

# Flexible Visual Summaries of High-Dimensional Data using Hierarchical Aggregated Data-Subsets

David Pfahler\*

*Supervised by: Harald Piringer†*

VRVis Research Center  
Wien / Austria

## Abstract

The number of installed sensors to acquire data, for example electricity meters in smart grids, is increasing rapidly. This huge amount of collected data needs to be analyzed and monitored by transmission system operators. This task is supported by visual analytics techniques, but traditional multi-dimensional data visualization techniques do not scale very well for high-dimensional data. The main contribution of this paper is a framework to efficiently inspect and compare such high-dimensional data. The key idea is to partition the data by the semantics of the underlying data dimensions into groups. Domain experts are familiar with the meta-information of the data and are able to structure these groups into a hierarchy. The proposed system visualizes the subsets of the data by appropriate means. These visual summaries can then be used to support the explorative overview tasks of the user.

**Keywords:** Visual Analytics, Hierarchical Aggregation, High-Dimensional Data

## 1 Introduction

Creating an overview of the data is the first task a user requires of an information visualization [17]. However, traditional multi-dimensional data visualization techniques, like parallel coordinates or scatterplot matrices, do not scale very well for high-dimensional data [5, 21]. This leads to the need of analyzing the unfiltered and raw high-dimensional data before the interesting information can be presented to the analyst [6].

Explorative overview tasks are relevant in various application domains. Examples include comparing outputs of multi-run simulations in the automotive sector or monitoring multiple quality indicators of products in advanced manufacturing.

This work is motivated by High-Dimensional Overview (HDO) tasks in the energy sector — a domain where the amount of acquired data is increasing rapidly. Power generation, power consumption, and meteorological quanti-

ties are constantly measured by the providers, creating a vast number of time series. The number of sensors will grow even further with the advent of smart meters. Until the year 2020 EU member states are required to equip at least 80% of their consumers with smart meter devices [13]. The transmission system operators need to analyze and process this acquired time series in regular intervals, e.g. weekly or monthly. It is impossible for them to inspect every single acquired time series. Thus, they need to get an overview of the data first and interactively explore the data further to get an insight.

The primary contribution of this paper is a framework for analyzing and comparing high-dimensional data. The key idea is to partition the data by meta-information and to visualize the resulting subsets by appropriate means. In the energy sector, various data sensors are placed at the same location or share the same type, for example, temperature sensors. This meta-information of the data is familiar to domain experts and allows them to analyze and compare the data in a more intuitive way. An example task would be the comparison of multiple time series of power consumptions of multiple locations, where only locations are compared and power consumptions of the different sensors within a location are combined. This scales better than comparing every single data dimension like at the Rank By Feature Framework (RBFF) [16] or comparing every data record like in parallel coordinates or scatterplot matrices. Still, the user is able to flexibly drill-down on demand in order to explore the details of the different dimensions.

## 2 Related Work

A vast amount of scientific and application areas are confronted with high-dimensional datasets. Interactive visual analysis is an effective way to understand and process the data [5]. Approaches to visualize multi-dimensional datasets are an important topic of research. Traditional visualization techniques like parallel coordinates or scatterplot matrices are well suited for targeting this problem for a small number of dimensions [12]. But when the number of dimensions increases, these techniques fail. This is because of the boundaries of our visual system, visual

---

\*pfahler@vrvis.at

†hp@vrvis.at

clutter, and technical challenges [3].

There are multiple surveys that focus on high-dimensional data visualization [9, 12]. One way to maintain the scalability of the data is to reduce the information of the data by the number of data records or the number of data dimensions [14]. Principle component analysis, multidimensional scaling and self organizing maps are common techniques to reduce dimensions in data visualization. The drawback of these methods is that they produce a subspace that has no intuitive meaning to the data analyst [4].

A technique that supports the generation of meaningful subspaces is called Visual Hierarchical Dimension Reduction [20]. It uses a similarity measure to hierarchically cluster the dimensions and allows the user of this framework to interactive explore and modify the created hierarchy. From this hierarchy clusters, a meaningful subset is selected. Representative dimensions of these selected clusters are then visualized as visual representations. The drawback of this method is that not all dimensions are used for the encoding of the visual representations.

The RBFF by Seo and Shneiderman [16] ranks small preview visualizations of one dimension or two dimensions by statistical properties which give a good initial overview of all dimensions. This approach scales well for the number of data records but has limitations regarding the number of dimensions. Especially for the comparison of dimension pairs (a scatterplot matrix), the number of simultaneously displayed pairs increases quadratically, but also in the case of one-dimensional statistics the limit of displayed visual representations that can be handled reasonably is a few hundred [15].

Stole and Hanrahan developed the system “Tableau” (former “Polaris”) [18]. It introduces a scalable pivoting algebra on meta-information on data records. This work applies this algebra on the partitioning of data dimensions into subsets using the categories of a data dimension as meta-information.

Elmqvist and Fekete presented a model for implementing hierarchically aggregated visualizations [2]. Their proposed guidelines are used in this work to create visual aggregates of the data subsets that are more scalable for the limited perceptual capabilities of a human viewer.

### 3 Tasks and Goals

This section characterizes tasks that are needed for an explorative overview of high-dimensional data in the energy sector. These are used to derive the goals of the design process of the Hierarchical Data Overview (HDO) framework.

#### 3.1 Task Analysis

Transmission system operators in the energy sector acquire time series data from different sensors on a regular basis. These are used for power control and risk manage-

ment. The inspection and analysis of newly acquired data is hence a frequent, recurring and important activity.

The time spent looking at the data can be shortened by identifying the recurring tasks a user needs to fulfill. This paper aims to focus on the following tasks:

**T1 - Finding structures:** A key task of a data analyst is to get insight into the data and to validate or discard an initial hypothesis. Hypotheses often refer to structures in the data. In contrast to the validation of expected structures, the discovery of new structures is also a user task. Important structures in time-dependent data are listed below and clarified by an example:

**Trends:** “Is the data increasing or decreasing?”, “Does it have recurring peaks, troughs or plateaus?”

**Modalities:** “Are the distributions of the underlying data uni- or multimodal?”

**Outliers:** “Are some of the data not fitting the general trend?”

**T2 - Rank by feature:** As described by the RBFF, a user is interested in the dimensions that match a specific feature. These features can be statistical properties like the median of the dimension. An example of this user task would be the exploration of the biggest electrical loads in the electric grid.

**T3 - Assessing the purity of groups:** By merging the data dimensions by meta-information into groups, it is necessary to identify whether or not the grouping is sufficient for the user. By characterizing the purity of these meta-information based groups, the user is able to make further decisions. An example for a follow-up question is: “How much does the given grouping coincide with the similarity of the data?”

**T4 - Exploration and Tuning:** After the user was able to get an initial overview, further questions concerning the data may arise. Those can be answered by exploring and tuning the created groups and receiving more detailed information. An example of this would be the question: “Is the purity of a group improving when I drill-down the group with another partitioner?”

#### 3.2 Design Goals

Based on the task analysis three design goals were established. These goals guided the design process of the HDO framework.

**G1 - Visual summaries of groups:** To support the user to find structures in the data (T1), efficient visual summaries of groups of (large numbers of) dimensions need to be displayed. The requirement on the summaries is that they give a good reproduction of statistical position, variance and distribution and also the trend of them over time or over categories (T3).

**G2 - Flexible drill-down and roll-up:** With respect to T4 the overview visualizations need to be explored in depth. The concept of drill-down and roll-up with respect to “any

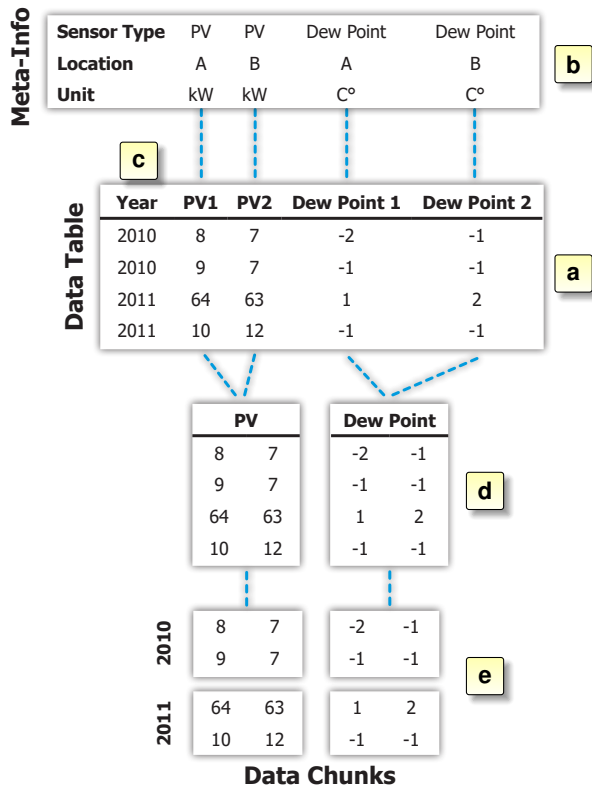


Figure 1: The subdivision of the data table (a) by utilizing the meta-information of data dimensions (b) and data records (c) to data chunks (d, e)

known structure of the feature space” [19] enables this fast change of the viewing granularity. Also, the ranking by a specific feature (T2) reduces the displayed visualizations to a set of interest.

**G3 - Scalability:** Like previously mentioned the number of data dimensions and records is rising. The framework should not be limited by any inherent upper limit of dimensions or data records. This also concerns the visual complexity of the used visualization. The goal of this framework is to support simple monitoring and reporting tasks (T1, T2), but also to allow detailed exploration tasks (T3, T4). This implies a trade-off between the simplicity and cognitive ease of the visualizations and the preserving of the distributions, modalities, and outliers of the underlying data.

## 4 Data Model

A raw dataset is most commonly available as a table, made up of columns and rows. In this work, a column is referred to as a data dimension and a row is referred to as a data record. Fig. 1 shows a data table, which is used as a guiding example for the data model.

In the energy sector, the data dimensions are numerical time series from various sensors. In Fig. 1a this time se-

ries are *Production Values* from photovoltaic plants (*PV*) and temperature measurements from *Dew Point* sensors. Further the data dimensions can contain categorical meta-information on the data records. One common categorization is the partitioning of the time value into time intervals, like the year of the data record. This is shown as the *Year*-column in Fig. 1c. Additionally, meta-information can be assigned to the data dimensions themselves. This could be the information of the location of a sensor, the measuring unit or the type of the sensor (see top table in Fig. 1b).

The numerical data dimensions can be of a similar type on a common scale (e.g.: Multiple power consumption sensors in watt). This enables the comparison for similar distributions or sequences, or the detection of outliers (T1). In contrast to the common scale, the dimensions can also have different units (e.g.: Weather time series with temperature, wind speed or wind direction) with no common scale. In Fig. 1a the *kW* values of the *PV* sensors are on a different scale than the *C°* values of the *Dew Point* time series.

The assigned meta-information can be present in a hierarchical manner. For data record based meta-information the categorical data dimensions can describe a level of detail (e.g. For time-categories the refinement over years to months to days). Also for data dimension based meta-information, a hierarchy can be derived. For example, if the location of a sensor is assigned as meta-information, the different levels of detail of the sensor hierarchy can be present from state to city to house to sensor type to the data dimension.

## 5 Description of the Hierarchical Data Overview

This section describes the visualization method for the HDO framework. Its design is engineered by the defined goals from Section 3.2 for an application to data described in Section 4.

By utilizing the meta-information, described in the previous section, the data table can be subdivided by data dimensions and by data records into smaller blocks of data. In this work, these blocks are referred to as *data chunks*. Fig. 1d shows the combination of the data dimensions by the meta-information *Sensor Type* to the data chunks *PV* and *Dew Point*. Not shown but also possible subdivisions would be the combination of the dimensions by the *Location* or the *Unit* meta-information.

The data chunks shown in Fig. 1e are created by the additional subdivision of the previously created data chunks by the combination of meta-information on data records (in this example by the *Year* of the time series). The key concept of this framework is to create a hierarchical relationship between these data chunks and to visualize a combination of them.

The overview visualization is designed using a hierarchical tabular layout (see Fig. 2). The design decision of

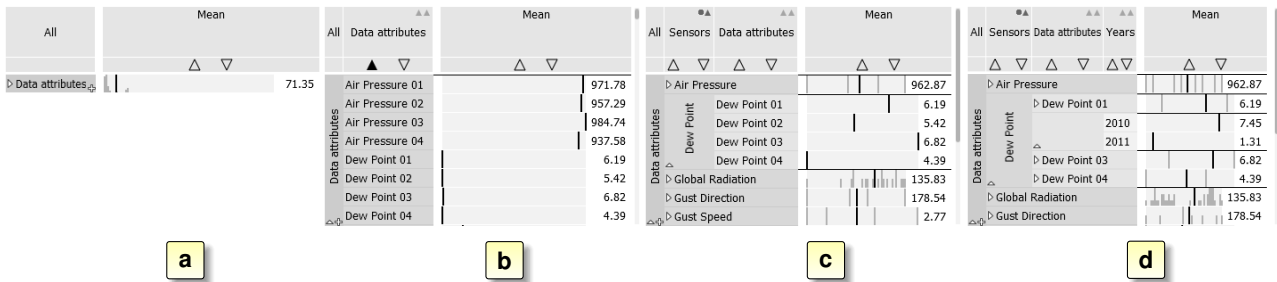


Figure 2: The HDO framework is visualized by a tabular layout. The left-hand side defines the rows of the visualization, which are defined by the hierarchy of the data model. The right-hand side visualizes the data chunks of the rows inside columns. The hierarchy can be collapsed (a), expanded by the data dimensions (b), combined by the sensor type (c) or refined by a categorical attribute (d). The refinement increases the number of visualized data chunks.

using a table-oriented display enables independent visual encodings of different aspects of the displayed data [8]. Additional previous work showed that the users are familiar with this kind of layout [1].

The tabular layout consists of two orthogonal parts: Rows (Section 5.1) and columns (Section 5.2).

A *Row* defines a combination of data chunks and a *Column* is responsible for creating a visual summary of a descriptive quantitative feature of the data chunks (G1).

## 5.1 Table Rows

The hierarchical structure is a key concept to ensure the scalability (G3) of the HDO framework [2]. The system offers the user an interactive hierarchical definition of the rows of the tabular layout.

Visually the first part of the tabular layout of the framework corresponds to the hierarchy of the data chunks. One row of this table is a node or a leaf of this hierarchy. In Fig. 2b the marked row shows a single data dimension *Dew Point 03* and in Fig. 2c the marked row shows all data dimensions with the same meta-information *Gust Direction*. The table headers of the left-hand part of Fig. 2d show the different hierarchy levels. The first column, marked with *All* in the example, is the root node of the hierarchy, the second column the sensor type, the third column the data dimensions and the last column a partitioning by the year of the data dimensions.

Fig. 2 visualizes the drill-down of the data table into a hierarchy of data chunks. The displayed column *Mean* plots the mean values of the underlying data chunks (see Section 5.2). Initially, the hierarchy consists only of one node, which contains all assigned data dimensions (Fig. 2a).

Additional hierarchy levels refine the data table into data chunks (G2). A usual refinement is the partitioning of the data table by data dimensions. Fig. 2b shows the table with all assigned data dimensions (in the figure called *Data attributes*) which displays a similar layout to the RBFF. However, too many dimensions are assigned to the visualization to display all of them on the limited screen area.

By combining data chunks, the number of rows can be reduced. In Fig. 2c the data dimensions are combined according to their sensor type. One can see that the number of displayed data chunks increases. For example, the sensor type *Air Pressure* contains four dimensions which are plotted in the mean column (see Section 5.2.1). In Fig. 2d the hierarchy level *Year* is assigned to the table. As previously described, this level partitions the data records of the data chunks by the categorical meta-information (see Fig. 1). As shown in the visualization of the dimension *Dew Point 01*, the number of chunks has increased.

The user is able to control the order and the level-of-detail of the levels (see Section 6). Additionally, the visible set of rows can be defined by collapsing and expanding the hierarchy nodes individually. Fig. 2d shows the expanded node *Dew Point* of the hierarchy level *Sensors* with an additional refinement for the dimension *Dew Point 02*.

The interactive refinement supports the identified G3 by enabling a visual scalability for a high number of dimensions and still allows the user to explore and tune (T4) the data.

## 5.2 Table Columns

The orthogonal part of the tabular layout is the visual encoding of different aspects of the data chunks as columns of the table. To maintain the scalability of the visual complexity (G3) of the visual representations of the data chunks the design decision for the visual encodings are simple and commonly used visualizations. For example in Fig. 2 the mean of a data chunk is visualized as lines inside the cells of the column.

### 5.2.1 Visual Aggregation Aspect

One column is responsible for calculating and visualizing a specific descriptive statistic for every row. These statistics are used to quantitatively describe and summarize different features of the defined data subsets (G1) [10]. The system differs three different classes of descriptive statistics: Central tendency, dispersion, and frequency distribution. Fig. 3 shows the visual encodings of these three classes.

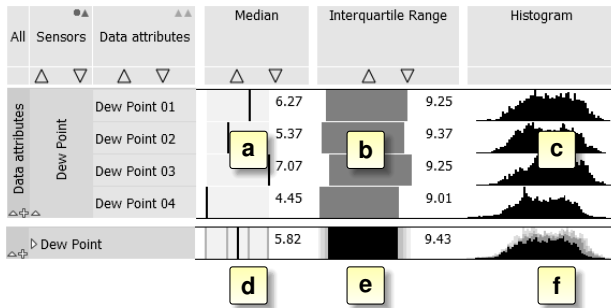


Figure 3: The data chunks are visually aggregated in the columns. A cell can visualize one leaf of the hierarchy (a, b, c) but is also able to visualize a node by combining the underlying leaves (d, e, f).

A simple way to visualize a statistic feature is by means of a textual representation. However, this representation does not give a good reproduction of the position, scattering, and distribution (G1) to the other displayed values. By plotting the features of the data chunks onto an axis a more descriptive representation can be achieved because one visual aggregate can be compared with another in an intuitive way.

**Central Tendency:** An univariate statistic can describe a point inside an axis. This includes the central tendency. Examples of this class of statistics are the average, the extreme values (minimum, maximum) or the percentile values (Median, Quantiles). Like other univariate statistics, they can be used for ranking (T2) but also for other tasks like finding outliers (T1). A textual representation of the aggregation, as shown in Fig. 3a on the right-hand side of the cell, enables the user to receive the actual value. To visualize the value of the aggregation, a line is positioned between the extents of the axis. This enables the user to compare the value of one data chunk with another (T1). When multiple data chunks have to be visualized within one cell, the same visual aggregate can be used. The location of the underlying features is plotted as gray lines, and the combined statistic is shown as a black line (see Fig. 3d). If more than one line is drawn at the same location, the visualization turns into a bar chart as shown in the *mean* cell of the *Global Radiation* row of Fig. 2d.

**Dispersion:** The second class of univariate statistics is the dispersion. It describes a positive range inside the axis. Examples for this class of descriptive statistics are the standard deviation or the Inter-Quartile Range (IQR). Similar to the previous encoding, an area is positioned around a dependent first moment inside the axis of the aggregation (e.g.: Area in the size of IQR around the median, Fig. 3b). When the dispersion of multiple data chunks has to be visualized within one cell, the range of the underlying ranges can be visualized in an aggregated cell by over-plotting the areas (see Fig. 3e).

**Frequency Distribution:** An example of a frequency distribution is the histogram. A textual representation is no

longer suitable because multiple values are visually encoded. To support details on demand (T4), detailed information of a bin value of a histogram can be displayed with tool-tip information. In contrast to the previous classes, it is not possible to rank the data chunks by a frequency distribution.

The binning of the histogram depends on all data chunks inside the axis (see Fig. 3c). When the frequency distribution of multiple data chunks has to be visualized within one cell, the histograms of the underlying data chunks are plotted over each other. The darker a part of a bin the more histograms overlap at this position (as shown in Fig. 3f).

## 5.2.2 Partitioning Aspect

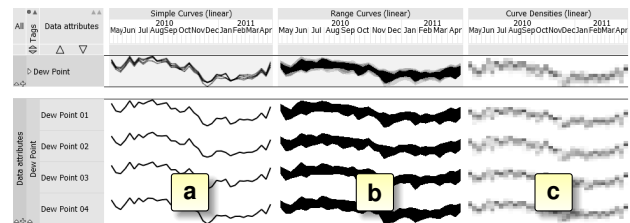


Figure 4: Central tendencies are shown as line graphs (a). Dispersions are shown as areas graphs (b). Frequency distributions are shown as heat maps (c).

Another important information is the trend of a descriptive statistic over time (T1). The system supports the partitioning of a column into sub-columns. Then the descriptive statistics for the data chunks are calculated for every partition of the column. This provides a global overview over local relationships of the statistic features (G1) [11].

Fig. 4 shows a partitioning of the column into time intervals (in this example thirds of months). For every row and every partition, the values from the three classes that were proposed in the previous section are computed and visualized:

**Central Tendency:** The calculated descriptive statistic of every partition is connected with a line, resulting in a line chart for every data chunk. Multiple lines are drawn inside one cell if multiple data chunks belong to one row (see Fig. 4a). This representation can be used to analyze the trend of the data or to find outliers (T1).

**Dispersion:** Similar to the simple visual aggregate, an area is drawn which is positioned around the dependent first moment for every partition. Its upper and lower boundaries are defined by the value of the dispersion. The upper and lower boundaries for all partitions are connected, and the area between them is filled (see Fig. 4b). This representation can be used to characterize the purity of the underlying data (T3).

**Frequency Distribution:** The frequency distribution of all data chunks inside a sub-column is calculated. A visualization called Curve Density Estimates [7] is used to encode the distributions of all sub-columns (see Fig. 4c).

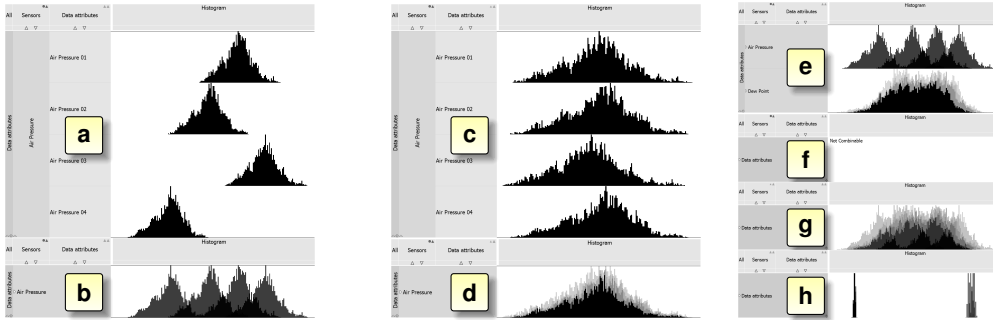


Figure 5: The axis of the visualized data chunks can be configured by the hierarchy levels. To compare the position of the data chunks, the nodes of the level are visualized on a common scale ( $\blacktriangle\blacktriangle$ ) (a, b). To compare the shape of the data chunks, every node receives an own scale ( $\blacktriangle\blacktriangle$ ) (c, d). If the nodes of levels are not combinable ( $\bullet\blacktriangle$ ) no visualization can be plotted. The rows (f, g, h) show the data chunks of two sensor types (e) for the scales  $\bullet\blacktriangle$ ,  $\blacktriangle\blacktriangle$  and  $\blacktriangle\blacktriangle$

This representation can be used to analyze the modalities of the underlying data chunks (T1).

### 5.2.3 Combination Aspect

All shown examples have used the same axis for creating the visual aggregate. A precondition is, however, that not all data chunks share a common scale or are combinable.

**Combination:** The created hierarchy levels, as described in Section 5.1, can have the property that they separate the data into distinctive nodes, in which the combination of the data chunks has no useful meaning. In the visualization, this is indicated by displaying the icon  $\bullet\blacktriangle$  in the header of the hierarchy level (see Fig. 5). An example of not combinable nodes is the partitioning of the data by the sensor type of the underlying dimension, in which it makes no sense to plot a temperature and a voltage value in the same coordinate frame (see Fig. 5f).

However, a level can also separate the data into chunks that are comparable with each other (for example the partitioning of power consumption time series by their location). This is indicated with the icon  $\blacktriangle\blacktriangle$  in the header (see Fig. 5a).

**Scaling:** If two nodes of a hierarchy are set to not comparable ( $\bullet\blacktriangle$ ), they cannot share a common axis. The common idiom “*comparing apples and oranges*” states that a non suitable comparison would indicate a false analogy (see Fig. 5f).

In contrast to this idiom, if two nodes are visualized that are set to be comparable, the scaling of the visualizations needs to be defined. One possibility is that the two nodes can share a common scale, which is indicated with  $\blacktriangle\blacktriangle$ . Another option is that every node may use its own scale and the visualizations are only displayed overlaid or in juxtaposition ( $\blacktriangle\blacktriangle$ ).

If the hierarchy level *Data attributes* is set to a common scale ( $\blacktriangle\blacktriangle$ ), it is possible to compare the position of the underlying data chunks (see Figs. 5a, 5b and 5h). As opposed to this, it is not as easy to compare the shape of the frequency distributions of the different data chunks (T1). The level can be set to no common scale ( $\blacktriangle\blacktriangle$ ), to be able

to compare the shape. Thereby, each node of the level receives its own axis. Figs. 5c and 5d show that the shapes of the frequency distributions of the different data chunks can be compared.

## 6 Exploring the Hierarchy

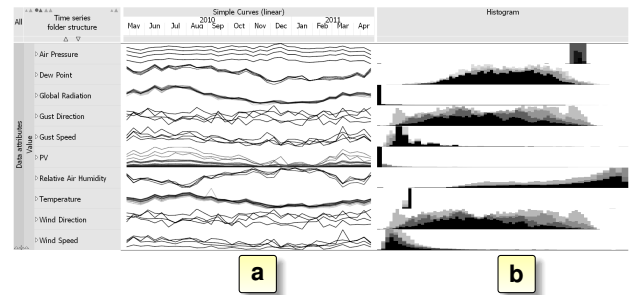


Figure 6: The HDO-framework visualizes the time curves (a) and the distributions (b) of 163 data dimensions combined by their sensor type.

This section describes how users are able to interact with the framework to be able to address their explorative overview tasks. As a precondition, it is assumed that the domain expert knows what data types, sensors, and meta-information one can expect from the data. Typically a user does not start his exploration task with no initial hierarchy like shown in Fig. 2a. A pre-defined set of hierarchy levels that are relevant in the domain is assigned. Fig. 6 visualizes 163 data dimensions of the sensor data set from the guiding example. The initial hierarchy is the sensor type, and the underlying data dimensions, which are collapsed so that the user is able to get an overview of the data.

To address the G1 the user is able to find structures by looking at the visual summaries (T1). Example structures that can be observed in Fig. 6 are: The yearly trend of meteorological quantities like the *Temperature* and the related *Production Values*. The modality of the *Gust Direction* which is different for every sensor. The outliers of the data

for example the null value in the *Air Pressure* distribution and the purity of some groups like the *Global Radiation* (T3).

The exploration process supports the extension of the complexity of the displayed visual representations to provide the user more information to get insight into the data. First, it is possible to increase the visual aggregation aspect of the data. This concept is described in Section 5.2.1 and shown in Fig. 3 where complexity increases from the central tendency (median values) to the dispersion (interquartile ranges) to the frequency distribution (histograms) of the data chunks. This increases the visual complexity but also enables the user to find structures, like the modality of the data chunks or outliers. The second complexity increase is the partitioning of the columns into sub-columns as described in Section 5.2.2 and shown in Fig. 4. This enables the user to analyze the trend and the modality of the data chunks.

To address the G2, the framework supports the drill-down and roll-up of the hierarchy nodes as described in Section 5.1 and shown in Fig. 2. Whilst exploring the visualizations, a user may observe that the variance in some displayed nodes is very high, and one wants to drill-down the node to see if the purity of the underlying nodes in the next hierarchy level increases (T3). This is achieved by clicking on the arrow inside of the hierarchy node. Fig. 2d shows this drill-down of multiple hierarchy nodes to a detailed representation.

Furthermore the user is able to refine the hierarchy further, by adding more hierarchy levels, by clicking on the plus sign in the lower left corner (Fig. 6). The further refining of the hierarchy may increase the purity of the newly partitioned data chunks.

Additionally the user may change the order of the hierarchy levels, by dragging the headers to another position. Fig. 7 shows the rearrangement of the hierarchy levels. The comparison of the individual *Temperature* sensors does not reveal specific structures in the data chunks (Fig. 7a). By dragging the columns *Sensors* and *Data attributes* to the end of the hierarchy, the comparison of *Years* of the underlying sensor data is possible and a new structure in the data is revealed (Fig. 7b).

To address the T4 further, every column can be configured individually. Possible interactions include the filtering of specific categories, the restriction of the displayed partitions of sub-columns, and further parameterizations of the visualizations of the data chunks.

## 7 Implementation Aspects

This framework has been implemented in C++ and uses OpenGL for rendering. The visual feedback of the visualization is in real-time to help the user to explore the data faster and thereby support him to make decisions faster. However, the computational costs of measurements on high-dimensional data are high. Several possibilities

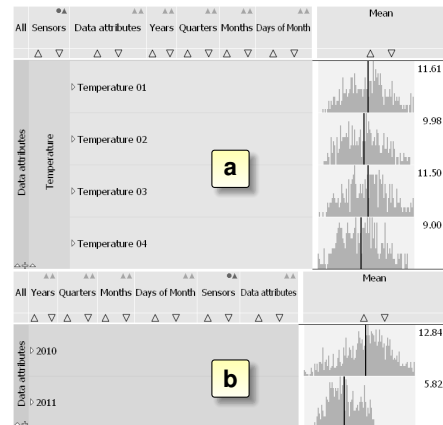


Figure 7: The rearrangement of the order of the hierarchy levels from the comparison of the *Data attributes* (a) to the comparison of the *Years* (b) enables the user to analyze different data chunks in the cells of the data columns.

to achieve real-time capability can be utilized: The hierarchy is used to compute the results of the measurements on higher levels of the hierarchy by a bottom-up approach. This is, for example, possible for the statistic means of values. Unfortunately, not all statistics can be calculated by reusing intermediate results of lower levels. For example, the quantiles of multiple chunks in an intermediate node of the hierarchy needs to be recalculated by using the data of the combined chunks.

To support these expensive tasks, the joint computation of results for multiple dimensions is considered. For example, a data dimension is sorted only once for all measures that need a sorting of subsets of the dimension.

## 8 Conclusion and Future Work

This work describes a Hierarchical Data Overview (HDO) framework to efficiently inspect and compare high-dimensional data. Motivated by the tasks of domain experts in the energy domain, three design goals are defined for this framework: Visual summaries, flexible interaction and the scalability for high dimensional data.

Based on these goals the HDO framework is described, which utilizes meta-information of the assigned data dimensions to partition the dimensions into data chunks. In a tabular layout, multiple descriptive statistics of these chunks can be visualized. The interactive refinement of the displayed rows and the flexible configuration of the columns of the tabular layout supports the interactive exploration tasks of domain experts.

The presented results show that this approach is able to address the identified goals. Future work could include hierarchy levels that partition the underlying data chunks automatically by a model or including columns that compute a bi-variate analysis that compares data chunks with another dimension.

## Acknowledgments

I would first like to thank my supervisor Dr. Harald Piringer of the VRVis research center in Vienna. He provided me with scientific guidance and always suggested helpful improvements. I would also like to acknowledge Dipl.-Ing. Thomas Mühlbacher of the VRVis research center for reviewing and supporting me during the development of this work.

## References

- [1] Clemens Arbesser, Florian Spechtenhauser, Thomas Mühlbacher, and Harald Piringer. Visplause: Visual data quality assessment of many time series using plausibility checks. *IEEE Transactions on Visualization and Computer Graphics*, 23(1):641–650, 2017.
- [2] Niklas Elmqvist and Jean-Daniel Fekete. Hierarchical aggregation for information visualization: Overview, techniques, and design guidelines. *IEEE Transactions on Visualization and Computer Graphics*, 16(3):439–454, 2010.
- [3] J-D Fekete and Catherine Plaisant. Interactive information visualization of a million items. In *Information Visualization, 2002. INFOVIS 2002. IEEE Symposium on*, pages 117–124. IEEE, 2002.
- [4] Arthur Flexer. On the use of self-organizing maps for clustering and visualization. *Intelligent Data Analysis*, 5(5):373–384, 2001.
- [5] Daniel A Keim, Jörn Kohlhammer, Geoffrey Ellis, and Florian Mansmann. *Mastering the information age-solving problems with visual analytics*. Florian Mansmann, 2010.
- [6] Daniel A Keim, Florian Mansmann, Jörn Schneidewind, and Hartmut Ziegler. Challenges in visual data analysis. In *Information Visualization, 2006. IV 2006. Tenth International Conference on*, pages 9–16. IEEE, 2006.
- [7] O Daae Lampe and Helwig Hauser. Curve density estimates. In *Computer Graphics Forum*, volume 30, pages 633–642. Wiley Online Library, 2011.
- [8] Alexander Lex, Nils Gehlenborg, Hendrik Strobel, Romain Vuillemot, and Hanspeter Pfister. Upset: visualization of intersecting sets. *IEEE transactions on visualization and computer graphics*, 20(12):1983–1992, 2014.
- [9] Shusen Liu, Dan Maljovec, Bei Wang, Peer-Timo Bremer, and Valerio Pascucci. Visualizing high-dimensional data: Advances in the past decade. *IEEE Transactions on Visualization and Computer Graphics*, 2016.
- [10] Prem S Mann. *Introductory statistics*. John Wiley & Sons, 2007.
- [11] Thomas Mühlbacher and Harald Piringer. A partition-based framework for building and validating regression models. *IEEE Transactions on Visualization and Computer Graphics*, 19(12):1962–1971, 2013.
- [12] Tamara Munzner. *Visualization Analysis and Design*. CRC Press, 2014.
- [13] Council of European Union. Council regulation (EU) no 189/2014, 2014. <http://eur-lex.europa.eu/legal-content/EN/TXT/?uri=SWD:2014:189:FIN>.
- [14] Harald Piringer. *Large data scalability in interactive visual analysis*. PhD thesis, Institute of Computer Graphics and Algorithms, Vienna University of Technology, 2011.
- [15] Harald Piringer, Wolfgang Berger, and Helwig Hauser. Quantifying and comparing features in high-dimensional datasets. In *2008 12th International Conference Information Visualisation*, pages 240–245. IEEE, 2008.
- [16] Jinwook Seo and Ben Shneiderman. A rank-by-feature framework for interactive exploration of multidimensional data. *Information Visualization*, 4(2):96–113, 2005.
- [17] Ben Shneiderman. The eyes have it: A task by data type taxonomy for information visualizations. In *Visual Languages, 1996. Proceedings., IEEE Symposium on*, pages 336–343. IEEE, 1996.
- [18] Chris Stolte, Diane Tang, and Pat Hanrahan. Polaris: A system for query, analysis, and visualization of multidimensional relational databases. *IEEE Transactions on Visualization and Computer Graphics*, 8(1):52–65, 2002.
- [19] Cagatay Turkay, Arvid Lundervold, Astri Johansen Lundervold, and Helwig Hauser. Representative factor generation for the interactive visual analysis of high-dimensional data. *IEEE Transactions on Visualization and Computer Graphics*, 18(12):2621–2630, 2012.
- [20] Jing Yang, Wei Peng, Matthew O Ward, and Elke A Rundensteiner. Interactive hierarchical dimension ordering, spacing and filtering for exploration of high dimensional datasets. In *Information Visualization, 2003. INFOVIS 2003. IEEE Symposium on*, pages 105–112. IEEE, 2003.
- [21] Jing Yang, Matthew O Ward, and Elke A Rundensteiner. Visual hierarchical dimension reduction for exploration of high dimensional datasets. 2002.