# Interactively Exploring Bibliographical Data for Literature Analysis

Stefan Schlechtweg, Stefan Büder, Marcel Götze

Otto-von-Guericke-Universität Magdeburg, Institut für Simulation und Graphik

**Abstract**

This paper introduces techniques for interactive navigation within large sets of bibliographic data. The main conceptual idea is to use various relations between the entries to navigate within the information space. A visualisation that supports both the display of the data itself and the relations is introduced. Interaction techniques offer possibilities to follow relations and, thus, create new views onto the data. The proposed visualisation and interaction techniques improve literature analysis tasks that occur when writing a scientific document or when reviewing and exploring a scientific topic based on literature.

## 1        Introduction

Writing a scientific paper, a thesis, or any other work in the scientific community requires a thorough understanding and analysis of other people's written work. Literature review and analysis is therefore an integral part of successful scientific work. Working with literature not only includes to know about the contents of papers and books but also to know relations between certain works, topics and authors who work in an area and how these are interrelated. Especially these relations are a valuable basis to derive new ideas.

Keeping track of the literature in a field is rarely supported by visual tools for exploring the above mentioned relations. Often, a database containing bibliographic information, possibly personal annotations, and a reference to an electronic version of the article is the only source for a literature survey or analysis. Relations are established and maintained mainly mentally by the person doing the survey, or, for example, taken down in (handdrawn) diagrams.

We propose techniques to support the literature analysis by means of visualization and interaction of and with such relations. Based on a database containing bibliograpic and additional information, several kinds of data are visualized. This includes information for single papers, a whole body of papers, and relations between papers and/or authors. We further propose interactive tools to explore the information space, and possibly to add new information.

# 2        Related Work

There exist a wealth of projects that deal with a certain visualization of bibliographic data. Most interfaces to digital libraries, however, do not use the full power that can be gained when exploiting relational connections among the documents in the library. A standard search engine interface is not enough if it comes to analyzing a field of research or following a scientific topic. The problem here is that the queries to the digital library have to be created manually from the contents of the read paper or the information presented within the digital library. The exploitation of relational information would help to create such queries automatically and, thus, provide an easier way to navigate within the library.

Almost all digital libraries today use a form based interface to their search engine and provide the query results in list form on a web page. Different information about the documents are given ranging from author, title and abstract to a list of citations or a list of documents that are related in some other way. Navigation is almost exclusively based on user formulated queries. A recent example that uses relations for navigation in a digital library is CITEULIKE (www.citeulike.org). Here, a user can collect documents and relate them to each other via the assignment of "tags". CITEULIKE offers a node-link-diagram to show which documents are related to each other and to navigate through the library. The relations are solely computed based on the assigned tags, no citation or authoring relations are considered.

Hsu et al. (2004) present with MONKELLIPSE an interactive visualization that builds on an elliptical design. All documents are chronologically laid out in an elliptical shape, where additional indicators show a grouping by years. Research topics are shown in the inner area of the ellipse with each area in the relative center of the documents that belong to it. Selecting a document leads to an emphasis of all cited documents as well as of the respective research area. A selection of the research area emphasises all comprising documents. While the pleasing layout, a good use of screen space and the elaborate interaction concept make MONKELLIPSE a nice tool to explore a body of literature, many features that are needed for a deeper literature analysis are missing. Also, the unchangeable layout poses some problems, especially if other criteria than the year of publication are needed to sort the documents.

A different approach was taken by Wong et al. (2004) with their IN-SPIRE system. Here, the most important point is the distribution of documents among various topic areas. A set of documents is visualized as topic map or galaxy view showing clusters of documents belonging together. Additional tools allow, for instance, filtering by time of publication. The outlier tool allows further filtering and a more detailed examination of the document set. Both the topic map and the galaxy view are not connected to each other. This makes the exploration of the information space rather difficult. Also, the visualizations can not be adapted to other information needs, as for instance, authors or document-author relationships.

With a focus on visualization and analysis of large complex networks, WILMASCOPE was presented by Ahmed et al. (2004). To visualize bibliographic data, relations between authors, documents, and other relations are extracted from the data and a network is built. This is then visualized in various ways, allowing to see, for example, the (co-)citations, authors and their relations, or relations within and in between research topics. Considering all such relations yields a complex network which can be visualized differently. The used layout techniques make the most prominent nodes in the network stand out. Also, the layered layout offers the

advantage of less edge crossings as with normal 2D graph layouts. WILMASCOPE leads to very comprehensive visualizations which are, however, not interactively explorable.

Most applications for handling bibliographical data restrict themselves to browsing the database and showing statistical data. BIBRELEX (Brueggemann 1999) uses standard graph drawing techniques to reveal and show document relationships, e.g. the citation network. It is, however, restricted to these kinds of relations and the chosen spring mass based graph layout changes if new nodes were added, possibly leading to a complete layout change. The DBL-BROWSER (Klink 2004) is primarily an interface for browsing in an online bibliographic database. Some visualization tools are added to explore various aspects. The timeline graph gives an overview of time-oriented aspects, e.g., the distribution of published papers of an author over a certain time period. Other more network based visualizations show relations between authors (co-authorship) or between documents (citations). While these visualizations are very useful, a global overview of the complete data set is missing.

Citation and Co-Citation analysis is another area to be considered. Often citation chains and co-citations lead to different documents that are of use for the task at hand. HISTCITE (Garfield 2002), for example, uses node link diagrams to visualize citation and co-citation graphs. While such visualizations help to get an overview, they might become rather com-plex and contain too much information for an actual literature analysis task. This can be seen in the CITESPACE system (Chen 2004; Chen 2006). Based on co-citation networks, Chen builds a visualization that supports the identification of intellectually significant articles based on a visual inspection of the graph. Even though the goal behind Chen's work is some-what different from ours, it shows that a visualization of relational data helps to get new information from a set of bibliographic data.

# 3          Requirements

For an overview of the complete data set, a comprehensive visualization of the *complete* data should be provided. Almost all of the mentioned systems offer such a feature. Searching as well as filtering needs to be provided to support navigation and to reduce the amount of data with which a user is working. WILMASCOPE does not offer searching while even more of the previously mentioned systems, namely WILMASCOPE, MONKELLIPSE, and BIBRELEX, do not contain filtering. As far as the navigation is concerned, a smooth transition between the views is required in order to keep the context. Detail-on-demand techniques can be used to get specific information about one selected entity, narrow a selection, or working as initiator for a new filtering. An example is to get information about the authors of a publication and then use the co-authorship relation to get to a different set of documents. To enable this filtering, relational information have to be extracted from the initial data and tools have to be provided for navigating with these relations. Astonishingly enough, only WILMASCOPE makes extensive use of such relations. User defined relations are a further information source and even more support navigation. Such relations need to be defined for or computed from the given data. Also, user defined annotations and possibly changes of the data (addenda, corrections) will help in supporting literature analysis in an interactive and visual way.

# 4        The Data

The original data which were used to form the information space consist of an extended BiBTeX database of approximately 600 documents from the area of computer graphics written by about 800 authors. We have added a field to each entry that holds the keys of all documents which are cited by the respective document, a field for topic keywords and a field containing an abstract of the document. Each document in the database is characterized by a set of attributes, which are derived from the given BiBTeX fields. Other input formats, e.g., XML, can also be used if a respective import filter exists. In order to provide the necessary data for a rich visualization and interaction, we derive three sets of data entities, a set of documents: $\mathbf{D}$, a set of authors: $\mathbf{A}$, and a set of research topics: $\mathbf{T}$.

There exist various relations between these three sets which can be calculated from the data base. The following are the most important of them:

- $\mathbf{D} \rightarrow \mathbf{A^n}$: all authors of a document

- $\mathbf{A} \rightarrow \mathbf{D^n}$: all documents written by an author

- $\mathbf{A} \rightarrow \mathbf{A^n}$: all co-authors of an author

- $\mathbf{D} \rightarrow \mathbf{D^n}$: all documents being cited by a document or all documents citing a document

- $\mathbf{D} \rightarrow \mathbf{T^n}$: all topic areas of a document

- $\mathbf{T} \rightarrow \mathbf{D^n}$: all documents in a topic area

In general, the following visualization tasks need to be realized for the task at hand:

- visualizing all members of a dataset (overview)

- visualizing the attributes of a single document (attribute relations)

- visualizing relations within one dataset (internal relations)

- visualizing relations between two different datasets (external relation)

# 5        Visualization and Interactive Literature Analysis

To visualize the literature data and relations, we borrow the elliptic design from TEXTARC (Paley 2002) and MONKELLIPSE (Hsu 2004). TEXTARC is an alternative way of displaying a continuous text which is arranged line by line in an elliptical form. In the inner area of the ellipse relevant words are placed according to their position in the text. Relations between words and the whole text are displayed via lines drawn when moving the pointer over a certain word. We build on this basic idea and generalize it for certain types of relations and data sets. The main power of our approach lies in the combination of those visualizations via interaction. Thus, the presented visualization together with interactive techniques allows a target-oriented navigation within the information space. In the following we will present tools for navigation using actual screenshots from the implemented system. The system was

realized as a prototype using OpenGL for visualization and a MySQL database to store the bibliographic information. The principles behind can, however, easily be adopted to other graphics libraries and to data handling via network.

The elliptic design offers two primary advantages. First, due to its aspect ratio, screen space is better used in comparison to a circular layout. Second, an elliptical visualization can be devided into an inner (the ellipse's surface) and an outer part (the ellipse's line and the surface outside of the ellipse), so two different data sets can be easily combined visually.

## 5.1        Visualization of Bibliographical Data

The most general visualization, tool is the *browse ellipse* which displays all entities of a certain set. Figure 1 shows an implementation of this concept. All entries are arranged around the circumference of the ellipse where a single entry is displayed as label.
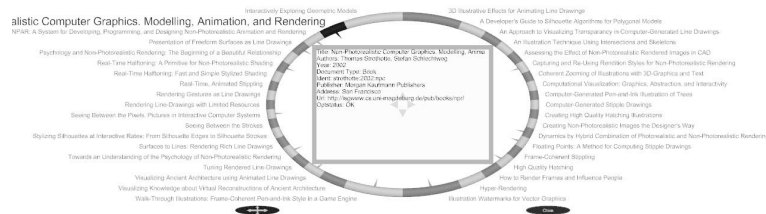


*Figure 1: Concept of the browse ellipse, showing all documents of one author*

Such visualization gives a first overview of the information space and offers various navigation tools to gain a deeper insight. The ellipse itself in subdivided into segments, one for each label, which are colored alternatively to make the items more distinguishable. If a very large number of items is present, the labels are not readable. Therefore, Focus+Context techniques according to the *Information Seeking Mantra* (Shneiderman 1996) should be provided. In a sense, this layout with labels made unreadable to fit the whole information space in the visualization, resembles the Information Murals as presented by Jerding and Stasko (Jerding & Stasko 1998). Moving the pointer over the labels is used to enlarge the focused label, a direct selection triggers the attribute relation and reveals detailed information in the inner area of the ellipse. Selecting more than one item gives the possibility to group these or open a new browse ellipse with just these items. Offering search functionality yields a selection of entries that match the search and which are then displayed.

The browse ellipse is also the basis for the visualization of relations. Figure 1 has already shown the display of attribute relations: selecting exactly one entry activates the attribute relation and the entry's attributes are shown in the inner surface of the ellipse. If internal rela-tions like $A \rightarrow A^n$ or $D \rightarrow D^n$ are to be displayed, the originator of the relation is selected and all related items will automatically become selected. Between both, lines or arrows, depending on the kind of relation, are drawn as in Figure 2. To enhance the visibility of the lines, especially when they reach to nearby items, they start from markers drawn in-

ward from the position of the respective labels. All selected labels are enlarged and made readable.
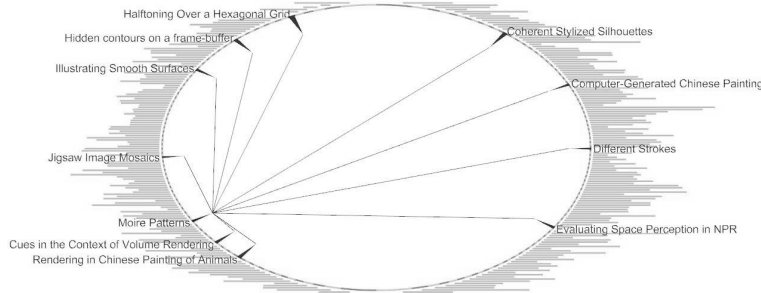


*Figure 2: Showing internal relations as arrows between the selected origin and all related entries*

External relations, i.e., relations between data items from two different sets use the inner area of the ellipse to display the second data set while the first is laid out as shown for the browse ellipse. This layout is especially useful for relations between authors and documents, i.e., $D \rightarrow A^n$ or $A \rightarrow D^n$. A single data item from the second set is placed in the inner area while lines emanate from there to the related items from the first data set on the circumference of the ellipse. The position of the inner data item is either central or can be computed using a force based algorithm that positions the inner item with respect to the positions of the related items in their relative center.
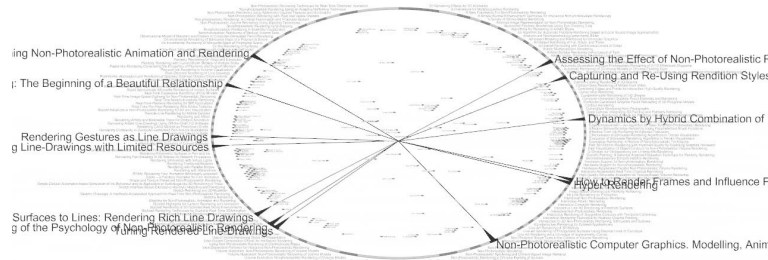


*Figure 3: A single item is related to some data items from another data set.*

Such a force based model is also needed if external relations between several data items from both sets are to be displayed. The items from the first set are, again, distributed on the ellipse's circumference. All items from the second set are laid out following the force based approach in the inner area. In this way, items having similar relations to items from the first set come close together. A direct link via lines and arrows is established when moving the pointer over the respective item (cf. Figure 3). This visualization seems especially valuable

for relations like $D \rightarrow A^n$ or $A \rightarrow D^n$. However, both data sets can also be identical, so that this concept can also be used to show internal relations like $A \rightarrow A^n$ or $D \rightarrow D^n$.

To visualize groupings, the browse ellipse is segmented according to the number of groups and these groups are visually emphasized by drawing wedges in the inner area of the ellipse. Each wedge can then hold certain information if space permits. The most simple kind of information being displayed is a label as in Figure 4. Moreover, all group members can be displayed as points within the area of the wedge belonging to the group so that a display of relations becomes possible. When the pointer is moved over a data item, all relations to other items from other groups are shown as direct lines between the respective points.
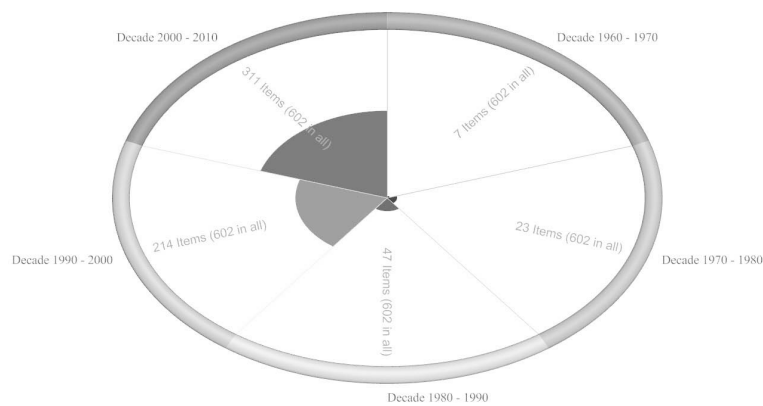


*Figure 4: Grouping of entries with general information about the groups shown as labels. Here, overview of the data grouped by decade of publication*

Some of the relations which are of interest in literature analysis lead to relation chains. The most prominent example here is the citation relation which yields citation networks. A standard way of visualizing such citation networks as graphs soon becomes problematic since the resulting graphs are rather complex. For a literature analysis in connection with a writing task, citation networks are often needed level by level, or, alternatively, single paths need to be followed, so that an interactive exploration of such networks is more appropriate. Starting with a browse ellipse, all citation relations (for example) can be treated like external relations, so that the ellipse itself with the selected item is moved into the center of a second ellipse which shows the same data set. All relations are drawn as direct links. For the next level of the network, a third ellipse is created in the same way. Moving those ellipses in 3D on separate levels even enhances the recognition of the links due to the possibility to view the network from various angles (see Figure 5).

The visualization designs so far concentrate mainly on one particular relation. When working with bibliographical data, there is often the need to follow different relations in order to gain deeper insight into the topic or to establish connections between different documents that are not obvious on a first glance. Therefore, interaction techniques are needed to explore the information space in various ways.
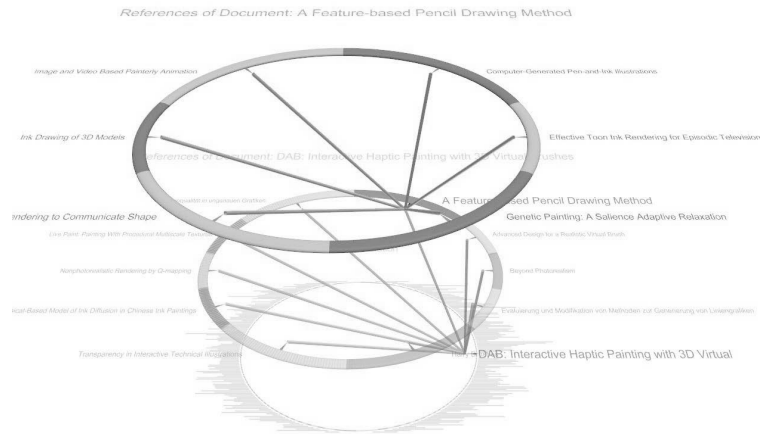
*Figure 5: Three levels of the citation network shown in three ellipses laid out in 3D*

## 5.2      Interactive Literature Analysis

The usual starting point for a literature analysis is a search in the database for either a specific author or a keyword that reveals a set of papers which contain that keyword in the title or even the document text. It would also be possible to start with an overview visualization of the complete set of documents (e.g., as shown in Figure 1). This is a good idea if someone who starts an analysis within a relatively unknown area and wishes to actually search for a paper or author. In our example, we shall start with a different kind of overview visualization that gives the temporal distribution and grouping of the papers related to computer graphics (our example data set). Figure 4 shows this overview where all papers are grouped by decade of publication. This overview is augmented by a histogram like indicator in the ellipse's area to quickly get a comparison of the number of documents in each decade.

We then select the biggest chunk of data, the decade starting 2000 and call a browse ellipse for all documents within this decade. Now the user can browse through these data by moving the mouse over the titles. Note that due to the number of elements the titles are reduced in size so that they are not readable at the moment. However, the title currently under the mouse will be enlarged and highlighted. Moreover, for each selected document, detail information is shown in the middle of the ellipse, as can be seen in Figure 1.

After the user has found a paper "show references" was selected yielding Figure 2 which shows all cited documents as internal relation of the kind $\mathbf{D} \rightarrow \mathbf{D^n}$. From this relation mainly alternative or previous research approaches can be derived. Taking this further and selecting one of the cited documents and again calling up all references leads to the second level of the citation network. Going another step yields Figure 5 where three levels of the citation network are shown. Note that only the subset of the documents that is actually cited is shown on higher levels to avoid too much overlapping and visual clutter.

After working with the citation network, we return to the document overview in Figure 1 and go a different way which is also often gone when analyzing scientific literature. Assuming that the author of a specific paper is working in a specific topic area, it is a good idea to find other papers by the same author. This will yield an overview of the author's research. There are two possible ways to get this information. First, a new browse ellipse can be opened which shows just the documents by the respective author (see Figure 1). Second, the external relation $A \rightarrow D^n$ can be used to relate author and documents visually to each other as can be seen in Figure 3. In contrast to the first option, here the displayed information space is not limited and can be explored further. However, the visualization is also more complex. The inner visualization is made up of the authors which are placed in the relative center of their documents. Selecting one author highlights the documents he or she has written.

Co-authors of an author are likely to work on the same or similar research topics, so that the co-authorship relation can be used to get a more focussed overview of a research area. Figure 2 is an example of a visualization of the co-author relation $A \rightarrow A^n$, which is an internal relation on the set of authors.

All visualizations as shown in this section are interlinked with each other and reachable via simple interactions, i.e., selecting an item and determining the new visualization via popup menu selection. The transition between one visualization to the other is smooth and an undo functionality helps in going back to earlier stages in the exploration process.

# 6        Conclusion and Future Work

The proposed visualization and interaction techniques offer a way to browse bibliographic data while being focussed on relational connections between documents, authors and their attributes. In comparison to standard form based interfaces to digital libraries, such a relation based exploration eases the way to find new documents or to mentally connect information. Form based interfaces to digital libraries are perfectly suited for specific queries if the user exactly knows what he or she is looking for. In comparison, queries that may accure from reading a document are rather unspecific, as for example "all documents form an author and his co-authors", or do not (yet) belong to standard queries for digital libraries.

We argue that for such queries an exploratory approach using relational information and directly visualizing relations is better suited as a standard form based input. Furthermore, the time to find matching results and the mental load is supposed to be smaller. The user is no longer required to build the query outside of the interface and, possibly to mentally integrate several earlier query results. First informal tests have supported this hypothesis, however, a detailed user study will be necessary to evaluate the amount of improvement. Such a study will, nonetheless, require a reimplementation of the techniques so that it could be integrated in existing digital libraries.

The core argument of this paper is that an explicit visualization of relational information within bibliographical data and the use of such visualizations in navigating bibliographical information helps in performing literature analysis tasks. So far we have only considered those relations that have a strong bibliographical background. Adding user-centered relations

as it is sometimes done in digital libraries ("users who have read A also read B") or even user defined relations will open new possibilities and directions. Also, a closer combination with standard query-based retrieval will offer advantages by combining the positive sides of both concepts. One application area that is interesting to investigate is the use of our techniques in thematic communities related to scientific literature. One such community, CITEULIKE has already been mentioned in Section 2. The discussion of literature will also be enhanced by offering navigation aids through the steadily increasing amount of publications.

Aside from bibliographical data, the presented techniques can also be applied to other areas where relational information between data items play an important role. It might be worthwhile to investigate the use of our techniques in the context of online communities and social networks where members are connected by a wealth of different relations.

## References

Ahmed A.; Dwyer T.; Murray C.; Song L.; Wu Y. X. (2004): Wilmascope graph visualisation. In: Proc. of INFOVIS'04, IEEE Computer Society, p. 216.4.

Brüggemann-Klein A.; Klein R.; Landgraf B. (1999): Bibrelex: Exploring bibliographic databases by visualization of annotated contents-based relations. D-Lib Magazine 5(11): (1999).

Chen, C. (2004): Searching for intellectual turning points: Progressive knowledge domain visualization. In: Proc. of the National Academy of Sciences of the United States of America (Washington, 2004), vol. 101, National Academy of Sciences, pp. 5303-5310.

Chen, C. (2006): CiteSpace II: Detecting and Visualizing Emerging Trends and Transient Patterns in Scientific Literature. Journal of the American Society for Information Science & Technology.

Garfield, E.; Pudovkin, A. I.; Istomin, V. S. (2002): Algorithmic Citation-Linked Historiography – Mapping the Literature of Science. In: Proc. of the 65th Annual Meeting of the American Society for Information Science & Technology, vol. 39, pp. 14-24.

Hsu, T.-W.; Inman, L.; Mccolgin, D.; Stamper, K. (2004): MonkEllipse: Visualizing the History of Information Visualization. In: Proc. of INFOVIS'04, IEEE Computer Society, p. 216.9.

Jerding, D. F.; Stasko, J. T. (1998): The Information Mural: A Technique for Displaying and Navigating Large Information Spaces. IEEE Transactions on Visualization and Computer Graphics 4(3):(1998), pp. 257-271.

Klink, S.; Ley, M.; Rabbidge, E.; Reuther, P.; Walter, B.; Weber, A. (2004): Browsing and visualizing digital bibliographic data. In: Proc. of the 2004 Eurographics/IEEE TVCG Workshop on Visualization., Eurographics Association, pp. 237-242.

Paley, W. B. (2002): TextArc: Revealing Word Associations, Distributions and Frequency. Interactive Poster at the IEEE INFOVIS'02.

Shneiderman, B. (1996): The Eyes Have it: A Task by Data Type Taxonomy for Information Visualization. In: Proc. of VL'96 Symposium on Visual Languages, IEEE Press, pp. 336-343.

Wong, P. C.; Hetzler, B.; Posse, C.; Whiting, M.; Havre, S.; Cramer, N.; Shah, A.; Singhal, M.; Turner A.; Thomas, J. (2004): In-spire infovis 2004 contest entry. In: Proc. of INFOVIS'04, IEEE Computer Society, p. 216.2.