShapeQuery as a functional primitive as part of the Process column in ZQL. For instance, Table 6 depicts an integrated query

for the above example for finding states with decreasing listing prices trends but increasing soldprice trends. Combining ZQL and ShapeQuery also adds to expressiveness and efficiency of ZQL—functional primitives are currently treated as black boxes and thus not optimized in Zenvisage. By adding support for ShapeQuery, Zenvisage can leverage the shape matching algorithms discussed earlier for efficient processing of visualizations.

### 6.2 In-Database Support for Fuzzy Matching

ShapeSearch performs shape matching outside relational databases; consequently as the size of the dataset increases, the data transfer and serialization/deserialization overheads tend to dominate, resulting in an increase in latency. On the other hand, recognizing patterns in a sequence of rows in relational databases has been widely desired but only supported by a few vendors. For instance, Oracle Database 12c supports a MATCH RECOGNIZE [3] clause for pattern matching in native SQL. SQL-TS (Simple Query Language for Time Series) [21] is another proposal on SQL extensions for pattern queries. Nevertheless, none of these extensions support fuzzy matching capabilities, instead they require users to define the patterns (e.g, up, down) using values of matching columns—making the specification quite tedious and verbose.

In order to support fuzzy shape queries, we envision developing new database extensions that take as input a ShapeQuery as part of the SQL query and leverage shape matching algorithms for efficiently executing the ShapeQuery within the database kernel. For instance, the following query depicts potential extensions for supporting ShapeQuery within the SQL syntax.

```
SELECT *
FROM Ticker T,
(MATCH BY symbol ON price
USING PATTERN [p=*]⊗p=down]⊗[p=up]⊗[p=*) AS score
ORDER BY score DESC
LIMIT 1) S
WHERE T.symbol = S.symbol
```

Given a table Ticker, the above query finds a stock symbol (specified via MATCH BY clause) with closest matching V-shaped trend on the values of column price (specified via ON clause) and outputs its corresponding tuples.

### 6.3 Supporting Context-Aware Search

During exploratory data analysis, users do not always have a precise pattern query to start with, instead they compose queries as the exploration evolves. Often, they break their pattern queries into a sequence of simpler queries that build upon prior ones. Thus, for a more fluid user-experience, there are interesting avenues for future work that capture context, suggest next steps, and support incremental composability. One option is to make meaningful query suggestions by mining query patterns from past search logs and match them with immediately preceding queries (i.e., the context) in the same session. Since there can be a large number of patterns, efficiently searching for prefixes while effectively capturing user’s intent is an interesting challenge.

For incremental composability, the context of user exploration can be represented using a state machine consisting of partial queries as different possible states. The state can then be incrementally updated as new queries arrive. For novice analysts, interactive querying interfaces, similar to systems such as GestureDB [16] and DataPlay [5] can help make query specification even easier. In addition, like in Zenvisage, ShapeSearch can be extended to automatically find typical and outliers patterns to help the users get started quickly.

### 6.4 Mixing Code and Interaction

To accelerate the process of data exploration, another important next step is to integrate visualization search abstractions supported



Name	X	Y	Z	Process
f1	year	listprice	z1 <- 'state'.*	z2 <- argany <sub>z1</sub> [p=down](f1)
f2	year	soldprice	z1	z3 <- argany <sub>z1</sub> [p=up](f2)
*f3	year	{listprice, soldprice }	z2 && z3	

Table 6: Example of a ZQL query using ShapeQuery as a functional primitive within Process column. The query finds states with decreasing listprice but increasing sold price over trends of houses.

via ZQL and ShapeQuery with existing data science libraries such as Pandas. This will allow users to seamlessly transition between writing code (e.g., for data edits, cleaning, and transformation); getting recommendations via search specifications, and performing interactions on visualizations—all in one place. As a step in this direction, LUX [1], a recent Python library, combines partial user-specifications with best practices from visual data analysis to recommend interesting visualizations for guiding users towards next steps. It further displays visualizations as a widget in-situ within a Jupyter notebook to support easier transitions between code and interaction. While specification in LUX is inspired from ZQL, adding natural-language or regex-based pattern searching functionalities, as supported in ShapeQuery, can further enhance the power of such libraries.

## 7. RELATED WORK

Our work draws on prior work in visual querying, as well as symbolic pattern mining. Visual querying tools [14, 15, 20, 23, 27] help users search for visualizations with a desired shape by taking as input a sketch of that shape. Most of these tools perform precise point-wise matching using measures such as Euclidean distance or DTW. A few tools such as TimeSearcher [7] let users apply soft or hard constraints on the  $x$  and  $y$  range values via boxes or query envelopes, but do not support mechanisms for specifying shape primitives beyond location constraints. ShapeSearch introduces a novel algebra that improves extensibility by acting as a common “substrate” for various input mechanisms, along with an optimization engine that efficiently matches patterns against a large collection of trendlines.

Symbolic sequence matching papers approach the problem of pattern matching by employing offline computation to chunk trendlines into fixed length blocks, encoding each block with a symbol that describes the pattern in that block [8, 9, 12, 17, 22]. Since these work operates on pre-chunked-and-labeled trendlines, the problem is one of matching regular expressions against string sequences (one per pre-labeled trendline). Most of these papers only return a boolean score for whether the pattern matches the string sequence. Moreover, since the trendlines are pre-labeled and indexed, they do not support on-the-fly pattern matching where the same trendline can change shapes based on filters or aggregation constraints. ShapeSearch, on the other hand, adopts a more online query-aware ranking of trendlines without requiring precomputation, and is thus more suited for ad-hoc data exploration scenarios.

## 8. CONCLUSION

In this work, we described ShapeSearch, a pattern matching system that complements our prior system Zenvisage by providing expressive and flexible mechanisms for domain experts to effortlessly and efficiently search for trendline visualizations. We described the ShapeQuery algebra that forms the core of ShapeSearch, and helps express a large variety of patterns with a minimal set of primitives and operators. The algebra is backed by a shape matching engine that enables on-the-fly and scalable pattern matching. Overall, together with Zenvisage, ShapeSearch offers a promising first step towards substantially simplifying and improving the process of interactive data exploration for novice and expert analysts alike.

## 9. REFERENCES

- [1] <https://github.com/lux-org/lux>. <https://github.com/lux-org/lux>.
- [2] Investopedia. <https://www.investopedia.com/terms/t/tripletop.asp>.
- [3] Match recognize. <https://bit.ly/3bWwUhs>. Online; accessed 17-Aug-2015].
- [4] Technical report. <https://arxiv.org/abs/1811.07977>.
- [5] A. Abouzied, J. M. Hellerstein, and A. Silberschatz. Playful query specification with dataplay. *PVLDB*, 5(12):1938–1941, 2012.
- [6] C. J. Bult, J. T. Eppig, J. A. Kadin, J. E. Richardson, J. A. Blake, and M. G. D. Group. The mouse genome database (mgd): mouse biology and model systems. *Nucleic acids research*, 36(suppl\_1):D724–D728, 2008.
- [7] P. Buono, A. Aris, C. Plaisant, A. Khella, and B. Shneiderman. Interactive pattern search in time series. In *Visualization and Data Analysis 2005*, volume 5669, pages 175–187. International Society for Optics and Photonics, 2005.
- [8] C. Faloutsos, M. Ranganathan, and Y. Manolopoulos. *Fast subsequence matching in time-series databases*, volume 23. ACM, 1994.
- [9] M. N. Garofalakis, R. Rastogi, and K. Shim. Spirit: Sequential pattern mining with regular expression constraints. In *VLDB*, volume 99, pages 7–10, 1999.
- [10] J. Lafferty, A. McCallum, and F. C. Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. 2001.
- [11] D. J.-L. Lee, J. Lee, T. Siddiqui, J. Kim, K. Karahalios, and A. Parameswaran. You can’t always sketch what you want: Understanding sensemaking in visual query systems. *IEEE transactions on visualization and computer graphics*, 2019.
- [12] R. A. K.-I. Lin and H. S. S. K. Shim. Fast similarity search in the presence of noise, scaling, and translation in time-series databases. In *Proceeding of the 21th International Conference on Very Large Data Bases*, pages 490–501. Citeseer, 1995.
- [13] M. Mannino and A. Abouzied. Expressive time series querying with hand-drawn scale-free sketches. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, page 388. ACM, 2018.
- [14] M. Mohebbi, D. Vanderkam, J. Kodys, R. Schonberger, H. Choi, and S. Kumar. Google correlate whitepaper. 2011.
- [15] P. K. e. a. Muthumanickam. Shape grammar extraction for efficient query-by-sketch pattern matching in long time series. In *Visual Analytics Science and Technology (VAST), 2016 IEEE Conference on*, pages 121–130. IEEE, 2016.
- [16] A. Nandi, L. Jiang, and M. Mandel. Gestural query specification. *PVLDB*, 7(4):289–300, 2013.
- [17] R. A. G. Psaila and E. L. Wimmers Mohamed & It. Querying shapes of histories. *Very Large Data Bases. Zurich, Switzerland: IEEE*, 1995.
- [18] L. Rabiner, A. Rosenberg, and S. Levinson. Considerations in dynamic time warping algorithms for discrete word recognition. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 26(6):575–582, 1978.
- [19] R. P. Roetter, C. T. Hoanh, A. G. Laborte, H. Van Keulen, M. K. Van Ittersum, C. Dreiser, C. A. Van Diepen, N. De Ridder, and H. Van Laar. Integration of systems network (sysnet) tools for regional land use scenario analysis in asia. *Environmental Modelling & Software*, 20(3):291–307, 2005.
- [20] K. Ryall, N. Lesh, T. Lanning, D. Leigh, H. Miyashita, and S. Makino. Querylines: approximate query for visual browsing. In *CHI’05 Extended Abstracts*, pages 1765–1768, 2005.
- [21] R. Sadri, C. Zaniolo, A. Zarkesh, and J. Adibi. Expressing and optimizing sequence queries in database systems. *ACM Transactions on Database Systems (TODS)*, 29(2):282–318, 2004.
- [22] H. Shatkey and S. B. Zdonik. Approximate queries and representations for large data sequences. In *Data Engineering, 1996. Proceedings of the Twelfth International Conference on*, pages 536–545. IEEE, 1996.
- [23] T. Siddiqui, A. Kim, J. Lee, K. Karahalios, and A. Parameswaran. Effortless data exploration with zenvisage: an expressive and interactive visual analytics system. *Proceedings of the VLDB Endowment*, 10(4):457–468, 2016.
- [24] T. Siddiqui, J. Lee, A. Kim, E. Xue, X. Yu, S. Zou, L. Guo, C. Liu, C. Wang, K. Karahalios, et al. Fast-forwarding to desired visualizations with zenvisage. In *CIDR*, 2017.
- [25] T. Siddiqui, P. Luh, Z. Wang, K. Karahalios, and A. Parameswaran. Shapesearch: flexible pattern-based querying of trend line visualizations. *Proceedings of the VLDB Endowment*, 11(12):1962–1965, 2018.
- [26] T. Siddiqui, P. Luh, Z. Wang, K. Karahalios, and A. Parameswaran. Shapesearch: A flexible and efficient system for shape-based exploration of trendlines. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, pages 51–65, 2020.
- [27] M. Wattenberg. Sketching a graph to query a time-series database. In *CHI’01 Extended Abstracts on Human factors in Computing Systems*, pages 381–382. ACM, 2001.