

Out-of-Core Dimensionality Reduction for Large Data via Out-of-Sample Extensions

Luca Reichmann*
University of Stuttgart

David Hägele†
University of Stuttgart

Daniel Weiskopf‡
University of Stuttgart

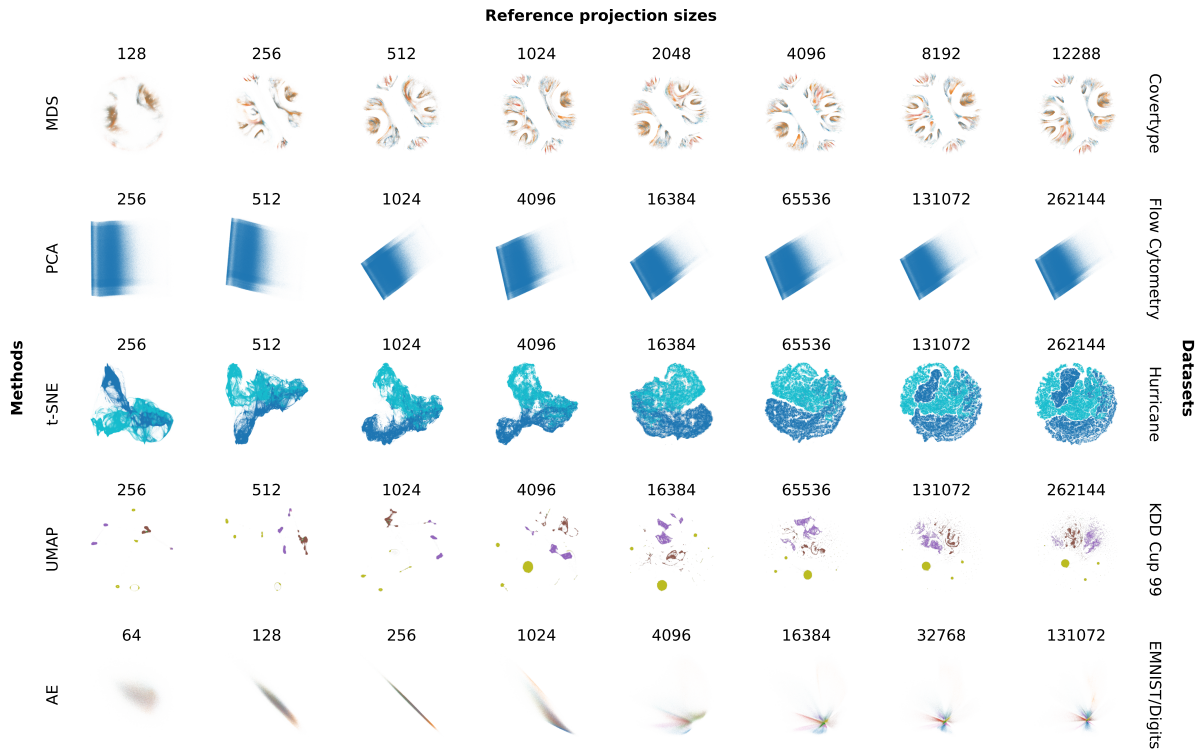


Figure 1: Visualizations of data sets with up to 50 million data points using increasing reference set sizes. We show the results for popular dimensionality reduction techniques: MDS, PCA, t-SNE, UMAP, and autoencoder. The number above each plot shows the reference set size used for creating the initial reference projection.

ABSTRACT

Dimensionality reduction (DR) is a well-established approach for the visualization of high-dimensional data sets. While DR methods are often applied to typical DR benchmark data sets in the literature, they might suffer from high runtime complexity and memory requirements, making them unsuitable for large data visualization especially in environments outside of high-performance computing. To perform DR on large data sets, we propose the use of out-of-sample extensions. Such extensions allow inserting new data into existing projections, which we leverage to iteratively project data into a reference projection that consists only of a small manageable subset. This process makes it possible to perform DR out-of-core on large data, which would otherwise not be possible due to memory and runtime limitations. For metric multidimensional scaling (MDS), we contribute an implementation with out-of-sample projection capability since typical software libraries do not support it.

*e-mail: st169765@stud.uni-stuttgart.de

†e-mail: david.haegle@visus.uni-stuttgart.de

‡e-mail: daniel.weiskopf@visus.uni-stuttgart.de

We provide an evaluation of the projection quality of five common DR algorithms (MDS, PCA, t-SNE, UMAP, and autoencoders) using quality metrics from the literature and analyze the trade-off between the size of the reference set and projection quality. The runtime behavior of the algorithms is also quantified with respect to reference set size, out-of-sample batch size, and dimensionality of the data sets. Furthermore, we compare the out-of-sample approach to other recently introduced DR methods, such as PaCMAP and TriMAP, which claim to handle larger data sets than traditional approaches. To showcase the usefulness of DR on this large scale, we contribute a use case where we analyze ensembles of streamlines amounting to one billion projected instances.

Index Terms: Dimensionality reduction, out-of-core, out-of-sample, evaluation.

1 INTRODUCTION

Dimensionality reduction (DR) is a popular visualization approach to gain insight into inherent structures of high-dimensional data in different domains. However, applying DR to large data sets consisting of millions of instances is challenging or even infeasible for some methods, such as multidimensional scaling (MDS) that scales quadratically in computational effort with the number of data points.

In this paper, we address the problem of DR for large data, such as high-dimensional biomedical data consisting of millions of individual observations. There are many DR techniques, with principal component analysis (PCA), uniform manifold approximation and projection (UMAP), and t-distributed stochastic neighbor embedding (t-SNE) being popular examples. Although there are some DR techniques that were developed to overcome the performance issues, such as TriMap [1] and pairwise controlled manifold approximation projection (PaCMAP) [47], they are still often limited by the available memory. Existing methods have been evaluated in the literature with data sets consisting of up to 20 million instances [4, 41], usually taking hours to compute and some requiring hardware exceeding consumer standards, e.g., 128+ GB of RAM.

In this paper, we use similarly sized data sets for comparability but also show that our approach is feasible for data sets comprised of 1 billion instances using around 240 GiB of memory. We resolve the memory problems by using out-of-sample (OOS) extensions of the DR methods. OOS extensions allow us to successively supplement an existing DR with new samples that are processed with respect to the existing projection. Depending on the DR technique, the implementation of OOS extensions can vary widely. While it is trivial to project new data points with PCA, the process is more complex with techniques like UMAP and t-SNE. Using OOS extensions usually not only lowers memory requirements but also reduces runtime because fewer data point comparisons become possible. However, this also leads to a loss in DR quality since less data is available to create the initial low-dimensional representation. Often, this trade-off is primarily influenced by the size of the already existing low-dimensional representation, which is then employed for processing subsequent OOS batches. We primarily use the term *projection* for the low-dimensional representation in our work, yet other terms, such as *embedding*, are also commonly used in the literature.

We provide a framework description enabling DR to be performed out-of-core by leveraging OOS extensions and making it viable to process data that does not fit into memory. Our main contribution is an extensive evaluation of OOS extensions for several DR methods using this framework: (metric) MDS [23], PCA, UMAP [28], t-SNE [44], and autoencoder [14], spanning DR techniques of different categories used for various purposes. The evaluation includes popular quantitative and qualitative measures of the projection quality of the DR methods. Figure 1 shows examples of the DR techniques applied to various data sets. We also conducted a runtime evaluation of the decrease in computational effort when using OOS extensions compared to the default use of DR techniques. In addition, we provide a description of a generalized framework for performing DR with OOS extensions to process large data out-of-core. Our approach can thus be applied to other DR techniques supporting OOS extensions.

Therefore, this paper aims to provide guidance when choosing the number of data points for the initial DR computation to balance the necessary computational effort and the resulting quality when employing OOS extensions with large data sets.

Furthermore, we provide supplemental material [37] containing additional figures, code to reproduce figures of this paper, and our own implementation of MDS with OOS extension.

2 RELATED WORK

We discuss related previous work that addresses the challenges of large data DR in general and specific approaches relying on and evaluating OOS extensions.

DR for Large Data There are several efforts to make DR applicable to large data. While we specifically propose using OOS extensions, others developed DR algorithms with a focus on performance or provided optimized implementations of existing DR algorithms, such as ports to the GPU. One of those techniques is Pair-

wise Controlled Manifold Approximation Projection (PaCMAP) by Wang et al. [47]. They applied PaCMAP to data sets such as KDD Cup '99, which contains over 4,000,000 data points. While conventional DR methods such as UMAP or t-SNE either failed because they ran out of memory or took longer than 24 hours to finish, PaCMAP could transform large data sets in a reasonable amount of time.

Another method is TriMap [1], which is also able to process large data sets with slightly higher runtimes. The authors of the TriMap publication also tested the method with data sets of up to 11 million data points. In both publications, UMAP and t-SNE tended to fail or exceed the 12-hour time limit at around a data set size of 1 million data points. While both papers used the UMAP implementation of McInnes et al. [29], TriMap relied on Multicore-TSNE [42] and PaCMAP on the Scikit-learn t-SNE implementation [34] for the comparison. LargeVis [41] was specifically designed to reduce the computational cost of t-SNE. Similar to the other two methods mentioned, LargeVis was evaluated using data sets with millions of data points. Zhu et al. [49] improved LargeVis, achieving higher performance with similar visualization quality.

One could also use parametric DR methods, which are intrinsically able to perform OOS projections by learning an explicit function to map points from the high- to low-dimensional space, to perform DR with large quantities of data. There have been proposals of parametric extensions for multiple non-parametric DR methods, such as t-SNE [43] and UMAP [39], which use neural nets to learn the function. Hinterreiter et al. [19] introduced a framework making parametric extensions generally available for DR techniques, as such extensions are not formulated for every DR method. While the parametric variants of t-SNE and UMAP provide high-quality projections, it is also noted that parametric versions of DR methods also come with more hyperparameters that have to be adjusted. In our paper, we use and evaluate OOS extensions instead of parametric extensions of DR methods to project OOS data points.

As implementation details also dramatically influence the feasibility of DR methods, there were different attempts to gain performance enhancements. A way of achieving this is GPU porting of the techniques, which exist for t-SNE [9, 30] and UMAP [32]. The authors could achieve a speedup of up to 700 times compared to the single-core Scikit-learn implementation and up to 100 times compared to the original multi-core UMAP implementation. We also use parallelization for the MDS algorithm by implementing its gradient descent scheme on the GPU. Typical quality metrics of DR algorithms, such as stress or trustworthiness, often suffer from poor algorithmic scalability [38] as they rely on pairwise comparisons. Therefore, Nolet et al. [32] also ported the trustworthiness computation to the GPU. The implementation works in batches to save memory. We continue this approach with the GPU-based implementations of the metrics described in Section 3.2.3. An alternative approach to dealing with the slow runtime of metrics is to only evaluate a subset of the whole DR computation. This approach was used by Kobak and Berens [22]. In contrast, we evaluate the complete DR computation with the GPU-based computation of the metrics to obtain more accurate and reliable insights.

Use and Evaluation of OOS Extensions OOS extensions are available for many DR methods that do not intrinsically support the transformation of new data points to the lower dimension with respect to the existing projection. While it is often mentioned in the literature that the use of OOS extensions might lead to lower quality [31, 33], a comprehensive study of the effects of using OOS extensions is lacking to the best of our knowledge.

However, there are proposals and corresponding evaluations of individual OOS extensions. One such extension was presented by Bengio et al. [6], where they introduced an OOS framework by learning the corresponding eigenfunctions. Similar to our evaluation, they evaluated the framework with different training set sizes.

Gisbrecht et al. [17] proposed kernel-based methods to obtain OOS extensions and demonstrated these for t-SNE. They presented three different methods and evaluated each of them with different training set sizes. Yet, their evaluation was limited, as they only tested three different training set sizes. Zhang et al. [48] extended the approach of kernel-based OOS extensions, as they identified weaknesses when processing outliers by introducing so-called bi-kernel t-SNE. They performed a more extensive evaluation of their approach, as they used multiple metrics to measure the quality of their approach and compare it to other state-of-the-art DR methods. In addition, they conducted a runtime evaluation.

We base our evaluation on these previous works. Many different quality metrics are available for quantifying projection quality, focusing on different DR aspects. Comprehensive studies of such measures are provided by Espadoto et al. [13], Gracia et al. [18], and Nonato and Aupetit [33]. These works comprise over 20 quality metrics categorized into local and global measures. We utilize a selection of them to quantify the projection quality of our proposed algorithm with metrics from both categories.

3 METHODOLOGY

Our methodology consists of two parts: a framework with a high-level algorithm to compute DR with OOS extensions for any DR method and our evaluation approach.

3.1 Computational Framework

We provide a generalized formulation of the OOS projection framework to perform DR on large data sets of many observations, allowing us to plug in any DR method that has an OOS extension. The pseudo code is listed under [Algorithm 1](#).

Algorithm 1 Projection with random subset reference

```

1: var  $n_{\text{ref}}$  ▷ reference size
2: var  $n_{\text{batch}}$  ▷ batch size
3: procedure PROJECT( $X, \Phi$ ) ▷  $X$ : data set,  $\Phi$ : DR method
4:    $X_a \leftarrow n_{\text{ref}}$  random points of  $X$ 
5:    $Y_a, \beta \leftarrow \Phi(X_a)$  ▷  $\beta$ : learned parameters
6:    $X_r \leftarrow X \setminus X_a$  ▷  $X_r$ : remaining data
7:   for  $i \in \{1 \dots \lceil \text{len}(X_r)/n_{\text{batch}} \rceil\}$  do ▷ project batches
8:      $X_{b(i)} \leftarrow i^{\text{th}}$  subset of  $X_r$ 
9:      $Y_{b(i)} \leftarrow \Phi_{\beta}(X_{b(i)})$  ▷ parameters  $\beta$  stay fixed
10:  end for
11:  return  $Y_a \cup Y_{b(1)} \cup \dots \cup Y_{b(n_{\text{batch}})}$ 
12: end procedure

```

A DR method in this framework is a procedure that takes a high-dimensional data set X as input and maps it to a low-dimensional representation Y while learning the set of parameters β for the mapping: $Y, \beta \leftarrow \Phi(X)$. The details and semantics of the parameters depend on the DR method. For PCA, β is the set of eigenvectors; for an autoencoder, it is the set of weights and biases; and for MDS, β is the discrete mapping $X \rightarrow Y$ itself. When the mapping has been learned, the parameters can be reused to map more data points, which we call the OOS projection: $Y' \leftarrow \Phi_{\beta}(X')$. This results in a mapping $X' \rightarrow Y'$ that is consistent with $X \rightarrow Y$.

The parameter β determines the amount of information available about the underlying high-dimensional space, providing a way of controlling the trade-off between the projection quality and time necessary to map OOS data. A set β of large cardinality may capture more information but come at higher computational costs for subsequent OOS projections. Note that the parameters β do not change when performing the OOS projection, even though it may be possible in general. Also, the points in the OOS batches cannot “see” each other, i.e., relationships between points of the OOS sets

are not taken into account. This is an important detail that provides deterministic behavior independent of batch size and order.

The algorithm starts by creating a small random subset X_a of the data points. The subset is then used to learn the mapping to low-dimensional space. The rest of the data X_r is then split up into batches $X_{b(i)}$ and mapped via OOS projection with the learned parameters. This means that the mapping $X_a \rightarrow Y_a$ is used as a reference for the projection. Hence, we call the initial subset the *reference set* similar to the *training set* in machine learning contexts. While splitting into batches is unnecessary from a theoretical point of view, practically, this is key to arriving at an out-of-core algorithm. The batch processing keeps the memory requirements of the individual DR methods low, and it allows for loading data in parts subsequently from external memory.

3.2 Evaluation

To evaluate the proposed algorithm, we test it with several DR methods that have an OOS extension. Our selection of DR methods includes some of the most popular algorithms and covers different data qualities that are preserved or optimized for, as outlined in [Table 1](#).

Table 1: DR methods used in the evaluation of the OOS framework.

Method	Optimizes for	Linear	β
PCA	reconstruction err.	yes	eigenvectors
MDS	global distances	no	$X_a \rightarrow Y_a$
t-SNE	local neighborhood	no	$X_a \rightarrow Y_a$, kNN
UMAP	local nb., global dst.	no	$X_a \rightarrow Y_a$, kNN
Autoencoder	reconstruction err.	no	weights, biases

In the evaluation, we use a systematic approach to gain insight into the quality, runtime, and memory usage of the OOS framework. A key factor is the size of the reference set X_a used to learn β . Thus, we report results for varying sizes of X_a . We also compare to the baseline without OOS extensions, i.e., where all data points are projected at once.

3.2.1 Out-of-Sample Extensions

The selected DR algorithms use different mechanisms to support OOS projection. PCA and the autoencoder learn a parametric mapping and give an explicit function to map the data, which serves as $\Phi_{\beta}(\cdot)$. For metric MDS, we use the process of performing stress minimization for a single point while keeping all others fixed, referred to as *single scaling* in the literature [3, 26], to perform the OOS projection of a single point. The gradient for an OOS point x' with respect to the points x_i of the reference set X_a and corresponding low-dimensional points y_i is given as $\tilde{\delta} = \sum_i (1 - d(x', x_i)/\|y' - y_i\|) \cdot (y' - y_i)$, with dissimilarity function $d(\cdot, \cdot)$. A similar mechanism has been proposed for t-SNE [7], which leverages informed initialization based on the OOS point’s k-nearest neighbors (kNN). The implementation we employed approximates kNN and uses interpolation-based t-SNE [27] to boost performance on large data sets. To the best of our knowledge, the OOS extension for UMAP is not described explicitly in the literature. However, similar approximate nearest neighbor mechanics can be found in UMAP’s source code [29].

3.2.2 Data Sets

We conduct the evaluation with large data sets of up to several millions of data points; see [Table 2](#). While these data sets are not particularly large in the sense of large data analysis and the visualization community in general, they are large in terms of what modern DR

Table 2: Data sets used in the evaluation of the OOS framework.

Name	Data points	Class count	Dim.	Source
EMNIST/Digits	280,000	10	784	[11]
Coverttype	581,012	7	54	[8]
Tornado	2,097,152	1	3	[12]
Flow Cytometry	3,176,162	1	23	[5]
KDD Cup '99	4,898,431	23	41	[40]
Higgs	11,000,000	2	28	[2]
Hurricane Isabel	50,000,000	2	13	1,2

methods can handle within hours of computation and limited memory. Although even much larger data sets would be possible with our OOS framework, we chose to stick to these data sets because they have already been covered in the DR literature. The metric-based evaluation also becomes challenging with respect to runtime with increasing data set size, which is another reason for us to stick to these sizes.

The evaluation data sets were chosen to cover a wide range of data sources, application areas, and data characteristics. The *EMNIST/Digits* data set contains 28×28 pixel images of handwritten character digits split into ten classes. *Coverttype* shows forest data in the US, where each data point represents a $30 \text{ m} \times 30 \text{ m}$ patch describing the cover type. The *Flow Cytometry* data set consists of fluorescence information acquired via the process of flow cytometry. *Tornado* is a uniformly sampled version of a synthetic 3D vector field representing a tornado. *KDD Cup '99* contains network data observations, including properties such as the protocol type and duration, which are categorized into normal and different types of intrusions. The data points of the *Higgs* data set represent two classes of signal processes: one that produces Higgs bosons and one that does not (background). *Hurricane Isabel* from the IEEE Visualization Contest 2004 consists of time-varying multidimensional data on a 3D uniform grid that represents simulated environmental variables as a hurricane travels across the land. We only use two time steps out of 48, amounting to 50 million data points.

3.2.3 Metrics

We use common quality metrics to measure the quality of OOS extension techniques. The metrics are often categorized as local or global ones [33, 18]. As the OOS framework is evaluated with different DR algorithms optimizing for different goals, we use metrics of both categories to provide comprehensive insights. Moreover, we use qualitative result inspection [20] to evaluate the types of quality of the projections that metrics cannot capture.

Global Metrics The *stress* metric [24] is defined as $\sqrt{\frac{\sum_{i \neq j} (d_{ij} - \|y_i - y_j\|)^2}{\sum_{i \neq j} (d_{ij}^2)}}$, where d_{ij} is the measure of dissimilarity of the high-dimensional points x_i and x_j that is compared to the Euclidean distance of corresponding low-dimensional points y_i and y_j . In our case, $d_{ij} = \|x_i - x_j\|$. Thus, stress measures the preservation of pairwise distances, where small stress indicates good preservation, with 0 being the ideal value. The metric is not bounded by a maximum value as it is scale-dependent. The stress calculation suffers from quadratic time and memory complexities due to its dependency on the pairwise distance matrix. Therefore, we perform the calculation in batches on the GPU.

Similar to Geng et al. [15], the *Pearson correlation coefficient* is used to calculate the correlation between d_i and \hat{d}_i , which represents the i -th distance entry of the flattened low- and high-dimensional distance vectors. A higher value indicates a higher correlation,

¹vis.computer.org/vis2004contest/

²vets.ucar.edu/vg/isabeldata/

with 1 being the best possible and 0 being the worst possible value. We chose the Pearson correlation coefficient instead of the Spearman rank correlation coefficient due to performance considerations when evaluating large projections. As the stress metric is scale-dependent and the Pearson correlation coefficient is not, we can compare how far clusters spread out in the projection with the stress metric. With the Pearson correlation coefficient, it is then possible to compare the projections independently of how spread out the projection is.

Local Metrics The following local metrics provide a parameter k , determining how many neighbors are taken into account when computing the metric value. As we apply the quality measures to larger data sets than in most other DR publications, we chose to consider larger neighborhoods (i.e., $k = 100$) than most other publications and metric implementations. This results in measuring a neighborhood containing 0.005% of all points of a data set with 2 million points, separating it from global measures.

The *KNN precision* [22] metric measures the accordance between the neighborhoods of a point in the low- and high-dimensional space. It is defined as $\frac{\sum_i \frac{1}{k} |H_{k_i} \cap L_{k_i}|}{n}$, where k denotes the number of neighbors considered, n the number of data points, H_{k_i} the indices of the k -nearest neighbors of the high-dimensional point x_i , and L_{k_i} the indices for the corresponding low-dimensional point y_i . The metric is again bounded to $[0, 1]$, where higher values indicate an increased neighborhood overlap between the low- and high-dimensional representation.

The *trustworthiness* metric [46] is defined as $1 - \frac{2}{nk(2n-3k-1)} \sum_{i=1}^n \sum_{j \in U_k(i)} (r(i, j) - k)$. The set $U_k(i)$ contains the low-dimensional elements that are in the set of k closest neighbors of y_i but not in the set of the closest neighbors of the same point in the higher-dimensional space. $r(i, j)$ yields the index of the low-dimensional point y_j in the nearest-neighbor ordering of the neighbors of the point y_i . As a result, the metric measures how well the local neighborhoods are preserved, where higher values indicate better local neighborhood preservation. Similar to the stress metric, the trustworthiness metric is calculated in batches on the GPU using the cuML library [36].

4 COMPUTATIONAL COMPLEXITIES

The DR techniques evaluated have different runtime and memory complexities, which could possibly be improved using the OOS extension framework described above. In this section, we discuss asymptotic complexity as a measure of the algorithmic scalability of the visualization techniques [38]. Regarding scalability with the number of data points n , MDS is in $\mathcal{O}(n^2)$ for both time and memory, due to its dependency on pairwise distances [16]. Similar to MDS, a time and memory complexity of $\mathcal{O}(n^2)$ is reported for t-SNE [44]. According to McInnes et al. [28], UMAP has an approximate runtime complexity of $\mathcal{O}(n^{1.14})$. With respect to n , PCA scales linearly for the computation of the covariance matrix (when dominated by n), and the eigendecomposition is independent of n . Its memory complexity is only dependent on the number of dimensions. The time for training an autoencoder scales linearly with n [45]. The memory complexity depends on the encoder’s architecture, which is also independent of the number of data points.

Our framework is designed to generically bring the computational effort down, making it possible to run demanding DR methods and large data sets on consumer-grade computers. The algorithmic time complexity of our framework can be described as

$$\mathcal{O}(\Phi^{\text{RT}}(n_{\text{ref}})) + \mathcal{O}(\Phi_{\beta}^{\text{RT}}(n_{\text{batch}}) \cdot \text{batch_count}), \quad (1)$$

where Φ^{RT} is the function that calculates the DR method’s runtime and Φ_{β}^{RT} represents the runtime of the method’s OOS projection operation. This is because the reference projection has to be created

first, and then each batch is projected sequentially using the information from the reference projection. The time complexities are dependent on the specific DR method used, yet the time complexity of each batch projection is equal due to constant batch sizes. This also applies to the memory requirements, as we do not need to compute on the whole data set at the same time. Consequently, the memory requirements are constantly $\mathcal{O}(\Phi^M(n_{\text{ref}})) + \mathcal{O}(\Phi_{\beta}^M(n_{\text{batch}}))$ for each individual batch projection, where Φ^M and Φ_{β}^M are the functions to calculate the memory requirement of the respective DR method. Please note that complexities for Φ and Φ_{β} are possibly unequal. For PCA and the autoencoder, the OOS projection via Φ_{β} is approximately in $\mathcal{O}(n_{\text{batch}})$, since the data only needs to be transformed via their explicit projection function. In MDS, the points in the OOS set are only compared to the reference set but not among themselves, resulting in a time complexity of $\mathcal{O}(n_{\text{ref}} \cdot n_{\text{batch}})$. Similarly, the t-SNE and UMAP implementations ignore relations within the OOS set.

5 EXPERIMENTS

In the following, we demonstrate how the framework performs with the described data sets and various DR techniques. We highlight both qualitative and quantitative results by inspecting the projections and including the results of the metric values. Furthermore, we show how the runtime changes with increasing reference set sizes and how the techniques compare to each other.

5.1 Setup and Implementation

Our benchmarks were run on a workstation equipped with an AMD Ryzen Threadripper PRO 3995WX with 64 CPU cores, an NVIDIA RTX A6000 GPU with 48 GiB of VRAM, and 251 GiB of RAM. This is not a consumer-grade computer, and the OOS framework requires much less computing resources. However, a potent machine was required for large projections with regular DR methods and the computation of quality metrics. The software environment used for running and evaluating the DR techniques was implemented in Python with Numba [25] compiled functions for performance-critical sections. We used the UMAP implementation of the umap-learn library [29], the openTSNE library [35] for t-SNE, scikit-learn [34] for PCA, and keras [10] for the autoencoder.

Since an MDS implementation with an OOS extension is missing in popular DR libraries, we provide our own implementation that can be applied in an out-of-core fashion with GPU-accelerated gradient descent, which can be found in the supplemental material [37]. We used standard hyperparameter settings if provided by the corresponding DR libraries. We set the number of iterations to 350 for UMAP, which is the maximum it defaults to for large data sets. For MDS, we chose a fixed number of 500 iterations with a step size of 10^{-4} to strike a balance between good quality results and reasonable runtime. The autoencoder models used in the evaluation were inspired by the models from the literature [14], where ReLU activation functions are used in the encoder- and decoder layers, linear and a sigmoid activation function in the last layers.

5.2 Reference Set Size Quality Trade-off

In this subsection, we evaluate the projection quality with respect to different reference set sizes. Figure 1 shows the impact of reference set size for combinations of different DR methods and data sets. The number in each cell of the figure refers to the reference set size used to create the corresponding projection. To restrict the runtime of MDS to a reasonable range, we decided to limit the set size to $2^{13} + 2^{12}$ data points. To substantiate the qualitative observations, we include metric values in Figure 2. Since the metric computations become increasingly slow with large data sets, we present the metric values for the smaller data sets of our collection (Table 2):

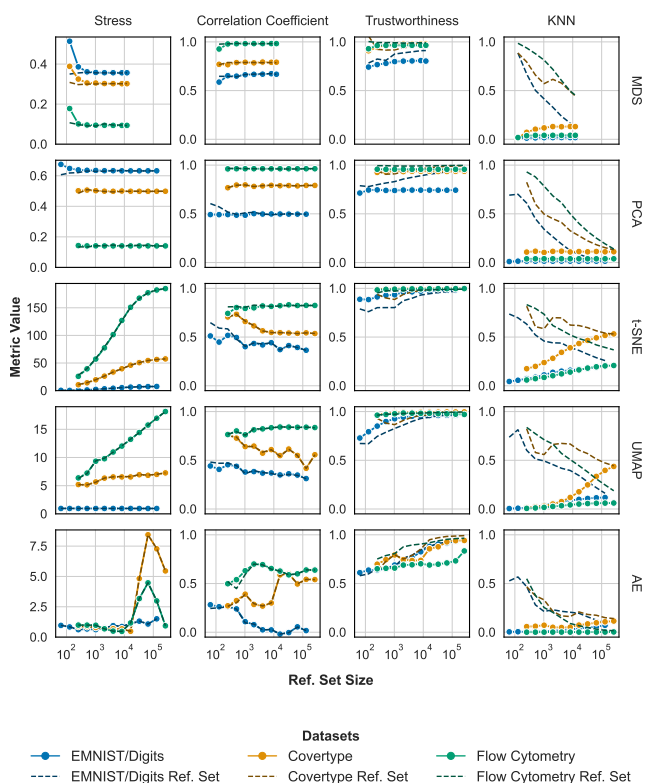


Figure 2: Metric values for increasing reference set size for the EMNIST, Covertypes, and Flow Cytometry data sets. The x-axis uses a logarithmic scale for the reference set sizes, and the y-axis displays the corresponding metric values.

EMNIST, Covertypes, and Flow Cytometry. We include metric values for the reference and whole projections, providing insights into OOS projection performance.

In Figure 1, it can be seen that the projections, in general, become more consistent with increasing reference set size. MDS appears consistent pretty early on. Even with 128 data points, the global overall point arrangement is already visible. This is also observable in the stress metric results, where the values stabilize at a size of 1,024. Both the stress and correlation coefficient metrics show a larger discrepancy between the projection of the reference set and the whole projection when using small reference set sizes. The local measures indicate that the projections also tend to have higher accordance between the low- and high-dimensional neighborhoods with larger reference set sizes. However, local measures are less meaningful for MDS than global measures.

For PCA, the structure of the Flow Cytometry data set projection is already well visible for the smallest reference set. However, different rotations of the data set are noticeable, showing that the principal vectors still vary for sizes up to 16,384. The smaller changes in the projections of PCA, when compared to the other methods, are also reflected in the corresponding metric values, where almost no changes are visible. Nonetheless, with very small reference sets, there are some minor weaknesses, as can be seen in the stress and trustworthiness values of the EMNIST data set (Figure 2).

The projections created with t-SNE and UMAP show the most variation throughout the different reference set sizes, which we think is an effect of overfitting on small reference sets. Due to their objective of preserving local neighborhoods, they typically also show higher values of trustworthiness and KNN accuracy with

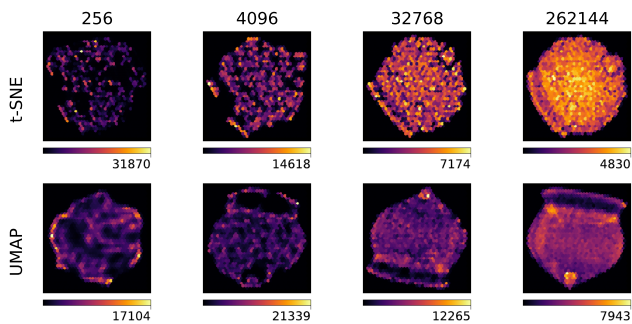


Figure 3: Heat maps of the projections of the Flow Cytometry data set with t-SNE and UMAP. The color bar label shows the number of data points in the most populated area.

small reference sets compared to the other methods (see Figure 2). The KNN metric increases more drastically with larger reference set sizes, which is especially noticeable with the Covertype data set. The plot also shows that the KNN metric for the reference set and the whole projection converge, which is expected since the reference set size approaches the actual data set size. The correlation coefficient measure shows larger differences between the reference set sizes, getting worse for the EMNIST and Covertype data sets with both techniques. These more drastic changes are also reflected in the projections itself, which continuously change with increasing reference set size. Both techniques tend to form more distinct, compressed clusters with smaller reference sets. This can be seen in the purple clusters of the KDD Cup '99 data set, both clusters of the Hurricane data set, and in Figure 3, where we include heat maps of the Flow Cytometry data set created with UMAP and t-SNE. It can be seen that with both methods, the point distribution in the projections becomes increasingly spread out. This is also apparent in the heat maps, where the densest regions drop from over 30,000 samples per tile to under 5,000 and from over 17,000 to under 8,000, respectively. The heat maps also indicate that t-SNE seems to generate more compressed clusters than UMAP, which is likely the consequence of UMAP optimizing more for global relationships than t-SNE. While it may seem that larger reference sets lead to better global distance relationship preservation with both techniques, the global stress and correlation coefficient measures contradict this observation.

The results of the autoencoder are more inconsistent than those of the other techniques. With the smallest reference set size, the data points in the projection of the EMNIST data set are arranged in a point cloud, then become line-shaped, and starting from 4,096 data points, the expected distinct separation of classes becomes visible. This behavior is also reflected in the metric values of the other data sets. While the other techniques tend to show a slow and steady improvement or decline, there are larger fluctuations in the autoencoder metrics. This leads to the conclusion that the training process of the autoencoder is influenced more significantly by the size of the reference set than those of the other DR methods. While the global quality metrics might imply that the projections become worse with larger reference sets, both local metrics show better results.

5.3 Trade-off Between Reference Set Size and Runtime

We measure the time consumed by the DR methods using different reference set sizes. We use a synthetic data set consisting of four isotropic Gaussian blobs with 1,000,000 instances in total, which is varied in the number of dimensions to capture the relation of dimensionality to runtime. In our experiment, we use dimensionalities of 64, 512, and 2,048, which are chosen based on the dimensionality ranges of typical data sets used for DR. We use reference set sizes

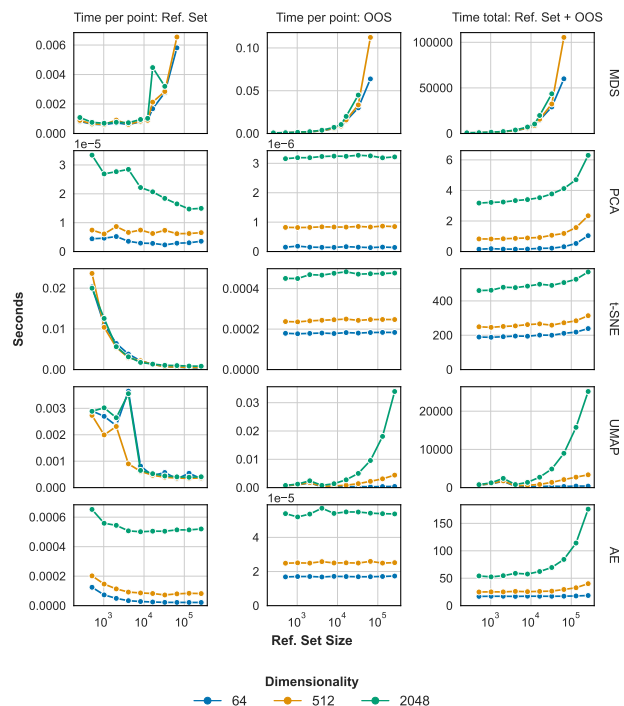


Figure 4: Runtime (in seconds) needed for the projections for varying reference set sizes. The x-axis uses a logarithmic scale for the reference set sizes. The line colors refer to the dimensionality of the projected data set.

similar to those in the quality trade-off experiments. The OOS projections were performed in batches of 100,000 data points, except for MDS, where we used 1,000 points per batch due to GPU memory limitations. We include the results of the runtime evaluation in Figure 4 and Figure 5. The left column of Figure 4 shows the time necessary to compute the reference projection, and the middle column the time needed for the projection of the OOS batches. The right column shows the total time necessary for the whole projection. We divide the times by the number of points projected in the left and middle columns to get a comparable unit of time per point.

In the runtime graph of MDS, it can be seen that the time to project subsequent batches increases up to over 0.1 seconds per point. The runtime of the reference set projection shows a similar effect, where the time per point rises from under 2 milliseconds to over 6. The projection runtime increases with dimensionality. Yet, the influence is negligible, as the dimensionality only affects the computation of the high-dimensional distance matrix, which is pre-computed once before the optimization process of each OOS batch.

The per-point projection runtime of OOS batches with PCA does not significantly increase with larger reference sets, which results from only needing to apply a single matrix multiplication on the OOS batches. We can see that the total runtime slightly increases to over 6 seconds with the largest reference set and the 2048-dimensional data set, which can be attributed to the overhead of processing a larger reference set.

The openTSNE library varies hyperparameter settings depending on the input data, influencing the runtime behavior of the algorithm. Thus, we decided to always use the approximate neighbor search and employ interpolation-based t-SNE in this experiment to obtain more consistent results. The runtime behavior shows that with increasing reference set sizes, the mean time to project an individual point steadily and quickly decreases. The OOS batch pro-

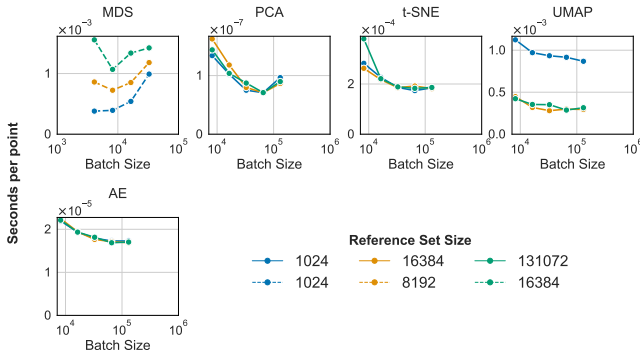


Figure 5: The time necessary to project a point with batches of varying size. The time is measured with the three different reference set sizes per technique (see the legend).

jection runtime is not influenced by the size of the reference set, as the runtimes per point are almost constant with all data set dimensionalities. Despite that, the overall runtimes increase with larger reference sets.

The runtime behavior of UMAP in the reference set projection process is hardly influenced by the dimensionality. In contrast, the times necessary for the OOS projections are affected more substantially by the dimensionality of the synthetic data set. The per-point runtimes for the 64-dimensional data set are low for all reference set sizes, while they increase to over 0.03 seconds for the 2048-dimensional data set with the largest reference set. The total runtime of the projection is, therefore, also sensitive to the size of the reference set, showing an increase from under 1 second to over 20,000 seconds for the 2048-dimensional data set.

The autoencoder runtimes of the OOS batch projection are comparable to those of PCA and t-SNE. While the autoencoder itself is slower than PCA and faster than t-SNE, the overall runtime progression is similar, as the times do not increase with larger reference sets. This is expected since transforming new samples into the low-dimensional space is independent of the size of the reference set once the network is trained. Yet, the training process runtime is dependent on the size and dimensionality of the reference. We can see that the training runtime per point decreases marginally in the first five projections before stabilizing.

In Figure 5, we show the time necessary to perform DR with an individual batch of varying sizes. We use the same synthetic data set with 32 dimensions. It can be seen that the time to project one point decreases with the size of the batch with each technique, except for PCA in the last step and MDS with increasing times for larger batches. The runtime results of MDS and UMAP indicate that they are more sensitive to the reference set size when it comes to the projection of subsequent OOS batches. This is in accordance with the results shown in Figure 4, where the per-point times of the projection are almost constant with PCA, t-SNE, and the autoencoder with increasing reference set sizes. While the batch projection times increase with larger reference set sizes using MDS, UMAP generally took the longest with the small reference sets.

6 COMPARISON TO LARGE-SCALE DR METHODS

In this section, we compare the OOS framework with UMAP to the large-scale DR methods TriMAP [1] and PaCMAP [47], both with respect to quality measures and runtime, using the KDD Cup '99 and Tornado data sets. In Table 3, we include the metric results of the comparison, in Table 4 the corresponding runtimes, and in Figure 6 the projections of the Tornado and KDD Cup '99 data sets. The figure shows that the projections of UMAP and PaCMAP

Table 3: Results of comparing UMAP and the OOS framework with other state-of-the-art DR methods for large data with the KNN (local) and correlation coefficient (global) quality measures. The best two results per data set and metric are marked in bold. Additionally, the results of running UMAP with 1,000 iterations are included in parentheses.

Method (Ref. Set)	Tornado		KDD Cup '99	
	KNN	Corr. Coeff.	KNN	Corr. Coeff.
UMAP (4,096)	0.036 (0.27)	0.165 (0.14)	0.017 (0.02)	0.508 (0.54)
UMAP (65,536)	0.275 (0.48)	0.097 (0.13)	0.094 (0.11)	0.579 (0.40)
UMAP (262,144)	0.363 (0.51)	0.074 (0.11)	0.146 (0.16)	0.523 (0.17)
UMAP	0.370 (0.35)	0.386 (0.007)	0.153 (0.18)	0.374 (0.35)
TriMAP	0.271	0.098	0.001	0.696
PaCMAP	0.522	0.879	0.207	0.768

Table 4: Results of comparing the runtimes of UMAP and the OOS framework with other state-of-the-art DR methods for large data, which includes both the projection of the reference set and the OOS batches, if the OOS framework is used.

Method (Ref. Set)	Runtime [in sec.]		
	Tornado	KDD Cup '99	Higgs
UMAP (4,096)	508.82	1,473.15	3,026.46
UMAP (65,536)	514.76	1,829.00	2,693.36
UMAP (262,144)	732.46	3,591.94	2,925.83
UMAP (none)	1,344.07	13,046.95	5,194.33
TriMAP (none)	13,037.25	33,155.39	76,889.32
PaCMAP (none)	9,199.30	20,661.68	65,754.48

are more similar to each other than the projections created with TriMAP, especially visible in the projection of the Tornado data set. Wang et al. [47] also noticed a similar effect with the KDD Cup '99 data set, stating that this is a result of TriMAP primarily preserving global structure. We can also see a difference in the projection of the Tornado data set between the UMAP projection with OOS projection and pure UMAP, where the latter is missing a hole in the center. A similar result is obtained for the KDD Cup '99 data set, where the purple points in the projection without an OOS set are arranged around the three yellow dots. There still seem to be changes with UMAP when using larger reference sets than 262,144 with the given data sets. It is also conceivable that the global structure of the projections of UMAP with OOS is more similar to the projections of PaCMAP. The metric values show that PaCMAP tends to perform the best both with regard to global and local quality measures. UMAP, without using the OOS framework, has the second-best values for the local KNN measure with 0.37 and 0.15. It can also be seen that in all cases, the KNN metric values increase when the reference set size is increased, where the best values are scored when not using the OOS framework. However, the correlation coefficient metric does not systematically increase when using larger reference sets, as already observed in previous sections.

TriMAP performs substantially worse with respect to the KNN measure, especially with the KDD Cup '99 data set, where a value of 0.001 is achieved. With the Tornado data set, the KNN metric value is higher at 0.271, yet it is still lower than the values reached with UMAP and PaCMAP. Amid and Warmuth [1] also observed lower values of local neighborhood accuracy with TriMAP compared to UMAP and t-SNE with several state-of-the-art DR benchmark data sets. Similar to their findings, TriMAP also performs well with respect to the global correlation coefficient quality measure in our experiment with the KDD Cup '99 data set, where a value of 0.696 is achieved compared to 0.768 with PaCMAP, which

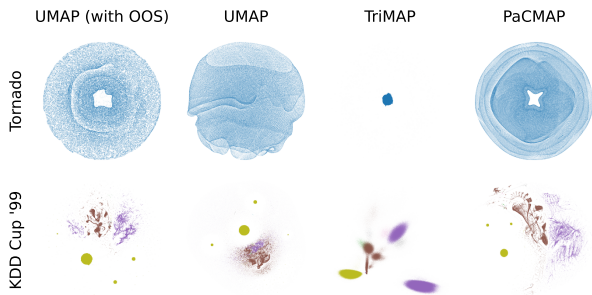


Figure 6: Projections of the Tornado and KDD Cup '99 data sets in the first and second row, respectively. The projections in the first column were created with UMAP and a reference set size of 262,144.

is the highest value with this configuration. However, TriMAP yields lower correlation coefficient values with the Tornado data set, where it performs slightly worse than UMAP and significantly worse than PaCMAP (i.e., 0.098 vs. 0.879), which correlates with the unproportional positioning of outlier data points visible in the projection. The comparison between the three DR methods shows that PaCMAP generally performs the best when compared to TriMAP and UMAP. TriMAP tends to be inconsistent with its results, while UMAP performs as assumed, where steady improvements to projection quality can be expected. We can also see that the quality of the projection usually increases for larger reference sets with UMAP, especially with local quality measures. This shows that with our OOS approach, the reference set size should be chosen as large as possible while considering memory and runtime requirements, as the projection is usually of higher quality.

7 USE CASE

In this section, we apply the OOS framework to an example of a large-scale data set. Our test data set consists of 1 billion streamlines that we generated on the *fluid simulation ensemble for machine learning* [21]. From the 8,000 flow fields, we used the first 1,000. For each flow field, we generated 1M streamlines using fourth-order Runge-Kutta (RK4) integration with step sizes $\{0.1, 0.2, 0.4, 0.8\}$ and 32 integration steps each. Our intention was to use DR to get insights into flow behavior and to be able to characterize the flow fields. To make the streamlines comparable, we translated their origins to zero and rotated them to consistently start in the positive x direction before applying DR. To build a reference set, 100,000 streamlines were generated at random locations from a small portion of 50 flow fields of the ensemble. The reference projection was then obtained using t-SNE, where we selected the perplexity parameter so that many clear clusters were forming. We used the Euclidean distance between the vectors consisting of the streamlines' x - and y -coordinates $(x_1, y_1, \dots, x_{32}, y_{32})$ as the similarity measure between two streamlines. Since the 1 billion streamlines that we aimed to OOS project would take approximately 240 GiB of memory, we chose to generate them on-the-fly and directly pass them in per-flow-field batches to t-SNE for projection. Thus, we did not need to sacrifice any disk space and could also simulate an in-situ application of the OOS framework.

Figure 7 shows a heat map of all projected streamlines with log-scaled color mapping. The projection exhibits many little clusters that correspond to similarly shaped streamlines. A few of them are illustrated below the heat map. From examining these patterns, we get an idea of the projection's semantics, where very elongated trajectories can be found on the left, and trajectories with little extent where flow is almost stationary are located on the right. In between, we can find medium-length oscillating streamlines, where clusters

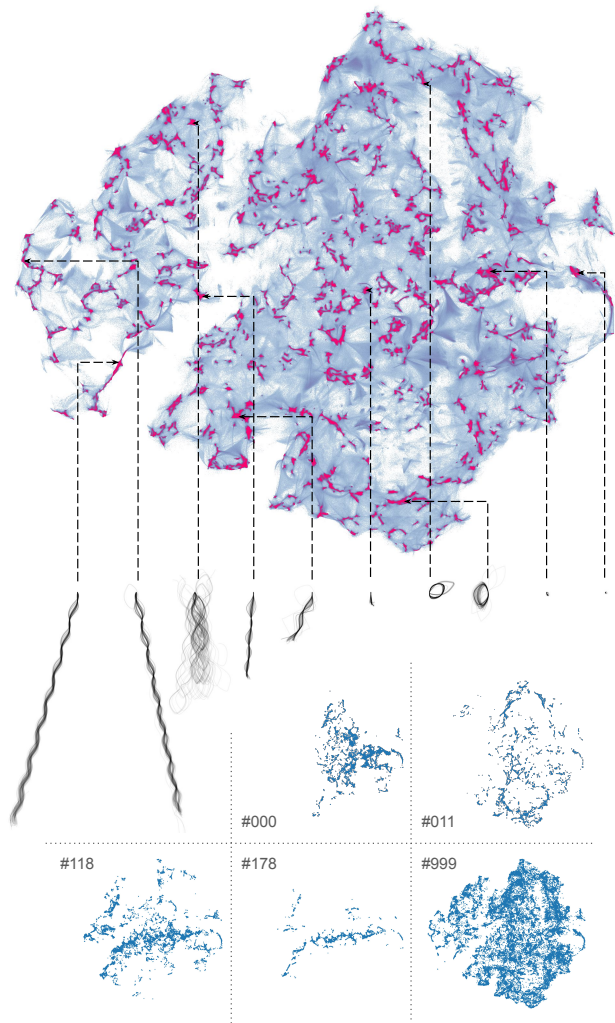


Figure 7: (Top) Heat map of all t-SNE-projected streamlines, where blue indicates low density and red high density. (Middle) For different clusters of the projection, 100 of the corresponding streamlines. (Bottom) Scatterplots of subsets of the projection that correspond to different flow fields of the ensemble, showing different coverage of the clusters.

differ in orientation and amplitude of the lines, and also elliptical patterns with clockwise and counterclockwise winding. At the bottom of the figure, there are scatterplots of a selection of flow fields that only contain their respective projected streamlines. What can be observed is that they differ in coverage of the space, i.e., different flow fields exhibit different kinds of streamlines and allow us to characterize them based on the areas of the projection that they cover. The projection could be effectively used to explore the flow ensemble when integrated into a visual analytics system, giving an overview from a streamline perspective and serving as a starting point for analysis.

The OOS projection of a batch of 1M streamlines took, on average, 11.8 minutes, of which around 7.4 minutes were used for actual optimization, and the rest of the time was spent on nearest neighbor search and computation of affinities between reference and OOS set. The whole process of reading each flow field from a file, generating the streamlines, and performing the OOS projection needed 200 hours of computation time, but since the batches are easily parallelizable, we ran four concurrently, resulting in 2 days and 2 hours

to obtain the whole projection. However, upon examining the logs, we observed that improvements in the Kullback-Leibler divergence were stalling early during the optimizations, meaning that the quality gain of using many iterations was negligible, and we could have used only 400 instead of 1000 iterations, which would have saved us 1/3 of the time.

8 DISCUSSION

In this section, the results of the evaluation of the OOS approach are discussed, along with limitations and possible future work.

8.1 Discussion of Results

We found that the OOS extensions deliver varying results depending on the DR technique. While techniques like MDS or UMAP are sensitive to the size of the reference set both in runtime and quality, other methods like PCA are not. Despite these differences, we observed a general increase in quality with respect to metrics measuring local neighborhood preservation. Depending on the technique, the results of global quality measures were mixed. While improvements with DR methods optimizing for preserving global relationships can be seen, both the metric values of t-SNE and UMAP dropped for larger reference sets. This does not necessarily indicate worse quality as the techniques work better toward their goal to visibly separate clusters, as can be seen in the t-SNE projection of the Hurricane data set. We also noticed that especially the projections of UMAP still change substantially with some data sets when the reference set size is larger than 262,144, which has to be taken into account when using the OOS extension.

We observed improvements in runtime with all techniques when using smaller reference sets. Again, the results differ depending on which DR algorithm (i.e., what is learned with β) and data set are used. Parametric methods, such as the autoencoder or PCA, can typically perform OOS projections quickly, as the learned function can just be applied to the OOS batches again. Substantial improvements in runtime were achieved when using the OOS approach with the implementations of MDS and UMAP, especially with very high-dimensional data sets.

Another benefit of using OOS extensions is that memory limitations are no longer an issue, as only the reference set and the individual OOS batch have to be stored in memory simultaneously. While we did not conduct a systematic study on the memory requirements of our approach, we were able to project the KDD Cup '99 data set with 4,898,431 data points with 41 dimensions using our GPU implementation of MDS on an NVIDIA RTX A6000 GPU, which would not easily be possible without an OOS extension due to memory constraints. Failures with this data set due to memory restrictions with techniques like UMAP and LargeVis were reported by Wang et al. [47], which we also were able to circumvent. While PaCMAP typically performs the best, UMAP with OOS typically yields better, more reliable results than TriMAP. Noteworthy is also that the projections of both data sets using the UMAP OOS extension are arguably more similar in global structure to those of PaCMAP than UMAP without OOS extension.

We have also seen that speed is not an issue with UMAP with the hardware used in our state-of-the-art comparison, as the UMAP implementation, in fact, performed better than both TriMAP and PaCMAP. Yet, the runtime was even further reduced using the OOS approach. Based on our evaluation and experience from the use case, we recommend that the reference projection should be tested before applying the OOS projection to the whole data set. By testing, we mean examining the emerging structures and checking if the analysis you aim to do is possible on the small reference subset. We recommend trying to find out which data characteristics are represented by the different areas of the projection to get an idea of what insights you can expect to find in the overall projection

later, so you only spend time and compute resources from promising starting points.

8.2 Limitations

The evaluation of the proposed OOS framework is quite time-consuming, mainly due to the computation of metrics but also due to the repeated experiments for different reference set sizes, so we needed to cut corners. For instance, we did not tune hyperparameters of the DR methods for the specific data sets, e.g., we did not test different perplexities for t-SNE or the number of neighbors for UMAP, but we used default parameters most of the time. Also, for the autoencoder, we used off-the-shelf architectures suggested in the literature. To keep runtimes foreseeable, we used a fixed number of iterations and did not employ convergence criteria for terminating optimizations. Thus, the results we are showing are not benchmark results to compare methods against each other because there was no mechanism in place to achieve the best quality with each DR method. Therefore, the produced projections and timings are not on a competitive level.

Regarding the limitations of the approach, we think that the reference set's projection quality determines the overall quality that is achievable. We did not evaluate this explicitly, but Figure 2 supports this hypothesis. When comparing the metric of the complete projection with the reference projection, it can be seen that they are on par or converging. In our experiments, the reference sets were selected randomly, so it is possible that certain patterns of the data may not be included in the reference and, therefore, cannot be uncovered for the OOS data. A more informed selection of the reference set could lead to improved results if the reference, therefore, becomes more representative of the data. Another limitation by design is that relationships between OOS points are ignored and that the underlying projection model (β parameter) is not updated with OOS points.

8.3 Future Work

There are several other aspects that we want to examine in detail in the future. To get a comprehensive evaluation of the OOS framework, we need to expand our selection of data sets to cover an even wider range of domains, sizes, and dimensions. An extensive benchmark suite for DR methods with OOS extension would be a valuable contribution for the community to get detailed insights into achievable quality and performance for various data sets. We also seek opportunities to apply the OOS framework as part of a visual analytics system and perform a design study to help solve real-world large data problems.

9 CONCLUSION

We formulated a generic algorithm for performing out-of-core DR by leveraging OOS extensions and showed that the proposed approach makes it feasible to project large data sets. The runtime performance and resulting quality of selected DR methods were examined for varying reference set sizes to fathom the trade-off between speed and quality. We found that some methods were more sensitive to reference set size and needed larger sets to produce consistent projections, while others produced similar projections with small reference sets already. With a use case, we showcased the feasibility of the OOS framework to produce usable results for a dataset as big as having one billion instances.

ACKNOWLEDGMENTS

This work was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation)—Project ID 251654672—TRR 161 (Project A01).

REFERENCES

- [1] E. Amid and M. K. Warmuth. TriMap: Large-scale dimensionality reduction using triplets. *arXiv preprint arXiv:1910.00204*, 2019. 2, 7
- [2] P. Baldi, P. Sadowski, and D. Whiteson. Searching for exotic particles in high-energy physics with deep learning. *Nature Communications*, 5:4308, 2014. doi: 10.1038/ncomms5308 4
- [3] W. Basalaj. Incremental multidimensional scaling method for database visualization. In R. F. Erbacher, P. C. Chen, and C. M. Wittenbrink, eds., *Visual Data Exploration and Analysis VI*, vol. 3643, pp. 149–158. International Society for Optics and Photonics, SPIE, 1999. doi: 10.1117/12.342830 3
- [4] A. C. Belkina, C. O. Ciccolella, R. Anno, R. Halpert, J. Spidlen, and J. E. Snyder-Cappione. Automated optimized parameters for T-distributed stochastic neighbor embedding improve visualization and analysis of large datasets. *Nature Communications*, 10(1):5415, Nov 2019. doi: 10.1038/s41467-019-13055-y 2
- [5] A. C. Belkina and J. E. Snyder-Cappione. OMIP-037: 16-color panel to measure inhibitory receptor signatures from multiple human immune cell subsets. *Cytometry Part A*, 91(2):175–179, 2017. doi: 10.1002/cyto.a.22983 4
- [6] Y. Bengio, J.-F. Paiement, P. Vincent, O. Delalleau, N. L. Roux, and M. Ouimet. Out-of-sample extensions for LLE, Isomap, MDS, Eigenmaps, and Spectral Clustering. In *Proceedings of the 16th International Conference on Neural Information Processing Systems, NIPS’03*, p. 177–184. MIT Press, Cambridge, MA, USA, 2003. 2
- [7] G. J. Berman, D. M. Choi, W. Bialek, and J. W. Shaevitz. Mapping the stereotyped behaviour of freely moving fruit flies. *Journal of The Royal Society Interface*, 11(99), 2014. doi: 10.1098/rsif.2014.0672 3
- [8] J. Blackard. Covertype. UCI Machine Learning Repository, 1998. doi: 10.24432/C50K5N. 4
- [9] D. M. Chan, R. Rao, F. Huang, and J. F. Canny. T-SNE-CUDA: GPU-accelerated t-SNE and its applications to modern data. In *2018 30th International Symposium on Computer Architecture and High Performance Computing (SBAC-PAD)*, pp. 330–338, Sep. 2018. doi: 10.1109/CAHPC.2018.8645912 2
- [10] F. Chollet et al. Keras. <https://keras.io>, 2015. 5
- [11] G. Cohen, S. Afshar, J. Tapson, and A. van Schaik. EMNIST: Extending MNIST to handwritten letters. In *2017 International Joint Conference on Neural Networks (IJCNN)*, pp. 2921–2926, May 2017. doi: 10.1109/IJCNN.2017.7966217 4
- [12] R. Crawfis. Tornado data set generator, 2003. 4
- [13] M. Espadoto, R. M. Martins, A. Kerren, N. S. T. Hirata, and A. C. Telea. Toward a quantitative survey of dimension reduction techniques. *IEEE Transactions on Visualization and Computer Graphics*, 27(3):2153–2173, March 2021. doi: 10.1109/TVCG.2019.2944182 3
- [14] Q. Fournier and D. Aloise. Empirical comparison between autoencoders and traditional dimensionality reduction methods. In *2019 IEEE Second International Conference on Artificial Intelligence and Knowledge Engineering (AIKE)*, pp. 211–214, June 2019. doi: 10.1109/AIKE.2019.00044 2, 5
- [15] X. Geng, D.-C. Zhan, and Z.-H. Zhou. Supervised nonlinear dimensionality reduction for visualization and classification. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 35(6):1098–1107, 2005. doi: 10.1109/TSMCB.2005.850151 4
- [16] A. Gisbrecht and B. Hammer. Data visualization by nonlinear dimensionality reduction. *WIREs Data Mining and Knowledge Discovery*, 5(2):51–73, 2015. doi: 10.1002/widm.1147 4
- [17] A. Gisbrecht, W. Lueks, B. Mokbel, B. Hammer, et al. Out-of-sample kernel extensions for nonparametric dimensionality reduction. In *ESANN 2012 Proceedings, 20th European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, pp. 531–536, 2012. 3
- [18] A. Gracia, S. González, V. Robles, and E. Menasalvas. A methodology to compare dimensionality reduction algorithms in terms of loss of quality. *Information Sciences*, 270:1–27, 2014. doi: 10.1016/j.ins.2014.02.068 3, 4
- [19] A. Hinterreiter, C. Humer, B. Kainz, and M. Streit. ParaDime: a framework for parametric dimensionality reduction. *Computer Graphics Forum*, 42(3):337–348, 2023. doi: 10.1111/cgf.14834 2
- [20] T. Isenberg, P. Isenberg, J. Chen, M. Sedlmair, and T. Möller. A systematic review on the practice of evaluating visualization. *IEEE Transactions on Visualization and Computer Graphics*, 19(12):2818–2827, 2013. doi: 10.1109/TVCG.2013.126 4
- [21] J. Jakob, M. Gross, and T. Günther. A fluid flow data set for machine learning and its application to neural flow map interpolation. *IEEE Transactions on Visualization and Computer Graphics*, 27(2):1279–1289, 2021. doi: 10.1109/TVCG.2020.3028947 8
- [22] D. Kobak and P. Berens. The art of using t-SNE for single-cell transcriptomics. *Nature Communications*, 10(1):5416, Nov 2019. doi: 10.1038/s41467-019-13056-x 2, 4
- [23] J. Kruskal and M. Wish. *Multidimensional Scaling*. SAGE Publications, Inc., 1978. doi: 10.4135/9781412985130 2
- [24] J. B. Kruskal. Nonmetric multidimensional scaling: A numerical method. *Psychometrika*, 29(2):115–129, Jun 1964. doi: 10.1007/BF02289694 4
- [25] S. K. Lam, A. Pitrou, and S. Seibert. Numba: a LLVM-based Python JIT compiler. In *Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC, LLVM ’15*. Association for Computing Machinery, New York, NY, USA, 2015. doi: 10.1145/2833157.2833162 5
- [26] D. Leon, A. Podgurski, and W. Dickinson. Visualizing similarity between program executions. In *16th IEEE International Symposium on Software Reliability Engineering (ISSRE’05)*, 2005. doi: 10.1109/ISSRE.2005.45 3
- [27] G. C. Linderman, M. Rachh, J. G. Hoskins, S. Steinerberger, and Y. Kluger. Fast interpolation-based t-SNE for improved visualization of single-cell RNA-seq data. *Nature Methods*, 16(3):243–245, Mar 2019. doi: 10.1038/s41592-018-0308-4 3
- [28] L. McInnes, J. Healy, and J. Melville. UMAP: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv: 1802.03426*, 2018. 2, 4
- [29] L. McInnes, J. Healy, N. Saul, and L. Grossberger. UMAP: Uniform manifold approximation and projection. *The Journal of Open Source Software*, 3(29):861, 2018. doi: 10.21105/joss.00861 2, 3, 5
- [30] B. H. Meyer, A. T. R. Pozo, and W. M. Nunan Zola. Global and local structure preserving GPU t-SNE methods for large-scale applications. *Expert Systems with Applications*, 201, Sept. 2022. doi: 10.1016/j.eswa.2022.116918 2
- [31] T. T. d. A. T. Neves, R. M. Martins, D. B. Coimbra, K. Kucher, A. Kerren, and F. V. Paulovich. Fast and reliable incremental dimensionality reduction for streaming data. *Computers & Graphics*, 102:233–244, Feb. 2022. doi: 10.1016/j.cag.2021.08.009 2
- [32] C. J. Nolet, V. Lafargue, E. Raff, T. Nanditale, T. Oates, J. Zedlewski, and J. Patterson. Bringing UMAP closer to the speed of light with GPU acceleration. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(1):418–426, May 2021. doi: 10.1609/aaai.v35i1.16118 2
- [33] L. G. Nonato and M. Aupetit. Multidimensional projection for visual analytics: Linking techniques with distortions, tasks, and layout enrichment. *IEEE Transactions on Visualization and Computer Graphics*, 25(8):2650–2673, 2019. doi: 10.1109/TVCG.2018.2846735 2, 3, 4
- [34] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011. 2, 5
- [35] P. G. Poličar, M. Stražar, and B. Zupan. openTSNE: A modular Python library for t-SNE dimensionality reduction and embedding. *Journal of Statistical Software*, 109(3):1–30, May 2024. 5
- [36] S. Raschka, J. Patterson, and C. Nolet. Machine learning in Python: Main developments and technology trends in data science, machine learning, and artificial intelligence. *Information*, 11(4), 2020. doi: 10.3390/info11040193 4
- [37] L. Reichmann, D. Hägele, and D. Weiskopf. Supplemental material for out-of-core dimensionality reduction for large data via out-of-sample extensions, 2024. doi: 10.18419/darus-4441 2, 5
- [38] G. Richer, A. Pister, M. Abdelaal, J.-D. Fekete, M. Sedlmair, and D. Weiskopf. Scalability in visualization. *IEEE Transactions on Vi-*

- ualization and Computer Graphics, 30(7):3314–3330, 2024. doi: 10.1109/TVCG.2022.3231230 2, 4
- [39] T. Sainburg, L. McInnes, and T. Q. Gentner. Parametric UMAP embeddings for representation and semisupervised learning. *Neural Computation*, 33(11):2881–2907, 10 2021. doi: 10.1162/neco.a.01434 2
- [40] S. Stolfo, W. Fan, W. Lee, A. Prodromidis, and P. Chan. KDD Cup 1999 data. UCI Machine Learning Repository, 1999. doi: 10.24432/C51C7N. 4
- [41] J. Tang, J. Liu, M. Zhang, and Q. Mei. Visualizing large-scale and high-dimensional data. In *Proceedings of the 25th International Conference on World Wide Web*, pp. 287–297. International World Wide Web Conferences Steering Committee, 2016. doi: 10.1145/2872427.2883041 2
- [42] D. Ulyanov. Multicore-TSNE. <https://github.com/DmitryUlyanov/Multicore-TSNE>, 2016. 2
- [43] L. van der Maaten. Learning a parametric embedding by preserving local structure. In D. van Dyk and M. Welling, eds., *Proceedings of the Twelfth International Conference on Artificial Intelligence and Statistics*, vol. 5 of *Proceedings of Machine Learning Research*, pp. 384–391. PMLR, Hilton Clearwater Beach Resort, Clearwater Beach, Florida USA, 16–18 Apr 2009. 2
- [44] L. van der Maaten and G. Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9(86):2579–2605, 2008. 2, 4
- [45] L. van der Maaten, E. O. Postma, H. J. van den Herik, et al. Dimensionality reduction: A comparative review. Technical Report TiCC-TR 2009-005, Tilburg University, 2009. 4
- [46] J. Venna and S. Kaski. Visualizing gene interaction graphs with local multidimensional scaling. In *ESANN 2006 Proceedings, 14th European Symposium on Artificial Neural Networks*, pp. 557–562, 2006. 4
- [47] Y. Wang, H. Huang, C. Rudin, and Y. Shaposhnik. Understanding how dimension reduction tools work: An empirical approach to deciphering t-sne, umap, trimap, and pacmap for data visualization. *Journal of Machine Learning Research*, 22(201):1–73, 2021. 2, 7, 9
- [48] H. Zhang, P. Wang, X. Gao, Y. Qi, and H. Gao. Out-of-sample data visualization using bi-kernel t-sne. *Information Visualization*, 20(1):20–34, 2021. doi: 10.1177/1473871620978209 3
- [49] H. Zhu, M. Zhu, Y. Feng, D. Cai, Y. Hu, S. Wu, X. Wu, and W. Chen. Visualizing large-scale high-dimensional data via hierarchical embedding of KNN graphs. *Visual Informatics*, 5(2):51–59, 2021. doi: 10.1016/j.visinf.2021.06.002 2