

# CrossVis: A Visual Analytics System for Exploring Heterogeneous Multivariate Data with Applications to Materials and Climate Sciences

Chad A. Steed<sup>a,\*</sup>, John R. Goodall<sup>b</sup>, Junghoon Chae<sup>c</sup>, Artem Trofimov<sup>d</sup>

<sup>a</sup>Computer Science and Mathematics Division, Oak Ridge National Laboratory, Oak Ridge, TN 37831, USA

<sup>b</sup>Cyber and Applied Data Analytics Division, Oak Ridge National Laboratory, Oak Ridge, TN 37831, USA

<sup>c</sup>Computer Science and Mathematics Division, Oak Ridge National Laboratory, Oak Ridge, TN 37831, USA

<sup>d</sup>Materials Science and Technology Division, Oak Ridge National Laboratory, Oak Ridge, TN 37831, USA

---

## ARTICLE INFO

### Article history:

Received March 29, 2020

---

Information Visualization, Multivariate Data Analysis, Statistical Analytics, Materials Science, Climate Science, Explainable Neural Networks, Parallel Coordinates

---

## ABSTRACT

We present a new visual analytics system, called CrossVis, that allows flexible exploration of multivariate data with heterogeneous data types. After presenting the design requirements, which were derived from prior collaborations with domain experts, we introduce key features of CrossVis beginning with a tabular data model that coordinates multiple linked views and performance enhancements that enable scalable exploration of complex data. Next, we introduce extensions to the parallel coordinates plot, which include new axis representations for numerical, temporal, categorical, and image data, an embedded bivariate axis option, dynamic selections, focus+context axis scaling, and graphical indicators of key statistical values. We demonstrate the practical effectiveness of CrossVis through two scientific use cases; one focused on understanding neural network image classifications from a genetic engineering project and another involving general exploration of a large and complex data set of historical hurricane observations. We conclude with discussions regarding domain expert feedback, future enhancements to address limitations, and the interdisciplinary process used to design CrossVis.

© 2020 Elsevier B.V. All rights reserved.

---

## 1. Introduction

Forming a comprehensive understanding of patterns and relationships in multivariate data, where the phenomena under investigation are influenced by multiple factors, is integral to unlocking the full potential of today's vast data sets, especially in scientific domains. Whether interpreting the output of deep learning algorithms or exploring historical climate observations, scientists require interactive tools to develop a comprehensive understanding of large, multivariate data sets.

Developing new and improved multivariate visualization techniques for data exploration has captured the attention of

many researchers as evidenced by a recent survey from Liu et al. [1]. However, real world analysis of such data remains a significant challenge for several reasons. One challenge is rooted in the technical difficulties of exploring increasingly large volumes of data. Another lies in equipping scientists with effective sense-making techniques (e.g., visual representations, interac-

12  
13  
14  
15  
16  
17

---

(Artem Trofimov)

This manuscript has been authored by UT-Battelle, LLC under Contract No. DE-AC05-00OR22725 with the U.S. Department of Energy. The United States Government retains and the publisher, by accepting the article for publication, acknowledges that the United States Government retains a non-exclusive, paid-up, irrevocable, worldwide license to publish or reproduce the published form of this manuscript, or allow others to do so, for United States Government purposes. The Department of Energy will provide public access to these results of federally sponsored research in accordance with the DOE Public Access Plan (<http://energy.gov/downloads/doe-public-access-plan>).

---

\*Corresponding author: Tel.: +1-865-574-7168

e-mail: [csteed@acm.org](mailto:csteed@acm.org) (Chad A. Steed), [jgoodall@ornl.gov](mailto:jgoodall@ornl.gov) (John R. Goodall), [chaej@ornl.gov](mailto:chaej@ornl.gov) (Junghoon Chae), [trofimovaa@ornl.gov](mailto:trofimovaa@ornl.gov)

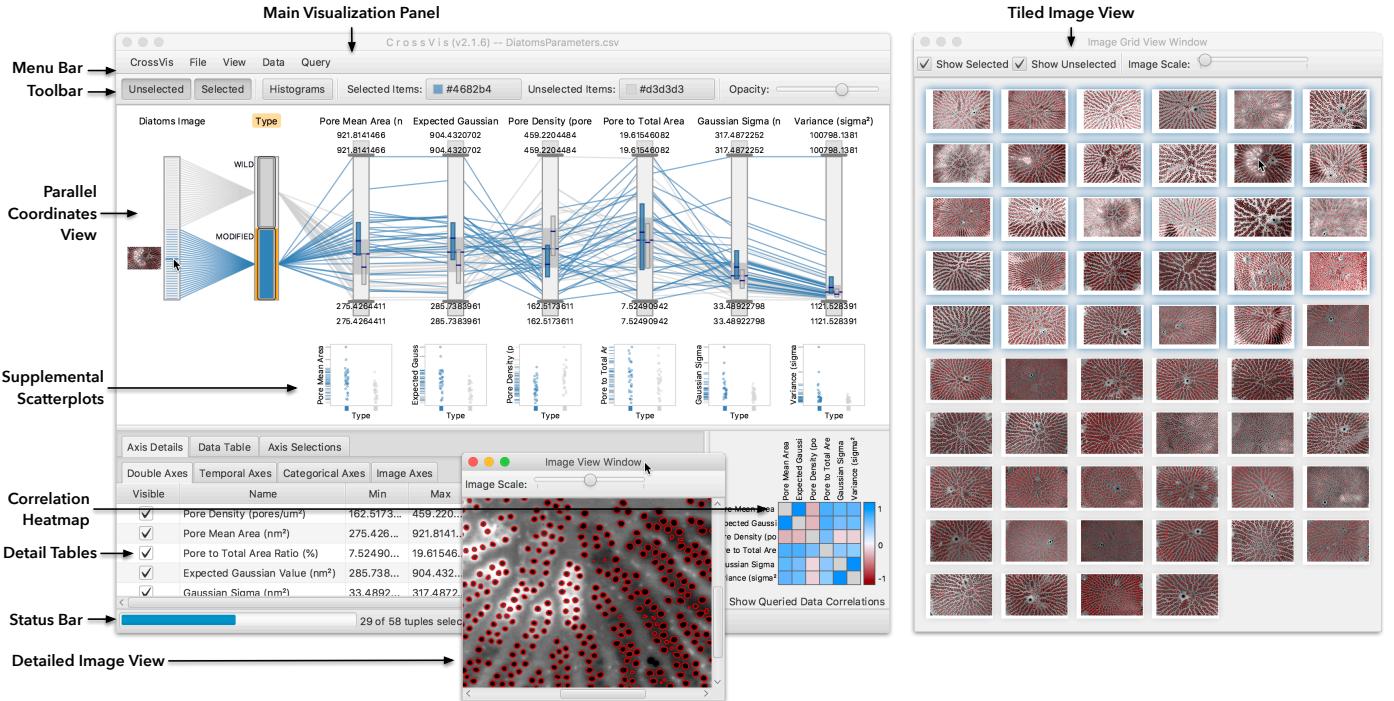


Fig. 1: The CrossVis system presents several interactive views of multivariate data that are coordinated using a single tabular data model. The main view, an extended version of the parallel coordinates plot, is supplemented with scatterplots, correlation visualizations, and image views to foster creative exploratory data analysis. CrossVis's interactive query capabilities help uncover key patterns, which are revealed through the various visual representations.

1 tive queries) for multivariate patterns. In practice, multivariate  
 2 data sets often contain heterogeneous data types, missing values,  
 3 and quality issues, which exacerbate the problem. Even  
 4 with moderately sized data sets, heterogeneous data types repre-  
 5 sent a significant barrier to thorough data exploration. But when  
 6 data sets are large and contain a mixture of data types (e.g.,  
 7 numeric, categorical, temporal, and image data) with problem-  
 8 atic values, developing an adequate understanding of the data  
 9 with the tools at hand becomes increasingly difficult, especially  
 10 when coupled with the scientist's desire to freely navigate alter-  
 11 nate paths during their analytical journey. To develop a com-  
 12 prehensive understanding, scientists need comprehensive solutions  
 13 that address these challenges.

14 Visual analytics offers a viable approach for coping with  
 15 these issues [2, 3]. By blending interactive visualizations with  
 16 computational guidance, well-designed visual analytics sys-  
 17 tems harness human and computational strengths to improve the  
 18 outcomes of data-driven studies. Yet, the number of visual ana-  
 19 lytics techniques that are readily available to non-visualization  
 20 experts for practical use with large, heterogeneous, and multi-  
 21 variate data, which are ubiquitous in modern scientific studies,  
 22 is low, especially when a combination of capabilities is desired.

23 In light of these and other practical challenges, we have de-  
 24 veloped the CrossVis visual analytics system (shown in Fig-  
 25 ure 1) in collaboration with domain experts. CrossVis ex-  
 26 pands the parallel coordinates plot (PCP) [4] to support new  
 27 axis representations for several non-numeric data types, em-  
 28 bedded bivariate PCP axes, linked supplemental visualizations,  
 29 focus+context axis scaling, and views that vary the level of de-  
 30 tail. A progressive rendering algorithm and an optimized data

model provide support for large data sets. These features are  
 31 motivated by feedback from domain experts in multiple sci-  
 32 entific domains. This paper presents the design and integration of  
 33 these methods into a flexible system that enables compre-  
 34 hensive multivariate data exploration. The main contributions of  
 35 the current work include the following:

- 36 • New visual representations of numerical, categorical, tem-  
 37 poral, image, and bivariate PCP axes that, when com-  
 38 bined with linked supplemental visualizations and inter-  
 39 active queries methods, reveal trends and patterns in het-  
 40 erogeneous, multivariate data
- 41 • Functional design considerations related to the visual rep-  
 42 resentations provided by CrossVis including discussion of  
 43 alternative approaches
- 44 • An overview of CrossVis describing the incorporation of  
 45 several PCP extensions (visual representations and inter-  
 46 action techniques) into a comprehensive system that is  
 47 greater than the sum of its parts
- 48 • Feedback from domain experts following their application  
 49 of CrossVis to practical scientific data analysis scenarios

## 51 2. Related Work

52 CrossVis employs multiple visualization methods, such as  
 53 scatterplots, tiled image views, and correlation heatmaps, but  
 54 the focal point is an extension of the classic PCP technique.  
 55 Inselberg [4] initially popularized the PCP as a method for vi-  
 56 sualizing hyper-dimensional geometries and later Wegman [5]

1 applied it to the analysis of multivariate data. The standard  
 2 PCP method yields a compact two-dimensional representation  
 3 of multidimensional data sets by mapping the  $N$ -dimensional  
 4 data tuple  $C$  with coordinates  $(c_1, c_2, \dots, c_N)$  to points on  $N$   
 5 parallel axes, which are joined using a polyline [6]. The PCP  
 6 is attractive for exploratory data analysis because it transforms  
 7 high-dimensional data sets into a two-dimensional plot without  
 8 dimensionality reduction. Although the number of variables  
 9 that can be shown is only geometrically restricted by the reso-  
 10 lution of the display, axes that are located next to one another  
 11 yield the most obvious insight. To analyze relationships be-  
 12 tween variables that are separated by one or more axes, interac-  
 13 tions and representations of information derived from analytical  
 14 algorithms are necessary.

15 As recent surveys of PCP methods [7, 8] and a book on  
 16 PCPs by Inselberg [6] demonstrate, the drive to improve and  
 17 apply PCPs has attracted considerable attention. In addition to  
 18 extensions of the technique, a significant portion of the prior  
 19 work involves the application of PCPs to a wide variety of do-  
 20 mains, such as climate science [9, 10, 11], cybersecurity [12],  
 21 computer forensics [13], genetics [14], biomedical [15], health-  
 22 care [16], and environmental pollution [17].

23 CrossVis implements several common PCP interactions  
 24 building on the direct interaction techniques described in Si-  
 25 iirtola [18] and visual data mining methods reviewed by Insel-  
 26 berg [6]. These interactions enhance exploratory data analy-  
 27 sis and include polyline selections, reorderable axes, and de-  
 28 tails on demand. CrossVis also extends the standard PCP axis  
 29 with graphical indicators of various summary statistics follow-  
 30 ing earlier designs by Hauser et al. [19]. The current work de-  
 31 scribes extensions to these interactions as well as linkages to  
 32 new visual representations.

33 PCPs and scatterplots [20] are two of the most popular mul-  
 34 tivariate data visualization techniques. Due to complemen-  
 35 tary characteristics, some prior work has combined the two  
 36 into a single layout [17, 21] using a coordinated multiple view  
 37 (CMV) strategy [22]. Yuan et al. [23] introduced the scatter-  
 38 ing points technique to embed scatterplot points between PCP  
 39 axes. Both PCPs and scatterplots excel at showing correlation  
 40 relationships between variables [24]. Some previous work di-  
 41 rectly augmented standard PCPs with graphical indicators that  
 42 encode correlation metrics, such as the Pearson correlation co-  
 43 efficient, to guide users to potentially significant trends [18].  
 44 Zhou and Weiskopf [25] delved deeper into correlation analysis  
 45 using PCPs by introducing an indexed point representation of  
 46 multivariate correlations as opposed to most PCP systems that  
 47 focus on the relationships between pairwise variable combi-  
 48 nations. In addition to supplemental scatterplots and correlation  
 49 indicators, CrossVis includes the ability to interactively embed  
 50 a scatterplot between axes in the main PCP view to investigate  
 51 pairwise correlations.

52 Although the vast majority of PCP methods focus on nu-  
 53 merical data, the desire to represent other data types has in-  
 54 spired PCP extensions. Kosara et al. [26] introduced the Par-  
 55 allelSets technique to allow interactive representations of cat-  
 56 egorical data. Fernstad and Johansson [27] demonstrated that  
 57 the ParallelSets method is superior to common quantitative en-

58 codings of categorical data in frequency related tasks. More  
 59 recently, Vosough et al. [28] described the parallel hierarchies  
 60 technique, which used parallel Icicle Plots to display hierarchi-  
 61 cal categorical data. CrossVis includes a variation of the Par-  
 62 allelSets approach for representing categorical data as well as  
 63 new representations for temporal and image data. To the best of  
 64 our knowledge, CrossVis represents the first PCP system with  
 65 support for numerical, temporal, categorical, and image-based  
 66 data in a single PCP framework.

67 When dealing with some moderate and most large scale data  
 68 sets, PCPs are prone to polyline overplotting and occlusion is-  
 69 sues [24, 29]. Clustering [29, 30] and binning methods [31]  
 70 alleviate these issues by reducing the number of polylines that  
 71 are rendered. Other approaches include alpha blending and  
 72 displaying statistical representations (e.g., summary statistics,  
 73 histograms), in lieu of or in combination with polylines, to  
 74 represent the data at a higher level of detail [19, 18, 32, 33].  
 75 Finally, both graphical processing units (GPUs) [34] and dis-  
 76 tributed computing infrastructure [35] have been harnessed to  
 77 improve the rendering speed and scalability of PCPs. CrossVis  
 78 uses a progressive rendering algorithm that leverages system  
 79 GPUs for improved performance. In addition, CrossVis uses  
 80 statistical representations of raw data to provide summarized  
 81 levels of detail, thereby reducing the need to represent every in-  
 82 dividual polyline for large data sets. CrossVis also integrates a  
 83 focus+context technique for PCPs, which allows users to zoom  
 84 into ranges of interest on numerical and temporal axes with  
 85 dense clusters of polylines while maintaining contextual aware-  
 86 ness, similar to work by Novotný et al. [31] and more recently  
 87 Richer et al. [36] where a focus+context approach is formalized  
 88 with abstract PCPs.

### 3. Design Requirements

89 For over a decade, we have collaborated closely with sci-  
 90 entists from climate, materials science, manufacturing, and other  
 91 fields that engage in multivariate data analysis. Through these  
 92 engagements, we have observed two fundamental limitations.  
 93 At the onset of these collaborations, scientists often state that  
 94 *they fail to examine enough of their data*. This issue is par-  
 95 tially due to large data volumes, but other factors come into play  
 96 such as inadequate visualization support for heterogeneous data  
 97 types and cumbersome query support. Scientists also state that  
 98 although they are good at finding patterns they already know,  
 99 *new discoveries are slow to occur*. This issue can be tied to an  
 100 inability to interactively explore the full data set as well as the  
 101 application of automated methods or workflows that focus on  
 102 known relationships, which can lead to anchoring bias.

103 We postulate that more flexible human-centered interactions,  
 104 scalable visualizations, and comparative techniques that are tai-  
 105 lored to specific data types are viable solutions to these issues.  
 106 In the remainder of this section, we consider these issues in  
 107 greater detail and present the design requirements that have  
 108 guided the development of CrossVis.

109 **R1: The visualizations should show the distinguishing**  
 110 **characteristics of heterogeneous data types.** Visualizations  
 111 of multivariate data often transform temporal and categorica-

1 data into numerical values because a wider range of numerical  
 2 representations are available. However, this process often  
 3 leads to misleading statistical summaries and informative char-  
 4 acteristics of the native data types are discarded. Therefore,  
 5 CrossVis includes visual representations that are tailored to key  
 6 data types (e.g., temporal, categorical, images) to enable more  
 7 comprehensive analysis.

8 **R2: The system should support flexible comparative**  
 9 **analysis of variables and subsets.** The ability to compare dif-  
 10 ferent variables and/or subsets of values empowers scientists to  
 11 freely ask questions and explore more of the data. This capa-  
 12 bility is challenging for multivariate data sets, especially het-  
 13 erogeneous data, due to the number of ways items can be com-  
 14 compared. To meet these demands, scientists require the ability to  
 15 quickly select data. Furthermore, visual representations must  
 16 clearly highlight variations between selections. Often both di-  
 17 rect and indirect selections are necessary; precise adjustment  
 18 capabilities for selecting parameters through indirect controls  
 19 can complement fast and approximate direct selections.

20 **R3: The system should clearly link selections in separate**  
 21 **views of the data through highlights in the visual represen-**  
 22 **tations.** Separate visualizations should be linked through inter-  
 23 actions so that changes in one view are propagated to all views.  
 24 Each view should be designed to clearly communicate specific  
 25 aspects of the data to increase the probability of finding new in-  
 26 sights. A viable way to achieve this requirement is to provide  
 27 coordinated multiple views where interactions are managed us-  
 28 ing a separate data model that shares selection state through an  
 29 event-based listener interface.

30 **R4: The system should maintain responsive interactions**  
 31 **and visualizations.** Exploratory data analysis techniques must  
 32 maintain responsive interactions to avoid disrupting cognitive  
 33 flow. By spawning threads to handle different aspects of inter-  
 34 action, summarization, and rendering, the system can orches-  
 35 trate processing demands and prioritize tasks vital to interactive  
 36 performance. Rendering threads are launched to render subsets  
 37 of the updated data and progressively refine the visualization.  
 38 The result is an increase in perceived scalability with large data  
 39 sets. In addition, computational optimization, preprocessing,  
 40 and caching strategies help sustain interactive performance.

41 **R5: The visualizations should support views at multiple**  
 42 **scales.** A complementary approach to progressive rendering is  
 43 to summarize the raw data at different levels of detail to reduce  
 44 the number of graphical elements that are needed to display the  
 45 essence of the data. In addition to reducing rendering times,  
 46 summarized views reduce occlusion and clutter. In CrossVis,  
 47 variable detail displays are implemented using a hierarchy of  
 48 statistical summaries that drill down to raw data views.

49 **R6: The visualizations should display magnified views**  
 50 **with context.** Focus+context techniques that allow users to ex-  
 51 pand or zoom into specific regions of interest in the overall data  
 52 space are helpful for selectively probing higher levels of de-  
 53 tail on demand. Focus regions magnify data within a particular  
 54 range and provide contextual awareness by preserving the prox-  
 55 imity of the focused region within the whole. These capabilities  
 56 make dense clusters of shapes more legible.

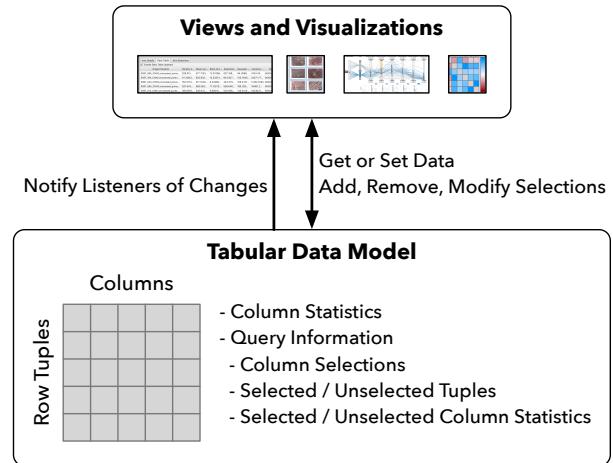


Fig. 2: A tabular data model supports multiple data types and provides access to both raw data and statistical summary information. The data model also coordinates user queries across the linked visualizations.

#### 4. Introducing the CrossVis System

57  
 58 As shown in Figure 1, CrossVis is comprised of a main visual-  
 59 ization panel that is supplemented by other linked views. In  
 60 addition to detailed table views, the main panel is augmented  
 61 by a correlation matrix, scatterplots, and image-based visual-  
 62 izations. The main visualization panel is an extension of the  
 63 PCP featuring a unique combination of axis representations and  
 64 embedded visualizations of statistical information.

65 CrossVis is an open source application<sup>1</sup> with performance  
 66 enhancements that address the large scale analysis needs men-  
 67 tioned in **R4**. The JavaFX graphics library is used to render  
 68 geometric shapes and the JavaFX user interface library is used  
 69 for layouts, menus, and windows. These libraries automatically  
 70 utilize system GPUs to boost rendering performance in a plat-  
 71 form independent manner. In addition, the rendering algorithms  
 72 use parallel threads to prioritize the display of more salient vi-  
 73 sual features and progressively reveal more details. In the re-  
 74 mainder of this section, the data model and key data visualiza-  
 75 tion techniques are described.

##### 4.1. Tabular Data Model

76 CrossVis is supported by a custom-developed, tabular data  
 77 model (see Figure 2) that stores raw data, derived statistics, and  
 78 selection criteria using collections of row and column objects.  
 79 This data model is a critical component in CrossVis, and it is in-  
 80 tegral to fulfilling all of the previously mentioned requirements  
 81 (**R1–R6**). Data structures are allocated in working memory as  
 82 files are loaded, but these structures are not serialized to disk  
 83 like a typical database system. The data model provides opti-  
 84 mized statistical summaries, fast data access, and modular sup-  
 85 port for columns of numeric, categorical, temporal, and image  
 86 data. Internally, row data are stored in native data types using a  
 87 generic object array. A row-based data structure is implemented

<sup>1</sup>CrossVis is available at <https://github.com/ORNL/CrossVis>

1 to match the access patterns of the PCP rendering algorithms,  
 2 which display row tuples as polylines. The column objects store  
 3 metadata and summary statistics, and they provide convenience  
 4 methods for accessing data elements as native values. Caching  
 5 mechanisms are also integrated for improved performance.

6 The data model manages subset selections using column se-  
 7 lection criteria (e.g., value ranges, value sets) and an event-  
 8 based listener interface to propagate changes. Data views reg-  
 9 ister as listeners with the data model and implement a set of  
 10 interface methods to respond to changes. Data views transmit  
 11 user interactions to the data model, which notifies registered lis-  
 12 teners. For quick access, the data model also maintains a query  
 13 object that holds column selection criteria and references to the  
 14 currently selected and unselected rows.

15 Column objects store summary statistics for the overall data  
 16 distribution. The query object also stores column summary  
 17 statistics for both the selected and the unselected rows. As se-  
 18 lection criteria change in the visualizations, the data model de-  
 19 tects the changes and updates the statistical summaries, which  
 20 triggers the event-listener interface and forces other views to re-  
 21 draw. The summaries supply the visualizations with a fast level  
 22 of detail hierarchy that enables multiple scale views. Standard  
 23 descriptive statistics (e.g., mean, median, standard deviation,  
 24 interquartile ranges) are calculated for numerical columns, and  
 25 frequency-based statistics (e.g., histograms) are calculated for  
 26 numerical, categorical, temporal, and image data.

27 The CrossVis data model performance is highly dependent  
 28 on the host system's processor, graphics cards, and memory  
 29 configurations. Most development occurred on a MacBook Pro  
 30 with 16GB of random access memory, a 3.1GHz Intel Core  
 31 i7 processor, and a 4GB AMD Radeon Pro 5600 GPU. With  
 32 data files containing less than 10,000 rows, CrossVis maintains  
 33 responsive interactions and rendering. As row counts are in-  
 34 creased, responsiveness tends to suffer. However, drawing such  
 35 a large number of lines in PCPs and points in scatterplots is  
 36 often not useful due to overplotting side effects that make it  
 37 difficult to see patterns. In these situations, CrossVis is de-  
 38 signed to show histograms and / or summary statistics. These  
 39 higher level views are capable of revealing patterns for larger  
 40 segments of data and the user can select subsets using the in-  
 41 teractive query techniques to draw polylines on demand. This  
 42 scheme leverages a level of detail hierarchy to allow exploration  
 43 of data sets containing upwards of 100,000 rows and a dozen  
 44 or more columns. Further performance enhancements can be  
 45 achieved, but these data scales represent a sweet spot for our  
 46 targeted users.

#### 4.2. New Axis Representations for Parallel Coordinates

47 The main visualization panel extends the PCP method to in-  
 48 clude new representations for specific data types and statisti-  
 49 cal information. PCP polylines correspond to row tuples in the  
 50 data model and vertical axes to columns. In addition to subtle  
 51 refinements (e.g., incremental line rendering, automated axis  
 52 layouts), the CrossVis PCP design includes new axis represen-  
 53 tations supporting additional data types (addressing **R1**), sup-  
 54 plemental displays of statistical information (addressing **R5**), and  
 55 focus+context scaling (addressing **R6**). In the remainder of  
 56 this subsection, these extensions are presented.

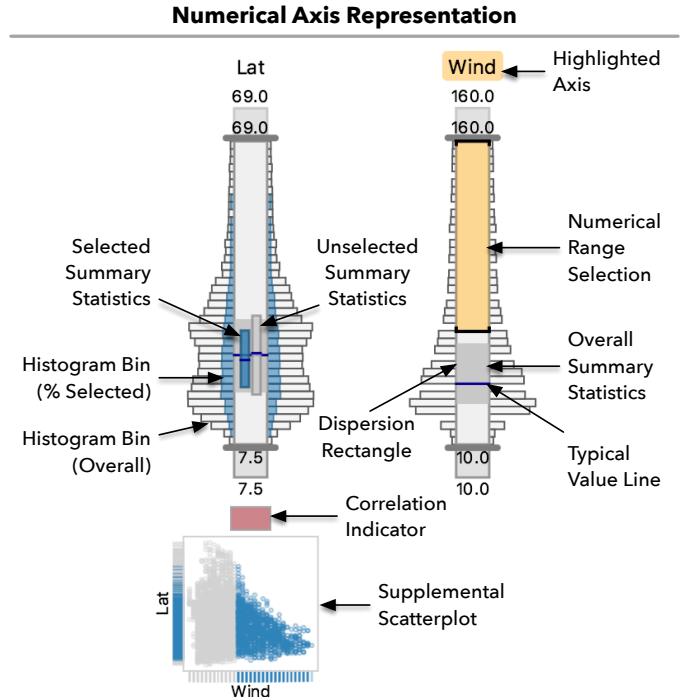


Fig. 3: CrossVis numerical axes are augmented with graphical statistical summaries using both descriptive and frequency statistics. Overall summaries, selected / unselected summaries, and scatterplots are visually represented providing an aggregated view to complement the detailed PCP polylines.

##### 4.2.1. Numerical Axis Representation

Numerical PCP axes are augmented with graphical statistical summaries (see Figure 3) and enhanced with the focus+context scaling technique described in Section 4.2.5. Both the typical value (mean or median) and dispersion range (two times the standard deviation range centered on the mean or the interquartile range) are represented in the axis bar interior. The user can toggle between mean- or median-based statistics using the application menu. These statistics are calculated for the overall distribution of values, the selected data, and the unselected data. For the overall distribution (see the 'Wind' axis in Figure 3), the height of the gray rectangle spanning the full axis bar width encodes the dispersion value and a blue line encodes the typical value. The overall statistical indicators serve as a baseline for comparisons against smaller subsets.

Narrow versions of the overall statistical indicators (see 'Lat' axis in Figure 3) summarize the selected and unselected data. Statistics for the selected data are shown on the left and unselected on the right. The selected and unselected statistical indicators are visually linked to the selected and unselected polylines using color; the two dispersion rectangles are filled with the current selected or unselected polyline colors, which the user can modify using buttons on the toolbar above the PCP panel (see Figure 1). These indicators are drawn over the overall indicators with a semi-transparent fill color to avoid completely masking the overall statistical information.

Numerical axes can also show frequency-based statistics as vertical histograms on the exterior of the axis bar (see Figure 3). A standard histogram is computed for both the overall distri-

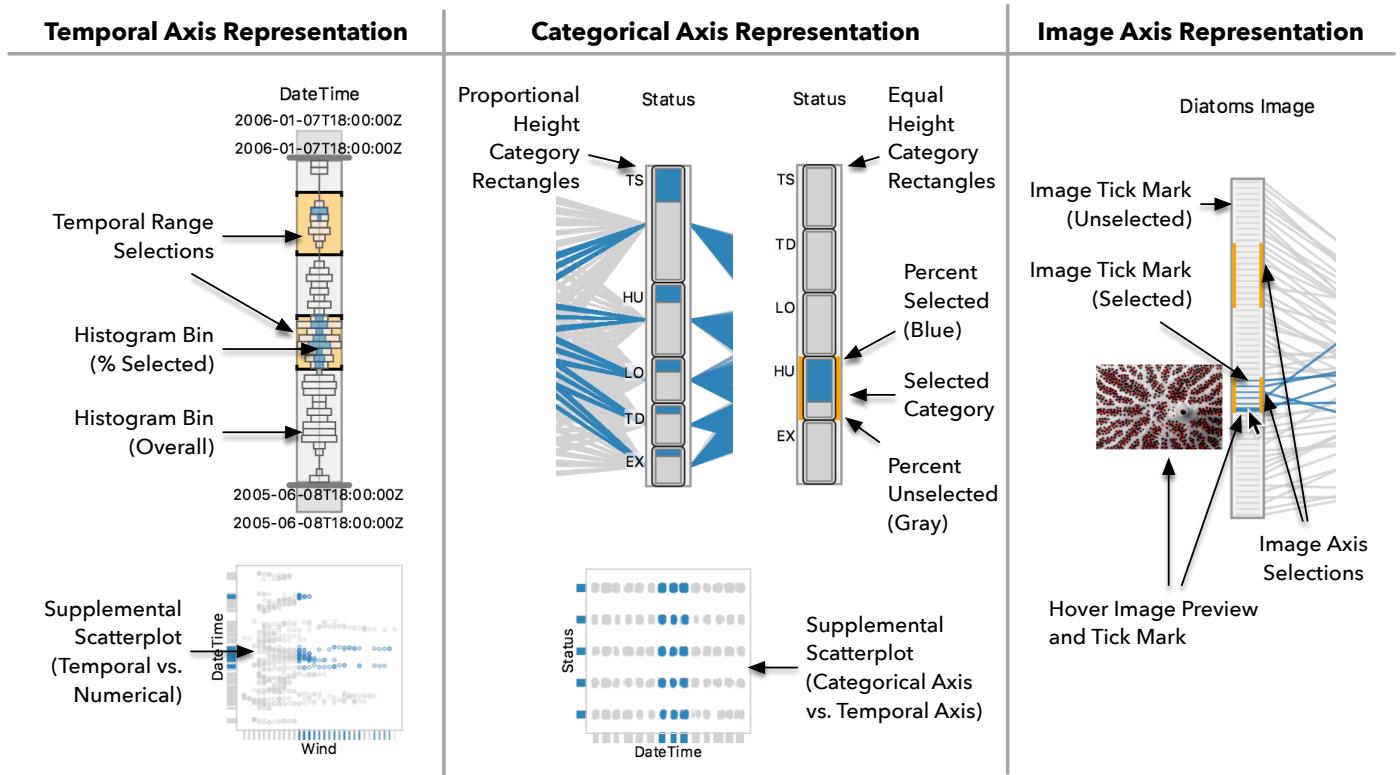


Fig. 4: CrossVis includes new axis representations for temporal, categorical, and image data. The unique characteristics of each data type inform the designs, which are augmented with statistical graphics and scatterplots. All the axes support interactive visual queries.

button and the selected data using bins covering equally sized value intervals. The number of bins is initially set to  $k = \lceil \sqrt{n} \rceil$ , where  $n$  is the number of rows in the data model, but the user can adjust the bin count through the application menu. A symmetrical layout is used with histogram bin rectangles shown on both sides of the axis. The rectangle width encodes the number of values falling within the bin's range. The percentage of values that are currently selected for each bin is encoded as the width of a semi-transparent rectangle drawn over the overall bin rectangle. The selected bin rectangle is filled with the current selected polyline color for consistency.

When both histograms and polylines are shown (see Figure 9b), dense polyline clusters can negatively affect the decoding of bin counts. To deal with this situation, a thin white line is added to the bin rectangle's outer edge and the bin outline stroke color has high value contrast with the fill color. These color effects are intended to make the histogram silhouette more salient. Some clutter is introduced due to the inability of the semi-transparent histogram bins to completely mask polylines connecting to the axis behind them, but the polyline colors are muted and often the ability to see individual polyline axis intersections is useful.

Before settling on the symmetrical histogram design, we experimented with an alternative approach where bins were drawn on only one side of the axis. This unbalanced design made it difficult to see trends when an axis was positioned between two other axes. Furthermore, the unbalanced design was visually inconsistent with the other visual elements, which are mostly symmetrical. We also experimented with mapping bin counts

to the fill colors of a band of smaller rectangles on the edge of the axis bar resembling a vertical heat map. Although using color required less space and avoided polyline occlusion, it was more difficult to compare relative bin counts, especially subtle differences. The superiority of positional to color encoding techniques is reported by Mackinlay [37].

Histograms can be shown instead of the polylines to improve performance with large data. Histograms provide more detail than the summary statistics, but less than individual lines. Thus, a series of increasingly detailed views is formed by showing statistical graphics, histograms, and then individual polylines providing a level of detail scheme that addresses **R5**.

#### 4.2.2. Temporal Axis Representation

Temporal axes (see Figure 4) are similar to numerical axes. Because the data model stores the values as time instants, labels, hover values, and range selections are shown using date-time formatted strings addressing **R1**. Temporal axis bars show a continuous value range and use the focus+context axis scaling technique (see Section 4.2.5). Instead of descriptive statistics, the axis interior displays a vertical temporal histogram where bin rectangles are centered horizontally. To provide more space for the histogram, temporal axis bars are wider than those of numerical axes. Locating the histogram inside the axis bar helps avoid the occlusion of polyline intersections. Similar to histograms on the numerical axes, the overall histogram bins are augmented with rectangles showing the percentage of selected values.

The symmetry of the temporal histogram visually unifies

it with histograms on numerical axes. Unlike numerical histograms, temporal histograms are always displayed. In earlier designs, the temporal histograms were drawn outside the axis bar. This alternative design had a stronger correspondence to the numerical histograms, but it left an empty axis interior. We compensated for the empty space by making the axis bar more narrow, but this approach disrupted the overall consistency of axis bar treatments. We settled on the interior representation to capitalize on the opportunity to avoid occluding polyline intersections while sacrificing some visual unity. In the future, we plan to revisit the temporal axis design to explore methods that encode additional statistical information in the axis bar interior (e.g., dynamic time warping similarity metrics [38]).

#### 4.2.3. Categorical Axis Representation

The categorical axis representation (see Figure 4) emphasizes the relative frequency of categories. Categories are represented as rectangles inside the axis bar. In Figure 4, the ‘Status’ axis has five categories. This figure also shows the two ways that category rectangles are displayed. On the left, the height of a category rectangle is mapped to the percentage of rows associated with its category. In this mode, the overall frequency trends are emphasized to support comparisons, but categories with a small percentage of values can be hard to see. The right ‘Status’ axis in Figure 4 shows the equal height mode, which divides the axis bar height by the number of categories making smaller categories more visible. Users can also enable the display of category names as labels drawn to the left of the axis.

When polylines are selected, category rectangles are split into two smaller rectangles with heights that encode the ratio of selected (on the top and filled with the current selected polyline color) to unselected (on the bottom and filled with the current unselected polyline color) polylines. For example, the ‘HU’ category of the right ‘Status’ axis in Figure 4 shows that about 75% of polylines associated with the ‘HU’ category are selected. To select polylines associated with a category, the user clicks on the category rectangle. On the right ‘Status’ axis in Figure 4, the ‘HU’ category is selected. A category is removed from a selection by clicking the category rectangle a second time. When the user hovers over a rectangle, a tooltip reveals detailed information about the category (see Figure 8).

Polylines are connected to the vertical centers of the overall category rectangles. In earlier designs, we evaluated representations that used polygonal shapes covering the full height of the category rectangles in a manner similar to the ParallelSets [26] design. However, during developer testing with data sets that had a large number of rows the polygonal shapes were difficult to read, especially between numerical (or temporal) and categorical columns. By adopting the polyline representation over the polygonal approach, we maintain visual consistency with representations between the other axis types and avoid clutter. However, we recognize that there is room for future improvements of the polyline representations on categorical axes.

#### 4.2.4. Image Axis Representation

For columns consisting of images, CrossVis features a unique axis representation. As shown in Figure 4, images are visually

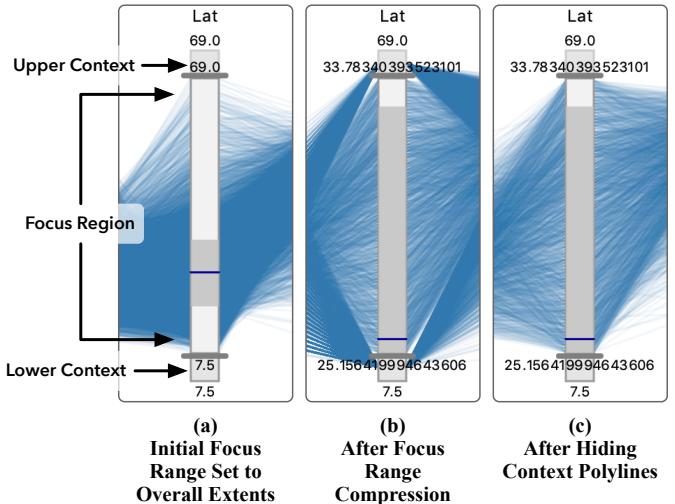


Fig. 5: Focus+context axis scaling is provided for temporal and numerical axes. In (a), the focus range, marked by the thick gray lines between the focus region and the two outer context regions, is set to the overall data extent. In (b), the focus range lines are moved in magnifying data between the 50% and 75% percentile. In (c), polylines in the context region are hidden to reduce clutter.

represented by horizontal tick marks inside the axis bar. The ordering of the image tick marks is determined by the file name. Tick marks are colored using either the selected or unselected polyline color. When the user hovers over a tick mark, its line's thickness increases and a small thumbnail copy of the corresponding image is shown to the left of the axis bar.

Images are selected by clicking on tick marks or dragging a selection range. To deselect images, the user presses the control key modifier while clicking or dragging. Selected images are indicated by orange highlights on the left and right of the axis (see Figure 4). Because selections are combined using an OR operation between different axes, it is possible that images can be included in selections on the image axis, but those images are not associated with the current overall set of selected polylines (see the upper axis selection on the ‘Diatoms Image’ axis in Figure 4). In such cases, the highlighting of selected images prevents selections on image axes from becoming invisible.

#### 4.2.5. Focus+Context Axis Scaling

Dense PCP polyline clusters make it difficult to decipher patterns due to overplotting. Adjusting the opacity of polylines helps, but the problem is not completely eliminated, especially with large data sets. To cope with this issue and address **R6**, CrossVis builds on the dynamic axis scaling technique introduced in MDX [11]. The CrossVis implementation provides more focus range control and adds support for temporal axes.

As shown in Figure 5, upper and lower context regions are located above and below the main focus region. The focus region extents are adjusted by dragging the thick gray lines at the focus range edges. Moving the maximum value boundary line down decreases the extent value and upward movement increases it. After the extents are modified, the display is redrawn, which spreads polylines in the focus region out and pushes some polylines into the context. Especially when lines connect from an

extreme edge of a neighbor axis to the opposite edge of another axis (e.g., top of one axis to the bottom of another), the context polylines can occlude the display after axis scaling. As shown in Figure 5c, the user can choose to hide polylines in either of the context regions to further reduce clutter.

#### 4.2.6. Bivariate Axis Representation

Bivariate PCP axis representations are supported because some variables are more easily understood in relationship to another variable. For example, geographic patterns are more apparent when latitude and longitude values are shown in a scatterplot (see Figure 9a). Univariate PCP axes make it hard to analyze such variables.

The bivariate axis is represented as a scatterplot using the same design as the supplemental scatterplots described in Section 4.3. However, the bivariate axis is embedded in the PCP with polylines of neighboring axes connecting to the y-axis of the scatterplot. By embedding the bivariate axis in the PCP, we alleviate potential perceptual issues associated with separated views, such as change blindness [39]. Any univariate axis can be combined with another to form a bivariate axis. The user can add bivariate axes manually by specifying the x and y columns from the data model, or the user can drag and release the x-axis on the target y-axis.

#### 4.2.7. Additional Interactive Axis Selection Considerations

The ability to query and filter data is essential for efficient exploratory data analysis. The selection capabilities in CrossVis meet comparative analysis needs in R2. As shown in Figure 3, users can select polylines that fall within value ranges on numerical and temporal axes. Multiple selections on multiple axes are supported (see Figure 4) by directly dragging a range selection using the mouse over an axis bar. The user can select a category by clicking on its associated rectangle on categorical axes and an image by clicking on its associated tick mark. Users can set bivariate selections by dragging rectangles within the scatterplots. For precise control, the user can also manually add selections for all axis types using the Axis Selections tab (see bottom of Figure 1).

Axis selections are visually unified using an orange highlight color. For range selections, the fill color of the selection rectangle uses the highlight color. For categorical and image selections, the rectangle or tick mark is augmented with an orange halo. The user can directly interact with the selection indicators to remove items or adjust the extents.

To compute the subset of selected polylines when multiple selections are present, a disjunction (OR) operation is first applied to selections on individual axes and then a conjunction (AND) operation is applied to the selected values between axes. This logic allows users to consider values spread out over non-contiguous value ranges for individual axes as well as relationships between different axes.

#### 4.3. Supplemental Scatterplots

In addition to the ability to embed scatterplots directly in the PCP, small scatterplots are shown below the PCP axes (addressing R3). When an axis is highlighted (see Figure 8), the

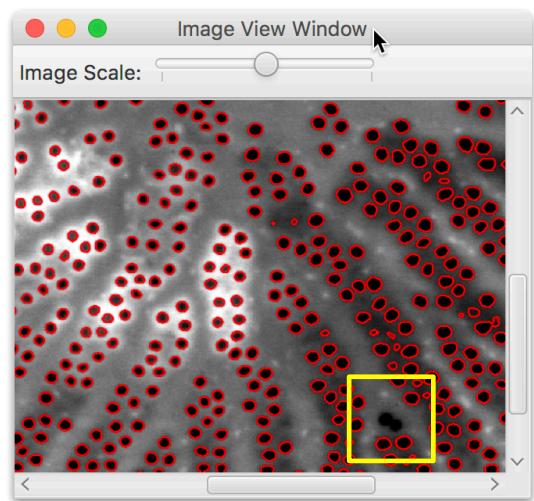


Fig. 6: The detailed image view allows scientists to visually examine images associated with polylines in the PCP. Here the view shows a diatom image where a pore detection algorithm failed to label two pores, outlined by the yellow box annotation, possibly because the boundary between the pores is blurred.

supplemental scatterplots show the pairwise relationship of the highlighted axis with all other axes. That is, the x-axis of each scatterplot is mapped to the highlighted axis and the y-axis is mapped to the axis above the scatterplot. When no axis is highlighted (see Figure 9b), the scatterplots are shown in the space between the axes where the x-axis is the left axis and the y-axis is the right axis. Both the supplemental scatterplots and the embedded bivariate axis scatterplots can be configured to show tick marks on the axis boundary and convey univariate distributions.

Scatterplots excel at conveying non-linear trends and clusters. Furthermore, scientists are usually familiar with scatterplots and displaying these in conjunction with the PCP, which is often new to them, can increase understanding. Thus, the combination of these two techniques is more valuable than showing either in isolation.

#### 4.4. Axis Correlation Coefficient Representations

Supporting R2 and R3 requirements, correlations between numerical axes are shown in several ways. The user may glean correlations from polyline configurations in the PCP (e.g., 'X' shaped crossings indicate negative correlations and more horizontal crossings indicate positive) and point configurations in the scatterplots. In addition, direct encodings of the Pearson correlation coefficient,  $r$ , are shown above the supplemental scatterplot as color-filled rectangles (see Figure 4). The  $r$  values are mapped to a color scale where the most saturated blue represents a perfect positive correlation, the most saturated red represents a perfect negative correlation, and white represents no correlation. The user can hover the mouse over a cell to see the exact  $r$  value. As shown in Figure 1, CrossVis also shows a linked correlation matrix using the same color coding scheme as the indicators in the parallel coordinates plot. To the right of the matrix heatmap, the  $r$  value color scale is shown.

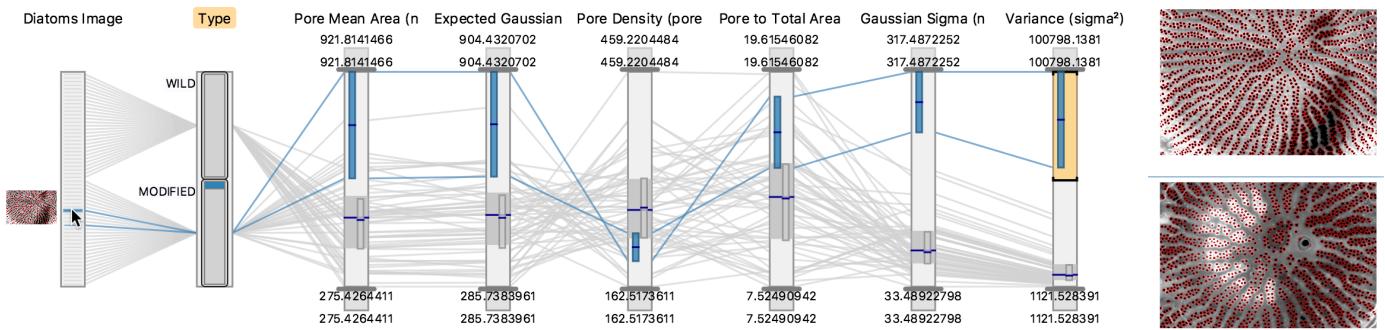


Fig. 7: Two outlier images are selected on the ‘Variance’ axis. The pixel variance of the images, both classified as ‘MODIFIED’, is evident in the thumbnails on the right. These images are outliers on all but the ‘Pore Density’ axis.

#### 1 4.5. Image Set and Individual Image Views

2 CrossVis provides two separate image views (addressing  
3 **R3**): a tiled image set view and a detailed image view (see  
4 labels in Figure 1). A slider at the top of both views allows  
5 adjustments of the image(s) size. The image set view is linked  
6 to the other visualizations in the main window. When polylines  
7 are selected, the images in the selection set are haloed with the  
8 selected polyline color. Likewise, the unselected images are  
9 haloed with the unselected color. Furthermore, selected images  
10 are located at the top of the image set view. In Figure 1, the  
11 cutoff between selected and unselected images is visible after  
12 the fifth image of the fifth row.

13 When the user double-clicks on an image in the tiled view, a  
14 detailed view appears (see Figure 1 at bottom) making it easier  
15 to focus on a specific image. In Figure 6, the detailed image  
16 view shows a magnified region of the image where some of  
17 the dark circular shapes (see yellow box) of the microscopic  
18 view are not enclosed by red outlines that were added by a pore  
19 detection algorithm.

#### 20 5. Two Practical Scientific Use Cases

21 In this section, we present two scientific use cases of  
22 CrossVis to demonstrate its data exploration capabilities. The  
23 first focuses on understanding the results of an artificial neural  
24 network (ANN) designed to classify microscopic imagery; the  
25 motivating scenario that inspired the development of CrossVis.  
26 The second describes exploration of a historical hurricane ob-  
27 servation data set, and it illustrates CrossVis’s capacity to dis-  
28 cover, investigate, and validate patterns in a larger and more  
29 complex data set as well as its suitability as a general purpose  
30 exploration system.

##### 31 5.1. Use Case 1: Understanding Neural Network Image Clas- 32 sifications in Genetic Engineering

33 Scientists at Oak Ridge National Laboratory’s Center for  
34 Nanophase Materials Science (CNMS), one of whom is a co-  
35 author of this paper, used an ANN to automatically classify  
36 scanning electron microscope (SEM) images of diatoms, some  
37 of which were genetically modified. The ANN predicts whether  
38 images correspond to genetically modified (‘MODIFIED’) di-  
39 atoms or not (‘WILD’). A diatom is a unicell alga with a sil-  
40 ica cell wall. Diatoms are attractive candidates for functional

41 systems of materials with applications ranging from photonics,  
42 sensing, filtration, and drug delivery [40, 41]. The scientists  
43 used CrossVis to analyze the ANN results and the following  
44 narrative captures some of their findings.

45 In addition to the ANN classification of ‘WILD’ or ‘MOD-  
46 IED’, the scientists computed a number of parameters that  
47 quantify characteristics of pores detected in the images: den-  
48 sity of pores, mean area of pores, and the percentage of area  
49 occupied by pores relative to the total area of the valve captured  
50 in an image. Additionally, pore area distribution was extracted  
51 and fitted with a Gaussian distribution to yield two more param-  
52 eters: Gaussian value and Gaussian sigma. These parameters,  
53 in combination with a variance metric, yield a total of six quan-  
54 titative values that supplement the categorical value output from  
55 the neural network classification. The goal of this study was to  
56 understand the significance of these parameters for distinguish-  
57 ing between modified and unmodified diatom images.

58 We begin by selecting the ‘MODIFIED’ category on the  
59 ‘Type’ axis (see orange highlight on ‘MODIFIED’ category  
60 in Figure 1). This action associates the selected polylines and  
61 summary statistics (shown in median/IQR mode) with images  
62 marked as ‘MODIFIED’ and the unselected with ‘WILD’ im-  
63 ages to allow comparative analysis. In the tiled image view, the  
64 ‘MODIFIED’ classified images at top are haloed with the selec-  
65 tion color and the ‘WILD’ images are below. This view reveals  
66 that ‘WILD’ images appear more uniform and have more pores  
67 (pores are outlined in red by a separate process) as compared  
68 to the ‘MODIFIED’ images, which exhibit more pixel variance  
69 and wider gaps between pores.

70 The variance is particularly evident in the hover image on  
71 the ‘Diatoms Image’ PCP axis in Figure 1. The hover image  
72 is located at the intersection of row 2 and column 5 in the tiled  
73 image view. Figure 6 shows a magnified view of the detailed  
74 image window at the bottom of Figure 1. The image shows that  
75 two pores with fuzzy edges at the lower right corner (see yellow  
76 highlight box) were missed by the pore detection step. Several  
77 other missed pores are apparent in top detail image view on  
78 the right in Figure 7. Visual inspections using CrossVis give  
79 scientists the ability to drill-down and find such subtle patterns,  
80 which in this case provides an opportunity to improve the pore  
81 detection process. Furthermore, CrossVis’s ability to display  
82 the images at multiple scales helps scientists see that the two  
83 image sets are visually distinct.

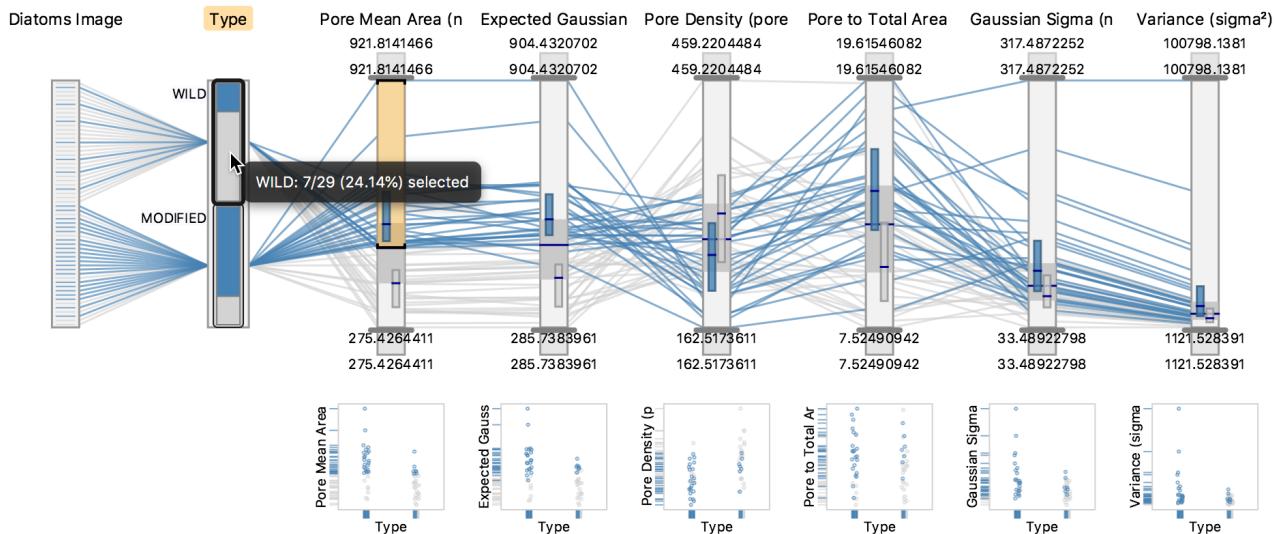


Fig. 8: A numerical range selection is used to compare images with high and low ‘Pore Mean Area’ values. The separation of the IQR rectangles on the ‘Pore Mean Area’ axis suggests the importance of these axes for distinguishing between the two ‘Type’ categories (‘WILD’ versus ‘MODIFIED’). The ‘Type’ axis is highlighted, resulting in the display of scatterplots below the other axes.

We shift our attention to the PCP to explore quantitative trends. In Figure 1, the IQR rectangles for the ‘Pore Mean Area’, ‘Expected Gaussian Value’, and ‘Pore Density’ axes show the most separation between the two image categories, which suggests that these are distinguishing features. The ‘Pore Mean Area’ and ‘Expected Gaussian’ axes exhibit a strong positive correlation ( $r = 0.97$ ), as indicated by a saturated blue square on the heatmap and nearly horizontal polylines between the two PCP axes. This finding suggests that one of the two variables can be removed since together they fail to add additional value. Having more overlap between the selected and unselected IQR rectangles, the ‘Expected Gaussian’ axis is a good candidate for removal. Both ‘Pore Mean Area’ and ‘Pore Density’ show clear separation between the two image categories, with ‘Pore Mean Area’ showing less overlap and fewer polylines crossing the median line.

In Figure 8, the ‘MODIFIED’ category selection on the ‘Type’ axis is removed and a range selection on ‘Pore Mean Area’ for values greater than the median is added. The selection captures 29 of the 59 images, most of which are of the ‘MODIFIED’ category. However, some ‘MODIFIED’ images are excluded and some ‘WILD’ images are included. The tooltip shows 7 of the 29 ‘WILD’ images are selected. The tooltip for the ‘MODIFIED’ category (not shown) shows that 22 of the 29 images are selected. This finding confirms the significance of pore sizes and density that we observed from the earlier visual inspections and it reinforces the importance of a multivariate process to accurately classify the images. As the materials scientist who led this analysis stated: “Such information suggests that not every single diatom in the ‘MODIFIED’ set underwent genetic modification, which was clearly revealed using CrossVis analysis.”

A glance at the ‘Variance’ axis in Figure 8 reveals two severe outliers in the upper range. These two polylines are selected in Figure 7 with the two associated images shown on the right. The ‘MODIFIED’ image discussed previously is shown and the

wide range of pixel fluctuations in both images confirms the pixel variance.

## 5.2. Use Case 2: Analyzing Historical Tropical Cyclone Observations

The National Oceanic and Atmospheric Administration (NOAA) maintains the Atlantic Hurricane Database (HURDAT2)<sup>1</sup>, which contains information on the location, winds, central pressure, and size (since 2004) of all known tropical and subtropical cyclones in the Atlantic basin between the years 1851 and 2017 [42]. HURDAT2 is important for understanding historical tropical cyclone trends, but the number of records (over 50,303 rows), number of variables (21 columns), and heterogeneous data types (categorical, temporal, and numerical) make it a challenge to analyze.

In Figure 9a, the full HURDAT2 data set is shown. Despite the size and number variables, CrossVis maintains interactive performance (< 1sec) during visual investigations. In the figure, scales for the 12 wind radii axes (located on the right side of the figure) are synchronized to a common range and several values with ‘no data’ wind radii values (wind radii fields were omitted for storms prior to 2004) are pushed into the context regions to achieve clearer views. The wind radii values provide wind swath size information for the 34 knot (34kt), 50 knot (50kt), and 64 knot (64kt) maximum wind ranges. For each range, four radii values are provided in nautical miles (nm) for the four quadrants: northeast (NE), southeast (SE), southwest (SW), and northwest (NW). The view reveals that wind swaths grow tighter (less dispersed) for fields with higher wind speed since the distance values decrease and associated wind speeds increase from left to right.

The resulting view highlights five records (selected in Figure 9a) with remarkably large ‘SE\_64kt’ values (see lines cap-

<sup>1</sup>HURDAT2 is available at <https://www.nhc.noaa.gov/data/>

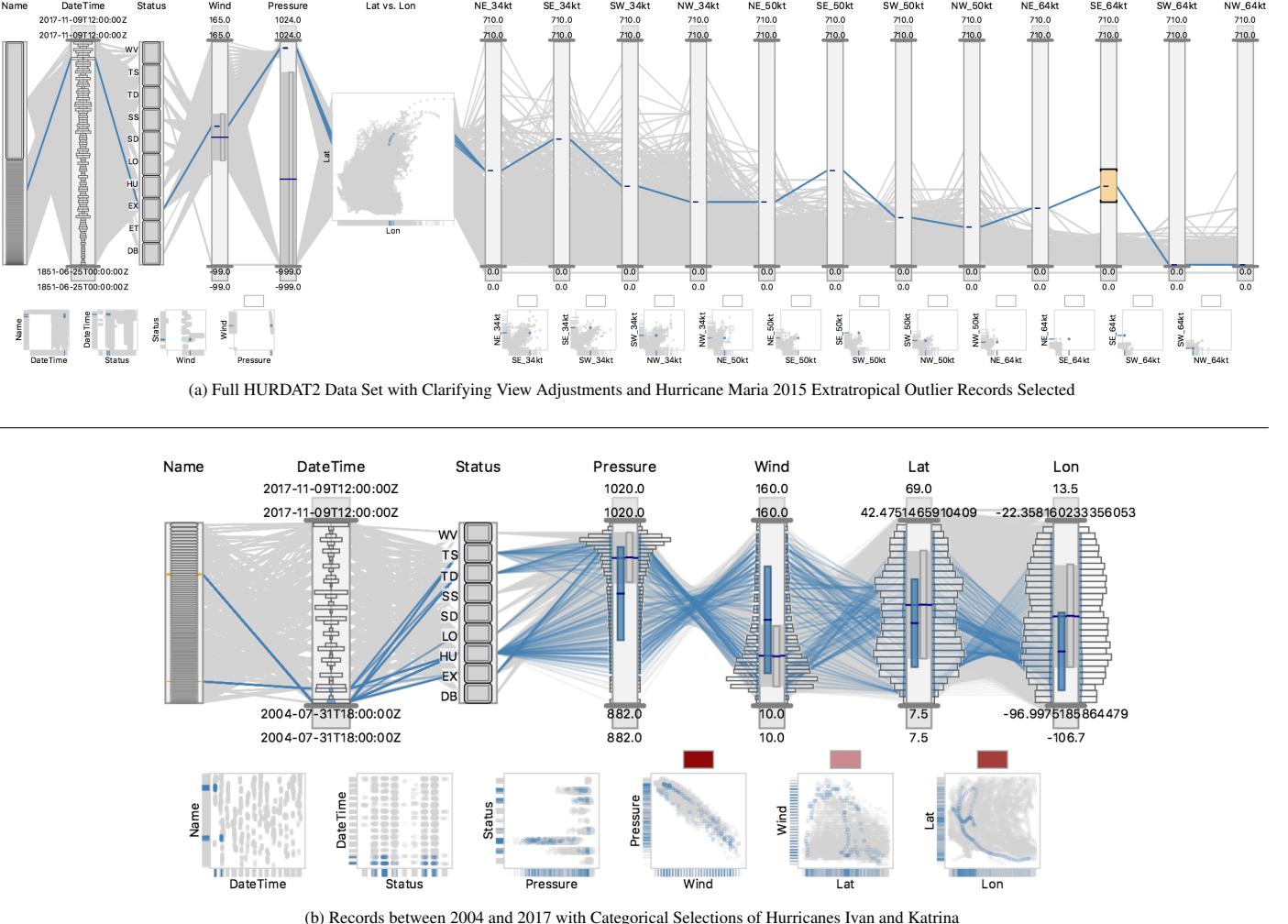


Fig. 9: CrossVis provides interactive techniques to cope with large multivariate data and the quality issues, such as missing data, that are common in scientific applications. In (a), the full NOAA HURDAT2 tropical cyclone data set is shown after we push ‘no data’ flagged values into the lower context region (see lines below 0.0 on the wind radii axes) and synchronize the 12 wind radii axis scales. Five outlier values on the ‘SE\_64kt’ axis are selected originating from extratropical records of Hurricane Maria 2015. In (b), a view of records between 2004 and 2017 is shown with a selection on the ‘Name’ axis for Hurricanes Ivan and Katrina.

tured by the range selection on this axis). Using the numeric data table view (not shown), we find that these records are from Hurricane Maria in 2005. With these records selected, it is evident that the southeast quadrant is larger than normal on the other wind region axes (‘SE\_34kt’ and ‘SE\_50kt’). The selection also shows that the western side of the storm swath has collapsed in the most intense wind region (0 nm for the ‘SW\_64kt’ and ‘NW\_64kt’ axes).

On the ‘Status’ axis, we see that the records for Hurricane Maria 2005 are assigned the ‘EX’ category, which denotes an extratropical system in the middle latitudes (between 30° and 60° latitude) of Earth. The latitudinal location is visible in the ‘Lat vs. Lon’ axis. These observations are consistent with the tendency of extratropical storms to become less symmetric in the middle latitudes [43] and the straightforward way these complex relationships are explored with CrossVis demonstrates its effectiveness in data verification and validation tasks.

Figure 9b focuses on storm records between 2004 and 2017, which offer the most detail. We select the ‘IVAN’ and ‘KAT-

RINA’ categories on the ‘Name’ axis highlighting records for Hurricane Ivan (2004) and Hurricane Katrina (2005), which were both major hurricanes that caused tremendous destruction. A high correlation between ‘Pressure’ and ‘Wind’ ( $r = -0.94$ ) is indicated by the highly saturated red correlation indicator, scatterplot point configuration, and PCP polyline crossings. The statistical indicators also show that pressure values for these storms were lower than normal and wind values were higher. Both storms progressed through different stages of development, which accounts for the variation of values on the ‘Status’ axis. After adjusting the focus range of the ‘Lat’ and ‘Lon’ axes to drill into on the affected geographic range, we observe that on average these storms were more southern and western in latitude and longitude, respectively. To put it another way, these storms lingered longer in the Caribbean and Gulf of Mexico regions, as can be seen in the supplemental scatterplot below the ‘Lat’ and ‘Lon’ axes.

## 1 6. Discussion

2 CrossVis increases scientists' capacity to develop a more  
 3 comprehensive understanding of multivariate data, particularly  
 4 when heterogeneous data types are present, by providing flexi-  
 5 ble techniques for quickly querying data and visually represent-  
 6 ing results. The system offers a unique collection of analysis ca-  
 7 pabilities that correspond to specific challenges scientists face.  
 8 These challenges were uncovered through prior collaborations  
 9 and the process of developing techniques to address them was  
 10 executed in close collaborations with the scientists who now use  
 11 the system. As we developed CrossVis and evaluated iterations  
 12 with domain experts, we observed indicators of its effectiveness  
 13 and limitations. In this section, we discuss these observations  
 14 and reflect on the interdisciplinary design process.

### 15 6.1. Domain Expert Feedback and Observations

16 The effectiveness of CrossVis is perhaps most apparent in  
 17 a more comprehensive discussion of the results of the ANN  
 18 diatom image classification project that our materials science  
 19 collaborators recently published [44]. That work focused on  
 20 the ANN design and scientific implications to genetic engineering  
 21 as opposed to a detailed description of CrossVis, which the  
 22 current work provides. Some findings related to explaining  
 23 the diatom images classifications were revealed with an ear-  
 24 lier version of CrossVis and communicated in the publication  
 25 with figures from the main PCP visualization panel. The sci-  
 26 entists stated that CrossVis significantly increased their under-  
 27 standing of the complex ANN process and, as the quote in Sec-  
 28 tion 5.1 states, “clearly revealed” complex relationships in the  
 29 data. Prior to using CrossVis, scientists were forced to flip be-  
 30 tween static plots (e.g., scatterplots, histograms) and file system  
 31 image viewers to compare features. Querying and filtering the  
 32 data involved running scripts to regenerate static plots, which  
 33 slowed investigations and limited the number of different com-  
 34 binations of conditions that could be realistically viewed. Di-  
 35 mensionality reduction processes were also employed, but re-  
 36 sults were difficult to translate back to the original parameters.

37 Scientists also noted that CrossVis helped them think in a  
 38 more multivariate manner during analysis. For example, it be-  
 39 came clear that multiple pore measures were instrumental in  
 40 classifying the images. Although some variables were more  
 41 correlated than others, it was clear that no single variable could  
 42 be tied to the results of the ANN process. The expanded  
 43 PCP visualization efficiently conveyed this condition through  
 44 the interactive representations of images and categorical val-  
 45 ues with quantitative metrics. One domain expert mentioned  
 46 that CrossVis “enabled faster pairwise comparative variable  
 47 comparisons because we don't have to pick through and cycle  
 48 between a set of scatterplots.” Referencing the diatom im-  
 49 age classification example in the previous section, this same ex-  
 50 pert noted that the correlation between ‘Pore Mean Area’ and  
 51 ‘Pore Density’ was not apparent until the data was viewed in  
 52 CrossVis. Moreover, the PCP revealed that the ‘Pore to Total  
 53 Area’ variable does not change between ‘WILD’ and ‘MODI-  
 54 FIED’ sets, which led to additional investigations. In the words  
 55 of one expert: “Only by using CrossVis were we able to find that  
 56 an increase in mean area of pores caused a decrease in density

57 of pores for the ‘MODIFIED’ set, and vice versa for the ‘WILD’  
 58 set, causing both sets to maintain approximately the same ratio  
 59 of pore area to total area.”

60 The scientists also found the ability to view categorical and  
 61 numerical data in a single system, especially with the image  
 62 representations, particularly helpful in distinguishing between  
 63 the two different image categories. The ability to select images  
 64 of one category and see the separation in values on the other nu-  
 65 matical axes was key to their identification of the most sensitive  
 66 parameters for the ANN classification process.

67 The image visualization capabilities in CrossVis were added  
 68 after the previously mentioned publication on the ANN diatom  
 69 image classification project. By integrating an image-based  
 70 PCP axis and the linked image view panel, scientists noted in-  
 71 creased productivity because they didn't have to rely on a sepa-  
 72 rate image viewer and manually link lines in the PCP to indi-  
 73 vidual images. They felt that the image views helped to supple-  
 74 ment the mostly quantitative analysis, especially for investigat-  
 75 ing outliers and visually exploring specific image features (e.g.,  
 76 pores, skeleton structures).

77 By viewing their data in new visual representations, scientists  
 78 were encouraged to consider their data from fresh perspectives,  
 79 which led to more creative analytical discourse. After an initial  
 80 training session to help them understand the views, they were  
 81 free to explore the data using direct interaction techniques. Be-  
 82 cause it was no longer a requirement to manually run scripts to  
 83 query the data, CrossVis increased the number of combinations  
 84 they could consider and expanded their depth of understanding  
 85 of the data. Although the total time they invested in data anal-  
 86 ysis may have been about the same as before, the direct query  
 87 capabilities allowed them to consider more of the data and find  
 88 unforeseen patterns; two main issues we have observed in most  
 89 of our collaborations with domain scientists.

### 90 6.2. Limitations and Future Enhancements

91 Although CrossVis works well in several scenarios, limita-  
 92 tions and opportunities for improvement remain. In the remain-  
 93 der of this section, we describe the most salient observations.

#### 94 6.2.1. Providing More Active Computational Guidance

95 CrossVis relies on the user to form hypotheses and drive most  
 96 of the analysis. Statistical analytics provide hints at potentially  
 97 significant trends, but these are predominately passive requiring  
 98 the user to see a visual indicator and follow the lead to deeper  
 99 investigation. In the current application, the statistical analytics  
 100 primarily support the level of detail rendering scheme.

101 To improve CrossVis, we envision the integration of auto-  
 102 mated machine learning techniques that suggest key patterns  
 103 and more actively guide the user during the analysis process.  
 104 We have already explored the integration of multiple linear re-  
 105 gression in our prior work with the MDX system [11]. We are  
 106 planning new methods to couple CrossVis visualizations with  
 107 other approaches, such as dynamic time warping to highlight  
 108 similar time series trends and other anomaly detection routines.  
 109 These techniques could be designed to capture user interactions,  
 110 either implicitly or explicitly, and feed the examples as labels to  
 111 the automated algorithms for highlighting similar patterns in the  
 112 full data volume, especially in unseen sections of the data.

### 1    6.2.2. Increasing Scalability

2    CrossVis is a standalone application designed to run on lap-  
 3    tops and workstations, which works well for typical data sets  
 4    we have encountered but doesn't scale smoothly to data sets that  
 5    exceed 10s of gigabytes. To support larger data sets, like those  
 6    generated from computer simulations in climate science, we are  
 7    investigating a distributed approach where the full data set re-  
 8    sides on remote high performance systems. Analytical process-  
 9    ing and query operations will be executed on the remote system  
 10   and only the results will be transmitted to the client, where the  
 11   interactive data visualization components operate. Such a sys-  
 12   tem could enable web-based visualizations of larger and more  
 13   complex data sets, thereby improving accessibility and main-  
 14   tenance.

### 15    6.2.3. Refining the PCP Axis Representations

16    We evaluated several variations of the overall CrossVis de-  
 17   sign before settling on the current version. One objective was  
 18   to pack as much useful information as possible into the axis  
 19   representations without overloading the user. For example, the  
 20   symmetrical presentation of histograms on the numerical axes  
 21   involved evaluations of alternative solutions using various vi-  
 22   sual feature mappings. We recognize that there is additional  
 23   room for exploring enhanced axis representations. Some of  
 24   these (e.g., reclaiming the axis bar interior for time series plots  
 25   or other statistical metrics, representing polyline connections  
 26   for categories) are mentioned in the previous sections. Other  
 27   future work involves supporting vertical PCP axis orientations,  
 28   as opposed to strictly horizontal, which may be easier for deci-  
 29   phering histogram and summary views. We are also interested  
 30   in new ways to represent multiple focus regions using the fo-  
 31   cus+context axis scaling technique. This expansion could en-  
 32   able exploration of multiple focus ranges on the same axis as  
 33   well as cascading axes that drill down into specific ranges.

34    We are actively expanding the concept of embedding addi-  
 35   tional multivariate views into PCPs in a manner similar to the  
 36   bivariate axis scatterplots. Alternatives include leveraging ad-  
 37   dditional visual features (e.g., color, size, shape) to encode addi-  
 38   tional variables, providing three-dimensional views of volumes  
 39   with plane slicing capabilities, and embedding additional visu-  
 40   alization techniques such as miniature treemaps or graphs. Us-  
 41   ing a modular axis design, a wide range of possibilities exist.

### 42    6.3. Reflections on the Interdisciplinary Design Process

43    CrossVis was developed in close collaboration with mate-  
 44   rials scientists following a participatory design process where  
 45   our collaborators were co-designers and primary users of the  
 46   system. We met with them regularly to discuss new develop-  
 47   ments, mock-ups, and evaluated use of the tool with their data  
 48   sets. The design was also influenced by prior collaborations  
 49   with scientists in climate, cybersecurity, health care, and other  
 50   domains. Such projects benefit from clear objectives related  
 51   to domain specific challenges, which ground feature develop-  
 52   ments and help fulfill the central promise of data visualization;  
 53   bringing the latest data visualization advances to data rich do-  
 54   mains where they are needed.

55    By integrating experts from data visualization and other do-  
 56   mains, interdisciplinary projects also benefit from a diversity of  
 57   ideas. Domain experts teach data visualization experts about  
 58   their analysis procedures, which often stimulates new interac-  
 59   tive visualization designs. On the other hand, data visualization  
 60   experts enlighten domain experts about new data visualization  
 61   techniques and trends, thereby bridging the gap between theo-  
 62   retical data visualization and practical applications. Our expe-  
 63   rience is that both sides welcome the engagements and the results  
 64   are almost always positive.

65    Data visualization experts often strive to generalize their  
 66   techniques to maximize impacts on broader endeavors. At the  
 67   same time, it can be hard to avoid degrading performance in the  
 68   motivating domain specific scenario. For example, CrossVis  
 69   supports reading data from CSV files while also supporting cus-  
 70   tom file formats for our domain collaborations to enhance per-  
 71   formance with their data sets. If generalization impacts domain  
 72   specific performance, the payoff must be carefully weighed.  
 73   Perhaps by sacrificing some performance, domain experts will  
 74   recognize the opportunity to use the technique with other data  
 75   sets and the sacrifice will be welcomed. In such cases where per-  
 76   formance degradation is unacceptable, the team may consider  
 77   deploying specific builds of the tools; one for general purpose  
 78   use and others for specific applications.

79    Interdisciplinary projects are challenging for data visualiza-  
 80   tion experts because they must exhibit agility in learning new  
 81   domains and concepts. By including the experts in the develop-  
 82   ment team to validate assumptions and algorithmic decisions,  
 83   the knowledge gap can be spanned and in the process data visu-  
 84   alization researchers gradually gain a more comprehensive un-  
 85   derstanding of the domain. Furthermore, this process engages  
 86   domain experts and can increase adoption rates as they share  
 87   the techniques with others in their field. Perhaps the greatest  
 88   reward from such an activity is when the experts publish results  
 89   found with the new techniques, an outcome that essentially val-  
 90   idates the effectiveness of the approach and one that we experi-  
 91   enced the development of CrossVis with our materials science  
 92   collaborators.

## 93    7. Conclusion

94    CrossVis extends the PCP concept by adding a range of new  
 95   axis representation techniques, interactions, and scalability ex-  
 96   tensions to enable large scale, multivariate exploration of het-  
 97   erogeneous data. The overall design requirements were derived  
 98   from key challenges uncovered during close interactions with  
 99   experts in a variety of fields using prior multivariate visualiza-  
 100   tion tools. The resulting system helps scientists look at more  
 101   of their data and find new, and often unexpected, insights; two  
 102   needs that are common in most scientific domains.

103   By working close with materials scientists to develop  
 104   CrossVis, we observed how it improved the depth of their anal-  
 105   ysis as well as certain limitations and opportunities for future  
 106   improvement. Important patterns were uncovered by domain  
 107   experts who used CrossVis to interpret a neural network pro-  
 108   cess for classifying microscopic images in a genetic engineer-  
 109   ing study. In the process of developing and applying the sys-

tem, we gained additional insight into interdisciplinary collaborations to develop data science systems, particularly visual analysis systems. Another scientific use case described in the current work involves exploring a complex hurricane observation data set and demonstrates the suitability of CrossVis in general purpose data exploration.

CrossVis is a general purpose visual analytics framework that has also proven useful in other fields, such as cybersecurity, earth system modeling, health care, and algorithmic performance analysis. In the future, we will continue to expand and apply CrossVis to other data rich domains while continuing our interdisciplinary development strategy.

## Acknowledgments

The authors wish to express our appreciation to the reviewers and editor for their excellent feedback. The authors also thank Dr. Alison Pawlicki of Oak Ridge National Laboratory's (ORNL) Center for Nanophase Materials Sciences (CNMS) for acquiring the raw Scanning Electron Microscope (SEM) diatom imagery used in the CrossVis evaluation. The HURDAT2 data set is a product of the NOAA National Hurricane Center.

This research was sponsored by the U.S. Department of Energy (DOE) under the Scientific Discovery through Advanced Computing (SciDAC) RAPIDS project and the Laboratory Directed Research and Development Program of ORNL, managed by UT-Battelle, LLC, for the U.S. DOE. The SEM analysis was partially conducted at ORNL's CNMS, which is a DOE Office of Science User Facility.

## References

- [1] Liu, S, Maljovec, D, Wang, B, Bremer, P, Pascucci, V. Visualizing high-dimensional data: Advances in the past decade. *IEEE Transactions on Visualization and Computer Graphics* 2017;23(3):1249–1268. doi:10.1109/TVCG.2016.2640960.
- [2] Thomas, JJ, Cook, KA. A visual analytics agenda. *IEEE Computer Graphics and Applications* 2006;26(1):10–13. doi:10.1109/MCG.2006.5.
- [3] Keim, D, Kohlhammer, J, Ellis, G, Mansmann, F, editors. Mastering the Information Age: Solving Problems with Visual Analytics. Goslar, Germany: Eurographics Association; 2010.
- [4] Inselberg, A. The plane with parallel coordinates. *The Visual Computer* 1985;1(4):69–91. doi:10.1007/BF01898350.
- [5] Wegman, EJ. Hyperdimensional data analysis using parallel coordinates. *Journal of the American Statistical Association* 1990;85(411):664–675. doi:10.2307/2290001.
- [6] Inselberg, A. Parallel coordinates: Interactive visualization for high dimensions. In: Zudilova-Seinstra, E, Adriaansen, T, Liere, R, editors. Trends in Interactive Visualization. London, UK: Springer-Verlag; 2009, p. 49–78. doi:10.1007/978-1-84800-269-2\_3.
- [7] Heinrich, J, Weiskopf, D. State of the art of parallel coordinates. In: Proceedings of Eurographics 2013 - State of the Art Reports. The Eurographics Association; 2013, doi:10.2312/conf/EG2013/stars/095–116.
- [8] Johansson, J, Forsell, C. Evaluation of parallel coordinates: Overview, categorization and guidelines for future research. *IEEE Transactions on Visualization and Computer Graphics* 2016;22(1):579–588. doi:10.1109/TVCG.2015.2466992.
- [9] Wang, J, Liu, X, Shen, H, Lin, G. Multi-resolution climate ensemble parameter analysis with nested parallel coordinates plots. *IEEE Transactions on Visualization and Computer Graphics* 2017;23(1):81–90. doi:10.1109/TVCG.2016.2598830.
- [10] Steed, CA, Ricciuto, DM, Shipman, G, Smith, B, Thornton, PE, Wang, D, et al. Big data visual analytics for exploratory earth system simulation analysis. *Computers & Geosciences* 2013;61:71–82. doi:10.1016/j.cageo.2013.07.025.
- [11] Steed, CA, Swan, JE, Jankun-Kelly, TJ, Fitzpatrick, PJ. Guided analysis of hurricane trends using statistical processes integrated with interactive parallel coordinates. In: IEEE Symposium on Visual Analytics Science and Technology 2009, p. 19–26. doi:10.1109/VAST.2009.5332586.
- [12] Choi, H, Lee, H, Kim, H. Fast detection and visualization of network attacks on parallel coordinates. *Computers & Security* 2009;28(5):276–288. doi:10.1016/j.cose.2008.12.003.
- [13] Wang, WB, Huang, ML, Lu, L, Zhang, J. Improving performance of forensics investigation with parallel coordinates visual analytics. In: Proceedings of the IEEE International Conference on Computational Science and Engineering. 2014, p. 1838–1843. doi:10.1109/CSE.2014.337.
- [14] Boogaerts, T, Tranchevent, L, Pavlopoulos, GA, Aerts, J, Vandewalle, J. Visualizing high dimensional datasets using parallel coordinates: Application to gene prioritization. In: Proceedings of the IEEE International Conference on Bioinformatics Bioengineering. 2012, p. 52–57. doi:10.1109/BIBE.2012.6399706.
- [15] Keefe, D, Ewert, M, Ribarsky, W, Chang, R. Interactive coordinated multiple-view visualization of biomechanical motion data. *IEEE Transactions on Visualization and Computer Graphics* 2009;15(6):1383–1390. doi:10.1109/TVCG.2009.152.
- [16] Caat, MT, Maurits, NM, Roerdink, JBTM. Design and evaluation of tiled parallel coordinate visualization of multichannel EEG data. *IEEE Transactions on Visualization and Computer Graphics* 2007;13(1):70–79. doi:10.1109/TVCG.2007.9.
- [17] Qu, H, Chan, W, Xu, A, Chung, K, Lau, K, Guo, P. Visual analysis of the air pollution problem in hong kong. *IEEE Transactions on Visualization and Computer Graphics* 2007;13(6):1408–1415. doi:10.1109/TVCG.2007.70613.
- [18] Siirtola, H. Direct manipulation of parallel coordinates. In: Proceedings of the IEEE Conference on Information Visualization. 2000, p. 373–378. doi:10.1109/IV.2000.859784.
- [19] Hauser, H, Ledermann, F, Doleisch, H. Angular brushing of extended parallel coordinates. In: Proceedings of IEEE Symposium on Information Visualization. 2002, p. 127–130. doi:10.1109/INFVIS.2002.1173157.
- [20] Cleveland, WC, McGill, ME. Dynamic Graphics for Statistics. Boca Raton, FL, USA: CRC Press, Inc.; 1988. doi:10.1214/ss/1177013104.
- [21] Claessen, JHT, van Wijk, JJ. Flexible linked axes for multivariate data visualization. *IEEE Transactions on Visualization and Computer Graphics* 2011;17(12):2310–2316. doi:10.1109/TVCG.2011.201.
- [22] Roberts, JC. Exploratory visualization with multiple linked views. In: Exploring Geovisualization. Elsevier; 2005, p. 159–180. doi:10.1016/B978-008044531-1/50426-7.
- [23] Yuan, X, Guo, P, Xiao, H, Zhou, H, Qu, H. Scattering points in parallel coordinates. *IEEE Transactions on Visualization and Computer Graphics* 2009;15(6):1001–1008. doi:10.1109/TVCG.2009.179.
- [24] Cuzzocrea, A, Zall, D. Parallel coordinates technique in visual data mining: Advantages, disadvantages and combinations. In: Proceedings of the International Conference on Information Visualisation. 2013, p. 278–284. doi:10.1109/IV.2013.96.
- [25] Zhou, L, Weiskopf, D. Indexed-points parallel coordinates visualization of multivariate correlations. *IEEE Transactions on Visualization and Computer Graphics* 2018;24(6):1997–2010. doi:10.1109/TVCG.2017.2698041.
- [26] Kosara, R, Bendix, F, Hauser, H. Parallel sets: interactive exploration and visual analysis of categorical data. *IEEE Transactions on Visualization and Computer Graphics* 2006;12(4):558–568. doi:10.1109/TVCG.2006.76.
- [27] Fernstad, SJ, Johansson, J. A task based performance evaluation of visualization approaches for categorical data analysis. In: Proceedings of the International Conference on Information Visualisation. 2011, p. 80–89. doi:10.1109/IV.2011.92.
- [28] Vosough, Z, Hogräfer, M, Royer, LA, Groh, R, Schulz, HJ. Parallel hierarchies: A visualization for cross-tabulating hierarchical categories. *Computers & Graphics* 2018;76:1–17. doi:10.1016/j.cag.2018.07.009.
- [29] Johansson, J, Ljung, P, Jern, M, Cooper, M. Revealing structure within clustered parallel coordinates displays. In: Proceedings of the IEEE Symposium on Information Visualization, 2005. INFOVIS 2005.

- 1 2005, p. 125–132. doi:10.1109/INFVIS.2005.1532138.
- 2 [30] Palmas, G, Bachynskyi, M, Oulasvirta, A, Seidel, HP, Weinkauf, T.  
3 An edge-bundling layout for interactive parallel coordinates. In: 2014  
4 IEEE Pacific Visualization Symposium. 2014, p. 57–64. doi:10.1109/  
5 PacificVis.2014.40.
- 6 [31] Novotný, M, Hauser, H. Outlier-preserving focus+context visualization  
7 in parallel coordinates. *IEEE Transactions on Visualization and Computer  
8 Graphics* 2006;12(5):893–900. doi:10.1109/TVCG.2006.170.
- 9 [32] Andrienko, G, Andrienko, N. Blending aggregation and se-  
10 lection: Adapting parallel coordinates for the visualization of large  
11 datasets. *The Cartographic Journal* 2005;42(1):49–60. doi:10.1179/  
12 000870405X57284.
- 13 [33] Janetzko, H, Stein, M, Sacha, D. Enhancing parallel coordinates: Sta-  
14 tistical visualizations for analyzing soccer data. In: Proceedings of the  
15 Visualization and Data Analysis Conference. 2016, p. 1–8. doi:10.2352/  
16 ISSN.2470-1173.2016.1.VDA-486.
- 17 [34] Blaas, J, Botha, C, Post, F. Extensions of parallel coordinates for  
18 interactive exploration of large multi-timepoint data sets. *IEEE Trans-  
19 actions on Visualization and Computer Graphics* 2008;14(6):1436–1451.  
20 doi:10.1109/TVCG.2008.131.
- 21 [35] Sansen, J, Richer, G, Jourde, T, Lanlanne, F, Auber, D, Bourqui, R.  
22 Visual exploration of large multidimensional data using parallel coordi-  
23 nates on big data infrastructure. *Informatics* 2017;4(3):1–21. doi:10.  
24 3390/informatics4030021.
- 25 [36] Richer, G, Sansen, J, Lanlanne, F, Auber, D, Bourqui, R. Enabling hi-  
26 erarchical exploration for large-scale multidimensional data with abstract  
27 parallel coordinates. In: Proceedings of the International Workshop on  
28 Big Data Visual Exploration and Analytics. 2018, p. 8 pp.
- 29 [37] Mackinlay, J. Automating the design of graphical presentations of  
30 relational information. *ACM Trans on Graphics* 1986;5(2):110–141.  
31 doi:10.1145/22949.22950.
- 32 [38] Salvador, S, Chan, P. Fastdtw: Toward accurate dynamic time warping  
33 in linear time and space. In: Proceedings of the ACM KDD Workshop on  
34 Mining Temporal and Sequential Data. 2004, p. 70–80.
- 35 [39] Rensink, RA. Change detection. *Annual Review of Psychol-  
36 ogy* 2002;53:245–577. doi:10.1146/annurev.psych.53.100901.  
37 135125.
- 38 [40] Hamm, CE, Merkel, R, Springer, O, Jurkoje, P, Maier, C, Prechtel,  
39 K, et al. Architecture and material properties of diatom shells provide ef-  
40 fective mechanical protection. *Nature* 2003;421:841–843. doi:10.1038/  
41 nature01416.
- 42 [41] Delalat, B, Sheppard, VC, Ghaemi, SR, Rao, S, Prestidge, CA,  
43 McPhee, G, et al. Targeted drug delivery using genetically engineered  
44 diatom biosilica. *Nature Communications* 2015;6:8791. doi:10.1038/  
45 ncomms9791.
- 46 [42] Landsea, CW, Franklin, JL. Atlantic hurricane database uncertainty  
47 and presentation of a new database format. *Monthly Weather Review*  
48 2013;141(10):3576–3592. doi:10.1175/MWR-D-12-00254.1.
- 49 [43] Evans, JL, Hart, RE. Objective indicators of the life cycle evolution of extratropical transition for atlantic tropical cy-  
50 clones. *Monthly Weather Review* 2003;131(5):909–925. doi:10.1175/  
51 1520-0493(2003)131<0909:OITLC>2.0.CO;2.
- 53 [44] Trofimov, AA, Pawlicki, AA, Borodinov, N, Mandal, S, Mathews,  
54 TJ, Hildebrand, M, et al. Deep data analytics for genetic engineering of  
55 diatoms linking genotype to phenotype via machine learning. *npj Com-  
56 putational Materials* 2019;5(1):67. doi:10.1038/s41524-019-0202-3.