

# A User-Centered Methodology for the Evaluation of (Semantic) Web Service Discovery and Selection

Friederike Klan  
Institute of Computer Science  
Friedrich-Schiller-University of Jena  
friederike.klan@uni-jena.de

Birgitta König-Ries  
Institute of Computer Science  
Friedrich-Schiller-University of Jena  
birgitta.koenig-ries@uni-jena.de

## ABSTRACT

Recently, a new breed of user-centric solutions to Web Service discovery and selection that applies Semantic Web Service technology in B2C settings such as e-Commerce has evolved. They significantly differ from traditional Web Service frameworks and have to cope with new challenges such as assisting consumers in specifying service requirements and providing effective decision support for service selection. Existing evaluation efforts within the scope of Semantic Web Services do not account for these user-specific requirements and hence are not appropriate for assessing the quality of solutions that are dedicated to end-users, i.e. service consumers. In this paper, we address this issue by proposing a user-centered methodology for the evaluation of Semantic Web Service retrieval and demonstrating its feasibility.

## Categories and Subject Descriptors

H.3.3 [Information Storage And Retrieval]: Information Search and Retrieval; H.3.4 [Information Storage And Retrieval]: Systems and Software

## General Terms

Algorithms, Human Factors, Measurement

## Keywords

Semantic Web Service selection, user-centered evaluation

## 1. INTRODUCTION

Service Oriented Architectures have become increasingly popular in enterprises as a means to integrate heterogeneous business applications into agile business processes based on loosely coupled Web Services. The focus of research efforts in this area was primarily on automatizing the process of service discovery, composition, binding, mediation and invocation. Since business process design and implementation is typically done by experienced application developers, tool

support for these kinds of tasks was primarily targeted to this audience.

Taking up the changing role of end-users, i.e. service consumers, as active contributors and emancipated users of the Internet, business services now break through the company boundaries. Leading IT-companies such as Amazon.com or Google.com open up their internal functionality and make it accessible to be used by end-users. Researchers and practitioners in the field forecast a shift from the current Internet of Information to an Internet of Services [12]. Semantic technologies are expected to play a central role in the realization of this vision, since they enable effective service retrieval, sophisticated, knowledge-based end-user support and convenient usage by semi-automatic configuration, composition and on-demand invocation of services [5].

These developments gave rise to a new breed of user-centric solutions, that apply Semantic Web Service (SWS) technology in B2C settings such as in e-Commerce [4, 9]. They significantly differ from traditional Web Service frameworks and have to cope with new research challenges such as assisting consumers in specifying service requirements and providing effective decision support for selecting appropriate services from a large set of offers with a variety of possibly conflicting characteristics. Existing evaluation efforts within the scope of SWSs mainly focus on retrieval performance and correctness. Though these are still valid evaluation targets, they are not sufficient to evaluate service retrieval solutions with respect to the specific requirements of end-users. Classical techniques for usability testing as employed by some researchers can elucidate whether a consumer used a retrieval system in the intended way, but cannot assess whether he actually made a good selection decision in terms of an objective measure. In this paper, we address this issue by

- identifying the specific requirements to service discovery and selection, if targeted to end-users, and
- by characterizing the general structure and operating mode of solutions that meet those requirements.
- Based on these considerations, we propose a user-centered evaluation methodology for SWS retrieval systems that is capable of assessing a system's quality with respect to the identified requirements by means of objective measures.
- To demonstrate the feasibility and appropriateness of the suggested approach, we provide details on its implementation and on the evaluation results that emerged when applying it to our own conversational approach to SWS retrieval [6].

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

WIMS'14 June 2-4, 2014 Thessaloniki, Greece

Copyright is held by the owner/author(s). Publication rights licensed to ACM.

ACM 978-1-4503-2538-7/14/06 ...\$15.00.

<http://dx.doi.org/10.1145/2611040.2611069>

## 2. REQUIREMENTS

Making "good" decisions when selecting a service is important, since a bad choice often results in some kind of (not necessarily monetary) loss or inconvenience for the decision maker. Think for instance about a consumer having booked a flight that has an intermediate hop, while a non-stop flight having roughly the same price would have been available. However, what qualifies a decision as "good"? According to Payne et al. [10], a "good" decision, i.e., in our scenario a "good" service selection decision, is characterized by the fact, that it is *well-informed*, i.e. taken in consideration of relevant service alternatives and their properties, *balanced*, i.e. taken after deliberate resolution of conflicting service requirements, and *consistent*, i.e. made in consciousness of a consumer's service needs and optimal or close to optimal with respect to these needs and the available service alternatives.

In the traditional B2B setting, where service requirements were known in advance and decisions were taken by a program, making well-informed, balanced and consistent choices was more or less a matter of effective service retrieval and of choosing proper decision rules. This is no longer true in selection scenarios, where humans, i.e. service consumers, are involved. Known weaknesses of human decision making impose additional requirements on solutions to service selection, particularly related to the provision of support for making "good" service selection decisions. In the remainder of this section, we will discuss those requirements addressing the most important limitations of human decision making as identified by research in behavioral psychology.

Existing approaches to SWS selection typically assume that service consumers have a clear goal in mind when looking for service functionality. However, as research results from behavioral decision theory indicate [10], this is often not true. Particularly, in service domains that are complex and unfamiliar, such as in an e-Commerce setting, consumers have no clear-cut requirements and preferences. People rather construct them instantaneously when facing choices to be made. As a consequence, a system that aims at providing decision support for service selection cannot merely elicit existing requirements and preferences from the user. It rather has to interactively acquire and incrementally refine those information as they are constructed by the user based on the characteristics of matching service alternatives (REQUIREMENT R1).

As argued, thoughtful selection decisions are characterized by the fact that they are made in awareness and consideration of the service requirements and preferences that are important to the decision maker. They should result from a process, where conflicts between requirements are resolved by making explicit tradeoffs instead of applying non-compensatory, i.e. non-rational, heuristics as often performed by decision makers. Hence, as argued in [10] and empirically shown in [3], the effectiveness of service selection is largely determined by the supporting system's ability to provide incentives for thinking about preferences and requirements and by its ability to encourage decision makers to make tradeoffs (REQUIREMENT R2).

Payne et al. [10] also emphasize that, in order to enable effective decisions making, it is crucial to educate consumers about relevant service alternatives and their characteristics and to motivate them to consider this knowledge when making a selection (REQUIREMENT R3). This is required, since

consumers tend to base their decisions on a narrow range of options (myopic decision frame) without considering relevant and potentially more desirable service alternatives.

Finally, it might happen that although a service consumer has well-constructed preferences and requirements in mind, he fails in picking a service that is optimal with respect to these and the available service alternatives. According to [10], the major causes for such inconsistent decisions are cognitive biases in scale usage such as anchoring effects (tendency of humans to overly rely on a certain, "anchored" aspect). A system that assists users in making thoughtful selection decisions should prevent inconsistent selection decisions (REQUIREMENT R4).

## 3. RELATED WORK

Besides evaluation efforts within the scope of single research projects, there have been several community initiatives such as the Semantic Web Service Challenge<sup>1</sup> or the S3 Contest<sup>2</sup>, which are dedicated to the comparative evaluation of SWS technologies. The EU-funded project SEALS<sup>3</sup> that aims at providing an infrastructure for benchmarking semantic tools also involves a campaign on Semantic Web Service tools.

Typical evaluation targets concern the expressiveness of the semantic description language underlying an approach as well as the functional coverage of the semantic match-maker, the mediation and the service composition capabilities of a solution. These are assessed by organizing challenges that provide sets of scenario-based discovery, mediation and/or choreography problems that have to be solved by the participants. The solutions are verified and certified by the challenge staff. Available test data collections for benchmarking purposes such as OWLS-TC<sup>4</sup> mainly target at assessing retrieval correctness and completeness, e.g. by means of Precision and Recall or graded relevance [7], and retrieval efficiency in terms of response time and resource consumption.

Though retrieval quality is a valid evaluation target, it does not assess whether an approach is actually helpful for its users. Therefore, additional measures have been taken to evaluate the effectiveness and efficiency of provided tool support, e.g. for creating service requests or for discovering and composing services. Several evaluation efforts (e.g. [1]) do this by employing the Thinking Aloud Protocol [8], an approved technique for usability evaluation. To this end, test users are asked to perform a given set of tasks using the tool to be evaluated and comment on their thoughts while trying to perform this task. By observing this process, the testers can gain insights on how users perceived the interface of the provided tool, whether they were able to accomplish the given tasks and, if not, why they failed. As a measure of efficiency, often the time required for completing a given task is determined. A supplementary questionnaire is typically used to assess subjective user satisfaction. Another approach to usability evaluation has been taken in [2]. The developed SWS technologies and their associated tools have been tested within a company, i.e. in a real world environment. Information regarding the usability and applicability

<sup>1</sup><http://sws-challenge.org>

<sup>2</sup><http://www-ags.dfki.uni-sb.de/klusch/s3>

<sup>3</sup><http://www.seals-project.eu>

<sup>4</sup><http://projects.semwebcentral.org/projects/owls-tc>

of the provided tools have been acquired using structured interviews.

Usability evaluations as those mentioned above measure whether the user was able to easily setup, maintain and use a provided tool in the intended way and whether he was satisfied with the tools capabilities. In terms of SWS retrieval, it is typically validated whether the user was able to specify his service requirements using a provided tool and whether he was able to find matching services. However, it is not assessed whether he expressed his actual requirements and whether he made a "good" selection decision in terms of an objective measure.

Finally, to the best of our knowledge, some projects that particularly address end-users (e.g. [4, 9]) did not publish evaluation results at all.

## 4. SYSTEM MODEL

The peculiarities of human decision making as discussed in the previous section shape the general structure and operating mode of systems that can provide effective and efficient support for end-user mediated Web Service retrieval and selection (Fig. 1).

To be able to effectively support consumers in making well-informed, balanced and consistent service selection decisions, those systems require a model of a consumer's service requirements and preferences. While the former allow to identify appropriate service offers, the latter enable the comparison and ranking of those service alternatives according to their suitability. To provide meaningful responses and personalized assistance with service selection, this model needs to be continuously and appropriately updated based on the user's input to the system to correctly reproduce his incrementally emerging requirements (REQUIREMENT R5).

As a consequence of REQUIREMENT R1, the two processes of requirements elicitation and service selection cannot be separated, they rather have to be interwoven into a process of incremental requirements elicitation and service selection that alternates phases of intermediate service recommendation based on partially known requirements and requirements refinement based on the presented service alternatives. Existing mechanisms for SWS discovery and match-making can be adopted to effectively retrieve and rank suitable service offers based on the consumer's known requirements.

## 5. EVALUATION METHODOLOGY

We propose a user-centered methodology for the evaluation of SWS selection. It can be adopted to assess whether an approach implementing the system model in Sect. 4 meets the requirements identified in the Sects. 2 and 4, i.e. effectively supports service consumers in making well-informed, balanced and consistent service selection decisions.

### 5.1 Evaluation Setting

The suggested methodology (Fig. 2) extends a user study design proposed by Chen et al. [3] and proceeds as follows. Participating test users are first asked to think about a kind of service they are interested in. They might not be completely free in their choice, but have to choose from certain service categories for which semantic offer descriptions are available. The participants are then asked to indicate their initial service requirements with respect to the chosen

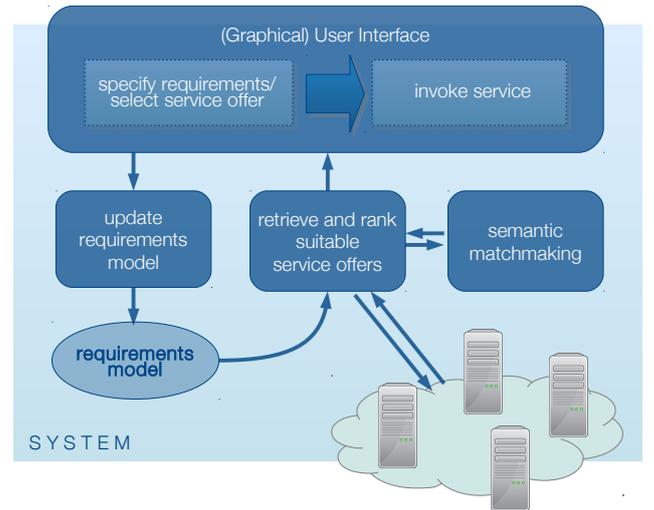


Figure 1: System model

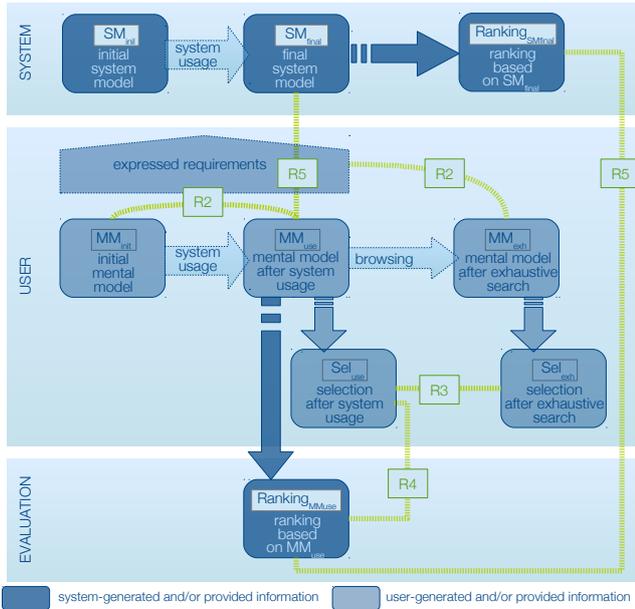
service category (the mental model  $MM_{init}$  of their requirements, cf. Fig. 2). After a (brief) introduction to the system to be tested, the users are asked to autonomously employ the considered solution to learn more about their requirements with respect to the considered service category, to specify these requirements and to finally identify the service offer ( $Sel_{use}$  in Fig. 2) out of a given collection of available service offers that best suits to their requirements. To make the task more challenging, all offers that are presented to the user may be taken from the selected service category. Once a user has chosen an offer, he is asked to state his (updated) service requirements ( $MM_{use}$  in Fig. 2).

In order to verify the quality of the elicited service requirements and the service selection decision that has been made by the user, the study participant is asked to look through a list comprising all available service offers associated with their properties and to check whether he prefers another offer over the selected one ( $Sel_{exh}$  in Fig. 2). To facilitate this task, the list should be scrollable and should be sortable by service property. In case of a switch, the user is asked about the reason for this. This can be either due to a yet unconsidered service alternative or due to a yet unconsidered service requirement that the user became aware of by browsing through the provided list. In the latter case, the user is also asked about his updated service requirements ( $MM_{exh}$  in Fig. 2).

During the test, the user's interactions with the system as well as the states taken by the system-maintained requirements model ( $SM_{init}$  to  $SM_{final}$  in Fig. 2) are logged.

### 5.2 Verifying Requirements Fulfillment

By comparing pairs of the recorded requirements models or the rankings of the available service alternatives that result from them, we can objectively verify different properties of the tested system and thus can decide, if the requirements that have been stated in the Sects. 2 and 4 are fulfilled. These objective findings have to be accompanied by the subjective impressions of the test users assessed e.g. by using a supplementary questionnaire. Note, that REQUIREMENT R1 is automatically fulfilled, if a proposed system is in compliance with the system model provided in Sect. 4.



**Figure 2: Evaluation methodology - the green lines indicate which information have to be compared to verify the indicated requirement**

### *Comprehension and awareness of requirements (R2).*

An inspection of the logged session interactions will reveal whether a proposed system succeeded in stimulating the test users to construct and specify service requirements via the system. However, we also want to find out, if the test users not just express service requirements, but indicate correct requirements. To verify that, we have to compare the service requirements that are indicated by the study participants after having inspected all available service alternatives based on all of their known properties ( $MM_{exh}$  in Fig. 2), and the requirements expressed during the system interaction. Thereby, we have to restrict ourselves to those participants that do not change their requirements after having looked through the entire list of available services. This is to ensure that the model to which the user’s expressed requirements are compared to is correct and complete. A comparison of the initial requirements specified by the participants ( $MM_{init}$  in Fig. 2) and those they provide after having chosen a service by using the tool ( $MM_{use}$  in Fig. 2) will reveal whether a proposed system successfully contributes to the test users awareness of their service requirements.

### *Consumer education (R3).*

To verify whether the test participants are effectively educated about available service alternatives, we have to check whether, after having viewed all available service alternatives, they switched to another service than the one originally selected (comparison of  $Sel_{use}$  and  $Sel_{exh}$  in Fig. 2).

### *Selection consistency (R4).*

To find out, whether a test participant’s selection decision is in fact consistent with his service requirements, we have to determine the rank of the offer that he selected by using the system ( $Sel_{use}$  in Fig. 2) within the service ranking resulting from the user’s actual service requirements, i.e. those

indicated by the user after having used the system ( $MM_{use}$  in Fig. 2).

### *Model consistency (R5).*

The degree of conformance between the user-specified requirements indicated after system usage ( $MM_{use}$  in Fig. 2) and the system-maintained requirements model at selection time ( $SM_{final}$  in Fig. 2) can be assessed by directly comparing both models. However, much more interesting than the discrepancies between the models themselves is the effect these inconsistencies have on the ranking of available service alternatives, i.e. its effect on the recommendation quality of the system. To evaluate this, we have to compare the service ranking which results from the user-indicated requirements model ( $MM_{use}$  in Fig. 2) to the ranking resulting from the system-maintained model ( $SM_{final}$  in Fig. 2).

## 6. CASE STUDY

To demonstrate the feasibility and to validate the appropriateness of the presented evaluation methodology, we show how we implemented it for the evaluation of our own conversational approach to SWS retrieval [6] and present selected evaluation results. We start with a brief summary of the evaluated conversational approach to SWS discovery and selection.

### 6.1 Conversational Service Selection

The suggested solution is in compliance with the system model outlined in Sect. 4, i.e. implements requirements elicitation and service selection as a unified, incremental and interactive process that alternates phases of intermediate service recommendation and requirements refinement. During that process, the user incrementally develops his service requirements and preferences and finally makes a selection decision. We illustrate that process using Fig. 3, which shows the user interface of the suggested retrieval solution. Besides viewing service offers matching to his known requirements, the user may indicate desirable service characteristics, i.e. additional requirements and preferences, based on the presented service alternatives. As an example, think of searching for a service to book a means of transport from the city of Jena to Paris. The system supports three ways of stating requirements: (1) by adding a not yet specified aspect such as *trip duration* to the requirements model, (2) by refining, i.e. subtyping, an aspect’s type, e.g. restricting acceptable carriers to trains, and (3) by critiquing one of the listed service offers. To implement the first two interaction opportunities, the system provides the user with a list of potential aspects, that have not yet been included into the requirements model, but might be added to it, as well as with a list of subtypes that can replace existing aspects’ types (Fig. 3 right). Potential aspects and subtypes are retrieved from the set of matching service offers.

In addition to these interaction opportunities, the user may select a service offer from the presented list, that fits reasonably well to his requirements. Extending previous work in the area of example critiquing recommender systems [3], we also implemented a feature, that allows users to indicate desirable service properties relative to this offer by critiquing it. For example, the user might indicate that the offer is fine, but too expensive (see Fig. 3 left). Based on the indicated property and the properties of the available service alternatives that fulfill this requirement, the system pro-



function over the entire range of values of the considered attribute as occurring in the available offers. Thereby, a fulfillment degree of 0 was assigned to the range minimum and a degree of 1 was assigned to the range maximum. The fulfillment degrees resulting from several constraints/preferences referring to a single service aspect were multiplied.

### 6.2.2 Selected Results and Lessons Learned

Employing the suggested evaluation procedure, we were able to verify that the proposed system effectively stimulates the construction of correct service requirements (REQUIREMENT R2). Thereby, the interaction log revealed on what service aspects the users expressed requirements, how often they refined the type of service aspects and how often they made use of the opportunity to critique available service alternatives. By comparing  $MM_{exh}$  to the requirements expressed during the system interaction, we were able to assess the fraction of the specified service aspects that were in fact relevant to the user as well as the portion of specified constraints and tradeoff opportunities that was in conformance with the user's actual service requirements. A comparison of the initial requirements specified by the participants ( $MM_{init}$  in Fig. 2) and those they provided after having chosen a service by using the tool ( $MM_{use}$  in Fig. 2) revealed that the proposed system successfully contributed to the test users awareness of their service requirements (REQUIREMENT R2). In particular, after having made their final choice, the respondents had revised constraints on service aspects, indicated additional service aspects they have not been aware of before and abandoned requirements on service aspects, that turned out to be of marginal relevance in light of the available service alternatives. The participants also changed the relative importance of their requirements related to the values of the considered service aspects.

After having viewed all available service alternatives, just one of the test persons switched to another service due to having found a service alternative that better fitted to his requirements than the one originally selected. This strongly indicates, that the test participants have been effectively educated about available service alternatives by the system (REQUIREMENT R3).

As it turned out, 90% of the study participants had no hidden requirements they discovered after having viewed all available service offers and their characteristics. The mean rank of the selected service alternative ( $Sel_{use}$ ) with respect to the user's actual requirements ( $MM_{use}$ ) was  $8.33 \pm 5.33$  (out of 200). Hence, the test users' selection decisions were in fact consistent with their service requirements (REQUIREMENT R4). To assess the effect that inconsistencies of the requirements model had on the recommendation quality of the system, we determined the top ten offers with respect to the user-indicated requirements model ( $MM_{use}$ ) and calculated the mean difference between these offers' rank with respect to the user-indicated model and the system-maintained model ( $SM_{final}$ ). As we found, the mean rank difference was  $30.98 \pm 8.50$ , i.e. about 15% of the maximal possible rank difference (200). Investigating the origin of the discrepancy in the rankings revealed, that it partially lies in the inconsistency of the relative weights that the test users assigned to the single service aspects of interest. In particular, the mean rank difference evaluated separately for the 5 fairly consistent weightings (measured in terms of the consistency ratio proposed by [11]) was lower ( $23.04 \pm 7.88$ ).

Hence, though fairly consistent, the weights elicited using the Analytic Hierarchy Process are still inconsistent to some degree. This reduces the significance of the results. However, the quality of the elicited weights is still sufficient to verify the fulfillment of REQUIREMENT R5. Altogether, we have shown that the suggested evaluation method qualifies as a means to verify the requirements discussed in Sect. 2.

## 7. CONCLUSION

We presented a user-centered methodology for the evaluation of SWS selection, that can be adopted to assess whether an approach effectively supports service consumers in making well-informed, balanced and consistent service selection decisions. We demonstrated the feasibility and appropriateness of the suggested approach using a case study.

## References

- [1] Deliverable 2.5.3. Final Summative Evaluation of SOA4All Studio, 2011.
- [2] insemives Deliverable 6.3.2 Test and Validation, 2012.
- [3] L. Chen and P. Pu. Hybrid Critiquing-based Recommender Systems. In *IUI*, pages 22–31, 2007.
- [4] S. Colucci, T. D. Noia, E. D. Sciascio, F. M. Donini, A. Ragone, and R. Rizzi. A Semantic-based Fully Visual Application for Matchmaking and Query Refinement in B2C e-Marketplaces. In *ICEC*, pages 174–184, 2006.
- [5] D. Fensel, F. M. Facca, E. P. B. Simperl, and I. Toma. *Semantic Web Services*. Springer, 2011.
- [6] F. Klan and B. König-Ries. A Conversational Approach to Semantic Web Service Selection. In *EC-Web*, pages 1–12, 2011.
- [7] U. Küster and B. König-Ries. Evaluating Semantic Web Service Matchmaking Effectiveness Based on Graded Relevance. In *Proc. of the 2nd Intl. Workshop SMR<sup>2</sup> at ISWC08*, 2008.
- [8] C. H. Lewis. Using the "Thinking Aloud" Method In Cognitive Interface Design. Technical Report RC-9265, IBM, 1882.
- [9] O. Noppens, M. Luther, T. Liebig, M. Wagner, and M. Paolucci. Ontology-supported Preference Handling for Mobile Music Selection. In *Proc. of the Multidisciplinary Workshop MPREF*, 2006.
- [10] J. W. Payne, J. R. Bettman, and D. A. Schkade. Measuring Constructed Preferences: Towards a Building Code. *Journal of Risk and Uncertainty*, 19(1-3):243–70, December 1999.
- [11] T. L. Saaty. Relative Measurement and its Generalization in Decision Making: Why Pairwise Comparisons are Central in Mathematics for the Measurement of Intangible Factors - The Analytic Hierarchy Process. *RACSAM*, 102(2):251–318, 2008.
- [12] C. Schroth and T. Janner. Web 2.0 and SOA: Converging Concepts Enabling the Internet of Services. *IT Professional*, 9(3):36–41, May 2007.