

Visual Analysis of Dimensionality Reduction Quality for Parameterized Projections

Rafael Messias Martins^{a,b}, Danilo Coimbra^{a,b}, Rosane Minghim^a, A. C. Telea^{b,c}

^a*ICMC, University of São Paulo, São Carlos, Brazil, 13566-590*

^b*Institute Johann Bernoulli, University of Groningen, The Netherlands*

^c*University of Medicine and Pharmacy ‘C. Davila’, Bucharest, Romania*

Abstract

In recent years, many dimensionality reduction (DR) algorithms have been proposed for visual analysis of multidimensional data. Given a set of n -dimensional observations, such algorithms create a 2D or 3D projection thereof that preserves relative distances or neighborhoods. The quality of resulting projections is strongly influenced by many choices, such as the DR techniques used and their various parameter settings. Users find it challenging to judge the effectiveness of a projection in maintaining features from the original space and to understand the effect of parameter settings on these results, as well as performing related tasks such as comparing two projections. We present a set of interactive visualizations that aim to help users with these tasks by revealing the quality of a projection and thus allowing inspection of parameter choices for DR algorithms, by observing the effects of these choices on the resulting projection. Our visualizations target questions regarding neighborhoods, such as finding false and missing neighbors and showing how such projection errors depend on algorithm or parameter choices. By using several space-filling techniques, our visualizations scale to large datasets. We apply our visualizations on several recent DR techniques and high-dimensional datasets, showing how they easily offer local detail on point and group neighborhood preservation while relieving users from having to understand technical details of projections.

Keywords: Visual Analytics, Dimensionality Reduction, Parameterization, Projection Errors, Image-based, Large Data

1. Introduction

Dimensionality reduction (DR) techniques are an increasingly popular and pervasive part of visual analytics solutions. Their key value is the ability to transform, or project, high-dimensional datasets into low-dimensional datasets which keep the underlying structure of the data similar. The results can be visualized by scatterplots [1], treemaps, timelines, and parallel coordinates [2]. DR methods have been used for the visual analysis of text documents [3, 4, 5], multimedia [6], text mining [7, 8], vector fields [9], and biomedical data [10, 11].

Although DR techniques have become increasingly more robust and computationally scalable, several major usability challenges still exist. One such challenge involves the quality analysis of DR algorithms. Currently, tens of DR algorithms exist, each with several parameters, whose values strongly influence the projection result. Changes of a single parameter can produce different projections, casting doubt on the correctness or meaning of the resulting projection. However, such parameters are typically quite technical and non-intuitive for the average end-user. Our question is, thus: How to provide insight into the *quality* of DR algorithms, and how to explore their parameter settings, so that users understand how these settings affect the shape, structure, and quality of the resulting projections? In this paper, we present a set of visualization techniques that help users with exploring the link between DR algorithm parameter settings and the quality of the resulting projections. Our visualizations target the following questions:

- 28 • How is the projection error spread over the 2D space?
- 29 • How to find points which are close in 2D but far in n D?
- 30 • How to find points which are close in n D but far in 2D?
- 31 • How do DR algorithm choice and parameter settings affect
32 the above quality aspects?

For this, we propose several space-filling techniques that visually scale to large datasets, offer a multiscale (or level-of-detail) view on the projection behavior, and do not require users to understand the internal formulation of DR algorithm. We illustrate our visualizations by exploring the parameters of five state-of-the-art DR techniques for several real-world datasets.

This paper is structured as follows. Section 2 presents related work on DR algorithm quality analysis. Section 3 presents our analysis goals. Section 4 describes our proposed visualizations. Section 5 uses these methods to explore the quality, as function of DR method parameters, of several DR techniques. Section 6 discusses our results. Section 7 concludes the paper.

2. Related Work

2.1. Dimensionality reduction

For a dataset $D^n = \{\mathbf{p}_i \in \mathbb{R}^n\}_{1 \leq i \leq N}$ of N n -dimensional points, dimensionality reduction (DR) can be seen as a function

$$f : \mathbb{R}^n \times P \rightarrow \mathbb{R}^m \quad (1)$$

which maps each point $\mathbf{p}_i \in D^n$ to a point $\mathbf{q}_i \in D^m$. Here, n is typically large (tens up to thousands of dimensions), and m is typically 2 or 3. P denotes the *parameter space* of f , *i.e.* the various settings that control the projection algorithm, including the algorithm type itself. f is designed to keep the so-called *structure* of the data as similar as possible in \mathbb{R}^n and \mathbb{R}^m . One way for this is to let f minimize the normalized stress function

$$\sigma = \frac{\sum_{1 \leq i \leq N, 1 \leq j \leq N} (d^n(\mathbf{p}_i, \mathbf{p}_j) - d^m(\mathbf{q}_i, \mathbf{q}_j))^2}{\sum_{1 \leq i \leq N, 1 \leq j \leq N} (d^n(\mathbf{p}_i, \mathbf{p}_j))^2} \quad (2)$$

where $d^n : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}^+$ and $d^m : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}^+$ are distance metrics for D^n and D^m respectively. Other ways to compute f are to optimize for having the k -nearest neighbors for a point $\mathbf{q}_i \in D^m$ be the same as the k -nearest neighbors of $\mathbf{p}_i \in D^n$.

Many DR methods are special cases of a wider class of techniques called Multidimensional Scaling (MDS). MDS methods compute f using only pairwise point distances. This avoids having to access the full nD coordinate data. However, computing distances creates additional costs ($O(N^2)$ for N points). The PLMP algorithm avoids this by using distances only for a small set of representative points and using nD coordinates for the other points [10].

DR methods can be classified by the techniques used to compute f [10]. *Spectral decomposition* techniques project points along the largest-eigenvalue eigenvectors of the pointwise distance matrix [12]. LLE [13] and ISOMAP [14, 15] use efficient numerical methods tailored to solve sparse eigenproblems. Landmarks MDS [16] and Pivot MDS [17] book further speed-ups by using classical MDS on a subset of representative points and projecting remaining points by local interpolation. Fastmap achieves linear complexity in the input point count but has a worse stress minimization [6].

Nonlinear optimization methods iteratively search the parameter space P to minimize the stress σ [18, 19]. Besides naive gradient descent, multigrid numerical solvers can be used to speed searching [20]. Pekalska *et al.* propose a speed-up that projects a representative subset (by gradient descent) and fits remaining points by local interpolation [21]. Force-based methods are a special class of nonlinear optimization with many uses in graph drawing [22]. Chalmers speeds this up by using the representative subset idea outlined earlier [23]. Further speed-ups are achieved by multilevel solvers and GPU techniques [24, 25], and by recursively selecting representatives via a multilevel approach [26]. Tejada *et al.* use a heuristic to embed instances by force-based relaxation [27]. LSP positions the representative subset by a force-based scheme and fits the remaining points by Laplacian smoothing [4]. LAMP also uses a representative subset to locally construct affine projections, and allows users to interactively place these points to optimize the overall projection layout [3]. More details on LSP, LAMP, and ISOMAP are given further in Section 5.2.

2.2. Visualizing projection quality

Although projection quality is acknowledged as important, most DR literature considers mainly aggregated quality metrics

such as the stress function (Eqn. 2), correlation [28], neighborhood preservation average plots [4], and distance scatterplots [3], which are distance and neighborhood based metrics, or cluster segregation metrics [29]. 2D scatterplots can show the correlation of D^n with D^m [3]. Such metrics capture the *overall* quality of a projection, but do not help finding local quality variations. In other words, they do not show projection problems for *any* point i vs *all* points $j \neq i$ in the input dataset.

Local metrics can be used to highlight where (in a projection) errors happen. Shreck *et al.* compute, for each $\mathbf{p} \in D^n$, the projection precision score (*pps*) defined as the normalized distance between the two k -dimensional vectors having as components the Euclidean distances between \mathbf{p} and its k nearest neighbors in D^n , respectively D^2 [30]. Visualizing *pps* as a color map shows areas where neighborhoods are not preserved. However, a neighborhood cannot be preserved for two distinct reasons: true neighbors (in D^n) are missing (in D^2), or neighbors (in D^2) are actually false neighbors (in D^n). The *pps* metric does not differentiate between such situations, and can also be sensitive to permutations of points that do not change distances.

Recognizing that DR methods can create distance approximation errors, Van der Maaten *et al.* extend the t-SNE technique [31] to output a set $\{M_i\}$ of 2D projections rather than a single one [32]. All points appear in all projections M_i , with potentially different weights and at different locations. This allows better modeling non-metric similarities. Yet, correlating points over the several M_i is done manually by the user, and can be challenging for large datasets and many projections M_i .

Several quality metrics for continuous DR techniques are proposed by Aupetit [33]. Point-based stretching and compression metrics measure, for each $\mathbf{p}_i \in D^n$, the aggregated increase, respectively decrease, of the distances of its projection $\mathbf{q}_i \in D^2$ to all other projections $\mathbf{q}_{j \neq i}$ vs the distances of \mathbf{p}_i to all other points $\mathbf{p}_{j \neq i}$. Segment stretching and compression measures the variation of distances of close point pairs (i, j) between \mathbb{R}^n and \mathbb{R}^2 . For a selected \mathbf{p}_i , the proximity metric maps distances in \mathbb{R}^n from \mathbf{p}_i to all other points $\mathbf{p}_{j \neq i}$ to the corresponding points $\mathbf{q}_i \in \mathbb{R}^2$ and thereby helps understanding how (and where) the projection may have distorted the structure of the data. These metrics are visualized with piecewise-constant interpolation of the point, respectively segment, data using Voronoi diagrams. Our proposed techniques in Secs. 4.2, 4.3, and 4.4 adapt and extend these visualizations in several directions.

Still using colored Voronoi cells, Lespinats and Aupetit show, at the same time, point stretching and compression by using a 2D color map [34]. The proposed color map encodes stretching as green, compression as purple, low-error points as white, and points with high stretching and compression as black, respectively. While this color map can show local error types (or the absence thereof), it cannot explicitly show the point-pairs which cause stretching and compression. Besides, as the authors also note, Voronoi cells can lead to visualization bias due to the cells' sizes and shapes being heavily dependent on the D^2 point density, and the fact that cells cover the entire \mathbb{R}^2 space, even in areas where no projected points exist.

To assist the task of navigating projections while also considering distortions, Heulot *et al.* present an interactive semantic

lens that filters points projected too closely to a user-selected focus point in \mathbb{R}^2 [35]. Such points, also called false neighbors, are pushed towards the lens border, so they do not attract the user's attention. Separately, points are colored by the distance in D^n to the focus point, to help users navigate to the so-called missing neighbors of the focus point. Instead of Voronoi cells of [33, 34], points are colored using Shepard interpolation, which yields a smoother, and arguably less distracting, image. However, in contrast to [33, 34], this method can only show errors related to a selected focus point.

3. Analysis goals

A projection f should preserve the *structure* of the original space \mathbf{R}^n . This implies a mix of distance and neighborhood preservations at various scales and happens at different rates for different datasets, projection algorithms and parameter values. For users, the projection's *precision* [30] is not clear unless they can interpret projected neighborhoods adequately [33]. Thus, given any DR algorithm (Eqn. 1), we aim to show how neighborhood preservation is affected by choices of parameter values in P , highlighting aspects that can adversely affect the interpretation of the projected point set in D^m . To simplify the discourse, we next consider $m = 2$, and that projections are drawn as scatterplots (the most common option for DR visualization). We identify the following aspects of interest:

A. False neighbors: Take a point $\mathbf{p}_i \in D^n$ and its 2D projection $\mathbf{q}_i = f(\mathbf{p}_i)$. A necessary condition for neighborhood preservation is that *all* points \mathbf{q}_j which are close to \mathbf{q}_i (in 2D) should be projections of points \mathbf{p}_j which are close to \mathbf{p}_i (in D^n). If not, *i.e.* we have a \mathbf{q}_j close to \mathbf{q}_i for which \mathbf{p}_j is not close to \mathbf{p}_i , the user wrongly infers from the projection that \mathbf{p}_j is close to \mathbf{p}_i . We call such a point j a *false neighbor* of i .

B. Missing neighbors: The second necessary condition for neighborhood preservation is that *all* \mathbf{p}_j which are close to \mathbf{p}_i (in D^n) project to points \mathbf{q}_j which are close to \mathbf{q}_i (in 2D). If not, *i.e.* we have a \mathbf{p}_j close to \mathbf{p}_i for which \mathbf{q}_j is not close to \mathbf{q}_i , the user will underestimate the set of points similar to point i . We call such a point j a *missing neighbor* of i .

C. Groups: A main goal of DR is to help users find groups of similar points, *e.g.* topics in a document set [3, 4] or classes of images in a database [6]. False and missing point neighbors generalize, for groups, to *false members* and *missing members* respectively. Given a group Γ of closely projected points, we aim to find if *all* points in Γ truly belong there (no false members), and if *all* points that belong to the topic described by Γ do indeed project in Γ (no missing members).

D. Detail: Aggregated local metrics such as [30, 33, 34, 35] can show, up to various extents, where missing or false neighbors occur. However, they do not directly show which are all such neighbors, for each projected point. Also, they do not explicitly address locating false and missing group members. We aim to provide interactive visual mechanisms to support these tasks on several levels of detail.

4. Visualization methods

We next propose several visualization methods to address the analysis goals outlined in Sec. 3. As a running example, we use LAMP as projection method, with the default parameter settings given in [3], and as input the well-known 19-dimensional Segmentation dataset with 2300 points from [36, 3, 37, 10]. Herein, each point describes a randomly drawn 3x3 pixel-block from a set of 7 manually segmented outdoor images, by means of 19 statistical image attributes, such as color mean, standard deviation, and horizontal and vertical contrast.

4.1. Preliminaries

To quantify the neighborhood preservation issues in Sec. 3, we first define the projection error of point i vs a point $j \neq i$ as

$$e_{ij} = \frac{d^m(\mathbf{q}_i, \mathbf{q}_j)}{\max_{i,j} d^m(\mathbf{q}_i, \mathbf{q}_j)} - \frac{d^n(\mathbf{p}_i, \mathbf{p}_j)}{\max_{i,j} d^n(\mathbf{p}_i, \mathbf{p}_j)}. \quad (3)$$

We see that $e_{ij} \in [-1, 1]$. Negative errors indicate points whose projections are too close (thus, false neighbors). Positive errors indicate points whose projections are too far apart (thus, missing neighbors). Zero values indicate 'good' projections, which approximate optimally the distances in D^n .

4.2. Aggregated error view

We first provide an overview of how the projection error spreads over an entire dataset, by computing for each point i the aggregate error

$$e_i^{agg} = \sum_{j \neq i} |e_{ij}|. \quad (4)$$

The value of e_i^{agg} gives the projection error of point i with respect to *all* other points. Low values of e_i^{agg} show points whose projections can be reliably compared with most other projections in terms of assessing similarity. These are good candidates for representatives in multilevel projection methods [6, 21, 23, 4]. Large values of e_i^{agg} show points which are badly placed with respect to most other points. These are good candidates for manual projection optimization [38, 37].

Fig. 1 (a) shows e_i^{agg} by color mapping its value on the 2D projected points, using a blue-yellow-red diverging colormap [39]. Brushing and zooming this image allows inspecting e_i^{agg} for individual points. However, given our goal of providing an overview first, we are actually not interested in *all* individual e_i^{agg} values, but rather to (a) find compact areas in the projection having similar e_i^{agg} values, (b) find outlier e_i^{agg} values in these areas (if any), and (c) see how e_i^{agg} globally varies across the projection. For this, we propose an image-based, space-filling visualization, as follows. Denote by $DT(\mathbf{x} \in \mathbb{R}^2) = \min_{\mathbf{q} \in D^m} \|\mathbf{q} - \mathbf{x}\|$ the so-called distance transform of the 2D point cloud D^m delivering, for any screen pixel \mathbf{x} , its distance to the closest point in D^m . We then compute e_i^{agg} at every screen pixel \mathbf{x} as

$$e_i^{agg}(\mathbf{x}) = \frac{\sum_{\mathbf{q} \in N_\epsilon(\mathbf{x})} \exp\left(-\frac{\|\mathbf{x} - \mathbf{q}\|^2}{\epsilon^2}\right) e_i^{agg}}{\sum_{\mathbf{q} \in N_\epsilon(\mathbf{x})} \exp\left(-\frac{\|\mathbf{x} - \mathbf{q}\|^2}{\epsilon^2}\right)} \quad (5)$$

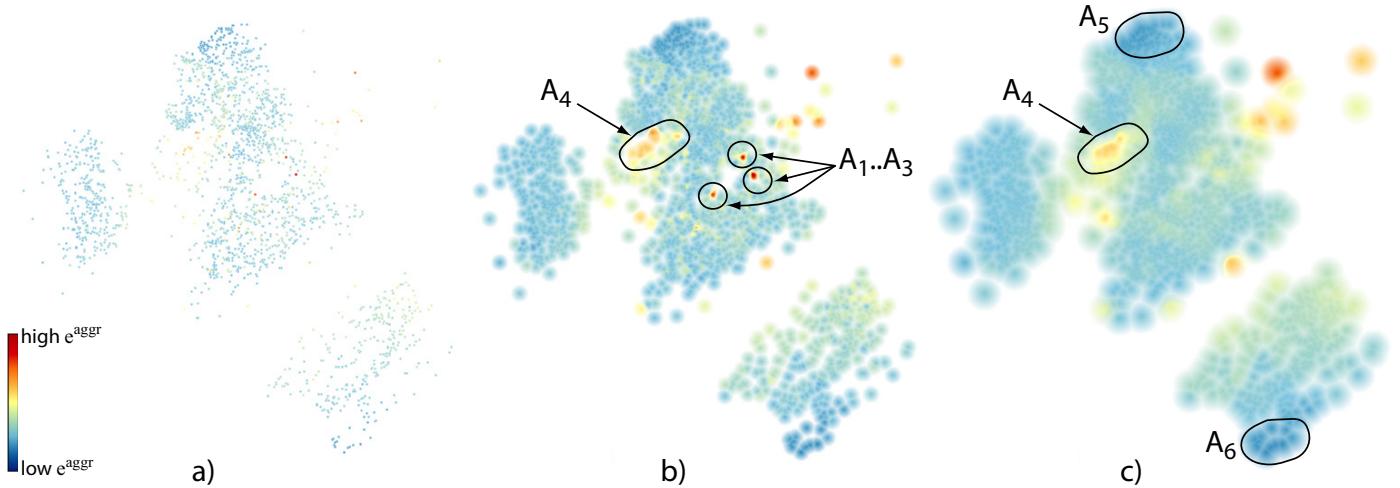


Figure 1: Aggregate error view, several levels of detail: (a) $\alpha = 1, \beta = 1$. (b) $\alpha = 5, \beta = 5$. (c) $\alpha = 20, \beta = 20$ pixels (see Sec. 4.2).

with

$$\epsilon = DT(\mathbf{x}) + \alpha. \quad (6)$$

Here, $N_\epsilon(\mathbf{x})$ contains all projections in D^m located within a radius ϵ from \mathbf{x} . We next draw $e^{agg}(\mathbf{x})$ as a RGBA texture, where the color components encode $e^{agg}(\mathbf{x})$ mapped via a suitable color map, and the transparency A is set to

$$A^{agg}(\mathbf{x}) = \begin{cases} 1 - \frac{DT(\mathbf{x})}{\alpha}, & \text{if } DT(\mathbf{x}) < \beta \\ 0, & \text{otherwise} \end{cases} \quad (7)$$

For $\alpha = 1, \beta = 1$, we obtain the classical colored scatterplot (Fig. 1 (a)). For $\alpha = 1, \beta > 1$, the space between projections is filled, up to a distance β , by the e^{agg} value of the closest data point. For $\alpha = 1, \beta = \infty$, we obtain a Voronoi diagram of the projections with cells colored by their e^{agg} values. This does not change the e^{agg} data values, but just displays them on larger spatial extents than individual pixels, making them easier to see. This creates visualizations identical to those obtained by drawing scatterplots with point radii equal to β , without having the issues created by overlapping points. For $\alpha > 1, \beta > 1$, the result is similar to Shepard interpolation where the kernel size ϵ is given by the *local* point density. The parameter $\alpha \geq 0$ controls the *global* level-of-detail at which we visualize e^{agg} : Small values show more detail in dense point zones, but also emphasize small-scale signal variations which are less interesting. Larger α values create a smoother signal where coarse-scale error patterns are more easily visible.

Figs. 1 (b,c) show the aggregate error for the Segmentation dataset for various values of the parameters α and β . Here, $e_{ij} \in [-0.67, 0.35]$. The error range already tells that we have poorly projected points, but does not tell where these are. In Fig. 1 (b), with low values for both α and β , we see that e^{agg} is relatively smoothly distributed over the entire projection. However, we see three small red spots $A_1..A_3$. These are high-error outlier areas, which indicate points that are badly placed with respect to *most* other points. We also see a relatively high error area A_4 of larger spatial extent. Increasing both α and β produces a simplified visualization (Fig. 1 (c)). Larger β values fill

in the gaps between points. Larger α values eliminate outlier regions whose spatial extent is smaller than α , such as the three small outlier areas $A_1..A_3$, but A_4 remains visible, since it is larger than α . We now also notice, better than in Fig. 1 (b), that the bottom and top areas (A_5, A_6) in the projection have dark blue values, with a significantly lower error than the rest of the projection.

Our image-based results are slightly reminiscent of the dense *pps* maps of Schreck *et al.* [30] (see Sec. 2.2). Differences exist, however. First, our e_i^{agg} is a *global* metric, that tells how point i is placed with respect to all other points, whereas the *pps* metric characterizes *local* neighborhoods. Interpolation-wise, our technique (used with $\alpha = 1, \beta = \infty$) delivers the same Voronoi diagram as Schreck *et al.*, which is also identical to the space partitioning of the point-based Voronoi diagrams in [33, 34]. The data being mapped is, however, different: Our e^{agg} shows the sum of distance compression and stretching, whereas [33, 34] treat these two quantities separately. In the next sections, we show how we split our aggregated insight into separate insights. Further on, both Schreck *et al.* and our method use smoothing to remove small-scale noise from such maps. However, whereas Schreck *et al.* uses a constant-radius smoothing kernel, which blurs the image equally strong everywhere, we use, as explained, a variable-radius kernel controlled by local density, which preserves better detail in non-uniform point clouds.

4.3. False neighbors view

However useful to assess the error distribution and find badly vs well-projected point groups, the aggregate error view does not tell us if the error is due to false neighbors, missing neighbors, or both. Let us first consider the false neighbors (case A, Sec. 3). To visualize these, we create a Delaunay triangulation of the projected point cloud that gives us the closest neighbors of each projected point in all directions, *i.e.*, the most important false-neighbor candidates for that point. To each edge E_k , $1 \leq k \leq 3$ of each triangle T of this triangulation, with vertices being the points \mathbf{q}_i and \mathbf{q}_j of D^m , we assign a weight

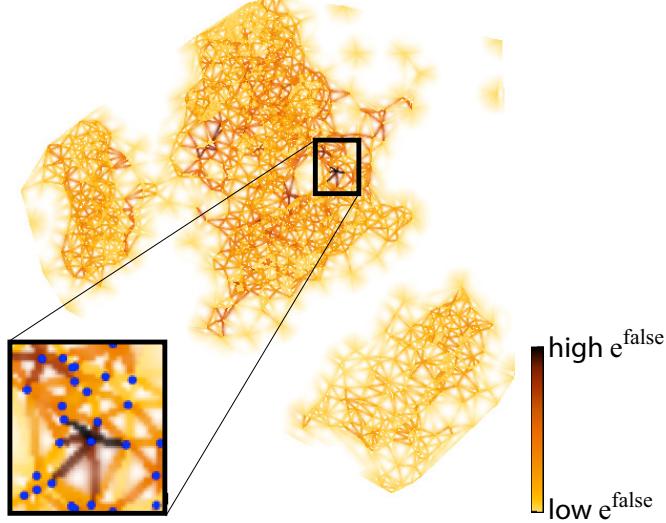


Figure 2: False neighbors view (see Sec. 4.3).

$e_k^{false} = |\min(e_{ij}, 0)|$, i.e., consider only errors created by false neighbors. Next, we interpolate e^{false} over all pixels \mathbf{x} of T by using

$$e^{false}(\mathbf{x}) = \frac{\sum_{1 \leq k \leq 3} \frac{1}{d(\mathbf{x}, E_k)} \|E_k\| e_k^{false}}{\sum_{1 \leq k \leq 3} \frac{1}{d(\mathbf{x}, E_k)} \|E_k\|} \quad (8)$$

where $d(\mathbf{x}, E)$ is the distance from \mathbf{x} to the edge E and $\|E\|$ is the length of the edge. Similarly to the aggregated error, we construct and render an image-based view for e^{false} as a RGBA texture. In contrast to the aggregated error, we use here a heated body colormap [39], with light hues showing low e^{false} values and dark hues showing high e^{false} values. This attracts the attention to the latter values, while pushing the former ones into the background. The transparency A is given by

$$A^{false}(\mathbf{x}) = A^{agg}(\mathbf{x}) \left(1 - \frac{1}{2} \left(\min \left(\frac{DT_T(\mathbf{x})}{DT_C(\mathbf{x})}, 1 \right) + \max \left(1 - \frac{DT_C(\mathbf{x})}{DT_T(\mathbf{x})}, 0 \right) \right) \right) \quad (9)$$

where $DT_T(\mathbf{x}) = \min(d(\mathbf{x}, E_1), d(\mathbf{x}, E_2), d(\mathbf{x}, E_3))$ is the distance transform of T at \mathbf{x} , $DT_C(\mathbf{x})$ is the distance from \mathbf{x} to the barycenter of T , and A^{agg} is given by Eqn. 7. The same technique is used in a different context to smoothly interpolate between two 2D nested shapes [40], where we refer for further implementation details. The combined effect of Eqns. 8 and 9 is to slightly thicken, or smooth out, the rendering of the Delaunay triangulation. Note that this interpolation does not change the actual values e_k^{false} rendered on the triangulation edges. The distance-dependent transparency ensures that data is shown only close to the projection points.

Fig. 2 shows the false neighbors for the Segmentation dataset. Several things are apparent here. First, the rendering is similar to a blurred rendering of the Delaunay triangulation of the 2D projections colored by e^{false} , showing how each point relates to its immediate neighbors. Light-colored edges show true neighbors, while dark edges show false neighbors. Since

edges are individually visible, due to the transparency modulation (Eqn. 9), we can see both the true and false neighbors of a point separately. The smooth transition between opaque points (on the Delaunay edges) and fully transparent points (at the triangles' barycenters) ensures that the resulting image is continuous and easier to follow at various screen resolutions than a Delaunay triangulation rendered with pixel-thin edges, as our edges appear slightly thicker.

In Fig. 2, two error-related aspects are visible. First, we see an overall trend from light to dark colors as we go further from the projection's border towards the projection center. This confirms the known observation on DR methods that projections on the border tend to be more accurate, since there is more freedom (and space) to place these. In contrast, projections falling deep inside the resulting point cloud tend to have more false neighbors, because the DR algorithm has there less space to shift points around to accommodate all existing distance constraints. Intuitively, we can think of this phenomenon as a ‘pressure’ which builds up within the projected point set from its border inwards. We shall see more examples of this phenomenon in Sec. 5. Secondly, we see a few small-scale dark outliers. Zooming in Fig. 2, we see that these are points connected by dark edges to most of their closest neighbors in a star-like pattern. Clearly, false neighbors exist here. These can be either the star ‘center’ or the tips of its branches. However, we also see that these tips have only one dark edge. Hence, they are too closely positioned to the star center only, and not to their other neighbors. Since the tip points are all positioned well with respect to their neighbors (except the star center), and the center point is positioned too closely with respect to all its direct neighbors, we can conclude that too little space was offered in the projection to the center point, or in other words that the center point is a false neighbor of its surrounding points.

The false neighbors view is related to Aupetit’s segment compression view, where the shortening of inter-point distances due to projection is visualized [33]. The underlying metrics, i.e. our e_{ij} (Eqn. 3) and m_{ij}^{distor} ([33], Sec. 3.2) are similar, up to

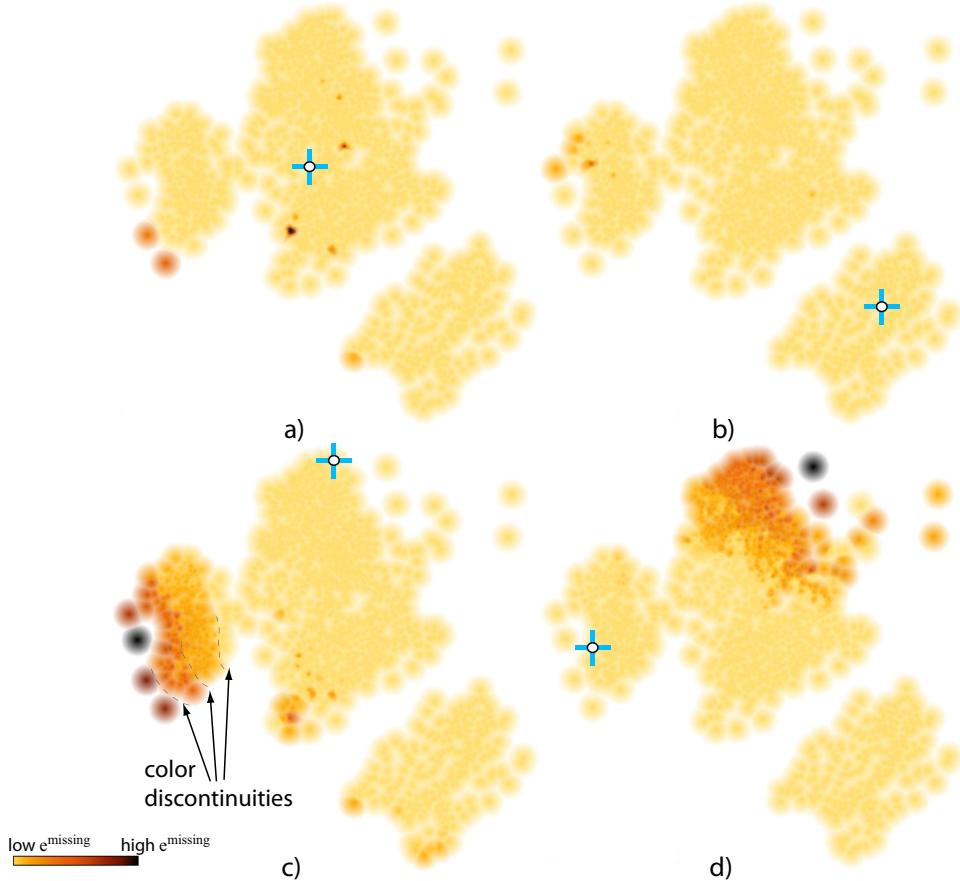


Figure 3: Missing neighbors view for different selected points. Selections are indicated by markers (see Sec. 4.4).

338 different normalizations. However, the proposed visualizations
 339 are quite different. Aupetit uses so-called ‘segment Voronoi
 340 cells’ (SVCs). SVCs essentially achieve piecewise-constant
 341 interpolation of the values e_k^{false} , defined on the edges E_k of
 342 each Delaunay triangle T , over T ’s area, by splitting T in three
 343 sub-triangles using its barycenter. In contrast, our interpola-
 344 tion (Eqn. 8) is C^∞ over T . Also, our triangles are increasingly
 345 transparent far away from their edges (Eqn. 9). Comparing our
 346 results (*e.g.* Figs. 2, 9 (a,d,g)) with SVCs (*e.g.* Figs. 7 (d), 12 (c)
 347 in [33]), we observe that SVCs exhibit several spurious elon-
 348 gated Voronoi cells that do not convey any information. Such
 349 cells do not exist in our visualization due to the transparency
 350 blending. Also, we argue that the artificial SVC edges linking
 351 projected points with Delaunay triangulation barycenters do not
 352 convey any information, but only make the visualization more
 353 complex. Such edges do not exist in our visualization due to
 354 our continuous interpolation.

355 4.4. Missing neighbors view

356 Besides false neighbors, projection errors (and subsequent
 357 misinterpretations) can also be caused by missing neighbors
 358 (case B, Sec. 3). Visualizing this by a space-filling method like
 359 for the aggregate error or false neighbors is, however, less easy.
 360 Given a projected point \mathbf{q} , its missing neighbors can be any-
 361 where in the projection, and are actually by definition far away
 362 from \mathbf{q} . To locate such neighbors, we would need to visualize a
 363 many-to-many relation between far-away projected points.

364 We first address this goal by restraining the question’s scope:
 365 Given a *single* point \mathbf{q}_i , show which of the other points $D^m \setminus \mathbf{q}_i$
 366 are missing neighbors for \mathbf{q}_i . For this, we first let the user select
 367 \mathbf{q}_i by means of direct brushing in the visualization. Next, we
 368 compute the error $e_i^{missing} = \max_{j \neq i}(e_{ij}, 0)$, *i.e.*, the degree to
 369 which \mathbf{q}_j is a missing neighbor for \mathbf{q}_i , and visualize $e^{missing}$ by
 370 the same technique as for the aggregated error (Sec. 4.2).

371 Fig. 3 shows this for the Segmentation dataset, using the
 372 same heat colormap as in Fig. 2. In Figs 3 (a,b), we selected two
 373 points deep inside the central, respectively the lower-right point
 374 groups in the image. Since Figs. 3 (a,b) are nearly entirely light-
 375 colored, it means that these points have few missing neighbors.
 376 Hence, the 2D neighbors of the selected points are truly *all* the
 377 neighbors that these points have in n D. In Figs. 3 (c,d), we next
 378 select two points located close to the upper border of the large
 379 central group and the left border of the left group respectively.
 380 In contrast to Figs. 3 (a,b), we see now an increasingly darker
 381 color gradient as we go further from the selected points. This
 382 shows that points far away from these selections are actually
 383 projected *too* far, as they are actually more similar than the
 384 projection suggests. This is a known (but never visualized as
 385 such) issue of many DR methods, which have trouble in em-
 386 bedding high-dimensional manifolds in 2D: points close to the
 387 embedding’s *border* are too far away from other points in the
 388 projection. Another interesting finding is that the color-coded
 389 Figs. 3 (c,d) do not show a *smooth* color gradient: We see, es-

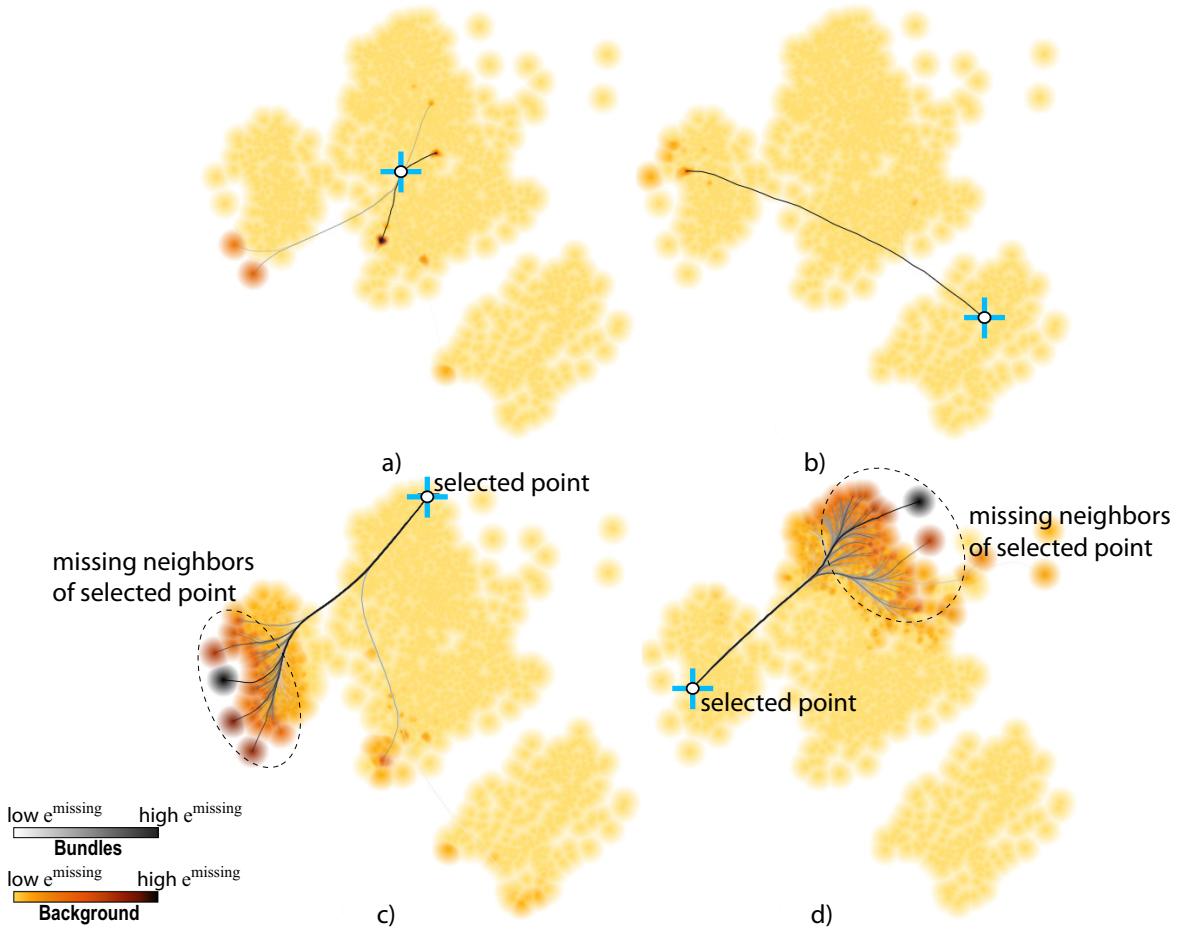


Figure 4: Missing neighbors finder view for four selected points. Selections are indicated by markers (see Sec. 4.5).

pecially in Fig. 3 (c) that the colors appear grouped in several ‘bands’, separated by discontinuities. In other words, the projection method suddenly increases the error as we get over a certain maximal 2D distance.

The missing neighbors view is related to the proximity view of Aupetit [33]. In both views, a point i is selected and a scalar value, related to this selection, is plotted at all other points $j \neq i$. For Aupetit, this is the distance $m_j^{prox} = d(\mathbf{p}_i - \mathbf{p}_j)$ (normalized by its maximum). For us, it is the error $e_j^{missing}$. Both the distance and $e^{missing}$ have, in general, the tendency to be small at points j close in 2D to the selected point i , and increase farther off from point i . However, the two quantities are different and serve different purposes. Visualizing m^{prox} is useful in finding points located within some distance to the selection i . Finding projection errors is only *implicitly* supported, as these appear as non-monotonic variations in the m^{prox} signal. In contrast, $e^{missing}$ specifically emphasizes points projected too far, rather than conveying the absolute distance. Thus, our visualization helps locating projection errors rather than assessing proximity.

4.5. Missing neighbors finder

Although providing details for single points, the views in Sec. 4.4 cannot show missing neighbors for an entire dataset. We address this goal by a different method, as follows. Consider all positive values of e_{ij} . By definition, these give all

point-pairs which are projected too far away. We sort these values decreasingly, and select the largest ϕ percent of them, where ϕ is a user-provided value. The selected values give the point pairs which are worst placed in terms of overestimating their true similarity. We next construct a graph $G = (V, E)$ whose nodes V are the projected points \mathbf{q}_i present in such point pairs, and edges E indicate the pairs, with e_{ij} added as edge weights. Next, we draw G using the KDEEB edge bundling technique [41], which provides robust, easy to use, and real-time bundling of graphs with tens of thousands of edges on a modern GPU. We color the bundled edges based on their weight using a grayscale colormap (with white mapping low and black mapping high weights), and draw them sorted back-to-front on weight and with an opacity proportional to the same weight. The most important edges thus appear atop and opaque, and the least important ones are at the bottom and transparent.

Fig. 4 shows this visualization, which we call the missing neighbors finder, with bundles that connect a single selected point with its most important missing neighbors (bundles connecting multiple points are discussed later on). The background images show $e^{missing}$ (Sec. 4.4). Dark bundle edges attract attention to the most important missing neighbors. For the selected points in images (a) and (b), we see that there are only very few and unimportant missing neighbors (few half-transparent edges). For the selected points in images (c) and (d), the situa-

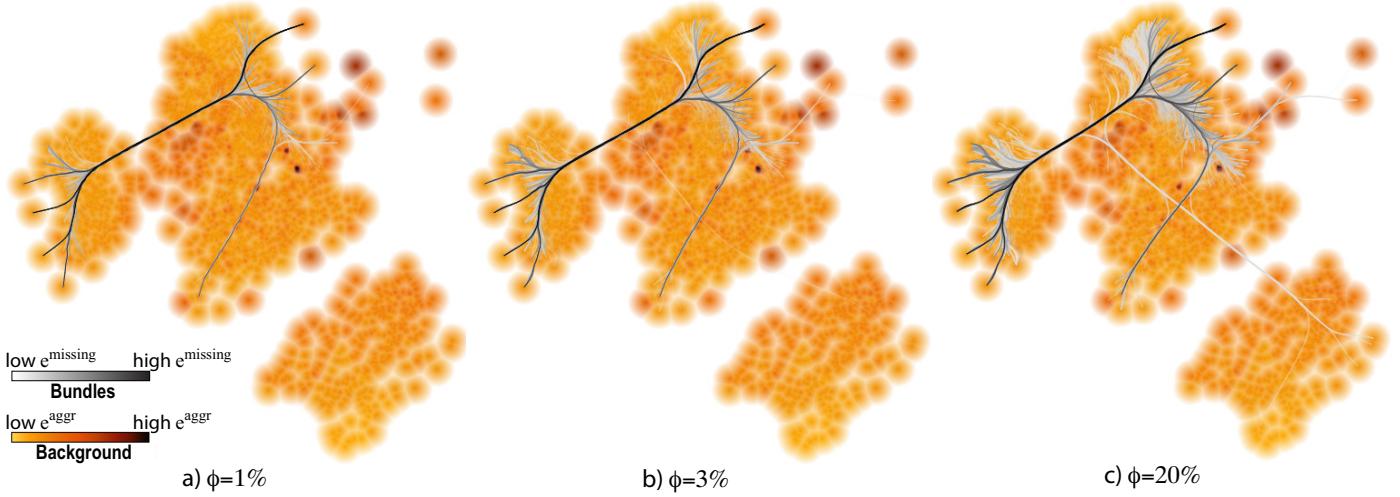


Figure 5: Missing neighbors finder view, all point pairs, for different ϕ values (see Sec. 4.5).

tion is different, as the bundles are thicker and darker. Bundle fanning shows the *spread* of missing neighbors for the selected points: In image (c), these are found mainly in the left point group, with a few also present in the lower part of the central group. In contrast, all missing neighbors of the point selected in image (d) are at the top of the central group.

The main added value of the missing neighbors finder appears when we visualize the many-to-many relations given by all projected points. Fig. 5 shows this result for three values of ϕ for the Segmentation dataset. The background shows now the aggregated error (e^{agg} , Sec. 4.2). We color bundles from black for largest error e_{ij} to white for largest error above the user-provided parameter ϕ . Image (a) shows the $\phi = 1\%$ worst missing-neighbor point-pairs. These link the top-right area of the central group with the left frontier of the left group. Adding more missing neighbor pairs to the view (image (b), $\phi = 3\%$) strengthens this impression. Adding even more missing neighbor pairs (image (c), $\phi = 20\%$) reveals additional missing neighbor pairs between the two areas indicated above (light gray parts of thick top bundle), and also brings in a few missing neighbors between these areas and the lower-right point group (light gray thin bundle going to this group). Nearly all bundles appear to connect point pairs located on the *borders* of the projection. This strengthens our hypothesis that such point pairs are challenging for the LAMP projection, which we noticed using the interactive missing neighbors view (Sec. 4.4). However, as compared to that view, the bundled view shows all such point pairs in a single go, without requiring user interaction.

4.6. Group analysis views

As outlined in Sec. 3, the false and missing neighbors issues for individual points become, at group level, the problems of false and missing group members respectively. We next propose two visualizations that assist in finding such issues.

First, let us refine the notion of a group. Given the tasks in Sec. 3 (C), a group $\Gamma \subset D^m$ is a set of projected points which form a visually well-separated entity. When users see points in a group, they understand that these share some commonality,

but are different from points in other groups. In the LAMP projection of our Segmentation dataset, we see three such groups (Figs. 1-5). Group perception is, obviously, subject to many factors such as user preferences and level-of-detail at which one focuses. However, once a user has established which are the groups (s)he sees in a visualization, the false and missing membership issues become relevant.

We allow users to select groups in a given projection by several mechanisms: direct interactive selection, mean-shift clustering [42], and upper thresholding of the point density [43]. Other user-controlled methods can be used if desired, e.g., K -means or hierarchical agglomerative clustering e.g. [44, 45]. The actual group selection mechanism is further of no importance to our visualization method. We next render each obtained group $\Gamma = \{\mathbf{q}_i\}$ by the shaded cushion technique in [46] as follows. First, we compute a density map $\rho(\mathbf{x}) = \sum_{\mathbf{q} \in \Gamma} K(\mathbf{x} - \mathbf{y})$, where K is an Epanechnikov kernel of width equal to the average inter-point distance δ in Γ , following [42]. Next, we compute a threshold-set Γ_δ of ρ at level δ , and its distance transform DT_{Γ_δ} . Finally, we render a RGBA texture over Γ_δ , where we set the color a fixed hue (light blue in our case) and the transparency A to $\sqrt{DT_{\Gamma_\delta}}$.

Having now groups both as a data structure and also shown in the visualization, we adapt the missing neighbors and finder techniques (Secs. 4.4, 4.5) to show missing group members. For this, we compute a value

$$e_{\Gamma}^{missing}(\mathbf{q}_i) = \begin{cases} \min_{\mathbf{q}_j \in \Gamma}(e_{ij}) & \text{if } \mathbf{q}_i \notin \Gamma \\ 0 & \text{otherwise} \end{cases} \quad (10)$$

at each projected point \mathbf{q}_i , and visualize $e_{\Gamma}^{missing}$ using the same technique as for missing neighbors.

Fig. 6 (a,b) show two missing group members views. The shaded cushions show the three groups identified in our Segmentation dataset. Several points fall outside of all groups. This is normal, in general, e.g. when the user cannot decide to which group to associate a point. In image (a), we select the bottom group Γ_{bottom} . The underlying color map shows now

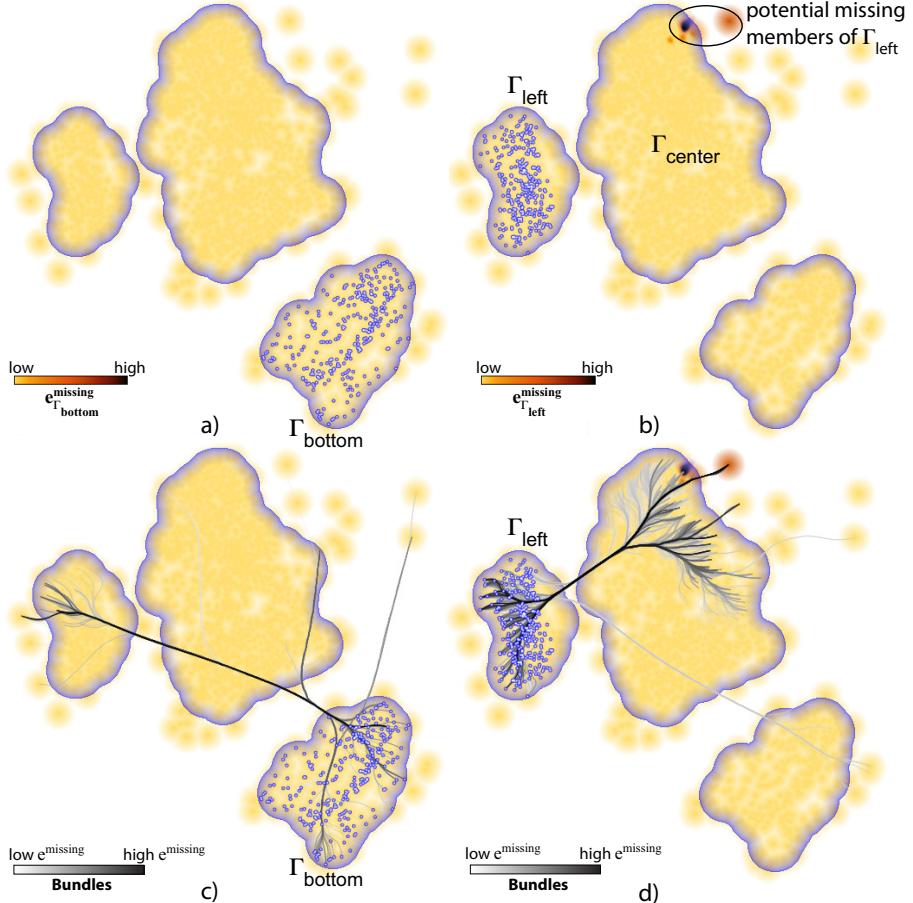


Figure 6: Missing members for two point groups. Points in the selected groups are drawn as marked (see Sec. 4.6).

⁵⁰⁶ $e_{\Gamma_{bottom}}^{missing}$, (Eqn. 10). All points appear light yellow. This means
⁵⁰⁷ that, with respect to Γ_{bottom} seen as a whole, no points are pro-
⁵⁰⁸jected too far, so Γ_{bottom} has no missing members. In image
⁵⁰⁹ (b), we do the same for the left group Γ_{left} . The image now
⁵¹⁰ appears overall light yellow, except for a small dark-red spot
⁵¹¹ in the upper-right corner of the central group Γ_{center} . Here are
⁵¹² a few points which are placed too far from any point in Γ_{left} .
⁵¹³ These are highly likely to be missing members of Γ_{left} . To ob-
⁵¹⁴tain more insight, we now use the bundle view in Sec. 4.5, with
⁵¹⁵ two changes. First, we build only bundles that have an endpoint
⁵¹⁶ in the selected group. Secondly, we consider all edges rather
⁵¹⁷ than showing only the most important ones. Image (c) shows
⁵¹⁸ the bundle view for Γ_{bottom} . We see only a few bundled edges,
⁵¹⁹ ending at a small subset of the points in Γ_{bottom} . This strength-
⁵²⁰ens our hypothesis that there are no points outside Γ_{bottom} which
⁵²¹should be placed closer to *all* points in Γ_{bottom} – or, in other
⁵²²words, that Γ_{bottom} has no missing members. Image (d) shows
⁵²³the bundled view for Γ_{left} . The bundle structure tells us that
⁵²⁴the top-right part of Γ_{center} contains many missing neighbors of
⁵²⁵ Γ_{left} . In particular, we see dark bundle edges that connect to
⁵²⁶dark-red points. This is a strong indication that these points can
⁵²⁷indeed be missing members of Γ_{left} . For a final assessment, the
⁵²⁸user can interactively query the discovered points’ details (at-
⁵²⁹tribute values) and, depending on these, finally decide if these
⁵³⁰points are missing group members or not.

531 4.7. Projection comparison view

⁵³² Consider running the same DR algorithm with two different
⁵³³parameter sets, or projecting a dataset by two different DR al-
⁵³⁴gorithms. How to compare the results from the viewpoint of
⁵³⁵neighborhood preservation? Subsequent questions are: Which
⁵³⁶points that were (correctly) placed close to each other in one
⁵³⁷projection are now ‘pulled apart’ in the other projection? Do
⁵³⁸the two projections deliver the same groups of points?

To answer such questions, we propose the *projection com-
⁵³⁹parison view*. The view reads two projections D_1^m and D_2^m of the
⁵⁴⁰same input dataset D^n . For each point-pair $(\mathbf{q}_i^1 \in D_1^m, \mathbf{q}_i^2 \in D_2^m)$,
⁵⁴¹we compute a displacement

$$e_i^{disp} = \frac{\|\mathbf{q}_i^1 - \mathbf{q}_i^2\|}{\max_i \|\mathbf{q}_i^1 - \mathbf{q}_i^2\|}. \quad (11)$$

⁵⁴² We next build a graph whose nodes are points in $D_1^m \cup D_2^m$.
⁵⁴³ Edges relate point pairs $(\mathbf{q}_i^1 \in D_1^m, \mathbf{q}_i^2 \in D_2^m)$, and have the val-
⁵⁴⁴ues e^{disp} as weights. We visualize this graph via edge bundling,
⁵⁴⁵as for the missing neighbors finder (Sec. 4.5).

⁵⁴⁶ Fig. 7 (a) shows a view where we compare the Segmentation
⁵⁴⁷dataset projected via LAMP (red points, D_1^m) and LSP (green
⁵⁴⁸points, D_2^m). The two projections are quite similar, since red and
⁵⁴⁹green points occur together in most cases. However, this image
⁵⁵⁰does not tell if the two projections create the same *groups* of
⁵⁵¹points, since we do not know how red points match the green

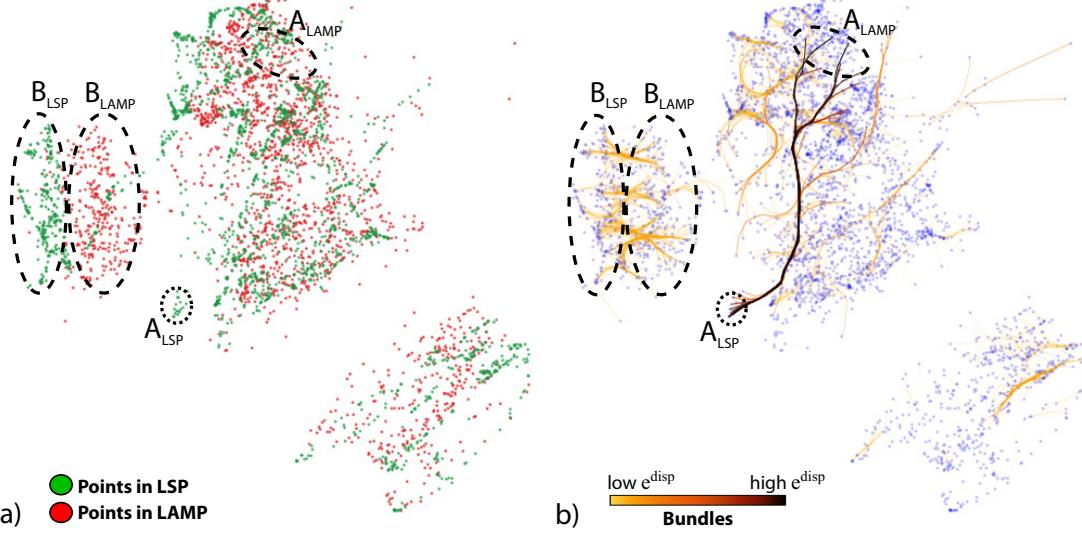


Figure 7: Comparison of two projections. (a) LAMP (blue) and LSP (red) points. (b) Bundles show corresponding point groups in the two projections (see Sec. 4.7).

549 ones. Fig. 7 (b) shows the projection comparison view for this
550 case. We immediately see a thin dark bundle in the center:
551 This links corresponding points which differ the most in the
552 two projections. Correlating this with image (a), we see that
553 LSP decided to place the respective points at the bottom (A_{LSP})
554 of the central group, while LAMP moved and also spread out
555 these points to the top (A_{LAMP}). However, points around the
556 locations A_{LSP} and A_{LAMP} do not move much between the two
557 projections, as we see only light-colored bundles around these
558 locations, apart from the dark bundle already discussed. Hence,
559 the motion of these points indicates a neighborhood problem
560 in one or both of the projections. Indeed, if e.g. the points in
561 A were correctly placed by LAMP (into A_{LAMP}), then the de-
562 cision of LSP to move the point-group A all the way up in the
563 visualization (to A_{LSP}) should also have moved the neighbors
564 of A_{LAMP} . Since this does not happen, A_{LSP} cannot be close to
565 the same points that A_{LAMP} was. A similar reasoning applies
566 if we consider that A_{LSP} is correct – it then follows that A_{LAMP}
567 cannot be correctly placed with respect to its neighbors.

568 Apart from this salient dark-colored bundle, we see many
569 shorter and light-colored bundles. These show smaller-scale
570 displacements between the two projections. For instance, we
571 see how the red points at the right of the left group (B_{LAMP}) are
572 moved to the left (B_{LSP}) of the same group. As these bundles
573 fan out relatively little, do not have many crossings, and they
574 are short, it means that B_{LSP} is almost a *translation* to the left
575 of B_{LAMP} , so the two projections depict the same structure of
576 the left group. Also, we do not see any bundle exiting this left
577 group. This means that both LAMP and LSP keep all points
578 in this group together. Finally, in the bottom-right group we
579 see just a very few short light-colored bundles. Most points in
580 this group do not have any bundles connected to them. This
581 means that e^{disp} for these points is very small (yielding thus
582 very short, nearly transparent, bundles). From this, we infer
583 that LAMP and LSP produce very similar layouts for this group.
584 If users are interested only to spot the most salient differences
585 between two projections, and want to ignore such small-scale

586 changes, this can be easily obtained by mapping e_i^{disp} to bundle-
587 edge transparency.

588 4.8. Usage scenario

589 Considering that the user is offered quite a few different
590 views to analyze projection errors, each with specific features
591 and goals, the next question arises: How to put all these views
592 together to form a coherent *usage scenario* for a common anal-
593 ysis task? Below we propose such a usage scenario. The view
594 names herein refer to the respective techniques presented earlier
595 in this section.

596 **Step 1.** Start with the *Aggregated Error* view. This shows an
597 overview of the error at all points, without a distinction be-
598 tween false or missing neighbors. Next, check if (a) there are
599 regions or groups with substantial errors or (b) the overall error
600 is low. Case (b) indicates that the projection is quite good and
601 that nothing else needs to be improved. In case (a), continue
602 with steps 2, 3, and 4.

603 **Step 2.** The *Missing Neighbors Finder* view can be enabled
604 and disabled freely over the *Aggregated Error* view to show the
605 most important missing neighbors between all points. The user
606 should notice now whether this view shows bundles having high
607 error values (i.e. dark-colored). If so, there are important miss-
608 ing neighbors between the groups connected by such bundles.
609 These groups must be further analyzed with the *Group Analy-
610 sis Views*. If not, i.e. the bundles are colored (light) gray, this
611 tells that the projection is good and, although there are missing
612 neighbors, they are in a low error range and should not threaten
613 the projection interpretation.

614 **Step 3.** Points, groups or regions found problematic in steps 1
615 and 2 are now analyzed in more detail using the *False Neigh-
616 bors* and *Missing Neighbors* views. For groups detected in step
617 1 the most important thing is to find out exactly what kind of er-
618 ror is present: Are they (a) wrongly placed with respect to each

619 other and other close points (false neighbors) or (b) in relation
620 to far away points that should be closer (missing neighbors)?
621 For groups detected in step 2, the error is already identified
622 from the beginning: They have a high rate of missing neigh-
623 bors. In this case, the question to be answered is: Which points
624 are exactly the problematic ones inside the detected groups, or
625 where exactly do the relations (bundle edges) with the highest
626 errors start and end from? By using these two views, the user
627 should be able to establish exactly which are the more problem-
628 atic points (or groups), and what kind of error these have.

629 *Step 4.* Knowing now where exactly errors occur, we consider
630 the next questions: (1) Are such errors really a problem? (2)
631 Do they show unexpected results related to how the projection
632 should work with the provided data? (3) Are the problematic
633 points important for the analysis task at hand? If questions (1-
634 3) all answer ‘no’, then we have a good projection for our data
635 and analysis task, and our analysis stops. If any question (1-
636 3) answers yes, then the user must improve the projection of
637 problematic points, as follows. If the user is a projection de-
638 signer testing the accuracy of a new method, (s)he should go
639 back to the algorithm and use the new insight gotten from this
640 analysis to improve that algorithm. If the user has no access to
641 the projection implementation, the solution is to re-execute the
642 analysis from step 1 with either (i) a new projection algorithm
643 that might better fit the specific data and task; or (ii) a new set
644 of parameters for the same algorithm. The new results can be
645 compared with the old ones to determine if the errors have de-
646 creased or if the errors moved into a new region where they are
647 not as important for the task at hand. For the second task, the
648 *Projection Comparison View* can be used.

649 5. Applications

650 We now use our views to study several projections for several
651 parameter settings – thus, to explore the space P that controls
652 the creation of a DR projection. First, we present the datasets
653 used (Sec. 5.1), the studied projection algorithms (Sec. 5.2), and
654 their parameters (Sec. 5.3). Next, we use our views to explore
655 the considered parameter settings (Secs. 5.4, 5.5).

656 5.1. Description of Datasets

657 Apart from the Segmentation dataset used so far, we consider
658 the following datasets:

659 **Freefoto:** contains 3462 images grouped into 9 unbalanced
660 classes [47]. For each image, we extract 130 BIC (border-
662 interior pixel classification) features. Such features are widely
663 used in image classification tasks [48].

664 **Corel:** composed of 1000 photographs that cover 10 specific
665 subjects. Similarly to the Freefoto dataset, we extract for each
667 image a vector of 150 SIFT descriptors [49].

668 **News:** contains 1771 RSS news feeds from BBC, CNN,
669 Reuters and Associated Press, collected between June and
671 July 2011. The 3731 dimensions were created by removing

672 stopwords, employing stemming and using term-frequency-
673 inverse-document-frequency counts. We manually classified
674 the data points based on the perceived main topic of the news
675 feed resulting in 23 labels. Given the imprecision of the manual
676 classification and the restriction to have one topic per point, the
677 labels are unbalanced for a number of points. Also, for other
678 points (with different labels), we can still have a high similarity
679 of content.

680 **Sourceforge:** This publicly available dataset contains 24 soft-
681 ware metrics computed on 6773 open-source C++ software
682 projects from the sourceforge.net website [50]. Metrics include
683 classical object-oriented quality indicators such as coupling, co-
684 hesion, inheritance depth, size, complexity, and comment den-
685 sity [51], averaged for all source code files within a project.

687 5.2. Description of Projections

688 We detail next the projection algorithms whose parameter
689 spaces we will next study. We chose these particular algo-
690 rithms based on their availability of documented parameters,
691 scalability, genericity, presence in the literature, and last but
692 not least availability of a good implementation.

693 **LSP:** The Least Squares Projection [4] uses a force-based
694 scheme to first position a subset of the input points, called
695 control points. The remaining points in the neighborhood
696 of the control points are positioned using a local Laplace-
698 like operator. Overall, LSP creates a large linear system
699 that is strong in local feature definition. LSP is very precise
700 in preserving neighborhoods from the n D space to the 2D space.

701 **PLMP:** The Part-Linear Multidimensional Projection
702 (PLMP) [10] addresses computational scalability for large
703 datasets by first constructing a linear mapping of the control
704 points using the initially force-placed control points. Next,
705 this linear mapping is used to place the remaining points, by a
706 simple and fast matrix multiplication of the feature matrix with
708 the linear mapping matrix.

709 **LAMP:** Aiming to allow more user control over the final lay-
710 out, the Local Affine Multidimensional Projection (LAMP) [3]
712 provides a user-controlled redefinition of the mapping matrix
713 over a first mapping of control points. LAMP also works by
714 defining control points, which are used to build a family of
715 orthogonal affine mappings, one for each point to project.
716 LAMP has restrictions regarding the number of dimensions
717 against the number of points. Also, LAMP cannot directly
718 work with distance relations, *i.e.*, it needs to access the n D
719 point coordinates. However, LAMP is very fast, without com-
720 promising the precision reached, for instance, by LSP. Both
721 LSP and LAMP can be controlled by a number of parameters,
722 such as the control point set.

723 **Pekalska:** Another class of projection techniques works with
724 optimization strategies. These are, in general, quite expensive
725 computationally. To improve speed, Pekalska *et al.* [21] first
727 embeds a subset of points in 2D by optimizing a stress function.
728 Remaining points are placed using a global linear mapping,

729 much like LAMP and LSP.

730 **ISOMAP:** The ISOMAP technique [14] is an extension of
731 classical Multidimensional Scaling (MDS) that aims to capture
732 nonlinear relationships in the dataset. ISOMAP replaces the
733 input distance between point pairs by an approximation of the
734 geodesic distance given by the shortest path on a graph created
735 connecting neighbor points in the original space with the origi-
736 nal distance as weight. The final 2D coordinates are computed
737 via a conventional MDS embedding with calculations of eigen-
738 values over the distance relations of the previous step.

740 5.3. Description of parameters to analyze

741 Most techniques that initially project control points use a
742 simplified iterative force-based algorithm, such as the one of
743 Tejada *et al.* [27]. The number of *iterations* of force-based
744 placement influences the control points' positions, and is, thus,
745 a relevant parameter. LSP control points are typically the cen-
746 troids of clusters obtained from a clustering of the input dataset.
747 The *number of control points* is thus a second relevant param-
748 eter for LSP. To position points in the neighborhood of a given
749 control point, LSP solves a linear system for that neighborhood.
750 The neighborhood size (*number of neighbors*) is a third relevant
751 parameter.

752 In LAMP, the affine mappings are built from a neighborhood
753 of control points. The size of the control point set used to build
754 the mapping, expressed as a *percentage* of the size of the con-
755 trol point set, is the main parameter here. The choice of con-
756 trol points and the choice of the initial projection of the con-
757 trol points are also parameterizable, just as for LSP, PLMP, and
758 Pekalska. However, in LAMP, these parameters are mainly in-
759 teractively controlled by the user, and thus of a lesser interest to
760 our analysis.

761 ISOMAP, just as the previous methods, also requires the ex-
762 pression of neighborhoods. The main, and frequently only, ex-
763 posed parameter of ISOMAP is the number of *nearest neigh-
764 bors* that defines a neighborhood.

765 5.4. Overview comparison of algorithms

766 To form an impression about how the goals outlined in Sec. 3
767 are better, or less well, satisfied by LAMP, LSP, PLMP, and
768 Pekalska, we start with an overview comparison.

769 Figure 8 shows the false neighbors, aggregated error, and
770 most important $\phi = 5\%$ missing neighbors for the Segmen-
771 tation dataset. To ease comparison, color mapping is normal-
772 ized so that the same colors indicate the same absolute values
773 in corresponding views. The aggregate error (top row) is quite
774 similar in both absolute values and spread for all projections,
775 *i.e.*, lower at the plot borders and higher inside, with a few dark
776 (maximum) islands indicating the worse-placed points. Over-
777 all, thus, all studied projections are quite similar in terms of dis-
778 tance preservation quality. The false neighbors views (middle
779 row) show a similar insight: Border points have few false neigh-
780 bors (light colors), and the density of false neighbors increases
781 gradually towards the projections' centers. Although local vari-
782 ations exist, these are quite small, meaning that all studied pro-
783 jections are equally good from the perspective of (not) creating

784 false neighbors. The missing neighbors view (bottom row) is
785 however quite different: By looking at the size and color of
786 the depicted bundles, we see that LSP and Pekalska have much
787 more important missing neighbors than PLMP, while LAMP
788 has the fewest missing neighbors. In all cases, we see bundles
789 that connect borders of the projected point-set. This confirms
790 that all studied projections optimize placement of close points
791 than far-away points. We also see that the missing neighbors are
792 spread differently over the data: For LAMP, there are no bun-
793 dles going to the bottom-right point cluster, showing that this
794 cluster is indeed well separated in the projection, as it should be
795 in relation to the n D data. In contrast, LSP, PLMP, and Pekalska
796 all have bundles going to this cluster, indicating that they place
797 these points too close to the remaining projected points.

798 5.5. Parameter analysis

799 We next refine our overview analysis by selecting two of the
800 studied algorithms: LAMP and LSP. We next vary several of
801 their parameters, and evaluate the resulting projections' quality
802 with respect to this variation.

803 **LAMP - Different control point percentages:** Fig. 9 shows
804 the results of LAMP for the Freefoto dataset with three different
805 values for the *percentage* parameter: 10%, 30% and 50%. The
806 error has been normalized on each view type (column in the
807 figure).

808 First, we see that the final layout of the point cloud does
809 not change drastically while varying the *percentage* parame-
810 ter, only showing a 90 degree clockwise rotation for the value
811 of 30%. While analyzing the false neighbors view, we also see
812 that, while the light brown areas are large – meaning that a mod-
813 erate amount of error can be expected on the whole layout – the
814 dark-colored spots are found nearer to the center. This suggests
815 that LAMP positions the most problematic points in the center,
816 surrounded by the rest of the points. By focusing on the dark
817 spots (points with the largest false neighbor errors) throughout
818 the parameter variation we can see that the value of the largest
819 errors on each result remain similar – no view has many more,
820 or much darker-colored, areas.

821 For the missing neighbors view, we selected a point near the
822 upper border of the layout, marked by a cross in Figs. 9 (b), (e)
823 and (h), since missing neighbors occur mainly on the borders
824 of the projection, as we have already observed in Section 4.4.
825 The dark spot in Fig. 9 (h) is where the largest error occurs over
826 these three views. While in Fig. 9 (b) there are a few orange
827 spots showing moderate error, in Fig. 9 (e) the error decreases
828 considerably, and then increases again in Fig. 9 (h). This sug-
829 gests that using about 30% of neighbors is a good value for
830 avoiding large numbers of missing neighbors. We confirmed
831 this hypothesis on several other datasets (not shown here for
832 brevity). Finally, the aggregated error view shows results very
833 similar to the false neighbors view: More problematic points
834 (dark spots) are pushed to the center, and moderate error is
835 found spread evenly over the entire layout. This shows that,
836 for LAMP, most errors come from false neighbors rather than
837 from missing neighbors.

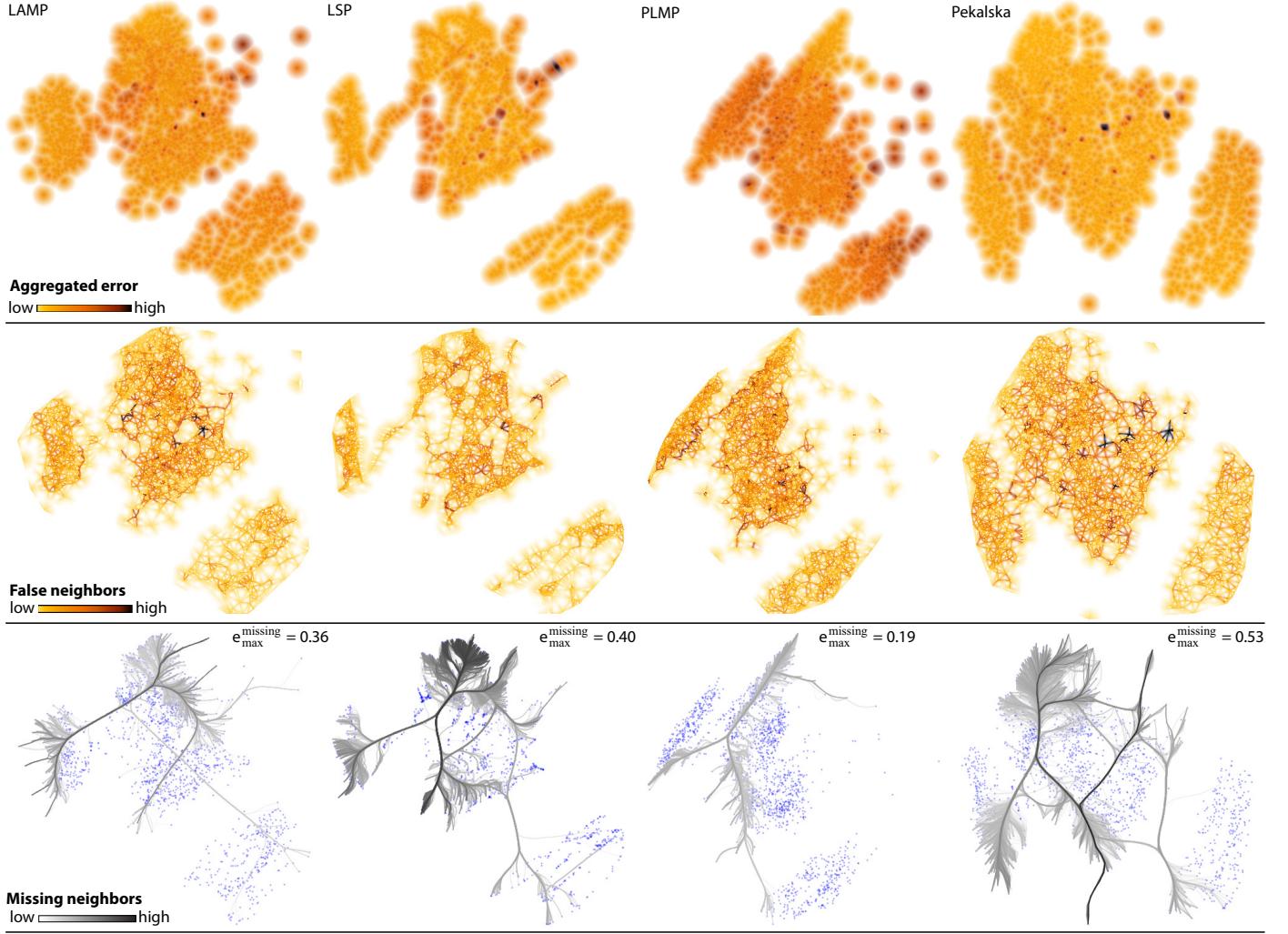


Figure 8: Comparison of LAMP, LSP, PLMP, and Pekalska projections for the Segmentation dataset (see Sec. 5.4)

LSP - Different numbers of control points: Figure 10 shows the same dataset (Freefoto) projected with LSP. The varying parameter is the *number of control points*. We use here the same views as in Fig. 9, and normalized the error in each column. By looking at the false neighbors views, we see a spatial interleaving of light-yellow and orange-brown colored areas in the projection. This contrasts with LAMP (Fig. 9) where the larger missing neighbor errors are consistently located away from the projection border. As the *number of control points* increases, the large error areas get more compact and closer to the projection center, but we see no increase in error severity (the amount of the orange and dark-red spots stays the same). In the missing neighbors views, the dark-colored areas in Fig. 10 (b) disappear largely in images (e) and (h), which means that the missing neighbors severity decreases when our control parameter increases. Comparing this with LAMP (Fig. 9 b,e,h), this shows that LAMP and LSP behave in opposite ways when dealing with missing neighbors. Finally, like for LAMP, the aggregate error views show the worst errors (dark spots) located in the center: The most problematic points are pushed inside by the other points which surround them,

creating a mix of both false neighbors and missing neighbors. The severity of the errors, however, does not change visibly between the three parameter values.

LSP - Different numbers of neighbors: We next examine the effect of a second parameter of LSP: *number of neighbors*. For the Freefoto dataset, we fix 250 control points and vary the number of neighbors to 10, 50 and 100. Fig. 11 shows the results with the missing neighbors finder view. We see that the most significant errors are initially concentrated between groups A, B and C, with C being essentially too far placed from both A and B. Increasing our parameter reduces has a positive impact on solving the missing neighbors problem between groups A and C, bringing them together into the group marked AC. The main missing neighbors are now concentrated in the relationship between groups AC and B. The ‘concentration’ of error given by the parameter increase is, upon further analysis, explainable by the working of LSP: Given a neighborhood N , LSP’s Laplace technique positions all points in N close to each other in the final layout. However, the position of the neighborhoods N_i themselves is given only by the control

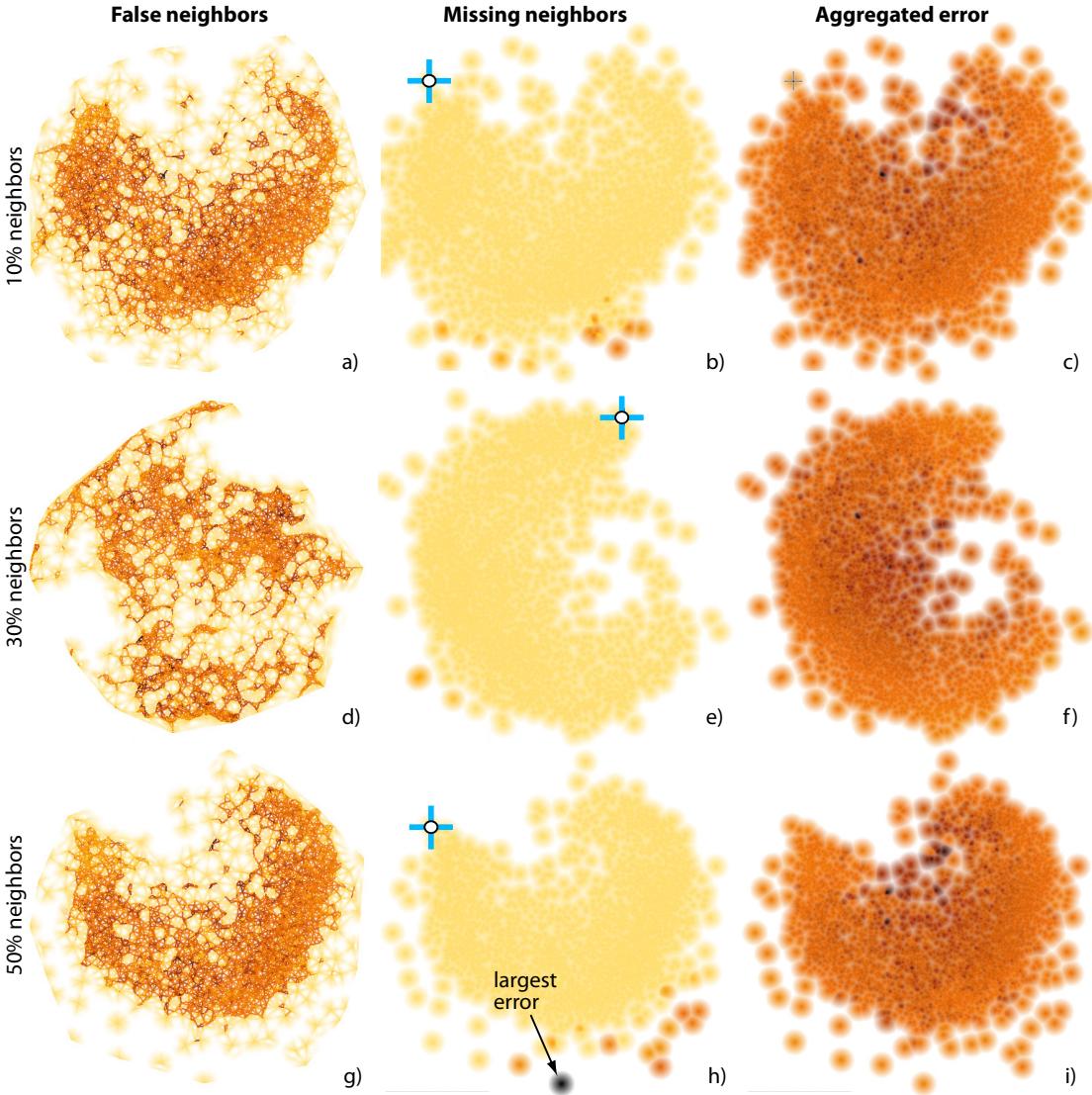


Figure 9: Applications – LAMP algorithm, Freefoto dataset, different neighbor *percentages* per row (see also Fig. 10).

points, which are determined by the initial force-based layout. If this layout suboptimally places two control points i and j too far away from each other, then *all* points within the neighborhoods N_i and N_j end up being too far away from each other. Hence, as the neighborhood size increases, the likelihood to see fewer thick high-error bundles increases. This insight we found is interesting since it was not reported in the LSP literature so far, and it can be explained (once we are aware of it) by the algorithmics of LSP.

LAMP - Different datasets: We next analyze the LAMP technique applied to three different datasets: Corel (1000 elements), Freefoto (3462 elements), and Sourceforge (6773 elements).

The varying parameter is now the input *dataset* itself. The aim is to see whether (and how) errors are affected by the nature of the input data, *e.g.* distribution of similarity, number of dimensions, and number of points. Figure 12 top row shows the false neighbors views. We see here that, while for the first two datasets the behavior of false neighbors is similar to earlier results, for the largest dataset (Sourceforge) there are much fewer

false neighbors. These are located close to the intersection area of the two apparent groups in the image, and on the borders of these groups. This, and the low errors (light colors) inside the groups may indicate that both groups have a high degree of cohesion between their inner elements. The large errors on close to the intersection areas and borders can indicate elements that could be in either group, respectively very different from all other elements. Figure 12 (a) shows a similar pattern: Most false neighbors are located at the ‘star’ shape’s center, while the arms of the star contain elements that are more cohesive. This may indicate that the dataset contains a number of cohesive groups equal to the number of start arms, and elements in the center belong equally to all groups.

While analyzing the missing neighbors for several points selected on the periphery of the projections, we see that the errors are smaller for Figs. 12 (d) and (e), and considerably larger for Fig. 12 (f). For the last image, we selected a point close to the intersection area of the perceived groups. Image (f) shows that this point is *equally* too far placed from most points

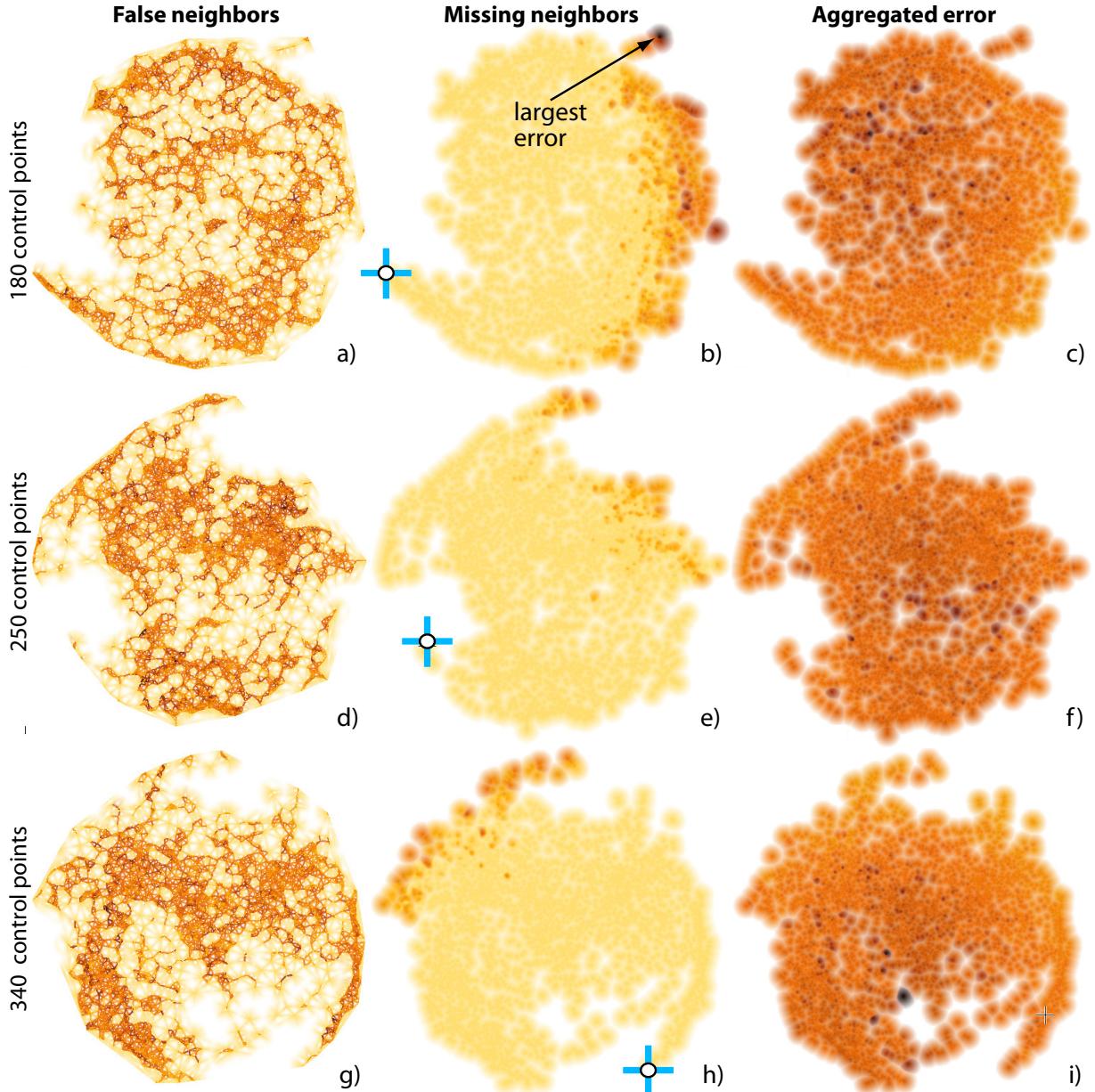


Figure 10: Applications – LSP technique, Freephoto dataset, different numbers of *control points* per row (compare with Fig. 9)

in both perceived clusters. The size and speed of increase of the error (as we get further from this point in the projection space) strongly suggests that the selected point belongs stronger to both perceived groups than the projection indicates. This strengthens our initial hypothesis that the area separating the two groups belongs equally to these groups.

ISOMAP - Different numbers of neighbors: To illustrate a different type of analysis made possible by our work, Fig. 13 shows the effect of changing the number of *neighbors* in ISOMAP on missing group members. Our group Γ of interest, shown first on Fig. 13 (a), is highlighted in images (b-d) by a shaded cushion. Besides the fact that Γ moves from the left of the projection to the right, images (b-d) show how its missing members behave as we change our parameter. At first, in Fig. 13 (b), we see that the most important missing neighbors

are found in two other areas A_1 and A_2 on the far side of the layout. We also notice many black edges, which means that the points in A_1 and A_2 are indeed too far away from all points in the selected group. The relatively large fan-out of the bundles show that the group misses many members, and these are scattered widely over the projection. As the parameter increases, we see in image (c) that the missing members spread out even more, but the severity of the errors decreases (as shown by the lighter colors of $e_{\Gamma}^{missing}$ background). The inner fanning of the edges, inside Γ , is still large, which shows that many group members miss neighbors. Finally, in Fig. 13 (d), issues decrease significantly: We see thinner bundles, which imply less error; the bundle fanning inside Γ is relatively small, meaning that most of Γ 's points do not miss neighbors; and the fan-out of the bundles is smaller, showing that the missing group members are now more concentrated than for the first

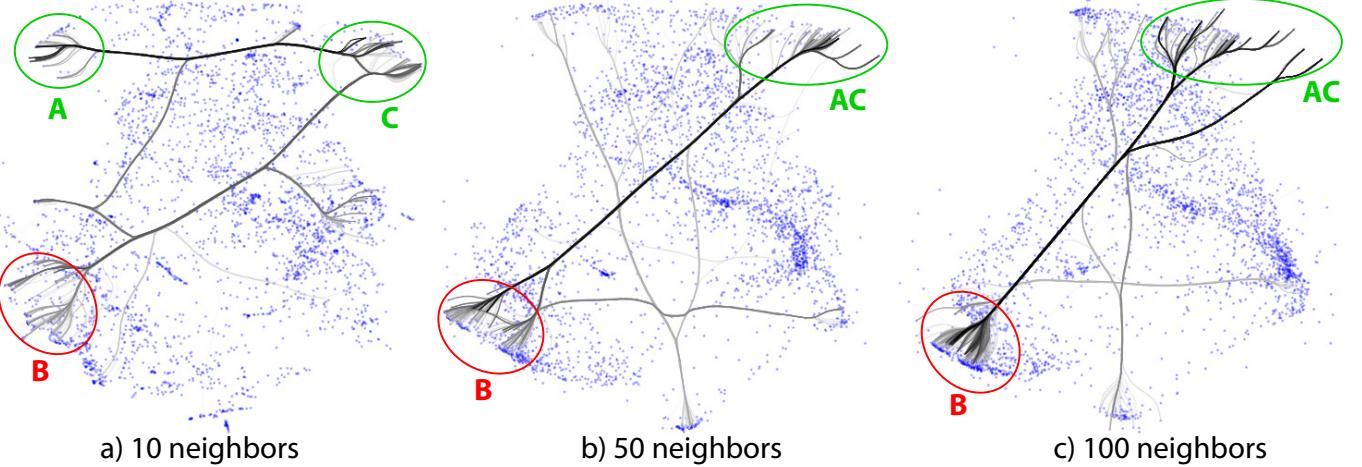


Figure 11: Applications – LSP technique, Freefoto dataset, different numbers of *neighbors*. Bundles show most important missing neighbors.

952 two parameter values. This leads to the conclusion that, for
953 the analyzed group, the increase of the number of neighbors
954 parameter has a positive impact on the final projection quality.
955

956 **LSP - Different numbers of iterations:** The final analysis we
957 present compares two different LSP projections of the same
958 dataset (News), computed using values of 50, respectively 100
959 for the *number of iterations* parameter of the control-point
960 force-directed placement.

961 Figures 14 (a) and (b) show the two LSP projections. In
962 each of them, several high-density groups are visible. These
963 are strongly related news feeds, *i.e.*, which likely share the same
964 topic (see Sec. 5.1). However, without extra help, we cannot *re-*
965 *late* the two projections, *e.g.*, find out (a) if points significantly
966 change places due to the parameter change; (b) which groups
967 in one projection map to groups in the other projection; and (c)
968 whether points in a group in one projection are also grouped in
969 the second projection.

970 To answer question (a), we use the projection comparison
971 view (Sec. 4.7). The result (Fig. 14 (c)) shows that there are
972 many large point shifts; the bundle criss-crossing also shows
973 that groups change places in the projection. This is a first indi-
974 cation that LSP is not visually stable with respect to its num-
975 ber of iterations parameter. Next, we manually select three
976 of the most apparent point groups in one projection, shown in
977 Fig. 14 (a) by the shaded cushions *A,B,C*. We examine these
978 in turn. In Fig. 14 (d), we show how points in group *A* shifted,
979 in the second projection, to a group *A*₁. Virtually all bundled
980 edges exiting *A* end in *A*₁, so the parameter change preserves
981 the cohesion of group *A* (though, not its position in the layout).
982 The same occurs for group *B* (Fig. 14 (e)). However, the pa-
983 rameter change spreads *B* more than *A* – in image (e), we see
984 that *B* maps to three groups, *B*_{1..B}₃. These visualizations thus
985 answer question (b). Group *C* behaves differently (Fig. 14 (f)):–
986 This group is split into two smaller groups *C*₁ and *C*₂ when we
987 change our parameter. For question (c), thus, the answer is par-
988 tially negative: not all groups are preserved in terms of spatial
989 coherence upon parameter change.

990 6. Discussion

991 We have implemented our visualization techniques in C++
992 using OpenGL 1.1, and tested them on Linux, Windows, and
993 OSX. Below we discuss several aspects of our method.

994 **Computational scalability:** For Delaunay triangulation and
995 nearest-neighbor searches, we use the Triangle [52] and
996 ANN [53] libraries. Both can handle over 100K points in
997 subsecond time on a commodity PC. Further, we accelerate
998 imaging operations using GPU techniques. For distance trans-
999 forms, we use [54]. On an Nvidia GT 330M, this allows us to
1000 compute shaded cushions and perform our Shepard interpola-
1001 tion at interactive frame rates for views of 1024² pixels. For
1002 edge bundling, we implemented KDEEB [41] fully on Nvidia’s
1003 CUDA platform. This yields a speed-up of over 30 times
1004 (on average) as compared to the C# implementation in [41]
1005 and allows bundling graphs of tens of thousands of edges in
1006 roughly one second. All in all, we achieve interactive query-
1007 ing and rendering of our views for projections up to 10K points.

1009 **Visual scalability:** Our image-based approaches scale well to
1010 thousands of data points or more, even when little screen space
1011 is available. Moreover, all our techniques have a multiscale
1012 aspect: The parameters α and β (Eqns. 6, 7) effectively control
1013 the visual *scale* at which we want to see false neighbors,
1014 missing neighbors, and the aggregate error. Increasing these
1015 values eliminates spatial outliers smaller than a given size,
1016 thereby emphasizing only coarse-scale patterns (see *e.g.*
1017 Fig. 1). The bundled views (Sec. 4.5) also naturally scales
1018 to large datasets given the inherent property of bundled edge
1019 layouts to emphasize coarse-scale connectivity patterns.

1021 **Genericity:** Our visualizations are applicable to any DR
1022 algorithm, as long as one can compute an error distance
1023 matrix encoding how much 2D distances deviate from their nD
1024 counterparts (Eqn. 2). No internal knowledge of, or access to,
1025 the DR algorithms is needed – these can be employed as black
1026 boxes. This allows us to easily compare widely different DR

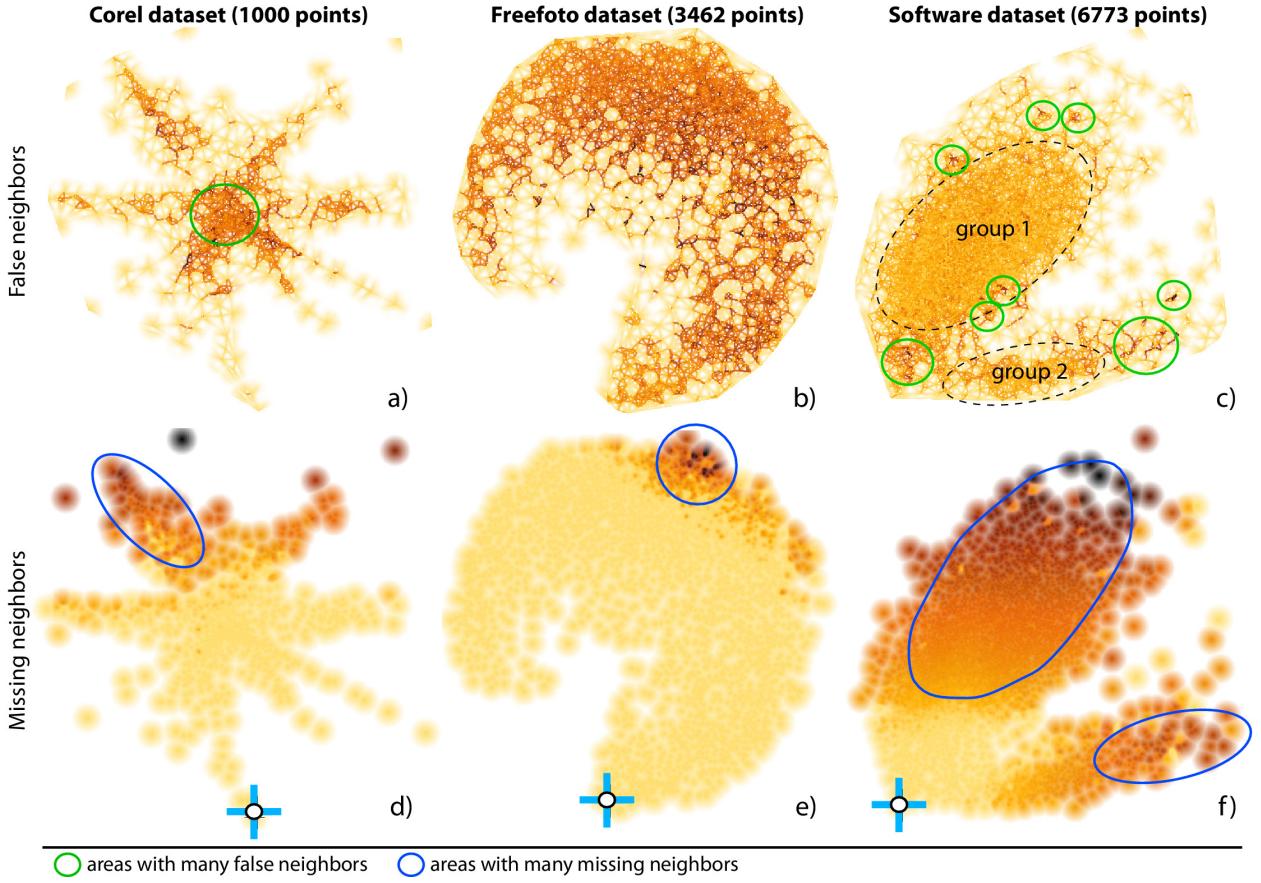


Figure 12: Applications – One algorithm (LAMP), different datasets. Top row: false neighbors. Bottom row: missing neighbors.

1028 algorithms, *e.g.* based on representatives, based on distance
1029 matrices, or based on direct use of the n D coordinates.

1030
1031 **Ease of use:** Our views are controlled by three parameters: α
1032 sets the scale of the visual outliers we want to show; β sets the
1033 radius around a point in which we want to display information,
1034 *i.e.*, controls the degree of space-filling of the resulting images;
1035 ϕ sets the percentage of most important missing neighbors we
1036 want to show. These parameters, as well as the interaction
1037 for selecting point groups (Sec. 4.6) are freely controllable by
1038 users by means of sliders and point-and-click operations.
1039

1040 **Comparison:** Similarly to Van der Maaten *et al.* [32], we use
1041 multiple views showing the same data points to explain a pro-
1042 jection, *e.g.*, the false neighbors, missing neighbors view, miss-
1043 ing neighbors finder, and group-related maps. However, the
1044 multiple maps in [32] are used to actually convey the projection,
1045 so the same point can have different locations and/or weights in
1046 different maps. In contrast, we use multiple views to convey
1047 different quality metrics atop of the same 2D projection. Sim-
1048 ilar to Aupetit [33], our error metrics encode discrepancies in
1049 distances in \mathbb{R}^n vs \mathbb{R}^2 . However, our error metrics are different.
1050 More importantly, our visualizations are different: Our false
1051 neighbors view does not show (a) spurious Voronoi cell edges
1052 far away from data points or (b) cell subdivision edges whose
1053 locations does not convey any information, since we (a) use
1054 distance-based blending and (b) continuous rather than constant

1055 per-cell interpolation (Sec. 4.3). Secondly, our missing neig-
1056 bors finder (Sec. 4.5) can show one-to-many and many-to-many
1057 error relationships, whereas all other methods are constrained to
1058 one-to-one relationships. Finally, we can show errors at group
1059 level, whereas the other studied techniques confine themselves
1060 to showing errors at point level only.

1061 Our projection comparison view is technically related to the
1062 method of Turkay *et al.*, which connects two 2D scatterplots
1063 to each other by lines linking their corresponding points [55].
1064 However, Turkay *et al.* stress that line correspondences only
1065 work for a *small* number of points. In contrast, we use bundles
1066 to (a) show up to thousands of correspondences, and coloring
1067 and blending to encode correspondence importance.
1068

1069 **Findings:** It can be argued that our results are limited, as we
1070 did not decide, using our method, which of the studied DR
1071 algorithms are best. However, this was not the aim of our
1072 work. Rather, our goal was to present a set of visual techniques
1073 that help analyze the effect of parameters on projection quality
1074 for several DR techniques of interest. Deciding whether a
1075 certain degree of quality, *e.g.* in terms of false neighbors,
1076 missing neighbors, grouping problems, or projection stability
1077 is a highly context, dataset, and application-dependent task.
1078 Having such a context, our tools can be then used to assess
1079 (a) which are the quality problems, (b) how parameter settings
1080 affect them, and (c) whether these problems are acceptable
1081 for the task at hand. The same observation applies to the

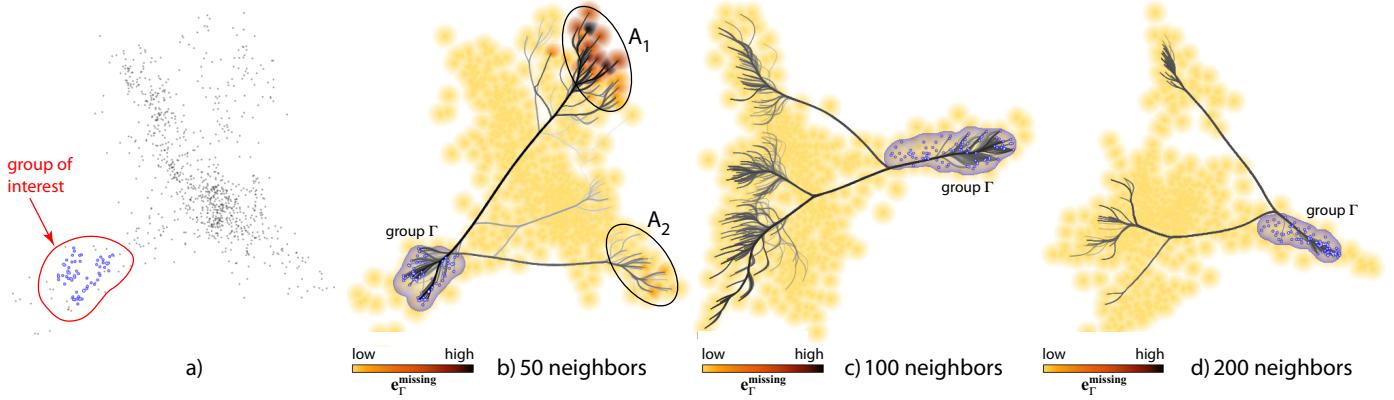


Figure 13: Applications – ISOMAP projection, finding missing group members for different numbers of *neighbors*.

1082 datasets used here. Our analyzes involving these should be
 1083 seen purely as test cases for assessing the quality problems of
 1084 DR projections, and not as findings that affect the underlying
 1085 problems captured by these datasets.

1086 **Limitations:** As outlined by our examples, our visualizations
 1087 can show (a) which projection areas suffer from low quality;
 1088 and (b) how two projections differ in terms of neighborhood
 1089 preservation. However, we cannot directly explain (c) *why* a
 1090 certain DR algorithm decided to place a certain point in some
 1091 position; and (d) how the user should *tune* (if possible) the al-
 1092 gorithm’s parameters to avoid errors in a given area. In other
 1093 words, we can explain the function $f : P$ (Eqn. 1) and its first
 1094 derivatives over P , but not the inverse f^{-1} . This is a much more
 1095 challenging task – currently not solved by any technique we
 1096 know of. Further explaining such second-order effects to help
 1097 users locally fine-tune a projection is subject to future work.
 1098 Secondly, the parameter space P of some DR algorithms can
 1099 be high-dimensional. So far, we can only analyze the variation
 1100 of one or two parameters at a time. Extending this to several
 1101 parameters is a second challenging next topic.

1103

1104 7. Conclusions

1105 We have presented a set of visualization methods for the anal-
 1106 ysis of quality of dimensionality-reduction (DR) algorithms by
 1107 exploration of their parameter settings. By generically mod-
 1108 eling such algorithms as functions from nD to $2D$ in terms of
 1109 their distance-preservation error, we propose several views for
 1110 assessing the distribution of false neighbors, missing neighbors,
 1111 and aggregated projection error at both individual point and
 1112 point-group level. We use several dense-pixel, visually scal-
 1113 able, techniques such as multi-scale scattered point interpola-
 1114 tion and bundled edges to make our methods visually and com-
 1115 putationally scalable to large datasets and also work in a mul-
 1116 tiscale mode. We demonstrate our techniques by analyzing the
 1117 parameters of five state-of-the-art DR techniques. In contrast
 1118 to existing assessments of DR projections by aggregate figures,
 1119 that can only infer overall precision, we offer more local tools to
 1120 examine how neighborhoods and groups are mapped in the fi-
 1121 nal projection. The usage of our techniques is simple and, most

1122 importantly, allows users of DR techniques to study their qual-
 1123 ity without needing to understand complex internal processes
 1124 or the exact role of each parameter in the projections.

1125 Future work can target several directions. First, we plan to
 1126 support ‘what if’ scenarios, *i.e.*, help users to decide how they
 1127 could correct local projection problems by shifting wrongly-
 1128 placed points while dynamically assessing the ensuing overall
 1129 projection errors. Secondly, we plan to explicitly visualize the
 1130 reasons that determine point placement, *i.e.*, depict the nD vari-
 1131 able values which cause points to be placed close to, or far away
 1132 from, each other. Additionally, we intend to provide tools for
 1133 local evaluation of projections customized for specific target au-
 1134 diences. By this, we hope to make the operation of DR algo-
 1135 rithms more transparent and understandable for users ranging
 1136 from algorithm designers to end-users.

1137 Acknowledgements

1138 This research was financially supported by the grant PN-
 1139 II-RU-TE-2011-3-2049 offered by ANCS, Romania, and by
 1140 the research project CAPES/NUFFIC 028/11. Authors also
 1141 wish to acknowledge the Brazilian financial agencies CNPq and
 1142 FAPESP for their support.

1143 References

- [1] Elmqvist N, Dragicevic P, Fekete JD. Rolling the dice: Multidimensional visual exploration using scatterplot matrix navigation. *IEEE TVCG* 2008;14:1141–8.
- [2] Heinrich J, Weiskopf D. Continuous parallel coordinates. *IEEE TVCG* 2009;15(6):1531–8.
- [3] Joia P, Paulovich FV, Coimbra D, Cuminato JA, Nonato LG. Local affine multidimensional projection. *IEEE TVCG* 2011;17:2563–71.
- [4] Paulovich FV, Nonato LG, Minghim R, Levkowitz H. Least square projection: A fast high-precision multidimensional projection technique and its application to document mapping. *IEEE TVCG* 2008;14(3):564–75.
- [5] Cui W, Wu Y, Liu S, Wei F, Zhou MX, Qu H. Context-preserving, dynamic word cloud visualization. *IEEE CG&A* 2010;42–53.
- [6] Faloutsos C, Lin KI. FastMap: a fast algorithm for indexing, data-mining and visualization of traditional and multimedia datasets. In: Proc. ACM SIGMOD Intl. Conf. on Management of Data. 1995, p. 163–74.
- [7] Chen Y, Wang L, Dong M, Hua J. Exemplar-based visualization of large document corpus. *IEEE TVCG* 2009;15:1161–8.
- [8] Paulovich FV, Nonato LG, Minghim R. Visual mapping of text collections through a fast high precision projection technique. In: Proc. IEEE Information Visualization. 2006, p. 282–90.

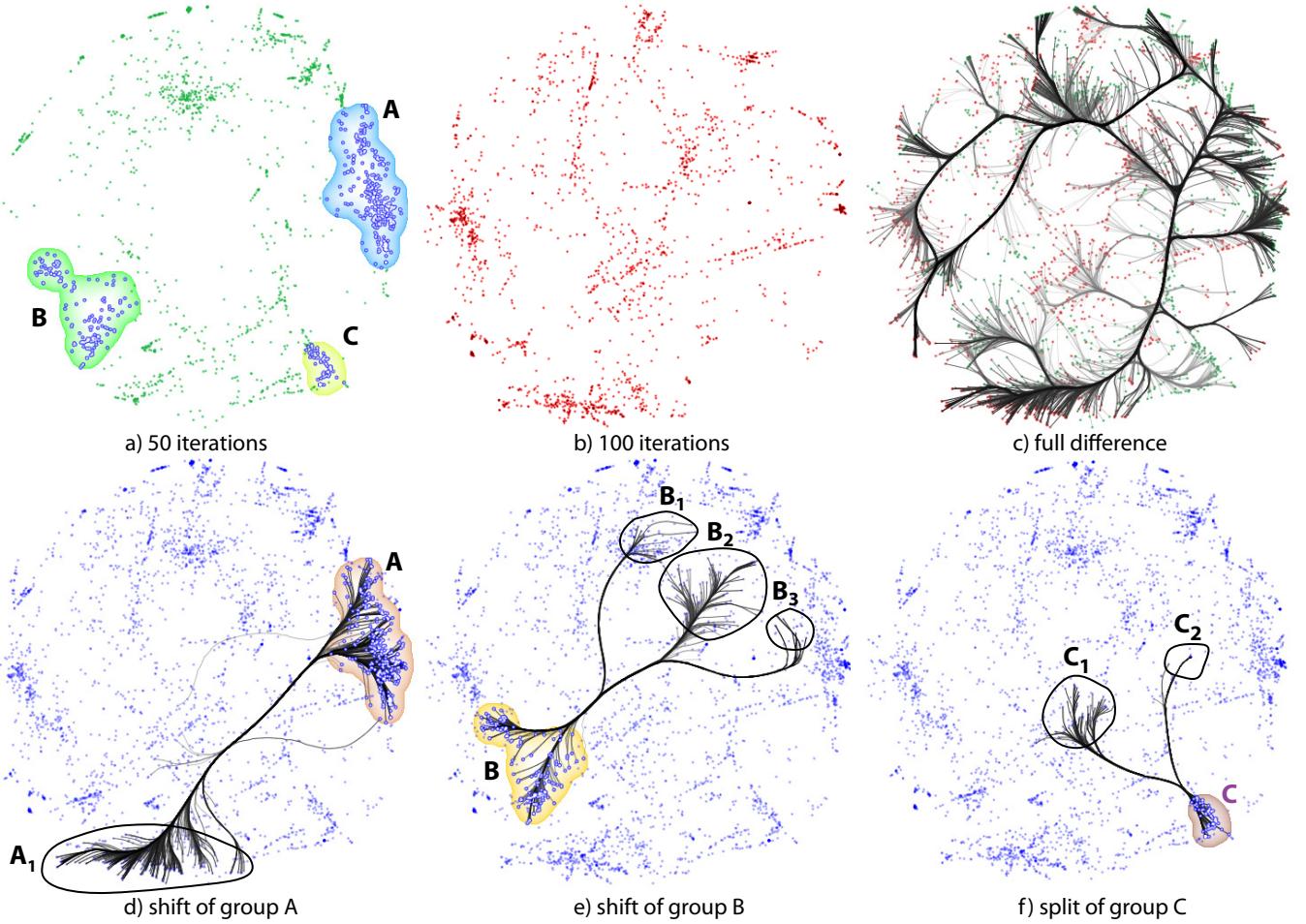


Figure 14: Applications – Shift between two LSP projections, for different numbers of force-directed *iterations*.

- 1164 [9] Daniels J, Anderson EW, Nonato LG, Silva CT. Interactive vector field
1165 feature identification. IEEE TVCG 2010;16:1560–8.
- 1166 [10] Paulovich FV, Silva C, Nonato LG. Two-phase mapping for projecting
1167 massive data sets. IEEE TVCG 2010;16:1281–90.
- 1168 [11] Poco J, Eler D, Paulovich FV, Minghim R. Employing 2D projections for
1169 fast visual exploration of large fiber tracking data. CGF 2012;31:1075–
1170 84.
- 1171 [12] Torgeson WS. Multidimensional scaling of similarity. Psychometrika
1172 1965;30:379–93.
- 1173 [13] Roweis ST, Saul LK. Nonlinear dimensionality reduction by locally lin-
1174 ear embedding. Science 2000;290(5500):2323–6.
- 1175 [14] Tenenbaum JB, de Silva V, Langford JC. A global geometric framework
1176 for nonlinear dimensionality reduction. Science 2000;290(5500):2319–
1177 23.
- 1178 [15] Silva VD, Tenenbaum JB. Global versus local methods in nonlinear di-
1179 mensionality reduction. In: Advances in Neural Information Processing
1180 Systems; vol. 15. MIT Press; 2003, p. 705–12.
- 1181 [16] de Silva V, Tenenbaum J. Sparse multidimensional scaling
1182 using landmark points. Tech. Rep.; Stanford Univ.; 2004.
1183 [window.stanford.edu/courses/cs468-05-winter/Papers/
1184 Landmarks/Silva_landmarks5.pdf](http://window.stanford.edu/courses/cs468-05-winter/Papers/Landmarks/Silva_landmarks5.pdf). Last access: 15/07/2013.
- 1185 [17] Brandes U, Pich C. Eigensolver methods for progressive multi-
1186 dimensional scaling of large data. In: Proc. Graph Drawing. Springer; 2007, p.
1187 42–53.
- 1188 [18] Gansner ER, Koren Y, North S. Graph drawing by stress majorization.
1189 In: Proc. Graph Drawing. Springer; 2004, p. 239–50.
- 1190 [19] Sammon JW. A nonlinear mapping for data structure analysis. ACM
1191 Trans Comp 1969;C-18:401–9.
- 1192 [20] Bronstein M, Bronstein A, Kimmel R, Yavneh I. Multigrid multidimen-
1193 sional scaling. Numer Linear Algebr 2006;13:149–71.
- 1194 [21] Pekalska E, de Ridder D, Duin R, Kraaijveld M. A new method of gen-
1195 eralizing Sammon mapping with application to algorithm speed-up. In:
1196 Proc. 5th Annual Conference of the Advanced School for Computing and
1197 Imaging (ASCI). Delft, NL; 1999, p. 221–8.
- 1198 [22] Eades PA. A heuristic for graph drawing. Congressus Numerantium
1199 1984;42:149–60.
- 1200 [23] Chalmers M. A linear iteration time layout algorithm for visualising high-
1201 dimensional data. In: Proc. IEEE Visualization. 1996, p. 127–31.
- 1202 [24] Frishman Y, Tal A. Multilevel graph layout on the gpu. IEEE TVCG
1203 2007;13:1310–9.
- 1204 [25] Ingram S, Munzner T, Olano M. Glimmer: Multilevel MDS on the GPU.
1205 IEEE TVCG 2009;15(2):249–61.
- 1206 [26] Jourdan F, Melangon G. Multiscale hybrid MDS. In: Proc. IEEE Infor-
1207 mation Visualization. 2004, p. 388–93.
- 1208 [27] Tejada E, Minghim R, Nonato LG. On improved projection techniques
1209 to support visual exploration of multidimensional data sets. Inf Vis
2003;2(4):218–31.
- 1210 [28] Geng X, Zhan DC, Zhou ZH. Supervised nonlinear dimensionality re-
1211 duction for visualization and classification. IEEE Trans Syst Man Cybern
2005;35(6):1098–107.
- 1212 [29] Sedlmair M, Tatu A, Munzner T, Tory M. A taxonomy of visual cluster
1213 separation factors. CGF 2012;31(3):1335–44.
- 1214 [30] Schreck T, von Landesberger T, Bremm S. Techniques for precision-
1215 based visual analysis of projected data. Inf Vis 2010;9(3):181–93.
- 1216 [31] van der Maaten L, Hinton G. Visualizing data using t-SNE. J Mach Learn
1217 Res 2008;9:2431–56.
- 1218 [32] van der Maaten L, Hinton G. Visualizing non-metric similarities in mul-
1219 tiple maps. Machine Learning 2012;87(1):33–5.

- 1222 [33] Aupetit M. Visualizing distortions and recovering topology in continuous
1223 projection techniques. Neurocomputing 2007;10(7-9):1304–30.
- 1224 [34] Lespinats S, Aupetit M. CheckViz: Sanity check and topological clues
1225 for linear and non-linear mappings. CGF 2011;30(1):113–25.
- 1226 [35] Heulot N, Aupetit M, Fekete JD. Proxilens: Interactive exploration of
1227 high-dimensional data using projections. In: EuroVis Workshop on Visual
1228 Analytics using Multidimensional Projections. 2013, p. 11–5.
- 1229 [36] Frank A, Asuncion A. UCI machine learning repository. 2013. www.ics.uci.edu/~mlearn. Last access: 15/07/2013.
- 1230 [37] Paulovich F, Eler D, Poco J, Botha C, Minghim R, Nonato LG. Piece
1231 wise Laplacian-based projection for interactive data exploration and or-
ganization. CGF 2011;30(3):1091–100.
- 1232 [38] Silva C, Paulovich F, Nonato LG. User-centered multidimensional pro-
1233 jection techniques. Comput Sci Eng 2012;14(4):74–81.
- 1234 [39] Brewer C, Harrower M. ColorBrewer. 2013. URL: <http://www.colorbrewer.org>.
- 1235 [40] Rumpf M, Telea A. A continuous skeletonization method based on level
1236 sets. In: Proc. VisSym. Eurographics; 2002, p. 151–9.
- 1237 [41] Hurter C, Ersoy O, Telea A. Graph bundling by kernel density estimation.
CGF 2012;31(3):865–74.
- 1238 [42] Comaniciu D, Meer P. Mean shift: A robust approach toward feature
1239 space analysis. IEEE TPAMI 2002;24(5):603–19.
- 1240 [43] Ester M, Kriegel HP, Sander J, Xu X. A density-based algorithm for
1241 discovering clusters in large spatial databases with noise. In: Proc. 2nd
1242 Intl. Conf. on Knowledge Discovery and Data Mining. 1996, p. 226–31.
- 1243 [44] Jain AK, Murthy M, Flynn P. Data clustering: A review. ACM Comput
1244 Surv 1999;31(3):264–323.
- 1245 [45] Jain AK. Data clustering: 50 years beyond k-means. Pattern Recognit
1246 Lett 2010;31(8):651–66.
- 1247 [46] Telea A, Ersoy O. Image-based edge bundles: simplified visualization of
1248 large graphs. CGF 2010;29(3):843–52.
- 1249 [47] Freefoto . Free stock photo collection. 2013. www.freefoto.com. Last
1250 access: 15/07/2013.
- 1251 [48] Stehling RO, Nascimento MA, Falcão AX. A compact and efficient image
1252 retrieval approach based on border/interior pixel classification. In: Proc.
1253 CIKM. ACM; 2002, p. 102–9.
- 1254 [49] Li J, Wang JZ. Automatic Linguistic Indexing of Pictures by a Statistical
1255 Modeling Approach. IEEE TPAMI 2003;25:1075–88.
- 1256 [50] CCSL/IME-USP . The complete data set used within the paper "a study
1257 of the relationships between source code metrics and attractiveness in free
1258 software projects". 2013. <http://ccsl.ime.usp.br/mangue/data>.
1259 Last access: 15/07/2013.
- 1260 [51] Lanza M, Marinescu R. Object-Oriented Metrics in Practice. New York:
1261 Springer; 2006.
- 1262 [52] Shewchuk JR. Delaunay refinement algorithms for triangular mesh gen-
1263 eration. Computational Geometry: Theory and Applications 2002;22(1–
1264 3):21–74.
- 1265 [53] Mount D, Arya S. Approximate nearest neighbor search software. 2011.
1266 www.cs.umd.edu/~mount/ANN. Last access: 15/07/2013.
- 1267 [54] Cao TT, Tang K, Mohamed A, Tan TS. Parallel banding algorithm to
1268 compute exact distance transform with the GPU. In: Proc. ACM I3D.
1269 2010, p. 83–90.
- 1270 [55] Turkay C, Filzmoser P, Hauser H. Brushing dimensions-a dual
1271 visual analysis model for high-dimensional data. IEEE TVCG
1272 2011;17(12):2591–9.