

A Systematic View on Data Descriptors for the Visual Analysis of Tabular Data

Information Visualization
XX(X):1–22 (to appear)
©The Author(s) 2016
Reprints and permission:
sagepub.co.uk/journalsPermissions.nav
DOI: 10.1177/1473871616667767
www.sagepub.com



Hans-Jörg Schulz¹, Thomas Nocke², Magnus Heitzler³, and Heidrun Schumann¹

Abstract

Visualization has become an important ingredient of data analysis, supporting users in exploring data and confirming hypotheses. At the beginning of a visual data analysis process, data characteristics are often assessed in an initial data profiling step. These include, for example, statistical properties of the data and information on the data's well-formedness, which can be used during the subsequent analysis to adequately parametrize views, and to highlight or exclude data items. We term this information *data descriptors*, which can span such diverse aspects as the data's provenance, its storage schema, or its uncertainties. Gathered descriptors encapsulate basic knowledge about the data and can thus be used as objective starting points for the visual analysis process. In this paper, we bring together these different aspects in a systematic form that describes the data itself (e.g., its content and context) and its relation to the larger data gathering and visual analysis process (e.g., its provenance and its utility). Once established in general, we further detail the concept of data descriptors specifically for tabular data as the most common form of structured data today. Finally, we utilize these data descriptors for tabular data to capture domain-specific data characteristics in the field of climate impact research. This procedure from the general concept via the concrete data type to the specific application domain effectively provides a blueprint for instantiating data descriptors for other data types and domains in the future.

Keywords

Metadata, data profiling, initial data analysis, climate impact research

Introduction

Over the last two decades, visualization has matured into an important tool for data analysis. The scientific literature encompasses a plethora of visualization techniques that support *exploratory analysis* (i.e., hypothesis generation) and *confirmatory analysis* (i.e., hypothesis testing). Yet, these visual analysis techniques require certain constraints to be met by the input data to ensure their applicability and usefulness – for example, a certain quantity and quality of data. As a given input dataset rarely carries information about these aspects beyond the raw data, comprehensive data analysis methodologies start with an *initial analysis* (Adèr 2008) or *data profiling* (Gschwandtner et al. 2014) that aims to assess these data characteristics before going into the exploratory or confirmatory phase. The outcome of such an assessment are what we term *data descriptors*.

We define a data descriptor as any objective data characterization that captures data properties with a particular focus on the data's subsequent visual analysis. The objectiveness of the descriptor is of importance to not bias this base information on which the remainder of the visual analysis workflow rests. The term “data descriptor” was chosen to reflect this objectiveness and to set it apart from interpretive information construing the data. It thus shares the intention of similar concepts, such as *metadata* and *semantic data*.

Data descriptors explicitly encode a dataset's characteristics, such as irregularities (e.g., format violations or extreme

values) and regularities (e.g., data types or constant data values). These make it possible for subsequent visual analysis techniques, to check the found irregularities against required quality constraints and to adapt their parametrization to the found regularities. A common example for the latter is the parametrization of a meaningful and effective color scale according to known properties of the data. These properties can range from simple information about the data's type (Silva et al. 2011) to its spatial frequency (Bergman et al. 1995) or background knowledge about its semantics (Lin et al. 2013; Setlur and Stone 2016).

Taking this knowledge about a dataset into account when visualizing the data can be vital. An impressive instance of how visualization fails, when the facts that are known about a dataset are not taken into account, was just recently given in the IEEE VIS 2016 tutorial on “Perception and Cognition for Visualization” by Bernice Rogowitz. Her example shows how a poorly chosen and inadequately parametrized color scale hides the known properties of the Higgs Boson dataset instead of showing them.*

¹University of Rostock, Germany

²Potsdam Institute for Climate Impact Research, Germany

³ETH Zurich, Switzerland

Corresponding author:

Hans-Jörg Schulz, University of Rostock, Department of Computer Science, Albert-Einstein-Str.22, 18059 Rostock, Germany.

Email: hjschulz@informatik.uni-rostock.de

*see <http://root.cern.ch/rainbow-color-map>

Yet, useful data properties and metrics are scattered across different levels of detail, subsuming various heterogeneous information about data and spanning different subdomains of analysis. For example, data descriptors can relate to such diverse aspects of a dataset as its provenance, its storage schema, its uncertainties, or its descriptive statistical measures. From these few examples, it is easily understandable that these aspects are rarely considered and described in concert and taken into account only as they become relevant for a particular computational analysis or visual mapping.

This paper aims on one hand to bring these scattered approaches for describing data together in a systematic form. And on the other hand, it aims to illustrate how these approaches can be used to support the visual analysis process. To form such a systematic understanding of data descriptors, the paper makes the following contributions:

- It gathers a wide variety of data descriptors from different fields of visual analysis in a **generic classification** that is not restricted to a particular data type or application domain.
- This classification is then **instantiated for tabular data**, which is one of the most common types of data across various application domains, and a pipeline for gathering descriptors from tabular data is presented.
- To exemplify its use, this instantiation is then **adapted for climate impact research**, which has a high demand for data descriptors due to the heterogeneity of the various involved disciplines and their diverse data standards and implementations.

This systematic view on data descriptors will provide a solid and comprehensive base for their application and further investigation. In this way, our paper gives a blueprint for how descriptors for other types of data – e.g., textual data, image/video data, or graph/network data – can be systematically established and adapted to their respective application domains. The structure of this paper follows this overall direction from the generic to the specific, starting with the definition and classification of data descriptors in the following section.

A Classification of Data Descriptors

To concretize our introductory remarks, we define a **data descriptor** as objective data about data that is available to a visual analysis system. *Objectivity* captures the important aspect of independence of any a-priori assumptions about the data and of any preconceived goal or path of analysis. Note that this does not imply that the data itself must be objective, if it ever can be (Gitelman 2013) – only its description. As it is hard to exactly delimit objectivity in technical terms, we use two indicators that a data description is objective: *invariance* (i.e., the same dataset stemming from the same data source will always result in the same description) and *independence* (i.e., the description only depends on the dataset and its source and no external parameters). If these indicators are fulfilled by a data description, we deem it sufficiently probable that the description is not distorted or biased by an outside influence. Finally, the description's *availability* to a visual analysis system is important, as it

means that the description is machine-readable, ruling out, for example, solely verbal or diagrammatic descriptions.

While other fields have already embraced the idea of leveraging “data about data” (Duval 2001) – be it out of convenience or out of necessity – the visualization community has just started to explore this direction. Examples include the Metadata Mapper (rog 2011) that utilizes data descriptors to map data between different visual analysis components, as well as the Knowledge-based Visual Analytics Reference Architecture (Flöring 2012) that captures analytical results (i.e., knowledge about the analyzed data and the analysis process) and feeds them back into future analyses. As applications like these make use of a few selected data descriptors to reach their particular goals, a general overview of data descriptors for visual analysis remains an open point of research.

For giving such an overview, we assume without loss of generality the described data to be *self-contained* and *homogeneous*. In cases in which this assumption does not hold – i.e., for datasets that link to external data of unknown properties (not self-contained) or that are combined of a number of individual datasets of different structure and content (heterogeneous dataset) or both – descriptors can hardly be applied across the whole dataset. In these cases, the dataset can be partitioned into self-contained homogeneous subsets, then to be characterized individually by data descriptors appropriate for each of them.

For such a self-contained homogeneous dataset, our classification shown in Figure 1 categorizes data descriptors according to the aspect of the data they are describing. As a first distinction, data can either be looked at from a temporal perspective, i.e., the *data flow*, or from a structural perspective, i.e., the *data space*. Current literature on metadata, data properties, data models, or any related term or notion hardly ever considers data flow descriptors and data space descriptors in concert. This is most certainly due to the fact that data flow descriptors are mainly used in database applications and information management scenarios, whereas data space descriptors are mainly used in data analysis and data mining approaches. Yet, describing these two aspects of data together is common in other fields – for example, for describing multimedia not only content-wise, but also in terms of who produced it and for which audience (Arens et al. 1993), or for describing web-based resources together with their utility if limited by legal or other conditions (Steinacker et al. 2001). Hence, the following sections give examples of data descriptors for both aspects, their common notations and models, as well as how they are used for visual analysis purposes.

Data Flow Descriptors

Data flow descriptors (DFD) give details about where the data came from (*data provenance*), where it is now (*data storage*), and where it can go from there (*data utility*). Data provenance information can range from a simple model number and firmware version of the device used to record the data to a full-fledged protocol of all analysis and data processing steps it has already undergone. Information about the current storage of the data captures mainly if and how the data can be retrieved and thus be used in its current state. Such information can include, for example, details on

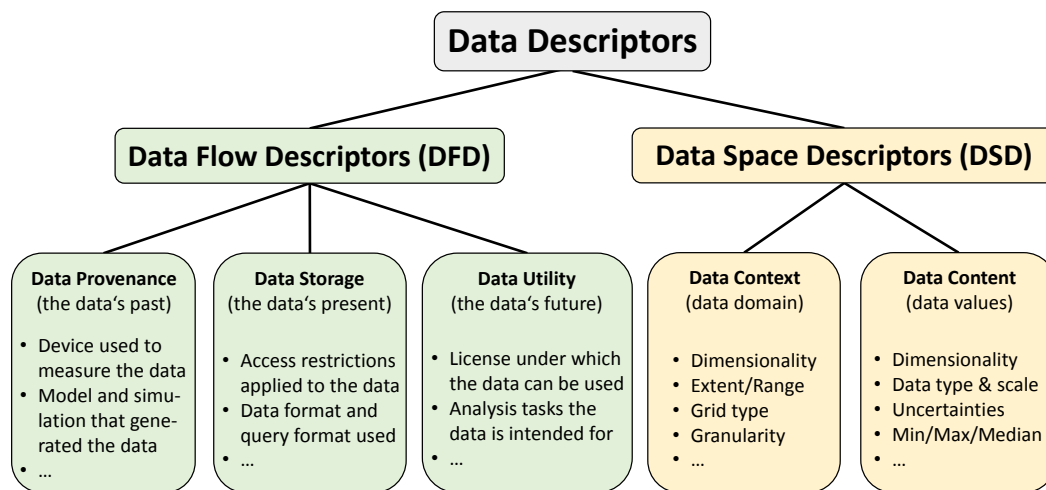


Figure 1. Our proposed classification of data descriptors by the data aspects they describe. The different classes of descriptors are arranged from general (top) to specific (bottom) and exemplified with a few common descriptors each.

the database schema (i.e., across how many tables does the data spread), or to how many queries per minutes the server restricts the access. Lastly, the data utility information can detail the uses for which a dataset is licensed, or it can specify for which purpose the data was collected.

Apart from a few papers contributing to the current research challenge of data provenance, literature from the visualization community is rather sparse on data flow descriptors. The reason for this may be that data flow descriptors are often less formal than data space descriptors or they make use of non-standard notations in which the description is given. Due to the diversity of data, standards for its description exist mainly for domain-specific cases, such as the ISO standard 19115-1:2014 (Organization for Standardization 2014) for geographic information or the Full-Metadata Format (FMF) (Riede et al. 2010) for data from scientific experiments. The following sections highlight data descriptors that have been proposed for the data aspects of provenance, storage, and utility.

Data Provenance Since data is ubiquitous in these days, it is not uncommon anymore to have at least minimal provenance information indicating who authored or curated a dataset and which version of the dataset one is looking at. With more knowledge about the origin and lineage of data, it may also be possible to judge the trustworthiness of the data and of the processes that generated it (Carata et al. 2014), or even to independently reproduce the data (Davison 2012).

A number of taxonomies for provenance descriptors exist. The taxonomy, which captures the widest range of provenance aspects is the one by Simmhan et al. (2005a,b). It contains not only such essential aspects as the data and the process that has produced it (*subject of provenance*), but also such diverse descriptors as the representation, the storage, the dissemination, and the use of the provenance information. Other taxonomies tend to focus more on the conceptual and technical aspects of collecting and managing provenance information – such as the ones by Glavic and Dittrich (2007) and by da Cruz et al. (2009). Regardless of their particular focus, each of these taxonomies can be thought of as a logical continuation and further subdivision of the class of provenance descriptors in Figure 1.

For recording and storing provenance information, either models from the database community (Buneman et al. 2001) or from the domain of scientific workflows (Cohen et al. 2006; Ludäscher et al. 2008) are more suitable – depending on whether the provenance information captures a series of data states or the sequence of processing steps that produced them. For interchangeability of provenance information beyond the data and software ecosystems of a particular domain, the Open Provenance Model (Moreau et al. 2011) has been developed.

In visual analysis, provenance is still largely an open research challenge (Keim et al. 2010, ch.3.3). Research results in this direction deal with capturing either the generation of the visualization – i.e., the visualization process itself (Freire et al. 2008; Silva et al. 2007), or the interaction history with the generated visualization – i.e., the knowledge discovery process (Groth and Streefkerk 2006; Heer et al. 2008). Along the lines of the current survey by Ragan et al. (2016), these two strategies can be called “follow the data” and “follow the user”, respectively. As an outcome of both, the generated visualization or the discovered knowledge can then be annotated with the information about how they were yielded. In particular in highly exploratory scenarios that are characterized by a frequent back and forth in which many different alternatives are tried, provenance information can become quite large and unwieldy. In these cases, the visual exploration of this provenance information poses a challenge in itself that is addressed by dedicated tools like *Tableau Behavior Graphs* (Heer et al. 2008) or *AVOCADO* (Stütz et al. 2016).

Data Storage For retrieving and querying data, information about its storage is needed. This information does not only entail descriptions of the data’s organization, such as the data structure or data schema, but also about the access mechanism with which to read and manipulate the data.

Descriptors that commonly hold such information can be given for different aspects of a data storage – e.g., for a logical or for a physical perspective (Vassiliadis 2009), as well as for stored data or for stored (business) rules/processes (Hoxmeier 2005; Vassiliadis 2009). The latter include descriptors that detail the behavior of a data

storage through usage statistics and security settings, which can help to optimize data access patterns or to understand access limitations, respectively.

Notations for information about the storage of a dataset are clearly centered on describing the logical model of data organization, with the relational model describing data as tables being one of the most prominent examples (Codd 1970, 1990). Other models include graph-based descriptions of data organization (Angles and Gutierrez 2008) or data being organized in a multi-dimensional space (Zhuge 2004; Zhuge et al. 2005). As most models come with their own notation, more abstract metamodels, like *information spaces* (Franklin et al. 2005), and a diverse set of standards, like the ISO/IEC 10027:1990 Information Resource Dictionary System (IRDS) (Organization for Standardization 1990) or the W3C Resource Description Framework (RDF) (World Wide Web Consortium (W3C) 2014) have been developed. They can be used to describe the specifics of different forms of data organization in a uniform way.

In visualization, information about the data storage is used in some cases for finding correspondences between data items and using them for visual highlighting (North and Shneiderman 2000) or for visual linking (Lieberman et al. 2011) in multiple coordinated views. Other authors utilize RDF-encoded information about the data organization to automatically establish mappings of data attributes to visual attributes without prior knowledge of the data sources and their schemas (Cammarano et al. 2007). While these approaches all work on data item level, others use information about datasets as a whole to visualize the landscape of all datasets in a particular data repository. For example, descriptors containing information about each dataset's size and server speeds can be used to graph an entire such *data landscape* to make an informed decision about which dataset to use for an analysis at hand (Tshagharyan and Schulz 2013).

Data Utility The utility of data, i.e., its intended and imaginable uses – which may not be the same, is rarely considered in the literature. Some of the data's utility (or lack thereof) may be inferred from its provenance, as for example outdated financial data cannot be used as a basis for day trading or stale patient records cannot be utilized to plan a medical procedure. Other utility aspects may be inferred from storage descriptors, as for example many database servers limit the number of queries per minute, which can severely hamper its usefulness for query-intensive analyses.

To the best of our knowledge, no list or taxonomy of data utility descriptors exist. Apart from legal constraints that might limit the use of a dataset (e.g., its license or confidentiality regulations), most of the literature on data utility relates to anonymized and/or obfuscated data (Russell 2008, ch.2.7). Depending on which methods were used for anonymization, the utility of the data may be limited to certain kinds of analyses. For statistical obfuscation methods (so called *statistical disclosure limitations*), metrics exist to measure the remaining utility of the data to find a reasonable trade-off between anonymity and usefulness (Karr et al. 2006). To the best of our knowledge, no such metrics exist for technical methods that have either introduced the utility limitation as an intended outcome (Grammer et al. 2012)

or as unintended byproduct – e.g., when data compression disrupts data properties (Bassiouni 1985). In these cases, the used method should be included with the provenance descriptor, so that while data utility cannot be automatically quantified, the user can at least be informed about them.

At this point, the notion of *data utility* is not widespread and the decision of which dataset to use for a particular analysis, or which analysis to perform on a given dataset is left to the analyst. Hence, a common notation or model for data utility to store such information alongside the data remains an open research challenge.

In visualization, privacy preservice is mostly reflected by methods that aim to provide a given level of anonymity in the resulting visualization, measured through screen-space privacy metrics (Dasgupta et al. 2013). Since there exist no standards for utility descriptors, visual analysis methods do not make use of them or adhere to them. The few visualization approaches, which aim at capturing some notion of utility, apply pragmatic models that link the available datasets with those visual and analytical techniques that are deemed appropriate for them by an expert user (Streit et al. 2012). This overall lack of concern with issues of data utility stands in contrast to the early observation in visualization that the *functional role of data* – i.e., its use – is a key data characteristic (Zhou and Feiner 1996).

Data Space Descriptors

Data space descriptors (DSD) give details about the data domain (*data context*) and the data values therein (*data content*). Properties of the data context describe aspects of the space in which the data was gathered or observed, as this is important, for example, to relate the data items to each other. Common instances are the observation space's dimensionality (e.g., 2D – Lat/Lon, 3D – Lat/Lon/Alt, 4D – Lat/Lon/Alt/Time), whether the data is scattered or gridded, and in case of the latter whether the grid type is structured or unstructured. Descriptors of the data content include properties, such as data type (e.g., scalar or vector) or min/max values, but also information about missing data or data that is affected by uncertainty.

Since the characteristics of the data space are of utmost importance for correct analysis and visualization of a dataset, they have been extensively investigated from the very beginning of visualization research. There exist a few descriptors that apply to both aspects of the data space – data context and data content – in the same manner. Examples of such descriptors are dimensionality of the data domain (context) and of the data values (content), as well as the scale type of each dimension (e.g., nominal, categorical, ordinal, or interval) (Stevens 1946; Roth and Mattis 1990). The following sections highlight data descriptors that are specific to either data context or data content.

Data Context The data context (often also called *data domain*, *independent variables*, or *primary key*) denotes the part of the data that specifies the frame of reference of the data values. Since the frame of reference is spanned via n axes in space, in time, or in some abstract space of identifiers, the data context can be understood as an n -dimensional space. A particular n -tuple specifying a point within that space forms the data context for its

associated data content, i.e., the gathered data values or data characteristics (Andrienko and Andrienko 2006). Knowledge about the data context is essential to determine, for example, the data's completeness – i.e., whether a data entry exists for all identifiers or locations.

One of the first characterizations of the data context was given by Zhou and Feiner (1996) under the term *data domain* and in particular *data domain entity*, which can be anything unique from a person to a point in time and space at which a measurement was taken. This generalizes other such notions, like the distinction between *coordinates* and *amounts* (Roth and Mattis 1990), or between the data types *1D*, *2D*, *3D*, *temporal* (Shneiderman 1996). These data domain entities can have a *point-wise*, *local*, or *global* extent for which they are deemed valid (Robertson 1991). For example, a given point in space could not only stand for this particular point, but for its local neighborhood as well. Finally, the characterization by Zhou and Feiner (1996) also defines *data relations* that can be used to describe a structure underlying the data domain, such as a grid or a multi-level topology. This is in line with the concept of *relations* by Roth and Mattis (1990), and with the data types *network* and *tree* from Shneiderman (1996). Some theories order these characteristics of the data context in layers that are hidden underneath the data content and only visible to the professional user (Lux 1998).

The characterization of the data context is probably the most influential in visualization research, as one prominent distinction between Information Visualization and Scientific Visualization is whether the data is spatially referenced (Tory and Möller 2004). Yet nowadays already 60%-80% of the data is geospatially referenced (Hahmann and Burghard 2013) – including document and image collections, whose depiction is usually considered to be part of Information Visualization. Hence, this common demarcation line is hard to uphold as most data is somehow spatially referenced. As a result, the distinction between scattered and gridded data becomes of increasing importance as a more meaningful characterization of Information Visualization and Scientific Visualization, respectively. Furthermore, the data context may indicate how to partition the data in a meaningful way, which can have repercussions all the way to the storage level (cp. OLAP).

Data Content The data content (often also called *attribute space* or *dependent variables*) denotes the part of the data that specifies the actual (gathered) data values. It is for this part of the data, for which probably the most descriptors exist and for which the border between descriptions from an initial analysis phase and analytical results from later analysis phases is the most blurred. For example, some literature considers clustering results as an inherent characteristic of the data content. Yet, requiring independence as part of the descriptors' objectivity forbids to consider it as such, as clustering depends on a number of subjective assumptions, such as a similarity threshold or even a predefined number of clusters (cf. *k*-Means clustering).

Purely descriptive properties of data content are, for example, the types of each data attribute. Abstract distinctions differentiate the data type merely as being *atomic* or *composite* (Zhou and Feiner 1996), while more

concrete descriptions in common software packages, such as OpenDX, distinguish further between *scalar*, *vector*, *matrix*, and *tensor* data. Another common data content descriptor relates to the quality of the data, which subsumes a whole range of dirty data properties, as they are surveyed by Kim et al. (2003), Oliveira et al. (2005), and Gschwandtner et al. (2012). This includes the overall *data quality* (Batini and Scannapieca 2006; Josko et al. 2016) and in particular the *data's uncertainty* (Ayyub and Klir 2006; Drogg 2009) that plays the most prominent role in visualization besides *missing data*, *unusable data*, or *undefined data*. On top of these given properties, it is common to derive further descriptors that can be computed without being biased by user parametrization – e.g., descriptive statistics (Cleary et al. 1996).

Notations for data content descriptors exist only partially in some data formats, such as the NetCDF format that will be discussed in further detail in the section on data from climate impact research. Only for the subset of data quality descriptors, specialized notations can be found. Among them are the ISO/IEC standard 25012:2008 (Organization for Standardization 2008), as well as a proposal for an extension to the Business Process Model and Notation (BPMN) to encode data quality requirements (Rodríguez et al. 2012).

A specific focus in the visualization community lies on representing data quality in general (Sulo et al. 2005; Josko and Ferreira 2016) and data uncertainty in particular, as communicating the data's trustworthiness is of essence when basing a visual analysis on it. Specifically for the challenge of uncertainty visualization, a number of extensive overview articles provide a good outline of the massive corpus of literature on this topic – see, for example, (Pang et al. 1997; Thomson et al. 2005; Correa et al. 2009; Skeels et al. 2010; Ward et al. 2011; Potter et al. 2012; Brodlie et al. 2012; Ristovski et al. 2014). Furthermore, the problem of visualizing missing data is frequently singled out as a particular challenge, since it is unclear how to show something that is not present. Notable approaches in this direction include *missing value charts* by Theus et al. (1997), *shadow plots* by Swayne and Buja (1998), *missingness profile plots* by Fernstad and Glen (2014), and *missingness maps* by Cheng et al. (2015). The opposite of missing data, namely duplicate data entries, are addressed by visual analytics tools, such as *D-Dupe* (Bilgic et al. 2006; Kang et al. 2008) and *GeoDDupe* (Kang et al. 2007).

Notations to describe the data space including data context and data content are, for example, the framework of Galhardas et al. (1998) that is based on first-order logic, as well as the *E-notation* by Brodlie (1992) and its extension into the *domino notation* (Brodlie and Noor 2007). Older variants are the fiber bundle-based notation by Butler and Pendley (1989) and the *L-notation* by Bergeron and Grinstein (1989). File formats with metadata capabilities are, for example, CDF, HDF, NetCDF, XDF, or XSIL. Data space descriptors like these are sometimes used to classify visual mappings – for example, as it was done by Rankin (1990) or Brodlie (1992).

Gathering Data Descriptors

The process of gathering data descriptors is not necessarily straightforward, as there exist at least three different sources for descriptors, which may differ in the required effort and their attainable reliability and objectivity. The first source for a descriptor is to *query* it from the data source, if it has been stored alongside the dataset, e.g., as annotation or supplemental material. The second source is to *derive* a descriptor by computing it from the dataset or by inferring it from other descriptors – i.e., inferring data utility from data provenance. The third source are the users with their background knowledge about the data, who can be prompted for *input* to specify a descriptor. On top of these basic mechanisms, combinations can be employed. A common combination is that a descriptor, which has been determined once through a computation or a user input, is then stored as an annotation to the dataset, so that it does not need to be recomputed or re-entered, but can be queried directly from the data source in the future.

Tool support for gathering generic data descriptors is rare. The gathering of data space descriptors is (if at all) only supported as a step in a larger process – e.g., for performing automated data quality assessment in *Profiler* (Kandel et al. 2012) or for generating automated previews of datasets in *AutoVis* (Wills and Wilkinson 2010). These tools understand the descriptors they compute as means towards a particular end and do not allow to uncouple them from the process in which they are embedded, even though they could be beneficial for other purposes as well. The only tools that are geared towards gathering generic data descriptors are designed for capturing data provenance and generally motivated by the goal of traceable and reproducible data analysis. Notable examples for such tools include the well-known frameworks for scientific workflow management *Karma* (Simmhan et al. 2008), *Kepler* (Ludäscher et al. 2006), and *Taverna* (Belhajjame et al. 2008). In the field of visual analysis, the most advanced provenance management is currently available from dedicated frameworks that capture the visualization process, such as *VisTrails* (Silva et al. 2007). Visualization tools can utilize *VisTrails*' features by integrating it through a common API (Callahan et al. 2008). Approaches that also aim to capture computational processing steps may be able to extract the provenance information from the analytical scripts being run on the data (Huq et al. 2013). The most comprehensive approach would be to use a system-wide capturing that spans different applications, as it is envisioned by *Glass Box* (Cowley et al. 2005, 2006).

After having established the fundamental notions of data descriptor classes and the different ways of gathering them, we make use of these concepts to compile data-type-specific descriptors for tabular data in the following section.

Data Descriptors for Tabular Data

This section concretizes our concept of data descriptors by taking a closer look at its concrete instantiation for tabular data. That includes the different descriptors such data entails, as well as methods to gather these descriptors if they are not supplied with the data.

Tabular data encompasses the overwhelming amount of data available in CSV files, spreadsheets, and relational databases. We assume tabular data to be given in the form of a single table (dataset) of rows (records), columns (variables), and cells (individual data items), which can be likened to Codd's 3rd normal form (Wickham 2014). Non-tabular data is often first transformed into a table, before being visualized. This is embodied in the first step of the visualization pipeline by Card et al. (1999) that performs a data transformation from raw data into data tables. More complex settings of multiple tables that are linked via foreign key relations can be broken down into this canonical form. As the different parts of the data correspond to sets of different cardinality – i.e., singleton (cell), tuple (row), multiset/bag (column), full dataset (table) – we call these different aspects of tabular data *granularities*.

Data Flow Descriptors for Tabular Data

As a first observation, we note that all data flow descriptors are applicable to all four granularities. For example, the data can have provenance information attached to individual values, to individual records, to individual variables, or to the entire table. While data provenance descriptors and data utility descriptors are conceptually independent of the kind of dataset they describe, technical particularities of the described dataset still require some adaptation. For example, there can be subtle differences depending on whether the provenance relates to relational databases and SQL queries (Glavic et al. 2013) or spreadsheets and embedded formulas (Asuncion 2011). These finer differences are usually captured by data storage descriptors that detail how to access the data, which is obviously different for relational databases and spreadsheets. For data having been stored from a spreadsheet in a CSV format, this could be whether the file is comma-separated or tab-separated, and how many lines of table header it contains. For data stored in relational database systems, this could be the schema of a table, as it is detailed by the `SHOW COLUMN FROM table` statement in SQL. The output of this statement also gives a number of data space descriptors, as they are discussed next.

Data Space Descriptors for Tabular Data

Data space descriptors are not only much more specific to tabular data than data flow descriptors, but they also apply mostly to its specific granularities. We list a number of typical data space descriptors for tabular data in Table 1. Note that this listing does not list uncommon usage of descriptors. For example, principally it is possible to have names for records or even individual cells and there certainly exist scenarios in which this is desirable – yet, it is not very common, so we list the descriptor “name” only for variables. Table 1 contains those general data space descriptors that also apply to tabular data and adds some data-type-specific descriptors. The descriptors in Table 1, which were not mentioned in the previous section, are:

Unit & valid range: Given a variable's unit of measurement, we can imply various other properties, such as the variable's semantics (e.g., degrees Celsius indicate

Data Space Descriptors Granularity	Data Context	Data Content
Value (Cell)	<ul style="list-style-type: none"> contextual uncertainty (e.g., uncertain position or time point) 	<ul style="list-style-type: none"> uncertainty of measurement, calculation, or simulation type of value (regular, missing, undefined)
Record (Row)	<ul style="list-style-type: none"> context outlier neighborhood (e.g., connected records via grid structure) 	<ul style="list-style-type: none"> content outlier topological feature (e.g., critical point)
Variable (Column)	<ul style="list-style-type: none"> variable name scale type (e.g., nominal, ordinal, interval) unit & valid range extent (point, local, global) type of dimensions (spatial, temporal, identifier, other) 	<ul style="list-style-type: none"> variable name scale type (e.g., nominal, ordinal, interval) data type (scalar, vector, matrix, tensor) unit & valid range univariate statistical measures (e.g., min, max, mean, skewedness) spatial/tempor. pattern (e.g., constancy, monotonicity, periodicity)
Dataset (Table)	<ul style="list-style-type: none"> kind of space (e.g., Euclidean) spatial dimensionality variable relations (e.g., day+month+year, first name+surname) grid type (structured, unstructured) variable combinations that form unique keys 	<ul style="list-style-type: none"> multivariate statistical measures (e.g., correlation, principal components)

Table 1. Data Space Descriptors for tabular data.

a temperature measurement) or its valid value range (e.g., values in degrees Celsius cannot be lower than $-273, 15^{\circ}\text{C}$).

Spatial/temporal continuity: If the data context provides a spatial and/or temporal frame of reference for the data content, certain continuities among the data values of a variable may emerge. For example, we can objectively evaluate if the data values of a variable are monotonically increasing/decreasing with time or spatial distance and express this as a descriptor of that variable.

Variable combinations that form unique keys: Besides the variables, which are explicitly denoted to be keys or IDs, a dataset may contain other combinations of variables that uniquely identify the records. Mechanisms, such as *primary key analysis* (Borek et al. 2011), can be used to find such alternative identifiers. Yet often the most interesting case is when a variable combination that is expected to be unique, turns out not to be, indicating inconsistencies in the dataset.

These descriptors are typically associated with a single granularity. Yet, this association is not necessarily exclusive. For example, univariate statistical measures, such as min/max/median, can be either descriptors of the variable (column) for which they have been computed, or descriptors of the individual value (cell) that constitutes the identified min/max/median. Furthermore, there exist descriptors that can be applied to all granularities and which we did not place in Table 1 for this reason. Important examples include:

Number/size: Everything in a dataset can be counted or measured in terms of its memory footprint. This can be of interest as an information about the data, but also point to errors in the dataset.

Duplications: Certain values and records, but also entire variables or datasets can be identical or close to identical, which can make them of higher or lower interest to a user. While a numerical value appearing twice somewhere within the dataset does not seem like a notable occurrence, this is certainly different if the value is a person's name or a supposedly unique identifier.

Inconsistencies/mismatches: Data descriptors can be erroneous as well. For example, the dataset may have changed since the descriptors were stored, or the descriptors describe the dataset in a prototypical ideal way,

but the actual data is messy and incomplete. In both cases, it is important to check if the descriptors match the data and to annotate those parts of the data that do not.

Out of this collection of descriptors for tabular data, only a few can be stored together with the data in the common file and database formats. The CSV format typically contains the variable name and sometimes also the scale type in the file header. Whereas relational databases keep these descriptors together with additional schema information in separate tables. Other than these basic descriptors are rarely given in the standard storage formats, even though it would be useful to have advanced descriptors available to bootstrap subsequent visual analysis steps. In this situation, we can gather further descriptors given appropriate tool support.

Gathering Data Descriptors for Tabular Data

When gathering data descriptors for variables, records, and the entire dataset, dependencies between them have to be taken into account and – if possible – to be resolved. This can hardly be achieved in a purely autonomous preprocess that runs without user intervention. Instead, this gathering process requires a well-defined workflow that combines computation for those descriptors that can be automatically derived with user interaction for those that rely on the background knowledge of the user. Such a semi-automated gathering of data descriptors poses the challenge that automated computations must be confined to acceptable runtimes so that users do not get frustrated waiting for intermediary prompts for their input.

To address these points, we propose a series of guidelines for gathering descriptors:

- The gathering process should **follow a step-wise gathering procedure**. This allows for running the individual gathering steps in a configurable order. That order can be defined so that it minimizes repeated calculations and user inputs. Cyclic dependencies between the data granularities – e.g., certain record descriptors requiring column descriptors and vice versa – can be resolved by breaking up the gathering of descriptors for a granularity into multiple steps.

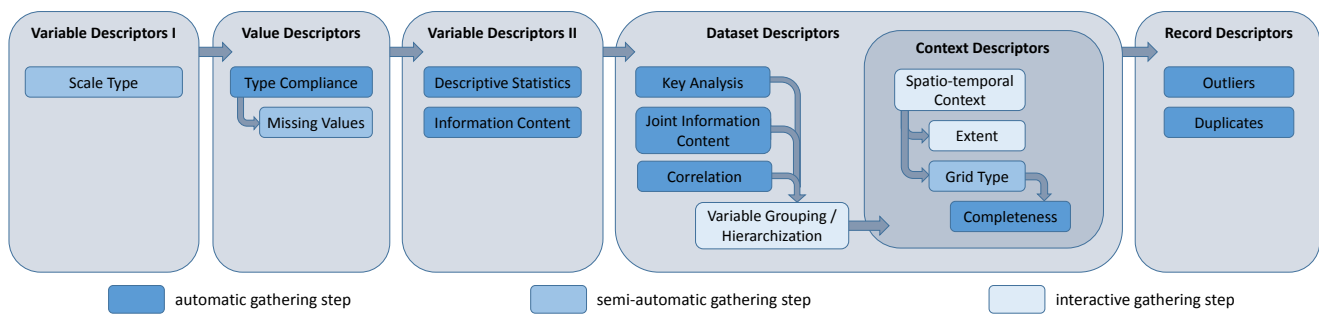


Figure 2. Steps for gathering descriptors from tabular data shown in the order in which our approach determines them. Arrows indicate dependencies between these steps and different shades of blue denote different levels of interactivity.

- The gathering process should allow for **different degrees of interaction**. By also providing the possibility for user inputs and manual adjustments, the process can on one hand adapt to “dirtier” data with many inconsistencies that might not lend itself to automatic descriptor acquisition. On the other hand, it can also benefit from the users’ additional background knowledge about the data. Note that together with a step-wise procedure, the degree of interaction can be different for the individual gathering steps.
- The gathering process should adhere to given **time frames / speed constraints**. As the computation of some descriptors is computationally expensive, it is important to be able to limit the necessary amount of time as needed. This limit can be observed by first querying all descriptors that are already available from the dataset itself, then computing only the very essential descriptors that are still missing, before finally gathering more “advanced” descriptors – e.g., running an automatic key analysis on a large number of variables where many possible variable combinations have to be checked. Due to the step-wise gathering procedure, the gathering can be stopped at each step along the process, when a predefined time limit is reached.

We have developed a software tool for gathering descriptors from tabular data that implements these guidelines. It is driven by the idea of providing explicit information about the dataset that allows for gaining first insights and deciding on visualization possibilities. In this sense, it is conceptually different from existing software tools for assessing the data quality and improving it through data cleaning. This is often termed *data wrangling* (Kandel et al. 2011a) and a number of software tools have been developed to help with it – e.g., AJAX (Galhardas et al. 2000), Potter’s wheel (Raman and Hellerstein 2001), Wrangler (Kandel et al. 2011b), or Profiler (Kandel et al. 2012). These tools aim to produce data that is consistently formatted and sufficiently complete for subsequent visual analysis steps. Whereas, our tool aims at gathering information about the data that can be used in the subsequent visual analysis to decide what to view (selection of interest) and how to view it (parametrization of the representation), as it is discussed in a later section. The procedure used by our tool is shown in Figure 2 and the corresponding user interface is depicted in the screenshot in Figure 3. It adopts the step-wise gathering approach, which is detailed in the following.

Variable Descriptors I. First and foremost, the process gathers information about the data types stored in each column as variable descriptors. This is the most basic information to gather, as it does not require any other information about the data. Our acquisition algorithm runs over all values per column and determines the type of the majority of entries using heuristics that match digits, delimiters, and characters and assigns fitting data types. Note that this step cannot be fully automated, as the algorithm can discern nominal variables from discrete numerical variables, but cannot detect ordinal data types. If there exists an ordering among nominal data values, it needs to be interactively specified by the user who has the appropriate domain knowledge and can thus redefine the variable into ordinal data type. This makes the gathering of the data type descriptor a semi-automatic step. As this first block of variable descriptors consists only of this single gathering step, our tool combines its interface (Figure 3A) with the interface for the following gathering steps of value descriptors (Figure 3B).

Value Descriptors. Once the type information is known, we can gather value descriptors. In particular, we aim to describe type compliance and missing values. The former can be done automatically by checking against the data type having been gathered for each variable in the previous step. If a value is not compliant, a corresponding descriptor will be added to that value’s table cell.

Determining the missing values requires some user input, as “missing” does not necessarily mean that the cell is empty, but it could also hold a placeholder value that is out of range, such as `UINT_MAX (4294967295, 0xffffffff)`. When such a placeholder value exists and it is type compliant, which we check first, only the users with their background knowledge about the data can decide whether this is a realistic value or a placeholder for a missing value. This makes this gathering step a semi-automatic one.

Variable Descriptors II. In this next step, we gather more variable descriptors, such as descriptive statistics and each variable’s information content. For these descriptors to be meaningful, it was important to mark down the non-compliant and missing values first, so that, for example, a `UINT_MAX` placeholder does not skew the computation of extreme values and distributions. These descriptive statistics can be computed automatically, since the variable types are known from the very first gathering step and appropriate statistics can thus be computed – e.g., the mean for

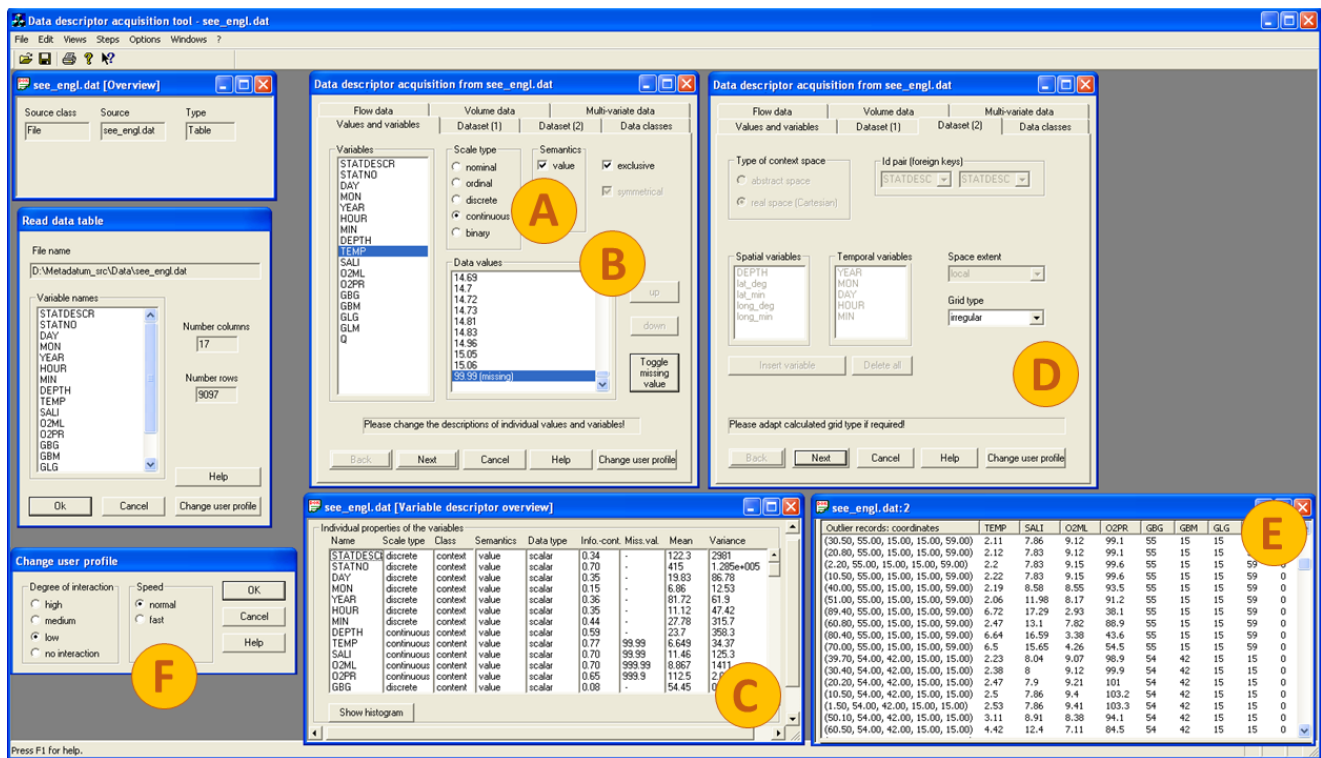


Figure 3. Screenshot from our data descriptor gathering tool for a tabular example dataset containing marine measurements conducted by the Leibniz Institute for Baltic Sea Research. Besides the generic dialogs for loading a dataset and getting a first overview in the top left, the screenshot features the interfaces for each of the individual gathering steps outlined in Figure 2: A – Variable Descriptors I (data type), B – Value Descriptors (type compliance, missing values), C – Variable Descriptors II (descriptive statistics, information content), D – Dataset Descriptors (variable groupings, spatial/temporal context, extent, grid type), E – Record Descriptors (outliers, duplicates). In addition, F shows the dialog that allows for adjusting the computational profile.

continuous numerical variables, the median for discrete numerical variables, and the most frequent term for nominal variables. As these statistical metrics are required by many of the following gathering steps, we placed their computation as early as possible in the process.

On top of these standard measures, we also determine each variable’s information content by computing the Shannon entropy over the set of all values per variable. These can be used to identify and eliminate variables that contain a constant value throughout, which is surprisingly common in practice. Both, statistical descriptors and information content are displayed for the users to inspect in Figure 3C. Note that for convenience reasons, this table also contains the data type from the first gathering step.

Dataset Descriptors. In this step, the first three descriptors have the common goal of providing information for discriminating between data context and data content. To this end, we run a primary key analysis, determine the joint information content of n -tuples of variables, and compute the bivariate Pearson correlation between all pairs of variables. These descriptors are then presented to the users who interactively specify data context and data content based on them. They can furthermore define hierarchies or groupings of associated variables, as it is suggested by Robertson (1990) – e.g., grouping two variables “first name” and “surname” into “name”.

Once specified, we gather additional descriptors for the data context. At first, the users are asked to interactively specify those context variables that denote a spatial and/or

temporal frame of reference. In addition, for each such reference, the users can define its extent – i.e., point, local, or global – which is important for a subsequent visualization, as it tells whether it is possible to interpolate between the data or not. The procedure then tries to establish the data’s grid type (structured or unstructured) by checking the spatial and temporal references for equidistance across records. Since this heuristic is easily misled by a single outlier or undetected header row, we follow a semi-automatic approach and present its results to the users for validation and correction. Lastly, if we have determined a regularly structured grid, we automatically check for completeness – i.e., whether there is a corresponding data record for all possible grid positions. While completeness would be a record descriptor by the data granularity it describes, it still is treated as a dataset descriptor, because the records it describes are missing and thus cannot be marked as such. This is different from missing individual values, as they were determined in the beginning, for which an empty table cell exists to which to attach the corresponding descriptor. The interface for these descriptors is spread over two tabs (see Figure 3D) – the first tab accommodates the descriptors that discern between data context and content, the second tab holds the dataset’s context descriptors.

Record Descriptors. This last step gathers information about outlier records and duplicate records in the dataset. Outlier records are determined automatically by statistical means (Blommestijn and Peerbolte 2012). As for duplicates, we automatically check for so-called *inconsistent duplicates*

that contain different information, but refer to the same entity – i.e., the same customer listed twice under different addresses (Oliveira et al. 2005). Figure 3E shows found outliers in a table view for the user to inspect.

In the standard configuration of our tool, we go through these steps trying to automatically gather as many descriptors as possible, while asking the user for input only as much as necessary. In accordance with our guidelines, we also provide other degrees of interaction from “no interaction” (automatically compute as many descriptors as possible and leave out the rest) to a “high degree of interaction” (report all automatically derived descriptors to the user for validation and readjustment). Furthermore, we also have two different computational profiles – “normal” and “fast”. The “normal” mode follows the default prioritization of automated gathering through computation for as many as descriptors as possible and an interactive gathering for the rest. In contrast, the “fast” mode tries first to query stored descriptors from the data source itself. For those descriptors that are not available from the data source, the system automatically performs computationally inexpensive gathering steps (e.g., descriptive statistics or correlations) and asks the user for input on computationally expensive ones (e.g., key analysis). Note that the fast mode comes with the price of possibly having incorrect descriptors, as the ones accompanying the dataset may be outdated and the ones entered by an inexperienced user may be incorrect. Hence, the fast mode is best used by an experienced user on trustworthy data sources that are known to provide valid descriptors. Both parametrizations of the gathering process are shown in Figure 3F.

Our software for gathering data descriptors was designed to be a general-purpose tool for tabular data about which nothing more than its tabular nature is known or assumed. The tool can be extended to discern more specific data types that follow known formats and value ranges, which can be exploited for their detection. For example, country codes could easily be detected and used accordingly in a geographic mapping, as it is done by recent versions of Tableau and MS Excel.

Leveraging Data Descriptors for Visualizing Tabular Data

On one hand, the gathered data descriptors can be visualized themselves to graphically communicate high-level information about the dataset, which is mainly used for subset selection. For example, dos Santos and Brodlić (2004) introduced visual displays for certain data descriptors, which were designed for providing easier access to the data filtering step in the visualization preprocess. Their *interaction graph* and *n-dimensional window* give an overview over the dimensionality of the attribute space – i.e., the data content. They allow for selecting subspaces of interest (reducing the dimensionality – i.e., columns) or variable ranges of interest (reducing the number of data records – i.e., rows), respectively. Other instances of descriptor visualizations include GeoVISTA’s display of the maximum conditional entropy and correlation values between data content dimensions (MacEachren 2003), as well as the

arrangement of data dimensions in a way that allows for exploring their interdependencies by Yang et al. (2007).

On the other hand, the description can be used to suitably parametrize visualizations of the described dataset. It is noteworthy that to this end, data descriptors are at least implicitly already part of each visualization system, as without knowledge about the distribution of data values no meaningful color-coding and no sensible scaling of axes is possible. The concept of data descriptors gives these already existing means of describing a dataset a formal framework and advocates for dedicated mechanisms for gathering and managing them. So it is not a question of whether to use data descriptors or not, but how to use (and re-use) them.

Figure 4 shows the marine dataset that was already used in Figure 3, as it would be depicted in the absence of any further information about it. By default, the coordinate axes are sorted alphabetically by variable name, which is certainly not optimal, but it at least allows users to quickly seek out a variable of interest. Ideally, one would want to see related variables placed on axes close to each other, so as to ease their combined inspection. One can furthermore observe some unreasonably high values, such as water temperatures (TEMP) of around the boiling point of 100°C or salinity measurements (SALI) around 100%. While these values can be easily spotted and identified as being invalid, probably placeholders for missing values, they nevertheless distort the axes and reduce the axis resolution for valid values. For example, the majority of values on the temperature and salinity axes are now being compressed into the lower part and are hardly discernible.

After gathering data descriptors, as depicted in Figure 3, we can leverage this additional information about the data to alleviate these problems. To first establish a sense of plausibility for a given dataset, we can use a network diagram depicting a bivariate correlation network of the dataset’s variables. This diagram, shown in Figure 5, exhibits a group of correlated data content variables on the right side of the figure. The correlations between salinity (SALI), relative oxygen levels (O2PR), absolute oxygen levels (O2ML), and water temperature (TEMP) are to be expected in a marine dataset. If they would not show up, the source of the dataset should be questioned. The credibility of the dataset is further underlined by the correlation between water temperature and month of the year (MON), which reflects a well-known seasonal pattern in the Baltic Sea. Merely the correlation between month, water temperature, and the numerical ID of the measurement station (STATNO) is an artifact that either happened by chance or is due to a particular numbering scheme for these stations that we are not aware of.

We can further use the bivariate correlations to improve the parallel coordinate plot from Figure 4 by auto-adjusting the order of the axes, so that correlated variables are placed in each other’s proximity (New 2009, ch.4). We can furthermore use identified placeholders for missing values to map them onto separate positions, so that the axes do not get distorted by them. Figure 6 shows the outcome of these descriptor-driven adaptations. The group of correlated data context variables from the left side of Figure 5 is clearly being placed together at the very left of the parallel coordinates, whereas the other observed group of correlated variables gets placed at the right side. The missing values

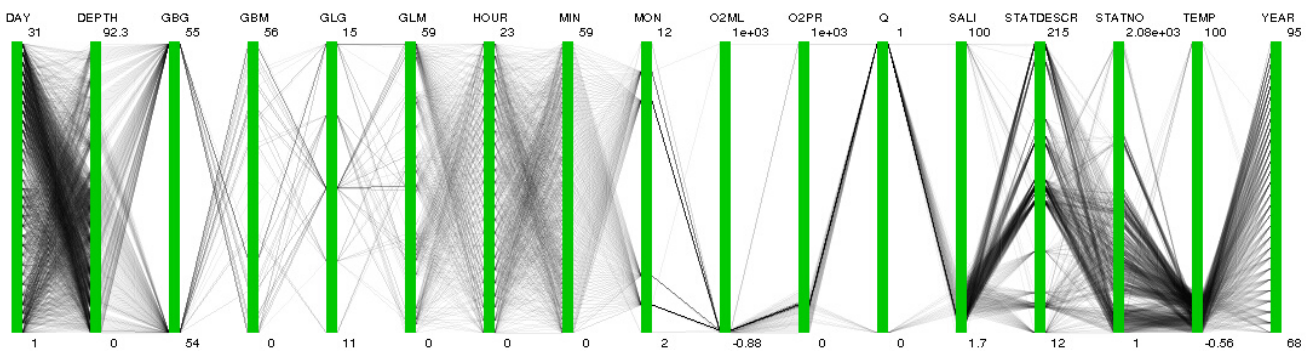


Figure 4. Parallel Coordinates view of the raw marine dataset from Figure 3 without taking descriptors into account.

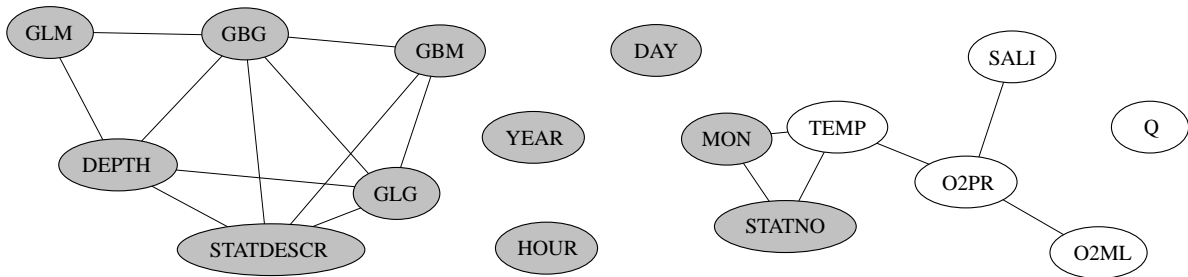


Figure 5. Network diagram of the unsigned bivariate correlation between different variables from the marine dataset described in Figure 3. Each node represents a variable, where gray nodes indicate data context variables and white nodes indicate data content variables. Links between two nodes denote a correlation of at least 0.15 – a threshold that is necessary to set, as all variables are minimally correlated, which would result in a fully connected and thus meaningless display.

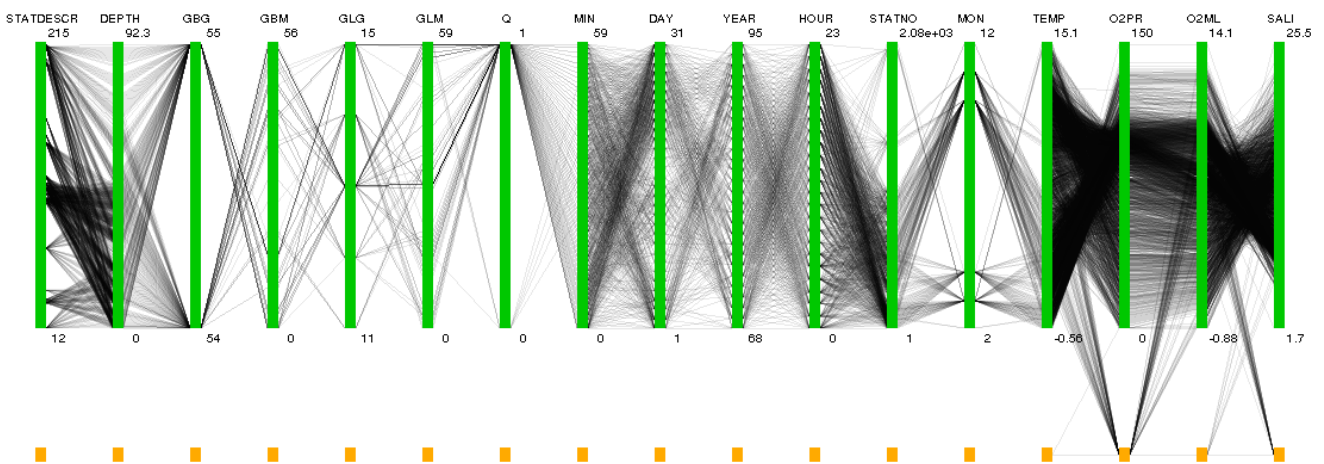


Figure 6. Improved parallel coordinates view from Figure 5. The axes have been re-ordered so that highly correlated variables are positioned close together. The orange marks below each axis single out the placeholder value for missing data.

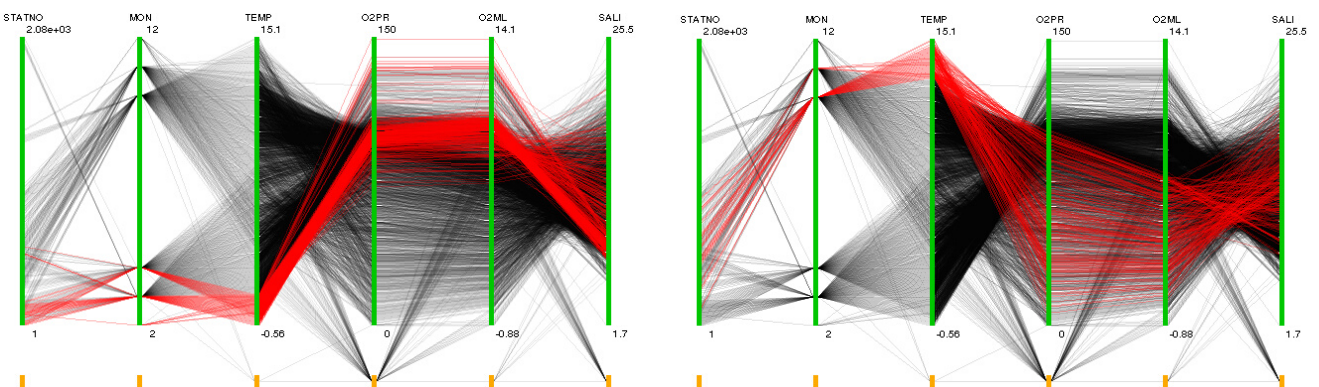


Figure 7. Looking at the variables from the six rightmost axes in Figure 6, we can investigate the correlations via brushing. The left figure highlights the low temperature records and we can observe that at these temperatures, we get mainly high oxygen levels and low salinity. Whereas in the right figure, the highlighted high temperature records exhibit the reverse relations.

are singled out below the axes and marked in orange. This does not only improve the visibility of the legitimate values as these are now spread across each axis, but also eases the application of interactive means, such as brushing. Figure 7 illustrates the correlations among a subset of variables by brushing low and high temperatures, respectively.

Data Descriptors for Climate Impact Research

Many data-intensive scientific domains have established domain-specific data descriptors for their datasets – for example, ecology (Jones et al. 2006), astronomy (Hurt et al. 2007), or systems biology (Stanford et al. 2015). Dedicated systems, so-called *metadata repositories*, that serve the purpose of storing, managing, and querying not the datasets themselves, but merely their DFDs and DSDs at dataset granularity have emerged from these fields (Jones et al. 2001; Berkley et al. 2009; Xiao et al. 2015).

For this article, we have chosen the field of climate impact research to instantiate our set of tabular data descriptors. It presents a most challenging scenario for data descriptors, as it encompasses measurement and simulation data from a multitude of different disciplines, including meteorology, climatology, ecology, agriculture, hydrology, economy, and sociology. Hence, the domain-specific requirements for data descriptors in each of these fields are to some level reflected in those established for climate impact research. Furthermore, working with data from such a multitude of disciplines makes it necessary to meticulously keep track of the different datasets and their various revisions for their cross-disciplinary analysis. Hence, it is not surprising that the field of climate impact research has already developed a number of different notations and standards for data descriptors. While having such standards is a promising direction in theory, in practice climate impact researchers have to deal with a number of caveats that more often than not prevent the use of given data descriptors. These include:

- **Competing standards:** The different standards for data descriptors each describe some data aspects and leave out others. Yet they are not complementary to each other, as they overlap in some parts and are disjoint in others.
- **Different versions of a standard:** Standards evolve to capture notions that arise over time. This creates incompatibilities among different versions of the same standard.
- **Flexibility in interpreting a standard:** The standards leave some flexibility to their realization, which is why two software tools that officially support the same standard may not be able to use each other's descriptors.
- **Incomplete implementations of the standard:** Standards are rarely implemented in full and most software tools work with some sensible subsets that are useful in their context, but hardly match across tools.
- **Standardized descriptors still require validation:** Even if standards are fully implemented and a full set of data descriptors is available for a dataset, that does not mean that the provided descriptors are correct.

This shows that by having such standards, the data description and the use of data descriptors can hardly be automated in the background. While necessary to keep track of the data, it becomes another aspect of the data to which the user has to attend. To do so, the user must have basic knowledge about them in order to resolve conflicting descriptors or to identify implausible ones. This section gives a list of data flow and data space descriptors that are specific for the field of climate impact research, including their availability in the most important standards. To aid the climate researchers in maintaining consistent and valid descriptors for their data despite the challenges outlined above, this section furthermore introduces a software module for gathering data descriptors within a climate data visualization support tool.

Data Flow Descriptors for Climate Data

We have catalogued data flow descriptors that are specifically suitable for climate-related data in Table 2 (top). In general, the rather generic ISO 19115-1 geodata standardization already includes data flow descriptors, such as *provenance information* (e.g., evaluation method for quality assessment), *storage information* (e.g., format, recommended decompression algorithm), and *utility information* (e.g., purpose). This standard is a good fit for geospatial data in general, but it does not explicitly support climate-specific descriptors. With explicitly, we mean that while the standard provides places to put climate-specific information in textual form for the user, there are no dedicated fields for this information that make it available to a visual analysis tool in the sense of our definition of data descriptors.

In particular data provenance descriptors play an important role in climate impact research, as the data can originate from a wealth of sources, including measurements, simulations, or further postprocessing steps, as well as from a variety of disparate fields that all contribute to climate impact research. This is reflected by the NetCDF-CF convention, where “CF” stands for “Climate Forecast” and denotes an extension to the NetCDF standard that is developed by the University Corporation for Atmospheric Research (UCAR). It includes a number of more specific data provenance descriptors, as indicated in the top part of Table 2. For data originating from simulations, such climate-specific descriptors include information about the climate model, such as type and version, and about the simulation run on that model, such as the used driver, which are of major importance to assess and reproduce climate simulations. For measured data, this includes information about the measurement device (accuracy, precision, resolution, and sensitivity) and other acquisition information when conducting weather observations, or collecting data from paleo-climatic ice cores or flowstones.

Data Space Descriptors for Climate Data

Data space descriptors that are relevant for the visual analysis of climate-related data are given in Table 2 (bottom). Here it is the data context that is best covered by existing standards. This is not surprising, as climate researchers measure and simulate very different aspects and processes, but always in the same geophysical space. Hence, there is a consensus

Data Provenance		Data Storage	Data Utility
NetCDF-CF: <ul style="list-style-type: none"> institution, date name of the climate (impact) model simulation experiment type (e.g., Monte-Carlo) data generation workflow / operator sequence Not standardized: <ul style="list-style-type: none"> author 		NetCDF-CF: <ul style="list-style-type: none"> climate/model-specific conventions Not standardized: <ul style="list-style-type: none"> storage format (e.g., GRIB, ASCII, binary) data partitioning scheme (e.g., all data in one file, each time step in a separate file) 	Not standardized: <ul style="list-style-type: none"> kind of analyses the dataset can be used for (e.g., hydrology simulations, storm track analysis)

Data Space Descriptors Granularity	Data Context	Data Content
Value	NetCDF-U: <ul style="list-style-type: none"> dating uncertainty (e.g., age dating for ice cores or flowstones) 	NetCDF-CF: <ul style="list-style-type: none"> domain-specific missing values NetCDF-U: <ul style="list-style-type: none"> domain-specific value uncertainties (e.g., of the emission scenario, global/regional climate model, impact model)
Record	NetCDF: <ul style="list-style-type: none"> grid values restricted to certain regions (e.g. ocean only, land surface only, with or without Greenland/Arctic) NetCDF-CF: <ul style="list-style-type: none"> measurement station measurement position change (e.g., balloon or ship) NetCDF-U: <ul style="list-style-type: none"> uncertainty 	NetCDF-CF: <ul style="list-style-type: none"> meteorological / climatic features: <ul style="list-style-type: none"> 1D: centers of pressure systems Not standardized: <ul style="list-style-type: none"> meteorological / climatic features: <ul style="list-style-type: none"> 2D: storm tracks, weather fronts, jet stream 3D: clouds, dust, circulation patterns
Variable	NetCDF: <ul style="list-style-type: none"> dimensions describing simulation ensemble factors NetCDF-CF: <ul style="list-style-type: none"> domain-specific type of spatial dimensions (longitude, latitude, pressure level) values defined for centers, edges, or vertices of grid cells? variable-specific properties of time (e.g., different time steps) kind of coordinate reference system (e.g., geographical, projected) and associated properties (e.g., rotated pole, conformal, equidistant) 	NetCDF-CF: <ul style="list-style-type: none"> holds certain climate quantity (e.g., temperature 2m above sea level, precipitation including snow and/or hail) Not standardized: <ul style="list-style-type: none"> fit of simulated variable distributions to reference data (e.g., measurements or reanalysis data) -> bias measures trends of mean values, variability, and extremes periodicities and time-delayed correlations
Dataset	NetCDF: <ul style="list-style-type: none"> number and kind of homogeneous subsets (e.g., earth surface and 3D atmospheric variable sets) NetCDF-CF: <ul style="list-style-type: none"> geospatial extent (global, regional, urban) 	NetCDF-CF: <ul style="list-style-type: none"> relation between variables (e.g., thickness and temperature of sea ice)

Table 2. Data descriptors for climate data as they are supported by the various standards.

about what climate researchers want to describe about the data context and this consensus has been formalized through standards. Whereas for the data content, possible descriptors are much more diverse and thus their set is much less standardized.

A few generic data context descriptors are part of the general NetCDF convention, such as the information about masked data – i.e., data that is only available for certain regions, such as land or sea area. Others that are more specific are captured by the NetCDF-CF extension. For example, for representing the dynamics of oceans, atmosphere, and ice shields, physical variables are provided as sub-models in different spatial dimensionalities (1D, 2D, 3D), with partly different, linked grid structures and varying temporal granularities (see, e.g., Petoukhov et al. (2000)) that define the *data context*. For a meaningful visualization of climate data, these structures and dependencies need to be known. In addition, typical *data content* descriptors, such as climate-related regions of interest in the data, range from centers of pressure systems (Wong et al. 2000), to weather fronts and storm tracks (Moorhead and Zhu 1993), and even to the 3D tracking of clouds, dust, and atmospheric pollutants (Ma and Smith 1993). Descriptors for paleo-climate analysis include periodicity and time-delayed correlations. While uncertainty information can be explicitly stored using the

NetCDF-U extension, in practice this is rarely done and it is more common to include an additional variable representing the uncertainty information.

Gathering Data Descriptors for Climate Data

In the previous section on tabular data, we concerned ourselves mostly with the computation and interactive specification of data descriptors as these are usually not provided alongside the data. For data from the domain of climate impact research, the situation is quite different, as the sections on climate-specific data descriptors and their standards have illustrated: a range of descriptors are usually already given – yet, they are possibly incomplete, incompatible, or inconsistent. Thus a gathering of data descriptors in climate impact research focuses on retrieving or querying those existing data descriptors, computing those that are missing or do not match the data, and converting them into the right standard and version for the visual analysis tool to be subsequently used.

Precisely for this purpose, we have developed a data descriptor module within a climate data visualization support tool (Nocke et al. 2007). This module helps climate scientists to bridge the gap between the data descriptors that are given and those that should be given for a subsequent visual analysis. In a first step, it *queries* the descriptors stored

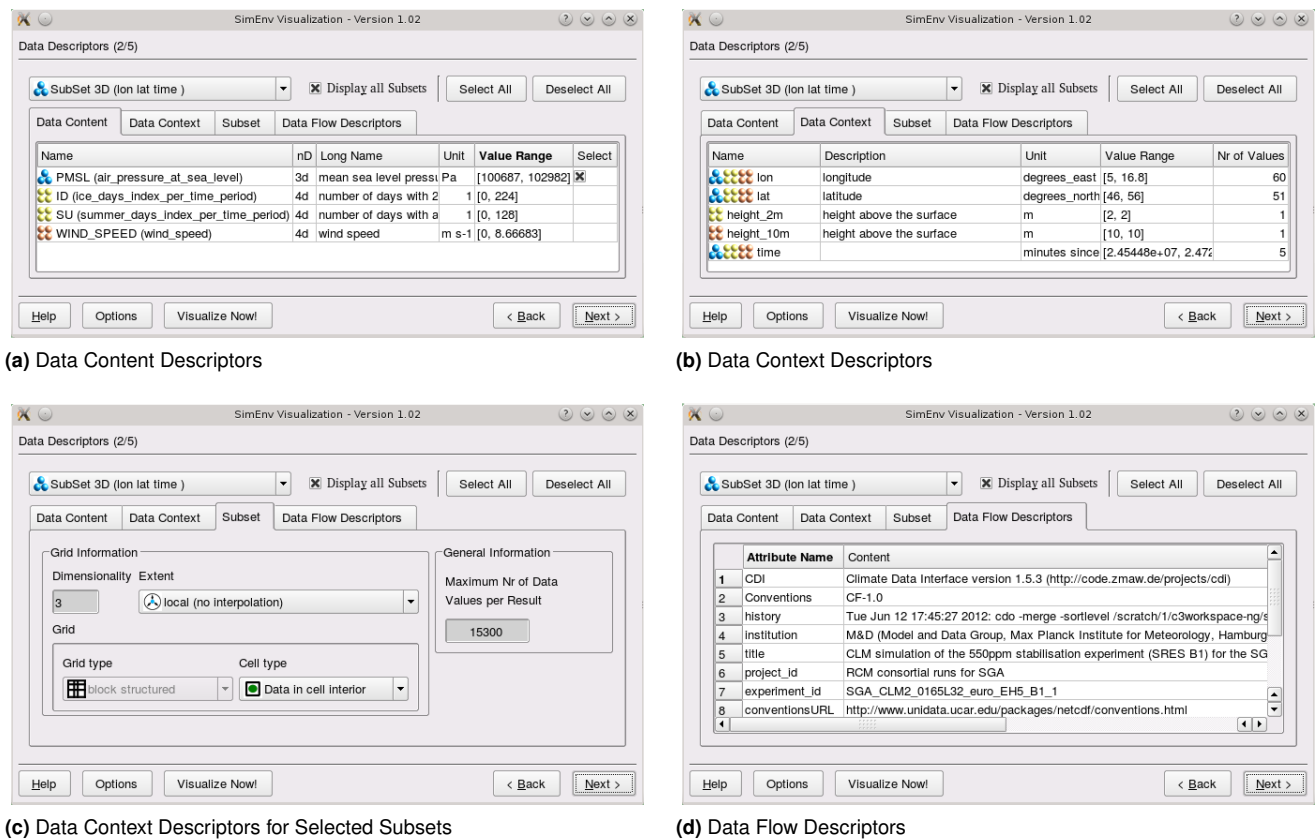


Figure 8. The gathering module allows a detailed examination of the descriptors for the data content (a), the data context (b), gridded subsets of the context (c), and the data flow (d). The shown example is from a dataset generated by a regional climate model simulation with the CLM model.

with the NetCDF/NetCDF-CF data to obtain any descriptors already provided by the data. If they are incomplete, it *derives* additional data descriptors from the dataset in a second step. Finally, in a third step, it presents the queried and derived descriptors to the user in order to prompt for *input* for the interactive adaptation of those that have been gathered and for the completion of those that cannot be computed automatically.

The first step queries data descriptors from the NetCDF data description. NetCDF is a particularly good source for dataset descriptors, as many descriptors are mandatory by the NetCDF convention and must be given. This includes information about the grid’s dimensionality and structure, which allows us to automatically distinguish between data context variables (i.e., longitude, latitude, and time) and data content variables (i.e., the measured or simulated values) as well as to relate different grids / ensemble members.

The second step gathers optional data descriptors that were not given with the data. For example, often missing within NetCDF are the data extent (point, local, global) or a variable’s unit. Where possible, our tool tries to infer these from domain knowledge and assumes, for example, that “lightning” is point information and that a variable “temperature” has the unit “Kelvin”. It also tries to find and fix common spelling errors and replaces, for example, the unknown identifier “units” with the standard-conform identifier “unit”. Only if no descriptor is found – neither directly nor by the described inference – our tool aims to gather it using the pipeline from Figure 2.

The third step presents the results of the previous two steps to the user, as it is illustrated in Figure 8. The shown interface provides an overview of the data content variables, each being assigned a colored glyph that denotes the context, which provides the frame of reference for that variable (Figure 8a). The number of nodes that comprise these glyphs indicate the dimensionality of the corresponding data context, with differently colored glyphs marking different contexts. In the shown example, one can see two 4-dimensional data contexts colored yellow for ice days and summer days, and orange for wind speed. To find out concretely which variables comprise these contexts, one can switch to the data context tab (Figure 8b). It becomes evident that the two 4-dimensional contexts are both composed of latitude, longitude, and time, and that they only differ in the included altitude variable. One can further notice that these altitude variables are somewhat peculiar, as the “Nr of Values” descriptor shows that they only contain a single data value each: 2 meters and 10 meters, respectively, as a quick glance at the “Value Range” descriptor reveals. So, the altitude is in both cases not a variable (as there is no variability), but a constant that simply means that wind speeds were simulated at an altitude of 10 meters, while ice days and summer days were defined at 2 meters. Trying to visualize the data in this “pseudo” 4-dimensional form would most certainly lead to an ill-configured 3D visualization, as the data is actually flat and should also be visualized as such. The interface further allows the users to inspect context descriptors for selected subsets (Figure 8c) and data flow

descriptors (Figure 8d) and to interactively readjust them if necessary.

The software module has been designed in collaboration with researchers from the Potsdam Institute of Climate Impact Research and is in active use. User interviews have highlighted two principal usage scenarios: before and after the visualization. Using the module before the visualization reflects the idea of an initial analysis step that gives first insight into the data to decide on later analysis steps. Whereas the usage afterwards reflects a debugging process where users found their data misrepresented and are looking for the reason. The above example of the “pseudo” 4-dimensional dataset illustrates how easily a visualization system can misinterpret a dataset and thus how important it is to be able to go back from an improper visualization to inspect and adjust the data descriptors.

Leveraging Data Descriptors for Visualizing Climate-related Data

One of the main problems in climate impact research is neither a lack of data descriptors, nor a lack of standards to convey them, but that they are only selectively and inconsistently used by different visualization tools, as can be seen in Figure 9. The figure shows a dataset that contains the initial state for a COSMO/CLM (CCLM) climate simulation run (Rockel et al. 2008). The depicted variable quantifies the height of the snow cover over Europe, as indicated by the following NetCDF description:

```
float W_SNOW(time, rlat, rlon);
  W_SNOW:standard_name =
    "lwe_thickness_of_surface_snow_amount";
  W_SNOW:long_name =
    "surface snow amount";
  W_SNOW:units = "m";
  W_SNOW:grid_mapping = "rotated_pole";
  W_SNOW:coordinates = "lon lat";
  W_SNOW:_FillValue = -1.e+20f;
```

This description contains valuable information for generating a proper visualization for the dataset: that the values are given in meters is important for providing a legend, that missing entries are indicated by a value of -1.0×10^{20} is helpful for masking these entries in the resulting view, and that the given coordinates are rotated is relevant for a correct spatial mapping to the map or globe. Such a rotation of the coordinate grid is often used to yield evenly-sized grid cells for regions that are not close to the equator, to ensure stable numerical simulation. If some of these descriptors were missing or incorrect – e.g., units or missing values – the researcher would have been able to add them using the gathering module. So, we assume that our data descriptors are correct, complete, and conform to the NetCDF standard, so that a visualization tool should theoretically be able to render a correct view of the snow cover dataset. The grid adjustment for the spatial mapping can be determined from the specifics of the rotation, which are also given in NetCDF:

```
char rotated_pole;
  rotated_pole:grid_mapping_name =
    "rotated_latitude_longitude";
  rotated_pole:grid_north_pole_latitude =
    39.25f;
  rotated_pole:grid_north_pole_longitude =
    -162.f;
float rlon(rlon);
  rlon:axis = "X";
  rlon:standard_name = "grid_longitude";
  rlon:long_name = "rotated longitude";
  rlon:units = "degrees";
float rlat(rlat);
  rlat:axis = "Y";
  rlat:standard_name = "grid_latitude";
  rlat:long_name = "rotated latitude";
  rlat:units = "degrees";
```

Yet, Figure 9 illustrates that given the same dataset and the same data description in NetCDF, different visualization tools generate very different outcomes. We have tested the ability of four different visualization tools that are commonly used in climate research to visualize the NetCDF-described snow cover dataset using their default settings. *Avizo*[†] properly recognizes the rotated grid from the NetCDF descriptors and automatically adjusts for the rotation, so that the data is correctly mapped onto Europe. Yet, in its standard configuration *Avizo* does not correctly mask the missing values, which basically coincide with sea surfaces, and it also does not display the unit of measurement in the legend. Whereas *OpenDX* (Thompson et al. 2004) properly masks the missing values and display the units correctly, but does not correctly translate the spatial mapping to adjust for the rotated grid. This leads to the data clearly showing the outline of Europe being overlaid on the map of Africa. Similarly, *Ferret* (Hankin et al. 1996) also handles the missing values correctly and gives an indication of the unit of measurement, yet neglects the rotation of the coordinates, as can be seen from the latitude/longitude labels on the axes. Finally *NCL*[‡], the NCAR Command Language, can actually leverage all three of the highlighted descriptors – units, missing values, and rotation – using one of its standard example scripts.

To be fair, we have to note that these visualizations were generated with the respective tools using standard parametrizations and default options without any further user intervention. Some of their deficits could be alleviated by further manual fine-tuning. For example, *Avizo* allows users to adjust the color mapping in the *Colormap Editor*, which can be used to mask the missing values by hand and thus to eventually produce a proper visualization of that dataset. Yet even for that, the users themselves must go over the NetCDF code in a text editor to find out which value is used to indicate missing data in order to adjust the colors accordingly. Domain knowledge about which visualization tool supports which kinds of NetCDF descriptors could potentially help to identify tools that are suitable for a dataset at hand in future work.

[†]see <http://www.fei.com/software/avizo3d>

[‡]see <http://www.ncl.ucar.edu>

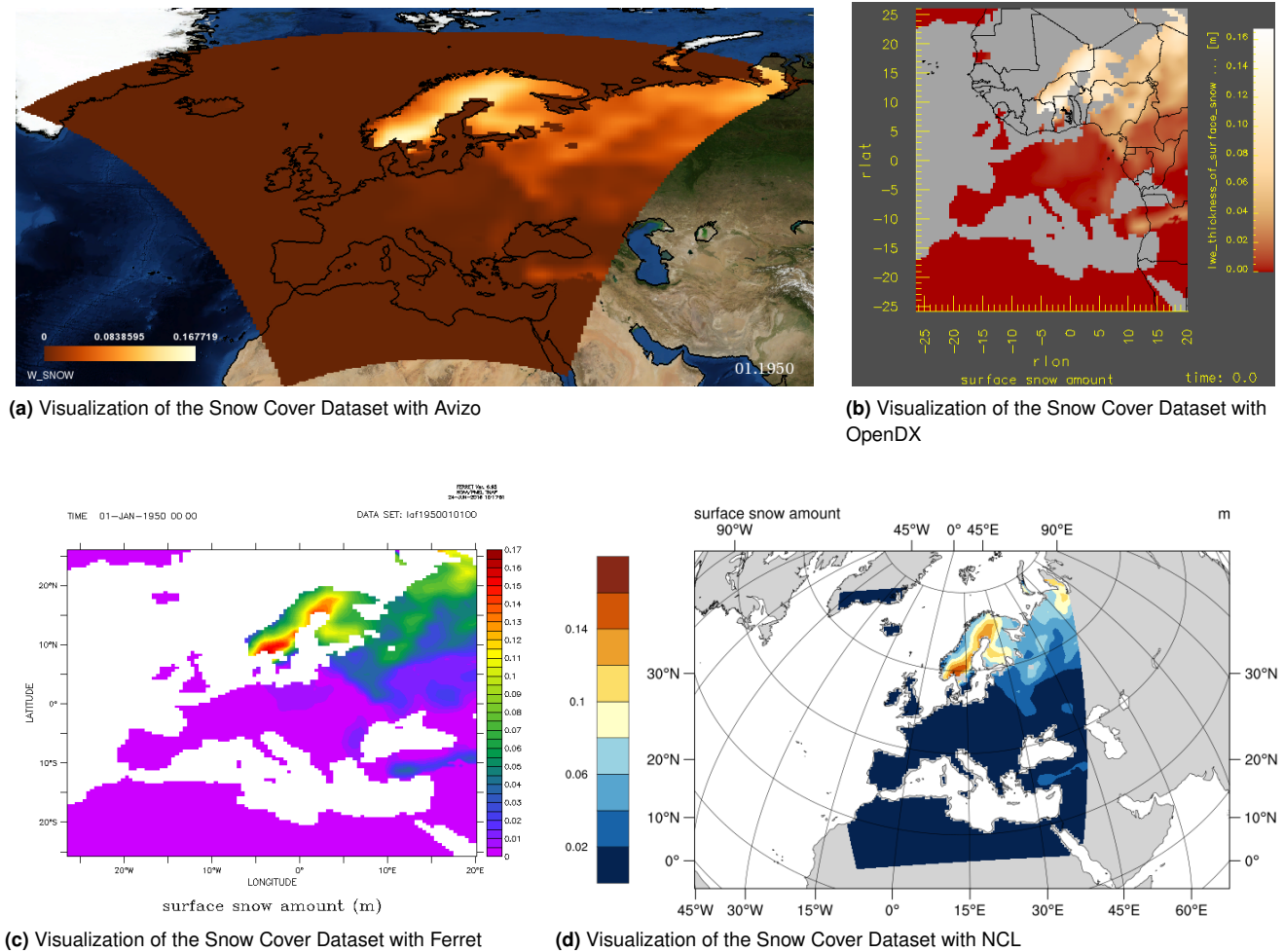


Figure 9. Visualizations of the Snow Cover Dataset with four different standard tools commonly used in climate science. It can be observed that Avizo, OpenDX, and Ferret make only selective use of the given NetCDF descriptors, which results in faulty visualizations. Only NCL was able to actually leverage all of the given descriptors for generating a proper visualization.

Concluding Remarks

With the proposed concept of data descriptors, we do not present yet another taxonomy of data types and structures, but instead a unifying view on those that already exist. From its scope, our concept is more in line with approaches, such as *meta-metadata* (Kerne et al. 2010), that aim to bring together different data descriptor standards under a unified meta standard. Yet, we differ in the path taken towards such a unifying view: While the existing approaches take a top-down perspective and define new standards for researchers to adhere to, our concept takes a bottom-up perspective by filling in the gaps and resolving inconsistencies.

Implications From our own experience in working with climate researchers, this bottom-up approach is the more practical one as it yields concrete results that we can already start using while waiting for the definite data standard to arrive. Yet this makes it also the more involved approach as compared to defining a suitable data description standard and further assuming that any given analysis input adheres to it. It takes computational effort and the involvement of the users with their background knowledge to provide information about data that is

- accurate (matches underlying data),
- complete (covers all relevant data aspects),

- consistent (does not contradict itself),
- current (reflects data changes), and
- conforms to whichever conventions are used by the variety of existing analysis tools.

Our concept of data descriptors lifts this information from a few scattered auxiliary measures to being data entities in themselves. This allows us to centralize the required effort for providing reliable information about data in a dedicated software module – our gathering pipeline. The pipeline serves as the principal access point to this information for users to adjust them and for analysis methods to utilize them.

Limitations Providing data descriptors is only one side of the equation. On the other side, we have no influence over whether and to which degree subsequent visual analysis steps actually use given data descriptors. This was illustrated by the examples in the previous section. To some degree, this lies in the nature of such a bottom-up approach that rather offers descriptors for use, than to enforce their observance. After all, it is hard to determine in general, which descriptors *must* be given and used, and which *can* be given and used – i.e., which are ultimately necessary and which would be helpful, but one could do without. This depends largely on the concrete visual analysis task at hand. While we do not know the analysis task in beforehand, data type (e.g., tabular) and application domain (e.g., climatology)

can help to limit the set of all possible data descriptors to those that are suitable for the data type and typical in the domain. In the same way, as we have tailored our general descriptor framework to tabular data in a first step, and then to the concrete requirements of the climate research domain, adaptations are necessary for other data types and domains. This highlights once again that the descriptors listed in this paper and the pipeline for gathering them are not a one-stop solution for all possible data/domain combinations, but rather a blueprint to be adapted to the specifics of other data types and to be refined for other domains. Providing a set of descriptors that is well adapted to the data and tasks of a particular domain is certainly more likely to be picked up on by the tools and users in that domain, than some unwieldy all-encompassing generic metadata solution.

Generalization It is noteworthy that common visual analysis strategies do not explicitly include such an initial analysis step of gathering information about the data. Following the commonly applied strategies, an analyst can initiate a visual analysis either with an *overview-first step* (Shneiderman 1996) or with an *analyze-first step* (Keim et al. 2006). Yet in both cases, it remains challenging to decide which overview or analysis method shall be invoked, respectively, on which data subset and with which parameter settings. Hence our initial gathering can be understood as a *describe-first step* that precedes overview or analysis and collects information about the data. This information can either be visualized directly so that users get a meta-view of their data, or it can be used indirectly for an informed descriptor-driven selection and parametrization of appropriate computational or visual methods.

Continuation For the future, we anticipate that descriptive information about data will play an increasingly important role in visual data analysis for two reasons. On one hand, the push for Big Data increases the amount of data, which in turn has to be described in a meaningful way to select the right subset for an analysis at hand. On the other hand, more often than not these days, the data provider is different from the data analyst and thus information about datasets needs to be passed on (Patro et al. 2003). As this trend grows with current movements, such as Open Data and Open Science, the visualization and visual analytics community needs to develop concepts and tools to deal with it. We strongly believe that the concept of data descriptors will serve as an important foundation for these developments.

Funding

This work originated from a long-standing collaboration between the University of Rostock and the Potsdam Institute of Climate Impact Research. The authors acknowledge financial support by the Federal Ministry for Education and Research via the Potsdam Research Cluster for Georisk Analysis, Environmental Change and Sustainability (PROGRESS).

References

(2011) Metadata mapper: a web service for mapping data between independent visual analysis components, guided by perceptual rules. In: Wong PC, Park J, Hao MC, Chen C, Börner K, Kao DL and Roberts JC (eds.) *VDA'11: Proceedings of the*

- Conference on Visualization and Data Analysis*. SPIE. ISBN 9780819484055, pp. 78650I–1–13. DOI:10.1117/12.881734.
- Adèr HJ (2008) Phases and initial steps in data analysis. In: Adèr HJ and Mellenbergh GJ (eds.) *Advising on Research Methods: A Consultant's Companion*. Johannes van Kessel Publishing. ISBN 9789079418022, pp. 333–355.
- Andrienko N and Andrienko G (2006) *Exploratory Analysis of Spatial and Temporal Data – A Systematic Approach*. Springer. ISBN 3540259945. DOI:10.1007/3-540-31190-4.
- Angles R and Gutierrez C (2008) Survey of graph database models. *ACM Computing Surveys* 40(1): 1–39. DOI:10.1145/1322432.1322433.
- Arens Y, Hovy EH and Vossers M (1993) On the knowledge underlying multimedia presentations. In: Maybury MT (ed.) *Intelligent Multimedia Interfaces*. AAAI Press. ISBN 0262631504, pp. 280–306.
- Asuncion HU (2011) In situ data provenance capture in spreadsheets. In: *Proceedings of the IEEE Conference on eScience*. IEEE. ISBN 9781457721632, pp. 240–247. DOI: 10.1109/eScience.2011.41.
- Ayyub BM and Klir GJ (2006) *Uncertainty Modeling and Analysis in Engineering and the Sciences*. CRC Press. ISBN 9781584886440. DOI:10.1201/9781420011456.
- Bassiouni MA (1985) Data compression in scientific and statistical databases. *IEEE Transactions on Software Engineering* SE-11(10): 1047–1058. DOI:10.1109/tse.1985.231852.
- Batini C and Scannapieca M (2006) *Data Quality: Concepts, Methodologies and Techniques*. Springer. ISBN 9783540331728. DOI:10.1007/3-540-33173-5.
- Belhajjame K, Wolstencroft K, Corcho O, Oinn T, Tanoh F, William A and Goble C (2008) Metadata management in the Taverna workflow system. In: Priol T, Lefevre L and Buyya R (eds.) *CCGRID'08: Proceedings of the IEEE International Symposium on Cluster Computing and the Grid*. IEEE, pp. 651–656. DOI:10.1109/CCGRID.2008.17.
- Bergeron RD and Grinstein GG (1989) A reference model for the visualisation of multi-dimensional data. In: Hansmann W, Hopgood FRA and Strasser W (eds.) *EG'89: Proceedings of the European Computer Graphics Conference and Exhibition*. Eurographics Association. ISBN 0444880135, pp. 393–399.
- Bergman LD, Rogowitz BE and Treinish LA (1995) A rule-based tool for assisting colormap selection. In: *VIS'95: Proceedings of the IEEE Conference on Visualization*. IEEE. ISBN 0818671874, pp. 118–125. DOI:10.1109/VISUAL.1995.480803.
- Berkley C, Bowers S, Jones MB, Madin JS and Schildhauer M (2009) Improving data discovery for metadata repositories through semantic search. In: Barolli L, Xhafa F and Hsu HH (eds.) *CISIS'09: Proceedings of the International Conference on Complex, Intelligent and Software Intensive Systems*. IEEE, pp. 1152–1159. DOI:10.1109/CISIS.2009.122.
- Bilgic M, Licamele L, Getoor L and Shneiderman B (2006) D-Dupe: An interactive tool for entity resolution in social networks. In: Wong PC and Keim D (eds.) *VAST'06: Proceedings of the IEEE Symposium on Visual Analytics Science and Technology*. IEEE, pp. 43–50. DOI:10.1109/VAST.2006.261429.
- Blommestijn SQ and Peerbolte EAL (2012) Outliers and extreme observations: What are they and how to handle them? In: Adèr HJ and Mellenbergh GJ (eds.) *Advising on Research Methods*

- *Selected Topics 2012*. Johannes van Kessel Publishing. ISBN 9789079418213, pp. 81–105.
- Borek A, Woodall P, Oberhofer M and Parlikad AK (2011) A classification of data quality assessment methods. In: *ICIQ'11: Proceedings of the International Conference on Information Quality*. pp. 189–203.
- Brodlie K, Osorio RA and Lopes A (2012) A review of uncertainty in data visualization. In: Dill J, Earnshaw R, Kasik D, Vince J and Wong PC (eds.) *Expanding the Frontiers of Visual Analytics and Visualization*. Springer. ISBN 9781447128038, pp. 81–109. DOI:10.1007/978-1-4471-2804-5_6.
- Brodlie KW (1992) Visualization techniques. In: Brodlie KW, Carpenter LA, Earnshaw RA, Gallop JR, Hubbard RJ, Mumford AM, Osland CD and Quarendon P (eds.) *Scientific Visualization: Techniques and Applications*. Springer. ISBN 3540545654, pp. 37–85. DOI:10.1007/978-3-642-76942-9_3.
- Brodlie KW and Noor NM (2007) Visualization notations, models and taxonomies. In: Lim IS and Duce D (eds.) *TPCG'07: Proceedings of the Theory and Practice of Computer Graphics Conference*. Eurographics Association, pp. 207–212. DOI: 10.2312/LocalChapterEvents/TPCG/TPCG07/207-212.
- Buneman P, Khanna S and Wang-Chiew T (2001) Why and where: A characterization of data provenance. In: van den Bussche J and Vianu V (eds.) *ICDT'01: Proceedings of the International Conference on Database Theory*, number 1973 in Lecture Notes in Computer Science. Springer. ISBN 9783540414568, pp. 316–330. DOI:10.1007/3-540-44503-X_20.
- Butler DM and Pendley MH (1989) A visualization model based on the mathematics of fiber bundles. *Computers in Physics* 3(5): 45. DOI:10.1063/1.168345.
- Callahan SP, Freire J, Scheidegger CE, Silva CT and Vo HT (2008) Towards provenance-enabling ParaView. In: Freire J, Koop D and Moreau L (eds.) *IPAW'08: Proceedings of the Provenance and Annotation Workshop*, number 5272 in Lecture Notes in Computer Science. Springer. ISBN 9783540899648, pp. 120–127. DOI:10.1007/978-3-540-89965-5_13.
- Cammarano M, Dong X, Chan B, Klingner J, Talbot J, Halevy A and Hanrahan P (2007) Visualization of heterogeneous data. *IEEE Transactions on Visualization and Computer Graphics* 13(6): 1200–1207. DOI:10.1109/TVCG.2007.70617.
- Carata L, Akoush S, Balakrishnan N, Bytheway T, Sohan R, Seltzer M and Hopper A (2014) A primer on provenance. *Communications of the ACM* 57(5): 52–60. DOI:10.1145/2596628.
- Card SK, Mackinlay J and Shneiderman B (1999) *Readings in Information Visualization: Using Vision to Think*. Morgan Kaufmann. ISBN 1558605339.
- Cheng X, Cook D and Hofmann H (2015) Visually exploring missing values in multivariable data using a graphical user interface. *Journal of Statistical Software* 68(6): 1–23. DOI: 10.18637/jss.v068.i06.
- Cleary J, Holmes G, Cunningham SJ and Witten IH (1996) Metadata for database mining. In: *Proceedings of the IEEE Conference on Metadata*.
- Codd EF (1970) A relational model of data for large shared data banks. *Communications of the ACM* 13(6): 377–387. DOI: 10.1145/362384.362685.
- Codd EF (1990) *The Relational Model for Database Management: Version 2*. Addison Wesley. ISBN 0201141922.
- Cohen S, Cohen-Boulakia S and Davidson S (2006) Towards a model of provenance and user views in scientific workflows. In: Leser U, Naumann F and Eckman B (eds.) *DILS'06: Proceedings of the Workshop on Data Integration in the Life Sciences*, number 4075 in Lecture Notes in Computer Science. Springer. ISBN 9783540365938, pp. 264–279. DOI:10.1007/11799511_24.
- Correa CD, Chan YH and Ma KL (2009) A framework for uncertainty-aware visual analytics. In: Stasko J and van Wijk JJ (eds.) *VAST'09: Proceedings of the IEEE Symposium on Visual Analytics Science and Technology*. IEEE. ISBN 9781424452835, pp. 51–58. DOI:10.1109/VAST.2009.5332611.
- Cowley P, Haack J, Littlefield R and Hampson E (2006) Glass Box: Capturing, archiving, and retrieving workstation activities. In: *CARPE'06: Proceedings of the ACM Workshop on Continuous Archival and Retrieval of Personal Experiences*. ACM. ISBN 1595934987, pp. 13–18. DOI:10.1145/1178657.1178662.
- Cowley P, Nowell L and Scholtz J (2005) Glass Box: An instrumented infrastructure for supporting human interaction with information. In: *HICSS'05: Proceedings of the Hawaii Conference on System Sciences*. IEEE. ISBN 0769522688, p. 296c. DOI:10.1109/HICSS.2005.286.
- da Cruz SMS, Campos MLM and Mattoso M (2009) Towards a taxonomy of provenance in scientific workflow management systems. In: Zhang LJ (ed.) *Proceedings of the World Conference on Services*. IEEE. ISBN 9780769537085, pp. 259–266. DOI:10.1109/SERVICES-I.2009.18.
- Dasgupta A, Chen M and Kosara R (2013) Measuring privacy and utility in privacy-preserving visualization. *Computer Graphics Forum* 32(8): 35–47. DOI:10.1111/cgf.12142.
- Davison A (2012) Automated capture of experiment context for easier reproducibility in computational research. *Computing in Science & Engineering* 14(4): 48–56. DOI:10.1109/MCSE.2012.41.
- dos Santos S and Brodlie K (2004) Gaining understanding of multivariate and multidimensional data through visualization. *Computers & Graphics* 28(3): 311–325. DOI:10.1016/j.cag.2004.03.013.
- Drosg M (2009) *Dealing with Uncertainties: A Guide to Error Analysis*. 2nd edition. Springer. ISBN 9783642013836. DOI: 10.1007/978-3-642-01384-3.
- Duval E (2001) Metadata standards: What, who & why. *Journal of Universal Computer Science* 7(7): 591–601. DOI:10.3217/jucs-007-07-0591.
- Fernstad SJ and Glen RC (2014) Visual analysis of missing data – To see what isn't there. In: Chen M, Ebert D and North C (eds.) *VAST'14: Proceedings of the IEEE Conference on Visual Analytics Science and Technology*. IEEE. ISBN 9781479962273, pp. 249–250. DOI:10.1109/VAST.2014.7042514.
- Flöring S (2012) *KnoVA: A Reference Architecture for Knowledge-based Visual Analytics*. PhD Thesis, University Oldenburg.
- Franklin M, Halevy A and Maier D (2005) From databases to dataspace: A new abstraction for information management. *ACM SIGMOD Record* 34(4): 27–33. DOI:10.1145/1107499.1107502.
- Freire J, Koop D, Santos E and Silva C (2008) Provenance for computational tasks: A survey. *Computing in Science & Engineering* 10(3): 11–21. DOI:10.1109/mcse.2008.79.

- Galhardas H, Florescu D, Shasha D and Simon E (2000) AJAX: An extensible data cleaning tool. *ACM SIGMOD Record* 29(2): 590–596. DOI:10.1145/335191.336568.
- Galhardas H, Simon E and Tomasic A (1998) A framework for classifying scientific metadata. In: *Proceedings of the AAAI Workshop on Artificial Intelligence and Information Integration*. ISBN 9781577350675, pp. 106–113.
- Gitelman L (ed.) (2013) *“Raw data” is an oxymoron*. MIT Press. ISBN 9780262518284.
- Glavic B and Dittrich KR (2007) Data provenance: A categorization of existing approaches. In: Kemper A, Schöning H, Rose T, Jarke M, Seidl T, Quix C and Brochhaus C (eds.) *BTW'07: Proceedings of the GI-Fachtagung Datenbanksysteme in Business, Technologie und Web*, number 103 in Lecture Notes in Informatics. Bonner Köllen Verlag. ISBN 9783885791973, pp. 227–241.
- Glavic B, Miller RJ and Alonso G (2013) Using SQL for efficient generation and querying of provenance information. In: Tannen V, Wong L, Libkin L, Fan W, Tan WC and Fourman M (eds.) *In Search of Elegance in the Theory and Practice of Computation, Lecture Notes in Computer Science*, volume 8000. Springer. ISBN 9783642416590, pp. 291–320. DOI: 10.1007/978-3-642-41660-6.16.
- Grammer G, Joshi S, Kroeschel W, Kumar S, Sathi A and Viswanathan M (2012) Obfuscating sensitive data while preserving data usability. United States Patent Application US 20120272329.
- Groth DP and Streefkerk K (2006) Provenance and annotation for visual exploration systems. *IEEE Transactions on Visualization and Computer Graphics* 12(6): 1500–1510. DOI:10.1109/tvcg.2006.101.
- Gschwandtner T, Aigner W, Miksch S, Gärtner J, Kriglstein S, Pohl M and Suchy N (2014) TimeCleanser: A visual analytics approach for data cleansing of time-oriented data. In: Lindstaedt S, Granitzer M and Sack H (eds.) *i-Know'14: Proceedings of the Conference on Knowledge Technologies and Data-driven Business*. ACM. ISBN 9781450327695, pp. 18:1–18:8. DOI:10.1145/2637748.2638423.
- Gschwandtner T, Gärtner J, Aigner W and Miksch S (2012) A taxonomy of dirty time-oriented data. In: Quirchmayr G, Basl J, You I, Xu L and Weippl E (eds.) *CD-ARES'12: Proceedings of the Cross Domain Conference and Workshop on Availability, Reliability, and Security*, number 7465 in Lecture Notes in Computer Science. Springer. ISBN 9783642324970, pp. 58–72. DOI:10.1007/978-3-642-32498-7_5.
- Hahmann S and Burghardt D (2013) How much information is geospatially referenced? networks and cognition. *International Journal of Geographical Information Science* 27(6): 1171–1189. DOI:10.1080/13658816.2012.743664.
- Hankin S, Harrison E, Osborne J, Davison J and O'Brien K (1996) A strategy and a tool, Ferret, for closely integrated visualization and analysis. *The Journal of Visualization and Computer Animation* 7(3): 149–157. DOI:10.1002/(SICI)1099-1778(199607)7:3<149::AID-VIS148>3.0.CO;2-X.
- Heer J, Mackinlay J, Stolte C and Agrawala M (2008) Graphical histories for visualization: Supporting analysis, communication, and evaluation. *IEEE Transactions on Visualization and Computer Graphics* 14(6): 1189–1196. DOI: 10.1109/tvcg.2008.137.
- Hoxmeier JA (2005) Dimensions of database quality. In: Khosrow-Pour M (ed.) *Encyclopedia of Information Science and Technology*. Idea Group. ISBN 9781591405535, pp. 886–891. DOI:10.4018/978-1-59140-553-5.ch155.
- Huq MR, Apers PMG and Wombacher A (2013) Provenance-Curious: A tool to infer data provenance from scripts. In: Paton NW, Guerrini G, Catania B, Castellanos M, Atzeni P, Fraternali P and Gounaris A (eds.) *EDBT'13: Proceedings of the Conference on Extending Database Technology*. ACM. ISBN 9781450315975, pp. 765–768. DOI:10.1145/2452376.2452475.
- Hurt RL, Gauthier AJ, Christensen LL and Wyatt R (2007) Sharing images intelligently: The astronomy visualization metadata standard. In: Christensen LL, Zoulias M and Robson I (eds.) *CAP'07: Proceedings of the Conference on Communicating Astronomy with the Public*. Eugenides Foundation, pp. 450–453.
- Jones MB, Berkley C, Bojilova J and Schildhauer M (2001) Managing scientific metadata. *IEEE Internet Computing* 5(5): 59–68. DOI:10.1109/4236.957896.
- Jones MB, Schildhauer MP, Reichman OJ and Bowers S (2006) The new bioinformatics: Integrating ecological data from the gene to the biosphere. *Annual Review of Ecology, Evolution, and Systematics* 37: 519–544. DOI:10.1146/annurev.ecolsys.37.091305.110031.
- Josko JMB and Ferreira JE (2016) Visualization properties for data quality visual assessment: An exploratory case study. *Information Visualization* DOI:10.1177/1473871616629516. To appear.
- Josko JMB, Oikawa MK and Ferreira JE (2016) A formal taxonomy to improve data defect description. In: Gao H, Kim J and Sakurai Y (eds.) *DASFAA'16: Proceedings of the Workshops on Database Systems for Advanced Applications*. Springer, pp. 307–320. DOI:10.1007/978-3-319-32055-7_25.
- Kandel S, Heer J, Plaisant C, Kennedy J, van Ham F, Riche NH, Weaver C, Lee B, Brodbeck D and Buono P (2011a) Research directions in data wrangling: Visualizations and transformations for usable and credible data. *Information Visualization* 10(4): 271–288. DOI:10.1177/1473871611415994.
- Kandel S, Paepcke A, Hellerstein J and Heer J (2011b) Wrangler: Interactive visual specification of data transformation scripts. In: *CHI'11: Proceedings of the International Conference on Human Factors in Computing Systems*. ACM. ISBN 9781450302289, pp. 3363–3372. DOI:10.1145/1978942.1979444.
- Kandel S, Parikh R, Paepcke A, Hellerstein JM and Heer J (2012) Profiler: Integrated statistical analysis and visualization for data quality assessment. In: Tortora G, Levialdi S and Tucci M (eds.) *AVI'12: Proceedings of the Working Conference on Advanced Visual Interfaces*. ACM. ISBN 9781450312875, pp. 547–554. DOI:10.1145/2254556.2254659.
- Kang H, Getoor L, Shneiderman B, Bilgic M and Licamele L (2008) Interactive entity resolution in relational data: A visual analytic tool and its evaluation. *IEEE Transactions on Visualization and Computer Graphics* 14(5): 999–1014. DOI:10.1109/TVCG.2008.55.
- Kang H, Sehgal V and Getoor L (2007) GeoDDupe: a novel interface for interactive entity resolution in geospatial data. In: Banissi E, Burkhard RA, Grinstein G, Cvek U, Trutschl M, Stuart L, Wyeld TG, Andrienko G, Dykes J, Jern M, Groth

- D and Ursyn A (eds.) *IV'07: Proceedings of the International Conference Information Visualization*. IEEE, pp. 489–496. DOI:10.1109/IV.2007.55.
- Karr AF, Kohnen CN, Oganian A, Reiter JP and Sanil AP (2006) A framework for evaluating the utility of data altered to protect confidentiality. *The American Statistician* 60(3): 224–232. DOI:10.1198/000313006x124640.
- Keim DA, Kohlhammer J, Ellis G and Mansmann F (eds.) (2010) *Mastering the Information Age*. Eurographics Association. ISBN 9783905673777.
- Keim DA, Mansmann F, Schneidewind J and Ziegler H (2006) Challenges in visual data analysis. In: Banissi E, Burkhard RA, Ursyn A, Zhang JJ, Bannatyne M, Maple C, Cowell AJ, Tian GY and Hou M (eds.) *IV'06: Proceedings of the International Conference on Information Visualisation*. IEEE. ISBN 0769526020, pp. 9–16. DOI:10.1109/IV.2006.31.
- Kerne A, Qu Y, Webb AM, Damaraju S, Lupfer N and Mathur A (2010) Meta-metadata: A metadata semantics language for collection representation applications. In: Huang XJ, Jones G, Koudas N, Wu X and Collins-Thompson K (eds.) *CIKM'10: Proceedings of the ACM International Conference on Information and Knowledge Management*. ACM. ISBN 9781450300995, pp. 1129–1138. DOI:10.1145/1871437.1871580.
- Kim W, Choi BJ, Hong EK, Kim SK and Lee D (2003) A taxonomy of dirty data. *Data Mining and Knowledge Discovery* 7(1): 81–99. DOI:10.1023/a:1021564703268.
- Lieberman MD, Taheri S, Guo H, Mirrashed F, Yahav I, Aris A and Shneiderman B (2011) Visual exploration across biomedical databases. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* 8(2): 536–550. DOI:10.1109/tcbb.2010.1.
- Lin S, Fortuna J, Kulkarni C, Stone M and Heer J (2013) Selecting semantically-resonant colors for data visualization. *Computer Graphics Forum* 32(3pt4): 401–410. DOI:10.1111/cgf.12127.
- Ludäscher B, Altintas I, Berkley C, Higgins D, Jaeger E, Jones M, Lee EA, Tao J and Zhao Y (2006) Scientific workflow management and the Kepler system. *Concurrency & Computation: Practice & Experience* 18(10): 1039–1065. DOI:10.1002/cpe.994.
- Ludäscher B, Podhorski N, Altintas I, Bowers S and McPhillips T (2008) From computation models to models of provenance: The RWS approach. *Concurrency and Computation* 20(5): 507–518. DOI:10.1002/cpe.1234.
- Lux M (1998) Level of data – a concept for knowledge discovery in information spaces. In: Banissi E, Khosrowshahi F and Sarfraz M (eds.) *IV'98: Proceedings of the Conference on Information Visualization*. IEEE. ISBN 0818685093, pp. 131–136. DOI:10.1109/IV.1998.694210.
- Ma KL and Smith PJ (1993) Cloud tracing in convection-diffusion system. In: Nielson GM and Bergeron D (eds.) *VIS'93: Proceedings of the IEEE Conference on Visualization*. IEEE. ISBN 0818639407, pp. 253–260. DOI:10.1109/VISUAL.1993.398876.
- MacEachren AM (2003) Exploring high-D spaces with multiform matrices and small multiples. In: Munzner T and North S (eds.) *InfoVis'03: Proceedings of the IEEE Symposium on Information Visualization*. IEEE. ISBN 0780381548, pp. 31–38. DOI:10.1109/INFVIS.2003.1249006.
- Moorhead RJ and Zhu Z (1993) Feature extraction for oceanographic data using a 3D edge operator. In: Nielson GM and Bergeron D (eds.) *VIS'93: Proceedings of the IEEE Conference on Visualization*. IEEE. ISBN 0818639407, pp. 402–405. DOI:10.1109/VISUAL.1993.398901.
- Moreau L, Clifford B, Freire J, Futrelle J, Gil Y, Groth P, Kwasnikowska N, Miles S, Missier P, Myers J, Plale B, Simmhan Y, Stephan E and van den Bussche J (2011) The open provenance model core specification (v1.1). *Future Generation Computer Systems* 27(6): 743–756. DOI:10.1016/j.future.2010.07.005.
- New JR (2009) *Visual Analytics for Relationships in Scientific Data*. PhD Thesis, University of Tennessee, Knoxville.
- Nocke T, Flechsig M and Böhm U (2007) Visual exploration and evaluation of climate-related simulation data. In: Henderson SG, Biller B, Hsieh MH, Shortle J, Tew JD and Barton RR (eds.) *WSC'07: Proceedings of the Winter Simulation Conference*. IEEE. ISBN 1424413060, pp. 703–711. DOI:10.1109/WSC.2007.4419664.
- North C and Shneiderman B (2000) Snap-together visualization: A user interface for coordinating visualizations via relational schemata. In: Gesú VD, Levialdi S and Tarantino L (eds.) *AVI'00: Proceedings of the Working Conference on Advanced Visual Interfaces*. ACM. ISBN 1581132522, pp. 128–135. DOI:10.1145/345513.345282.
- Oliveira P, Rodrigues F and Henriques P (2005) A formal definition of data quality problems. In: *ICIQ'05: Proceedings of the International Conference on Information Quality*.
- Organization for Standardization (1990) ISO/IEC 10027:1990 – Information Resource Dictionary System (IRDS) framework.
- Organization for Standardization (2008) ISO/IEC 25012:2008 – Software product Quality Requirements and Evaluation (SQuaRE) data quality model.
- Organization for Standardization (2014) ISO 19115-1:2014 – Geographic Information (Metadata).
- Pang AT, Wittenbrink CM and Lodha SK (1997) Approaches to uncertainty visualization. *The Visual Computer* 13(8): 370–390. DOI:10.1007/s003710050111.
- Patro A, Ward MO and Rundensteiner EA (2003) Seamless integration of diverse data types into exploratory visualization systems. Technical Report WPI-CS-TR-03-12, Worcester Polytechnic Institute.
- Petoukhov V, Ganopolski A, Brovkin V, Claussen M, Eliseev A, Kubatzki C and Rahmstorf S (2000) CLIMBER-2: A climate system model of intermediate complexity. Part I: Model description and performance for present climate. *Climate Dynamics* 16(1): 1–17. DOI:10.1007/PL00007919.
- Potter K, Rosen P and Johnson CR (2012) From quantification to visualization: A taxonomy of uncertainty visualization approaches. In: Dienstfrey AM and Boisvert RF (eds.) *Proceedings of the IFIP Working Conference on Uncertainty Quantification in Scientific Computing, IFIP Advances in Information and Communication Technology*, volume 377. Springer. ISBN 9783642326769, pp. 226–249. DOI:10.1007/978-3-642-32677-6_15.
- Ragan ED, Endert A, Sanyal J and Chen J (2016) Characterizing provenance in visualization and data analysis: An organizational framework of provenance types and purposes. *IEEE Transactions on Visualization and Computer Graphics* 22(1): 31–40. DOI:10.1109/TVCG.2015.2467551.
- Raman V and Hellerstein JM (2001) Potter's Wheel: An interactive data cleaning system. In: *VLDB'01: Proceedings of the*

- Conference on Very Large Data Bases. Morgan Kaufmann Publishers. ISBN 1558608044, pp. 381–390.
- Rankin R (1990) A taxonomy of graph types. *Information Design Journal* 6(2): 147–159. DOI:10.1075/idj.6.2.03ran.
- Riede M, Schueppel R, Sylvester-Hvid KO, Kühne M, Röttger MC, Zimmermann K and Liehr AW (2010) On the communication of scientific data: The full-metadata format. *Computer Physics Communications* 181(3): 651–662. DOI:10.1016/j.cpc.2009.11.014.
- Ristovski G, Preusser T, Hahn HK and Linsen L (2014) Uncertainty in medical visualization: Towards a taxonomy. *Computers & Graphics* 39: 60–73. DOI:10.1016/j.cag.2013.10.015.
- Robertson PK (1990) A methodology for scientific data visualisation: Choosing representations based on a natural scene paradigm. In: Kaufman A (ed.) *Visualization'90: Proceedings of the IEEE Conference on Visualization*. IEEE. ISBN 0818620838, pp. 114–123. DOI:10.1109/VISUAL.1990.146372.
- Robertson PK (1991) A methodology for choosing data representations. *IEEE Computer Graphics and Applications* 11(3): 56–67. DOI:10.1109/38.79454.
- Rockel B, Will A and Hense A (2008) The regional climate model COSMO-CLM (CCLM). *Meteorologische Zeitschrift* 17(4): 347–348. DOI:10.1127/0941-2948/2008/0309.
- Rodríguez A, Caro A, Cappiello C and Caballero I (2012) A BPMN extension for including data quality requirements in business process modeling. In: Mendling J and Weidlich M (eds.) *BPMN'12: Proceedings of the International Workshop on Business Process Model and Notation, Lecture Notes in Business Information Processing*, volume 125. Springer. ISBN 9783642331541, pp. 116–125. DOI:10.1007/978-3-642-33155-8_10.
- Roth SF and Mattis J (1990) Data characterization for intelligent graphics presentation. In: Chew JC and Whiteside J (eds.) *CHI'90: Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM. ISBN 0201509326, pp. 193–200. DOI:10.1145/97243.97273.
- Russell S (2008) *Resource Availability Awareness and Data Utility: The Foundation for a DSS Framework in a Pervasive Computing Environment*. PhD Thesis, University of Maryland.
- Setlur V and Stone MC (2016) A linguistic approach to categorical color assignment for data visualization. *IEEE Transactions on Visualization and Computer Graphics* 22(1): 698–707. DOI:10.1109/TVCG.2015.2467471.
- Shneiderman B (1996) The eyes have it: A task by data type taxonomy for information visualizations. In: *VL'96: Proceedings of the IEEE Symposium on Visual Languages*. IEEE. ISBN 081867508X, pp. 336–343. DOI:10.1109/VL.1996.545307.
- Silva CT, Freire J and Callahan SP (2007) Provenance for visualizations: Reproducibility and beyond. *Computing in Science and Engineering* 9(5): 82–89. DOI:10.1109/mcse.2007.106.
- Silva S, Santos BS and Madeira J (2011) Using color in visualization: A survey. *Computers & Graphics* 35(2): 320–333. DOI:10.1016/j.cag.2010.11.015.
- Simmhan YL, Plale B and Gannon D (2005a) A survey of data provenance in e-science. *ACM SIGMOD Record* 34(3): 31–36. DOI:10.1145/1084805.1084812.
- Simmhan YL, Plale B and Gannon D (2005b) A survey of data provenance techniques. Technical Report IUB-CS-TR618, Computer Science Department, Indiana University, Bloomington.
- Simmhan YL, Plale B and Gannon D (2008) Karma2: Provenance management for data driven workflows. In: Zhang LJ (ed.) *Web Services Research for Emerging Applications: Discoveries and Trends*. IGI Global. ISBN 9781615206841, pp. 317–339. DOI:10.4018/978-1-61520-684-1.ch014.
- Skeels M, Lee B, Smith G and Robertson GG (2010) Revealing uncertainty for information visualization. *Information Visualization* 9(1): 70–81. DOI:10.1057/ivs.2009.1.
- Stanford NJ, Wolstencroft K, Golebiewski M, Kania R, Juty N, Tomlinson C, Owen S, Butcher S, Hermjakob H, Novère NL, Mueller W, Snoep J and Goble C (2015) The evolution of standards and data management practices in systems biology. *Molecular Systems Biology* 11(12): 851. DOI:10.15252/msb.20156053.
- Steinacker A, Ghavam A and Steinmetz R (2001) Metadata standards for web-based resources. *IEEE Multimedia* 8(1): 70–76. DOI:10.1109/93.923956.
- Stevens SS (1946) On the theory of scales of measurement. *Science* 103(2684): 677–680. DOI:10.1126/science.103.2684.677.
- Stitz H, Luger S, Streit M and Gehlenborg N (2016) AVOCADO: Visualization of workflow-derived data provenance for reproducible biomedical research. *Computer Graphics Forum* 35(3): 481–490. DOI:10.1111/cgf.12924.
- Streit M, Schulz HJ, Lex A, Schmalstieg D and Schumann H (2012) Model-driven design for the visual analysis of heterogeneous data. *IEEE Transactions on Visualization and Computer Graphics* 18(6): 998–1010. DOI:10.1109/TVCG.2011.108.
- Sulo R, Eick S and Grossman RL (2005) DaVis: A tool for visualizing data quality. In: *InfoVis'05: Poster Compendium of the IEEE Symposium on Information Visualization*. pp. 45–46.
- Swayne DF and Buja A (1998) Missing data in interactive high-dimensional data visualization. *Computational Statistics* 13(1): 15–26.
- Theus M, Hofmann H, Siegl B and Unwin A (1997) MANET – extensions to interactive statistical graphics for missing values. In: *New Techniques and Technologies for Statistics II: Proceedings of the Second Bonn Seminar*. IOS Press, pp. 247–259.
- Thompson DL, Braun JA and Ford R (2004) *OpenDX – Paths to Visualization*. 2nd edition. Visualization and Imagery Solutions, Inc.
- Thomson J, Hetzler E, MacEachren A, Gahegan M and Pavel M (2005) A typology for visualizing uncertainty. In: Erbacher RF, Roberts JC, Gröhn MT and Börner K (eds.) *VDA'05: Proceedings of the Conference on Visualization and Data Analysis*. SPIE. ISBN 9780819456427. DOI:10.1117/12.587254.
- Tory M and Möller T (2004) Human factors in visualization research. *IEEE Transactions on Visualization and Computer Graphics* 10(1): 72–84. DOI:10.1109/TVCG.2004.1260759.
- Tshagharyan G and Schulz HJ (2013) A graph-based overview visualization for data landscapes. *Computer Science and Information Technology* 1(3): 225–232.
- Vassiliadis P (2009) Data warehouse metadata. In: Liu L and Özsu MT (eds.) *Encyclopedia of Database Systems*. Springer. ISBN 9780387355443, pp. 669–675. DOI:10.1007/

- 978-0-387-39940-9_912.
- Ward M, Xie Z, Yang D and Rundensteiner E (2011) Quality-aware visual data analysis. *Computational Statistics* 26(4): 567–584. DOI:10.1007/s00180-010-0226-0.
- Wickham H (2014) Tidy data. *Journal of Statistical Software* 59(10): 1–23. DOI:10.18637/jss.v059.i10.
- Wills G and Wilkinson L (2010) AutoVis: Automatic visualization. *Information Visualization* 9(1): 47–69. DOI:10.1057/ivs.2008.27.
- Wong PC, Foote H, Leung R, Jurrus E, Adams D and Thomas J (2000) Vector fields simplification – a case study of visualizing climate modeling and simulation data sets. In: Ertl T, Hamann B and Varshney A (eds.) *VIS'00: Proceedings of the IEEE Conference on Visualization*. IEEE. ISBN 0780364783, pp. 485–488. DOI:10.1109/VISUAL.2000.885738.
- World Wide Web Consortium (W3C) (2014) Resource Description Framework (RDF) 1.1 concepts and abstract syntax. URL <http://www.w3.org/TR/rdf-concepts/>.
- Xiao B, Zhang C, Mao Y and Qian G (2015) Review and exploration of metadata management in data warehouse. In: *ICIEA'15: Proceedings of the IEEE Conference on Industrial Electronics and Applications*. IEEE, pp. 928–933. DOI:10.1109/ICIEA.2015.7334243.
- Yang J, Hubball D, Ward MO, Rundensteiner EA and Ribarsky W (2007) Value and relation display: Interactive visual exploration of large data sets with hundreds of dimensions. *IEEE Transactions on Visualization and Computer Graphics* 13(3): 494–507. DOI:10.1109/TVCG.2007.1010.
- Zhou MX and Feiner SK (1996) Data characterization for automatically visualizing heterogeneous information. In: Gershon ND, Card S and Eick SG (eds.) *InfoVis'96: Proceedings of the IEEE Symposium on Information Visualization*. IEEE. ISBN 081867668X, pp. 13–20. DOI:10.1109/INFVIS.1996.559211.
- Zhuge H (2004) Resource space model, its design method and applications. *Journal of Systems and Software* 72(1): 71–81. DOI:10.1016/s0164-1212(03)00058-x.
- Zhuge H, Yao E, Xing Y and Liu J (2005) Extended resource space model. *Future Generation Computer Systems* 21(1): 189–198. DOI:10.1016/j.future.2004.09.016.