

This is an author-prepared version of a journal article that has been published by Elsevier.
The original article can be found at doi:10.1016/j.jss.2007.12.815

Architectural Knowledge Discovery with Latent Semantic Analysis: Constructing a Reading Guide for Software Product Audits

Remco C. de Boer and Hans van Vliet

The Journal of Systems and Software, Vol. 81, Issue 9, September 2008, pp. 544–550

Architectural Knowledge Discovery with Latent Semantic Analysis: Constructing a Reading Guide for Software Product Audits¹

Remco C. de Boer*, Hans van Vliet

*VU University Amsterdam, Dept. of Computer Science, De Boelelaan 1081a,
1081HV Amsterdam, the Netherlands*

Abstract

Architectural knowledge is reflected in various artifacts of a software product. In a software product audit this architectural knowledge needs to be uncovered and its effects assessed in order to evaluate the quality of the software product. A particular problem is to find and comprehend the architectural knowledge that resides in the software product documentation. In this article, we discuss how the use of a technique called Latent Semantic Analysis can guide auditors through the documentation to the architectural knowledge they need. We validate the use of Latent Semantic Analysis for discovering architectural knowledge by comparing the resulting vector-space model with the mental model of documentation that auditors possess.

Key words: Software architecture, architectural knowledge, knowledge discovery, latent semantic analysis, software product audit.

1 Introduction

The architectural design of a software product and the architectural design decisions taken play a key role in software product audits. Architectural design decisions and their rationale provide, for instance, insight into the trade-offs

* Corresponding author. Tel.: +31 20 59 87767; fax: +31 20 59 87728

Email addresses: remco@cs.vu.nl (Remco C. de Boer), hans@cs.vu.nl (Hans van Vliet).

¹ This article has been based on earlier work by the authors, presented at the 6th Working IEEE/IFIP Conference on Software Architecture in January 2007 in Mumbai, India (de Boer and van Vliet (2007)).

that were considered, the forces that influenced the decisions, and the constraints that were in place. The architectural design that is the result of these decisions allows for comprehension of such matters as the structure of the software product, its interactions with external systems, and the enterprise environment in which the software product is to be deployed. Following a recent trend in software architecture research (e.g., (Bosch, 2004; Jansen and Bosch, 2005; Kruchten et al., 2006; van der Ven et al., 2006)) we refer to the collection of architectural design decisions and the resulting architectural design as ‘architectural knowledge’.

For a given software product there is no single source that contains or provides all relevant architectural knowledge. Instead, architectural knowledge is reflected in various artifacts such as source code, data models, and documentation. A complicating factor in distilling relevant architectural knowledge from software product documentation is the fact that there are often many different documents. Each of these documents is tailored to specific stakeholders and different documents can therefore reflect architectural knowledge at different levels of abstraction. A high-level project management summary, for instance, will reflect architectural design decisions and their effects differently than a document describing detailed technical design.

The ISO/IEC 14598-1 international standard (ISO/IEC, 1999) defines a software product as ‘the set of computer programs, procedures, and possibly associated documentation and data’. Quality is defined as ‘the totality of characteristics of an entity that bear on its ability to satisfy stated and implied needs’, while quality evaluation is ‘a systematic examination of the extent to which an entity is capable of fulfilling specified requirements’. Consequently, when we refer in this article to a software product quality audit - i.e., an audit in which the quality of a software product is evaluated - we refer to ‘the systematic examination of the extent to which a set of computer programs, procedures, and possibly associated documentation and data are capable of fulfilling specified requirements’.

We have conducted a study at a company that has broad experience in performing software product audits. This company conducts independent quality audits of software products. Its customers range from large private companies to governmental institutions. In this study we have investigated the use of architectural knowledge in software product audits. To this end we observed an audit that was being conducted for one of the company’s customers. We attended and observed the audit team meetings and had discussions with the audit team members on their use of architectural knowledge in the audit. In addition, we held more general interviews on this topic with five employees who had been involved in various audits, two of whom were also directly involved in the observed audit. The interviewed employees possess different levels of experience and have different focal points when conducting an audit.

The problem of finding relevant architectural knowledge sketched above corresponds to a problem that is perceived by all auditors as being difficult to deal with. In short, the auditors need a reading guide that guides them through the documentation.

In this article we outline the problem of discovering architectural knowledge in software product documentation and present a technique that can be used to alleviate this problem. This technique, Latent Semantic Analysis, uses a mathematical technique called Singular Value Decomposition to discover the semantic structure underlying a set of documents. We employ this latent semantic structure to guide the auditors through the documentation to the architectural knowledge needed. A comparison of the discovered semantic structure with the ideas auditors have of software product documentation shows that Latent Semantic Analysis produces a good approximation of the auditors' mental models.

The remainder of this article is organized as follows. The next section discusses the use of architectural knowledge in software product audits based on our observations in the case study we conducted. Section 3 presents Latent Semantic Analysis (LSA) and its mathematical background. Section 4 discusses the application of LSA to a set of documents that contain software product documentation and shows how we can employ the semantic structure uncovered by LSA to guide the auditor to relevant architectural knowledge. In Section 5 we validate the LSA results through a comparison with auditors' mental models of software product documentation. Section 6 contains a discussion on related work regarding the application of LSA to similar problems as well as related work in the area of research into architectural knowledge. Section 7 outlines research areas that are still open for further study. In Section 8 we sketch the use of Architectural Knowledge Discovery in a broader scope, and Section 9 contains concluding remarks on this article.

2 Architectural Knowledge in a Software Product Audit

In a software product audit, two types of architectural knowledge can be distinguished. On the one hand there is architectural knowledge pertaining to the *current state* of the software product; this knowledge reflects the architectural decisions *made*. On the other hand there is architectural knowledge pertaining to the *desired state* of the software product; this knowledge reflects the architectural decisions *demand*ed (or expected). It is the auditor's job to compare the current state with the desired state.

In order to perform a comparison of current state and desired state, the auditor has to have a firm grasp on both types of architectural knowledge. A

common method to structure the architectural knowledge of the desired state is to define a number of review criteria. These criteria can be phrased as (architectural) decisions, and are a combination of the wishes of the customer and the expertise of the auditor. An example of such a criterion might be ‘All errors in the software are written to a log. Each log entry contains enough information to determine the cause of the error.’. A software product audit consists of a comparison of these review criteria against the current state of the software product.

The ‘current state’ architectural knowledge of the software product is reflected in different artifacts, in particular in source code and the accompanying documentation. Some architectural knowledge, for instance alternative solutions that were considered but have been rejected, might not be explicitly captured in these artifacts at all. This architectural knowledge is left implicit and lives only in the heads of its originators. Particular methods that are used to distill the architectural knowledge needed from these three sources - source code, documentation, and people - are:

- scenario analysis,
- interviews,
- source code analysis, and
- document inspection.

Both interviews and scenario analysis are techniques to elicit implicit architectural knowledge from people’s minds, and consequently require extensive interaction with the software product supplier. Source code analysis and document inspection, however, are performed using only the artifacts that have been delivered as part of the software product. In terms of availability of resources, the latter two are hence to be preferred. In the remainder of this article we will focus on document inspection in particular. A typical first use of the architectural knowledge reflected in the documentation is for auditors to familiarize themselves with a software product. Once a certain level of comprehension has been attained, the documents are used as a source of evidence for findings regarding the software product quality.

While document inspection is an important method in a software product audit, it can also be a difficult method to use. The difficulty of performing document inspection lies in the sheer amount of documentation that accompanies most software products. Auditors are swamped with documentation, and there is no single document that contains all architectural knowledge needed. Moreover, a ‘reading guide’, which tells the auditors which information can be found where, is usually not available up front. Auditors need to fall back on interviews, a resource-intensive technique, to gain an initial impression of the organization of architectural knowledge in the documentation.

In general, from the interviews held we learned that auditors have three major questions regarding software product documentation and the architectural knowledge contained in it. These three questions are:

- (1) Where should I start reading?
- (2) Which documents should I consult for more information on a particular architectural topic?
- (3) How should I progress reading? In other words, what is a useful ‘route’ through the documentation to gain a sufficient level of architectural knowledge?

From the above it should be clear that the auditors who perform a software product audit would greatly benefit from tools and techniques that can direct them to relevant architectural knowledge. We refer to the goal of such tools and techniques as ‘Architectural Knowledge Discovery’ (de Boer, 2006). A core capability of Architectural Knowledge Discovery is the ability to grasp the semantic structure, or meaning, of the software product documentation. Employing this structure transforms the set of individual texts into a collection that contains architectural knowledge elements and the intrinsic relations between them. A technique that can be deployed to support the discovery of directions to relevant architectural knowledge is Latent Semantic Analysis.

3 Latent Semantic Analysis

A method that can be used to capture the meaning of a collection of documents is the construction of a vector-space model. Vector-space models are based on the assumption that the meaning of a document can be derived from the terms that are used in that document. In a vector-space model, a document d is represented as a vector of terms $d = (t_1, t_2, \dots, t_n)$, with t_i ($i = 1, 2, \dots, n$) being the number of occurrences of term i in document d (Letsche and Berry, 1997).

Figure 1 depicts a matrix based on the vector-space model constructed for three texts that were taken from the documentation of a software product. The three texts used are representative selections from a use case definition (UC), a service specification (SVC), and an architecture description (ARCH). Together, the three document vectors corresponding to these three texts contain approximately 90 distinct terms, excluding stopwords. This so-called term-document frequency matrix represents the number of occurrences of each of these terms in each of the three documents. The original document vectors are hence extended with terms that did not occur in the document itself, but do occur in one of the other texts. In these extended document vectors t_i is set to 0 if term i does not occur in the document. The cutout shows the

	UC	SVC	ARCH
<i>[redacted]</i>	0	0	0
<i>[redacted]</i>	0	0	0
<i>[redacted]</i>	0	0	0
domain entity	0	2	0
service	0	3	1
SOA	0	0	3
system	0	0	4
use case	1	0	0
business object	8	3	0
<i>[redacted]</i>	0	0	0
<i>[redacted]</i>	0	0	0
<i>[redacted]</i>	0	0	0

domain entity	0	2	0
service	0	3	1
SOA	0	0	3
system	0	0	4
use case	1	0	0
business object	8	3	0

Fig. 1. Term-document frequency matrix based on the vector-space model for three software product documentation excerpts.

exact number of occurrences of six terms in the respective texts. For reasons of non-disclosure, the terms ‘domain entity’, ‘use case’, and ‘business object’ have been substituted for the product-specific terminology.

Although the vector-space model in Fig. 1 captures some of the semantics of the three texts, parts of the underlying semantic relationships are not represented very well. Based on Fig. 1 we can, for instance, only conclude that in theory neither the use case definition nor the service specification has anything to do with the term ‘SOA’ (an abbreviation for ‘Service Oriented Architecture’). In practice, however, we would expect at least some relevance of the term ‘SOA’ in the context of a *service* specification. Latent Semantic Analysis allows us to find and exploit such underlying, or latent, semantic relationships.

Latent Semantic Analysis (LSA) relies on a mathematical technique called Singular Value Decomposition (SVD). SVD decomposes a rectangular m -by- n matrix A into the product of three other matrices: $A = U\Sigma V^T$. The matrix Σ is a r -by- r diagonal matrix, in which the diagonal entries $(\sigma_1, \sigma_2, \dots, \sigma_r)$ are singular values and r is the rank of A . As explained in (Deerwester et al., 1990), SVD is closely related to standard eigenvalue-eigenvector decomposition of a square symmetric matrix. In fact, U is the matrix of eigenvectors of the square symmetric matrix AA^T , while V is the matrix of eigenvectors of $A^T A$. Σ^2 is the matrix of eigenvalues for both AA^T and $A^T A$. The interested reader can find more technical details on SVD in advanced linear algebra literature such as (Golub and Loan, 1996).

Since SVD can be applied to any rectangular matrix, it can also be used to decompose a term-document frequency matrix such as the one depicted in Fig. 1. After such a decomposition, depicted in Fig. 2, the matrices U and V contain vectors that specify the locations of the terms and documents in a

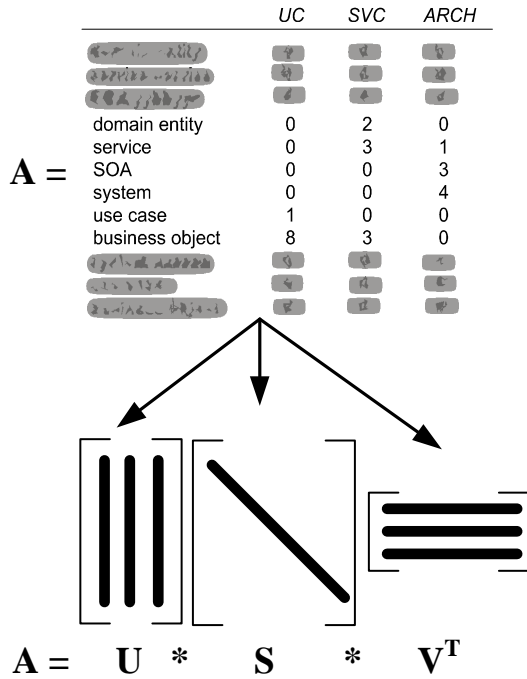


Fig. 2. Singular value decomposition of a term-document frequency matrix.

term-document space, respectively. The r orthogonal dimensions in this space can be interpreted as representations of r abstract concepts (cf. (Landauer et al., 1998)). The left-singular and right-singular vectors u_i and v_j indicate how much of each of these abstract concepts is present in term i and document j .

As outlined above, the original matrix A can be reconstructed by calculating the product of $U\Sigma V^T$. Instead of a reconstruction, a rank- k *approximation* of A can be calculated by setting all but the highest k singular values in Σ to 0. This approximation, A_k , is the closest rank- k approximation to A (Berry et al., 1994). Calculating A_k for a term-document space, such as the one depicted in Fig. 1, results in the closest k -dimensional approximation to the original term-document space (Letsche and Berry, 1997). In other words, by using SVD it is possible to reduce the number of dimensions in a term-document space. It is exactly this capability of SVD, depicted in Fig. 3, that is employed by LSA.

By using only k dimensions to reconstruct a term-document space, LSA no longer recalculates the exact number of occurrences of terms in documents. Instead, LSA estimates the number of occurrences based on the dimensions that have been retained. The result is that terms that originally did not appear in a document might now be estimated to appear, and that other words that did appear in a document might now have a lower estimated frequency (Landauer et al., 1998). This is the way in which LSA infers the latent semantic structure underlying the term-document space, and the way in which the deficiencies in the semantics captured in a vector-space model are overcome.

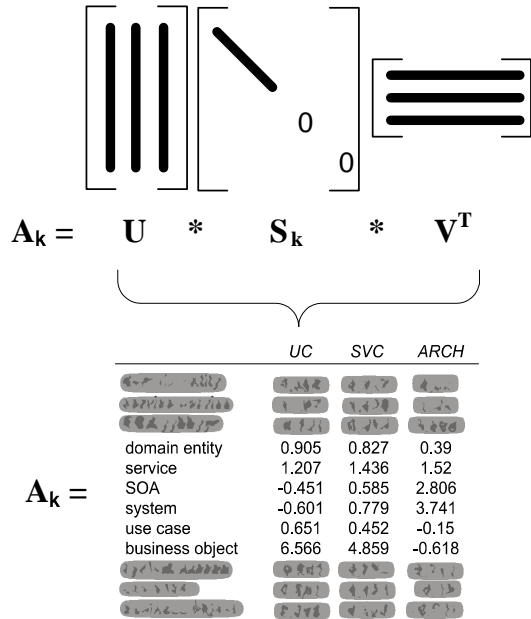


Fig. 3. Calculation of the closest rank- k approximation to the original term-document space.

In the reduced dimensional reconstruction of the term-document space, the meaning of individual words is inferred from the context in which they occur. This means that LSA largely avoids problems of synonymy, for instance introduced because two different authors of documentation for the same software product use two different terms to denote the same concept. One of the authors might for instance use the full product name in the documentation, while the other author prefers to use an acronym. Since the contexts in which these different terms are used will often be similar, LSA will expect the product acronym to occur with relatively high frequency in texts where the full product name is used and vice versa. However, it should probably be stressed here that we cannot expect LSA to improve the documentation other than making it more accessible. LSA will happily accept wrong, superfluous, or obsolesced documentation and guide anyone interested to ‘relevant’ parts of that documentation. Nonetheless, for reasonably well-written documentation the latent semantic structure LSA infers can be very well exploited to guide the reader.

Figure 4 shows the result of the application of LSA to the term-document frequency matrix from Fig. 1. The cutout shows the same six terms that are shown in the cutout in Fig. 1, but this time the numbers correspond to the *estimated* term frequencies based on retaining only 2 dimensions. Upon inspection of this result, interesting patterns appear. For starters, the term SOA is now expected to be present in the service specification as well, albeit at a lower frequency than in the architecture description. This corresponds to our intuitive notion that we would expect at least some relevance of SOA to a

	UC	SVC	ARCH
domain entity	0.905	0.827	0.39
service	1.207	1.436	1.52
SOA	-0.451	0.585	2.806
system	-0.601	0.779	3.741
use case	0.651	0.452	-0.15
business object	6.566	4.859	-0.618

domain entity	0.905	0.827	0.39
service	1.207	1.436	1.52
SOA	-0.451	0.585	2.806
system	-0.601	0.779	3.741
use case	0.651	0.452	-0.15
business object	6.566	4.859	-0.618

Fig. 4. Estimated term-document frequencies after the application of LSA to the matrix in Fig. 1.

service specification. The negative expected frequency of SOA in the use case specification is somewhat awkward to interpret mathematically, but might perhaps best be regarded as a kind of ‘surprise factor’. In a sense, LSA tells us not only that it does not expect the term SOA to crop up in the use case specification (estimated number of occurrences = 0), but that indeed it would be quite surprised to encounter this term there.

In general, a pattern seems to emerge in Fig. 4. If we regard the use case specification as the lowest level of abstraction text, the architecture description as the highest level, and the service definition somewhere in between, we see that low-level concepts (such as ‘business object’ and ‘use case’) have a diminishing level of association as the level of abstraction of the text increases and vice versa. LSA also seems to indicate that the term ‘service’ is a central concept in the documentation: its estimated frequency is almost equal for all three documents. Those patterns stem from the semantic structure in the documents. We can employ this uncovered semantic structure to guide an auditor to the information needed.

4 Constructing a Reading Guide: A Case Study

The LSA technique introduced in Section 3 forms the basis of a detailed case study in which we examine how the semantic structure discovered by LSA can be employed to guide the auditors through the documentation. This section presents the results of this case study.

Figure 5 depicts the interactive process by which an auditor is guided through the documentation. Initially, auditors start with a set of unread documents. Although the content of these documents is still unknown, the auditors have a

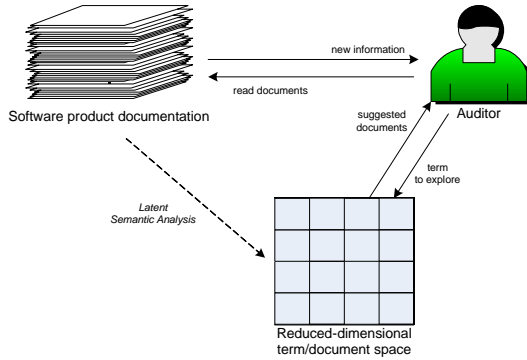


Fig. 5. Schematic overview of the construction of a reading guide using the reduced-dimensional term-document space calculated by LSA.

goal that needs to be satisfied by reading (part of) the documentation. Examples of such goals are obtaining a global understanding of the software product, investigating certain quality attributes, or locating (further) evidence for certain findings. The reduced-dimensional term-document space, which results from LSA, can be inspected to locate documents that are highly associated with a term that corresponds to the auditor’s goal. For instance - and this example will be worked out in more detail below - the term ‘architecture’ could be used to find documents that provide high-level information about the software product. From reading the suggested documents, an auditor learns new information including new - potentially product-specific - terms that can be used to locate documents that provide more detail on the new terms. In short, reading guidance consists of an iterative process of selecting and reading documentation, in which the auditor can use the architectural knowledge gained from reading suggested documents to steer the selection of new documents.

We applied LSA to a total of 80 documents that were subject to the audit that has been described earlier. The term-document frequency matrix that was constructed for these documents contained a total of 3290 terms found in the 80 documents. These 3290 terms did not contain very common words (‘stopwords’) such as ‘a’, ‘the’, or ‘is’. It is common practice to disregard these stopwords, since they tend to be evenly spread over all documents and hence do not bear any distinctive meaning. The length of the document vectors that make up the term-document frequency matrix had been normalized using cosine normalization (Salton and Buckley, 1988) before LSA was applied. This normalization reduces the effect of document size (i.e., the number of terms in the document); without normalization, longer documents tend to be favored over shorter documents when a document selection is made.

Using the technique described in Section 4, we calculated a 5-dimensional approximation of the constructed term-document frequency matrix. The selection of the number of dimensions to retain is an empirical issue (Landauer et al., 1998), although some heuristics exist (Berry et al., 1999). The rank-

5 approximation chosen here requires a 49% change relative to the original term-document frequency matrix. Although this might appear as a rather large change, the results obtained with this approximation suit our needs; they can be effectively used to construct a reading guide.

The case study presented here reconstructs the early phase of the software product audit, in which the auditors need to attain a global understanding of the software product in order to further assess its quality. As in the previous section, for reasons of non-disclosure the results presented here have been anonymized.

In general, when auditors commence a software product audit they want to gain an initial, high-level understanding of the software product. This global understanding is necessary to successfully perform the audit, since it is a prerequisite for subsequent audit activities. For instance, in scenario analyses the supplier of the software product is asked how the product reacts to certain change scenarios or failure scenarios. In order to judge the answer the supplier provides, an auditor needs to have a thorough understanding of the software product. Otherwise, the supplier might provide an answer that is incomplete or inconsistent with the real state of the product, without this being noticed.

Auditors who want to attain overall comprehension of the software product can be guided through the documentation using the semantic structure discovered by LSA. A route that is preferred by all auditors we interviewed is to start with high-level, global information and gradually descend to texts that contain more detailed and fine-grained information. A single term that can be expected to cover the high-level information about the software product well is the term ‘architecture’.

We can inspect the reduced 5-dimensional approximation of the original term-document frequency matrix that LSA has calculated to find the documents that best match the term ‘architecture’. In order to do so, it suffices to rank the documents by their respective values in the row that coincides with the term ‘architecture’ (the ‘architecture’ term vector). Documents that have a high value in the ‘architecture’ term vector correspond to the documents in which LSA expects the highest number of occurrences of the term ‘architecture’. Recall that the highest-ranking documents do not necessarily include the literal term ‘architecture’, but LSA inferred that it would be likely to encounter the term ‘architecture’ in these documents; they are semantically close to the meaning of ‘architecture’. In other words, these documents talk *about* architecture, perhaps without mentioning the word ‘architecture’ itself.

The list in Table 1 shows the 10 highest ranked documents for the term ‘architecture’, together with the actual number of occurrences of ‘architecture’ in these documents. Given this list, an auditor can simply start reading top-down,

Table 1

Top-10 documents that match the term ‘architecture’, with the number of occurrences of ‘architecture’ in the document.

<i>Rank</i>	<i>Document ID</i>	<i># ‘architecture’</i>
1	79	44
2	39	1
3	44	3
4	41	2
5	78	10
6	46	0
7	45	1
8	42	1
9	40	2
10	49	0

in this case starting with document 79. However, some of these documents are fairly large while others are rather small. In fact, documents 46 and 45 both consist of only 2 pages. If an understanding of the software product can be attained by either reading a (very) small document or by ploughing through a large number of pages, the former is obviously preferred by the auditors. Table 2 lists the same top-10 documents for ‘architecture’ also shown in Table 1. In this table, however, the documents have been categorized according to their size. The size categories have been defined as: very small (< 5 pages), small (< 10 pages), medium (< 30 pages) and large (≥ 30 pages). The rank according to Table 1 is given in parentheses, to illustrate the differences.

Table 2 shows that, given a preference for smaller documents, an auditor looking for information about the architecture of the software product should first read document 46. Note that this document does not contain the term ‘architecture’ at all (see Table 1). Nevertheless, upon inspection this document indeed contains high-level ‘architectural’ information.

From document 46, the auditor learns for instance that the software product consists of three high-level components, which we will call X, Y, and Z. Furthermore, the document identifies two external systems that interact with the software product as well as an organizational unit that will handle certain types of operational problems that might occur. Finally, the document contains a list of intended uses of the software product.

Now that the auditor knows a little more about the software product, the next document has to be selected. Since the auditor still has not read all

Table 2

Top-10 documents that match the term ‘architecture’, grouped by size.

	<i>Rank</i>	<i>Doc. ID</i>	<i># pages</i>
<i>Very small</i>	1 (6)	46	2
	2 (7)	45	2
<i>Small</i>	3 (5)	78	6
	4 (10)	49	9
<i>Medium</i>	5 (1)	79	24
	6 (2)	39	13
	7 (3)	44	21
<i>Large</i>	8 (4)	41	30
	9 (8)	42	48
	10 (9)	40	31

‘architecture’ documents, there are in principle two options: either remain with the ‘architecture’ documents and select a document from that list (e.g., document 45) or use the architectural knowledge obtained to delve into a particular topic.

Good candidates for further exploration of the documentation are the components X, Y, and Z. Since these components conceptually divide the software product in three distinct parts, auditors might want to examine each of these parts to further their global understanding. In its current form, the selection of the right terms for exploration is a matter of experience. It is from experience that the auditor knows that the term ‘architecture’ is likely to be related to high-level software product documentation. It is from experience that the auditor suspects that the three components are good candidates for further exploration.

Given the fact that the auditors want to gradually progress through the documentation, the degree of deviation of each of the components from the meaning of the term ‘architecture’ is an indication of the route that should be followed through the documents. In order to assess the deviation, a calculation can be performed of the similarity between the terms ‘architecture’ and ‘componentX’, ‘componentY’, and ‘componentZ’ respectively. This enables us to identify how much the texts for which LSA infers a high association with each of the components deviate from the text in document 46, which closely resembled the meaning of the term ‘architecture’.

A common measure of similarity between terms (or ‘term-term similarity’) is the cosine of the angle between the two term vectors (Berry et al., 1999). Let

t_i and t_j be the term vectors for terms ‘i’ and ‘j’ respectively, i.e., the rows from A_k that correspond to the terms ‘i’ and ‘j’. Then the cosine of the angle θ between these term vectors is $\cos \theta = \frac{t_i \cdot t_j}{\|t_i\| \|t_j\|}$.

Table 3 shows the similarity of each of the terms ‘componentX’, ‘componentY’, and ‘componentZ’ with the term ‘architecture’, calculated as the cosine between their respective term vectors. It becomes clear from these three values that ‘componentX’ is semantically closest to ‘architecture’ followed by ‘componentZ’, and that ‘componentY’ is the least similar to ‘architecture’.

An interesting observation is that the relations between the three components are not readily apparent from the text in the document itself, nor from the names given to the components. Although the document does contain a picture that seems to suggest a layered ordering of the components, the text in document 46 does not mention or reflect such a layered approach at all. Here, by using LSA we have truly discovered architectural knowledge that the auditor could not have distilled from reading document 46 alone. This discovered knowledge can be used to further explore the documentation.

Based on the similarity of ‘architecture’ and each of the three components, a logical next step to read the documentation seems to first read the top-ranking documents for ‘componentX’, then for ‘componentZ’, and finally for ‘componentY’. By following this route, the semantic distance between the document just read and the newly selected documents increases gradually.

Table 3

Cosine term-term similarity of ‘architecture’ and the high-level components

<i>componentX</i>	<i>componentY</i>	<i>componentZ</i>
0.9814	0.4096	0.7900

Analogous to the selection of the top-10 documents for the term ‘architecture’, we selected the top-10 documents for each of the terms ‘componentX’, ‘componentY’, and ‘componentZ’. For each of the components, we analyzed the top-10 documents found. The results of this analysis show an interesting pattern in the selection of documents.

Due to its close semantic resemblance of ‘architecture’, shown in Table 3, the top-10 documents that were found for the term ‘componentX’ are in fact the same 10 documents that were found for ‘architecture’. The only difference is a small change in the ranking of the documents. The top-10 documents found for the term ‘componentZ’ (the next closest term to ‘architecture’) comprises a mix of service specifications and architectural design descriptions, with a clear focus on service specifications (the first four documents in the list are service specifications). The top-10 documents found for ‘componentY’ are all use case definitions.

The route through the documentation found by analyzing the result of LSA suggests that, using ‘architecture’ as a starting point, the auditors should first read the architecture descriptions and similar high-level documentation, then proceed with service specifications, and finally read the use case definitions. Although LSA does not in and of itself know of the distinction between high-level documents (i.e., architecture descriptions) and low-level documents (i.e., use case definitions), the documents that it suggests to read are grouped along this axis. Moreover, from interviews with the auditors we learned that this is indeed a route they prefer to follow to familiarize themselves with a software product.

5 Validation of the use of LSA

The previous section shows how the application of LSA delivers results that support auditors in finding a route through the documentation. The auditors indicate that the results show correspondence to their preferences for selecting and reading documents. In this section we empirically validate this correspondence.

The knowledge discovered by using LSA can only be regarded valid if it fits the expectations of the auditor. In other words, the discovered semantic structure must conform to the auditor’s mental model of software product documentation and the architectural knowledge contained within. This means that the validity of the result is in principle highly subjective.

Although an auditor’s mental model is subjective, the auditors who were interviewed in the light of our case study on architectural knowledge use appear to have commonalities in their mental model of software product documentation. Perhaps the most obvious commonality is the categorization of documentation according to its level of abstraction. However, when confronted with the question how they find the architectural knowledge needed without having a reading guide available, most auditors reply that it is largely a matter of experience. This means that, even if we can employ the commonalities in the auditors’ perception of software product audits, we still need to validate our results against tacit knowledge.

A method that can be used to elicit tacit knowledge is the repertory grid technique. This technique stems from personal construct theory, and was devised by George Kelly (Kelly, 1955). The repertory grid technique can be used for exploring so-called ‘personal construct systems’; the collection of implicit personal theories. The repertory grid technique is “an attempt to stand in others’ shoes, to see their world as they see it, to understand their situation, their concerns” (Fransella and Bannister, 1977).

The repertory grid technique investigates bipolar constructs that form someone's mental model of (part of) the world. Examples of such constructs are *past-future*, *good-bad*, and *everything-nothing*. A particular method to elicit these constructs is by presenting a person with triads of elements from the domain under investigation. In our case, we could regard the individual documents as elements. When three of these elements are presented, the person is asked in which way two of them are alike and how the third one is different. This identifies a bipolar construct or axis that is apparently part of the person's mental model.

Based on elicited constructs, a so-called rank order grid can be constructed. To do so, the person is asked to rank all elements according to how well they fit the construct poles (e.g., from 'past' to 'future'). This rank order grid can then be further analyzed, for instance to determine the distances between different elements or constructs (Jankowicz and Thomas, 1982). In our case, we derive the distance matrix for the distances between documents as perceived by an auditor.

Ideally, we would apply the repertory grid technique to the same documents that were used in the case study described in Section 4. Unfortunately, using 80 elements in a repertory grid experiment is infeasible. Because of the exponential growth of possible combinations (triads) of elements, the upper limit for the application of the repertory grid technique is somewhere around 20 elements. For a validation of the LSA results, we therefore selected a second audit project in which a much smaller number of documents were involved.

We understand that it is dangerous to attribute the validity of one project to the validation of another, especially so if the projects differ in size. Nevertheless, we believe since auditors already acknowledge that the LSA results are sensible validation within a small project at the very least significantly adds to the credibility of our application of LSA in larger projects. This is even more so because LSA results are generally thought to improve for larger corpora (Landauer et al., 1998).

For the audit project in which we applied the repertory grid technique, the number of documents that had to be assessed was limited to 10. The audit team consisted of three members: one project manager and two auditors. The repertory grid technique was used in two experiments, one for each auditor. Each experiment lasted for approximately two hours, until the auditors felt they could not think of any more sensible (and distinctive) constructs. At the time we conducted the experiment, the audit project itself had been finished approximately two months earlier. To prevent distortions of the experiment results because of the two month gap between audit and experiment – part of the mental model might have been forgotten in the meantime – the triads presented to the auditors consisted of the actual (physical) documents, which the

auditors could freely browse and read during the experiment. Two exceptions were documents AS and DB (see below), which were no longer physically or electronically available. Those two documents were represented by a proxy: a single sheet of paper with the document's title.

The 10 documents will be referred to by the following abbreviations:

- FM contains a functional data model.
- FD describes the functional design.
- XX is an addendum to FD. This is not immediately clear from the title of the document, nor from the document's layout.
- PD contains the process design.
- TP contains a test plan.
- UM is the user manual.
- RN is a set of release notes.
- IM is the installation manual.
- AS is an administration manual for the application server.
- DB is an administration manual for the database server.

Tables 4 and 5 depict the rank order grids for auditors 1 and 2 respectively. The auditors were asked to rank the documents on a 1 to 5 scale. Hence, the first row in Table 4 should be read as follows: auditor 1 considers the contents of documents FM, FD, UM, and DB to be invariable over releases, whereas RN varies per release. PD, TP, and AS are considered to be more invariable than variable (but not completely invariable), while the opposite is true for IM. XX is exactly in the middle. The constructs should be interpreted from the auditor's point of view in the context of the performed audit. In other words, the construct 'used by me'/'not used by me' in Table 4 means '(not) used by auditor 1 in the audit'.

While many things could be said about the constructs that were elicited from both auditors – for instance regarding the commonalities and differences between the two grids – we will not perform such an analysis here. Within the grids themselves, clearly not all constructs (or dimensions) are orthogonal. This, too, does not matter for the discussion at hand. It suffices to take the two grids simply for what they are: a representation of the auditor's own mental model of the 10 documents used in the audit.

We analyzed the auditors' rank order grids to calculate the distances between the documents as perceived by the auditors. The resulting distance matrices can be further analyzed, for instance to determine clusters of documents that contain similar documents.

Figures 6 and 7 depict the document clusters according to auditor 1 and 2 respectively. The clusters have been determined with the single linkage hierarchical clustering method (Johnson, 1967). The axis denotes the difference

Table 4
Rank Order Grid Auditor 1

1	<i>XX</i>	<i>PD</i>	<i>FM</i>	<i>FD</i>	<i>TP</i>	<i>UM</i>	<i>RN</i>	<i>IM</i>	<i>AS</i>	<i>DB</i>	5
invariable	3	2	1	1	2	1	5	4	2	1	varies per release
used by me	5	1	1	1	2	4	5	5	3	3	not used by me
input for development	1	1	1	1	3	5	5	5	4	4	output of development
prescriptive	1	1	1	1	3	5	5	4	5	5	descriptive
development	1	1	1	1	2	2	3	4	5	5	deployment
functionality	2	2	1	1	2	2	4	5	5	5	no functionality
diagrams expected	3	2	1	1	5	2	5	5	5	5	no diagrams expected
whole application	5	1	1	1	1	1	1	5	5	5	part of application
use	4	2	2	2	3	1	3	5	5	5	deployment
test	1	2	3	2	1	1	1	5	5	5	deployment
data flow	1	1	5	3	1	2	3	3	3	4	data entity
too global	5	5	5	4	1	2	3	5	3	3	too detailed
good size	5	5	4	1	1	1	3	1	3	3	overwhelming

Table 5
Rank Order Grid Auditor 2

1	<i>XX</i>	<i>PD</i>	<i>FD</i>	<i>FM</i>	<i>TP</i>	<i>RN</i>	<i>UM</i>	<i>IM</i>	<i>DB</i>	<i>AS</i>	5
abstract	2	1	1	2	4	5	5	5	5	5	concrete
content	1	1	1	1	3	5	5	5	4	4	packing
input for development	1	1	1	1	3	4	5	5	4	4	output of development
descriptive / static	1	1	1	1	3	3	5	5	4	4	use / time dimension
app. functionality	1	1	1	1	3	3	3	4	5	5	system administration
design	1	1	1	1	3	4	3	5	5	5	deployment
conceptual	1	1	1	2	2	4	3	4	5	5	concrete/instance
high level	5	3	1	1	2	5	3	4	4	4	detailed
absolute	1	1	1	1	3	5	2	2	2	2	relative (wrt prev. version)
application	1	1	1	1	5	2	2	2	2	2	organisation
used by me	4	2	1	2	5	2	1	1	3	3	not used by me

between the documents, calculated as 1 minus the similarity. For instance, for auditor 1 the similarity between documents FD and FM has been calculated as 0.87, therefore the difference between the two equals 0.13, as shown in Fig. 6.

Although there are some differences between the two cluster configurations, both auditors seem to agree that there are two large document clusters. One cluster contains documents FD, FM, PD, and XX. The other documents are grouped in the second cluster. Note that we left AS and DB out of the figures to allow for a fair comparison with the effect of LSA. Recall that those two documents were no longer available, due to which LSA was unable to process them. Had we included them in the cluster figures, they would have shown as a small sub-cluster of two very similar documents (similarity according to auditor 1: 0.96; auditor 2: 1.00). For both auditors this sub-cluster is most similar to document IM (auditor 1: 0.79, auditor 2: 0.84).

To illustrate the effect of LSA on the document vector-space model, we applied LSA to the 8 documents from the audit that were still available. We determined the distance matrices for both the term-document frequency matrix and the LSA result. The distance measure used to calculate distances between two documents is the cosine of the angle between the two corresponding document vectors.

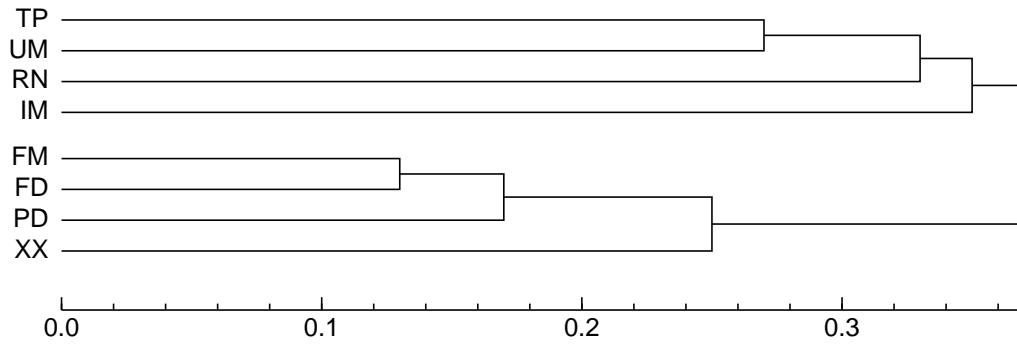


Fig. 6. Auditor 1: Hierarchical documentation clusters.

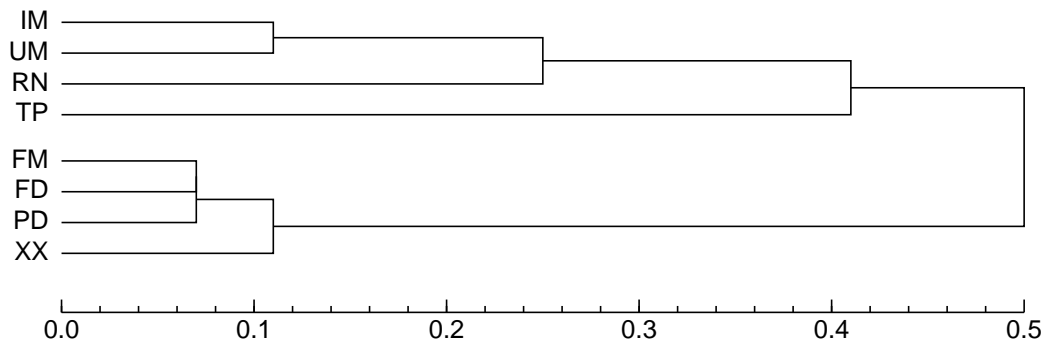


Fig. 7. Auditor 2: Hierarchical documentation clusters.

Figure 8 depicts the clusters for the original term-document frequency matrix. Those clusters already show some correspondence to the two clusters which the auditors perceive. However, documents XX and RN are two clear outliers. Moreover, the distinction between the two clusters (the cut around 0.5 similarity) is not very obvious.

The distinction between the two clusters is much clearer in Figure 9. This figure shows the clusters after LSA has been applied. The only outlier remaining is the set of release notes (RN). A possible (but unverified) explanation might be that the release notes follow a style different from the other documents. The release notes consist mainly of a list of short, somewhat stenographic messages that describe the changes from one version to another while the other documents contain more extensive text.

The hierarchical clusters leave the impression that the application of LSA transforms the vector space model (initially represented by the term-document frequency matrix) to better resemble the auditors' mental models. A quantitative expression of this 'better resemblance' can be given by calculating the correlation between the various distance matrices.

To calculate the correlation between the distance matrices we have performed a simple Mantel test using the *zt* software tool (Bonnet and Peer, 2002). In this test, the null hypothesis is that the distances in one distance matrix

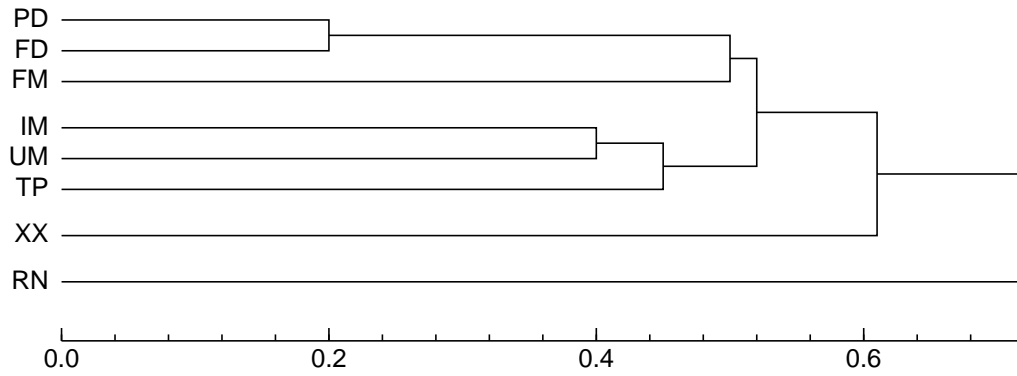


Fig. 8. Term-frequency: Hierarchical documentation clusters.

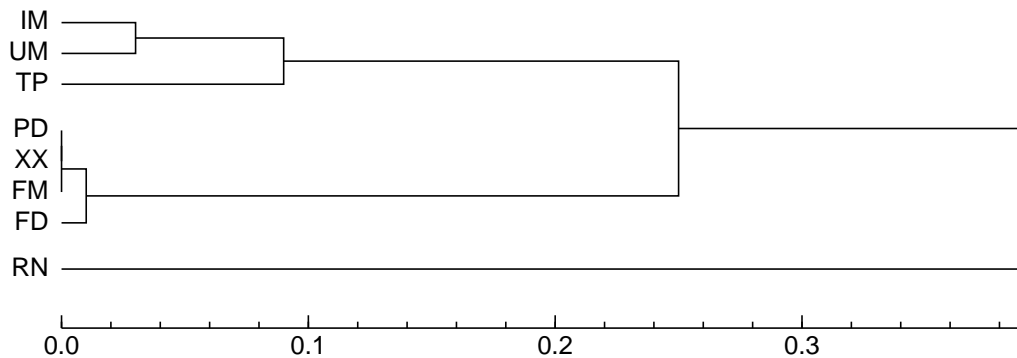


Fig. 9. LSA: Hierarchical documentation clusters.

are independent from those in another. The results of this test are shown in Table 6.

Table 6

Simple Mantel Test r and p -values

	<i>Auditor 1</i>	<i>Auditor 2</i>
Term-document frequency matrix	$r = 0.627, p = 0.00184$	$r = 0.666, p = 0.00017$
LSA	$r = 0.730, p = 0.00179$	$r = 0.871, p = 0.00057$

The results of the simple Mantel test show that there is already a significant correlation between the distances according to the term-document frequency matrix and the auditors' mental models. This corresponds with how the two clusters already appear in Figure 8. However, the correlation between the auditors' rank order grids and the result of LSA are clearly higher. As a matter of fact, the lowest correlation with LSA ($r = 0.730$ for auditor 1) is comparable to the correlation between the rank order grids of the auditors themselves. A simple Mantel test shows that the correlation between the two auditors' rank order grids is 0.738 ($p = 0.00020$).

In summary, the correlation coefficients calculated by the simple Mantel test indicate that application of LSA yields a vector space model more similar to an auditor's mental model than merely counting the words in the documents.

Visually, this shows as a better correspondence of document clusters with the clusters according to the auditors' perception. This empirical validation supports our initial finding that the reading guide constructed by applying LSA intuitively made sense for the auditors.

One thing that should again be stressed here, however, is that the result of the application of LSA is very much dependent on the selected number of dimensions to reduce to. Unfortunately we currently have no better guideline than trial-and-error heuristics. Since there was only a small number of documents used in this particular case, we were able to calculate the results for all possible number of reduced dimensions – the original term-document frequency matrix has a maximum number of orthogonal dimensions equal to the number of documents. This allowed us to select the number of dimensions – 4, for the record – that looked most promising. When larger numbers of documents are involved, this approach becomes infeasible quite soon.

6 Related Work

The application of Latent Semantic Analysis to architectural knowledge discovery discussed in this article bears some relation to other work, both within and outside of the software engineering research domain. The origin of LSA lies in information retrieval. LSA was presented in 1990 by Deerwester et al. as 'a new method for automatic indexing and retrieval' of documents (Deerwester et al., 1990). Later research also focused on the psycholinguistic significance of LSA. Landauer and Dumais, for instance, use LSA to simulate the acquisition of vocabulary from text, and present LSA as a theory of acquisition, induction, and representation of knowledge (Landauer and Dumais, 1997).

Over the years, LSA has seen various application domains, including software engineering. For instance, Maletic et al. applied LSA to source code of software components in order to support program comprehension (Maletic and Marcus, 2000; Maletic and Valluri, 1999). Another approach is taken by Hayes et al., who use LSA to support the construction of requirements traceability matrices (Hayes et al., 2005).

Our approach adds a new item to the list of LSA applications in software engineering. Although construction of a reading guide for software product documentation also contributes to better program understanding, our approach differs from the use of LSA by Maletic et al., who apply LSA to source code artifacts. The effect of a reading guide is also not limited to better product comprehension, but further supports the auditors in locating evidence for their findings. Since architectural knowledge can be reflected differently in source code and documentation, some of the evidence and knowledge that can be

found in the documentation might not be available from the software product's source code.

7 Future Work

The work presented in this article gives rise to a number of issues that warrant further research. An overall issue that remains to be investigated is the scalability of our approach. LSA proved to be feasible for a corpus of 80 documents, but in practice software product documentation might comprise many more documents. Document sets of several hundreds of documents are not uncommon.

Furthermore, the selection of the right number of reduced dimensions is still difficult. In this area, a comparison of the auditor's mental model with the result of LSA – as shown in Section 5 – could possibly provide guidance to selection of the right number of reduced dimensions, since the optimal number of dimensions yields the best match of the LSA result with the auditor's perception. This of course implies that the auditor's mental model is known. Although determining this model is feasible by means of the repertory grid technique, this is by no means a trivial task to perform. Especially not if this has to be done for every new project in which LSA is to be used.

Besides these global issues, we have identified three main areas that are to be further explored in the present context: enhancement of the workflow, the use of background knowledge, and user interaction. This section describes each of these areas in more detail.

7.1 *Workflow Enhancement*

The 'workflow' presented in this article, i.e., the selection of terms to explore the documentation, is still rather ad hoc and depends heavily on the auditor's experience. One could wonder whether the same result would have been obtained had the process been started with another 'high-level' term, such as the name of the software product instead of the generic term 'architecture'.

As a matter of fact, using the name of the product, or the high-level term 'system' instead of 'architecture', would have yielded a different result. The documents that are suggested for these terms resemble the documents that were suggested for 'componentZ', i.e., with an emphasis on service specifications. Document 46 would not have been suggested for any of these terms. Depending on one's opinion this may or may not come as a surprise. Some

might argue that ‘system’ indeed carries more of a notion of implementation than ‘architecture’. It does show, however, the importance of the selection of the (initial) terms to explore. It also shows that the auditor would benefit from assistance in this selection instead of having to rely on experience alone. We would like to capture this kind of experience to enhance the workflow and aid the auditor in selecting new terms to explore.

We believe that we can capture relevant experience and enhance the workflow by introducing a feedback loop. By keeping track of terms that worked well in earlier projects, the auditor can be presented with suggestions as to which terms to explore in a new project. This helps the auditor to get the project started (e.g., by suggesting initial terms such as ‘architecture’), but can also circumvent potential dead-ends in the exploration by explicitly ignoring terms that are known to have led to dead-ends in previous audits. Such suggestions could perhaps also be mined from the documentation itself, using techniques such as frequency profiling to locate uncommon words (with respect to a standard corpus) that are likely to be part of a domain-specific vocabulary. Sawyer et al. report successful application of frequency profiling in extraction of domain terms from requirements engineering documents (Sawyer et al., 2005).

While a feedback loop could relatively easily handle common generic terms such as ‘architecture’, additional research is needed in particular to determine how to cope with product-specific terminology such as ‘componentX’. The reasons that auditors choose certain (product-specific) terms for further exploration need to be analyzed and translated to more generic ‘heuristics’ and best practices that are applicable to other projects as well. An example heuristic might be that, given the auditor’s goal of overall product comprehension, terms that signify components are better candidates for further exploration than terms that signify intended uses. This heuristic can change when the auditor’s goal changes. If the auditor is looking for the satisfaction of certain requirements, intended uses might be preferred over components.

7.2 Background Knowledge Incorporation

The ‘queries’ that are used in this article to explore the software product documentation are logical from an auditor’s point of view. The auditor starts with a high-level exploration of the software product’s architecture, gradually zooming in to reveal more detailed architectural knowledge. Through Latent Semantic Analysis, documents with diminishing levels of abstraction are identified: from architecture descriptions at the highest level through service definitions to use case specifications at the lowest level. However, this analysis sequence still requires extensive human interpretation. As remarked earlier,

LSA itself has no notion of ‘high-level’ or ‘low-level’ documentation, nor of any other domain-specific knowledge.

To further enhance the support for auditors reading the software product documentation we intend to investigate the incorporation of relatively static background knowledge in the automated analysis of the documentation. The words ‘relatively static’ signify domain knowledge that does not change for each audit. Apart from often used classifications such as high-level vs. low-level documentation, examples of such background knowledge are:

- generic models, such as quality models (e.g., (ISO/IEC, 2001)) and process models (e.g., (ISO/IEC, 1998));
- ontologies, for instance of architectural patterns (e.g., the Handbook of Software Architecture by Booch (<http://www.booch.com/architecture/>)) and their known relations with for example quality attributes, generic software engineering ontologies (e.g., the ontology by the SEONTOLOGY project team (<http://www.seontology.org/>)) and application-generic software architecture ontologies (e.g., (Babu T. et al., 2007));
- ‘heuristics’, such as an auditor’s general preference for smaller documents (see also Section 4).

Background knowledge can be incorporated when constructing a reading guide by using it in the selection of suggested documents to read. Models and ontologies can for instance be used to broaden the scope when exploring a certain term; they can be used to formulate rules of the form ‘if auditors are interested in X (e.g., ‘maintainability’) they are (probably) also interested in Y (e.g., ‘changeability’, see also (ISO/IEC, 2001))’. Heuristics can for example be applied to change the suggested order in which the documents should be read, as demonstrated in Section 4.

Note that there is also some overlap of background knowledge incorporation with the workflow enhancements described in Section 7.1. The heuristics (or best practices) described in that section can in fact be regarded as background knowledge. The translation of product-specific terminology to generic terms used in these heuristics could very well be based on an ontology structure.

Background knowledge can hence be employed at two levels: to suggest terms to explore, steering the workflow; and to suggest documents to read, steering the analysis. Both affect the outcome of the process, the reading guide.

Using the correlation coefficient calculated by the simple Mantel test, we can closely monitor whether adjustment of our method improves the result. We could, for instance, assess the effect of the incorporation of background knowledge in the analysis of the documentation. Such an assessment consists of a comparison of the distances perceived by the auditors with the distances before and after incorporating background knowledge. If the result of the analysis

corresponds more to the auditor's mental model when background knowledge is taken into account, this means that using this background knowledge is indeed an improvement. Finally, the effect of using other techniques instead of, or together with, LSA could be easily judged analogous to the assessment of the incorporation of background knowledge. Techniques that could further improve architectural knowledge discovery results include techniques that, complementary to LSA, exploit certain more structured properties of the documentation. If, for example, a set of documents is structured according to a particular template – which is not uncommon – knowledge of this template could be used to guide the reader to particular parts of the documentation.

7.3 User Interaction and Tool Support

A final area in which further research is needed is the area of user interaction. The results presented in this article all show direct operations on the reduced-rank approximation of the original term-document frequency matrix. This matrix is probably not the best form of presentation for the end users, i.e., the auditors.

In order to be useful and used in an auditor's everyday practice, the techniques discussed in this article should be implemented in an interactive environment that abstracts away from the underlying estimated term frequencies. This environment should provide intuitive support for the workflow discussed in Section 7.1.

A particular area that requires further research is visualization of the reduced-dimensional term-document space. A desirable visualization supports both locating terms to explore and locating documents to read. Ideally, this would be presented to the auditors as a space through which they could navigate in search of the architectural knowledge they need. In this space, distances between terms and/or documents have actual meaning (cf. the three terms in Table 3). Such a visualization requires a projection of the reduced-dimensional term-document space to two - or at most three - dimensions. Through such a visualization, auditors can obtain quick visual clues as to which documents are closely related and how to proceed reading the document set.

8 Architectural Knowledge Discovery in a Broader Scope

This article considers Architectural Knowledge Discovery (AKD) as a means to construct a reading guide for software product audits. Although this application is undoubtedly valuable, we believe AKD has merit in a broader

scope.

We envision AKD as one particular technique used in a broad range of architectural knowledge management tools and methods. The role of AKD would mainly be to refine existing (codified) architectural knowledge from such diverse sources as documents, email, meeting minutes, and source code. Those sources contain mainly unstructured (documents, etc.) or at best semi-structured (source code) information. Refinement by AKD would therefore mainly consist of adding structure to this information.

AKD could play a major role in the transition from personalization to codification (cf. (Ali Babar et al., 2007)). In the early phases of architecting, often little attention is paid to a structured description (codification) of the architectural knowledge that is created, such as the architectural design decisions taken and the alternatives considered. More emphasis is usually put on personalization, i.e., knowing who knows what and discussing options and issues directly with peers. It is usually only after the fact – if ever – that the architectural knowledge created during this process is captured in architectural descriptions and other documents. However, this early architecting phase often leaves many traces in for example minutes, emails, or forum discussions. AKD could be employed to mine and structure those traces of unstructured information, thereby providing semi-automated support for effective codification of the architectural knowledge.

As a final remark, our current use of AKD with LSA concerns a final set of documents. We may also envision using LSA in a forward-engineering sense, to judge the quality of the evolving (architectural or otherwise) documentation of a system, and giving guidelines as to where the documentation needs attention.

9 Conclusion

Document inspection is a method used in software product audits to distill architectural knowledge from the software product documentation. Unfortunately, document inspection is often hard to perform. Auditors are in need of a reading guide that tells them where to start reading, how to progress reading, and which documents to consult for more detail on a particular topic.

We have demonstrated how auditors can be guided through the documentation in a case study in which we reconstructed the early phase of a software product audit. In this phase, the auditor has not read any documents yet and needs to attain a certain level of understanding of the software product.

To construct a reading guide, we have employed the semantic structure dis-

covered by Latent Semantic Analysis. This semantic structure is used as the basis for an interactive process in which auditors indicate terms that they want to explore and are subsequently given reading suggestions for documents containing information about these terms. The knowledge obtained from the suggested documents can give rise to new terms to explore, and the discovered semantic structure can be used to determine the order in which the terms - and corresponding documents - should be explored.

Using the repertory grid technique, we were able to elicit the mental model of documentation from two auditors. A comparison of the auditors' mental models with the result of the application of LSA shows that LSA's vector-space model is closer to the auditors' idea than a naive term-document frequency matrix.

We have identified three areas of future work: workflow enhancement, use of background knowledge, and user interaction. We intend to direct research efforts toward each of these areas in order to further improve the work presented in this article. Apart from the construction of a reading guide, we envision the application of architectural knowledge in a broader scope. Architectural knowledge discovery could be employed to refine unstructured traces of the early architecting phase, thereby aiding effective codification of architectural knowledge.

Acknowledgement

This research has been partially sponsored by the Dutch Joint Academic and Commercial Quality Research & Development (Jacquard) program on Software Engineering Research via contract 638.001.406 GRIFFIN: a GRId For inFormatIoN about architectural knowledge. The authors would like to thank Eefje Cuppen for helpful discussions on the repertory grid technique

References

- Ali Babar, M., de Boer, R. C., Dingsøyr, T., Farenhorst, R., 2007. Architectural Knowledge Management Strategies: Approaches in Research and Industry. In: Second Workshop on SHaring and Reusing architectural Knowledge / Architecture, Rationale, and Design Intent (SHARK/ADI). Minneapolis, MN, USA.
- Babu T., L., Seetha Ramaiah, M., Prabhakar, T., Rambabu, D., 2007. ArchVoc—Towards an Ontology for Software Architecture . In: Second Workshop on SHaring and Reusing architectural Knowledge / Architecture, Rationale, and Design Intent (SHARK/ADI). Minneapolis, MN, USA.

- Berry, M. W., Drmac, Z., Jessup, E. R., 1999. Matrices, Vector Spaces, and Information Retrieval. *SIAM Review* 41 (2), 335–362.
- Berry, M. W., Dumais, S. T., O’Brien, G. W., 1994. Using Linear Algebra for Intelligent Information Retrieval. Tech. Rep. CS-94-270.
- Bonnet, E., Peer, Y. V. d., 2002. *zt*: A Software Tool for Simple and Partial Mantel Tests. *Journal of Statistical Software* 7 (10).
- Booch, G., <http://www.booch.com/architecture/>. Handbook of Software Architecture.
- Bosch, J., 2004. Software Architecture: The Next Step. In: Oquendo, F., Warboys, B., Morrison, R. (Eds.), *Software Architecture: First European Workshop (EWSA)*. Vol. 3047 of Lecture Notes in Computer Science. Springer-Verlag GmbH, St. Andrews, UK, pp. 194–199.
- de Boer, R. C., 2006. Architectural Knowledge Discovery: Why and How? In: *First Workshop on SHaring and Reusing architectural Knowledge (SHARK)*. ACM Press, Torino, Italy.
- de Boer, R. C., van Vliet, H., 2007. Constructing a Reading Guide for Software Product Audits. In: *Sixth Working IEEE/IFIP Conference on Software Architecture (WICSA)*. IEEE Computer Society, Mumbai, India.
- Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., Harshman, R., 1990. Indexing by Latent Semantic Analysis. *Journal of the American Society for Information Science (JASIS)* 41 (6), 391–407.
- Fransella, F., Bannister, D., 1977. *A Manual for Repertory Grid Technique*. Academic Press.
- Golub, G. H., Loan, C. F. V., 1996. *Matrix Computations*, 3rd Edition. The John Hopkins University Press.
- Hayes, J. H., Dekhtyar, A., Sundaram, S. K., 2005. Improving After-the-Fact Tracing and Mapping: Supporting Software Quality Predictions. *IEEE Software* 22 (6), 30–37.
- ISO/IEC, 1998. Information technology - Software product evaluation - Part 5: Process for evaluators. Tech. Rep. ISO/IEC 14598-5.
- ISO/IEC, 1999. Information technology - Software product evaluation - Part 1: General overview. Tech. Rep. ISO/IEC 14598-1.
- ISO/IEC, 2001. Software engineering - Product quality - Part 1: Quality model. Tech. Rep. ISO/IEC 9126-1.
- Jankowicz, D., Thomas, L., 1982. An Algorithm for the Cluster Analysis of Repertory Grids in Human Resource Development. *Personnel Review* 11 (4), 15–22.
- Jansen, A., Bosch, J., 2005. Software Architecture as a Set of Architectural Design Decisions. In: *5th Working IEEE/IFIP Conference on Software Architecture (WICSA)*. IEEE Computer Society, Pittsburgh, Pennsylvania, USA, pp. 109–118.
- Johnson, S. C., 1967. Hierarchical Clustering Schemes. *Psychometrika* 32 (3), 241–254.
- Kelly, G. A., 1955. *The Psychology of Personal Constructs*. Norton, New York.
- Kruchten, P., Lago, P., van Vliet, H., 2006. Building up and Reasoning

- about Architectural Knowledge. In: Hofmeister, C., Crnkovic, I., Reussner, R. (Eds.), 2nd International Conference on Quality of Software Architectures (QoSA). Vol. 4214 of Lecture Notes in Computer Science. Springer, Västerås, Sweden, pp. 43–58.
- Landauer, T. K., Dumais, S. T., 1997. A Solution to Plato’s Problem: The Latent Semantic Analysis Theory of Acquisition, Induction, and Representation of Knowledge. *Psychological Review* 104 (2), 211–240.
- Landauer, T. K., Foltz, P. W., Laham, D., 1998. An Introduction to Latent Semantic Analysis. *Discourse Processes* 25, 259–284.
- Letsche, T. A., Berry, M. W., 1997. Large-Scale Information Retrieval with Latent Semantic Indexing. *Information Sciences* 100 (1-4), 105–137.
- Maletic, J. I., Marcus, A., 2000. Using Latent Semantic Analysis to Identify Similarities in Source Code to Support Program Understanding . In: 12th IEEE International Conference on Tools with Artificial Intelligence (ICTAI). IEEE Computer Society, Vancouver, BC, Canada, pp. 46–53.
- Maletic, J. I., Valluri, N., 1999. Automatic Software Clustering via Latent Semantic Analysis. In: 14th IEEE international conference on Automated Software Engineering (ASE). IEEE Computer Society, Cocoa Beach, FL, USA, pp. 251–254.
- Salton, G., Buckley, C., 1988. Term-Weighting Approaches in Automatic Text Retrieval. *Information Processing & Management* 24 (5), 513–523.
- Sawyer, P., Rayson, P., Cosh, K., 2005. Shallow knowledge as an aid to deep understanding in early phase requirements engineering. *IEEE Transactions on Software Engineering* 31 (11), 969–981.
- SEONTOLOGY project team, <http://www.seontology.org/>. The Software Engineering Ontology.
- van der Ven, J. S., Jansen, A. G. J., Nijhuis, J. A. G., Bosch, J., 2006. Design Decisions: The Bridge between Rationale and Architecture. In: Dutoit, A. H., McCall, R., Mistrik, I., Paech, B. (Eds.), *Rationale Management in Software Engineering*. Springer-Verlag, pp. 329–346.