

# Social Aspects of Web Page Contents

Miloš Kudělka, Václav Snášel, Zdeněk Horák  
VŠB – Technical University of Ostrava, Czech Republic  
Email: milos.kudelka@inflex.cz, vaclav.snaselvsb.cz, zdenek.horak.st4@vsb.cz

Ajith Abraham  
Machine Intelligence Research Labs (MIR Labs)  
Scientific Network for Innovation and Research Excellence, USA  
Email: ajith.abraham@ieee.org

## Abstract

*In this paper we try to consider a Web page as information with social aspects. Each Web page is the result of invisible social interaction. This interaction between different groups of people translates into a certain unification of Web page creation. External signs of this unification are the features of the Web page, that meets the user's expectations. Through analysis of the features, we can obtain information that can simply describe the Web page. This simple description contains strong information about the social group the page is intended for. If the user uses this information to refine the search, then he identifies himself as a member of a certain social group. For the description of the social aspects of Web pages we use the term MicroGenre. This paper describes the fundamental concepts of MicroGenre and also illustrate experiments for the detection and usage of MicroGenres.*

## 1 Introduction

One of the key problems of the current Internet is the problem of retrieving relevant Web pages. One possible way - which can improve user orientation in large sets of Web pages - is Web page description using so-called Genres [22]. A Genre is a kind of typology corresponding to the human view of Web page content. It is, therefore, clear that the specific Genre is associated with a specific social group of users, which find the pages of this Genre useful. A significant feature for the Genres usability in retrieval is their ability to be recognized by machine-applied algorithm. There are different approaches to detect the Web page Genre automatically. The problem is that the contemporary Web pages are often very complex and cannot be clearly classified to a single Genre (see [3]). The other side of the

Web page description problem is the so-called Web design patterns (see [29], [28]). They are used by the Web page developers to describe frequently used high-level page elements and methods for their design. It is similar to Genres, because each Web design pattern is associated with a particular social group, the Web developers. It is known fact that the design patterns can serve as a vocabulary between different groups of people (designers, analysts, developers, project managers, etc.).

From the view of Web page description, the Genres and Web design patterns are very close (see Table 1). This paper introduces the concept of MicroGenres. Under the notion of MicroGenres, we understand more or less independent Web page elements, which have some specific purpose and content. Utilizing the term "MicroGenre" we can describe the Web page as a collection of MicroGenres.

<b>Genres (Roussinov, Meyer zu Eissen, Boese, etc.)</b>
Homepage, Articles, News bulletin, Glossary, Course Lists, Instructional Materials, Geographical Location, Special Topics, Publications, Product Information, Product Lists, Ads, Order Forms, Ratings Help, Article, Discussion, Shop, Portrayal, Hub, Download, etc.
<b>Web Design Patterns (Page Type) – www.welie.com</b>
Article Page, Blog Page, Case Study, Contact Page, Event Calendar, Forum, Guest Book, Help Page, Homepage, Newsletter, Printer-friendly Page, Product Page, Tutorial.
<b>MicroGenres</b>
Price Information, Purchase Possibility, Special Offer, Product Catalogue, Product Technical Features, Discussion and Comments, Review, Customer Reviews, News, Author and Publications, Book Info, Job Advertisement, Personal Advertisement, etc.

**Table 1. Genres, Patterns and MicroGenres**

One of the contemporary trends is the analysis of the behavior of users and identification of the social group the user belongs to.

The results of this analysis may be useful for increasing the accuracy of search engine results. If the user belongs to some specific social group then he probably expects results

# Forum

## Problem

Users want to discuss a certain topic or react on a particular piece of content on the site

## Solution

Create a list of topics and allow users to place comments on the topic

The screenshot shows the Esato forum interface. At the top, there are navigation links for Home, Free Links, Colour Logos, Screensavers, Themes, Ringtones, News, Mobile Phones, Discussion Forum, Resources, FAQ, Site Map, and Contact. Below this, there are statistics: 34 New posts, 10788 Total Messages, and 21876 Members. A table lists forum topics with columns for Forum, Topics, Posts, and Last Post. The table includes categories like Esato news, Sony Ericsson, Accessories, Headset, Bluetooth, MP3-player problems and solutions, and General.

The screenshot shows a search results page for the Esato forum. It includes a search bar and a list of results. The results are organized into categories like Technical, Regional, and Americas. Each result shows the search terms, the number of messages and topics, and the date of the last post.

## Use when

You are designing a [Community Site](#) or other site for which you are interested in [Community Building](#). Many people can be tied into a site when there is ample opportunity to interact with the site. Such interaction with a site, or its content in particular can also be found on [News Site](#) or on a [Article Page](#).

## How

A forum is literally a discussions place. It is based on discussion topics and their

Figure 1. Sample of Forum pattern (www.welie.com)

selected by other users from the same social group.

An unpleasant feature of this kind of analysis is the need to store information about individual users. From our point of view, the information about a social group is contained in the page itself. Both approaches cannot be easily compared. In our paper we want to show that the user working with some specific Web page also shows – up to a certain degree – his aims and social group involvement. The important fact is that this aim can be obtained from the page itself, not from the concrete Web page content, but more likely from the context of this content. As an example we can consider discussion, review, product catalog, blogs, etc. Information context is a key element to us. It defines the purpose of this page, and this purpose determines the target group of this Web page and as a result it defines roughly a certain social group. Using the content and the purpose of the Web page we can fairly accurately estimate which group the user belongs to and what his aim is.

Rest of the paper is organized as follows. In the Section 2, we review the term Genre and Web design pattern. Section 3 introduces the term MicroGenre and discusses its detection. In Section 4, we describe the experiments dealing with Web site description. The last Section contains paper recapitulation and focuses on possible directions of further research.

## 2 Web Genres and Web Design Patterns

When we perceive the Web page as whole, the purpose is represented by a so-called Web genre. Similarly, the view of individual segments of the Web page is closely related to Web design patterns. A Web genre is a taxonomy that incorporates the style, form and content of a document which

is orthogonal to the topic, with fuzzy classification to multiple Web genres [3]. For classification, there are many approaches and also many methods for genre identification. Kennedy and Shepherd [10] analyzed home page genres (personal home page, corporate home page or organization home page). Chaker and Habib [5] proposed a flexible approach for Web page genre categorization. Flexibility means that the approach assigns a document to all predefined genres with different weights. Dong et al. [8] described a set of experiments to examine the effect of various attributes of Web genre for the automatic identification of the genre of Web pages. Four different genres were used in the data set (FAQ, News, E-Shopping and Personal Home Pages).

Rosso [22] explored the use of genre as a document descriptor in order to improve the effectiveness of Web searching. Author formulated three hypotheses: (1) Users of the system must possess sufficient knowledge of the genre. (2) Searchers must be able to relate the genres to their information needs. (3) Genre must be predictable by a machine applied algorithm because it is not typically explicitly contained in the document.

Design patterns and pattern languages came from architecture from the work of Christopher Alexander et al. [1]. From the mid sixties to mid seventies, Alexander et al. defined a new approach to architectural design. The new approach, centered on the concept of pattern languages, which is described in a series of books [1]. Alexander's definition of pattern is as follows: "Each pattern describes a problem, which occurs over and over again in our environment, and then describes the core of the solution to that problem, in such a way that you can use this solution a million times over, without ever doing it the same way twice."

According to Tidwell [27], patterns are structural and behavioral features that improve the applicability of software architecture, a user interface, a web site or something another in some domain. They make things more usable and easier to understand. Patterns are descriptions of best practices within a given design domain. They capture common and widely accepted solutions, their validity is empirically proved. Patterns are not novel, patterns are captured experiences and each their implementation is a little different.

Good examples are also the so-called “Web design patterns”, which are patterns for design related to the web. A typical example of a Web design pattern can be the Forum pattern (see Figure 1). This pattern is meant for designers who need to implement this element on an independent web page or as a part of another web page. The pattern describes key solution features without implementation details.

Generally, the design patterns describe a proven experience of repeated problem solving in the area of software solution design. From this point of view, the design patterns belong to key artifacts securing efficient reuse. While the design patterns have been proven in real projects, their usage increases the solution quality and reduces the time of their implementation.

There is a wide area of methods, which aim to detect objects related to patterns and extract their semantic details (e.g. Opinion extraction [15], News extraction [30], Web discussion extraction [14], Product detail extraction [20], Technical features extraction [25]).

### 3 MicroGenres

According to Wikipedia.org, the Genre is the division of concrete forms of art using the criteria relevant to the given form (e.g. film genre, music genre and literature genre). In all sectors of the arts, the Genres are vague categories without fixed boundaries and are especially formed by the sets of conventions.

Many artworks are cross-genre and employ and combine these conventions. Probably the most deeply theoretically studied Genres are the literary ones. It allows us to systematize the world of literature and consider it as a subject of scientific examination. We can find the term MicroGenre in this field. For example, in [19] the MicroGenre is seen as part of a combined text. This term has been introduced to identify the contribution of inserted Genres to the overall organization of text. The motivation of using the term “MicroGenre” is because it is used as a building block of the analytic descriptive system. On the other hand, Web design patterns are used more technically and provide means for good solution of Web pages. From our point of view the Web page is structured similarly to the literature text using parts which are relatively independent and have their own purpose (see fig. 2). For these parts we have chosen the

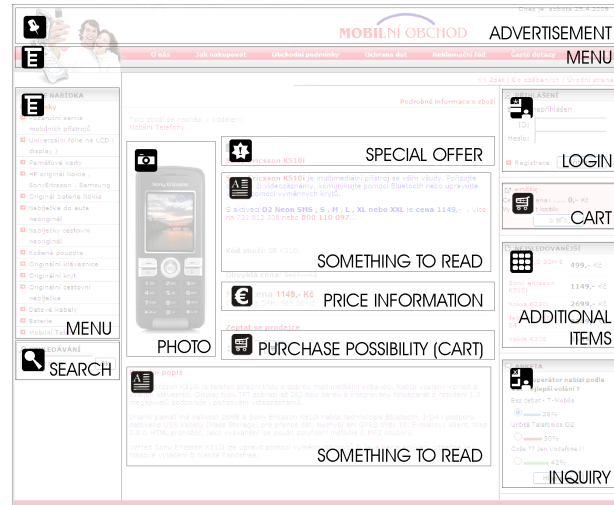


Figure 2. Structure of Web page

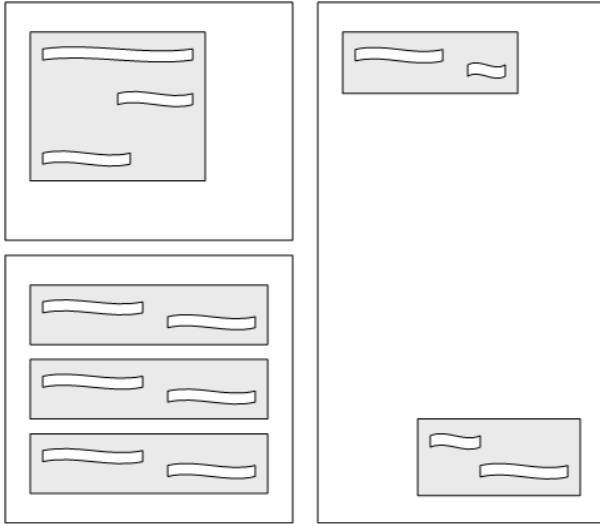
term “MicroGenre”. Contemporary Web pages are often very complex, nevertheless they can usually be described using several MicroGenres. This kind of description can be more flexible than the Genre description (which usually represents the whole page).

**Definition:** (Web) MicroGenre is a part of a Web page,

1. whose purpose is general and repeats frequently
2. which can be named intelligibly and more or less unambiguously so that the name is understandable for the Web page user (developer, designer, etc.)
3. which is detectable on a Web page using computer algorithm

The MicroGenre can, but does not have to, strictly relate to the structure of a Web page in a technical sense, e.g. it does not necessarily have to apply that it is represented by one subtree in the DOM tree of the page or by one block in the sense of the visual layout of the page. Rather it can be represented by one or more segments of a page, which form it together (see fig. 3).

**Remark:** MicroGenres are also context which encapsulates related information. In paper [13] we show the way we extract snippets from individual MicroGenres. We use these snippets in our Web application as an additional information for Web page description. The detection of MicroGenres can be considered in similar way as a first step for using Web information extraction methods (see [6]).



**Figure 3. MicroGenres formed by Web page segments**

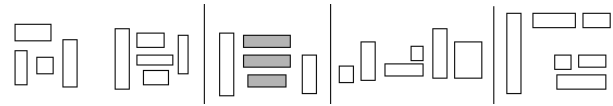
### 3.1 MicroGenre recognition

It follows the previous description that in order to be able to speak about the MicroGenre, this element has to be distinguishable by the user. From what attributes should the user recognize, if and what the MicroGenre there is in question? We work with up to three levels of view:

1. The first view is purely semantic in the sense of the textual content of a page. It does not always have to be a meaning in a sense of natural language such as sentences or paragraphs with a meaningful content. Logically coherent data blocks can still lack of grammars (see [31]).

For example, Price information can be only a group of words and symbols ('price', 'vat', symbol \$) of a datatype (price, number). For similar approach see [23].

2. The second view is visual in the sense of page perception as a whole. Here individual segments of perception or groups of segments of the page are in question. It is dependent on the use of colors, font and auxiliary elements (lines, horizontal and vertical gaps between the segments etc.) Approaches based on visual analysis of Web pages can be found in [4], [26].
3. The third view is a structural one in a technical sense. It is about the use of special structures, such as tables, links, navigation trees, etc. There are approaches based on analysis of the DOM tree and special structures as tables [16], [21].



**Figure 4. Gestalt principles (proximity, similarity, continuity, closure)**

The first view is dependent on the user's understanding of the text stated on a Web page. The second and third views are independent of this user ability. However, it can be expected that an Arabic or Chinese product page will be recognized also by an English-speaking user who does not have a command of those languages. It is determined by the fact that for the implementation of certain intentions there are habitual procedures which provide very similar results regardless of the language. On the other hand, if the user understands the page, he/she can focus more on the semantic content of the MicroGenre. For example, in the case of Product main info, the user can read what the product in question is, what its price is and on what conditions it can be purchased.

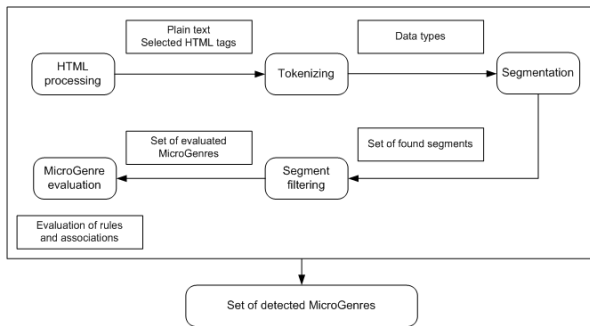
### 3.2 Pattrio Method

In our previous work, we have focused on Web design pattern detection. The research findings illustrated that it is useful to generalize the detection algorithm on all objects that provide the benefits to the users. These objects can be out of the scope of what is useful to the Web developers, but can be helpful in similar way like Genres. Therefore we have chosen the term "MicroGenre" for our further research.

In our approach, there are elements with semantic contents (words or simple phrases and data types) and elements with importance for the structure of the web page where the MicroGenre instance can be found (technical elements). The rules are the way that individual elements take part in the MicroGenre display. While defining these rules, we have been inspired by the Gestalt principles (see Figure 4 and [27]). We formulated four rules based on these principles. The first one (proximity) defines the acceptable measurable distances of individual elements from each other. The second one (closure) defines the way of creating of independent closed segments containing the elements. One or more segments then create the MicroGenre instance on the web page. The third one (similarity) defines that the MicroGenre includes more related similar segments. The fourth one (continuity) defines that the MicroGenre contains more various segments that together create the Web pattern instance. The relations among MicroGenres can be on various levels similar as classes in OOP (especially simple association and aggregation).

**Table 2. HTML tags - classification for analysis**

Types	Tags
Headings	H1, H2, H3, H4, H5, H6
Text containers	P, PRE, BLOCKQUOTE, ADDRESS
Lists	UL, OL, LI, DL, DIR, MENU
Blocks	DIV, CENTER, FORM, HR, TABLE, BR
Tables	TR, TD, TH, DD, DT
Markups	A, IMG
Forms	LABEL, INPUT, OPTION



**Figure 5. Object detection process**

The basic algorithm for detection of MicroGenres then implements the pre-processing of the code of the HTML page (only selected elements are preserved – e.g. block elements as table, div, lines, etc., see Table 2), segmentation and evaluation of rules and associations. The result for the page is the score of MicroGenres that are present on the page. The score then says what is the probability of expecting the MicroGenre instance on the page for the user. The entire process, including MicroGenre detection, is displayed in Figure 5.

The accuracy of the proposed method is about 80% (see [11]). Figure 6 shows the accuracy of the Pattrio method for three selected products (Apple iPod Nano 1GB, Canon EOS 20D, Star Wars Trilogy film) and for the *Discussion* and the *Purchase possibility* MicroGenres. We used only the first 100 pages for each product. We manually, and using Pattrio method, evaluated the pages using a three-degree scale:

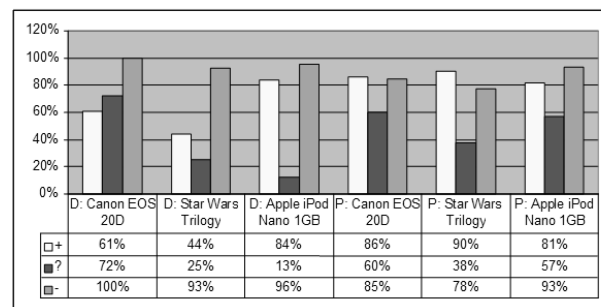
- + Page does not contain required MicroGenre.
- ? Unable to evaluate results.
- Page do not contain required MicroGenre.

Then we compared these evaluations. For example the first value 61% expresses the accuracy for the pages with Canon EOS 20D product where there was a discussion.

## 4 Experimentation Results and Discussions

We implemented a Web application with user interface connected to the API of different search engines (google.com, msn.com, yahoo.com and the Czech search engine jyx.cz). Users from a group of students and teachers of high schools and VŠB – Technical University Ostrava, Czech Republic were using this application for more than one year to search for everyday information. We have not influenced the process of searching in any way. The purpose of this part of the experiment was to view the World Wide Web using the perspective of users (as the search engines play key role in World Wide Web navigation). In the end we obtained dataset with more than 115,000 Web pages. After clean up, 77,850 unique Czech pages remained. For every single Web page we have performed the detection of sixteen MicroGenres. The page did not have to contain any MicroGenre, as well as it may have theoretically contained 16 MicroGenres (Price information, Purchase possibility, Special offer, Hire sale, Second hand, Discussion and comments, Review and opinion, Technical features, News, Enquire, Login, Something to read, Link group, Price per item, Date per item, Unit per item). The names of MicroGenres emerged from the discussions between us and students that took part in our experiments. They are therefore outcomes of social interaction.

We used such preprocessed dataset all the experiments. The Web pages were grouped by Web site. In the experiments, we attempted to visualize the structure and relations of Web sites referring to one specific topic. For visualization we have used Formal Concept Analysis.



**Figure 6. Accuracy of Pattrio method for detection of Discussion and Purchase Possibility MicroGenres - percentage of agreement between human and Pattrio method evaluation on sets of Web pages returned for different search queries**





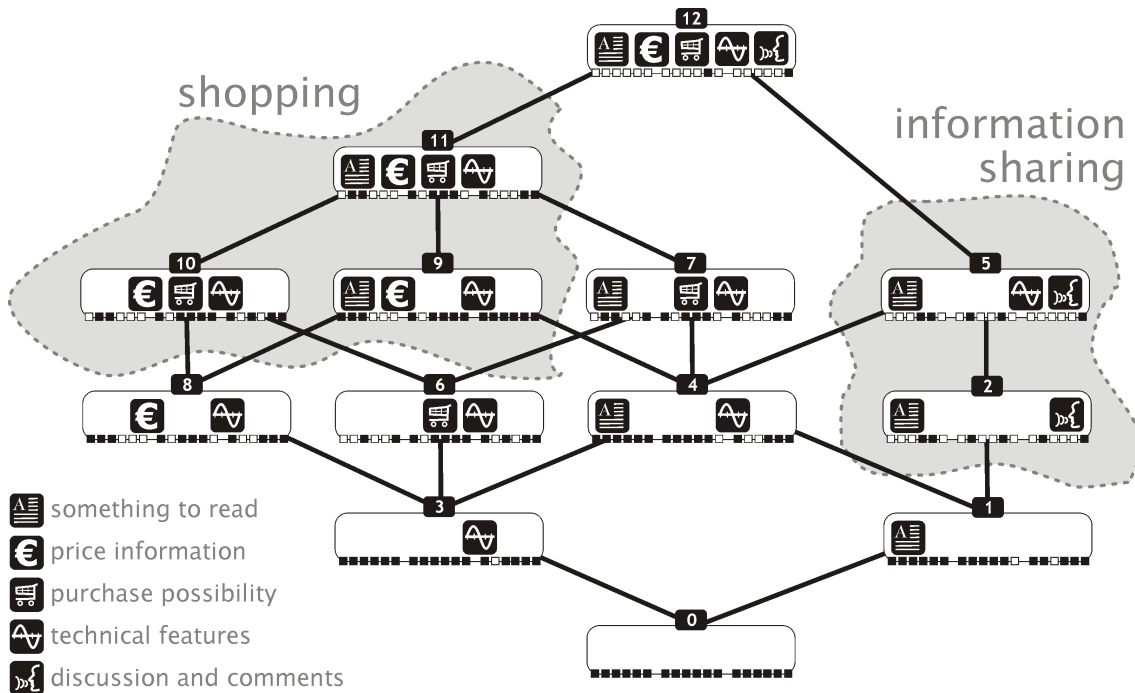


Figure 8. Part of lattice

ing for example the well known Latent Semantic Indexing. As a result of this paper, the Social group (and the association of Web site and Social group) can be also considered as another clustering criteria. Social groups are our approach to reduce the Web.

## 5 Conclusions and Future Work

This paper introduced the concept of MicroGenre in the context of Web page description. There is an interaction between users, developers and Web site owners hidden in different MicroGenres. For us, this is the social aspect of Web pages. It allows the different social groups to agree on the name and purpose of various Web page elements. By using the proposed methodology, it is possible to follow the evolution of communities and observe the expectancies, rules and behavior they share. From this point of view, Web 2.0 is only a result of the existence and interaction of these social groups. Our experiments illustrate that if we focus on Web sites and the Web page content they provide, we might come across a variety interesting questions. These questions may bear upon the Web sites' similarity and the similarity of social groups involved in these pages, which could formulate interesting future research directions. The identification of additional MicroGenres and design of specific algorithms for their detection will be the matter of our future research.

## References

- [1] Ch. Alexander: A Pattern Language: Towns, Buildings, Construction, Oxford University Press, New York (1977)
- [2] R. Belohlavek, V. Vychodil: What is a fuzzy concept lattice, Proceedings of the CLA, 3rd Int. Conference on Concept Lattices and Their Applications, pp. 34–45 (2005)
- [3] E. S. Boese: Stereotyping the web: Genre classification of Web documents, Colorado State University (2005)
- [4] D. Cai, S. Yu, J.-R. Wen, W.-Y. Ma: Extracting Content Structure for Web Pages based on Visual Representation, Fifth Asia Pacific Web Conference, pp. 406–417 (2003)
- [5] J. Chaker, O. Habib: Genre Categorization of Web Pages, Proceedings of the Seventh IEEE International Conference on Data Mining Workshops, pp. 455–464 (2007)
- [6] Ch. H. Chang, M. Kayed, M. R. Girgis, K. F. Shaalan: A Survey of Web Information Extraction Systems, IEEE Transactions on Knowledge and Data Engineering, vol. 18, pp. 1411–1428 (2006)

- [7] R. J. Cole, P. W. Eklund: Scalability in Formal Concept Analysis, *Computational Intelligence*, vol. 15, pp. 11–27 (1999)
- [8] L. Dong, C. Watters, J. Duffy, M. Shepherd: An Examination of Genre Attributes for Web Page Classification, *Proceedings of the Proceedings of the 41st Annual Hawaii International Conference on System Sciences*, pp. 133–143 (2008)
- [9] B. Ganter, R. Wille: *Formal Concept Analysis: Mathematical Foundations*, Springer-Verlag, New York (1997)
- [10] A. Kennedy, M. Shepherd: Automatic Identification of Home Pages on the Web, *Proceedings of the 38th Hawaii International Conference on System Sciences* (2005)
- [11] J. Kocibova, K. Klos, O. Lehecka, M. Kudelka, V. Snašel: Web Page Analysis: Experiments Based on Discussion and Purchase Web Patterns, *Web Intelligence and Intelligent Agent Technology Workshops*, pp. 221–225 (2007)
- [12] M. Kudelka, V. Snašel, O. Lehecka, E. El-Qawasmeh: Semantic Analysis of Web Pages Using Web Patterns, *Proceedings of the 2006 IEEE/WIC/ACM International Conference on Web Intelligence*, pp. 329–333 (2006)
- [13] M. Kudelka, V. Snašel, O. Lehecka, E. El-Qawasmeh, J. Pokorný: Web Pages Reordering and Clustering Based on Web Patterns, *SOFSEM 2008*, pp. 731–742 (2008)
- [14] H. Y. Limanto, N. N. Giang, V. T. Trung, J. Zhang, Q. He, N. Q. Huy: An information extraction engine for web discussion forums, *International World Wide Web Conference*, pp. 978–979 (2005)
- [15] D. Lee, O. R. Jeong, S. Lee: Opinion mining of customer feedback data on the web, *Proceedings of the 2nd international conference on Ubiquitous information management and communication*, pp. 230–235 (2008)
- [16] B. Liu, R. Grossman, Y. Zhai: Mining data records in Web pages, *KDD 2003*, pp. 601–606 (2003)
- [17] D. Lee, H. Seung.: Learning the parts of objects by non-negative matrix factorization, *Nature*, vol. 401, pp. 788–791 (1999)
- [18] T. Letsche, M. W. Berry, S. T. Dumais.: Computational methods for intelligent information access, *Proceedings of the 1995 ACM/IEEE Supercomputing Conference* (1995)
- [19] Martin, J. R.: *Text and clause: Fractal resonance*, *Text*, vol. 15, pp. 5–42 (1995)
- [20] Z. Nie, J. R. Wen, W. Y. Ma: Object-level Vertical Search, *Third Biennial Conference on Innovative Data Systems Research*, pp. 235–246 (2007)
- [21] A. Pivk, P. Cimiano, Y. Sure, M. Gams, V. Rajkovic, R. Studer: Transforming arbitrary tables into logical form with TARTAR. *Data & Knowledge Engineering*, vol. 60, pp. 567–595 (2007)
- [22] M. A. Rosso: User-based identification of Web genres. *JASIST (JASIS)* 59(7), pp. 1053–1072 (2008)
- [23] M. Santini: Description of 3 feature sets for automatic identification of genres in web pages, [www.nltg.brighton.ac.uk/home/Marina.Santini/three\\\_feature\\\_sets.pdf](http://www.nltg.brighton.ac.uk/home/Marina.Santini/three\_feature\_sets.pdf) (last access 2009-04-30)
- [24] V. Snašel, M. Polovincak, H. M. Dahwa, Z. Horak: On concept lattices and implication bases from reduced contexts, *Supplementary Proceedings of the 16th International Conference on Conceptual Structures, ICCS 2008*, pp. 83–90 (2008)
- [25] S. Schmidt, H. Stoyan: Web-based Extraction of Technical Features of Products, *Beiträge der 35. Jahrestagung der Gesellschaft für Informatik*, pp. 256–261 (2005)
- [26] Y. Takama, N. Mitsuhashi: Visual Similarity Comparison for Web Page Retrieval, *Web Intelligence*, pp. 301–304 (2005)
- [27] J. Tidwell: *Designing Interfaces: Patterns for Effective Interaction Design*, O’Reilly, pp. 0–596 (2005)
- [28] D. K. Van Duyne, J. A. Landay, J. I. Hong: *The Design of Sites: Patterns, Principles, and Processes for Crafting a Customer-Centered Web Experience*, Addison-Wesley Professional (2003)
- [29] M. Van Welie: *Pattern Library for Interaction Design*, [www.welie.com](http://www.welie.com), (last access 2009-04-30)
- [30] S. Zheng, R. Song, J. R. Wen: Template-independent news extraction based on visual consistency, In *Proceedings of the Twenty-Second AAAI Conference on Artificial Intelligence*, pp. 1507–1513 (2007)
- [31] J. Zhu, B. Zhang, Z. Nie, J. R. Wen, H. W. Hon: Web-page understanding: an integrated approach, *Conference on Knowledge Discovery in Data*, pp. 903–912. (2007)