

CiteSeer^x: a Cloud Perspective

Pradeep B. Teregowda
Pennsylvania State University

Bhuvan Urgaonkar
Pennsylvania State University

C. Lee Giles
Pennsylvania State University

Abstract

Information retrieval applications are good candidates for hosting in a cloud infrastructure. CiteSeer^x a digital library and search engine was built with the goal of efficiently disseminating scientific information and literature over the web. The framework for CiteSeer^x as an application of the SeerSuite software is a design built with extensibility and scalability as fundamental features. This loosely coupled architecture with service oriented interfaces allows the whole or parts of SeerSuite to readily be placed in the cloud. We discuss in brief the architecture, approaches, and advantages of hosting CiteSeer^x in the cloud. We present initial results on costs of migrating whole or parts of CiteSeer^x to two popular cloud offerings as well as discuss the effort involved.

1 Introduction

Digital library search engines have been a continuing topic of research and development for the past several years [4]. The growth in information available both on the Web and from rapid growth in electronic resources make information retrieval systems like CiteSeer^x [2, 16] important for access. At time of this publication the CiteSeer^x collection indexes more than 1.6 million documents and receives several hundred thousand unique visits per day.

The rate of growth of digital information is always a challenge to the effective design of information retrieval systems. Particularly, Web based digital library search engines such as CiteSeer^x can readily take advantage of the reduced maintenance, elasticity, and availability of infrastructure on demand provided by a cloud infrastructure [6].

SeerSuite includes components common to other information retrieval applications. Inspired by services provided by CiteSeer [9], SeerSuite provides among others full text indexing, autonomous citation indexing, and

a personal portal in the form of MyCiteSeer. It extracts and publishes extensive metadata for documents, authors and citations. Its design takes advantage of open source applications such as Tomcat, Solr/Lucene, Java Spring Framework and open source RDBM systems. Advances in and development of automatic metadata extraction for parsing header and citation information have also been important.

SeerSuite based applications share a common set of infrastructure challenges to support a growing set of documents. Although the CiteSeer^x architecture allows hosting of all components and services in the cloud, the size of the CiteSeer^x collection and the amount of data transferred make cloud hosting of CiteSeer^x challenging. There are several cost-effective approaches for solving this problem. We discuss some of these approaches in detail and identify the lessons learned from this analysis. The rest of the paper is arranged in the following manner. Background architecture and services of SeerSuite are discussed in Section 2. The issues of hosting are identified in Section 3. Various strategies for hosting services are discussed in Sections 4, 5 and 6 with future work in Section 7 and conclusions in Section 8.

2 Background and SeerSuite Architecture

2.1 Background

Cloud offerings have taken a number of forms, not limited to Infrastructure, such as Platform and Software as a service approaches. Recent research has focused on adoption, economics and applications. Armbrust et al. [6] explain and quantify benefits from the elasticity of a cloud. They argue that although costs for using cloud may appear higher than buying the hardware, elasticity and the ability to transfer the risk of under/provisioning outweighs the calculated costs. They show that cloud cost makes sense when factors such as cooling, power, and operational costs are taken into ac-

count. Campbell et al. [5] show using simple calculations that for OpenCirrus, the break-even point in terms of server utilization is 33%. A number of other recent papers present simple calculations showing the suitability (or lack thereof) of migrating a certain application to a cloud [20, 18, 15]. Cost, ROI calculators are available from several vendors and consulting groups.

Cloud computing infrastructure for information retrieval and scientific computing have focused on implementing particular features for cloud storage [14] or studying cost benefit trade-offs [8]. Recent research on the role of cloud infrastructure in information retrieval systems has focused primarily on its use for information extraction [12]. Furthermore, the focus has been on the computational costs with little attention to data storage costs with Learned applications pertinent to grid and distributed computing [11]. In contrast, we focus on the use of cloud infrastructure hosted on infrastructure already offered by various vendors.

2.2 Architecture

The following subsections are meant to provide a brief introduction to SeerSuite architecture and services supporting SeerSuite. We also provide a discussion of the feasibility of and refactoring required to migrate these services or components to the cloud. Figure 1 shows various components of SeerSuite. Service-oriented interfaces allow components to be distributed across physical systems. These components can be broadly grouped into those responsible for handling user requests and those handling acquisition and ingestion documents. Among those handling user requests is the Web application which provides presentation and personalization services. The focused crawler, document conversion and extraction, ingestion and maintenance services are responsible for acquiring and ingesting documents. These acquired documents and metadata are then stored in the data storage components for user access.

2.2.1 Web Application

User requests at the Web application are processed with the support of the database, index or the repository. SeerSuite supports interfaces such as the OAI [17] API to allow programmatic access to data stored in the collection. The Web application allows users to search for authors, documents, citations and view document metadata. Some services provided by the Web application require state based interactions with the user, particularly MyCiteSeer. For cloud based hosting minor refactoring will be required to support user authentication with MyCiteSeer. The Web application load depends on traffic, which varies throughout the day, making the Web

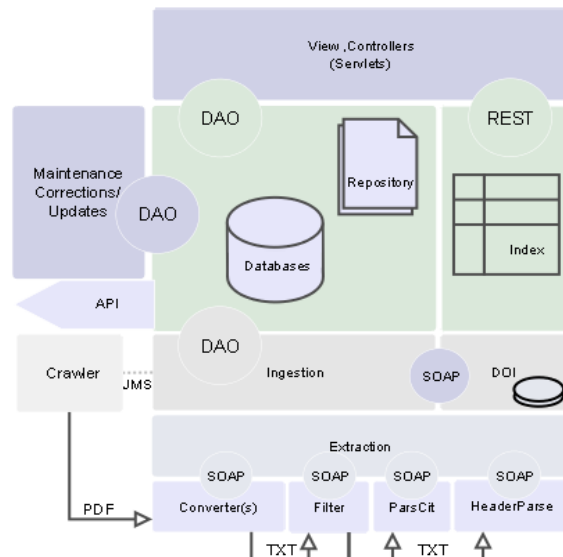


Figure 1: SeerSuite Architecture

application a strong candidate for such a hosting.

2.2.2 Focused Crawler

Document acquisition drives the growth of SeerSuite instances. In particular, focused crawlers [7] help acquire relevant documents efficiently. The focused crawler is a strong candidate for cloud hosting, since it can take advantage of the elasticity and on demand provisioning with efficiently scheduled crawls.

2.2.3 Document Conversion and Information Extraction

Before documents can be processed by the extraction system, the documents in PDF or PostScript formats are converted into text and filtered to remove documents not containing citations. Documents acquired from the Web are processed by multiple modules which extract extensive document, citation, and author metadata. These modules are based on state-of-the-art machine learning techniques. Prominent among these is the header parser, which extracts document and author information. The ParsCit module is utilized to extract citation information. The metadata extraction system is not a strong candidate for Platform-as-a-Service cloud offerings (e.g., Windows Azure), as extensive refactoring will be required.

2.2.4 Document Ingestion

The documents processed by the extraction and conversion service are ingested into the system. This includes adding the document and related metadata to the

database and to the repository. By comparing the checksum of the document to be ingested, the ingestion system avoids adding duplicates into the collection. Documents are assigned a unique document object identifier from the DOI service. The use of service oriented interfaces and the minimal code footprint allow the ingestion system to be easily hosted in the cloud.

2.2.5 Data Storage

Persistence of data extracted is achieved by the use of index, databases and file storage components (repository). The database is utilized by the web application to provide document summaries and metadata. The index allows users to query the full text and citation information. The repository caches documents crawled by the crawler and metadata extracted by the extractors.

2.2.6 Maintenance Service

Tasks not part of the ingestion system such as updates to the index, inference based metadata updates, charts and generation of citation charts and statistics are performed by the maintenance system. The maintenance systems generate very little data, and can be scheduled by the administrator. These services need to be closer to the data storage due to vast amount of information processed for each iteration of their operation. Their candidacy is hence dependent on the hosting of the data storage components.

2.2.7 Federated Services

SeerSuite provides several features that are not part of the main application. These features are supported by services which may not share the same framework or application components but share infrastructure. Many of these services are under development. Such components are strong candidates for migration into the cloud since they can take advantage of the pay-as-you-go charging offered by cloud products.

2.3 Deployment

The current deployment of SeerSuite as CiteSeer^x is on a group of heterogeneous server machines. Two Web application instances are hosted on the Apache Tomcat platform in a cluster. Each Web application instance is hosted on a machine with two dual core CPUs and 16GB of RAM. The Web traffic is load balanced through a software-based L4 load balancer cluster. The database and the repository are hosted on separate machines with large storage (> 15 TB), dual core dual CPUs, and 16GB of RAM. MySQL is used as the RDBMS for the system. Indices for document, tables, author names are hosted separately on machines with dual core dual CPU and

16GB of RAM. The repository is accessed by the Web servers using Global File System (GFS) over Global Network Block Device.

2.4 Terminology

A description of terms used in future discussions, relevant in the context of SeerSuite are provided below.

Request Types: User requests can be grouped into search, document views, MyCiteSeer and others or misc. The search request and document view requests involve the Web application, database and the index. MyCiteSeer requests involve the Web application and database. Others include requests for stylesheets, images.

Peak Load: This represents a set of requests observed at the web server, exceeding a set threshold of requests per second (90th percentile).

3 Problem Definition

SeerSuite as a whole can be hosted in Infrastructure as a Service platforms with minimal refactoring. Such a hosting, however, is expensive with current cloud offerings. This is due to the large collection size and the volume of data transferred between the application and the user. The key question we are interested in answering in this context is *moving which sections, components or subset of CiteSeer^x to a cloud would be most cost effective?*

To answer this question, we consider the entire application with particular focus on the Web application including its supporting components, the index, database and the repository. We present three different approaches by presenting first a hypothesis and discuss the cost and implications. We utilize the existing log monitoring data collected from CiteSeer^x. The Web application logs for a period of 15 days were analyzed to obtain data used in the following sections. Figure 2 shows the number of requests made during this period along with the type of request.

4 Component Hosting

From our discussion of SeerSuite, it is possible to identify specific components in SeerSuite to be hosted. Each component hosted in the cloud includes all of its associated modules and subsystems. For example, hosting of the SeerSuite index includes hosting of the storage for the index, application code, and interfaces.

We consider components of the system, choosing components based on size, data transfer, and feasibility of migration. From Table 1, we see the cost of hosting is dominated by the cost of computation (per instance cost). The other major cost components include the cost of data

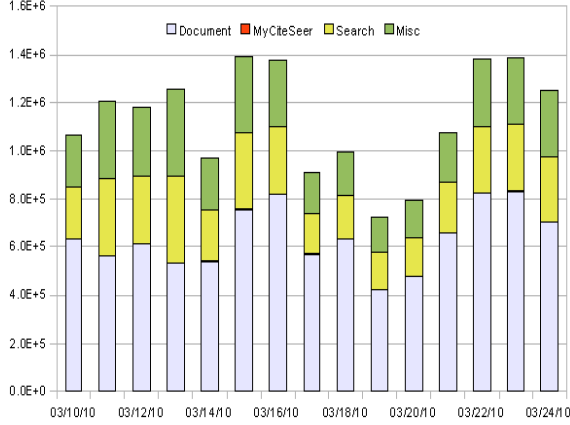


Figure 2: Traffic By Request Type

storage and data transferred in and out of the cloud. For hosting components the compute costs are a constant, since the components always operational. Therefore we now address the data storage and transfer costs.

Figure 3 shows the flow of data between components of CiteSeer^x. Data for creating this graph was obtained from log files and application specific monitors. In the case of crawlers, the information extraction data flow was assumed proportional to the number of documents acquired and processed. This graph is useful for determining candidates for hosting. For example, if CiteSeer^x were hosted entirely within the cloud infrastructure, the amount of data stored in the cloud would be 1.7 TB, with 3.2 TB of data transferred between the user, web and the application. Clearly, the repository is the largest com-

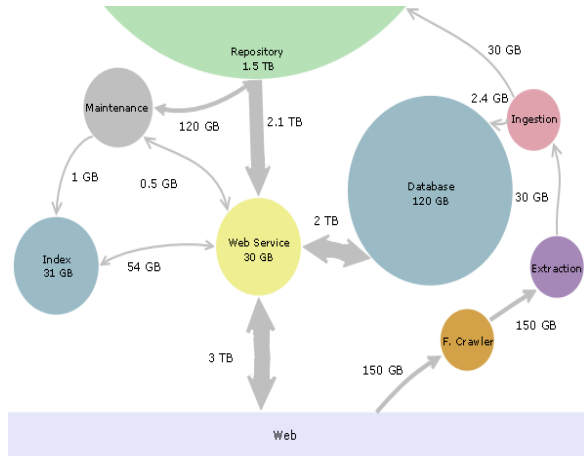


Figure 3: Data Flow - CiteSeer^x Components

ponent in size, while the web application has the largest data volume. We provide a cost estimate for the components based on cloud infrastructure services offered by

Amazon EC2 [1] and Google App Engine [3]. Cost estimates are based on a 30 day month.

Choice of vendors is a result of support in terms of environment and libraries offered or supported by these vendors. In the case of Amazon EC2, we consider a mapping of one to one to an extra large instance for hosting the database, application, index, repository, extraction and crawler services. We assume that additional instances are provided as required in Google App Engine with no additional cost.

Cost		Amazon	Google
Initial Setup	Data In	1820.4	0
Monthly	Stored	1820.4	182.04
	Data In	152	0
	Data Out	3072	460.8
	Trans.	368	190.77
	CPU	30*24	2937.6
Total Monthly		\$3771.21	\$800.9

Table 1: CiteSeer^x Hosting

Table 1 provides the cost of hosting CiteSeer^x in the cloud for a month. We now examine the cost of hosting individual components in the cloud, with all other services hosted locally. Estimates are provided in Table 2. CiteSeer^x migration costs include initial setup costs, as

Component	A. EC2		G. App Engine	
	Initial	Month	Initial	Month
Web Service	0	1448.18	0	942.53
Repository	0	1011.88	163.8	593.21
Database	0	858.89	12	348.05
Index	0	527.08	3.1	83.48
Extraction	0	499.02	0	90.6
Crawler	0	513.4	0	105

Table 2: Component Costs (USD) in the Cloud

a substantial collection already exists. New applications may not incur initial data transfer costs.

Note that Amazon currently provides free data transfers into the cloud. If this were not the case, hosting services on Amazon would be much more expensive and also incur initial setup costs. Calculating the cost of hosting the entire application leads to a figure of \$3771 for Amazon EC2 and \$800 for Google App Engine. The cost of hosting components also lends support to the conclusions drawn about data and access in [10].

Individual components hosted on the cloud have implications beyond the cost of hosting them in the cloud. Costs related to refactoring code for migration has not been accounted for in Table 2. In the case of Google App Engine, existing code written in languages not sup-

ported by App Engine will require significant refactoring. Along with components hosted in the cloud, components hosted locally may require refactoring. This refactoring is minimal if the service or component utilized a service oriented interface and significant when services are closely coupled.

Lessons Learned: If the cost of hosting an entire service is prohibitive, hosting components may be a reasonable approach to taking advantage of cloud infrastructure. The cost effectiveness of such an approach depends on data transferred through the service. Loosely coupled components are easier to migrate. For existing components and code, refactoring costs will provide a closer estimates of costs. This approach is suitable, when a fixed budget constrains the placement of services or components. By identifying components, data transfer and refactoring costs a hosting solution can be identified.

5 Content Hosting

Content particularly static images, stylesheets, javascript common to most web pages need not be hosted locally. An analysis of peak traffic at the web services provides an insight on how this can be achieved. From analysis of figure 4, we see that most requests for peak traffic are for such content. In this case the amount of computation required and data stored on the cloud is small, the cost of hosting is cost-effective. The total size of all files to be

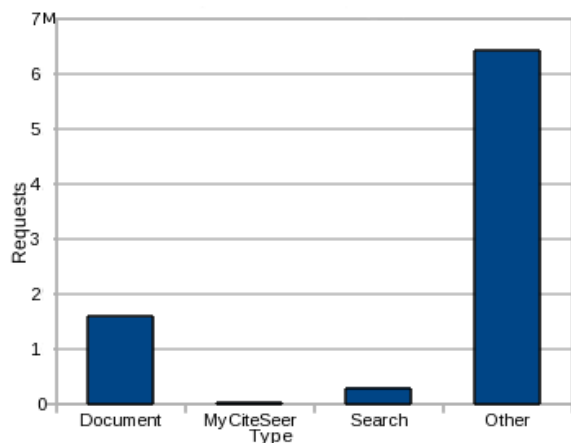


Figure 4: Request Types at Peak

placed on the cloud is 2.24 MB. By hosting these files in the cloud, the amount of data transferred for CiteSeer^x from the cloud is 390.26 GB costing less than \$142 per month (on both Amazon EC2 and Google App Engine services including a small instance cost). While this is a small part of more than 3 TB of data volume between the application and the user, it helps the system satisfy a significant number of peak load requests.

The same approach can be used to identify elements such as a subset of the repository to be placed in the cloud. Such an approach would involve identifying the most commonly accessed documents and placing them both locally and in the cloud. During peak loads, clients can be directed to the cloud for access.

Lesson Learned: Hosting specific content relevant to peak load scenarios in the cloud can be beneficial, and the simplest approach to hosting services in the cloud.

6 Load based partitioning

This approach is particularly important for supporting the growth in traffic, flash crowds providing users access to service.

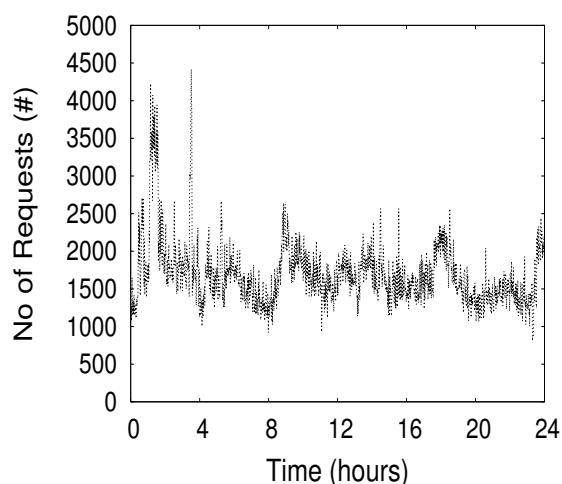


Figure 5: Number of Requests per Second

Figure 5 shows the requests received at the web server. From the graph we identify that the 90th percentile is represented by 60 requests per second. Most of these requests are for elements associated with presentation (javascript and stylesheets). Assuming that the traffic growth continues at the same pace and as more features (Algorithm and Figure search) are added, There is a need for provisioning more systems. Instead of procuring these systems, infrastructure at the cloud can be considered to fulfill this need.

Further examination provides evidence of self-similarity in the request arrival process, which has interesting implications for resource provisioning. The peak resource needs of several CiteSeer^x components are significantly higher than their average-case (or even a high percentile) needs.

Two strategies are possible in partitioning based on load. Of these, one strategy would be to host a copy of the entire application in the cloud, using load balancers

to identify and direct traffic during peak load conditions. Table 3 provides the costs of such a hosting solution for CiteSeer^x in Amazon EC2 and Google App Engine. All data measurements are in GB, and transaction measurements in transactions per second obtained via iostat.

Cost			Amazon	Google
Initial Setup	Data In	1820.4	0	182
Monthly	Stored	1820.4	182.04	273.06
	Data In	14.78	0	1.48
	Data Out	298.7	44.8	35.84
	Trans.	368	9.27	0
	CPU	70	285.6	7
Total Monthly			\$521.71	\$317.38

Table 3: CiteSeer^x Peak Load Hosting

These costs can be considered in comparison to the cost of procuring, maintaining systems. Savings by avoiding adoption of storage systems locally add to the attractiveness of cloud infrastructure.

An alternate approach would be to host only the component under stress in the cloud. For example, a database replica to support a locally hosted database could be deployed in the cloud. If this instance were used only during peak load conditions, the costs would decrease to \$385, since the instance would be in use for 70 hours.

Lessons Learned: By utilizing a replica or subset of the application for handling only peak loads, we can take advantage of cloud infrastructure in a cost-effective manner. This can resolve issues stemming from the growth of the collection and user traffic.

7 Future Work

We explored various strategies for hosting SeerSuite in the cloud. In Section 6 we briefly mentioned the temporal nature of traffic and user behavior. By identifying user patterns, the hosting solutions can be optimized to take advantage these patterns. While this discussion included the Amazon EC2 and Google App Engine for cost comparison, this work needs to be extended by examining in depth options offered by other cloud offerings, private clouds and virtualization solutions.

Products such as private clouds offered Eucalyptus [13] can be utilized to take advantage of hardware already existing as part of the system. Components related to user interaction with CiteSeer^x hosting with services like Amazon Virtual Private Clouds, and local clouds can be considered for these services.

Impact of including cloud hosted services on other services has not be considered in the current discussion. Inclusion of cloud services could require significant refactoring and changes to maintenance cycles. Is-

ues with how issues like latency, load balancing have not been addressed in this paper. Several opportunities exist within SeerSuite framework for adopting virtualization and cloud infrastructure. In particular, the repository can be restructured to take advantage of cloud based storage solutions in an effective manner. Hadoop [19] based metadata extraction and log analysis systems could enable faster document acquisition.

8 Conclusions

Preliminary costs for hosting SeerSuite instances such as CiteSeer^x in the cloud prove reasonable. We develop different approaches that can be adopted either for their cost-efficiency, simplicity, or handling peak loads. Cost estimation for each approach along with lessons learned from analysis provide a guideline for further exploration. Our future work would focus on adoption of virtualization and extraction systems suitable for hosting in the cloud. In addition to these goals, we would like to examine user behavior, issues in privacy, security for components hosted in the cloud that were not discussed in this work. As part of these discussions, we have presented a detailed examination of the existing deployment of SeerSuite in CiteSeer^x.

9 Acknowledgments

The authors acknowledge partial support from NSF CISE. Urgaonkar’s research was funded in part by NSF grants CCF-0811670 and CNS-0720456.

References

- [1] Amazon EC2. <http://aws.amazon.com/ec2/>.
- [2] CiteSeer^x. <http://citeseerx.ist.psu.edu/>.
- [3] Google App Engine. <http://code.google.com/appengine/>.
- [4] Greenstone. <http://www.greenstone.org/>.
- [5] Open cirrus: Cloud computing research testbed. <https://opencirrus.org/>.
- [6] ARMBRUST, M., FOX, A., GRIFFITH, R., JOSEPH, A., KATZ, R., KONWINSKI, A., LEE, G., PATTERSON, D., RABKIN, A., STOICA, I., ET AL. Above the clouds: A Berkeley view of cloud computing. *EECS Department, University of California, Berkeley, Tech. Rep. UCB/EECS-2009-28* (2009).
- [7] CHAKRABARTI, S., VAN DEN BERG, M., AND DOM, B. Focused crawling: a new approach to topic-specific Web resource discovery. *Computer Networks* 31, 11-16 (1999), 1623–1640.
- [8] DEELMAN, E., SINGH, G., LIVNY, M., BERRIMAN, B., AND GOOD, J. The cost of doing science on the cloud: the montage example. In *Proceedings of the 2008 ACM/IEEE conference on Supercomputing* (2008), pp. 1–12.
- [9] GILES, C., BOLLACKER, K., AND LAWRENCE, S. CiteSeer: An automatic citation indexing system. In *Proceedings of the third ACM conference on Digital libraries* (1998), ACM New York, NY, USA, pp. 89–98.

- [10] GRAY, J. Distributed computing economics. Tech. rep., Microsoft Research, 2003.
- [11] GROSSMAN, R., AND GU, Y. Data mining using high performance data clouds: Experimental studies using sector and sphere. In *Proceeding of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining* (2008), ACM, pp. 920–927.
- [12] MIKA, P., AND TUMMARELLO, G. Web semantics in the clouds. *IEEE Intelligent Systems* 23, 5 (2008), 82–87.
- [13] NURMI, D., WOLSKI, R., GRZEGORCZYK, C., OBERTELLI, G., SOMAN, S., YOUSEFF, L., AND ZAGORODNOV, D. The eucalyptus open-source cloud-computing system. In *Proceedings of the 2009 9th IEEE/ACM International Symposium on Cluster Computing and the Grid-Volume 00* (2009), IEEE Computer Society, pp. 124–131.
- [14] SINGH, A., SRIVATSA, M., AND LIU, L. Search-as-a-service: Outsourced search over outsourced storage. *ACM Trans. Web* 3, 4 (2009), 1–33.
- [15] TERELOWDA, P. URGANOKAR, B., AND GILES, C. Cloud computing: A digital libraries perspective. In *3rd IEEE 2010 International Conference on Cloud Computing* (2010).
- [16] TERELOWDA, P. B., COUNCILL, I. G., FERNANDEZ, J. P. R., KASBHA, M., ZHENG, S., AND GILES, L. C. Seersuite: Developing a scalable and reliable application framework for building digital libraries by crawling the web. In *USENIX Conference on Web Application Development* (2010).
- [17] VAN DE SOMPEL, H., NELSON, M., LAGOZE, C., AND WARNER, S. Resource harvesting within the OAI-PMH framework. *D-Lib Magazine* 10, 12 (2004), 1082–9873.
- [18] WALKER, E., BRISKEN, W., AND ROMNEY, J. To lease or not to lease from storage clouds. *Computer* 43 (2010), 44–50.
- [19] WEIGEL, F., PANDA, B., RIEDEWALD, M., GEHRKE, J., AND CALIMLIM, M. Large-scale collaborative analysis and extraction of web data. *Proc. VLDB Endow.* 1, 2 (2008), 1476–1479.
- [20] WOOD, T., CECCHET, E., RAMAKRISHNANY, K., SHENOY, P., VAN DER MERWEY, J., AND VENKATARAMANI, A. Disaster recovery as a cloud service: Economic benefits & deployment challenges. In *2nd USENIX Workshop on Hot Topics in Cloud Computing* (2010).