

Towards Quantitative Visual Analytics with Structured Brushing and Linked Statistics

S. Radoš¹, R. Splechtna¹, K. Matković¹, M. Duras², E. Gröller³, and H. Hauser⁴

¹VRVis Research Center in Vienna, Austria

²AVL Zagreb, Croatia

³Vienna University of Technology, Austria

⁴University of Bergen, Norway

Abstract

Until now a lot of visual analytics predominantly delivers qualitative results—based, for example, on a continuous color map or a detailed spatial encoding. Important target applications, however, such as medical diagnosis and decision making, clearly benefit from quantitative analysis results. In this paper we propose several specific extensions to the well-established concept of linking and brushing in order to make the analysis results more quantitative. We structure the brushing space in order to improve the reproducibility of the brushing operation, e.g., by introducing the percentile grid. We also enhance the linked visualization with overlaid descriptive statistics to enable a more quantitative reading of the resulting focus+context visualization. Additionally, we introduce two novel brushing techniques: the percentile brush and the Mahalanobis brush. Both use the underlying data to support statistically meaningful interactions with the data. We illustrate the use of the new techniques in the context of two case studies, one based on meteorological data and the other one focused on data from the automotive industry where we evaluate a shaft design in the context of mechanical power transmission in cars.

1. Introduction

Interactive visual data exploration and analysis has become an indispensable complement to automatic analysis techniques. Still, we see quite often that analysts prefer computational techniques for a few important reasons.

First, focus+context visualization is often only qualitative by nature. As compared to the context, the data subset(s) in focus are shown in a different color, or in another visualization style [Hau06], resulting in only approximate readings of such views. In certain application cases, including decision making, for example, “hard”, quantitative facts are often useful (think of a “no go”-decision, if the p-value of a statistical test is above a predefined threshold).

Another reason is that results from interactive procedures, like most of traditional visual analytics, often seem to lack a sufficiently good reproducibility. Redoing a visual analytics session, for example, where linking and brushing is used, will most likely not result in exactly the same result. This is due to small variations in the placement of the brushes, for example. A recent study by Kandogan et al. [KBHP14],

based on 34 in-depth interviews, documents this situation clearly in the context of business intelligence. It seems obvious that extensions to visual analytics, which enable reproducible and quantitative results, may become key to a further strengthened deployment of interactive visualization in analytics applications.

In this paper, we contribute several specific extensions to the well-established concept of linking and brushing in coordinated, multiple views. This amounts to the first major collection of techniques targeted specifically towards reproducible and quantitative visual analytics.

With respect to brushing, we describe several particular extensions, including percentile brushing and Mahalanobis brushing, i.e., two techniques that support reproducibility. In abstract terms, we discuss the brushing space and how it can be structured for improved reproducibility.

With respect to linking, we introduce further extensions, including the integration of descriptive statistics, which enables a quantitative reading of linked views with focus+context visualization. We also support the user during

the visual analysis by reducing mental load during brushing, by allowing him to record a brush path. The brush can then be animated, i.e., reproduced repeatedly, and the user can pay all attention to the linked views. Additionally, we provide animated transitions in linked views in combination with a descriptive statistics overview. We also introduce the relative difference plot as a novel way of describing the history of linked data statistics. We illustrate the use of the new techniques in the context of two case studies, one based on meteorological data and the other one focused on data from the automotive industry. Further we explain, in which way our results are reproducible and quantitative. We conclude by discussing benefits and limitations of the current approach and outlining selected plans for future work.

2. Related Work

The concept of linking and brushing is key to interactive visual analysis (IVA) [WH14, KH13]. It is modeled as an interactive and iterative method to reveal insight into large and multi-faceted datasets. The term *brushing* was defined by Becker and Cleveland [BC87] and different brush shapes were proposed, including rectangles and circles [CM88]. Martin and Ward researched N-dimensional, multiple, fuzzy, and composite brushes. They employed brushing for the analysis of multi-dimensional data in the XmdvTool [War94]. The user configures composite brushes by applying logical operations and expressions (e.g., with AND, OR, XOR, and NOT) [MW95]. Doleisch et al. [DGH03] introduced a feature definition language for the specification of multi-dimensional and/or complex features, using logical combinations of brushes in coordinated, multiple views. The concept of compound brushing, developed by Chen [Che03], helps in describing many existing brushing techniques and it is also useful for exploring new techniques. Animation is also used in interactive visual applications for helping users to follow changes in the visualization [HR07, ROC97, BPF14]. However, animation must be used with caution, since it could lead to perceptual errors, and can slow down the analysis [RFF*08].

Brushing techniques are commonly categorized into three groups, according to the space in which the selection is being performed: *screen*, *data* and *structure* brushes [FWR00]. While screen-space techniques traditionally limit the shape of a brush to two dimensions, data-space techniques permit brushes with dimensionality greater than two. For example, the N-dimensional brush [War94] provides insight into a spatial relationship over N dimensions. The third group extends the brush metaphor to structures. It encompasses structure-space techniques [FWR00] which are based on structural relationships between data points, such as clusterings, orderings, groupings, etc. Structure-space brushing techniques are particularly useful for datasets with natural and imposed structures. In this paper, we introduce the Mahalanobis brush as a new structure-based brushing technique.

It takes the underlying data distribution into account, while specifying the brush in screen space. Traditionally, brushing has been performed *unconstrained* – brushes can be created anywhere in the view and the analyst can move or resize them freely. As an addition to the free (unconstrained) brushing, and to support reproducibility, we now introduce an alternative mode that we term *constrained brushing*.

Visual analytics, especially the field of analytic provenance, has been interested in reproducible methods for several years. Examples include the work of Gotz et al. [GWL*10] on history keeping in the Harvest system and the work of Silva et al. [SFSA10] on provenance support in the VisTrails system. This application systematically maintains provenance in the data exploration process by capturing all the steps which have been taken during an interactive visualization session. Yang et al. [YXRW07] developed the Nugget Management System (NMS) for the housekeeping of user findings, called “nuggets”. They consequently manage nuggets, e.g., by organizing them in an intuitive manner. These approaches, in general, focus on the reproducibility of the whole analysis session. In our work we primarily focus on the reproducibility of the brushing operation itself, being an important part of the overall interactive visual analysis.

Up to now, not much related work is available on quantitative visual analytics. Chen [Che03] showed how to enable analytical filtering through the addition of the quantile range-filter for one variable to validate or filter data selections. In our work, we contribute constrained brushing using a percentile-derived grid as a related extension. This supports analytical tasks that are ranking-based (instead of value-based). Kehrer et al. [KFH10] integrated statistical aggregates along selected, independent data dimensions in a framework of coordinated, multiple views. Brushing particular statistics, the analyst can investigate data characteristics such as trends and outliers. Haslett et al. [HBC*91] introduced the ability to show the average of the points that are currently selected by the brush. Based on this idea, summarizations of the data are commonly used as a representative information for clusters in hierarchically organized large datasets [Shn92, FWR00]. We also use summarizations, in the context of brushing, and show several descriptive statistics in linked views, in a table, as overlay or in combination with traces from brushing.

3. Quantitative and Reproducible Linking&Brushing

In the following, we first discuss in which way standard linking and brushing is qualitative (as opposed to quantitative analytics) and why there are challenges with reproducibility. Then, we provide a detailed description of our contributions. In order to illustrate the new techniques, we visualize meteorological data from about 300 weather stations in California [NOA14]. This dataset contains geographic information and temperature and precipitation values.

The qualitative character is, in fact, a critical strength of

visual analytics, since it naturally harmonizes with the integration of a human in the analysis loop. After spotting a data subset of interest in the visualization, interactive brushing is used to mark up this data subset, directly and interactively in the view. All linked views get immediately updated and a consistent focus+context visualization is generated. However, linking&brushing provides mostly qualitative insight only, and it is not 100% reproducible.

Firstly, the brushed data subset is always highlighted in all linked views, while the rest of the dataset is shown as context (differently colored, smaller, accumulated, etc.). This results in only approximate readings of such views. The following example illustrates this situation: A typical result from standard IVA is something like “Using linking&brushing, we see that low values of dimension x [as brushed in view A] are correlated with high values of dimension y as apparent in the linked focus+context visualization [view B].” The meaning of “low” and “high” remains vague/relative. A computational data analysis would usually put a number on such a relation – maybe a Pearson correlation coefficient. Clearly, the brushed and linked visualization also provides additional information about the relation between x and y . It indicates if the relation is linear or not, for example, and this is highly useful. For decision making, however, “hard”, quantitative facts are often more valuable.

Secondly, it is typical in brushing that users select freely what they deem interesting. Considering a rectangular brush on a scatterplot as an example, the user chooses an arbitrary point as the top-left corner of the brush and then extends the brush-rectangle to the desired size. Due to the high-resolution of the visualization, and the corresponding interface technology it is highly improbable that an attempt to exactly recreate such a brush will succeed. This results in challenges with respect to the reproducibility of exploratory visualizations. Up to now, a possible way for repeating an exploratory task was to save the complete history, by using a provenance management system such as Vis-Trails [SFS10]. In our work, we think about the reproduction of IVA results after they have been documented, e.g., in a report. A typical example would be the following “We look at the screenshot of view B and we see that the highlighted data are linearly correlated. From the given textual description we know that the 25% lowest values of dimension x were selected in view A. After an update of the dataset (with additional data points, for example), we wish to swiftly reproduce the reported analysis, i.e., to brush the 25% lowest values of dimension x in view A and compare the updated linked view B with the screenshot in the report.” With standard IVA, this is only approximatively possible. Most automated, computational approaches, however, will score very well on reproducibility.

In the following, we describe how we structure the brushing space in order to make brushing more reproducible. Then we describe how we support the interpretation of linked

views by integrating descriptive statistics. Finally, after presenting the case study with domain experts, we discuss possibilities for further development.

3.1. Structured Brushing

In addition to standard brushing, which we call unconstrained (unstructured), we suggest as a complement constrained and automatic brushing. The brushing space is structured with respect to the anchoring of the brush, its extent, and the movement of the brush. Table 1 describes examples of possible solutions for (partially) constrained and (semi-)automatic brushing. Furthermore, two new brushes, the *percentile brush* the *Mahalanobis brush* are two concrete suggestions of how to realize an advanced brush, based on the structured/informed brushing space (see below).

Snap-to-Grid Brush. As in drawing programs, we can introduce a snap-to-grid option for brushing. This functionality is a useful mechanism to confine brushing to reproducible shapes that also can be interpreted quantitatively. A regular grid and the snap-to-grid functionality also works for categorical data. We can require that brushes are anchored at grid points and we can confine the extent of brushes to correspond to grid cells. For example, if we define a regular 4×4 grid, and we create a brush in the bottom-left grid cell, then we instantly know quantitatively that we have selected the $[0\%, 25\%]$ interval on the x axis, and the $[0\%, 25\%]$ interval on the y axis (Figure 1 brush (a)). If we constrain also the movement of the brush to allow only a vertical movement and activate the snap-to-grid functionality, only predefined intervals will be taken by the moving brush. Even an imprecise interaction in the brushed view will result in an exact, quantitative brush movement. This allows the user to concentrate on the linked view, knowing exactly which intervals are selected anyway, without the need to pay attention to value-accurate brushing.

Percentile Grid Brush. With the help of descriptive statistics, it is usual in (computational) data analysis to either do a value-based analysis, or a rank-based analysis. The latter could, for example, be enabled through quantile filters [Che03] or statistical estimators [KFH10]. Hence, we suggest to also provide brushing opportunities which match these analytics perspectives. Using a regular grid corresponds to a value-oriented perspective. Alternatively, often a rank-based analytics perspective is also very useful. An example here would be the Spearman correlation [Spe87]. Instead of selecting all items that correspond to a certain range of values, we are interested in a certain number of data items, e.g., the top 10% of all data items. If we define the grid so that each division on an axis contains a certain percentage of the items, we create a *percentile grid*. Each vertical and each horizontal strip of the scatterplot shown in Figure 1(B), for example, contains exactly 25% of the data. Brushing in the snap-to-grid mode has a different meaning then. Brushing all left-most cells, snapped to a 25% percentile grid, we know,

	<i>Brush Anchoring</i>	<i>Brush Extent</i>	<i>Brush Movement</i>
<i>unconstrained</i>	The user initiates the brush anywhere in the view, for example, on a scatterplot by specifying the top-left corner of a rectangular brush at an arbitrary position.	Any extent of the brush is possible and brush boundaries can be modified freely.	The brush can be moved freely.
<i>constrained</i>	A “snap-to-grid” functionality is used to constrain the anchoring of brushes to grid vertices.	The size of the brush can be adapted in discrete, predefined steps only.	If moved, the brush assumes only grid-aligned positions.
<i>automatic</i>	The user specifies a particular brush parameter, for example, a data-related property so that the brush is positioned automatically.	The brush resizes itself automatically due to certain constraints, for example, maintaining that a certain number of data points is selected.	The brush moves automatically, for example, following a user-defined animation procedure.

Table 1: Structuring the brushing space.

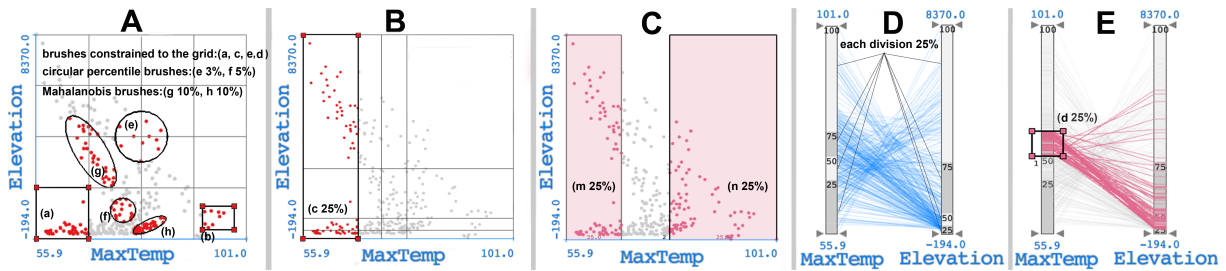


Figure 1: Overview of the extensions for structured brushing. **A:** A scatterplot with a regular 4×4 grid (value-based). **B:** A scatterplot with a percentile 4×4 grid (rank-based). **C:** A scatterplot with the 25% percentile brushes. **D and E:** Parallel coordinates with a 4×4 grid (rank-based) which show the data distributions in each dimension.

again quantitatively, that we have selected the 25% lowest values with respect to the dimension that is mapped to the horizontal axis (Figure 1 brush (c)). Moving the brush along the grid from left to right, then in each step the brush would select the next 25% of all items. The grid implicitly reveals also some insight into the data distribution as illustrated in Figure 1(D). The analyst may benefit from the grid even if the constrained brushing is not enabled. The grid can, e.g., assist the navigation of the brush over the presented data.

Percentile Brush. The *percentile brush* constrains the extent so that the brush always contains a predefined number of items, like 10%. The brush can be moved freely, or snapped to a conventional grid, or to a percentile grid also. When moved, the extent of the brush is adapted continuously so that it always selects the predefined number of items. In a scatterplot, we suggest two standard shapes for realizing percentile brushes, i.e., the rectangular and the circular percentile brush. The rectangular brush is easy to interpret in the scatterplot. When creating the brush, the user can decide whether the brush considers the data distribution in the horizontal or in the vertical dimension. Figure 1(C) shows two 25% percentile brushes (**m**, **n**) created over the horizontal dimension. The brush (**c**) in Figure 1(B) selects the lowest 25% using the snap-to-grid, which is equivalent to the brush (**m**) in this case. Note however, that the brush (**m**) can be moved

freely in the horizontal dimension, while the brush (**c**) can be moved only between grid positions. The circular percentile brush selects a specified number of items in the vicinity of a user-specified point, i.e., from the center of the brush, see Figure 1(A) and brushes (**e**) and (**f**). If the snapped circular percentile brush (**e**) is moved it jumps from one grid vertex to another one (with the center snapped to a grid vertex). In a parallel coordinates plot we use only the one dimensional percentile brush over the individual axes Figure 1(E).

Mahalanobis Brush. The percentile brush changes its extent, but keeps its shape, the circular brush changes its radius but remains circular. Dependent on the data distribution, this is sometimes not the most useful behavior. The Mahalanobis distance [Mah36] is a metric, which takes the data distribution into account. The Mahalanobis distance for two points \vec{x} and \vec{y} , both from the same distribution with covariance matrix \mathbf{C} , is given by $((\vec{x} - \vec{y})^T \mathbf{C}^{-1} (\vec{x} - \vec{y}))^{\frac{1}{2}}$. In a two-dimensional case (as in the scatterplot), equidistant lines around a point will be ellipses with the axes corresponding to the principal component directions of the data. If we compute the percentile brush using the Mahalanobis distance instead of the Euclidean distance we get the Mahalanobis brush and the brush accommodates itself to the underlying data distribution. Depending on the user preferences, the data distribution is calculated from the whole

dataset or from a local data subset D . The size of this reference subset is given as the percentage of data points from the whole data set, parameter m_d . Depending on this value, the Mahalanobis brush will be more or less sensitive to the distribution of the data near the center of the brush. The whole algorithm for computing the Mahalanobis brush is shown in Algorithm 1. Figure 1 (A) shows two 10% Mahalanobis brushes (g, h). Note that the shape adapts to the data distribution as the brush is moved. It is still a rank-based brush, selecting always a predefined number of points, of the closest possible number as noted for the percentile brush. Alternatively, we could transform the data space and use the previously explained percentile brush. Such an approach, however, would make the data interpretation more complicated.

```

Data: all data in the horizontal and vertical dimension,
 $\vec{p}$  – mouse position, percentage  $n$  – of all points
to be selected, percentage  $m_d$  – of all points
forming the basis for metric computations, vector
 $\vec{d}$  – of data points closest to point  $\vec{p}$ 
Result:  $\vec{m}$  – vector of all brushed data points
/* Step 1: computing the local
Mahalanobis metric. */
while percentage of points in the subset  $D < m_d$  do
| Increase size of subset  $D$  by adding nearby points;
end
Compute covariance matrix  $C$  using  $D$ ;
 $\vec{d}$  = compute Mahalanobis distances( $\vec{p}, D, C$ );
/* Step 2: aspect ratio and rotation
of the brush ellipse according to
an eigen-analysis of  $\vec{d}$ . */
while percentage of points selected by the brush  $< n$  do
| Increase the size of the brush and associate the
contained data items with  $\vec{m}$ ;
end

```

Algorithm 1: The pseudo-code of the Mahalanobis brush. The main steps are. **Step 1:** We use a rectangle-shaped or circle-shaped area for selecting a data subset D . The initial size of the area varies depending on the distribution around \vec{p} . We start with selecting 1% of all data points. **Step 2:** We use the eigenvalues to define the rotation angle for the ellipse which represents the Mahalanobis brush m_brush .

Animated Brush. Once the user knows, how the brush should be moved in order to analyze the data, an animated brush can be defined. For example, when the user is interested in observing changes in several linked views, the brush has to be moved over the same path repeatedly in order to study possible correlations. The animated brush can save a lot of time in this case. We enable path storing for different brushing techniques. This includes constrained and unconstrained brushing. Two types of a path recordings are considered in this paper. Firstly, the user can freely draw a brush path. As an example, the user creates a constrained brush, snapped to the first cell in the horizontal dimension

of a 10% percentile grid. While the brush is moved horizontally across all adjacent grid cells, the positions of the brush and the brushed data points are saved in each step. Secondly, the user defines the start and the end position for an unconstrained brush, and a number of frames to be generated in-between. The brush is then interpolated linearly. The user can also insert additional key frames and the brush is linearly animated between those. This is a complementary solution, when compared to completely free brushing. The brushing session can be automatically replayed, following exactly the same positions, extents, and brushed data. This allows the user to solely focus on the linked view(s). The scatterplot in Figure 2 (top-left) shows three key frames of the recorded animation. The start key frame and the end key frame are represented with dashed lines. This brush updates its position and moves along the created trajectory as the animation proceeds. The user can pause an animation and move the brush away from the defined path and/or continue the animation from the paused position.

3.2. Quantitative Linked Views

Interactive visual analysis is highly effective if information about relevant relations between different aspects of the data have to be revealed quickly. However, the qualitative insight by linking&brushing is not always the best help. Further, if the relations are complex, it is usually not easy to understand the trends and patterns. Even if we pay full attention to the linked view(s), we still need other methods to support the mental image creation and to quantify the analysis results. With a better understanding of what is happening on the brushing side (cf. extensions as described in section 3.1), we also aim at a better understanding of the linked side. As analysts need quantitative results, and statistics can provide these, a logical step is to enhance the linked views with additional descriptive statistics about the brushed data.

One of the first ideas to support IVA with statistics from the brushed data comes from Haslett et al. [HBC*91]. They computed the average on a local basis and showed the result as *Moving Average* trace added to the Trace View. The center of the data is certainly the most commonly used statistical measure in data analysis. We compute three different center points: the median, the mean, and the midrange (the value exactly between the minimum and the maximum). Additionally we determine the total spread, and the spread based on the standard deviation. Estimating the center and spread, we already have a first useful summarization of the data. Depending on the task, the user configures what is displayed in a view, i.e., she configures the descriptive statistics overview.

In addition to the *Moving Average* trace, we show traces for other common statistics, as shown in the Trace View Figure 2 (top-right, trace-view). The statistics are computed as the brush moves and new points are added to the trace on each position change. This can result in overplotting if unconstrained brushes are used. Additionally the user can con-

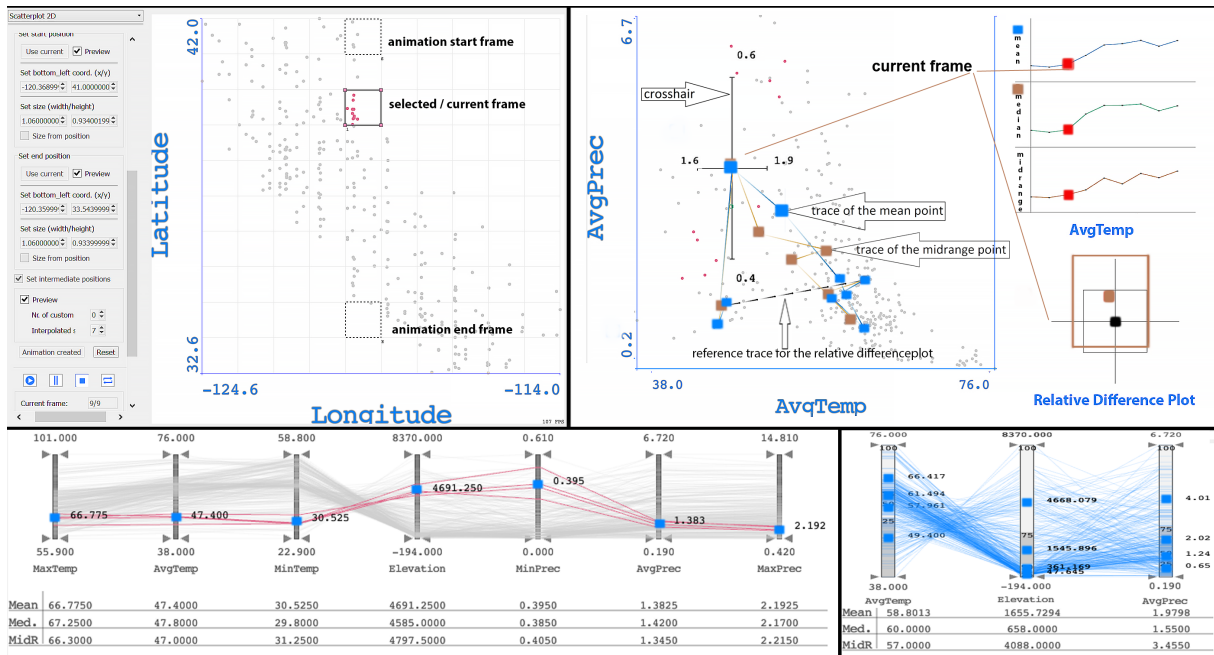


Figure 2: *Top-Left:* An animation brush moves automatically from the start to the end point. *Top-Right:* The user can stop or pause the animation and select any point on the path and the descriptive statistics overview moves and shows updated values. The relative difference plot shows the difference between the actual (brown) and the reference (black) center point and spread. *Bottom-Left:* The mean value for brushed data is shown in parallel coordinates, and the table below the view shows additional statistics. *Bottom-Right:* The mean values are shown for each cell of the three 25% percentile grids in parallel coordinates.

figure the size for the trace buffer. Optionally, we can add a new point to the trace only if its value is different to the previously saved value.

Traces shown in the Trace View are computed for selected dimension only, i.e., in a case of scatterplot either for the horizontal or vertical dimension. We also provide an option to draw the paths of the center points (we call it pathlines) in the view (Figure 2 (top-right, main-view)). To support the comprehension of a position change of the moving brush, we encode the direction in the paths from center points too. Further to support perception and cognition, we overlay a cross-hair depicting the spread in the linked scatterplot. Depending on the user preferences, the pathline is extended as the animation evolves, or the complete path is shown and the cross-hair moves along the path.

In Figure 2 (top-right, main-view) the blue squares indicate the mean center point and the brown squares indicate the midrange point. The cross-hair shows the one standard deviation spread in both directions from the mean center point. In this way, the user sees all points in the focus, but we help her in perceiving the data characteristics. This is exactly what we try to extract from the data when observing the trends.

In order to quantify changes of center points and spreads with respect to a moving brush, we depict them numerically

as well. We display the values for the current brush in a table. As the brush moves, the descriptive statistics overview and the table with the numerical values update accordingly. This is done also for all selected axes in the parallel coordinates, as shown in Figure 2 (bottom).

The psychologist Barbara Tversky [TMB02] found from reviewing nearly 100 studies of animation and visualization, that rich static diagrams are outperforming animations. Following this we provide the possibility to analyze the pathlines showing them in the linked view as a static overlay. The user can analyze statistics computed at different pathline positions by simply clicking on a center point, either in a pathline shown in a scatterplot or in a Trace View.

As the pathline of the mean point in Figure 2 (top-right, main-view) shows, the center points change significantly between the frames in the linked view. Such a change causes sudden jumps in the linked view, and distracts the user. This distraction exacerbates the mental image creation. In combination with the animated brush we propose to animate the crosshair transition in the linked view in order to prevent a distraction of the user. The cross-hair stays at a brush position for some predefined time “hold-time” (few seconds), then it smoothly travels to a new position within some predefined time (usually shorter than the “hold-time”). Visual-

ization of the transition does not only help in eliminating distraction, it also actively amplifies cognition of the trend evolution. However a case study is necessary to quantify the exact impact.

We extended the idea to support the mental image creation further, focusing on the change of the center and the spread. We design the relative difference plot in order to support the comprehension of data changes in a linked view. We need a reference for the relative difference plot. As we have an animation brush that moves linearly in the brush space, we establish a reference brush path as a linear path between the first and the last brush in the linked space. We interpolate center positions and horizontal and vertical spread values. These values represent the reference. Now, for each brush we compute the linked data center point and spread values and depict them relative to the reference values. The brown path is created by connecting the center points of each frame of the animation, just as shown in Figure 2 (top-right, main-view). The path in black is the reference path. The relative difference plot gives clear information about how average temperature and average precipitation are related non-linearly to the analyzed country region.

4. Demonstration

The newly proposed techniques are evaluated in the context of a vehicle simulation model which is representing a four wheel drive (4WD) power transmission vehicle. The model represents the engine, a manual gear box, the central differential, the front and rear shafts, the front and rear differential, and the axles. The transmission shafts are modelled as elastic components with different stiffness and damping parameters for each shaft. The stiffness and damping of the main shafts are varied through a wider range. Additionally, the central differential split ratio, representing the distribution of torque between the front and the rear axle, is varied from 0 (rear wheel drive RWD) through 0.5 (4WD) to 1 (front wheel drive FWD). The simulation is done for a full load acceleration test, where the maximum acceleration performance is checked. Under such conditions, power transmission elements are maximally stressed. Due to the elasticity of the power transmission elements, oscillations in the power transmission can occur. If they are large and at a low frequency, discomfort is caused. The target of the analysis is to check how performance and comfort parameters are sensitive to the stiffness and damping values of the main shafts in the modelled vehicle for various torque split regimes. The variability of stiffness and damping is influenced in a narrow range by imperfections in manufacturing, assembly and material. The differences in oscillations cause comfort effects (increased amplitude and frequency in vehicle acceleration), as well as performance issues marked in fuel consumption and acceleration. A data ensemble is computed, varying differential split ratio in the range 0 to 1, as well as damping and stiffness of front and rear shafts (in the range $\pm 30\%$ of the

nominal value). 2000 calculations are performed with five input variables varied with a Sobol sequence. In each case, we study fuel consumption, the maximum torque reached for specific gears, vehicle longitudinal acceleration, and maximum torques on front and rear shafts.

First we check how stiffness and damping influences the consumption and longitudinal acceleration (one of the comfort measures). The test has been split into two parts. First, a stiffness check has been done by changing the front-shaft stiffness. The results showed that there is no significant influence on longitudinal acceleration and fuel consumption. Mean values did not change much. By changing the stiffness of the rear shaft, it is found that a lower stiffness reduces longitudinal acceleration (more comfort). In both cases, spread in consumption is large for varying stiffness. So the consumption could be sensitive to variations in stiffness due to the manufacturing process, but only within a narrow range of less than 0.1% of the absolute value. The percentile grid proved very useful to accurately move the brush across the input data space. At each step the brush accommodates 10% of the observed data points for the stiffness of the rear shaft as shown in Figure 3 (A-brushed view), as this is the expected maximum variation due to errors in manufacturing and material. The calculated performance parameters are investigated concerning mean value and distribution. The target is to find a brush view position with the smallest distribution range. This brush position specifies a nominal damping and stiffness that will in case of a manufacturing/material error cause the smallest effect on performance/comfort. The goal looks like precisely defined, and we could calculate results automatically, but actually it is not that simple to automate this before establishing the analysis steps. IVA supported with the new extensions for linking and brushing is a great help in finding and defining the relevant analysis steps, as one of the domain experts stated. We do this by moving the brush across the entire rear-shaft stiffness-range in ten steps (with the snap-to-grid option enabled), and reading the spread value from the statistics table. In this way it was easy to move the brush forth and back, knowing that at each position change, the brush will select the next 10% of the data. We also used the “select and highlight” option in the brush trace, after the trace was created. This was done to easily select the point of interest, for example, the point with the lowest value for cumulated consumption, see Figure 3 (A-brush trace in the linked view). The relevant components are cross-referenced, including the brush and the brushed points in the brushed view, which is updated according to the selected position in the trace. The cross-hair was useful as a qualitative indicator for spread change but we need a quantitative value to confirm the visual insight, especially if the cross-hair changes its size only slightly.

Next, we check the dependency of the maximal torques in different gears for various stiffness parameters. We use a 10×10 percentile grid which provides a visual assistance for brushing the lowest and highest 10% of stiffnesses values for

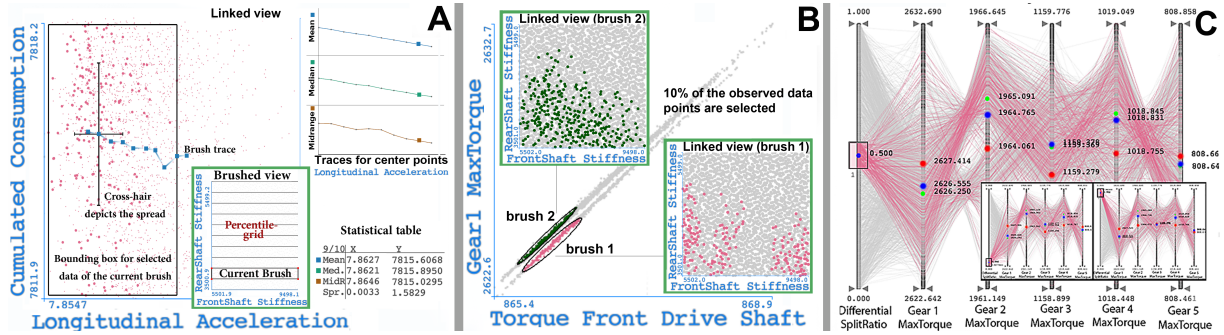


Figure 3: *A-Brushed view:* The view shows a scatterplot with a percentile 1x10 grid, and the current brush position. *A-Linked view:* The spread for y-dimension is 1.5829 as shown in the statistical table, and at the same time this is the smallest distribution range. *B:* Using the Mahalanobis brush to highlight structured data – note how this was not possible with a conventional, screen space techniques. *C:* Three different views of the same parallel coordinates plot, each showing the 10% percentile brush placed at different positions of the differential split ratio.

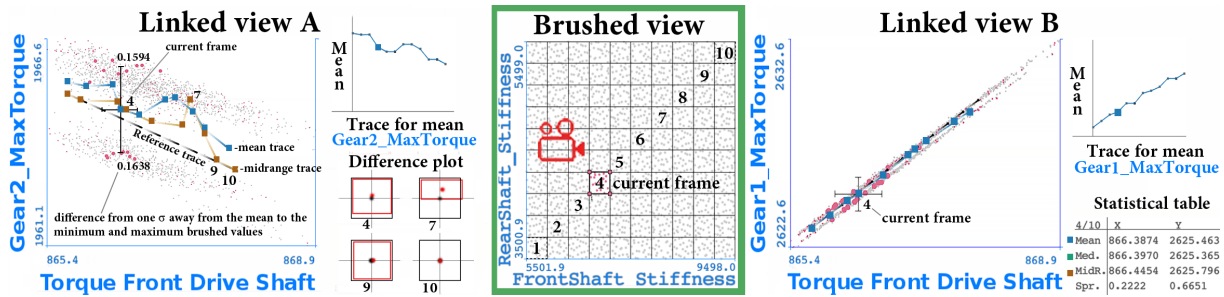


Figure 4: A screenshot from an analysis of the main shafts, torques and stiffness parameters.

the front and rear shaft. At these positions respectively we set the start and the end frame for the animation. We use eight inter-frames, as this constrains the brush to move uniformly across the data space. In this case we prefer to use animation instead of moving the brush with the mouse. Albeit the snap-to-grid option works great, moving the brush always diagonally while concentrating on several linked views is time consuming. While the animation is played in the brushed view, the linked scatterplots show changes for the first three gears. Due to space constraints we show only two linked views (Figure 4). Contrary to our expectation, the maximal values of torque rise with higher stiffness only for the first gear. The distribution of results shows that for example maximum torque of the 2nd gear can fluctuate significantly with changes in stiffness. However we see from the absolute numerical values that the fluctuation is in a range of less than one percent. This makes the selected stiffness range “robust” concerning on manufacturing imperfections.

In the next example, we use the Mahalanobis brush. We move the brush along the two visibly separated paths (Figure 3 (B)), which are parallel to each other and include always 10% of all data items in the brush. Such an exploration would be very complicated to do using conventional brush-

ing only. In our case it was a success right in the first attempt. Figure 3 (B) shows the results. It is very interesting to see, as shown in Figure 3 (B-linked view for the brush 1), how the linked parameter space splits. This happens once we return towards the lower values for both the torque front and torque rear drive-shaft parameters. But, this happens only for the upper path in (Figure 3 (B)), the one with constantly slightly higher front drive shaft values.

The last check is performed for maximum torques in different gears, for different driving regimes: FWD, RWD or 4WD. Parallel coordinates are used to show six data dimensions at once, and a statistics for the center points are enabled in the view (Figure 3 C). The first axis shows differential split ratio, and gears are mapped to the successive axes. We use the 10% percentile brush for selecting the differential split ratio at three different positions. The analysis shows that the maximum torques in gear two and four have relatively higher mean values than the torques in the other gears for three transmission cases. This is a hard to find design phenomenon which is determined by coupling an engine with its power characteristics and used gearbox.

Constraint brushing is an invaluable feature in a team-

work, if team members work on same types of datasets and need to (re)build an analysis step by step. As with this feature brushing becomes accurately defined, it is easy to step back in the analysis and try a different path, preserving all what has been done up to that point. Constraint brushing makes analysis steps recordable and easy to communicate. Our linking&brushing extensions proved to be useful for data analysis in the presented case. Linking quantitative and statistical parameters extended the boundaries of what can be recognized from raw data. One request that followed from this case study is to depict also the 'Spread' as a graph, next to Mean, Median, Midrange, in the statistical overview.

5. Discussion

In this paper we propose the use of *constrained brushing*, as an addition to the traditional (unconstrained) brushing, supporting the reproducibility of the analysis results. Specifically, our aim was to simplify the way how the user can repeatedly select the same data subset of interest, without the need to record an entire workflow. We show how to control the brushing interaction by introducing the concept of a structured brushing space, based on anchoring the brush, the extent of the brush, and the movement of the brush. The user can decide how to combine these constraints, for example, she can snap to a conventional grid for moving the brush and use a percentile brush for the brush extent. Although we exemplify the newly proposed techniques for scatterplot and parallel coordinates only, it is straightforward to extend them also to other views with quantitative axes.

Constrained brush movements provide benefits when doing a rank-based analysis, since at each step we can better control and interpret the brush. Extensions like the percentile grid and percentile brushes are powerful options for doing rank-based analysis. The results can be reproduced later very easily, for example, based on a textual description of the brushed data. The analyst can benefit from the structured visualization space even if constrained brushing is not enabled. An example is to depict the grid which assists to navigate the brush over the presented data. The constrained brushing can help the user to stay in the 'flow of analysis', while also providing quantitative precision. The user can quantitatively interpret the brush while moving it along the constrained direction. Indirect manipulation, e.g., through off-screen widgets, such as sliders, can compromise the user's focus on direct interaction to a certain degree. An example of indirect manipulation would be the Mahalanobis brush. The user sets with a slider the percentage of the points that should be selected by the brush. The brush adapts its size and shape automatically depending on the underlying data distribution. An alternative option could be to use a clustering algorithm to automatically calculate a meaningful percentage for the size of the Mahalanobis brush.

Grids proved to be very useful for structuring the brushing space. We provide some meaningful default values for the

grid size, e.g., we divide the data space into four quartiles, but we also allow the user to specify non-uniform grids. We also consider possibilities to use automatic methods for exploring the data space and divide the grid according to the data distribution. For now the user can manually set the grid, e.g., task driven, either rank-based or value-based.

Summarized statistics shown in linked views, in a table, or as an overview, present a natural way for adding quantitative information about the brushed data to other dimensions. Those can be added also for views which do not have quantitative axes. However quantitative extensions to show descriptive statistics for categorical data are not covered in our current work.

Traces from brushing can be used for analyzing data at different pathline positions. To follow the principles of IVA this should be interactive and cross-linked with other views. If a point on the trace is selected, the brush in the brushed view should also be updated. This way the user can go back to some point of the analysis and maybe explore it towards some other direction. Obviously, the extensions that we present in this paper are only a first step and we expect substantial future research towards quantitative and reproducible visual analytics.

6. Conclusions and Future Work

In this paper, we address two important limitations in current visual analytics, namely the lack of reproducibility and quantitative results. We present extensions to the well-established concept of linking&brushing including constrained brushing, animated brushing and percentile brushing. They can improve the reproducibility of visual analytics and provide the user with quantitative results. We discuss a possible structuring of the brushing space that is oriented towards an improved reproducibility of interactive brushing. The Mahalanobis brush takes the local data distribution into account and selects a predefined number of points. This brush is especially useful in areas with an elongated data distribution. Compared to the circular percentile brush and the standard rectangular brush it does not select outliers from the underlying data distribution. An advantage of integrating descriptive statistics is that it helps in creating a better mental image of changes in the linked views while the brush is moving. Animation is an example of how to structure the brushing space, such that in the brush view the selections remain simple and easy, while the user is free to concentrate on the interpretation of the linked view(s). As an addition to the animation, the relative difference plot adds to the comprehension of data changes in the linked view(s). In general, and in order to conquer important new application fields, we conclude that there is a need for visual analytics to (also) provide reproducible and quantitative results.

References

- [BC87] BECKER R. A., CLEVELAND W. S.: Brushing scatterplots. *Technometrics* 29, 2 (May 1987), 127–142. 2
- [BPF14] BACH B., PIETRIGA E., FEKETE J.-D.: Graphdiaries: Animated transitions and temporal navigation for dynamic networks. *Visualization and Computer Graphics, IEEE Transactions on* 20, 5 (May 2014), 740–754. 2
- [Che03] CHEN H.: Compound brushing [dynamic data visualization]. In *Information Visualization, 2003. INFOVIS 2003. IEEE Symposium on* (Oct 2003), pp. 181–188. 2, 3
- [CM88] CLEVELAND W. C., MCGILL M. E.: *Dynamic Graphics for Statistics*, 1st ed. CRC Press, Inc., Boca Raton, FL, USA, 1988. 2
- [DGH03] DOLEISCH H., GASSER M., HAUSER H.: Interactive feature specification for focus+context visualization of complex simulation data. In *Proc. of the 5th Joint IEEE TCVG - EUROGRAPHICS Symposium on Visualization (VisSym 2003)* (2003), pp. 239–248. 2
- [FWR00] FUA Y.-H., WARD M. O., RUNDENSTEINER E. A.: Structure-based brushes: A mechanism for navigating hierarchically organized data and information spaces. *IEEE Transactions on Visualization and Computer Graphics* 6, 2 (Apr. 2000), 150–159. 2
- [GWL*10] GOTZ D., WHEN Z., LU J., KISSA P., CAO N., QIAN W. H., LIU S. X., ZHOU M. X.: Harvest: An intelligent visual analytic tool for the masses. In *Proceedings of the First International Workshop on Intelligent Visual Interfaces for Text Analysis* (New York, NY, USA, 2010), IVITA '10, ACM, pp. 1–4. 2
- [Hau06] HAUSER H.: *Generalizing Focus+Context Visualization, in Scientific Visualization: The Visual Extraction of Knowledge from Data*. Springer, 2006, ch. Generalizing Focus+Context Visualization, pp. 305–327. 1
- [HBC*91] HASLETT J., BRADLEY R., CRAIG P., UNWIN A., WILLS G.: Dynamic graphics for exploring spatial data with application to locating global and local anomalies. *The American Statistician* 45, 3 (1991), 234–242. 2, 5
- [HR07] HEER J., ROBERTSON G.: Animated transitions in statistical data graphics. *IEEE Transactions on Visualization and Computer Graphics* 13, 6 (Nov. 2007), 1240–1247. 2
- [KBHP14] KANDOGAN E., BALAKRISHNAN A., HABER E., PIERCE J.: From data to insight: Work practices of analysts in the enterprise. *Computer Graphics and Applications, IEEE* 34, 5 (2014). 1
- [KFH10] KEHRER J., FILZMOSER P., HAUSER H.: Brushing moments in interactive visual analysis. In *Proceedings of the 12th Eurographics / IEEE - VGTC Conference on Visualization* (Aire-la-Ville, Switzerland, Switzerland, 2010), EuroVis'10, Eurographics Association, pp. 813–822. 2, 3
- [KH13] KEHRER J., HAUSER H.: Visualization and visual analysis of multifaceted scientific data: A survey. *Visualization and Computer Graphics, IEEE Transactions on* 19, 3 (March 2013), 495–513. doi:10.1109/TVCG.2012.110. 2
- [Mah36] MAHALANOBIS P. C.: On the generalized distance in statistics. *Proceedings of the National Institute of Sciences (Calcutta)* 2 (1936), 49–55. 4
- [MW95] MARTIN A., WARD M.: High dimensional brushing for interactive exploration of multivariate data. In *Visualization, 1995. Visualization '95. Proceedings., IEEE Conference on* (1995), pp. 271–. 2
- [NOA14] NOAA: National Climatic Data Center, 2014. 2
- [RFF*08] ROBERTSON G., FERNANDEZ R., FISHER D., LEE B., STASKO J.: Effectiveness of animation in trend visualization. *Visualization and Computer Graphics, IEEE Transactions on* 14, 6 (Nov 2008), 1325–1332. 2
- [ROC97] RENSINK R. A., O'REGAN J. K., CLARK J. J.: To see or not to see: The need for attention to perceive changes in scenes. *Psychological Science* 8, 5 (1997), 368–373. 2
- [SFSA10] SILVA C., FREIRE J., SANTOS E., ANDERSON E.: Provenance-enabled data exploration and visualization with vis-trails. In *Graphics, Patterns and Images Tutorials (SIBGRAPI-T), 2010 23rd SIBGRAPI Conference on* (Aug 2010), pp. 1–9. 2, 3
- [Shn92] SHNEIDERMAN B.: Tree visualization with tree-maps: 2-d space-filling approach. *ACM Trans. Graph.* 11, 1 (Jan. 1992), 92–99. URL: <http://doi.acm.org/10.1145/102377.115768>, doi:10.1145/102377.115768. 2
- [Spe87] SPEARMAN C.: The proof and measurement of association between two things. By C. Spearman, 1904. *The American journal of psychology* 100, 3-4 (1987), 441–471. URL: <http://view.ncbi.nlm.nih.gov/pubmed/3322052>. 3
- [TMB02] TVERSKY B., MORRISON J. B., BETRANCOURT M.: Animation: Can it facilitate? *Int. J. Hum.-Comput. Stud.* 57, 4 (Oct. 2002), 247–262. 6
- [War94] WARD M. O.: Xmdvtool: integrating multiple methods for visualizing multivariate data. In *VIS '94: Proceedings of the conference on Visualization '94* (Los Alamitos, CA, USA, 1994), IEEE Computer Society Press, pp. 326–333. 2
- [WH14] WEBER G. H., HAUSER H.: Interactive visual exploration and analysis. In *Scientific Visualization: Uncertainty, Multifield, Bio-Medical and Scalable Visualization*, Hansen C. D., Chen M., Johnson C. R., Kaufman A. E., Hagen H., (Eds.), Mathematics and Visualization. Springer-Verlag, 2014, pp. 161–174. LBNL-6655E. 2
- [YXRW07] YANG D., XIE Z., RUNDENSTEINER E. A., WARD M. O.: Managing discoveries in the visual analytics process. *SIGKDD Explor. Newsl.* 9, 2 (Dec. 2007), 22–29. 2