






# HiRegEx: Interactive Visual Query and Exploration of Multivariate Hierarchical Data

Guozheng Li , Haotian Mi , Chi Harold Liu , Takayuki Itoh , and Guoren Wang 

**Abstract**—When using exploratory visual analysis to examine multivariate hierarchical data, users often need to query data to narrow down the scope of analysis. However, formulating effective query expressions remains a challenge for multivariate hierarchical data, particularly when datasets become very large. To address this issue, we develop a declarative grammar, HiRegEx (Hierarchical data Regular Expression), for querying and exploring multivariate hierarchical data. Rooted in the extended multi-level task topology framework for tree visualizations (e-MLTT), HiRegEx delineates three query targets (node, path, and subtree) and two aspects for querying these targets (features and positions), and uses operators developed based on classical regular expressions for query construction. Based on the HiRegEx grammar, we develop an exploratory framework for querying and exploring multivariate hierarchical data and integrate it into the TreeQueryER prototype system. The exploratory framework includes three major components: top-down pattern specification, bottom-up data-driven inquiry, and context-creation data overview. We validate the expressiveness of HiRegEx with the tasks from the e-MLTT framework and showcase the utility and effectiveness of TreeQueryER system through a case study involving expert users in the analysis of a citation tree dataset.

**Index Terms**—Multivariate hierarchical data, declarative grammar, visual query

## 1 INTRODUCTION

Multivariate hierarchical data are ubiquitous in real-world applications and can be found in datasets like citation trees of publications [58], reposting trees in social media [16, 47, 69], and hierarchical tabular data [33–35]. One technique often used to analyze such hierarchical data is exploratory visual analysis (EVA), which involves examining data, extracting patterns, gaining insights, and refining hypotheses [3]. Visual analytics techniques, such as visual encoding and querying, can facilitate an EVA process by enabling rapid specification of data visualizations and transformations [3, 55]. While significant progress has been made in the visual encoding of hierarchical data visualizations [51], visual querying remains a challenge in the EVA of multivariate hierarchical datasets. Specifically, the unpredictable characteristic of an EVA process indicates that users often lack a clear idea of query targets and must continuously try different queries to reach a goal. However, the complexity of multivariate hierarchical data, in terms of topological structures and node attributes, makes constructing practical query expressions time-consuming and error-prone.

Take an example of analyzing a citation tree dataset with EVA. Each node in the citation tree represents a publication with multiple attributes, such as “topics” and “authors”, while the links among nodes signify their reference relationships. In this scenario, the tasks of a researcher include capturing important publications, identifying their various patterns, and comprehending the development of research trends. The researcher needs to formulate diverse query expressions frequently to accomplish tasks or validate hypotheses. For instance, the researcher may query publications within the “graph” topic from the past three years that have been cited by more than five publications within the “immersive” research topic to understand the recent intersection of disciplinary directions. The query expression is related to both node attributes, like topics, and topological structures, such as the number of children on specific topics. In particular, the topics and parameters in the above query might be changed frequently during EVA.

Existing techniques allow users to query multivariate hierarchical

data by programming and interactive filtering. The programming approach lets users craft general low-level imperative codes, which can be easily verified but requires clearly defined goals and query targets, a significant burden for non-programming users. Although some declarative languages for hierarchical data queries [4, 32, 56] are less challenging, they are usually domain- or task-specific (e.g., the syntax tree in natural language processing [32], or graph query languages [2, 18, 48, 61]) and lack adequate support for diverse visual analytics tasks [42]. With the interactive filtering approach, on the other hand, users can see specific topology structures and build dynamic queries on multiple node attributes. However, constructing flexible and diverse queries requires continuous data filtering and checking activities, a very time-consuming task, especially when the dataset is large. Therefore, we seek an efficient and expressive query approach that can be integrated into the EVA process for multivariate hierarchical data.

This work presents HiRegEx (**H**ierarchical data **R**egular **E**xpression), a novel declarative grammar for querying multivariate hierarchical data. HiRegEx builds upon a tree-specific extension to the multi-level task topology framework (e-MLTT) [42], a dataset with a collection of 213 tasks from existing studies. We develop a classification of three distinct targets (node, path, and subtree) from the task abstraction framework and then define two aspects of target querying, features, and positions. Specifically, HiRegEx sees a node as the elementary unit that is used to set constraints on various attributes and borrows operators from regular expressions to specify parent-child relations [52]. Some new operators are also introduced to specify sibling relations [36].

Furthermore, we introduce a query-based framework to support the EVA of multivariate hierarchical data based on the HiRegEx declarative grammar. The framework is designed to assist users in conducting exploratory analysis tasks based on the sense-making model [30, 46]. The framework includes top-down pattern specification, bottom-up data-driven inquiry, and context-creation data overview. We implement a prototype system, TreeQueryER, to integrate the exploratory framework. The top-down specification is supported by a visual editor based on HiRegEx, in which users can construct query expressions interactively. To conduct bottom-up inquiries, users can test different query expressions based on a visualized target dataset. The context-creation overview shows the hierarchical data collection through the dimension reduction technique and incorporates a graph contrastive learning model for the computation of similarities among distinct data.

We validate our techniques in two ways. First, we demonstrate the expressiveness of HiRegEx grammar based on the e-MLTT framework for hierarchical data. The results indicate that the HiRegEx grammar effectively supports 174 tasks of the entire task collection, consisting of 213 tasks. Second, we validate the utility of the TreeQueryER

- Guozheng Li, Haotian Mi, Chi Harold Liu, and Guoren Wang are with Beijing Institute of Technology. Chi Harold Liu is the corresponding author. E-mail: {guozheng.li, haotian.mi, chihliu, wanggr}@bit.edu.cn
- Takayuki Itoh is with Ochanomizu University. E-mail: itot@is.ocha.ac.jp.

Manuscript received xx xxx. 201x; accepted xx xxx. 201x. Date of Publication xx xxx. 201x; date of current version xx xxx. 201x. For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org. Digital Object Identifier: xx.xxx/TVCG.201x.xxxxxxx

prototype system in a case study involving a citation tree dataset. The results show that users can interactively construct query expressions for various tasks and confirm the capacity of the underlying framework to facilitate the EVA process of multivariate hierarchical data.

In summary, the main contributions of this paper are as follows:

- The HiRegEx declarative grammar to facilitate diverse tasks in querying multivariate hierarchical data;
- A query-based exploratory framework for exploring multivariate hierarchical data, supporting top-down, bottom-up, and context-creation data query modes;
- The TreeQueryER prototype system that integrates the exploratory framework for multivariate hierarchical data and is validated through a case study involving a citation tree dataset. Its source code is available at <https://github.com/bitvis2021/HiRegEx>.

## 2 RELATED WORK

In this section, we review literature in the areas of query languages of hierarchical data and interactive visual query techniques.

### 2.1 Query Languages of Hierarchical Data

Hierarchical data query has been extensively studied due to its wide applications in various domains. A typical application scenario is for querying syntax trees, which are hierarchical representations of different syntactic categories of a sentence, in natural language processing. Tregex [32] is a query language for syntax trees based on character matching and supports the specification of relative relationships between nodes. However, dealing with a large-scale syntax tree with this approach involves continuously integrative processes, which can be computationally complex and produce results with poor readability. In addition, unlike the traditional usage scenarios of hierarchical data queries, the design of Tregex emphasizes the order between sibling nodes, imposing some challenges for many users. Moreover, many databases contain data arranged in a hierarchical structure. Jaql [4] is a declarative scripting language most commonly used for querying and processing JSON data, a classic format used for hierarchical data. The hierarchical query language (HQL) [56] is a language for querying a hierarchical database. Based on HQL, users can specify transactions against a database with a hierarchical structure. Both Jaql and HQL are designed for users to set the attribute value for querying targets, such as single node attributes or aggregated node attributes. However, these techniques do not support the specification of the topological structure.

Given a tree as a connected acyclic undirected graph, querying languages for general graphs can be used in most situations. Cruz et al. [13] proposed a declarative language called G to support regular expressions in specifying the path between any two nodes in a graph. G+ [15] extends G with a summary graph to restructure the results obtained by a query statement. Furthermore, GraphLog [12] is a query language for hypertext based on G+, and adds negation and unifies the concept of a query graph. These graph query languages [12, 13, 15] use regular expressions to define the path between nodes and allow the concise specification of relationships between two non-adjacent nodes within a large-scale graph. In addition to the languages defined over a simple graph model, industry users also adopt other graph database query grammars, such as Gremlin [48], Cypher [18], PGQL [61], and G-CORE [2]. These languages have good expressiveness and support the specification of topological structures and node attributes. However, existing graph query languages cannot specify the sibling or same-level relation expressively, an essential feature of hierarchical data. In addition, most of these imperative query languages are not intuitive for users to complete query tasks for multivariate hierarchical data.

Despite the availability of many query languages for filtering hierarchies across various scenarios, they often lack the flexibility to support users in querying multivariate hierarchical data for EVA.

### 2.2 Interactive Visual Query Approach

In visualization research, efforts on interactive visual queries involve nearly all data types, such as tabular data, time-series data, event sequences, and graphs. The following discussions emphasized the closely

relevant data types with hierarchies: graph, time series, and event sequence. More specifically, hierarchy is a special graph data with hierarchical relationships and no circles. In addition, both time series and event sequences have a linear structure, similar to the decomposition paths of hierarchies, linked from the root node to the leaf nodes.

**Graph data query.** Many techniques have been developed for visual query of graph data. Methods like Graphite [11], VOGUE [5], and Visage [45], for example, allow users to construct query patterns for graph data intuitively and present the matching subgraphs of a large property graph. With Vertigo [14], users can build and recommend graph queries and delve into their findings within multi-layer networks. In comparison, VIGOR [44] emphasizes the efficient summarization of subgraph queries by grouping results according to node characteristics and structural similarity. VIMO [59] improves the intuitiveness and flexibility of graph queries by allowing users to sketch the targets in the query interface and enable users to define structural similarity constraints directly. Furthermore, to reduce the difficulty in querying graphs, researchers have paid more attention to data-driven graph query approaches and developed methods like selecting high-quality patterns (small subgraphs) automatically from the underlying graph database [6, 23, 67]. All these methods focus on how to specify the graph topology accurately by constructing subgraph patterns. However, different from graph data, hierarchical data consists of both parent-child and sibling relations. Because of the lack of the capability to distinguish these two types of relations, existing graph query methods fall short on tasks of multivariate hierarchical data.

**Time series query.** The most straightforward approach for the visual query of time series is to interactively select a partial time series as the query. Users can brush the partial time series of interest through timeboxes [8, 9, 21, 22] on a two-dimensional space. The extent of the timebox on the time axis ( $x$ -axis) specifies the time of interest, while the extent on the value axis ( $y$ -axis) specifies a constraint on the range of values of interest within the time. These methods allow users to draw multiple timeboxes interactively to query a significant amount of time series data simultaneously. Querying time series using timeboxes requires users to understand the data clearly. To solve this problem, QuerySketch [62], QueryLines [49], and Qetch [41] query time series data using sketching. Users can directly draw on the visualized time series to find similar instances, and Qetch [41] combines freehand drawing and regular expressions to query time series data. In addition to sketching patterns, COQUITO [27] provides a visual interface to define cohorts with temporal constraints through intuitive drag-and-drop operations, enabling the exploration and analysis of time series data with specific temporal patterns.

Data-driven approaches play an essential role in data query when users cannot precisely articulate targets. Peax [31] is the first deep-learning-based approach for an interactive visual pattern search in sequential data. It uses a convolutional autoencoder to capture more visual details of complex patterns. Zenvisage++ [30], which finds that sketch-based methods are not always efficient, develops a taxonomy of visual query system capabilities, including a top-down pattern search (translating a pattern “in-the-head” into a visual query), bottom-up data-driven inquiries (querying or recommending based on data), and context-creation (navigating across different collections of visualizations). Our effort to develop an exploratory framework for multivariate hierarchical data is partially inspired by this taxonomy.

**Event sequence query.** For the interactive visual query of event sequence data, (s)queries [68] is a multi-attribute event sequence data query method based on regular expressions. Users can interactively construct their query pattern by adding various constraints on a node-link diagram. EventPad [10] also uses regular expressions to query event sequence data. The difference between these two techniques is that EventPad defines query results as abstract behavior rather than for data extraction, and supports the aggregation, alignment, and combination of event sequence data. VESPa 2.0 [28] facilitates the discovery and validation of movement sequence patterns through interactive exploration and querying, supporting a bottom-up and data-driven exploration approach. Beyond the query task, MAQUI [29] offers a novel approach for recursive event sequence exploration, combining querying and pattern

1	<b>Target:=</b> (Tree   Path   Node) - EC	
2	<b>Tree:=</b> Node(Branch,...)	<b>EC:=</b> (EC <sub>a</sub>  EC <sub>e</sub> )EC
3	<b>Branch:=</b> [<Path> <sub>1</sub> <sup>(0..*)</sup> Branch,...] EPT	<b>EC<sub>e</sub>:=</b> ∃(Node) <sup>(0..*)</sup>   EPT
4	<b>Path:=</b> (Node !Node) <sup>(0..*)</sup> Path	<b>EC<sub>a</sub>:=</b> ∀(Node)   EPT
5	<b>Node:=</b> CustomNode   \$   ^   ·	<b>EPT:=</b> empty

"EC": Element Composition	" ": or	".": wildcard
"0..*": repetition	"∃": there exists	"^": root
"!": not	"∀": for all	"\$": leaf

Fig. 1: The formal specification of the HiRegEx declarative grammar. The first row introduces the overall structure of the HiRegEx specification. Rows 2 to 5 below present various query targets on the left, grounded in the e-MLTT framework, while the right side details element compositions that enable users to specify how these query targets are structured.

mining to help users analyze and understand complex event sequence data more effectively. Similar to event sequence data, orders exist between data items in the hierarchical data because of the parent-child relationships. Therefore, hierarchical data can be decomposed into a collection of multiple sequences, and the specification of hierarchical data can borrow ideas from regular expressions.

Visual queries play a crucial role in enabling EVA and offer flexible filtering capabilities. While visual query techniques for data types with relational and linear structures have been extensively studied, there remains a gap in addressing visual query for hierarchical data.

### 3 HIERARCHICAL DATA QUERY TASK SPACE

In this section, we define the task space related to querying multivariate hierarchical data. Pandey *et al.* comprehensively summarized analytical tasks for tree visualizations and introduced the e-MLTT framework [42], which enhances the specificity of task abstraction definitions tailored to tree visualizations. The e-MLTT framework decomposes tree visualization tasks along two dimensions: targets and actions (see Fig. 4 in Pandey’s study [42]). The target is defined as the object related to the tree visualization tasks, while the action is the operation users performed to accomplish the tasks. Our task space definition is also based on these two dimensions of the e-MLTT framework, but only keeps the query-related analytical tasks in the framework. In this section, we detail the targets and actions in our task space, using the citation tree as an example to explain relevant concepts and techniques.

#### 3.1 Query Targets

A query target is defined as the object that the query expression aims to match. The traditional multi-level task typology (MLTT) framework proposed by Brehmer and Munzner [7] categorizes the target of tasks into “topology” and “attributes”. The e-MLTT framework [42] further extends MLTT to include more specificity to support tree-specific tasks. It divides the specific targets of “topology” for tree visualization tasks into four categories: tree, subtree, path, and node, and the “attributes” into two categories: *node attributes* and *link attributes*. The query targets of the HiRegEx grammar fully consider the targets defined in the e-MLTT framework, except for link attributes. This is because links in hierarchical data are typically used to represent parent-child relationships, and the e-MLTT framework [42] also indicates that links are less frequently used for encoding data attributes in tree visualizations, and tasks related to links are rare in tree visualizations.

#### 3.2 Query Actions

A query action is an operation users perform to accomplish a task. The e-MLTT framework classifies query actions into three levels (high, mid, and low). Our focus here is on the mid-level actions—“search”, a type of operations a user must perform to find targets. Under e-MLTT, the “search” action is further divided into four subtypes—“lookup”, “locate”, “browse”, and “explore”—based on whether the knowledge of targets and locations is available. In particular, search in the “target

known” subtype refers to tasks with explicit knowledge of a target’s identity. For example, when users want to find some popular papers in certain directions that have attracted much attention, they can define a query expression to find publications with more than five iterative citations within a short period. Conversely, the “location known” search concerns tasks that clearly describe a target’s position within the tree. For instance, to understand the ongoing development of a technique, users can look for the inaugural paper that introduced the technique and then trace a subsequent path in the citation tree. As a result, all papers along this path are likely to be related to this technique.

## 4 QUERY GRAMMAR DESIGN

In this section, we first analyze the design requirements for query grammar based on the task space for querying multivariate hierarchical data and then present the detailed grammar specifications of HiRegEx.

### 4.1 Design Requirements

Existing studies [20,26,43,50] have significantly influenced and shaped the following two requirements for a declarative grammar designed for querying multivariate hierarchical data.

**R1: Expressiveness.** A visualization grammar should have the capability to articulate the entire design space [17,43]. Hierarchical data is a generic data type with broad applications in areas such as finance, biology, and computer science [38,39]. This implies that the grammar should support versatile data query functions applicable to various scenarios involving hierarchical data. The e-MLTT framework [42] has summarized 213 analysis tasks related to the hierarchical data from the literature. As mentioned in Sec. 3, query targets in the task space consist of topology and attribute, while query actions are divided into four types according to whether “target” and “location” are known. Consequently, the declarative query grammar should empower users to obtain different kinds of query targets based on various query actions, aligning with the requirements of most tree visualization tasks defined by the task abstraction model.

**R2: Conciseness.** EVA is a typical usage scenario of query languages for multivariate hierarchical data. In an EVA process, analysts often begin with an unclear goal and refine their objectives as they explore the data [64]. This process necessitates frequent data queries for effective exploration [19,54]. To improve the efficiency of EVA, analysts prefer query expressions that can be easily constructed, such as higher-level grammar like ggplot2 [63] and grammar-based systems like Tableau [57]. In addition, analysts need to understand the previous exploratory analysis process to make informed decisions about further exploration. Therefore, it is crucial to ensure the conciseness of the grammar, making it easy for users to understand and construct. The query expressions should align with users’ cognitive understanding of analysis tasks to enhance readability and construction efficiency.

### 4.2 Grammar Specification

Existing research has extensively explored semi-structured data querying [2,4,18,32,48,56]. However, these techniques either lack the expressive power to cover all tree analysis tasks or produce non-intuitive expressions (R1, R2), making them unsuitable for EVA scenarios. Since a tree can be viewed as a composition of multiple paths, and regular expressions provide an intuitive way to represent paths, we extend regular expressions to support various tasks in multivariate hierarchical data querying (R1). We present HiRegEx, a declarative grammar for querying multivariate hierarchical data. Regular expression uses ordinary characters (“a” to “z”) and special operators to define text sequences conforming to a pattern, commonly employed for text querying and replacement. Similarly, HiRegEx is designed to search multivariate hierarchical data by specifying patterns through nodes with attribute constraints and special operators. The operators in HiRegEx are borrowed from the regular expression and further extended according to the characteristics of the hierarchical data. The simplicity of operators ensures the conciseness of the query expression (R2).

The specifications of HiRegEx grammar (see Fig. 1) support all query targets in the e-MLTT framework, node, path, subtree, and tree, as explained in Sec. 3.1 (R1). Additionally, we also introduce element

compositions to allow users to specify the compositions of various query targets. Users need to specify different query patterns for various targets. For node queries, we introduce various constraints for node attributes. Querying paths necessitates the definition of parent-child relationships based on node specifications. Moreover, querying subtrees and trees entails specifying the complete topological structure, including parent-child and sibling relationships.

In addition to query targets, our query grammar is also designed to support different query actions defined in the task space (**R1**). For “*target known*” queries, the grammar assists users in specifying the inherent features of the target, encompassing topology structures, node attributes, and element compositions. Furthermore, to facilitate “*location known*” queries, the grammar enables users to precisely define the target’s location within the entire hierarchical dataset, thereby representing its positional features.

(1) **Node Query.** A node is the elementary unit of hierarchical data and has multiple attributes. The node expression within HiRegEx enables users to specify patterns for node attributes. We divide the node attributes into two categories: inherent and additional.

- **Inherent attributes**, often quantitative data like depth, help users in specifying target positions (location-known query). In addition, each node can also be seen as the root of a subtree. Therefore, HiRegEx allows users to specify the tree-specific attributes in the node, including size, height, and width. At the same time, the parent-child relationship defines a strict order between nodes, and each node is in a unique path starting from the root node. HiRegEx facilitates users in specifying node attributes according to their related nodes, with operators “&” and “#” for relative and absolute positions across levels, respectively. For instance, (*degree* = &-1) signifies that the degree of a node is the same as that of its parent, and (*degree* = #1) implies equivalence to the degree of the root node.
- **Additional attributes** are often quantitative and categorical node features, such as citation number and authors in the citation tree (target-known query). To articulate the constraints of node attributes, we introduce several predicate operators in the grammar specification, including “>”, “≥”, “<”, “≤”, “=” for quantitative data, and “∈” for categorical data. Nodes in multivariate hierarchical data are selected if their attributes satisfy all defined constraints within the expression. Query results of a node expression consist of individual nodes within multivariate hierarchical data. We also pre-define three special nodes: the wildcard node (●), root node (∧), and leaf node (⊙). The following defines the formal specification of *Node*:

$$Node := CustomNode | \$ | \wedge | \bullet$$

(2) **Path Query.** Based on the node specifications, we define the path-related syntax within hierarchical data by detailing the parent-child relationships. A path in hierarchical data exhibits a linear structure, wherein node sequence signifies parent-child relationships. Specifically, the preceding node serves as the parent of the subsequent nodes.

To query paths in hierarchical data, we adopt operators from regular expressions, including “or” (|), “not” (!), and “repetition” ( $\{min, max\}$ ). The repetition operator empowers users to specify an exact number or a range. For instance,  $node^2$  signifies that the *node* pattern repeats twice, and  $node^{\{2,5\}}$  indicates a repetition range from two to five. Users can specify only the minimum or maximum number of repetitions, such as  $node^{\{2,\}}$  matching at least two times and  $node^{\{,5\}}$  matching up to five times. HiRegEx consistently employs a lazy matching strategy to query paths according to the expression, which means the path that first satisfies the query expression is the result, concluding the matching process. An exception arises when the maximum number of repetitions is unspecified. In this case, the matching process concludes only if no nodes can be matched. The formal specification of *Path* is presented below. For instance, the query expression (*authors*=“Ben Shneiderman”)<sup>{3,}</sup> can search for a citation path that indicates Ben Shneiderman’s continuous iterative studies in a specific research topic. The formal specification of *Path* is presented below:

$$Path := (Node | !Node)^{\{min, max\}} Path$$

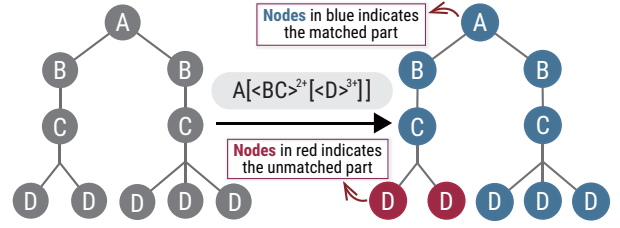


Fig. 2: The explanations of *Branch* operator in the HiRegEx. The nodes in blue indicate the matched part with the HiRegEx expression, while the nodes in red indicate the unmatched part.

(3) **Subtree/Tree Query.** The above path query expressions facilitate users in specifying parent-child relationships between nodes, yet the topological structures of hierarchical data also require the determination of sibling relations [36]. This section delves into the specification of sibling relationships that do not dictate the node sequence.

The relations between nodes’ siblings are implied by relations between paths. To specify the sibling relation, we introduce the *Branch* operator, which merges multiple paths and allows users to specify the repetition number of paths. Notably, an inner path within a *Branch* can be followed by another *Branch*, aligning with the recursive characteristics of hierarchical data. The formal specification of the *Branch* operator is expressed as:

$$Branch := [Path_1^{\{min, max\}} Branch, \dots, Path_n^{\{min, max\}} Branch] \quad (1)$$

When a path is connected with a *Branch*, the query results can only be determined if the matching process is finished. An illustrative example in Fig. 2 displays two paths (B-C) beneath node A. However, the nodes in red do not meet the requirement. More specifically, the repetition number should exceed three according to the expression. Only one path under node A satisfies the condition, falling short of the expression’s demand for more than two paths. Finally, the query expression cannot be matched with the hierarchical data. Based on the *Branch* operator, the formal specification of *Subtree* is shown below. For example, the query expression for citation tree (*authors*=“Ben Shneiderman”)[*(citation*≥200)<sup>{3,}</sup>] can search for a Shneiderman’s paper which has inspired more than three highly cited papers.

$$Subtree := Node Branch$$

(4) **Element Composition.** We find that specific query tasks necessitate consideration of the comprehensive element compositions of the query target (target-known query). For instance, analysts may seek to identify influential papers by querying citation trees published in 2019, comprising more than ten highly cited papers. To accommodate such query tasks, we introduce the *Element Composition* (*EC*) operator in HiRegEx, empowering analysts to specify compositions as an additional aspect of the query target. The specifications of *EC* can be categorized into two types: for all ( $\forall$ , denoted as  $EC_a$ ) and there exists ( $\exists$ , denoted as  $EC_e$ ). The formal specification of *EC* is presented below.

$$EC := (EC_a | EC_e) EC$$

$$EC_a := \forall \langle Path \rangle^{\{min, max\}} | EPT$$

$$EC_e := \exists \langle Path \rangle^{\{min, max\}} | EPT$$

Note that the *repetition* operator ( $min, max$ ) in *EC* refers to the occurrence number of *Path*, distinguishing it from the repetition ( $min, max$ ) used with the *Node*. With the *EC*, we can articulate the task above through the following expressions.

$$(year = 2019)[((\bullet)^{\{0,\}})^{\{0,\}}] - \exists (citation \geq 200)^{\{10,\}}$$

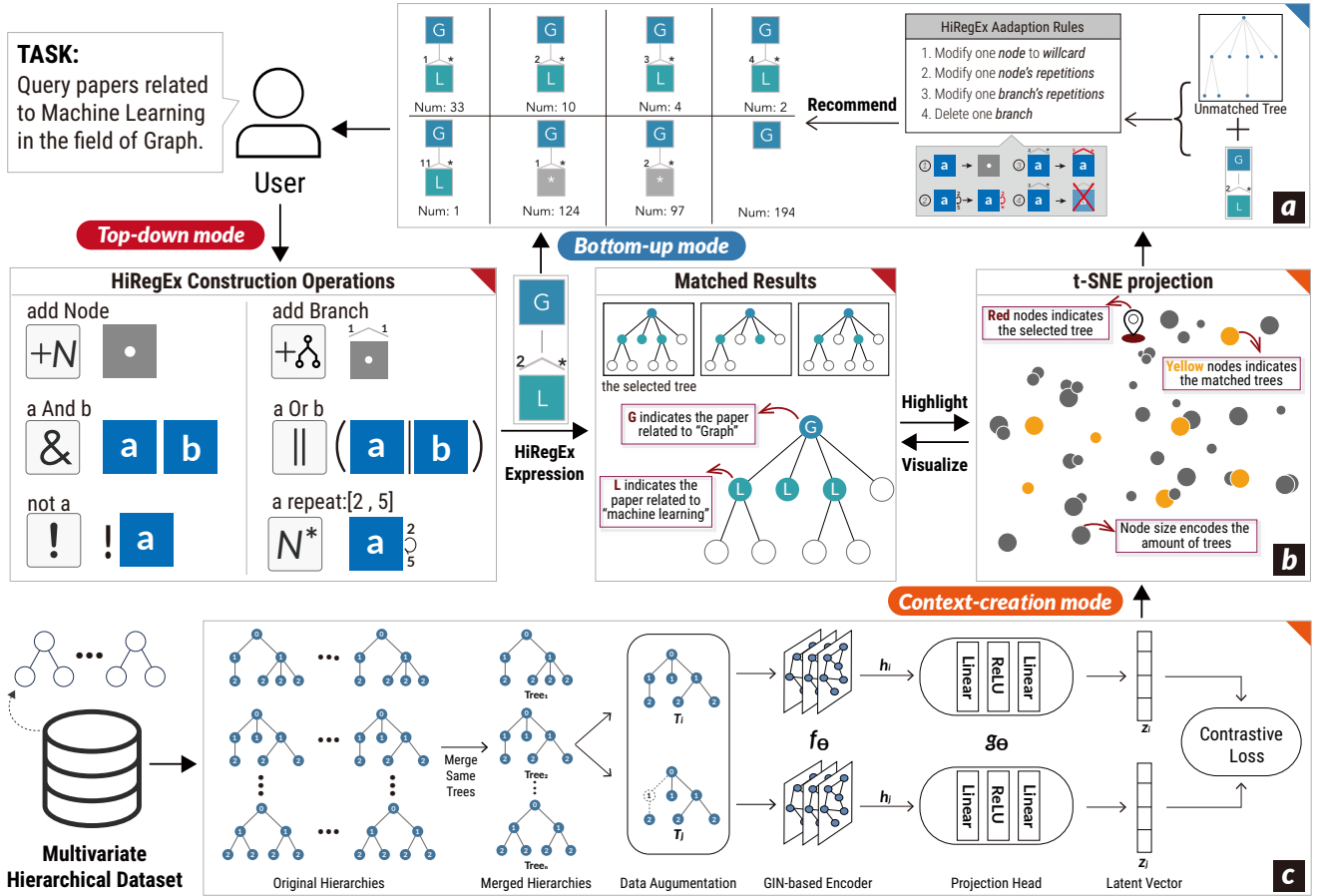


Fig. 3: The exploratory framework for querying multivariate hierarchical data comprises three modes: *top-down*, *bottom-up*, and *context-creation*. The top-down mode starts from a clear query task. Users construct the corresponding query expression through direct manipulations interactively. The bottom-up mode recommends related query expressions based on the initial expression and the multivariate hierarchical data collection. The context-creation mode offers users an overview of the entire hierarchical data collection. Modules associated with the top-down, bottom-up, and context creation modes in the framework are denoted by red, orange, and blue triangles.

With all the previously defined operators, we can formally specify the query target, denoted as *Target*, as follows.

$$Target := (Subtree|Path|Node)-EC$$

## 5 QUERY-BASED EXPLORATION FRAMEWORK

EVA constitutes an iterative process involving data presentation and interactive queries [3], aligning with the principles outlined in the visual information-seeking mantra [53]. This holds for multivariate hierarchical data as well. This section presents a query-based exploratory framework tailored for multivariate hierarchical data. The exploratory framework, rooted in the visual query sense-making model [30], is designed to explore multivariate hierarchical data comprehensively. Illustrated in Fig. 3, the framework provides analysts with three distinct modes tailored to address various requirements: top-down mode, bottom-up mode, and context-creation mode.

(1) **Top-down mode** is designed for a goal-oriented query process where users have a clear understanding of the target pattern and aim to search for data instances exhibiting the pattern. As explained in Sec. 4, HiRegEx is designed based on the task space for querying multivariate hierarchical data. Specifically, users can convert requirements into query expressions without low-level programming. After executing a query expression, users can obtain target data. The challenge in the top-down mode lies in translating the desired patterns by users into executable query expressions because constructing HiRegEx expressions in a textual format requires a steep learning curve. More specifically, users need to memorize operators and parameters in the specification, and text-based expressions lack cognitive consistency with the query

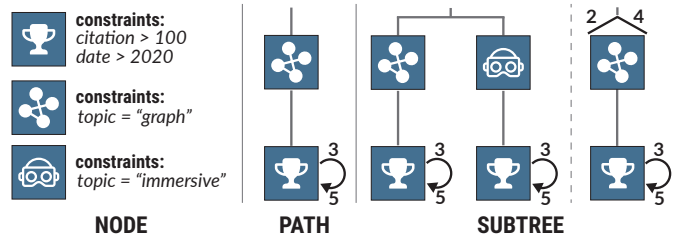


Fig. 4: Three visual operators in the query expression, Node, Path, and Branch, which are the basic components of query expressions.

results of hierarchical data. To address this limitation, we propose a visual operator for each component in a query expression, including the Node, Path, and Branch shown in Fig. 4.

(2) **Bottom-up mode** is a data-driven process enabling users to identify something of interest from the data collection. In this mode, analysts are initially unclear about query targets or unable to specify query patterns accurately. They need to determine the query expression or specific parameters according to the data collection. The challenge in this mode lies in generating an appropriate set of stimuli through recommendations that can prompt further data-driven inquiries.

We devise an expression recommendation algorithm that can derive several relevant query patterns based on a pre-determined expression and data collection. The recommendation algorithm checks each item in the data collection based on the initial HiRegEx expression. For

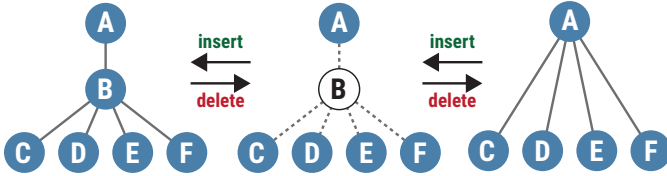


Fig. 5: The *delete* and *insert* operations for computing tree edit distance.

each data item that fails to match the initial query expression, the recommendation algorithm iteratively refines the expression through the following operations until the matching process is finished (Fig. 3a): (1) changing a node with several constraints to a wildcard; (2) modifying the repetition of one node; (3) modifying the repetition of one path; (4) deleting a path from a branch. These four operations are prioritized in descending order based on their impact on the matching results. After traversing all entries, the algorithm merges the adjusted expressions and provides them as recommendation results to users. Details of the algorithm can be found in the supplemental material.

(3) **Context-creation mode** aims to assist users in understanding the data distribution and offering relevant data as context for each query result, thereby assisting users in subsequent exploration. The design challenge of the context-creation mode is constructing an effective overview to reflect the similarities between hierarchical data. The tree edit distance is a typical metric to quantify similarities in tree-structured data. It is defined as the minimum-cost sequence of node operations (*e.g.*, insert) required to transform one tree into another. However, the tree edit distance lacks consistency with the query expression matching algorithm, which traverses hierarchical data from top to bottom, as explained in Sec. 4.2. While tree edit distance might yield a small value between two hierarchical data, their topological structures can vary significantly from top to bottom due to operations that allow inserting or deleting any nodes. For example, Fig. 5 illustrates the insertion and deletion operations between two hierarchical datasets. The edit distance between different hierarchies is small. However, their topological structures differ significantly; one has only one node at the second level, while the other has four nodes.

To address the above limitations, we construct a semantic overview based on the topological structures of hierarchical data using a graph embedding method. Our goal is to learn a low-dimensional representation that captures the structural information of the graph, ensuring that graphs with similar structures are adjacent in a two-dimensional space. To achieve this goal, we first merge hierarchical data with identical topological structures and then apply a graph contrastive learning method (GraphCL) [66] to map the structural information of graphs into high-dimensional vectors. GraphCL employs a contrastive loss function to maximize the consistency between positive pairs in comparison to negative pairs. Fig. 3c illustrates the detailed architecture of the framework. GraphCL augments each hierarchy in the dataset by randomly dropping a node and its subtree to construct positive pairs. To ensure that the hierarchical structure is not significantly changed, the height of the dropped node should be less than or equal to two. Considering the node matching process of HiRegEx is from top to bottom, we design node attributes as *depth* because the changes in the node depth can significantly influence the topological structure. Subsequently, we apply the t-SNE dimensionality reduction algorithm [60] to project the high-dimensional latent vectors of hierarchical data into a two-dimensional space for an overview. From the overview, analysts can understand the similarities between any pairs of hierarchical data and identify the patterns/anomalies.

## 6 TREEQUERYER PROTOTYPE SYSTEM

We have designed and implemented the TreeQueryER prototype system to facilitate the exploratory visual analysis for multivariate hierarchical data based on the HiRegEx grammar.

### 6.1 Design Consideration

**DC1: Reducing the cognitive burden for constructing query expressions based on the HiRegEx specification.** Constructing HiRegEx expressions in a textual format has a steep learning curve. Specifically, users need to memorize the operators and parameters in HiRegEx specifications. The prototype system should enable users to analyze multivariate hierarchical data effectively and efficiently. However, manually writing textual query expressions in textual format contradicts this goal. More specifically, the manual construction process is time-consuming and query expressions are not intuitive. Inspired by various visual query studies for graph data [59], event sequence [10], movement sequence [28], and temporal pattern [27], which support user interaction for constructing visual query patterns intuitively, TreeQueryER also aims to allow users to construct query expressions through direct manipulation and display them in a visual format.

**DC2: Enabling users to achieve the comprehensive analysis of multivariate hierarchical data.** Different users conduct data analysis with varying intentions: some may have explicit goals and tasks, some may lack specific objectives, and some may possess vague goals with a few initial tasks [3], leading to different data analysis requirements. During an analytical process, user interests may refine or evolve as they observe and discover new insights, ultimately seeking the desired information [24, 25, 64]. This process is a key aspect of exploratory analysis. Therefore, it is essential to implement various exploration modes, such as top-down, bottom-up, and context-creation (introduced in Sec. 5), to support diverse analysis needs [3]. These modes should be integrated into the system to realize comprehensive data exploration.

### 6.2 User Interface and Interaction

The user interface of the TreeQueryER system is shown in Fig. 6, and it includes a visual editor panel, a data overview panel, an expression recommendation panel, and a tree visualization panel.

The **data collection overview panel** (Fig. 6a) is tailored to fulfill the requirements of the context-creation exploration mode. This panel demonstrates the distributions of multivariate hierarchical data through a scatter plot. Each node in the scatter plot signifies a sub-collection with the same topological structure. The scatter plot visualization encodes the amount of data in the collection into node size. Distances between nodes indicate similarities in hierarchical data from the topological perspective. The scatter plot highlights query results, providing context for visual query results and aiding users in comprehending distributions across the entire data collection (**DC2**). The TreeQueryER system also enables data filtering by diverse attributes (*e.g.*, size, height, and width) in the distribution panel (Fig. 6b).

The **tree visualization panel** displays the results of visualizing multivariate hierarchical data that match the query expression. We employ the same color to associate nodes with their matched elements in query expressions. Users can click on each node within the tree visualization to inspect detailed node attributes. At the bottom of the tree visualization panel (Fig. 6c), the TreeQueryER system furnishes users with thumbnails of the tree visualizations, encompassing the entire collection matched with the expression.

The **visual editor panel** (Fig. 6d) facilitates the top-down exploration mode (**DC2**). It offers users an interface for constructing HiRegEx expressions through direct manipulations. The design of this panel adheres to the HiRegEx specification outlined in Sec. 4.2. The node corresponds to a rectangle, allowing users to define constraints for multiple attributes. The path corresponds to multiple sequentially connected rectangles, and users can specify the repetition number of nodes. Similarly, the branch consists of multiple paths, and users can specify a repetition number for each path. Components can be dragged and connected to construct a query expression. Users can connect the components to denote parent-child relationships. The visual representation of the expression in the visual editor panel employs the node-link tree visualization, enhancing consistency with users' cognitive understanding of targets and aiding in identifying matching relationships between components and query results (**DC1**).

The **expression recommendation panel** (Fig. 6e) is designed to support bottom-up, data-driven inquiries (**DC2**). The bottom-up mode is

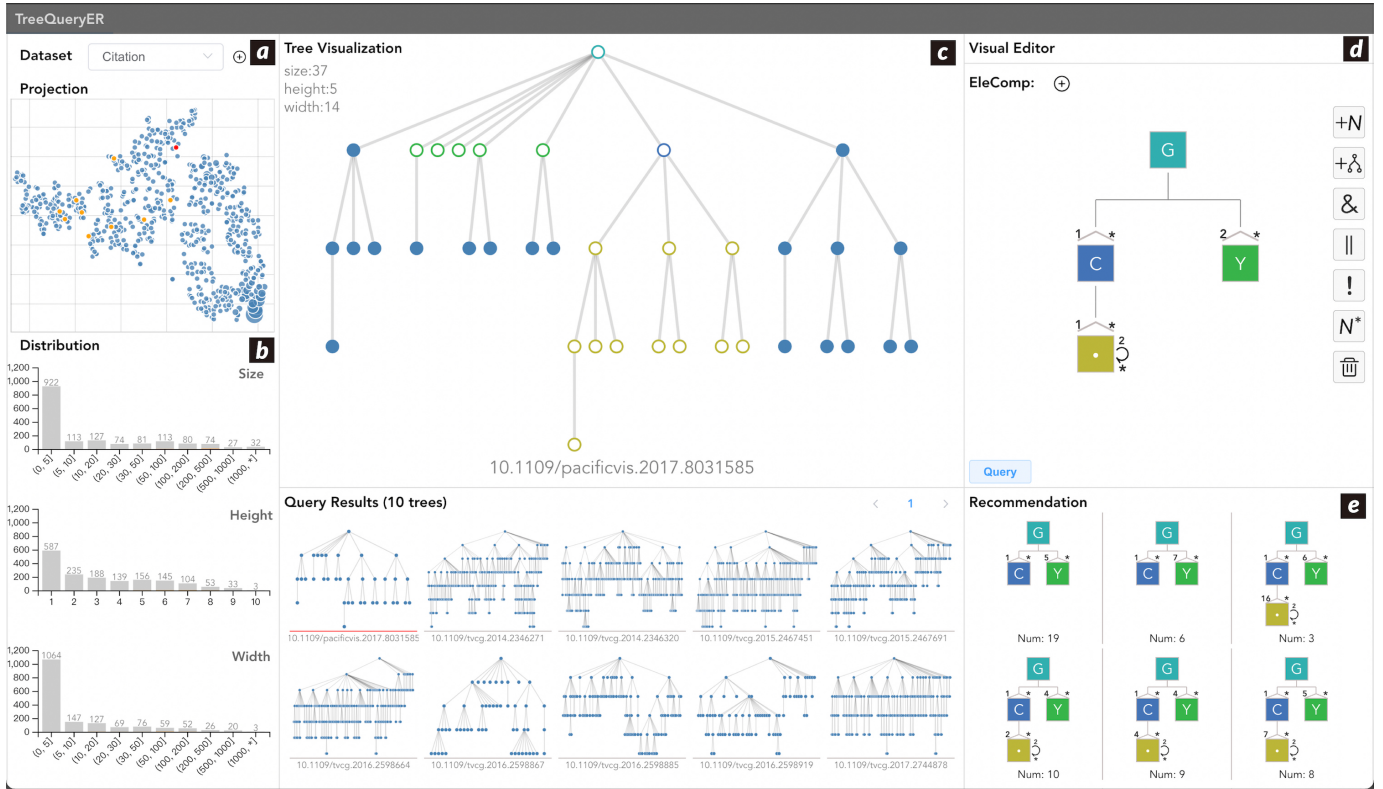


Fig. 6: The user interface of the TreeQueryER prototype system. (a) data collection overview panel. (b) data distribution panel. (c) tree visualization panel. (d) visual editor panel. (e) expression recommendation panel.

part of a browsing-oriented process where users lack a clear target or a specific expression for the target, needing to determine the expression based on the hierarchical data collection. The query expression constructed by users in the visual editor panel comprises multiple parameters that are challenging to determine. For instance, users seeking important papers may set a criterion that their children in the citation tree consist of many highly-cited papers, but determining the threshold is challenging. In such cases, the query expressions constructed by users represent a rough direction rather than a determined pattern for the targets. After users provide an initial query statement, TreeQueryER displays relevant expressions and the quantity of matching hierarchical data in the entire dataset, thereby accelerating the process of obtaining the desired expression and query results (DC1). Users can further select the recommended query expression to visualize their query results in the tree visualization panel (Fig.6c).

In summary, the visual editor panel, expression recommendation panel, and data collection overview panel can support the top-down, bottom-up, and context-creation modes, respectively. The tree visualization panel allows users to understand the query results to refine their query expressions. All the above panels can support users' analysis for multivariate hierarchical data.

## 7 EVALUATION

We validate the above techniques from two aspects. First, we demonstrate the expressiveness of HiRegEx based on the e-MLTT framework. Second, we validate the utility of the TreeQueryER system through a use case on the citation tree dataset in the visualization field.

### 7.1 Performance Evaluation

**Expressiveness of the HiRegEx grammar.** We validated the expressiveness of the HiRegEx grammar based on the e-MLTT framework. More specifically, we utilized the HiRegEx grammar to specify the targets in the 213 tasks underlying the e-MLTT framework. We carefully recorded the number of tasks that could be supported by the HiRegEx grammar, with the results detailed in Fig. 7. Out of the 213 tasks,

HiRegEx can support 174 of them. For each category, we selected a representative task and demonstrated the application of the HiRegEx grammar to define its target. Among the 39 tasks that are not supported, 7 of them involve no query operations. An example of such tasks is “comparison of different subtrees”. The remaining 32 tasks primarily concern computing “extreme” values, such as “find a node having the maximum attribute value of the second layer”. HiRegEx does not consider aggregation operations for query results. This limitation affects its effectiveness in these tasks related to extreme values. However, this shortfall is effectively mitigated by those tools in TreeQueryER that help users to interactively filtering extreme values within the results set. Moreover, the experiment results revealed that HiRegEx proficiently supports tasks related to topology, signifying its capability to represent the structural attributes of hierarchical data accurately.

### 7.2 Case Study

We validate the utility of TreeQueryER by collaborating with two expert users. This section presents a use case on a citation tree dataset and demonstrates how TreeQueryER can help users achieve a comprehensive exploratory visual analysis for the citation tree dataset and get insights into the development and intersection of research topics.

**Dataset.** To assemble a thorough and representative dataset for our study, we crawled a total of 1644 research papers in the field of visualization, covering the period from 2014 to 2020. This dataset encompasses a diverse range of publications from IEEE VIS/TVCG, EuroVis, and PacificVis. The attribute data of each paper include title, authors, publication date, affiliation, country, and citation list/count. We then processed the data collection to construct a citation tree for each paper. In a tree, a paper “b” that cited another paper “a” is a child node of “a”. This structure enabled us to visualize and analyze the intricate network of citations within the field of visualization.

**Expert Users and Tasks.** We invited two experienced visualization researchers, referred to as E1 and E2, to evaluate the TreeQueryER system. Both experts have over five years of research experience in the field. Typically, researchers build citation graphs using tools

Target	Specific Target	Target Attribute	Example Tasks	HiRegEx Results	Num	Available Num	
Topology	Tree	Balance	Find all balanced tree/s.	$\text{Node}[\text{T}_{\text{balanced}} = \text{True}] \{ < (\cdot)^{(0)} > \} > \{ 0 \}$	6	6	
		Height	Find the tree/s of a specific height N.	$\text{Node}[\text{T}_{\text{height}} = N] \{ < (\cdot)^{(0)} > \} > \{ 0 \}$	5	5	
		Size	Find the tree/s of a specific size N.	$\text{Node}[\text{T}_{\text{size}} = N] \{ < (\cdot)^{(0)} > \} > \{ 0 \}$	4	4	
		Fanout	—	—	0	0	
		Order	Find the tree/s in which all nodes are assigned a fixed ordering.	$\text{Node}[] \{ < (\text{Node}[\text{value} > \& - 1])^{(1)} > \} > \{ 1 \}$	1	1	
	Subtree	—	Find the tree/s containing only node 'A'.	$\text{Node}[] \{ < (\cdot)^{(1)} > \} > \{ 1 \}$	2	0	
		Size	Find the subtrees that have minimum N descendants at level 'Y'.	$\text{Node}[\text{T}_{\text{size}} \geq N; \text{level} = Y] \{ < (\cdot)^{(0)} > \} > \{ 0 \}$	20	18	
		Height	Identify the height of the subtree with root node 'A'.	$\text{Node}[\text{label} = A] \{ < (\cdot)^{(0)} > \} > \{ 0 \}$	4	3	
		Depth	Find the deepest subtree inside node 'A'.	—	1	1	
		—	Identify the existence of the subtree with root node 'A'.	$\text{Node}[\text{label} = A] \{ < (\cdot)^{(0)} > \} > \{ 0 \}$	6	3	
	Path	Length	Identify the distance between nodes 'A' and 'B'.	$\text{Node}[\text{label} == A] \{ < (\cdot)^{(0)} \} \text{Node}[\text{label} = B] > \{ 1, 1 \}$	1	1	
		—	Determine whether node 'A' and 'B' are on the same path.	$\text{Node}[\text{label} == A] \{ < (\cdot)^{(0)} \} \text{Node}[\text{label} = B] > \{ 1, 1 \}$	4	3	
		Degree	Find nodes with at least 10 children at level 'Y'.	$\text{Node}[\text{Degree} \geq 10; \text{level} = Y]$	22	21	
	Node	Ancestors	Identify the common ancestors of node 'A' and node 'B'.	$\text{Node}[] \{ < (\cdot)^{(0)} \} \text{Node}[\text{label} = A] > \{ 1, 1 \}, < (\cdot)^{(0)} \} \text{Node}[\text{label} = B] > \{ 1, 1 \}$	14	14	
		Descendants	Identify all remote descendants (8 generations) of node 'A'.	$\text{Node}[\text{label} == A] \{ < (\cdot)^{(8)} > \} > \{ 1 \}$	4	4	
		Ancestor/Descendant	Identify ancestors and descendants of node 'A'.	$\text{Node}[] \{ < (\cdot)^{(1)} \} \text{Node}[\text{label} == A] > \{ 1, 1 \} \{ < (\cdot)^{(0)} > \} > \{ 0 \}$	8	8	
		Depth	Find the depth of node 'A' in the tree	$\text{Node}[] \{ < (\cdot)^{(1)} \} \text{Node}[\text{label} = A] > \{ 1, 1 \}$	5	5	
		Siblings	Identify siblings of node 'A'.	$\text{Node}[] \{ < (\cdot)^{(1)} > \} > \{ 1 \} \{ < \text{Node}[\text{label} = A] > \{ 1, 1 \}, < (\cdot)^{(1)} > \}$	3	3	
		Attribute	Node	Categorical Value	Identify the node with label 'A'.	$\text{Node}[\text{label} == A]$	56
	Quantitative Value			Identify the nodes with a value higher than N.	$\text{Node}[\text{value} > N]$	47	22
SUM					213	174	

Fig. 7: The multi-level task topology of the e-MLTT framework, which consists of three levels, target (high-level), specific target (mid-level), and target attribute (low-level). Each category provides a representative task and the corresponding query expression using the HiRegEx grammar. The cells with a yellow background color indicate task categories that HiRegEx cannot fully achieve.

such as Connected Papers [1] to explore papers’ relationships. This involves manually navigating through citation graphs to understand the evolution. In our case study, E1 and E2 aimed to explore how research in graph visualizations, a traditional topic, intersects with deep learning using TreeQueryER. Before delving into the dataset, they received an introduction to the TreeQueryER system and the specifications of the HiRegEx grammar. Then, they used TreeQueryER for exploratory visual analysis.

**Step 1: Querying papers on graph visualizations and deep learning through the top-down mode.** Initially, the experts were given a brief introduction to the citation tree dataset and various tools provided by TreeQueryER. Then, they decided to explore papers related to deep learning techniques and graph visualization. To identify papers related to graph visualization, they initialized a node in the visual editor, requiring their keyword list to include “graph”. The node, along with the corresponding constraints of attribute values, was labeled as **G** (representing the topic “graph”), as shown in the  $exp_1$  of Fig. 8b. Employing this query expression, the experts retrieved 194 papers.

**Step 2: Querying papers cited by deep learning-related papers using the bottom-up mode.** Many publications may not employ deep learning techniques but still have an impact on this domain. These papers can provide researchers with valuable inspiration but cannot be obtained through the above query expression. To retrieve these papers, the experts constructed another expression to query those cited by more than five papers related to the “deep learning” topic. They added a branch with more than five repetitions in the visual editor, comprising a node labeled **D** with a constraint that the keywords included the term “deep learning”. Connecting this branch with the node **G**, they executed the query expression. The resulting query ( $exp_2$  in Fig. 8) had only one paper due to the relatively strict constraint of being cited by more than five deep learning-related papers. This result hindered a comprehensive understanding of the dataset. At this point, the bottom-up recommendations of the HiRegEx expression became pivotal. From the expression recommendation panel, experts ascertained that the number of papers cited by one or two deep learning-related papers was 33 and 10. Based on these recommendations, they adjusted the parameter of branch repetitions in the query expression to enhance the diversity of results. Fig. 8c presents the query results. The projection panel is updated accordingly, highlighting circles that match the query results in yellow, as depicted in Fig. 8a.

**Step 3: Identifying anomalies through the context-creation mode.** The researchers identified distinct clusters by examining the distributions in the projection panel. From the projection view, they learned that trees located in the lower right corner exhibited small sizes, in contrast to the larger trees represented by circles in the upper left corner, as shown in Fig. 8a. Based on this observation, they selected a cluster of

interest for in-depth analysis. Upon closer inspection, the researchers discovered that certain citation trees consisted of only a few levels, suggesting that the research topics of these papers were outdated and lacked continued exploration by researchers.

**Step 4: Refining the query expression interactively.** Furthermore, they introduced a new branch under the node **G**, complementing the constraint that the descendants must comprise more than five papers after 2019 (denoted as **Y**). These constraints were implemented to ensure that the papers in the query results remain within an active research area. Executing this refined query expression yielded 14 citation trees. Notably, four trees within the results lacked subsequent citations, indicating a limited impact. To measure a paper’s impact, they considered its citations and whether highly-cited papers referenced it. Accordingly, they set a constraint for the “citation” attribute of “node” to be greater than ten (denoted as **C**). Next, they introduced a new branch to restrict papers cited by at least one highly-cited paper ( $exp_6$  in Fig. 8). Note that the parameters in the above expression are adjusted by the recommendation algorithm based on the dataset without users’ manual specifications.

From the expression recommendation panel, experts identified an expression ( $exp_7$  in Fig. 8) corresponding to a specific query result (VIDX [65]), which emerged as the most relevant result for the given expression. This result pertained to “graph” and was cited by influential papers on “deep learning”. After retrieving the paper, they explored the projection view and found that many nearby trees contain another influential paper, “Analyzing the Training Process of Deep Generative Model” [40]. They learned that the outlier detection capability of visual analytics, such as the methods based on Marey’s Graph in the smart factories usage scenario, can be used to help users identify the outlier causing a failed training process.

**Expert Feedback.** We conducted one-on-one 30-minute interviews with the two experts to gather their feedback on our techniques after they finished their hand-on explorations. We encouraged them to freely share their thoughts on our methods as well as their impressions of the overall experience. Both experts expressed positive attitudes towards our method and agreed that the system could improve the efficiency of the paper-searching process. More specifically, they appreciated the design of HiRegEx. “HiRegEx allows me to flexibly define various conditions, such as citation counts and papers’ relations. While each condition alone may not be complex, combining them can be quite intricate. HiRegEx provides an intuitive and user-friendly way to describe these queries (E1)”. HiRegEx extends the basic usage of regular expressions, making it easy for them to understand the rules and use the grammar conveniently. Furthermore, they appreciated the ability to construct their query expressions effectively using the visual editor panel of TreeQueryER. “The visual representation of



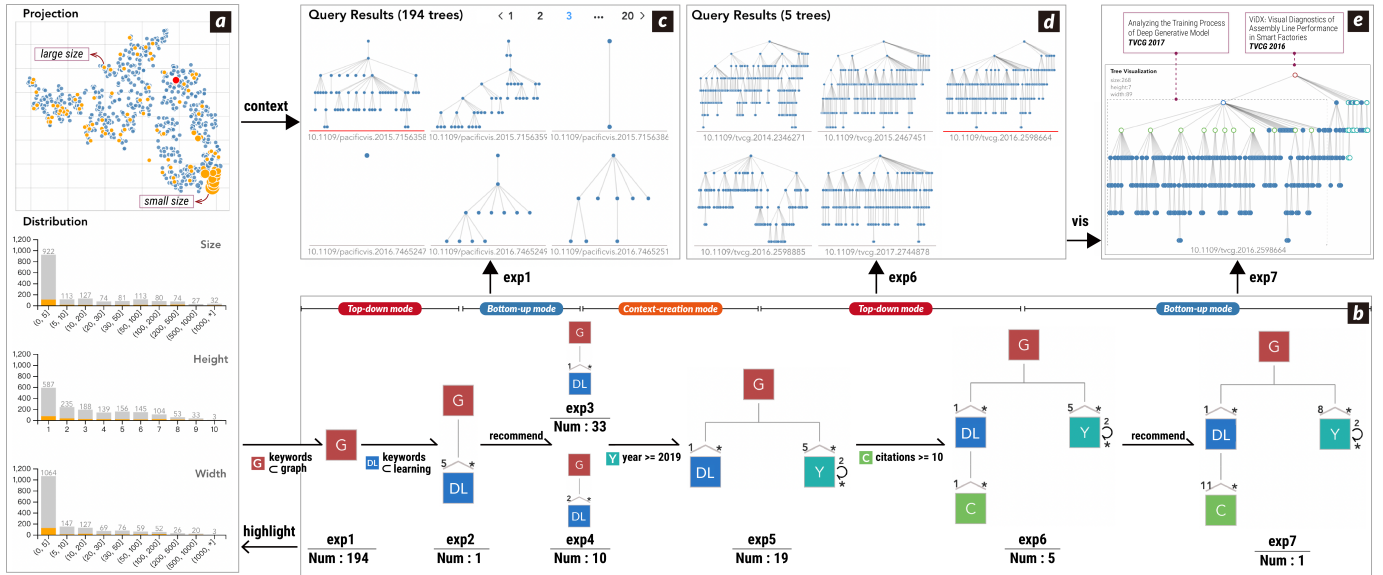


Fig. 8: The use case on the citation tree dataset. (a) The data collection overview panel offers the exploration context and highlights user query results. (b) The process of constructing query expressions comprises seven statements indicating the respective number of query results. The links between statements reveal the attributes and corresponding constraints. In particular, *exp4* and *exp7* are suggested by the recommendation algorithm; (c) The query results of the *exp1* statement; (d) The query results of the *exp6* statement; (e) The query results of the *exp7* statement.

*HiRegEx expressions is intuitive, allowing me to quickly turn an idea into an expression. Through interactive manipulation, I can easily construct the query expressions I need (E1, E2)*. The visual cues enabled them to recognize the characteristics of their query statements. In addition, they found the visual query expression recommendation approach beneficial for refining their queries and selecting appropriate parameters. *The recommendation mechanism always provides me with effective suggestions for the next-step exploration (E2)*. The semantic overview in the projection view can present the context of query results and facilitate the identification of patterns for further investigation.

## 8 DISCUSSION AND FUTURE WORK

**Comparison with Existing Query Grammars.** The techniques competing with HiRegEx for hierarchical data query include Tregex [32], Jaql [4], and HQL [56]. For graph data, our competitors are Gremlin [48], Cypher [18], PGQL [61], and G-CORE [2]. These existing techniques differ from our approach in motivation, expressiveness, available tutorials, and prototype systems. Hence, we did not conduct a quantitative study to evaluate the efficiency of HiRegEx grammar. Most existing studies predominantly focus on node attributes, and often overlooked tree-specific attributes such as size, height, and depth. This limitation hampers their capability to effectively describe large trees, as their query languages necessitate detailed node specifications to define a tree structure accurately. For instance, Tregex [32], tailored for syntax trees, adeptly specifies topological structures and node attributes but cannot define constraints from the perspective of tree compositions. In contrast, Jaql [4] and HQL [56], designed for semi-structured hierarchical database queries, exhibit predetermined topology in their query targets, restricting users from defining query patterns flexibly. Gremlin [48], Cypher [18], PGQL [61], and G-CORE [2] are graph query languages. However, these techniques only partially consider tree-specific characteristics. PGQL [61] and G-CORE [2] can support the specification of the sibling relationships of hierarchical data. However, they do not support queries for a large tree because they need to specify the query targets in a fine-grained manner.

**Expressiveness of the HiRegEx Grammar.** The e-MLTT framework, dedicated to tree visualizations, encompasses 213 analysis tasks. HiRegEx can support 174 of them, as shown in Fig. 7. We have categorized the unsupported analysis tasks into two distinct groups. Tasks falling within the first category necessitate user interactions, aggregation, and computation after querying. These tasks involve actions

related to finding the extreme value or aggregating data, exemplified by queries like “What is the maximum depth of the hierarchy” or “How many files are there in the directory”. Another common example entails determining the least common ancestor of two nodes. In contrast, tasks in the second category are incompatible with query tasks. An illustrative example includes tasks centered around assessing the balance of trees or subtrees in hierarchical data.

**Further Improvements of the Construction Efficiency.** Given the various operators and constraints for the node attributes, the query expression of HiRegEx cannot be completely expressed or constructed in the textual format like the traditional regular expressions. Therefore, we plan to develop a library to integrate the HiRegEx expression with popular programming languages, like Python, to improve the utility of the HiRegEx in textual format. Another future work is to improve user construction efficiency for query expressions by natural language processing techniques. We will explore natural language processing techniques to generate the corresponding query expression and expression description based on Large language models (LLMs) [37] in the future. To improve readability, we will also explore providing a short natural language description of the expressions constructed by users.

## 9 CONCLUSION

To support effective visual query on large, multivariate hierarchical datasets, we proposed HiRegEx, a declarative grammar designed for querying multivariate hierarchical data. HiRegEx borrows the operators from the classical regular expressions and further extends their expressiveness according to the characteristics of multivariate hierarchical data. Based on HiRegEx, we developed a query-based exploratory framework, which consists of *top-down* pattern specification, *bottom-up* data-driven inquiry, and *context-creation* data overview. We implemented a prototype system, TreeQueryER, to integrate our exploratory framework. We validate the expressiveness of HiRegEx based on the e-MLTT framework. We also demonstrate the effectiveness and utility of the exploratory framework and the TreeQueryER system through a case study involving a citation tree dataset.

## ACKNOWLEDGMENTS

This work is supported by National Key R&D Program of China (2021YFB3301500), NSFC (62302038, U2268205), Young Elite Scientists Sponsorship Program by CAST (2023QNRC001), and Tencent Rhino-Bird Focused Research Program.

## REFERENCES

- [1] Connected papers. <https://www.connectedpapers.com>. 8
- [2] R. Angles, M. Arenas, P. Barceló, P. Boncz, G. Fletcher, C. Gutierrez, T. Lindaaker, M. Paradies, S. Plantikow, J. Sequeda, et al. G-CORE: A core for future graph query languages. In *Proc. Int. Conf. Management of Data*, pp. 1421–1432, 2018. doi: 10.1145/3183713.3190654 1, 2, 3, 9
- [3] L. Battle and J. Heer. Characterizing exploratory visual analysis: A literature review and evaluation of analytic provenance in tableau. *Computer Graphics Forum*, 38(3):145–159, 2019. doi: 10.1111/CGF.13678 1, 5, 6
- [4] K. S. Beyer, V. Ercegovac, R. Gemulla, A. Balmin, M. Eltabakh, C.-C. Kanne, F. Ozcan, and E. J. Shekita. Jaql: A scripting language for large scale semistructured data analysis. *Proc. the VLDB Endowment*, 4(12):1272–1283, 2011. doi: 10.14778/3402755.3402761 1, 2, 3, 9
- [5] S. S. Bhowmick, B. Choi, and S. Zhou. VOGUE: Towards a visual interaction-aware graph query processing framework. In *Proc. Conf. Innovative Data Systems Research*. Citeseer, 2013. 2
- [6] S. S. Bhowmick, K. Huang, H. E. Chua, Z. Yuan, B. Choi, and S. Zhou. AURORA: Data-driven construction of visual graph query interfaces for graph databases. In *Proc. Int. ACM Conf. Management of Data (SIGMOD)*, pp. 2689–2692, 2020. doi: 10.1145/3318464.3384681 2
- [7] M. Brehmer and T. Munzner. A multi-level typology of abstract visualization tasks. *IEEE Transactions on Visualization and Computer Graphics*, 19(12):2376–2385, 2013. doi: 10.1109/TVCG.2013.124 3
- [8] P. Buono, A. Aris, C. Plaisant, A. Khella, and B. Shneiderman. Interactive pattern search in time series. *Visualization and Data Analysis*, 5669:175–186, 2005. doi: 10.1117/12.587537 2
- [9] P. Buono and A. L. Simeone. Interactive shape specification for pattern search in time series. In *Proc. Conf. Advanced Visual Interfaces*, pp. 480–481, 2008. doi: 10.1145/1385569.1385666 2
- [10] B. C. Cappers and J. J. van Wijk. Exploring multivariate event sequences using rules, aggregations, and selections. *IEEE Transactions on Visualization and Computer Graphics*, 24(1):532–541, 2018. doi: 10.1109/TVCG.2017.2745278 2, 6
- [11] D. H. Chau, C. Faloutsos, H. Tong, J. I. Hong, B. Gallagher, and T. Eliassi-Rad. Graphite: A visual query system for large graphs. In *Proc. IEEE Int. Conf. Data Mining Workshops*, pp. 963–966, 2008. doi: 10.1109/ICDMW.2008.99 2
- [12] M. P. Consens and A. O. Mendelzon. GraphLog: a visual formalism for real life recursion. In *Proc. ACM SIGACT-SIGMOD-SIGART Symp. Principles of Database Systems*, pp. 404–416, 1990. doi: 10.1145/298514.298591 2
- [13] I. F. Cruz, A. O. Mendelzon, and P. T. Wood. A graphical query language supporting recursion. *Proc. Int. ACM Conf. Management of Data (SIGMOD)*, 16(3):323–330, 1987. doi: 10.1145/38713.38749 2
- [14] E. Cuenca, A. Sallaberry, D. Ienco, and P. Poncelet. VERTIGO: A visual platform for querying and exploring large multilayer networks. *IEEE Transactions on Visualization and Computer Graphics*, 28(3):1634–1647, 2022. doi: 10.1109/TVCG.2021.3067820 2
- [15] I. Curz. G+: Recursive queries without recursion. In *Proc. Int. Conf. Expert Database Systems*, pp. 355–368, 1988. 2
- [16] T. Diefenbach and J. A. Sillince. Formal and informal hierarchy in different types of organization. *Organization studies*, 32(11):1515–1537, 2011. doi: 10.1177/0170840611421254 1
- [17] M. Fowler. *Domain-specific Languages*. Pearson Education, 2010. 3
- [18] N. Francis, A. Green, P. Guagliardo, L. Libkin, T. Lindaaker, V. Marsault, S. Plantikow, M. Rydberg, P. Selmer, and A. Taylor. Cypher: An evolving query language for property graphs. In *Proc. Int. Conf. Management of Data*, pp. 1433–1445, 2018. doi: 10.1145/3183713.3190657 1, 2, 3, 9
- [19] D. Gotz and M. X. Zhou. Characterizing users’ visual analytic activity for insight provenance. In *Proc. IEEE Symp. Visual Analytics Science And Technology (VAST)*, pp. 123–130, 2008. doi: 10.1109/VAST.2008.4677365 3
- [20] P. Hanrahan. VizQL: a language for query, analysis and visualization. In *Proc. Int. ACM Conf. Management of Data (SIGMOD)*, pp. 721–721, 2006. doi: 10.1145/1142473.1142560 3
- [21] H. Hochheiser and B. Shneiderman. Interactive exploration of time series data. In *Proc. The Craft of Information Visualization*, pp. 313–315. Elsevier, 2003. doi: 10.1007/3-540-45650-338 2
- [22] H. Hochheiser and B. Shneiderman. Dynamic query tools for time series data sets: textbox widgets for interactive exploration. *Information Visualization*, 3(1):1–18, 2004. doi: 10.1057/PALGRAVE.IVS.9500061 2
- [23] K. Huang, H. E. Chua, S. S. Bhowmick, B. Choi, and S. Zhou. MIDAS: towards efficient and effective maintenance of canned patterns in visual graph query interfaces. In *Proc. Int. Conf. Management of Data*, pp. 764–776, 2021. doi: 10.1145/3448016.3457251 2
- [24] S. Idreos, O. Papaemmanouil, and S. Chaudhuri. Overview of data exploration techniques. In *Proc. Int. ACM Conf. Management of Data (SIGMOD)*, pp. 277–281, 2015. doi: 10.1145/2723372.2731084 6
- [25] D. A. Keim. Visual exploration of large data sets. *Communications of the ACM*, 44(8):38–44, 2001. doi: 10.1145/381641.381656 6
- [26] Y. Kim and J. Heer. Gemini: A grammar and recommender system for animated transitions in statistical graphics. *IEEE Transactions on Visualization and Computer Graphics*, 27(2):485–494, 2021. doi: 10.1109/TVCG.2020.3030360 3
- [27] J. Krause, A. Perer, and H. Stavropoulos. Supporting iterative cohort construction with visual temporal queries. *IEEE Transactions on Visualization and Computer Graphics*, 22(1):91–100, 2016. doi: 10.1109/TVCG.2015.2467622 2, 6
- [28] R. Krueger, T. Tremel, and D. Thom. VESPa 2.0: data-driven behavior models for visual analytics of movement sequences. In *Proc. Int. Symp. Big Data Visual Analytics (BDVA)*, pp. 1–8, 2017. doi: 10.1109/BDVA.2017.8114626 2, 6
- [29] P.-M. Law, Z. Liu, S. Malik, and R. C. Basole. MAQUI: Interweaving queries and pattern mining for recursive event sequence exploration. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):396–406, 2019. doi: 10.1109/TVCG.2018.2864886 2
- [30] 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*, 26(1):1267–1277, 2020. doi: 10.1109/TVCG.2019.2934666 1, 2, 5
- [31] F. Lekschas, B. Peterson, D. Haehn, E. Ma, N. Gehlenborg, and H. Pfister. Peax: Interactive visual pattern search in sequential data using unsupervised deep representation learning. *Computer Graphics Forum*, 39(3):167–179, 2020. doi: 10.1111/CGF.13971 2
- [32] R. Levy and G. Andrew. Tregex and Tsurgeon: Tools for querying and manipulating tree data structures. In *Proc. Int. Conf. Language Resources and Evaluation (LREC)*, pp. 2231–2234, 2006. 1, 2, 3, 9
- [33] G. Li, P. He, X. Wang, R. Li, C. H. Liu, C. Ou, D. He, and G. Wang. InsightHTable: Insight-driven hierarchical table visualization with reinforcement learning. *IEEE Transactions on Visualization and Computer Graphics*, pp. 1–18, 2024. doi: 10.1109/TVCG.2024.3404454 1
- [34] G. Li, R. Li, Y. Feng, Y. Zhang, Y. Luo, and C. H. Liu. CoInsight: Visual storytelling for hierarchical tables with connected insights. *IEEE Transactions on Visualization and Computer Graphics*, 30(6):3049–3061, 2024. doi: 10.1109/TVCG.2024.3388553 1
- [35] G. Li, R. Li, Z. Wang, C. H. Liu, M. Lu, and G. Wang. HiTailor: Interactive transformation and visualization for hierarchical tabular data. *IEEE Transactions on Visualization and Computer Graphics*, 29(1):139–148, 2023. doi: 10.1109/TVCG.2022.3209354 1
- [36] G. Li, M. Tian, Q. Xu, M. J. McGuffin, and X. Yuan. GoTree: A grammar of tree visualizations. In *Proc. ACM Conf. Human Factors in Computing Systems (CHI)*, pp. 170:1–170:13, 2020. doi: 10.1145/3313831.3376297 1, 4
- [37] G. Li, X. Wang, G. Aodeng, S. Zheng, Y. Zhang, C. Ou, S. Wang, and H. C. Liu. Visualization generation with large language models: An evaluation. *arXiv preprint arXiv:2401.11255*, 2024. doi: 10.48550/arXiv.2401.11255 9
- [38] G. Li and X. Yuan. GoTreeScape: Navigate and explore the tree visualization design space. *IEEE Transactions on Visualization and Computer Graphics*, 29(12):5451–5467, 2023. doi: 10.1109/TVCG.2022.3215070 3
- [39] G. Li, Y. Zhang, Y. Dong, J. Liang, J. Zhang, J. Wang, M. J. McGuffin, and X. Yuan. BarcodeTree: Scalable comparison of multiple hierarchies. *IEEE Transactions on Visualization and Computer Graphics*, 26(1):1022–1032, 2020. doi: 10.1109/TVCG.2019.2934535 3
- [40] M. Liu, J. Shi, K. Cao, J. Zhu, and S. Liu. Analyzing the training processes of deep generative models. *IEEE Transactions on Visualization and Computer Graphics*, 24(1):77–87, 2017. doi: 10.1109/TVCG.2017.2744938 8
- [41] M. Mannino and A. Abouzied. Qetch: Time series querying with expressive sketches. In *Proc. Int. Conf. Management of Data*, pp. 1741–1744, 2018. doi: 10.1145/3183713.3193547 2
- [42] A. Pandey, U. Syeda, C. Shah, J. Guerra-Gomez, and M. Borkin. A state-of-the-art survey of tasks for tree design and evaluation with a curated task dataset. *IEEE Transactions on Visualization and Computer Graphics*,

- 28(10):3563–3584, 2022. doi: 10.1109/tvcg.2021.3064037 1, 3
- [43] D. Park, S. M. Drucker, R. Fernandez, and N. Elmqvist. ATOM: A grammar for unit visualizations. *IEEE Transactions on Visualization and Computer Graphics*, 24(12):3032–3043, 2018. doi: 10.1109/TVCG.2017.2785807 3
- [44] R. Pienta, F. Hohman, A. Endert, A. Tamersoy, K. Roundy, C. Gates, S. Navathe, and D. H. Chau. VIGOR: Interactive visual exploration of graph query results. *IEEE Transactions on Visualization and Computer Graphics*, 24(1):215–225, 2018. doi: 10.1109/TVCG.2017.2744898 2
- [45] R. Pienta, A. Tamersoy, A. Endert, S. Navathe, H. Tong, and D. H. Chau. Visage: Interactive visual graph querying. In *Proc. Int. Conf. Advanced Visual Interfaces*, pp. 272–279, 2016. doi: 10.1145/2909132.2909246 2
- [46] P. Pirolli and S. Card. The sensemaking process and leverage points for analyst technology as identified through cognitive task analysis. In *Proc. Int. Conf. Intelligence Analysis*, vol. 5, pp. 2–4, 2005. 1
- [47] D. Ren, X. Zhang, Z. Wang, J. Li, and X. Yuan. WeiboEvents: A crowd sourcing weibo visual analytic system. In *Proc. IEEE Pacific Visualization Symposium (PacificVis)*, pp. 330–334, 2014. doi: 10.1109/PACIFICVIS.2014.38 1
- [48] M. A. Rodriguez. The gremlin graph traversal machine and language. In *Proc. Symp. Database Programming Languages*, pp. 1–10, 2015. doi: 10.1145/2815072.2815073 1, 2, 3, 9
- [49] K. Ryall, N. Lesh, T. Lanning, D. Leigh, H. Miyashita, and S. Makino. Querylines: approximate query for visual browsing. In *Extended Abstracts on Human Factors in Computing Systems*, pp. 1765–1768, 2005. doi: 10.1145/1056808.1057017 2
- [50] A. Satyanarayan, D. Moritz, K. Wongsuphasawat, and J. Heer. Vega-lite: A grammar of interactive graphics. *IEEE Transactions on Visualization and Computer Graphics*, 23(1):341–350, 2017. doi: 10.1109/TVCG.2016.2599030 3
- [51] H.-J. Schulz. Treevis. net: A tree visualization reference. *IEEE Computer Graphics and Applications*, 31(6):11–15, 2011. doi: 10.1109/MCG.2011.103 1
- [52] H.-J. Schulz, S. Hadlak, and H. Schumann. The design space of implicit hierarchy visualization: A survey. *IEEE Transactions on Visualization and Computer Graphics*, 17(4):393–411, 2011. doi: 10.1109/TVCG.2010.79 1
- [53] B. Shneiderman. The eyes have it: a task by data type taxonomy for information visualizations. In *Proc. IEEE Symp. Visual Languages (VL)*, pp. 336–343, 1996. doi: 10.1109/VL.1996.545307 5
- [54] T. Siddiqui, A. Kim, J. Lee, K. Karahalios, and A. Parameswaran. Effortless data exploration with zenvisage: an expressive and interactive visual analytics system. *Proc. VLDB Endow.*, 10(4):457–468, 2016. doi: 10.14778/3025111.3025126 3
- [55] T. Siddiqui, P. Luh, Z. Wang, K. Karahalios, and A. G. Parameswaran. Expressive querying for accelerating visual analytics. *Communications of the ACM*, 65(7):85–94, 10 pages, 2022. doi: 10.1145/3535337 1
- [56] A. Slingsby, J. Dykes, and J. Wood. Configuring hierarchical layouts to address research questions. *IEEE Transactions on Visualization and Computer Graphics*, 15(6):977–984, 2009. doi: 10.1109/TVCG.2009.128 1, 2, 3, 9
- [57] C. Stolte, D. Tang, and P. Hanrahan. Polaris: a system for query, analysis, and visualization of multidimensional relational databases. *IEEE Transactions on Visualization and Computer Graphics*, 8(1):52–65, 2002. doi: 10.1109/2945.981851 3
- [58] M. Tian, G. Li, and X. Yuan. LitVis: a visual analytics approach for managing and exploring literature. *Journal of Visualization*, 26(6):1445–1458, 2023. doi: 10.1007/S12650-023-00941-3 1
- [59] J. Troidl, S. Warchol, J. Choi, J. Matelsky, N. Dhanyasi, X. Wang, B. Wester, D. Wei, J. W. Lichtman, H. Pfister, and J. Beyer. ViMO - visual analysis of neuronal connectivity motifs. *IEEE Transactions on Visualization and Computer Graphics*, 30(1):748–758, 2024. doi: 10.1109/TVCG.2023.3327388 2, 6
- [60] L. Van der Maaten and G. Hinton. Visualizing data using t-SNE. *Journal of machine learning research*, 9(86):2579–2605, 2008. 6
- [61] O. van Rest, S. Hong, J. Kim, X. Meng, and H. Chafi. PGQL: a property graph query language. In *Proc. Int. Conf. Graph Data Management Experiences and Systems*, pp. 1–6, 2016. doi: 10.1145/2960414.2960421 1, 2, 9
- [62] M. Wattenberg. Sketching a graph to query a time-series database. In *Extended Abstracts on Human factors in Computing Systems*, pp. 381–382, 2001. doi: 10.1145/634067.634292 2
- [63] H. Wickham. A layered grammar of graphics. *Journal of Computational and Graphical Statistics*, 19(1):3–28, 2010. doi: 10.1198/jcgs.2009.07098 3
- [64] K. Wongsuphasawat, D. Moritz, A. Anand, J. Mackinlay, B. Howe, and J. Heer. Voyager: Exploratory analysis via faceted browsing of visualization recommendations. *IEEE Transactions on Visualization and Computer Graphics*, 22(1):649–658, 2016. doi: 10.1109/TVCG.2015.2467191 3, 6
- [65] P. Xu, H. Mei, L. Ren, and W. Chen. ViDX: Visual diagnostics of assembly line performance in smart factories. *IEEE Transactions on Visualization and Computer Graphics*, 23(1):291–300, 2017. doi: 10.1109/TVCG.2016.2598664 8
- [66] Y. You, T. Chen, Y. Sui, T. Chen, Z. Wang, and Y. Shen. Graph contrastive learning with augmentations. In *Proc. Advances in Neural Information Processing Systems*, vol. 33, pp. 5812–5823, 2020. doi: 10.48550/arXiv.2010.13902 6
- [67] Z. Yuan, H. E. Chua, S. S. Bhowmick, Z. Ye, W.-S. Han, and B. Choi. Towards plug-and-play visual graph query interfaces: data-driven selection of canned patterns for large networks. *Proc. VLDB Endow.*, 14(11):1979–1991, 13 pages, 2021. doi: 10.14778/3476249.3476256 2
- [68] E. Zraggen, S. M. Drucker, D. Fisher, and R. Deline. (slq)eries: Visual regular expressions for querying and exploring event sequences. In *Proc. ACM Conf. Human Factors in Computing Systems (CHI)*, pp. 2683–2692, 2015. doi: 10.1145/2702123.2702262 2
- [69] B. Zheng and F. Sadlo. On the visualization of hierarchical multivariate data. In *Proc. IEEE Pacific Visualization Symposium (PacificVis)*, pp. 136–145, 2021. doi: 10.1109/pacificvis52677.2021.00026 1