# Sacrificing vs. Salvaging Coherence:
# An Issue for Adaptive Agents in Information Navigation

Kerstin Voigt
Department of Computer Science
California State University, San Bernardino
5500 University Parkway
San Bernardino, CA 92407
email: voigt@csci.csusb.edu

## Abstract

Currently numerous research efforts are under way with the objective of designing intelligent agents for information navigation which have the capability to adapt to the changing informational needs of their users. "Individualization" is to help users cope with the flood of information made available by large information repositories (e.g., the WWW). One challenge in the design of such agents for information navigation lies in striking a balance between the opposing goals of (1) low access cost for highly relevant information and (2) the preservation of coherence of the information. This paper proposes means of assessing access costs and "loss of coherence" in a manner that can help adaptive agents strike a viable compromise between both competing forces.

## 1 Introduction

Various efforts in research are currently under way in the area of adaptive software agents for information navigation according the informational needs of the individual user (e.g., see [Bow94, Mae94, Oos94, Voi95, Arm95, Lie94, Hol94]). Adaptive agents for navigating large information repositories such as the WWW are to help the user to better cope with the explosive amounts of information available on-line. Typically, a process of adaptation ("individualization") aims at facilitating the access of information that is perceived of particular relevance to the user. The discussion in the paper is biased towards information navigation and retrieval agents whose mode of operation is that of an information *browser*. In this context, the *objective of individualization the access and retrieval of information is to make more relevant items more cheaply accessible than items that are of little personal importance.* Thereby, the presentation of information is being adapted according to the information items' relevance to the individual user.

In the context of information browsing, cheaper access is equivalent to fewer options to scan and fewer clicks of hyperlinks or menu-buttons [DeM95]. Access cost, e.g., the number of required scans and clicks, is a function of the information structure. So reduction in access cost is typically accomplished by varying the structure. Therein lies a danger: access cost may be reduced for several relevant pieces of information, but the restructuring may have moved information out of its original context. Consequently, the information is cheaper to access, but out of context, its original coherence may be compromised (e.g., see [Nie90]).

```
[X] ai_textbook --> [ ] artificial_int
                    [ ] problem_solv  --> [ ] search
                    [ ] logic             [ ] informed search
                    [X] learning  --      [ ] game_play
                                  |
                           -->  [X] l_by_obs --> [ ] gen_mod_lea_ag
                                [ ] l_neural     [ ] induct_lea
                                                 [X] dec_tree_lea --> [ ] perf_elem
                                                                      [ ] expr_dec_tr
                                                                      [ ] dec_tr_from_ex
                                                                      [X] assess_perf
```
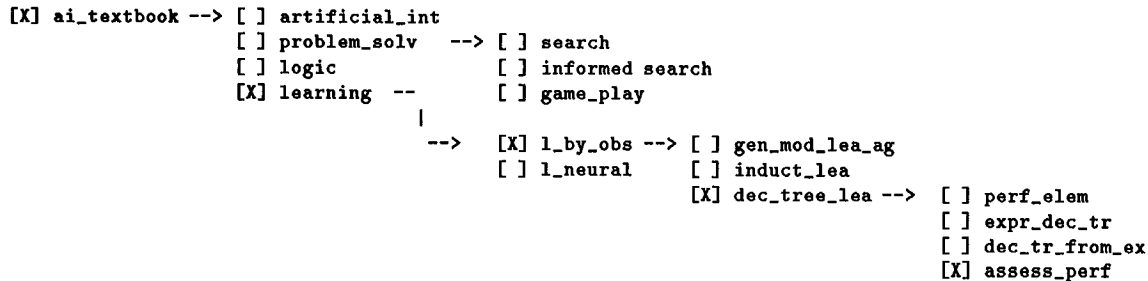
Figure 1: The original information structure, inspired by a textbook on artificial intelligence.

This paper elaborates on the *opposing forces of access cost reduction and preservation of coherence* when individualizing information navigation and retrieval. We will suggest measures of both forces and discuss insights gained from the study of an example. We will demonstrate how knowledge of the competing forces can help in trading off reduction in access cost of relevant information against severity of coherence loss.

**Information with "Expository" Structure.** In this paper we say that an organization of information items has an *expository* structure when the following is true: (1) the information items are linked hierarchically as a *tree*, (2) each "child item" elaborates specific details of the more general "parent" item, and (3) the children of each parent are arranged in a left-to-right order such that preceding pieces of information, in some sense or other, facilitate the understanding of the message to be conveyed in succeeding (sibling) pieces of information.

Prototypical examples of "expositorily structured" information are textbooks, technical articles, or (thoughtfully constructed) Web-pages. Much of our discussion will be illustrated with an example which simulates a textbook in the form it might appear in the World Wide Web information repository. The example is inspired by a textbook on artificial intelligence [Rus95].

**Example.** In Figure 1, the information structures rooted in ai_textbook are expository in the sense that information of the descendent levels elaborate details of concepts touched upon by the parent levels. For example, the chapter of problem solving (problem_solv) elaborates on issues of problem solving in sections on search (search), informed search (inf_search), and game playing (game_play). The subtopics on problem solving constitute part of the horizontal structure of Figure 1. Along the vertical (lateral) dimension we find, for instance, the information items ai_textbook (the common root), learning, l_by_obs ("learning by observation"), etc. It is obvious that considerable care has been taken in organizing the parts of the textbook in a fashion that presents the information as coherently as possible.

**Information structure and access cost.** "Browsing" for information is synonymous with traversing the information structure in order to locate individual informational items. Different locations of desired information account for varying traversal (or, access) costs. Large-scale general-purpose information repositories (e.g., the WWW) are intended for a

wide audience. Typically their design is not tailored to the informational needs of a particular user. Consequently, access costs tend to not correlate favorably with the user's personal interests.

It is the author of the book, article, or Web-site who determines the organization of the information within the larger context of a information repository. Within the repository, this structure is fixed, and without any further adaptive means, any interested user has to traverse information along the same original path. For users with pronounced interest in some piece of information, it may prove too cumbersome to retrace this fixed sequence of steps (e.g., traversing hyperlinks, flipping through and scanning the same book/article pages) each and every time.

In this paper we introduce our notion of "access cost" of information within an expositorily structured information repository. In addition, we suggest measures of "coherence loss " due to information restructuring. We discuss the antagonism between access cost and coherence with regard to the example introduced in Figure 1.

**Individualization through information restructuring.** Imagine that the above "book" is accessible via the WWW, and that the above hierarchical information structure allows the browsing of the book contents. Suppose a user wants to access information of performance assessment in decision-tree learning. It makes sense for the user to first read up on performance elements in decision tree learning, and it certainly helps to first familiarize oneself with decision-tree learning beforehand. Similarly, it would probably be best if the user read up on general notions of artificial intelligence and learning before going into the details of decision tree learning. Notice that the hierarchical structure shown in Figure 1 organizes the information in just the right fashion. Working one's way top-down and left-to-right, starting at the top-level item ai_textbook, the information appears in an order which presents the information coherently.

The user may hold a strong interest in assessing performance of decision trees, and may want to access this topic repeatedly. Fast access of this particular information item is desired. A fairly naive, but simple approach to facilitating retrieval of information on access_ perf could consist in simply placing the item first in the top-level menu. The organization which would result is shown below:

```
[X] ai_textbook --> [X] assess_perf
                    [ ] artificial_int
                    ...
                    [ ] learning
```

**Loss of Coherence.** Certainly, the access cost of the relevant piece of information is minimal. On the other hand, most would not advocate this form of adapted organization as the one to be preferred over all others. Why not? The most disconcerting aspect of the above structure is the fact that the information on performance assessment appears entirely out of context, without mention of any higher-level concepts that lead up to the more detailed pieces of information.

## 2 Reduction of Access Cost

Graphical information browsers have greatly simplified the procedures that need to be

| Information Item | Relevance Score |
|---|---|
| ai_textbook | (1.00) |
| artificial_int | 0.17 |
| prob_solv | 0.30 |
| ... | ... |
| gen_mod_lea_ag | 0.99 |
| dec_tree_lea | 0.31 |
| perf_elem | 0.46 |
| expr_dec_tr | 0.48 |
| dec_tr_from_ex | 0.95 |
| assess_perf | 0.98 |

Table 1: Information items and their relevance scores (for some imaginary user.)

carried out by the user in order to locate information. Point-and click are the predominant user actions. Yet even these can become too much effort for today's demanding information consumers [DeM95]. Users want to home in on their targets with the least number of clicks, where clicking should be preceded by the least possibly amount of scanning and selecting options. We take the position that the effort of scanning is the larger effort, and tends to cause the greater discomfort to the user. This view is reflected in our definition of *access cost for information item* $x$ below:

$$access\_cost(x) = \begin{cases} 0 & \text{if } x \text{ is root} \\ \frac{posit(x)/nsibs(x)+access\_cost(parent(x))}{maxcost} & \text{otherwise} \end{cases}$$

($nsib(x)$ is the number of siblings of $x$ (incl. $x$); $posit(x)$ is the position of $x$ among its siblings (viewed left to right); $maxcost$ is the maximal access cost among all items in the structure.)

In words, the access cost of an information item is the number of items that need to be scanned (called *lateral_cost*) within each substructure along the path from the root of the structure to the target item, divided by the cost of the most "inaccessible" item ($maxcost$). It is assumed that the cost of "clicking" an item is zero, and the user effort is solely due to having to scan (and possibly scroll) the presented options. This measure also reveals some of the shortcomings of flat hotlists. Surely, once the desired item is located in the list, a single click will access the information; however, the number of items that would need to be scanned may be cumbersomely large.

For the example in Figure 1 we obtain 0.88 as the access cost of item assess_perf. It is computed by $(4/4 + 3/3 + 1/2 + 4/4)/4.0 = 3.5/4.0 = 0.88$ (The largest possible access cost ($max\_cost$) of any item in the example is 4.) We see that the item assess_perf is the most expensive item to access in the given information structure.

**Restructuring Information by "Relevance Scores".** Let us assume that for each individual user, information items are assigned a value between 0.0 and 1.0 which indicates the degree of relevance of each item to the user (e.g., see [Voi95, Kra94, Jur94]). For the above example, assume that the scores in Table 1 have been established.

For the remaining discussion in this paper, notice that the item assess_perf is ranked as one of the most relevant items. Its relevance score of 0.98 is exceeded only by item

130

```
[X] ai_textbook --> [X] learning   --> [X] l_by_obs ==> [ ] gen_mod_lea_ag
                    [ ] logic          [ ] l_neural     [X] dec_tree_lea  --> [X] assess_perf
                    [ ] art_intell                      [ ] induct_lea        [ ] dec_tr_from_ex
                    [ ] prob_solv                                             [ ] expr_dec_tr
                                                                              [ ] perf_elem
```

Figure 2: Information structure after "sorting" by relevance score.

gen_mod_lea_ag with a score of 0.99. As before, we focuse on the fate of assess_perf in the course of individualization by restructuring.

Expository information structures can be reorganized by altering the presentation of their items along either the vertical or the lateral dimension. Since the originator of the information chose the structuring for good reason (one hopes), one can expect restructuring to be the cause a certain loss of coherence with regard to the overall information. Having examined various structured repositories of information, we have noticed that changes along the vertical dimension tends to lead to greater loss in coherence more consistently than changes to the lateral structure. In our current research we decided to refrain from altering vertical structure at all, and rather focus on how the adaptation of information navigation would fare with only lateral restructuring of information sites.

We have investigated the effects of lateral restructuring of the original expository information structure by *recursively sorting the subtrees rooted in each information item in ascending order of the largest relevance scores among the item's subtrees.* Sorting of the information structure in Figure 1 yields the new structure in Figure 2.

When comparing the both structures, one finds that the positions of items within their substructures have shifted. This is, for instance, reflected in the reduced access cost of item assess_perf. The cost before restructuring, was 0.88; the new improved cost is 0.42. The latter results from $1/4 + 2/3 + 1/2 + 1/4)/4.0 = 1.67/4.0 = 0.42$. Thus, a reduction by more than 50% has been achieved.

## 3 Salvaging Coherence

So far we have been able to successfully reduce the access cost of relevant items by restructuring the information. Would the user (whose informational needs are captured by our example) really prefer to navigate the structure provided in Figure 2 over the original structure in 1? A critical look at the former reveals at least two shortcomings. (1) The first information the user is confronted with is information on 'learning'; one wonders whether out of the blue, i.e. without the context of AI, and some general knowledge of issues in problem solving, the user will be able to comprehend the topic of learning sufficiently. Or, (2) will the user really be able to appreciate the information assess_perf ("performance assessment for decision trees") without at least being given the prior option to look up information on decision trees as performance elements. In short, while a relevant topic such as assess_perf is now retrieved more cheaply, the information is presented in an "out-of-context" fashion that threatens to make it less coherent overall.

*In order to be able to tell whether information restructuring has lead to genuine improve-*

*ment from the user's point of view, we need to not only look at access costs, but also at the degree of coherence that may have been sacrificed in the process.* One needs to be careful to not have reduction of access costs be outweighed by an unacceptable loss coherence. In the next section, we will suggest a method by which it is possible to quantify such loss of coherence. Note, that some objective measure of "coherence" is the prerequisite for any adaptive agent's attempt at striking a proper balance between both opposing forces.

**Measuring loss of coherence.** What causes a collection of information items to form a coherent whole? We have stated earlier, that much of the secret of coherence lies in the manner in which individual pieces of information are organized. Coherence is typically enhanced by first introducing general concepts and notions, and then moving on to the more specific cases ("vertical"). For example, it helps to be familiarized with the basic concerns in the field of machine learning, before moving on to the specifics of decision tree learning. Coherence is also furthered by organizing information such that details are presented in their "logic" order ("lateral"). For example, some background on AI and problem solving may want to be presented prior to topics on machine learning. Or, before you discuss the assessment of performance of some component, this component better be introduced before hand. Whenever such "natural" order of information is violated in its presentation, loss of coherence is likely.

**Penalty based on violation of precedence constraints among information items.** Who determines and realizes the "natural" order among pieces of information in structured information repositories? It is usually the *author* of the information (the book, the article, the web-page etc.) Under ideal circumstances, the author knows exactly the kind of message that is to be conveyed, and will structure the information in a fashion that presents the information in the most coherent way possible. Naturally, given a collection of information items, not all of the items are tightly linked, and some degree of freedom may exist with regard to their order of presentation. Yet, for certain pieces it may be absolutely mandatory that they appear within the context of other specific pieces of information. While it has not been customary for authors explicitly state such *precedence constraints* on information, providing such constraints would enable an agent to quantify the otherwise vague notion of coherence. Developing this idea further, let us assume that the author of the original structure in Figure 1 provided the following three precedence constraints:

> (PC1) Artificial intelligence (`artificial_int`) <u>before</u> `learning`.
> (PC2) Problem solving (`problem_solv`) <u>before</u> `learning`.
> (PC3) Performance elements (`perf_elem`) <u>before</u> assess performance (`assess_perf`).

We see that in Figure 1, all precedence constraints are satisfied. In contrast, Figure 2 violates all three of them. We suggest to quantify loss of coherence as a function of violations of precedence constraints. We not only count the violations, but also weigh them against the total number of theoretically possible violations among all siblings of the violating item. Given an information item with $n$ siblings, the current order of siblings could violate at most $n * (n - 1)/2$ precedence constraints. Thus, one violation of precedence contributes $1/(n * (n - 1)/2) = 2/(n * (n - 1))$ to the overall loss of coherence.

We measure loss of coherence by a *penalty score* (*penalty*1) which is defined as the following function of an information structure ($IS$):

$$penalty1(IS) = \sum_{x \in IS} \sum_{y \in siblings(x)} [viol\_0\_1(x,y) * 2/(nsibs(x) * (nsibs(x) - 1))]$$

The penalty for the original structure is zero. The restructured information in Figure 2 carries a penalty of 0.5 (because $2/12 + 2/12 + 2/12 = 6/12 = 0.5$).

In the above penalty measure, each violation of precedence contributes evenly (here 0.17), to the total of 0.5. This is equivalent to saying that all violations of precedence constraints are equally damaging to the coherence of the information. Yet, intuitively, the violation of precedence among the topics assess_perf and perf_elem seems more severe than the violation involving artificial_int and learning. It is conceivable that a reader understands at least the basics of learning without necessarily having been giving a general AI background (although the later would be advantageous). On the other hand, information about performance assessment in decision tree learning must inevitably appear rather incomprehensible without knowledge of decision trees as performance elements. Loss of coherence is not adequately captured by a measure which is not capable of distinguishing severely damaging violations of precedence constraints from violations that are only mildly damaging. A more discriminating measure of coherence loss will be introduced next.

**How to salvage coherence.** Why does it seem detrimental to the information when the topic of learning is presented before the topics of artificial intelligence (or, problem solving)? Why should the damage be greater when the discussion of performance assessment appears before the definition of the actual performance elements that are subject to assessment? Answers to these questions seem to lie in the level of generality/ specificity of the information items involved. In particular, we believe violation of precedence matters less among more general pieces of information than among information items which capture contents of more specific detail. Generalities can, at least to a certain degree, be understood in isolation. Issues about machine learning, when kept sufficiently general, do not require much prior knowledge of artificial intelligence or problem solving. Violations of precedence constraints (PC1) and (PC2) should be given lesser weight than the violation of (PC3). Notice that the information referred to in the latter constraint are at a level of greater specificity. Also realize that greater specificity is found at greater depths of the vertical dimension in the expository information structure. Consequently, a more discriminating penalty measure for coherence loss may be obtained *by weighing each violation by the information's level of specificity.* The less specific, the smaller a weight will be assigned. Accordingly, the more discriminating penalty measure (*penalty2*) is defined as:

$$penalty2(IS) = \sum_{x \in IS} \frac{\sum_{y \in siblings(x)} [viol(x,y)*2/(nsibs(x)*(nsibs(x)-1))]}{2^{(maxd-lev(x))}}$$

(The constant *maxd* is the maximal depth of the expository; lev(x) is the level at which item $x$ appears; the root is at level 0.)

This second measure of coherence loss differs from *penalty1* in that the violations counted for an item are divided by 2 to the power of the difference between the maximal level in the structure and the level of the item considered. The more specific the information item, the smaller this term becomes (smallest: $2^0 = 1$), and the greater value of *penalty2*. This new measure not only captures our intuition more appropriately, but is also consistent with standard information theory (e.g., [Cov91]). In a hierarchical information structure, items at greater depth tend to have greater "information content". The more information content

```
[X] ai_textbook --> [X] learning   --> [X] l_by_obs ==> [ ] gen_mod_lea_ag
                    [ ] logic          [ ] l_neural    [X] dec_tree_lea --> [ ] perf_elem
                    [ ] art_intell                     [ ] induct_lea       [ ] dec_tr_from_ex
                    [ ] prob_solv                                           [ ] expr_dec_tr
                                                                            [X] ass_perf
```

Figure 3: Information structure after "salvaging" some amount of coherence.

is presented in the wrong order, the greater one should expect the damage to the overall information to be.

For the restructured information in Figure 2 the revised penalty (*penalty*2) evaluates to 0.21. It is computed by $2/12*1/2^3 + 2/12*1/2^3 + 2/12*1/2^0 = 0.21$ What matters primarily, is not the fact that the new measure returns a value (0.21) that is lower than the value of *penalty*1 (0.5), but the fact that in *penalty*2 different violations contributed different amounts to the overall score. This knowledge can be exploited in striking a compromise between the desired reduction of access cost and the potential loss of coherence.

**Undo selected violations to salvage coherence.** Some of the coherence lost in the restructured information displayed in Figure 2 can be salvaged by selectively undoing some violations and thereby partially restoring the original state in Figure 1. *Coherence is restored where it matters most*, namely at the more specific levels of the information structure. Undoing these low-level violations will remove the largest contributors to the *penalty*2 measure. In our example, it is recommended to restore the original placements of items **perf_elem** and **assess_perf**. Low-level information on performance elements and and their assessments is now presented in the original coherent form. The entire information structure after coherence salvaging is shown in Figure 3. The value of *penalty*2 for this structure is now 0.04 (compared to the previous 0.21). The gain in coherence comes at the cost of a slightly increased cost of accessing **assess_perf**; the new access cost has increased to 0.69 (after having previously been reduced from 0.88 to 0.42). *Overall, noticeable reduction in access cost for a relevant information items has still been accomplished, yet it has been possible to keep the penalty for coherence loss low.*

## 4 Discussion

In this paper we have elaborated on two competing forces in individualizing the navigation of information repositories which are structured in an expository manner. We have introduced means for quantifying information access cost and loss of coherence due to altering information structure, and we have developed a measure of coherence loss which discriminates between intuitively more and less severe violations of coherence. We have shown, how the ability to quantify and trade off access cost and loss of coherence can help in selectively salvaging some of the lost coherence between restructured information items. In more general terms, the problem is one of *multi-objective optimization* (here: striking the "best" compromise between two opposing objectives), and it seems well-suited to the application of genetic algorithms (e.g., [Gol89]). Genetic algorithms (GAs) have the advantage of being general enough to, in principle, apply to any type of solution structure (e.g., to tree-like structures

of information sites, but also to more general graph-like networks of sites). Furthermore, from the view-point of computational cost, GAs appear to be viable in the context of the current system where information restructuring takes place *off-line*, that is, between user sessions. Browser idle-times can be exploited by GA's in order to evolve new, "fitter", more up-to-date structures of information sites.

**Source of coherence constraints.** Knowledge of precedence constraints between information items is key to our system's ability to judge coherence. Such constraints are not customarily provided by the originators of information. Requesting the explicit posting of such constraints may be of some burden to the authors. Nevertheless providing such information might pay off in view of how such knowledge can help an adaptive agent in individualizing information navigation. In the long run, however, the utility of an adaptive information agent should not rest on the shoulders of the information authors. We plan to gear some of our continuing research towards the *automatic elicitation of coherence constraints*.

Work on *automated text-understanding* could greatly contribute if techniques were developed to automatically determine precedence requirements among collections of information items. Yet, much less sophisticated approaches may suffice to provide some useful hints on intended coherence. *Syntactic analysis of hypertext* documents appears to be a candidate. For instance, items organized in an ordered list are indicative of the fact that orders matters. The extraction and posting of coherence constraints will serve to prevent loss of coherence due to arbitrary information ordering. *Analyses of browsing histories* across users may reveal coherence information as well. Imagine, for example, a scenario where users view some item $A$ (briefly), and then, fairly consistently, view a specific other item $B$, before returning to spend more time on the first item. This observation could lead to the justifiable conclusion that item $B$ holds information whose prior viewing is helpful or even necessary for the processing of item $A$.

**Current limitations.** The research presented in this papers is work in progress. It has has several obvious limitations which we plan to address in the near future. As mentioned before, one limitation is inherent in the assumption that information structures are tree-structured. The sites in the most popular information repository to date – the World Wide Web – are not necessarily organized as trees. Therefore, we need to generalize our techniques such that the more realistically encountered graph-like information structures can be accommodated. By adopting general techniques like genetic algorithms to compute trade-offs between competing objectives, we expect to overcome such limitations with relative ease.

We have further limited ourselves to altering information structure along the lateral, but not the vertical dimension. Continued work will show whether lifting this restriction is beneficial, and how this could be accomplished without unacceptable losses in coherence.

# References

[Cas95]  J. Castro. Just Click to Buy. *TIME Magazine, Special Issue*, Spring 1995.

[Bow94]  C.M. Bowman, et.al. Scalable Internet Resource Discovery: Research Problems and Approaches. *CACM*, August 1994.

[Gol89]  D. E. Goldberg. *Genetic Algorithms in Search, Optimization, and Machine Learning*. Addison-Wesley, 1989.

[Dou94]  D. Dougherty, et.al. *The Mosaic Handbook for the X Window System*. O'Reilly & Associates, 1994.

[Kra94]  D.H. Kraft and C. Barry. Relevance in Textual Retrieval. In: *Relevance*, Papers from the 1994 AAAI Fall Symposium, Technical Report FS-94-02, 1994.

[Lie94]  H. Lieberman. Demonstrational Techniques for Instructible User Interface Agents. In: *Software Agents*, Papers from the 1994 AAAI Spring Symposium, Technical Report SS-94-03, 1994.

[Ber96]  H. Berghel. The Client's Side of the World Wide Web. *CACM*, January 1996.

[Jur94]  I. Jurisica. How to Retrieve Relevant Information. In: *Relevance*, Papers from the 1994 AAAI Fall Symposium, Technical Report FS-94-02, 1994.

[Oos94]  K.A. Oostendorp, W.F. Punch, and R.W. Wiggins. A Tool for Individualizing the Web. In *Second International WWW Conference '94: Mosaic and the Web*, Chicago, 1994.

[Kro94]  E. Krol. *The Whole Internet, User's Guide and Catalog*. O'Reilly & Associates, 1994.

[Bal95]  M. Balabanovic and Y. Shoham. Learning Information Retrieval Agents: Experiments with Automated Web Browsing. In: *Information Gathering from Heterogeneous, Distributed Environments*, Papers from the 1995 AAAI Symposium, Technical Report SS-95-08, 1995.

[Mae94]  P. Maes. Agents that Reduce Work and Information Overload. *CACM*, July 1994.

[DeM95]  N. deMause. The Virtues of Restraint. *Web Week*, September 1995.

[Nie90]  J. Nielsen. The Art of Navigating Through Hypertext. *CACM*, March 1990.

[Hol94]  R. C. Holte and C. Drummond. A Learning Apprentics for Browsing. In: *Software Agents*, Papers from the 1994 AAAI Spring Symposium, Technical Report SS-94-03, 1994.

[Arm95]  R. Armstroing, D. Freitag, T. Joachims, and T. Mitchell. WebWatcher: A Learning Apprentice for the World Wide Web. In: *Information Gathering from Heterogeneous, Distributed Environments*, Papers from the 1995 AAAI Symposium, Technical Report SS-95-08, 1995.

[Rus95]  S. Russell and P. Norvig. *Artificial Intelligence: A Modern Approach*. Prentice Hall, 1995.

[Cov91]  T.M. Cover and J. A. Thomas. *Elements of Information Theory*. John Wiley, 1991.

[Voi95]  K. Voigt. Reasoning about Changes and Uncertainty in Browser Customization. AAAI Fall Symposium 1995, Working Notes of AI Applications to Knowledge Navigation and Retrieval, November 1995.