

SIGMOD Officers, Committees, and Awardees

Chair	Vice-Chair	Secretary/Treasurer
Yannis Ioannidis University of Athens Department of Informatics Panepistimioupolis, Informatics Bldg 157 84 Ilissia, Athens HELLAS +30 210 727 5224 <yannis AT di.uoa.gr>	Christian S. Jensen Department of Computer Science Aalborg University Selma Lagerlöfs Vej 300 DK-9220 Aalborg Øst DENMARK +45 99 40 89 00 <csj AT cs.aau.dk >	Alexandros Labrinidis Department of Computer Science University of Pittsburgh Pittsburgh, PA 15260-9161 USA +1 412 624 8843 <labrinid AT cs.pitt.edu>

SIGMOD Executive Committee:

Sihem Amer-Yahia, Curtis Dyreson, Christian S. Jensen, Yannis Ioannidis, Alexandros Labrinidis, Maurizio Lenzerini, Ioana Manolescu, Lisa Singh, Raghu Ramakrishnan, and Jeffrey Xu Yu.

Advisory Board:

Raghu Ramakrishnan (Chair), Yahoo! Research, <First8CharsOfLastName AT yahoo-inc.com>, Rakesh Agrawal, Phil Bernstein, Peter Buneman, David DeWitt, Hector Garcia-Molina, Masaru Kitsuregawa, Jiawei Han, Alberto Laender, Tamer Özsu, Krithi Ramamritham, Hans-Jörg Schek, Rick Snodgrass, and Gerhard Weikum.

Information Director:

Jeffrey Xu Yu, The Chinese University of Hong Kong, <yu AT se.cuhk.edu.hk>

Associate Information Directors:

Marcelo Arenas, Denilson Barbosa, Ugur Cetintemel, Manfred Jeusfeld, Dongwon Lee, Michael Ley, Rachel Pottinger, Altigran Soares da Silva, and Jun Yang.

SIGMOD Record Editor:

Ioana Manolescu, INRIA Saclay, <ioana.manolescu AT inria.fr>

SIGMOD Record Associate Editors:

Magdalena Balazinska, Denilson Barbosa, Chee Yong Chan, Ugur Çetintemel, Brian Cooper, Cesar Galindo-Legaria, Leonid Libkin, and Marianne Winslett.

SIGMOD DiSC and SIGMOD Anthology Editor:

Curtis Dyreson, Washington State University, <cdyreson AT eecs.wsu.edu>

SIGMOD Conference Coordinators:

Sihem Amer-Yahia, Yahoo! Research, <sihemameryahia AT acm.org> and Lisa Singh, Georgetown University, <singh AT cs.georgetown.edu>

PODS Executive: Maurizio Lenzerini (Chair), University of Roma 1, <lenzerini AT dis.uniroma1.it>, Georg Gottlob, Phokion G. Kolaitis, Leonid Libkin, Jan Paradaens, and Jianwen Su.

Sister Society Liaisons:

Raghu Ramakrishnan (SIGKDD), Yannis Ioannidis (EDBT Endowment).

Awards Committee:

Laura Haas (Chair), IBM Almaden Research Center, <laura AT almaden.ibm.com>, Rakesh Agrawal, Peter Buneman, and Masaru Kitsuregawa.

Jim Gray Doctoral Dissertation Award Committee:

Johannes Gehrke (Co-chair), Cornell Univ.; Beng Chin Ooi (Co-chair), National Univ. of Singapore, Alfons Kemper, Hank Korth, Alberto Laender, Boon Thau Loo, Timos Sellis, and Kyu-Young Whang

SIGMOD Officers, Committees, and Awardees (continued)

SIGMOD Edgar F. Codd Innovations Award

For innovative and highly significant contributions of enduring value to the development, understanding, or use of database systems and databases. Until 2003, this award was known as the "SIGMOD Innovations Award." In 2004, SIGMOD, with the unanimous approval of ACM Council, decided to rename the award to honor Dr. E.F. (Ted) Codd (1923 - 2003) who invented the relational data model and was responsible for the significant development of the database field as a scientific discipline. Recipients of the award are the following:

Michael Stonebraker (1992)	Jim Gray (1993)	Philip Bernstein (1994)
David DeWitt (1995)	C. Mohan (1996)	David Maier (1997)
Serge Abiteboul (1998)	Hector Garcia-Molina (1999)	Rakesh Agrawal (2000)
Rudolf Bayer (2001)	Patricia Selinger (2002)	Don Chamberlin (2003)
Ronald Fagin (2004)	Michael Carey (2005)	Jeffrey D. Ullman (2006)
Jennifer Widom (2007)	Moshe Y. Vardi (2008)	Masaru Kitsuregawa (2009)
Umeshwar Dayal (2010)		

SIGMOD Contributions Award

For significant contributions to the field of database systems through research funding, education, and professional services. Recipients of the award are the following:

Maria Zemankova (1992)	Gio Wiederhold (1995)	Yahiko Kambayashi (1995)
Jeffrey Ullman (1996)	Avi Silberschatz (1997)	Won Kim (1998)
Raghu Ramakrishnan (1999)	Michael Carey (2000)	Laura Haas (2000)
Daniel Rosenkrantz (2001)	Richard Snodgrass (2002)	Michael Ley (2003)
Surajit Chaudhuri (2004)	Hongjun Lu (2005)	Tamer Özsu (2006)
Hans-Jörg Schek (2007)	Klaus R. Dittrich (2008)	Beng Chin Ooi (2009)
David Lomet (2010)		

SIGMOD Jim Gray Doctoral Dissertation Award

SIGMOD has established the annual SIGMOD Jim Gray Doctoral Dissertation Award to *recognize excellent research by doctoral candidates in the database field.* This award, which was previously known as the SIGMOD Doctoral Dissertation Award, was renamed in 2008 with the unanimous approval of ACM Council in honor of Dr. Jim Gray. Recipients of the award are the following:

- **2006 Winner:** Gerome Miklau, University of Washington. *Runners-up:* Marcelo Arenas, Univ. of Toronto; Yanlei Diao, Univ. of California at Berkeley.
- **2007 Winner:** Boon Thau Loo, University of California at Berkeley. *Honorable Mentions:* Xifeng Yan, UIUC; Martin Theobald, Saarland University
- **2008 Winner:** Ariel Fuxman, University of Toronto. *Honorable Mentions:* Cong Yu, University of Michigan; Nilesh Dalvi, University of Washington.
- **2009 Winner:** Daniel Abadi, MIT. *Honorable Mentions:* Bee-Chung Chen, University of Wisconsin at Madison; Ashwin Machanavajjhala, Cornell University.
- **2010 Winner:** Christopher Ré (advisor: Dan Suciu), University of Washington
Honorable Mentions: Soumyadeb Mitra (advisor: Marianne Winslett), University of Illinois, Urbana-Champaign; Fabian Suchanek (advisor: Gerhard Weikum), Max-Planck Institute for Informatics

A complete listing of all SIGMOD Awards is available at: <http://www.sigmod.org/awards/>

Editor's Notes

Welcome to the December 2010 issue of the ACM SIGMOD Record.

The first contribution in this issue is a regular article by Martinez-Gil and Aldana-Montes, on the topic of “Reverse Ontology Matching”. Ontology matching is the well-established problem of identifying possible alignments between two existing ontologies; the process is typically performed with the help of specific matching functions, and often involves also a human expert. The article considers the process of reverse-engineering ontology matching, by finding the matching functions used in a specific ontology alignment execution. Reverse ontology matching helps in order to capture the reasoning of a human expert aligning the ontologies, to perform matching by example, or to reverse-engineer existing tools.

The issue continues with a very timely survey of “Scalable SQL and NoSQL Data Stores” by Cattell, from the perspective of their support of OLTP workloads involving many users performing concurrent reads and updates. The field is so dynamic that the author had already revised his report half a dozen times in the year and the half before submitting it, and you will find in the paper that some players in the field had bought others by the time the draft was completed! The survey is rich with clearly-identified dimensions for comparing and analyzing the systems considered; it is a valuable introduction to a rich and dynamic data management field.

Chris Ré is the guest of our “Distinguished Profiles in Databases” column. Chris is the recipient of the 2010 SIGMOD Jim Gray Dissertation Award, the runners-up being Soumyadeb Mitra and Fabian Suchanek. Chris talks about his first encounter with databases and being spirited away from mathematics upon hearing Jim Gray’s present his work on the World Wide Telescope as a guest lecturer in a database course! Chris talks about learning by watching his PhD supervisor Dan Suciu and later on other co-workers, about the many applications of probabilistic databases and about the difficulty of staying focused in the face of so many exciting topics.

The “Open Forum” column is represented with an article perfectly fitting its scope: “Paper Bricks: an Alternative to Complete-Story Peer Reviewing”, by Dittrich. The paper is the companion of the winner CIDR 2011 5-minutes gong show (check it out on YouTube!) and has been quite widely discussed in the short interval between CIDR and the publication in the Record. I am pleased that the Record gets the opportunity to publish such thought-provoking articles on how our scientific community could alternatively organize itself. I am convinced the discussion around Jens Dittrich’s idea will be fruitful in the short and long term, regardless of how much of it will actually be implemented. And no, the funny letters on top of the first page are not a typographic mistake!

The Record continues to welcome your contributions through the RECESS submission site (<http://db.cs.pitt.edu/recess>); be sure to peruse the new Editorial Policy on the SIGMOD Record’s Web site (<http://www.sigmod.org/publications/sigmod-record/sigmod-record-editorial-policy>).

Ioana Manolescu

April 2011

Past SIGMOD Record Editors:

Harrison R. Morse (1969)
Daniel O'Connell (1971 – 1973)
Randall Rustin (1975)
Thomas J. Cook (1981 – 1983)
Jon D. Clark (1984 – 1985)
Margaret H. Dunham (1986 – 1988)
Arie Segev (1989 – 1995)
Jennifer Widom (1995 – 1996)
Michael Franklin (1996 – 2000)
Ling Liu (2000 – 2004)
Mario Nascimento (2005 – 2007)
Alexandros Labrinidis (2007 – 2009)

Reverse Ontology Matching

Jorge Martinez-Gil
University of Malaga
Dept. of Computing Sciences
Boulevard Louis Pasteur 35, 29071 Malaga
jorgemar@lcc.uma.es

Jose F. Aldana-Montes
University of Malaga
Dept. of Computing Sciences
Boulevard Louis Pasteur 35, 29071 Malaga
jfam@lcc.uma.es

ABSTRACT

Ontology Matching aims to find the semantic correspondences between ontologies that belong to a single domain but that have been developed separately. However, there are still some problem areas to be solved, because experts are still needed to supervise the matching processes and an efficient way to reuse the alignments has not yet been found. We propose a novel technique named Reverse Ontology Matching, which aims to find the matching functions that were used in the original process. The use of these functions is very useful for aspects such as modeling behavior from experts, performing matching-by-example, reverse engineering existing ontology matching tools or compressing ontology alignment repositories. Moreover, the results obtained from a widely used benchmark dataset provide evidence of the effectiveness of this approach.

1. INTRODUCTION

In the new approaches to develop information systems, the use of a type of formal schema called ontology is usual. Ontologies are considered to be semantically richer than schemas in general, and therefore, techniques for schema matching can be easily adapted to ontologies but not vice versa [12].

There are many ontologies available on the web currently. These ontologies are usually developed for different collections of information, and different kinds of applications. Nowadays, the Swoogle search engine¹ has indexed thousands of ontologies. There are several reasons for the quick proliferation of ontologies, but we consider mainly two:

- It is often easier to construct a new ontology, than find an existing one which is appropriate for a given task.
- There is often a desire for direct control over the ontology for a particular domain, rather than having the structure dictated by external forces.

¹<http://swoogle.umbc.edu>

A direct consequence of having large numbers of ontologies available is that it is necessary to integrate knowledge which is represented in different ways. Ontology matching aims to produce alignments, that is, sets of semantic correspondences between elements from different ontologies. This task is very expensive in terms of time and resource consumption. The reason is that it is necessary a lot of work from domain experts to match ontologies or to supervise results from existing semiautomatic tools. Our approach is based on the extraction of the ontology matching functions used by the agents, experts or tools when matching ontologies, so it is a powerful way to reuse, store and understand their knowledge. Moreover, there are other collateral benefits as the ability to implement strategies for ontology matching by example, reverse engineering existing ontology matching tools or compress large ontology alignment repositories. In this way, we think that the main contributions of our work can be summarized as follow:

- We propose, for the first time to the best of our knowledge, a methodology for reverse engineering an ontology alignment which tries to find the matching function that have been used to generate an ontology alignment.
- We perform an empirical evaluation of our approach in order to show its practical viability in the real world.

The rest of this work is structured in the following way: Section 2 describes the problem statement related to Reverse Ontology Matching. Section 3 presents the related works regarding other reverse engineering proposals. Section 4 presents the core of our approach, a methodology for reverse engineer an ontology alignment, and some real examples. Section 5 contains an evaluation that shows the applicability of Reverse Ontology Matching in the practice. In Section 6, we describe the conclusions extracted from this work.

2. PROBLEM STATEMENT

An ontology is “a specification of a conceptualization” [9] that it is to say, an abstract representation of the world like a set of objects. In this work, we are going to use the intuitive notion of ontology as a set of classes with relations among them that serves primarily to represent formal knowledge in a way which is understandable by people and machines.

Definition 1 (Ontology Matching Function). *Ontology Matching Function is a function f where, given two input ontologies o and o' , an (optional) input incomplete alignment A , a set of configuration parameters p and, a set of external resources r , an alignment y is returned.*

$$y = f(o, o', A, p, r) \quad (1)$$

Definition 2 (Ontology Alignment). *An ontology alignment is a set of mappings, thus, a set of tuples in the form (id, e, e', n, R) . Where id is an unique identifier, e and e' are entities belonging to two different ontologies, R is the relation of correspondence between these entities and n is a real number between 0 and 1 representing the confidence for R .*

Definition 3 (Reverse Ontology Matching Function). *We define Reverse Ontology Matching as the function g that has been used for obtaining an alignment y using two input ontologies o and o' , an (optional) input incomplete alignment A , a set of configuration parameters p and, a set of external resources r according to the following equation:*

$$y = f(o, o', A, p, r) \rightarrow \exists g, g(o, o', A, p, r, y) = f \quad (2)$$

The computation of f is far from being trivial. There are hundreds of algorithms to match ontologies, and everything indicates that more algorithms will appear. Moreover, matching algorithms can be compose so we obtain a solution space populated by many millions of possibilities. Lastly, in alignments created by humans it is possible to find mappings that have been found using heterogeneous rules [13]. Therefore, we cannot be sure that we are going to obtain the original function, but a function that is equivalent to the original one.

Definition 4 (Equivalent Ontology Matching Function). *Let f be an ontology matching function, then we define an Equivalent Ontology Matching Function f' as a function which return the same result that f for the same input. More formally*

$$f' \equiv f \leftrightarrow f'(o, o', A, p, r) = f(o, o', A, p, r) \quad (3)$$

These two basic ideas behind the notions of Reverse Ontology Matching Function and Equivalent Ontology Matching Function allow us to formulate the definition Equivalent Reverse Ontology Matching Function as follows.

Definition 5 (Equivalent Reverse Ontology Matching Function). *Let g be a reverse ontology matching function, then we define a Equivalent Reverse Ontology Matching Function g' as an function which return the same result that g for the same input. More formally*

$$g' \equiv g \leftrightarrow g'(o, o', A, p, r, y) = g(o, o', A, p, r, y) \quad (4)$$

Take into account, that given an alignment between two ontologies, we know nothing about the heuristic used for the expert in order to provide the results. However, we know two main things: firstly, in Ontology Matching there are a limited amount of categories for grouping algorithms with similar behaviors, and secondly, we understand the notion of composite matchers, that it is to say, the idea behind to combine similarity values predicted by multiple algorithms to determine correspondences between ontology elements in order to benefit from the high degree of precision of some algorithms and at the same time the broader coverage of other algorithms [6]. For the rest of this work, we are working under the assumption that the agent, expert or tool which generate an initial alignment always try to maximize the precision and coverage of their solutions.

2.1 Use Cases

There are many applications where reverse ontology matching will be very useful. We are going to show here four of them: a) Capturing behavior from experts, b) Matching by example, c) Reverse engineer existing tools, and d) Compressing large volumes of ontology alignments.

2.1.1 Capturing behavior from experts.

If we ask an expert for creating mappings between two ontologies that belong to its area of expertise but that have been developed separately, we are going to obtain a few correspondences but nothing else. These correspondences are only useful for the current case and we are not going to be able to get profit from them in the future. But, if we reverse engineer these correspondences, we will be able to obtain the heuristic used by the expert and apply

it in a lot of additional scenarios. Moreover, we can become experts because we are not going to see only results but the way to reach these results. Finally, but not least, we can compare heuristics from a wide variety of experts and obtain easily a core heuristic, thus, a common way to solve problems.

2.1.2 *Matching by example.*

In many ontology matching scenarios is popular the use of a technique called matching by example. This technique consists of given two ontologies, try to find several samples correspondences in order to the system may learn how to find the rest of existing correspondences between the two ontologies automatically. In this way, the user only has to do a little part of the work manually. The existing techniques use methods from the machine learning field (e.g. genetic algorithms, neural networks, and so on). In this way, better the set of mappings provided by the user larger the quality of the automatic matching to be performed. One of the advantages of reverse ontology matching functions is that can be computed in real time, so it is possible to compute the equivalent reverse matching function for a little set of mappings in order to apply this function to the rest of the given ontologies. If the user is able to provide all possible cases initially, the automatic part of the matching process will be very good.

2.1.3 *Reverse engineer existing tools.*

Author of the initial set of mappings is not relevant for our reverse ontology matching approach. This means that is possible to detect, and therefore to simulate, an equivalent working mode for the most of deterministic ontology matching tools. Deterministic here means that for a given input, the same output is always provided. The reason is that our approach evaluates some sample inputs and outputs for these tools, and then, configures a deterministic black box which uses well-known techniques to generate the same results for the initially given sample inputs. This technique can be useful to analyze and categorize existing tools. Larger our knowledge of these tools larger the possibility to find errors or improve them.

2.1.4 *Compressing large ontology alignments.*

There are many repositories of ontology alignments available on the web. The problem when storing an ontology alignment is that it is necessary to store a lot of information which a) needs much disk space and b) is very difficult to reuse. The reason for the first fact is that it is necessary to store the mappings, information regarding to the initial ontologies, related overhead, and so on. Secondly,

knowledge contained in the alignment only can be reused when comparing the same correspondences. Storing only the function that was used to generate the alignment can save much disk space (only a function is stored), contains the same knowledge that the alignment (the alignment can be generated again using the function), and is very reusable (the function can be used in other scenarios).

3. RELATED WORKS

The importance of ontology matching is evidenced by the large number of related works that have been made. Unable to cite all these works, we reference the most important surveys in this field, [3, 5, 7, 11, 14, 16] where ontology matching methods and tools are described. There are several improvements like the possibility to match very large ontologies [10] or the capability to make predictions [15].

Many authors tend to categorize simple ontology matching algorithms in the groups defined by Ehrig [5], thus, they try to categorize basic matching algorithms in four categories corresponding to the ontology features to exploit, i.e. Linguistic Features, Structural Features, Constraint-based Features and Integration-Knowledge-based features.

On the other hand, we have not found works addressing the problem of the reverse ontology matching. However, the problem has been treated in adjacent fields such as data exchange. For example, Fagin et al [8] developed a framework for reverse data exchange that supports source instances that may contain nulls. This development required a careful reformulation of all the important notions, including the identity schema mapping, inverse, and maximum recovery. Like in our approach, operators originally introduced by Arenas et al.[1, 2], thus, the composition operator and the inverse operator have been recognized as two fundamental ones.

4. REVERSE ONTOLOGY MATCHING

It is possible to compute an equivalent reverse matching function for the alignments that have been created using several of techniques surveyed previously, either they have been combined in a parallel or in a sequential way. Algorithm combination means that algorithms are considered independently of each other, instead of algorithm composition which consists of using several algorithms in order to create a new one (hybrid algorithm). The way to obtain this equivalent reverse matching function requires four main steps that are going to be described now. It should be taken into account that, although engineering details are outside the aim of this work, these steps are susceptible to automation.

1. Choose the set of algorithms which are going to be used to obtain the equivalent matching function
2. Apply rules for computing the matching function and generating an intermediate equation
3. Obtain the equivalent matching function from the intermediate equation
4. Simplify (if possible) the equivalent reverse ontology matching function

On the other hand, if the mappings are given in a probabilistic form, firstly we have to decide a threshold in order to identify the valid ones. Then, we can proceed with point number one. A future improvement could consist of managing the uncertainty inherent in these mappings.

4.1 Choosing the set of algorithms

Choosing an appropriate set of algorithms to compose the equivalent matching function is very important in order to get success in the reverse matching process. Ideally, we need to use all existing matching algorithms, but this choice is not viable in practice, so we propose to choose, at least, a representative from each of the categories, although it is possible to choose hybrid algorithms too.

4.2 The equivalent matching function

In this step, it is necessary to apply several rules to know the algorithms that has been used originally to perform the alignment. There rules are:

Rule 1. *All algorithms which satisfy a mapping will be included in the intermediate equation.*

Rule 2. *All algorithms which satisfy a same input will be combined in a sequential way, thus, using the operator AND.*

Rule 3. *All algorithms, or set of algorithms, which satisfy two different inputs will be combined parallelly, thus, using the operator OR.*

After applying these rules, there might be mappings which cannot be found using the algorithms included in the previous step. For this reason it is necessary to define the concept of magic mapping.

Definition 5 (Magic mapping). *We define a magic mapping as the tuple (id, e, e', n, R) belonging to an alignment A so we do not know what algorithm was used to find it.*

The notion of magic mapping tell us that either we have not used an appropriate set of matching algorithms or the need to design a new matching algorithm that addresses this issue.

4.3 Extraction of the generalization pattern

This step consists of erasing duplicates expressions from the intermediate equation and making free linked variables. In this way, we obtain a clean of redundancies function.

4.4 Reduction Properties

We have borrowed several rules from the boolean algebra in order to reduce the length of the equivalent reverse ontology matching function. In fact, we have classified these rules in two different groups. For expressions with overlapped (set of) algorithms:

$$(a \wedge b) \wedge b \rightarrow (a \wedge b) \quad (5)$$

$$(a \wedge b) \vee b \rightarrow b \quad (6)$$

$$(a \vee b) \wedge b \rightarrow b \quad (7)$$

$$(a \vee b) \vee b \rightarrow (a \vee b) \quad (8)$$

For expressions without overlapped (set of) algorithms:

$$(a \wedge b) \wedge c \rightarrow (a \wedge b \wedge c) \quad (9)$$

$$(a \vee b) \wedge c \rightarrow (a \wedge c) \vee (b \wedge c) \quad (10)$$

4.5 Reverse Ontology Matching in practice

We show here two examples of how reverse ontology matching can be performed in the practice: (a) We extract the equivalent reverse matching function from a set of mappings (as example of capturing expert behavior), (b) we apply the obtained function that we have obtained to find mappings between two lemmaries (as example of matching-by-example).

Example 1. *Given the following set of mappings (Note the misspellings) $\{(Paris, Charles-de-Gaulle), (London, Heathrow), (Berlin, Schonenefeld), (Rome, Romans), (Madrid, Barajas), (Lisboa, Lisbon)\}$ compute the equivalent reverse matching function that has been used to generate them.*

1. We are going to choose this set of non-overlapped² ontology matching algorithms:

$\{Synonym, 3Grams, Stoilos, Wikip., Google\}$

2. If we follow the rules proposed, we are going to obtain the following intermediate equation:

$(Wikipedia (Paris, Charles-de-Gaulle) AND Google (Paris, Charles-de-Gaulle)) OR (Wikipedia (London, Heathrow) AND Google (London, Heathrow)) OR Google (Berlin, Schonenfeld) OR (3-Grams (Rome, Romans) AND Stoilos (Rome, Romans)) OR (Wikipedia (Madrid, Barajas) AND Google (Madrid, Barajas)) OR (3-Grams (Lisboa, Lisbon) AND Stoilos (Lisboa, Lisbon))$

3. Now, we obtain the generalization pattern:

$(Wikipedia (c1, c2) AND Google (c1, c2)) OR Google (c1, c2) OR (3-Grams (c1, c2) AND Stoilos (c1, c2))$

4. Finally, we apply the appropriate reduction properties ((6) in this case) in order to obtain the final equivalent matching function:

$Google (c1, c2) OR (3-Grams (c1, c2) AND Stoilos (c1, c2))$

What means that all input mappings that meet these conditions will be included in the final alignment. If more complex alignments are going to be analyzed, the two ontologies have to be accessible so that matching algorithms can detect structural similarities.

As it can be seen, we have captured the equivalent reverse matching function that was initially applied by the expert in order to match the concepts. If the function is applied to the input set of concepts, the mappings will be obtained again.

Example 2. Use the equivalent reverse matching function obtained in the Example 1 to match these two simple lemmaries $\{Canada, Asia, Boston, Mexico, New-York\}$ and $\{Celtics, Canadian, MexicoDF, Lakers, Manhattan\}$

²Two algorithms are overlapped if they aim to exploit the same ontology characteristics when looking for a correspondence

1. The equivalent matching function that we obtained in the Example 1 was:

$Google (c1, c2) OR (3-Grams (c1, c2) AND Stoilos (c1, c2))$

2. We generate the set of all possible correspondences between the two given lemmaries:

$\{(Canada, Celtics), (Canada, Canadian), (Canada, MexicoDF), (Canada, Lakers), (Canada, Manhattan), (Asia, Celtics), (Asia, Canadian), (Asia, MexicoDF), (Asia, Lakers), (Asia, Manhattan), (Boston, Celtics), (Boston, Canadian), (Boston, MexicoDF), (Boston, Lakers), (Boston, Manhattan), (Mexico, Celtics), (Mexico, Canadian), (Mexico, MexicoDF), (Mexico, Lakers), (Mexico, Manhattan), (New-York, Celtics), (New-York, Canadian), (New-York, MexicoDF), (New-York, Lakers), (New-York, Manhattan)\}$

3. We apply the equivalent matching function to the set of all correspondences and we have,

- *Google:* $(Boston, Celtics), (Boston, Lakers), (Mexico, MexicoDF), (New-York, Manhattan)$
- *3-Grams AND Stoilos:* $(Canada, Canadian), (Mexico, MexicoDF)$

4. The final set of mappings is

$\{(Canada, Canadian), (Boston, Celtics), (Boston, Lakers), (Mexico, Mexico DF), (New-York, Manhattan)\}$

5. We have that Boston belongs to two different mappings. This is because there are a lot of pages referring to NBA indexed by Google, so this algorithm generates a false positive. It is possible to implement the system in two ways: a) Allowing only 1:1 correspondences, in this case, only the mapping with a higher degree of confidence according to the algorithms will be added to the final alignment b) Allowing n:m correspondences, in this case, all mappings that meet the conditions will be included in the final results.

5. EVALUATION

We perform here an evaluation of our proposal. Firstly, we define that way to measure the quality of an equivalent matching function. Then, we describe and discuss the cases that we can find when evaluating this kind of functions, and lastly, we apply our technique in several real world scenarios in order to show that reverse ontology matching is viable in the practice.

Definition 6 (Equivalent reverse matching function evaluation). *An equivalent reverse matching function evaluation $ermfe$ is a function $ermfe : S \times S' \mapsto precision \in \mathbb{R} \in [0, 1] \times recall \in \mathbb{R} \in [0, 1]$ that associates an alignment S and a reference alignment S' to two real numbers stating the precision and recall of S in relation to S' .*

Precision states the fraction of retrieved mappings that were included in the original alignment S . Recall is the fraction of the correct mappings that are obtained successfully in comparison with the mappings belonging to S . In this way, precision is a measure of exactness and recall a measure of completeness. The problem here is that techniques can be optimized to obtain a high precision at the cost of the recall or, on the contrary, it is easy to optimize the recall at the cost of the precision. By this reason a F-measure is defined as a weighting factor between precision and recall. In this work, we use the most common configuration which consists of weighting precision and recall equally.

Let S the alignment provided initially, and let emf be the equivalent matching function obtained using reverse engineering and S' its output alignment. Then, we can face to these three cases:

- $S' = S$. We have a perfect equivalent matching function. The reason is that the equivalent matching function has been able to replicate exactly the results of the expert, technique or tool that created the original alignment.
- $S' \subset S$ ($S' = S - \mu$, where μ is the set of magic mappings). In this case, it has not been possible to find some of the algorithms used by the expert, technique or tool to generate some specific mappings. The final set of mappings provided by the equivalent matching function is a subset of the original set. The rest of mappings are magic mappings. A large number of magic mappings means that either we have not used an appropriate set of matching algorithms or that the alignment was generated using a hitherto unknown technique.

Ontology	#Map.	Pr.	Rc.	F-M.
Russia12	85	0.97	0.04	0.08
RussiaAB	117	1.00	0.07	0.13
TourismAB	226	0.96	0.12	0.21
Sports	150	0.99	0.02	0.04
AnimalsAB	24	0.92	0.14	0.24

Table 1: Quality for the Equivalent Reverse Ontology Matching functions using Web Knowledge algorithms

Ontology	#Map.	Pr.	Rc.	F-M.
Russia12	85	0.86	0.53	0.65
RussiaAB	117	1.00	0.51	0.68
TourismAB	226	0.93	0.47	0.62
Sports	150	0.94	0.39	0.55
AnimalsAB	24	0.75	0.80	0.77

Table 2: Quality for the Equivalent Reverse Ontology Matching functions not using Web Knowledge algorithms

- $S' \supset S$ ($S' = S + \lambda$, where λ is a set of new discovered mappings). In this case, the expert, technique or tool used an ambiguous strategy to create the final alignment. If the result provided by the equivalent matching function is a superset from the original one, then, we know that a strategy was used only for an arbitrary set of entities. For example, (Mexico, Mexican) was included in the final alignment but (Canada, Canadian) was not. Our technique is able to capture the strategy, but it cannot be applied in the same arbitrary way.

5.1 Empirical Results

In our experiment (see Table 1 and 2), we have noticed that algorithms which use Web Knowledge (Google and Wikipedia distance in this case) have a big impact in our results. The reason is that such kind of algorithms, which detect the co-occurrence of terms in the same websites of the Web, are able to find a lot of correspondences, and therefore the precision is increased when using them. But these algorithms generate a lot of false positives too, so the recall is decreased.

We have that values for the precision are good. This means that algorithms that we have used are able to capture the most of the mappings from the alignment. On the other hand, the value for the recall is lower, what means that these algorithms find more mappings than the alignment had originally. Therefore, F-measure, thus, the overall quality measure decreases.

6. CONCLUSIONS

We have presented a novel approach for reverse ontology matching. To the best of our knowledge, this approach is the first attempt to extract the functions that were originally used to create an alignment between ontologies. As we have shown, it is very difficult to obtain the original function, but it is possible to compute an equivalent reverse ontology matching function for all ontology matching functions that have been created using one or several of the techniques studied, either they have been combined in a parallel or in a sequential way.

Results show us that we have reached a reasonable quality when capturing the equivalent reverse matching functions in the most of cases. Moreover, we have introduced the notion of magic mapping as a way to deal with those mappings which we do not know how they were found. The notion of magic mapping tells us that either we have not used an appropriate set of matching algorithms or we need to design a new matching algorithm that addresses this issue. However, in practice, web knowledge algorithms limit the presence of magic mappings and have a great impact on the results. The reason is that such algorithms are able to find a lot of correspondences (even those that are not very frequent), and therefore precision is increased. But these algorithms generate a lot of false positives too, so the recall is decreased. For this reason, we propose to use web knowledge algorithms only in domains where a great precision is required.

As future work, we have to face some important challenges. Firstly, it is necessary to improve the recall of the equivalent reverse matching functions. We think that this can be achieved either by the incorporation of more efficient matching algorithms either the design of a new composition model more effective than the current model. Secondly, it is necessary to research faster ways to reverse engineer an alignment. Checking one by one the mappings is a time consuming strategy, so it is necessary to research ways to accelerate the process without loss of quality. One possible way to do this could be working only on a sample of mappings, for example.

Acknowledgements

We wish to thank to the anonymous reviewers for the suggestions which have helped to improve this work. We thank to Anne Doherty for proofreading this document. This work has been funded by the Spanish Ministry of Innovation & Science through the project TIN2008-04844 and by the Dept. of Innovation & Science from the Regional Government of Andalusia through the project P07-TIC-02978.

7. REFERENCES

- [1] Arenas, M, Perez, J., Riveros, C. (2008) The recovery of a schema mapping: bringing exchanged data back. *Proc. of PODS* 13–22.
- [2] Arenas, M, Perez, J., Reutter, C. (2009) Inverting Schema Mappings: Bridging the Gap between Theory and Practice. *PVLDB* 2(1) 1018–1029.
- [3] Choi, C., Song, I., Han, H. (2006) A Survey on Ontology Mapping. *ACM Sigmod Record* 35(3) 34–41.
- [4] Cilibrasi, R., Vitanyi, P. (2007) The Google Similarity Distance. *IEEE Trans. Knowl. Data Eng.* 19(3) 370–383.
- [5] Ehrig, M. (2006) *Ontology Alignment: Bridging the Semantic Gap*. Springer-Verlag.
- [6] Eckert, K., Meilicke, C., Stuckenschmidt, H. (2009) Improving Ontology Matching Using Meta-level Learning. *Proc. of European Semantic Web Conference ESWC 2009* 158–172.
- [7] Euzenat, J., Shvaiko, P. (2007) *Ontology Matching*. Springer-Verlag.
- [8] Fagin, R., Kolaitis, P.G., Popa, L., Tan, W.C. (2009) Reverse data exchange: coping with nulls. *Proc. of PODS* 23–32.
- [9] Gruber, T. (1993) A translation approach to portable ontology specifications. *Knowledge Acquisition* 5(2) 199–220.
- [10] Hu, W., Qu, Y., Cheng, G. (2008) Matching large ontologies: A divide-and-conquer approach. *Data Knowl. Eng.* 67(1) 140–160.
- [11] Kalfoglou, Y., Schorlemmer, M. (2003) Ontology mapping: the state of the art. *Knowledge Eng. Review* 18(1) 1–31.
- [12] Li, J., Tang, J., Li, Y., Luo, Q. (2009) RiMOM: A Dynamic Multistrategy Ontology Alignment Framework. *IEEE Trans. Knowl. Data Eng.* 21(8) 1218–1232.
- [13] Martinez-Gil, J., Aldana-Montes, J. (2011) Evaluation of two heuristic approaches to solve the ontology meta-matching problem *Knowl. Inf. Syst.* 26(2) 225–247.
- [14] Noy, N. (2004) Semantic Integration: A Survey Of Ontology-Based Approaches. *ACM Sigmod Record* 33(4) 65–70.
- [15] Pirró, G., Talia, D. (2010) UFOme: An ontology mapping system with strategy prediction capabilities. *Data Knowl. Eng.* 69(5) 444–471.
- [16] Rahm, E., Bernstein, P.A. (2001) A survey of approaches to automatic schema matching. *VLDB J.* 10(4) 334–350.

Scalable SQL and NoSQL Data Stores

Rick Cattell

Cattell.Net Software

Email: rick@cattell.net

ABSTRACT

In this paper, we examine a number of SQL and so-called “NoSQL” data stores designed to scale simple OLTP-style application loads over many servers. Originally motivated by Web 2.0 applications, these systems are designed to scale to thousands or millions of users doing updates as well as reads, in contrast to traditional DBMSs and data warehouses. We contrast the new systems on their data model, consistency mechanisms, storage mechanisms, durability guarantees, availability, query support, and other dimensions. These systems typically sacrifice some of these dimensions, e.g. database-wide transaction consistency, in order to achieve others, e.g. higher availability and scalability.

Note: Bibliographic references for systems are not listed, but URLs for more information can be found in the System References table at the end of this paper.

Caveat: Statements in this paper are based on sources and documentation that may not be reliable, and the systems described are “moving targets,” so some statements may be incorrect. Verify through other sources before depending on information here. Nevertheless, we hope this comprehensive survey is useful! Check for future corrections on the author’s web site cattell.net/datastores.

Disclosure: The author is on the technical advisory board of Schooner Technologies and has a consulting business advising on scalable databases.

1. OVERVIEW

In recent years a number of new systems have been designed to provide good horizontal scalability for simple read/write database operations distributed over many servers. In contrast, traditional database products have comparatively little or no ability to scale horizontally on these applications. This paper examines and compares the various new systems.

Many of the new systems are referred to as “NoSQL” data stores. The definition of NoSQL, which stands for “Not Only SQL” or “Not Relational”, is not entirely agreed upon. For the purposes of this paper, NoSQL systems generally have six key features:

1. the ability to horizontally scale “simple operation” throughput over many servers,
2. the ability to replicate and to distribute (partition) data over many servers,

3. a simple call level interface or protocol (in contrast to a SQL binding),
4. a weaker concurrency model than the ACID transactions of most relational (SQL) database systems,
5. efficient use of distributed indexes and RAM for data storage, and
6. the ability to dynamically add new attributes to data records.

The systems differ in other ways, and in this paper we contrast those differences. They range in functionality from the simplest distributed hashing, as supported by the popular memcached open source cache, to highly scalable partitioned tables, as supported by Google’s BigTable [1]. In fact, BigTable, memcached, and Amazon’s Dynamo [2] provided a “proof of concept” that inspired many of the data stores we describe here:

- Memcached demonstrated that in-memory indexes can be highly scalable, distributing and replicating objects over multiple nodes.
- Dynamo pioneered the idea of *eventual consistency* as a way to achieve higher availability and scalability: data fetched are not guaranteed to be up-to-date, but updates are guaranteed to be propagated to all nodes eventually.
- BigTable demonstrated that persistent record storage could be scaled to thousands of nodes, a feat that most of the other systems aspire to.

A key feature of NoSQL systems is “shared nothing” horizontal scaling – replicating and partitioning data over many servers. This allows them to support a large number of simple read/write operations per second. This simple operation load is traditionally called OLTP (online transaction processing), but it is also common in modern web applications

The NoSQL systems described here generally do not provide ACID transactional properties: updates are eventually propagated, but there are limited guarantees on the consistency of reads. Some authors suggest a “BASE” acronym in contrast to the “ACID” acronym:

- BASE = Basically Available, Soft state, Eventually consistent
- ACID = Atomicity, Consistency, Isolation, and Durability

The idea is that by giving up ACID constraints, one can achieve much higher performance and scalability.

However, the systems differ in how much they give up. For example, most of the systems call themselves “eventually consistent”, meaning that updates are eventually propagated to all nodes, but many of them provide mechanisms for some degree of consistency, such as multi-version concurrency control (MVCC).

Proponents of NoSQL often cite Eric Brewer’s CAP theorem [4], which states that a system can have only two out of three of the following properties: consistency, availability, and partition-tolerance. The NoSQL systems generally give up consistency. However, the trade-offs are complex, as we will see.

New relational DBMSs have also been introduced to provide better horizontal scaling for OLTP, when compared to traditional RDBMSs. After examining the NoSQL systems, we will look at these SQL systems and compare the strengths of the approaches. The SQL systems strive to provide horizontal scalability without abandoning SQL and ACID transactions. We will discuss the trade-offs here.

In this paper, we will refer to both the new SQL and NoSQL systems as *data stores*, since the term “database system” is widely used to refer to traditional DBMSs. However, we will still use the term “database” to refer to the stored data in these systems. All of the data stores have some administrative unit that you would call a database: data may be stored in one file, or in a directory, or via some other mechanism that defines the scope of data used by a group of applications. Each database is an island unto itself, even if the database is partitioned and distributed over multiple machines: there is no “federated database” concept in these systems (as with some relational and object-oriented databases), allowing multiple separately-administered databases to appear as one. Most of the systems allow horizontal partitioning of data, storing records on different servers according to some key; this is called “sharding”. Some of the systems also allow vertical partitioning, where parts of a single record are stored on different servers.

1.1 Scope of this Paper

Before proceeding, some clarification is needed in defining “horizontal scalability” and “simple operations”. These define the focus of this paper.

By “simple operations”, we refer to key lookups, reads and writes of one record or a small number of records. This is in contrast to complex queries or joins, read-mostly access, or other application loads. With the advent of the web, especially Web 2.0 sites where millions of users may both read and write data, scalability for simple database operations has become more important. For example, applications may search and update multi-server databases of electronic mail, personal profiles, web postings, wikis, customer

records, online dating records, classified ads, and many other kinds of data. These all generally fit the definition of “simple operation” applications: reading or writing a small number of related records in each operation.

The term “horizontal scalability” means the ability to distribute both the data and the load of these simple operations over many servers, with no RAM or disk shared among the servers. Horizontal scaling differs from “vertical” scaling, where a database system utilizes many cores and/or CPUs that share RAM and disks. Some of the systems we describe provide both vertical and horizontal scalability, and the effective use of multiple cores is important, but our main focus is on horizontal scalability, because the number of cores that can share memory is limited, and horizontal scaling generally proves less expensive, using commodity servers. Note that horizontal and vertical *partitioning* are not related to horizontal and vertical *scaling*, except that they are both useful for horizontal scaling.

1.2 Systems Beyond our Scope

Some authors have used a broad definition of NoSQL, including any database system that is not relational. Specifically, they include:

- *Graph database systems*: Neo4j and OrientDB provide efficient distributed storage and queries of a graph of nodes with references among them.
- *Object-oriented database systems*: Object-oriented DBMSs (e.g., Versant) also provide efficient distributed storage of a graph of objects, and materialize these objects as programming language objects.
- *Distributed object-oriented stores*: Very similar to object-oriented DBMSs, systems such as GemFire distribute object graphs in-memory on multiple servers.

These systems are a good choice for applications that must do fast and extensive reference-following, especially where data fits in memory. Programming language integration is also valuable. Unlike the NoSQL systems, these systems generally provide ACID transactions. Many of them provide horizontal scaling for reference-following and distributed query decomposition, as well. Due to space limitations, however, we have omitted these systems from our comparisons. The applications and the necessary optimizations for scaling for these systems differ from the systems we cover here, where key lookups and simple operations predominate over reference-following and complex object behavior. It is possible these systems can scale on simple operations as well, but that is a topic for a future paper, and proof through benchmarks.

Data warehousing database systems provide horizontal scaling, but are also beyond the scope of this paper. Data warehousing applications are different in important ways:

- They perform complex queries that collect and join information from many different tables.
- The ratio of reads to writes is high: that is, the database is read-only or read-mostly.

There are existing systems for data warehousing that scale well horizontally. Because the data is infrequently updated, it is possible to organize or replicate the database in ways that make scaling possible.

1.3 Data Model Terminology

Unlike relational (SQL) DBMSs, the terminology used by NoSQL data stores is often inconsistent. For the purposes of this paper, we need a consistent way to compare the data models and functionality.

All of the systems described here provide a way to store scalar values, like numbers and strings, as well as BLOBs. Some of them also provide a way to store more complex nested or reference values. The systems all store sets of attribute-value pairs, but use different data structures, specifically:

- A “tuple” is a row in a relational table, where attribute names are pre-defined in a schema, and the values must be scalar. The values are referenced by attribute name, as opposed to an array or list, where they are referenced by ordinal position.
- A “document” allows values to be nested documents or lists as well as scalar values, and the attribute names are dynamically defined for each document at runtime. A document differs from a tuple in that the attributes are not defined in a global schema, and this wider range of values are permitted.
- An “extensible record” is a hybrid between a tuple and a document, where families of attributes are defined in a schema, but new attributes can be added (within an attribute family) on a per-record basis. Attributes may be list-valued.
- An “object” is analogous to an object in programming languages, but without the procedural methods. Values may be references or nested objects.

1.4 Data Store Categories

In this paper, the data stores are grouped according to their data model:

- *Key-value Stores*: These systems store values and an index to find them, based on a programmer-defined key.

- *Document Stores*: These systems store documents, as just defined. The documents are indexed and a simple query mechanism is provided.
- *Extensible Record Stores*: These systems store extensible records that can be partitioned vertically and horizontally across nodes. Some papers call these “wide column stores”.
- *Relational Databases*: These systems store (and index and query) tuples. The new RDBMSs that provide horizontal scaling are covered in this paper.

Data stores in these four categories are covered in the next four sections, respectively. We will then summarize and compare the systems.

2. KEY-VALUE STORES

The simplest data stores use a data model similar to the popular memcached distributed in-memory cache, with a single key-value index for all the data. We’ll call these systems *key-value stores*. Unlike memcached, these systems generally provide a persistence mechanism and additional functionality as well: replication, versioning, locking, transactions, sorting, and/or other features. The client interface provides inserts, deletes, and index lookups. Like memcached, none of these systems offer secondary indices or keys.

2.1 Project Voldemort

Project Voldemort is an advanced key-value store, written in Java. It is open source, with substantial contributions from LinkedIn. Voldemort provides multi-version concurrency control (MVCC) for updates. It updates replicas asynchronously, so it does not guarantee consistent data. However, it can guarantee an up-to-date view if you read a majority of replicas.

Voldemort supports optimistic locking for consistent multi-record updates: if updates conflict with any other process, they can be backed out. Vector clocks, as used in Dynamo [3], provide an ordering on versions. You can also specify which version you want to update, for the put and delete operations.

Voldemort supports automatic sharding of data. *Consistent hashing* is used to distribute data around a ring of nodes: data hashed to node K is replicated on node K+1 ... K+n where n is the desired number of extra copies (often n=1). Using good sharding technique, there should be many more “virtual” nodes than physical nodes (servers). Once data partitioning is set up, its operation is transparent. Nodes can be added or removed from a database cluster, and the system adapts automatically. Voldemort automatically detects and recovers failed nodes.

Voldemort can store data in RAM, but it also permits plugging in a storage engine. In particular, it supports a Berkeley DB and Random Access File storage engine. Voldemort supports lists and records in addition to simple scalar values.

2.2 Riak

Riak is written in Erlang. It was open-sourced by Basho in mid-2009. Basho alternately describes Riak as a “key-value store” and “document store”. We will categorize it as an advanced key-value store here, because it lacks important features of document stores, but it (and Voldemort) have more functionality than the other key-value stores:

- Riak objects can be fetched and stored in JSON format, and thus can have multiple fields (like documents), and objects can be grouped into buckets, like the collections supported by document stores, with allowed/required fields defined on a per-bucket basis.
- Riak does not support indices on any fields except the primary key. The only thing you can do with the non-primary fields is fetch and store them as part of a JSON object. Riak lacks the query mechanisms of the document stores; the only lookup you can do is on primary key.

Riak supports replication of objects and sharding by hashing on the primary key. It allows replica values to be temporarily inconsistent. Consistency is tunable by specifying how many replicas (on different nodes) must respond for a successful read and how many must respond for a successful write. This is per-read and per-write, so different parts of an application can choose different trade-offs.

Like Voldemort, Riak uses a derivative of MVCC where vector clocks are assigned when values are updated. Vector clocks can be used to determine when objects are direct descendents of each other or a common parent, so Riak can often self-repair data that it discovers to be out of sync.

The Riak architecture is symmetric and simple. Like Voldemort, it uses consistent hashing. There is no distinguished node to track status of the system: the nodes use a gossip protocol to track who is alive and who has which data, and any node may service a client request. Riak also includes a map/reduce mechanism to split work over all the nodes in a cluster.

The client interface to Riak is based on RESTful HTTP requests. REST (REpresentational State Transfer) uses uniform, stateless, cacheable, client-server calls. There is also a programmatic interface for Erlang, Java, and other languages.

The storage part of Riak is “pluggable”: the key-value pairs may be in memory, in ETS tables, in DETS

tables, or in Osmos tables. ETS, DETS, and Osmos tables are all implemented in Erlang, with different performance and properties.

One unique feature of Riak is that it can store “links” between objects (documents), for example to link objects for authors to the objects for the books they wrote. Links reduce the need for secondary indices, but there is still no way to do range queries.

Here’s an example of a Riak object described in JSON:

```
{
  "bucket":"customers",
  "key":"12345",
  "object":{
    "name":"Mr. Smith",
    "phone":"415-555-6524" }
  "links":[
    ["sales","Mr. Salesguy","salesrep"],
    ["cust-orders","12345","orders"]]
  "vclock":"opaque-riak-vclock",
  "lastmod":"Mon, 03 Aug 2009 18:49:42 GMT"
}
```

Note that the primary key is distinguished, while other fields are part of an “object” portion. Also note that the bucket, vector clock, and modification date is specified as part of the object, and links to other objects are supported.

2.3 Redis

The Redis key-value data store started as a one-person project but now has multiple contributors as BSD-licensed open source. It is written in C.

A Redis server is accessed by a wire protocol implemented in various client libraries (which must be updated when the protocol changes). The client side does the distributed hashing over servers. The servers store data in RAM, but data can be copied to disk for backup or system shutdown. System shutdown may be needed to add more nodes.

Like the other key-value stores, Redis implements insert, delete and lookup operations. Like Voldemort, it allows lists and sets to be associated with a key, not just a blob or string. It also includes list and set operations.

Redis does atomic updates by locking, and does asynchronous replication. It is reported to support about 100K gets/sets per second on an 8-core server.

2.4 Scalaris

Scalaris is functionally similar to Redis. It was written in Erlang at the Zuse Institute in Berlin, and is open source. In distributing data over nodes, it allows key ranges to be assigned to nodes, rather than simply hashing to nodes. This means that a query on a range of values does not need to go to every node, and it also may allow better load balancing, depending on key distribution.

Like the other key-value stores, it supports insert, delete, and lookup. It does replication synchronously (copies must be updated before the operation is complete) so data is guaranteed to be consistent. Scalaris also supports transactions with ACID properties on multiple objects. Data is stored in memory, but replication and recovery from node failures provides durability of the updates. Nevertheless, a multi-node power failure would cause disastrous loss of data, and the virtual memory limit sets a maximum database size.

Scalaris reads and writes must go to a majority of the replicas before an operation completes. Scalaris uses a ring of nodes, an unusual distribution and replication strategy that requires $\log(N)$ hops to read/write a key-value pair.

2.5 Tokyo Cabinet

Tokyo Cabinet / Tokyo Tyrant was a sourceforge.net project, but is now licensed and maintained by FAL Labs. Tokyo Cabinet is the back-end server, Tokyo Tyrant is a client library for remote access. Both are written in C.

There are six different variations for the Tokyo Cabinet server: hash indexes in memory or on disk, B-trees in memory or on disk, fixed-size record tables, and variable-length record tables. The engines obviously differ in their performance characteristics, e.g. the fixed-length records allow quick lookups. There are slight variations on the API supported by these engines, but they all support common get/set/update operations. The documentation is a bit unclear, but they claim to support locking, ACID transactions, a binary array data type, and more complex update operations to atomically update a number or concatenate to a string. They support asynchronous replication with dual master or master/slave. Recovery of a failed node is manual, and there is no automatic sharding.

2.6 Memcached, Membrain, and Membase

The memcached open-source distributed in-memory indexing system has been enhanced by Schooner Tehnologies and Membase, to include features analogous to the other key-value stores: persistence, replication, high availability, dynamic growth, backup, and so on. Without persistence or replication, memcached does not really qualify as a “data store”. However, Membrain and Membase certainly do, and these systems are also compatible with existing memcached applications. This compatibility is an attractive feature, given that memcached is widely used; memcached users that require more advanced features can easily upgrade to Membase and Membrain.

The Membase system is open source, and is supported by the company Membase. Its most attractive feature is probably its ability to elastically add or remove servers in a running system, moving data and dynamically redirecting requests in the meantime. The elasticity in most of the other systems is not as convenient.

Membrain is licensed per server, and is supported by Schooner Technologies. Its most attractive feature is probably its excellent tuning for flash memory. The performance gains of flash memory will not be gained in other systems by treating flash as a faster hard disk; it is important that the system treat flash as a true “third tier”, different from RAM and disk. For example, many systems have substantial overhead in buffering and caching hard disk pages; this is unnecessary overhead with flash. The benchmark results on Schooner’s web site show many *times* better performance than a number of competitors, particularly when data overflows RAM.

2.7 Summary

All the key-value stores support insert, delete, and lookup operations. All of these systems provide scalability through key distribution over nodes.

Voldemort, Riak, Tokyo Cabinet, and enhanced memcached systems can store data in RAM or on disk, with storage add-ons. The others store data in RAM, and provide disk as backup, or rely on replication and recovery so that a backup is not needed.

Scalaris and enhanced memcached systems use synchronous replication, the rest use asynchronous.

Scalaris and Tokyo Cabinet implement transactions, while the others do not.

Voldemort and Riak use multi-version concurrency control (MVCC), the others use locks.

Membrain and Membase are built on the popular memcached system, adding persistence, replication, and other features. Backward compatibility with memcached give these products an advantage.

3. DOCUMENT STORES

As discussed in the first section, document stores support more complex data than the key-value stores. The term “document store” may be confusing: while these systems could store “documents” in the traditional sense (articles, Microsoft Word files, etc.), a document in these systems can be any kind of “pointerless object”, consistent with our definition in Section 1. Unlike the key-value stores, these systems generally support secondary indexes and multiple types of documents (objects) per database, and nested documents or lists. Like other NoSQL systems, the

document stores do not provide ACID transactional properties.

3.1 SimpleDB

SimpleDB is part of Amazon's proprietary cloud computing offering, along with their Elastic Compute Cloud (EC2) and their Simple Storage Service (S3) on which SimpleDB is based. SimpleDB has been around since 2007. As the name suggests, its model is simple: SimpleDB has Select, Delete, GetAttributes, and PutAttributes operations on documents. SimpleDB is simpler than other document stores, as it does not allow nested documents.

Like most of the systems we discuss, SimpleDB supports eventual consistency, not transactional consistency. Like most of the other systems, it does asynchronous replication.

Unlike key-value datastores, and like the other document stores, SimpleDB supports more than one grouping in one database: documents are put into domains, which support multiple indexes. You can enumerate domains and their metadata. Select operations are on one domain, and specify a conjunction of constraints on attributes, basically in the form:

```
select <attributes> from <domain> where  
  <list of attribute value constraints>
```

Different domains may be stored on different Amazon nodes.

Domain indexes are automatically updated when any document's attributes are modified. It is unclear from the documentation whether SimpleDB automatically selects which attributes to index, or if it indexes everything. In either case, the user has no choice, and the use of the indexes is automatic in SimpleDB query processing.

SimpleDB does not automatically partition data over servers. Some horizontal scaling can be achieved by reading any of the replicas, if you don't care about having the latest version. Writes do not scale, however, because they must go asynchronously to all copies of a domain. If customers want better scaling, they must do so manually by sharding themselves.

SimpleDB is a "pay as you go" proprietary solution from Amazon. There are currently built-in constraints, some of which are quite limiting: a 10 GB maximum domain size, a limit of 100 active domains, a 5 second limit on queries, and so on. Amazon doesn't license SimpleDB source or binary code to run on your own servers. SimpleDB does have the advantage of Amazon support and documentation.

3.2 CouchDB

CouchDB has been an Apache project since early 2008. It is written in Erlang.

A CouchDB "collection" of documents is similar to a SimpleDB domain, but the CouchDB data model is richer. Collections comprise the only schema in CouchDB, and secondary indexes must be explicitly created on fields in collections. A document has field values that can be scalar (text, numeric, or boolean) or compound (a document or list).

Queries are done with what CouchDB calls "views", which are defined with Javascript to specify field constraints. The indexes are B-trees, so the results of queries can be ordered or value ranges. Queries can be distributed in parallel over multiple nodes using a map-reduce mechanism. However, CouchDB's view mechanism puts more burden on programmers than a declarative query language.

Like SimpleDB, CouchDB achieves scalability through asynchronous replication, not through sharding. Reads can go to any server, if you don't care about having the latest values, and updates must be propagated to all the servers. However, a new project called CouchDB Lounge has been built to provide sharding on top of CouchDB, see:

<http://code.google.com/p/couchdb-lounge/>

Like SimpleDB, CouchDB does not guarantee consistency. Unlike SimpleDB, each client does see a self-consistent view of the database, with repeatable reads: CouchDB implements multi-version concurrency control on individual documents, with a Sequence ID that is automatically created for each version of a document. CouchDB will notify an application if someone else has updated the document since it was fetched. The application can then try to combine the updates, or can just retry its update and overwrite.

CouchDB also provides durability on system crash. All updates (documents and indexes) are flushed to disk on commit, by writing to the end of a file. (This means that periodic compaction is needed.) By default, it flushes to disk after every document update. Together with the MVCC mechanism, CouchDB's durability thus provides ACID semantics at the document level.

Clients call CouchDB through a RESTful interface. There are libraries for various languages (Java, C, PHP, Python, LISP, etc) that convert native API calls into the RESTful calls for you. CouchDB has some basic database administration functionality as well.

3.3 MongoDB

MongoDB is a GPL open source document store written in C++ and supported by 10gen. It has some similarities to CouchDB: it provides indexes on collections, it is lockless, and it provides a document query mechanism. However, there are important differences:

- MongoDB supports automatic sharding, distributing documents over servers.
- Replication in MongoDB is mostly used for failover, not for (dirty read) scalability as in CouchDB. MongoDB does not provide the global consistency of a traditional DBMS, but you can get local consistency on the up-to-date primary copy of a document.
- MongoDB supports dynamic queries with automatic use of indices, like RDBMSs. In CouchDB, data is indexed and searched by writing map-reduce views.
- CouchDB provides MVCC on documents, while MongoDB provides atomic operations on fields.

Atomic operations on fields are provided as follows:

- The update command supports “modifiers” that facilitate atomic changes to individual values: \$set sets a value, \$inc increments a value, \$push appends a value to an array, \$pushAll appends several values to an array, \$pull removes a value from an array, and \$pullAll removes several values from an array. Since these updates normally occur “in place”, they avoid the overhead of a return trip to the server.
- There is an “update if current” convention for changing a document only if field values match a given previous value.
- MongoDB supports a findAndModify command to perform an atomic update and immediately return the updated document. This is useful for implementing queues and other data structures requiring atomicity.

MongoDB indices are explicitly defined using an `ensureIndex` call, and any existing indices are automatically used for query processing. To find all products released last year costing under \$100 you could write:

```
db.products.find(
  {released: {$gte: new Date(2009, 1, 1)},
   price: {$lte: 100},})
```

If indexes are defined on the queried fields, MongoDB will automatically use them. MongoDB also supports map-reduce, which allows for complex aggregations across documents.

MongoDB stores data in a binary JSON-like format called BSON. BSON supports boolean, integer, float, date, string and binary types. Client drivers encode the local language’s document data structure (usually a dictionary or associative array) into BSON and send it over a socket connection to the MongoDB server (in contrast to CouchDB, which sends JSON as text over an HTTP REST interface). MongoDB also supports a GridFS specification for large binary objects, eg.

images and videos. These are stored in chunks that can be streamed back to the client for efficient delivery.

MongoDB supports master-slave replication with automatic failover and recovery. Replication (and recovery) is done at the level of shards. Collections are automatically sharded via a user-defined shard key. Replication is asynchronous for higher performance, so some updates may be lost on a crash.

3.4 Terrastore

Another recent document store is Terrastore, which is built on the Terracotta distributed Java VM clustering product. Like many of the other NoSQL systems, client access to Terrastore is built on HTTP operations to fetch and store data. Java and Python client APIs have also been implemented.

Terrastore automatically partitions data over server nodes, and can automatically redistribute data when servers are added or removed. Like MongoDB, it can perform queries based on a predicate, including range queries, and like CouchDB, it includes a map/reduce mechanism for more advanced selection and aggregation of data.

Like the other document databases, Terrastore is schema-less, and does not provide ACID transactions. Like MongoDB, it provides consistency on a per-document basis: a read will always fetch the latest version of a document.

Terrastore supports replication and failover to a hot standby.

3.5 Summary

The document stores are schema-less, except for attributes (which are simply a name, and are not pre-specified), collections (which are simply a grouping of documents), and the indexes defined on collections (explicitly defined, except with SimpleDB). There are some differences in their data models, e.g. SimpleDB does not allow nested documents.

The document stores are very similar but use different terminology. For example, a SimpleDB Domain = CouchDB Database = MongoDB Collection = Terrastore Bucket. SimpleDB calls documents “items”, and an attribute is a field in CouchDB, or a key in MongoDB or Terrastore.

Unlike the key-value stores, the document stores provide a mechanism to query collections based on multiple attribute value constraints. However, CouchDB does not support a non-procedural query language: it puts more work on the programmer and requires explicit utilization of indices.

The document stores generally do not provide explicit locks, and have weaker concurrency and atomicity properties than traditional ACID-compliant databases.

They differ in how much concurrency control they do provide.

Documents can be distributed over nodes in all of the systems, but scalability differs. All of the systems can achieve scalability by reading (potentially) out-of-date replicas. MongoDB and Terrastore can obtain scalability without that compromise, and can scale writes as well, through automatic sharding and atomic operations on documents. CouchDB might be able to achieve this write-scalability with the help of the new CouchDB Lounge code.

A last-minute addendum as this paper goes to press: the CouchDB and Membase companies have now merged, to form Couchbase. They plan to provide a “best of both” merge of their products, e.g. with CouchDB’s richer data model as well as the speed and elastic scalability of Membase. See Couchbase.com for more information.

4. EXTENSIBLE RECORD STORES

The extensible record stores seem to have been motivated by Google’s success with BigTable. Their basic data model is rows and columns, and their basic scalability model is splitting both rows and columns over multiple nodes:

- Rows are split across nodes through sharding on the primary key. They typically split by range rather than a hash function. This means that queries on ranges of values do not have to go to every node.
- Columns of a table are distributed over multiple nodes by using “column groups”. These may seem like a new complexity, but column groups are simply a way for the customer to indicate which columns are best stored together.

As noted earlier, these two partitionings (horizontal and vertical) can be used simultaneously on the same table. For example, if a customer table is partitioned into three column groups (say, separating the customer name/address from financial and login information), then each of the three column groups is treated as a separate table for the purposes of sharding the rows by customer ID: the column groups for one customer may or may not be on the same server.

The column groups must be pre-defined with the extensible record stores. However, that is not a big constraint, as new attributes can be defined at any time. Rows are analogous to documents: they can have a variable number of attributes (fields), the attribute names must be unique, rows are grouped into collections (tables), and an individual row’s attributes can be of any type. (However, note that CouchDB and MongoDB support nested objects, while the extensible record stores generally support only scalar types.)

Although most extensible record stores were patterned after BigTable, it appears that none of the extensible records stores come anywhere near to BigTable’s scalability at present. BigTable is used for many purposes (think of the many services Google provides, not just web search). It is worthwhile reading the BigTable paper [1] for background on the challenges with scaling.

4.1 HBase

HBase is an Apache project written in Java. It is patterned directly after BigTable:

- HBase uses the Hadoop distributed file system in place of the Google file system. It puts updates into memory and periodically writes them out to files on the disk.
- The updates go to the end of a data file, to avoid seeks. The files are periodically compacted. Updates also go to the end of a write ahead log, to perform recovery if a server crashes.
- Row operations are atomic, with row-level locking and transactions. There is optional support for transactions with wider scope. These use optimistic concurrency control, aborting if there is a conflict with other updates.
- Partitioning and distribution are transparent; there is no client-side hashing or fixed keyspace as in some NoSQL systems. There is multiple master support, to avoid a single point of failure. MapReduce support allows operations to be distributed efficiently.
- HBase’s B-trees allow fast range queries and sorting.
- There is a Java API, a Thrift API, and REST API. JDBC/ODBC support has recently been added.

The initial prototype of HBase released in February 2007. The support for transactions is attractive, and unusual for a NoSQL system.

4.2 HyperTable

HyperTable is written in C++. Its was open-sourced by Zvents. It doesn’t seem to have taken off in popularity yet, but Baidu became a project sponsor, that should help.

Hypertable is very similar to HBase and BigTable. It uses column families that can have any number of column “qualifiers”. It uses timestamps on data with MVCC. It requires an underlying distributed file system such as Hadoop, and a distributed lock manager. Tables are replicated and partitioned over servers by key ranges. Updates are done in memory and later flushed to disk.

Hypertable supports a number of programming language client interfaces. It uses a query language named HQL.

4.3 Cassandra

Cassandra is similar to the other extensible record stores in its data model and basic functionality. It has column groups, updates are cached in memory and then flushed to disk, and the disk representation is periodically compacted. It does partitioning and replication. Failure detection and recovery are fully automatic. However, Cassandra has a weaker concurrency model than some other systems: there is no locking mechanism, and replicas are updated asynchronously.

Like HBase, Cassandra is written in Java, and used under Apache licensing. It is supported by DataStax, and was originally open sourced by Facebook in 2008. It was designed by a Facebook engineer and a Dynamo engineer, and is described as a marriage of Dynamo and BigTable. Cassandra is used by Facebook as well as other companies, so the code is reasonably mature.

Client interfaces are created using Facebook's Thrift framework:

<http://incubator.apache.org/thrift/>

Cassandra automatically brings new available nodes into a cluster, uses the phi accrual algorithm to detect node failure, and determines cluster membership in a distributed fashion with a gossip-style algorithm.

Cassandra adds the concept of a "supercolumn" that provides another level of grouping within column groups. Databases (called keyspaces) contain column families. A column family contains either supercolumns or columns (not a mix of both). Supercolumns contain columns. As with the other systems, any row can have any combination of column values (i.e., rows are variable length and are not constrained by a table schema).

Cassandra uses an ordered hash index, which should give most of the benefit of both hash and B-tree indexes: you know which nodes could have a particular range of values instead of searching all nodes. However, sorting would still be slower than with B-trees.

Cassandra has reportedly scaled to about 150 machines in production at Facebook, perhaps more by now. Cassandra seems to be gaining a lot of momentum as an open source project, as well.

For applications where Cassandra's eventual-consistency model is not adequate, "quorum reads" of a majority of replicas provide a way to get the latest data. Cassandra writes are atomic within a column family. There is also some support for versioning and conflict resolution.

4.4 Other Systems

Yahoo's PNUTs system also belongs in the "extensible record store" category. However, it is not reviewed in this paper, as it is currently only used internally to Yahoo. We also have not reviewed BigTable, although its functionality is available indirectly through Google Apps. Both PNUTs and BigTable are included in the comparison table at the end of this paper.

4.5 Summary

The extensible record stores are mostly patterned after BigTable. They are all similar, but differ in concurrency mechanisms and other features.

Cassandra focuses on "weak" concurrency (via MVCC) and HBase and HyperTable on "strong" consistency (via locks and logging).

5. SCALABLE RELATIONAL SYSTEMS

Unlike the other data stores, relational DBMSs have a complete pre-defined schema, a SQL interface, and ACID transactions. Traditionally, RDBMSs have not achieved the scalability of some of the previously-described data stores. As of 5 years ago, MySQL Cluster appeared the most scalable, although not highly performant per node, compared to standard MySQL.

Recent developments are changing things. Further performance improvements have been made to MySQL Cluster, and several new products have come out, in particular VoltDB and Clustrix, that promise to have good per-node performance as well as scalability. It appears likely that some relational DBMSs will provide scalability comparable with NoSQL data stores, with two provisos:

- *Use small-scope operations:* As we've noted, operations that span many nodes, e.g. joins over many tables, will not scale well with sharding.
- *Use small-scope transactions:* Likewise, transactions that span many nodes are going to be very inefficient, with the communication and two-phase commit overhead.

Note that NoSQL systems avoid these two problems by making it difficult or impossible to perform larger-scope operations and transactions. In contrast, a scalable RDBMS does not need to *preclude* larger-scope operations and transactions: they simply penalize a customer for these operations *if they use them*. Scalable RDBMSs thus have an advantage over the NoSQL data stores, because you have the convenience of the higher-level SQL language and ACID properties, but you only pay a price for those

when they span nodes. Scalable RDBMSs are therefore included as a viable alternative in this paper.

5.1 MySQL Cluster

MySQL Cluster has been part of the MySQL release since 2004, and the code evolved from an even earlier project from Ericsson. MySQL Cluster works by replacing the InnoDB engine with a distributed layer called NDB. It is available from MySQL (now Oracle); it is not open source.

MySQL Cluster shards data over multiple database servers (a “shared nothing” architecture). Every shard is replicated, to support recovery. Bi-directional geographic replication is also supported.

MySQL Cluster supports in-memory as well as disk-based data. In-memory storage allows real-time responses.

Although MySQL Cluster seems to scale to more nodes than other RDBMSs to date, it reportedly runs into bottlenecks after a few dozen nodes. Work continues on MySQL Cluster, so this is likely to improve.

5.2 VoltDB

VoltDB is a new open-source RDBMS designed for high performance (per node) as well as scalability.

The scalability and availability features are competitive with MySQL Cluster and the NoSQL systems in this paper:

- Tables are partitioned over multiple servers, and clients can call any server. The distribution is transparent to SQL users, but the customer can choose the sharding attribute.
- Alternatively, selected tables can be *replicated* over servers, e.g. for fast access to read-mostly data.
- In any case, shards are replicated, so that data can be recovered in the event of a node crash. Database snapshots are also supported, continuous or scheduled.

Some features are still missing, e.g. online schema changes are currently limited, and asynchronous WAN replication and recovery are not yet implemented. However, VoltDB has some promising features that collectively may yield an order of magnitude advantage in single-node performance. VoltDB eliminates nearly all “waits” in SQL execution, allowing a very efficient implementation:

- The system is designed for a database that fits in (distributed) RAM on the servers, so that the system need never wait for the disk. Indexes and record structures are designed for RAM rather than disk, and the overhead of a disk cache/buffer is eliminated as well. Performance will be very

poor if virtual memory overflows RAM, but the gain with good RAM capacity planning is substantial.

- SQL execution is single-threaded for each shard, using a shared-nothing architecture, so there is no overhead for multi-thread latching.
- All SQL calls are made through stored procedures, with each stored procedure being one transaction. This means, if data is sharded to allow transactions to be executed on a single node, then no locks are required, and therefore no waits on locks. Transaction coordination is likewise avoided.
- Stored procedures are compiled to produce code comparable to the access level calls of NoSQL systems. They can be executed in the same order on a node and on replica node(s).

VoltDB argues that these optimizations greatly reduce the number of nodes needed to support a given application load, with modest constraints on the database design. They have already reported some impressive benchmark results on their web site. Of course, the highest performance requires that the database working set fits in distributed RAM, perhaps extended by SSDs. See [5] for some debate of the architectural issues on VoltDB and similar systems.

5.3 Clustrix

Clustrix offers a product with similarities to VoltDB and MySQL Cluster, but Clustrix nodes are sold as rack-mounted appliances. They claim scalability to hundreds of nodes, with automatic sharding and replication (with a 4:1 read/write ratio, they report 350K TPS on 20 nodes and 160M rows). Failover is automatic, and failed node recover is automatic. They also use solid state disks for additional performance (like the Schooner MySQL and NoSQL appliances).

As with the other relational products, Clustrix supports SQL with fully-ACID transactions. Data distribution and load balancing is transparent to the application programmer. Interestingly, they also designed their system to be seamlessly compatible with MySQL, supporting existing MySQL applications and front-end connectors. This could give them a big advantage in gaining adoption of proprietary hardware.

5.4 ScaleDB

ScaleDB is a new derivative of MySQL underway. Like MySQL Cluster, it replaces the InnoDB engine, and uses clustering of multiple servers to achieve scalability. ScaleDB differs in that it requires disks shared across nodes. Every server must have access to every disk. This architecture has not scaled very well for Oracle RAC, however.

ScaleDB's sharding is automatic: more servers can be added at any time. Server failure handling is also automatic. ScaleDB redistributes the load over existing servers.

ScaleDB supports ACID transactions and row-level locking. It has multi-table indexing (which is possible due to the shared disk).

5.5 ScaleBase

ScaleBase takes a novel approach, seeking to achieve the horizontal scaling with a layer entirely on top of MySQL, instead of modifying MySQL. ScaleBase includes a partial SQL parser and optimizer that shards tables over multiple single-node MySQL databases. Limited information is available about this new system at the time of this writing, however. It is currently a beta release of a commercial product, not open source.

Implementing sharding as a layer on top of MySQL introduces a problem, as transactions do not span MySQL databases. ScaleBase provides an option for distributed transaction coordination, but the higher-performance option provides ACID transactions only within a single shard/server.

5.6 NimbusDB

NimbusDB is another new relational system. It uses MVCC and distributed object based storage. SQL is the access language, with a row-oriented query optimizer and AVL tree indexes.

MVCC provides transaction isolation without the need for locks, allowing large scale parallel processing. Data is horizontally segmented row-by-row into distributed objects, allowing multi-site, dynamic distribution.

5.7 Other Systems

Google has recently created a layer on BigTable called Megastore. Megastore adds functionality that brings BigTable closer to a (scalable) relational DBMS in many ways: transactions that span nodes, a database schema defined in a SQL-like language, and hierarchical paths that allow some limited join capability. Google has also implemented a SQL processor that works on BigTable. There are still a lot of differences between Megastore / BigTable "NoSQL" and scalable relational systems, but the gap seems to be narrowing.

Microsoft's Azure product has some replication capabilities, but it does not directly support scaling through sharding. It allows tables to be stored "in the cloud" and can sync multiple databases. It supports SQL. We have not covered it in this paper.

The major RDBMSs (DB2, Oracle, SQL Server) also include some horizontal scaling features, either shared-nothing, or shared-disk.

5.8 Summary

MySQL Cluster uses a "shared nothing" architecture for scalability, as with most of the other solutions in this section, and it is the most mature solution here.

VoltDB looks promising because of its horizontal scaling as well as a bottom-up redesign to provide very high per-node performance. Clustrix looks promising as well, and supports solid state disks, but it is based on proprietary software and hardware.

Limited information is available about ScaleDB, NimbusDB, and ScaleBase at this point; they are at an early stage.

In theory, RDBMSs should be able to deliver scalability as long as applications avoid cross-node operations. If this proves true in practice, the simplicity of SQL and ACID transactions would give them an advantage over NoSQL for most applications.

6. USE CASES

No one of these data stores is best for all uses. A user's prioritization of features will be different depending on the application, as will the type of scalability required. A complete guide to choosing a data store is beyond the scope of this paper, but in this section we look at some examples of applications that fit well with the different data store categories.

6.1 Key-value Store Example

Key-value stores are generally good solutions if you have a simple application with only one kind of object, and you only need to look up objects up based on one attribute. The simple functionality of key-value stores may make them the simplest to use, especially if you're already familiar with memcached.

As an example, suppose you have a web application that does many RDBMS queries to create a tailored page when a user logs in. Suppose it takes several seconds to execute those queries, and the user's data is rarely changed, or you know when it changes because updates go through the same interface. Then you might want to store the user's tailored page as a single object in a key-value store, represented in a manner that's efficient to send in response to browser requests, and index these objects by user ID. If you store these objects persistently, then you may be able to avoid many RDBMS queries, reconstructing the objects only when a user's data is updated.

Even in the case of an application like Facebook, where a user's home page changes based on updates made by the user as well as updates made by others, it may be possible to execute RDBMS queries just once when the user logs in, and for the rest of that session show only the changes made by that user (not by other

users). Then, a simple key-value store could still be used as a relational database cache.

You could use key-value stores to do lookups based on multiple attributes, by creating additional key-value indexes that you maintain yourself. However, at that point you probably want to move to a document store.

6.2 Document Store Example

A good example application for a document store would be one with multiple different kinds of objects (say, in a Department of Motor Vehicles application, with vehicles and drivers), where you need to look up objects based on multiple fields (say, a driver's name, license number, owned vehicle, or birth date).

An important factor to consider is what level of concurrency guarantees you need. If you can tolerate an "eventually consistent" model with limited atomicity and isolation, the document stores should work well for you. That might be the case in the DMV application, e.g. you don't need to know if the driver has new traffic violations in the past minute, and it would be quite unlikely for two DMV offices to be updating the same driver's record at the same time. But if you require that data be up-to-date and atomically consistent, e.g. if you want to lock out logins after three incorrect attempts, then you need to consider other alternatives, or use a mechanism such as quorum-read to get the latest data.

6.3 Extensible Record Store Example

The use cases for extensible record stores are similar to those for document stores: multiple kinds of objects, with lookups based on any field. However, the extensible record store projects are generally aimed at higher throughput, and may provide stronger concurrency guarantees, at the cost of slightly more complexity than the document stores.

Suppose you are storing customer information for an eBay-style application, and you want to partition your data both horizontally and vertically:

- You might want to cluster customers by country, so that you can efficiently search all of the customers in one country.
- You might want to separate the rarely-changed "core" customer information such as customer addresses and email addresses in one place, and put certain frequently-updated customer information (such as current bids in progress) in a different place, to improve performance.

Although you could do this kind of horizontal/vertical partitioning yourself on top of a document store by creating multiple collections for multiple dimensions, the partitioning is most easily achieved with an extensible record store like HBase or HyperTable.

6.4 Scalable RDBMS Example

The advantages of relational DBMSs are well-known:

- If your application requires many tables with different types of data, a relational schema definition centralizes and simplifies your data definition, and SQL greatly simplifies the expression of operations that span tables.
- Many programmers are already familiar with SQL, and many would argue that the use of SQL is simpler than the lower-level commands provided by NoSQL systems.
- Transactions greatly simplify coding concurrent access. ACID semantics free the developer from dealing with locks, out-of-date data, update collisions, and consistency.
- Many more tools are currently available for relational DBMSs, for report generation, forms, and so on.

As a good example for relational, imagine a more complex DMV application, perhaps with a query interface for law enforcement that can interactively search on vehicle color, make, model, year, partial license plate numbers, and/or constraints on the owner such as the county of residence, hair color, and sex. ACID transactions could also prove valuable for a database being updated from many locations, and the aforementioned tools would be valuable as well. The definition of a common relational schema and administration tools can also be invaluable on a project with many programmers.

These advantages are dependent, of course, on a relational DBMS scaling to meet your application needs. Recently-reported benchmarks on VoltDB, Clustrix, and the latest version of MySQL Cluster suggest that scalability of relational DBMSs is greatly improving. Again, this assumes that your application does not demand updates or joins that span many nodes; the transaction coordination and data movement for that would be prohibitive. However, the NoSQL systems generally do not offer the possibility of transactions or query joins across nodes, so you are no worse off there.

7. CONCLUSIONS

We have covered over twenty scalable data stores in this paper. Almost all of them are moving targets, with limited documentation that is sometimes conflicting, so this paper is likely out-of-date if not already inaccurate at the time of this writing. However, we will attempt a snapshot summary, comparison, and predictions in this section. Consider this a starting point for further study.

7.1 Some Predictions

Here are some predictions of what will happen with the systems we've discussed, over the next few years:

- Many developers will be willing to abandon globally-ACID transactions in order to gain scalability, availability, and other advantages. The popularity of NoSQL systems has already demonstrated this. Customers tolerate airline over-booking, and orders that are rejected when items in an online shopping cart are sold out before the order is finalized. The world is not globally consistent.
- NoSQL data stores will not be a "passing fad". The simplicity, flexibility, and scalability of these systems fills a market niche, e.g. for web sites with millions of read/write users and relatively simple data schemas. Even with improved relational scalability, NoSQL systems maintain advantages for some applications.
- New relational DBMSs will also take a significant share of the scalable data storage market. If transactions and queries are generally limited to single nodes, these systems should be able to scale [5]. Where the desire for SQL or ACID transactions are important, these systems will be the preferred choice.
- Many of the scalable data stores will not prove "enterprise ready" for a while. Even though they fulfill a need, these systems are new and have not yet achieved the robustness, functionality, and maturity of database products that have been around for a decade or more. Early adopters have already seen web site outages with scalable data store failures, and many large sites continue to "roll their own" solution by sharding with existing RDBMS products. However, some of these new systems will mature quickly, given the great deal of energy directed at them.
- There will be major consolidation among the systems we've described. One or two systems will likely become the leaders in each of the categories. It seems unlikely that the market and open source community will be able to support the sheer number of products and projects we've studied here. Venture capital and support from key players will likely be a factor in this consolidation. For example, among the document stores, MongoDB has received substantial investment this year.

7.2 SQL vs NoSQL

SQL (relational) versus NoSQL scalability is a controversial topic. This paper argues against both extremes. Here is some more background to support this position.

The argument for relational over NoSQL goes something like this:

- If new relational systems can do everything that a NoSQL system can, with analogous performance and scalability, and with the convenience of transactions and SQL, why would you choose a NoSQL system?
- Relational DBMSs have taken and retained majority market share over other competitors in the past 30 years: network, object, and XML DBMSs.
- Successful relational DBMSs have been built to handle other specific application loads in the past: read-only or read-mostly data warehousing, OLTP on multi-core multi-disk CPUs, in-memory databases, distributed databases, and now horizontally scaled databases.
- While we don't see "one size fits all" in the SQL products themselves, we do see a common interface with SQL, transactions, and relational schema that give advantages in training, continuity, and data interchange.

The counter-argument for NoSQL goes something like this:

- We haven't yet seen good benchmarks showing that RDBMSs can achieve scaling comparable with NoSQL systems like Google's BigTable.
- If you only require a lookup of objects based on a single key, then a key-value store is adequate and probably easier to understand than a relational DBMS. Likewise for a document store on a simple application: you only pay the learning curve for the level of complexity you require.
- Some applications require a flexible schema, allowing each object in a collection to have different attributes. While some RDBMSs allow efficient "packing" of tuples with missing attributes, and some allow adding new attributes at runtime, this is uncommon.
- A relational DBMS makes "expensive" (multi-node multi-table) operations "too easy". NoSQL systems make them impossible or obviously expensive for programmers.
- While RDBMSs have maintained majority market share over the years, other products have established smaller but non-trivial markets in areas where there is a need for particular capabilities, e.g. indexed objects with products like BerkeleyDB, or graph-following operations with object-oriented DBMSs.

Both sides of this argument have merit.

7.3 Benchmarking

Given that scalability is the focus of this paper and of the systems we discuss, there is a “gaping hole” in our analysis: there is a scarcity of benchmarks to substantiate the many claims made for scalability. As we have noted, there are benchmark results reported on some of the systems, but almost none of the benchmarks are run on more than one system, and the results are generally reported by proponents of that one system, so there is always some question about their objectivity.

In this paper, we’ve tried to make the best comparisons possible based on architectural arguments alone. However, it would be highly desirable to get some useful objective data comparing the architectures:

- The trade-offs between the architectures are unclear. Are the bottlenecks in disk access, network communication, index operations, locking, or other components?
- Many people would like to see support or refutation of the argument that new relational systems can scale as well as NoSQL systems.
- A number of the systems are new, and may not live up to scalability claims without years of tuning. They also may be buggy. Which are truly mature?
- Which systems perform best on which loads? Are open source projects able to produce highly performant systems?

Perhaps the best benchmark to date is from Yahoo! Research [2], comparing PNUTS, HBASE, Cassandra, and sharded MySQL. Their benchmark, YCSB, is designed to be representative of web applications, and the code is available to others. Tier 1 of the benchmark measures raw performance, showing latency characteristics as the server load increases. Tier 2 measures scaling, showing how the benchmarked system scales as additional servers are added, and how quickly the system adapts to additional servers.

In this paper, I’d like to make a “call for scalability benchmarks,” suggesting YCSB as a good basis for the comparison. Even if the YCSB benchmark is run by different groups who may not duplicate the same hardware Yahoo specified, the results will be informative.

7.4 Some Comparisons

Given the quickly-changing landscape, this paper will not attempt to argue the merits of particular systems, beyond the comments already made. However, a comparison of the salient features may prove useful, so we finish with some comparisons.

Table 1 below compares the concurrency control, data storage medium, replication, and transaction mechanisms of the systems. These are difficult to summarize in a short table entry without oversimplifying, but we compare as follows.

For concurrency:

- Locks: some systems provide a mechanism to allow only one user at a time to read or modify an entity (an object, document, or row). In the case of MongoDB, a locking mechanism is provided at a field level.
- MVCC: some systems provide multi-version concurrency control, guaranteeing a read-consistent view of the database, but resulting in multiple conflicting versions of an entity if multiple users modify it at the same time.
- None: some systems do not provide atomicity, allowing different users to modify different parts of the same object in parallel, and giving no guarantee as to which version of data you will get when you read.
- ACID: the relational systems provide ACID transactions. Some of the more recent systems do this with no deadlocks and no waits on locks, by pre-analyzing transactions to avoid conflicts.

For data storage, some systems are designed for storage in RAM, perhaps with snapshots or replication to disk, while others are designed for disk storage, perhaps caching in RAM. RAM-based systems typically allow use of the operating system’s virtual memory, but performance appears to be very poor when they overflow physical RAM. A few systems have a pluggable back end allowing different data storage media, or they require a standardized underlying file system.

Replication can insure that mirror copies are always in sync (that is, they are updated lock-step and an operation is not completed until both replicas are modified). Alternatively, the mirror copy may be updated asynchronously in the background. Asynchronous replication allows faster operation, particular for remote replicas, but some updates may be lost on a crash. Some systems update local copies synchronously and geographically remote copies asynchronously (this is probably the only practical solution for remote data).

Transactions are supported in some systems, and not in others. Some NoSQL systems provide something in between, where “Local” transactions are supported only within a single object or shard.

Table 1 compares the systems on these four dimensions.

Table 1. System Comparison (grouped by category)

System	Conc Control	Data Storage	Repl-ication	Tx
Redis	Locks	RAM	Async	N
Scalaris	Locks	RAM	Sync	L
Tokyo	Locks	RAM or disk	Async	L
Voldemort	MVCC	RAM or BDB	Async	N
Riak	MVCC	Plug-in	Async	N
Membrain	Locks	Flash + Disk	Sync	L
Membase	Locks	Disk	Sync	L
Dynamo	MVCC	Plug-in	Async	N
SimpleDB	None	S3	Async	N
MongoDB	Locks	Disk	Async	N
Couch DB	MVCC	Disk	Async	N
Terrastore	Locks	RAM+	Sync	L
HBase	Locks	Hadoop	Async	L
HyperTable	Locks	Files	Sync	L
Cassandra	MVCC	Disk	Async	L
BigTable	Locks+s tamps	GFS	Sync+ Async	L
PNUTs	MVCC	Disk	Async	L
MySQL Cluster	ACID	Disk	Sync	Y
VoltDB	ACID, no lock	RAM	Sync	Y
Clustrix	ACID, no lock	Disk	Sync	Y
ScaleDB	ACID	Disk	Sync	Y
ScaleBase	ACID	Disk	Async	Y
NimbusDB	ACID, no lock	Disk	Sync	Y

Another factor to consider, but impossible to quantify objectively in a table, is code maturity. As noted earlier, many of the systems we discussed are only a couple of years old, and are likely to be unreliable. For this reason, existing database products are often a better choice if they can scale for your application's needs.

Probably the most important factor to consider is actual performance and scalability, as noted in the discussion of benchmarking. Benchmark references will be added to the author's website cattell.net/datastores as they become available.

Updates and corrections to this paper will be posted there as well. The landscape for scalable data stores is likely to change significantly over the next two years!

8. ACKNOWLEDGMENTS

I'd like to thank Len Shapiro, Jonathan Ellis, Dan DeMaggio, Kyle Banker, John Busch, Darpan Dinker, David Van Couvering, Peter Zaitsev, Steve Yen, and Scott Jarr for their input on earlier drafts of this paper. Any errors are my own, however! I'd also like to thank Schooner Technologies for their support on this paper.

9. REFERENCES

- [1] F. Chang et al, "BigTable: A Distributed Storage System for Structured Data", *Seventh Symposium on Operating System Design and Implementation*, November 2006.
- [2] B. Cooper et al, "Benchmarking Cloud Serving Systems with YCSB", *ACM Symposium on Cloud Computing (SoCC)*, Indianapolis, Indiana, June 2010.
- [3] B. DeCandia et al, "Dynamo: Amazon's Highly Available Key-Value Store", *Proceedings 21st ACM SIGOPS Symposium on Operating Systems Principles*, 2007.
- [4] S. Gilbert and N. Lynch, "Brewer's conjecture and the feasibility of consistent, available, and partition-tolerant web services", *ACM SIGACT News* 33, 2, pp 51-59, March 2002.
- [5] M. Stonebraker and R. Cattell, "Ten Rules for Scalable Performance in Simple Operation Datastores", *Communications of the ACM*, June 2011.

10. SYSTEM REFERENCES

The following table provides web information sources for all of the DBMSs and data stores covered in the paper, even those peripherally mentioned, alphabetized by system name. The table also lists the licensing model (proprietary, Apache, BSD, GPL), which may be important depending on your application.

System	License	Web site for more information
Berkeley DB	BSD	oss.oracle.com/berkeley-db.html
BigTable	Prop	labs.google.com/papers/bigtable.html
Cassandra	Apache	incubator.apache.org/cassandra
Clustrix	Prop	clustrix.com
CouchDB	Apache	couchdb.apache.org
Dynamo	Internal	portal.acm.org/citation.cfm?id=1294281

GemFire	Prop	gemstone.com/products/gemfire
HBase	Apache	hbase.apache.org
HyperTable	GPL	hypertable.org
Membase	Apache	membase.com
Membrain	Prop	schoonerinfotech.com/products/
Memcached	BSD	memcached.org
MongoDB	GPL	mongodb.org
MySQL Cluster	GPL	mysql.com/cluster
NimbusDB	Prop	nimbusdb.com
Neo4j	AGPL	neo4j.org
OrientDB	Apache	orienttechnologies.com

PNUTs	Internal	research.yahoo.com/node/2304
Redis	BSD	code.google.com/p/redis
Riak	Apache	riak.basho.com
Scalaris	Apache	code.google.com/p/scalaris
ScaleBase	Prop	scalebase.com
ScaleDB	GPL	scaledb.com
SimpleDB	Prop	amazon.com/simpledb
Terrastore	Apache	code.google.com/terrastore
Tokyo	GPL	tokyocabinet.sourceforge.net
Versant	Prop	versant.com
Voldemort	None	project-voldemort.com
VoltDB	GPL	voltdb.com

Christopher Ré Speaks Out on His ACM SIGMOD Jim Gray Dissertation Award, What He Wishes He Had Known As a Graduate Student, and More

by Marianne Winslett



Christopher Ré

<http://pages.cs.wisc.edu/~chrisre/>

Dear readers, you may have noticed that the printed versions of these interviews are not appearing in every issue now. I took a position as the director of UIUC's new research center in Singapore (<http://adsc.illinois.edu>), and that has reduced the amount of time I have available for other activities. Yet the video versions of the interviews still pop up on the SIGMOD web site regularly – so what is the problem?

The difficulty is that spoken English is almost a different language from written English, so the conversion of an interview transcript into a coherent document is a labor-intensive process. I am looking for a second editor who can take over this process, similar to the video editor who produces the video versions.

The new editor(s) of the written versions needs to be really good with English, because the written version must make it perfectly clear what the interviewee meant, while at the same time not changing the interviewee's unique verbal style. And there are delicate judgment calls to make, e.g., if the interviewee sounds extremely Italian, that should still show up in some way in the written version. To understand the task, try listening to a video version of an interview while reading the printed version. The reward for taking on this new position is to have your name in the byline above. So, if you are interested, please send email to winslett@illinois.edu !

Welcome to this installment of ACM SIGMOD Record's series of interviews with distinguished members of the database community. I'm Marianne Winslett, and today we are at the SIGMOD 2010 conference in Indianapolis. I have here with me Chris Ré, who is an assistant professor of computer science at the University of Wisconsin - Madison. Chris is the recipient of the 2010 ACM SIGMOD Jim Gray Dissertation Award, for his dissertation entitled Managing Large-Scale Probabilistic Databases. His PhD is from the University of Washington. So, Chris, welcome!

Our runners-up for the award this year were Soumyadeb Mitra (University of Illinois at Urbana-Champaign) and Fabian Suchanek (Max Planck Institute for Informatics).

So, Chris, what is the thesis of your thesis?

The thesis of my thesis is that it is possible to build large-scale probabilistic databases. The reason you would want to build such a crazy thing is that there are a lot of applications out there that are managing data that is increasingly large and increasingly imprecise.

Can you give me some examples?

Certainly. One example is RFID data, where you are trying to track objects and people as they move through space. Whenever you go to measure the physical world, you have issues like measurement error, so the resulting data is much less precise than the data we are used to having in traditional relational databases. And there are many other examples, such as information extraction, or data integration, where the data is somehow just less precise than what we are used to putting inside a database.

What are the research challenges in managing this type of data?

There are many challenges, and I was fortunate [to win this prize], because there are a lot of people working in this area. The main challenges I focused on were scalability and performance. In general, the way we model all these imprecise information sources is to use probability theory. Probability theory is great, because it lets you talk about many different kinds of imprecision in one unified way. But then, when you go to actually process a query, you have to consider all these different alternatives. And when you consider all these different alternatives, you are doing a lot more work than you would do in a traditional relational SQL style query processor. So I focused on performance and scale, and that is really the contribution of my dissertation.

It sounds rather impossible! What is the secret to be able to manage all that?

The secret sauce, I think, is a lot of other smart people's work. We took techniques from all over the place: we borrowed techniques from AI literature, machine learning literature, logic, probability. I had great co-authors around me, and we mixed all that stuff together. So some of the contributions that were in my dissertation were just observations, such as that in a lot of these applications, people are only interested in, say, the top five or ten answers. So when we are building a system that can handle imprecise data, we don't have to compute very many answers to each query. Once you have the idea that people want to focus on just the most likely answers, then you can zoom right in on that problem. We had some query processing techniques that allowed us to address that problem. The other thing we had in our bag of tricks was classical database tools, like materialized views, and some ideas from approximate query processing, popular old topics that we could bring to bear on probabilistic processing. So, the secret sauce was really other people's work. Ha!

Okay, so you claim! Although, since you got the award, I'm sure there are a lot of new results in your dissertation. In fact, the AI community and the logic and probability communities are not known for their attention to scalability.

That is one thing that I was really excited about in my work. I was applying almost a traditional database focus to these classical problems from these other fields. It was really a joy to pursue, and I hope that some of the people in those fields are excited by this work too.

What do you know now that you wish you had known as a graduate student?

The one thing I wish I had known was how much of the learning during the PhD was by observation. I had a great model when I was doing my PhD: Dan Suciu, my advisor. Just watching how he did things, and imitating him, was a great part of the process. But now that I have gone through the PhD process, I see that I didn't realize how many things I should imitate from the other people inside the field. Everything from how they communicate, to how they write their papers, to how they conduct their experiments, down to really fine details. There is really no way you can tell someone all those things, they just have to pick them up. I didn't realize how much of that was on the graduate student agenda until I got into my second or third year, and then all of a sudden I understood.

How did you get into the database field?

That is an interesting story. As an undergrad, I was a math and CS major. And then, on a lark with a friend, we decided to take the database course at Cornell, which is CS432. Jay Shanmugasundaram and Johannes Gehrke were team-teaching it that semester. And they had a guest lecturer, Jim Gray, and he was talking about the World Wide Telescope project. I remember being blown away by the simple answers he was giving to complex problems. As a mathematically minded undergrad, this was just amazing to me. I really fell in love with the topic as a result of that talk, so to me it is a real honor that his name is on this award.

I have a kind of similar story, only the teacher was Phil Bernstein. I was also a math major, and Phil made databases seem so interesting to me as an undergrad.

It just somehow came alive. They had the simple answers to complex problems that just blew me away.

If you could change one thing about yourself as a computer science researcher, what would it be?

Probably I would be a bit more disciplined. I tend to get very interested in a lot of different topics, which has been nice in some ways. For example, I have a thesis that touches a couple of different areas that I really like, and I am proud of that thesis. But I get distracted by topics pretty easily. There are a lot of fascinating topics in data management and if I could stick and hold on to one, I would be pretty happy about that.

Thank you for talking with us today!

Thank you so much!

Paper Bricks: an Alternative to Complete-Story Peer Reviewing¹

Jens Dittrich
Saarland University

ABSTRACT

The peer review system as used in several computer science communities has several flaws including long review times, overloaded reviewers, as well as fostering of niche topics. These flaws decrease quality, lower impact, slowdown the innovation process, and lead to frustration of authors, readers, and reviewers. In order to fix this, we propose a new peer review system termed *paper bricks*. Paper bricks has several advantages over the existing system including shorter publications, better competition for new ideas, as well as an accelerated innovation process. Furthermore, paper bricks may be implemented with minimal change to the existing peer review systems.

1. INTRODUCTION

The current peer review system is heavily criticized in a variety of scientific communities, e.g. [7, 11, 13, 6]. Examples of criticism include from the point of view of the *authors*: lack of fairness, intransparency, low quality or superficial reviews, biased reviewers, reviews based on half-read papers, decisions based on one or two reviews only, author feedback with zero impact, overfocus on getting details right, overformalized papers, overselling, as well as frustration — especially for Ph.D. students. In addition, *readers* of research papers criticize the flood of syntactically correct yet meaningless delta papers, fostering of niche topics, over-polished papers, suppress of dissent with mainstream ideas, crushing of unpolished yet interesting research ideas and directions, topic killing, missing re-experimentation, no publishing of negative results, biased experimentation, dataset and query picking, long review times, and a slow innovation process. Finally, *re-*

viewers criticize the review overload/burst at few times a year, missing reviewing standards and guidelines, as well as the huge investment to read a 12-page paper.

This criticism has led to a number of proposals in the past including *open peer review* [8], *post-publication reviews* [13, 4], and double-blind reviewing [10]. However, none of the many proposals has so far been able to substantially heal the many issues of our peer review systems.

In this paper we propose a new system for scientific peer reviews coined *paper bricks*. In a nutshell, paper bricks defines fine-granular *subpapers* that may be submitted and reviewed independently. Paper bricks is very simple to implement as it only requires minimal changes to the existing peer review system. Yet, paper bricks has many benefits for *authors*: less focus on selling details, less focus and paper polishing, less frustration. In addition, *readers* benefit from more high-impact papers, better confidence in experimental results, and an accelerated innovation process. Finally for *reviewers* there is less investment for reviews and the review load is spread out over the entire year.

2. PROBLEM STATEMENT

Design a publication culture and appropriate review system that solves the above issues. The system should guarantee quality, increase the impact of CS research in industry and other sciences, accelerate the innovation process, and provide a better experience for authors, reviewers, as well as readers of research publications.

3. HIGH-LEVEL SOLUTION IDEA

Observation: Almost all research papers have the same structure: Introduction, Problem Statement, High-Level Solution Idea, Details, Performance Evaluation.

Idea: allow people to publish pre-defined paper sections individually. The possible sections are termed

¹This publication is an extended version of the original presentation shown at the CIDR 2011 Gong Show: [The Bowyers. http://www.youtube.com/watch?v=4sorEcLjN04](http://www.youtube.com/watch?v=4sorEcLjN04)

Paper brick	Symbol
Introduction	I
Problem Statement	PS
High-Level Solution Idea	HLSI
Details	D
Performance Evaluation	PE

Table 1: List of possible paper bricks

paper bricks. Paper bricks may be submitted, reviewed, and published individually. In addition, combinations of paper bricks may be submitted, reviewed, and published as well. Thus, a *publication* consists of one or more paper brick(s). Table 1 shows a tentative list of paper bricks. Each publication must be clearly marked to signal the paper bricks it contains on its first page, e.g. this publication is marked with I+PS+HLSI.

We illustrate how this aims at solving the problem statement with a set of examples (Sections 3.1 and 3.2). After that we discuss further advantages (Section 4) and possible issues (Section 5).

3.1 Single Paper Brick Publication Examples

Example 1 (I): Just publish an Introduction paper brick to draw people’s attention to an interesting area. For instance, assume a material scientist specializing in bowery. An interesting new area for her might be “Steel Bows”. This type of paper brick must be understandable by a broad audience, e.g. an M.Sc. in bowery should be enough to understand it. This allows people from different communities to link to the same Introduction, e.g. different communities will approach the area with different tools, including possibly different Problem Statements. Introductions with potential interest to a broad community should be selected and re-published in magazines such as CACM.

Example 2 (PS): Just publish a concise Problem Statement paper brick to crisply define a research problem. This clearly defines the problem to be attacked by the community, e.g. “The Steel Bow Vibration Problem” or “The Steel Bow Arrow Problem”. PS-paper bricks must cite at least one I- or one PS-paper brick. For each I-paper brick there may be several PS-paper bricks, e.g. in Figure 1 PS-paper bricks 2 and 9 refer to I-paper brick 1. Likewise a PS-paper brick may **specialize** an existing PS-paper brick, e.g. in Figure 1 PS-paper brick 19 refers to PS-paper brick 9.

Example 3 (HLSI): Just publish a High-Level Solution Idea paper brick to sketch a possible solution. This sketches on a high-level how to solve a partic-

ular PS, e.g. “Towards Vibration-free Steel Bows” or “Steel Bows: Copper Arrows to the Rescue?”. HLSI-paper bricks must cite at least one PS- or one HLSI-paper brick. For each PS-paper brick there may be several HLSI-paper bricks, e.g. in Figure 1 HLSI-paper bricks 3 and 10 both refer to PS-paper brick 2.

Example 4 (D): Just publish a Details paper brick to precisely define the details of an HLSI-paper brick. A D-paper brick usually contains algorithms or detailed descriptions of a method or system, e.g. “How to Calm a Steel Bow” or “Two Copper Arrow Algorithms”. D-paper bricks must cite at least one HLSI- or one D-paper brick. For each HLSI-paper brick there may be several D-paper bricks, e.g. in Figure 1 D-paper bricks 21, 22, and 24 all refer to HLSI-paper brick 20; in contrast D-paper brick 25 refers to D-paper brick 24 thus specializing paper brick 24, i.e. providing ‘details on details’. Notice that the more sideways links you require to reach the root node, the higher the likelihood that the topic being treated is a niche.

Example 5 (PE): Just publish a Performance Evaluation paper brick to compare one or more algorithms. A PE-paper brick presents performance results for one or more D-paper bricks, e.g. in Figure 1 PE-paper brick 12 refers to D-paper bricks 5, 7, and 11. This is what some communities call an “experiments paper”. Valid performance evaluation techniques are Measurement, Simulation, and Analytical Modeling [12]. This type of publication allows researchers to validate the performance of different approaches, e.g. “Steel Bows: How Still are they Really?” or “On the Performance of Copper Arrows”. Again, for these types of publications it is often only vaguely described what exactly needs to be part of the publication other than experimental results [5], i.e. should the algorithm description be included as well? Is simulation a valid technique? It is [12]! We believe that only the performance evaluation setup as well as the results and their discussion should be contained in such type of publication. PE-paper bricks must cite at least one D-paper brick. For each D-paper brick there may be several PE-paper bricks.

3.2 Multiple Paper Brick Publication Examples

Example 6 (I+PS): Publish an Introduction plus Problem Statement. This is interesting for people from industry to make people in academia aware of their problems: “Here is a real problem, please solve it!”.

Example 7 (I+PS+HLSI or I+HLSI): This is what some communities call a “vision paper”. Some con-

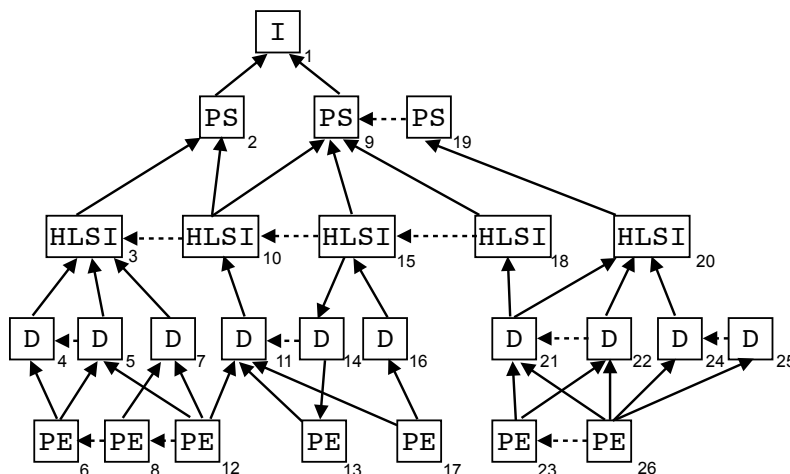


Figure 1: A PaperBrick Graph: Each node is a paper brick. Related work: Black edges are links to paper bricks on a different level. Dotted edges are links to paper bricks on the same level. Numbers show publication order.

ferences have recently started special tracks for vision papers (e.g. PLDI [3], CIDR 2011 [9], or VLDB 2011 [1]). However the calls for papers often do not clearly specify what exactly these publications should contain; some of these publications do not have a Problem Statement at all. We argue that a “vision paper” should not contain details. In addition, algorithms and experiments are not desired; they may only be used as illustration, i.e. to give directions or to support examples: “This is our vision on what we should be doing.”

Example 8 (PS+HLSI): This publication combines a crisp Problem Statement for an existing area with a High-Level Solution Idea. This type of publication allows you to identify new subproblems plus directions in an existing research area.

Example 9 (D+PE): This is basically what we often find in M.Sc. or Ph.D. theses. A D+PE publication allows students who initially can neither identify a new problem area, problem statement, nor a high-level solution idea to start with some ‘details’ for an existing PS or HLSI. Like that the scope of a student’s work is clearly defined without forcing the student to come up with a new HLSI-contribution; forcing students to come up with ‘something new’ often leads to niches. A D+PE publication also allows for **reverse order publications**: a Ph.D.-student may add other paper bricks in reverse order, e.g. he may start his work with a PE-paper brick (e.g., paper brick 13 in Figure 1). Then he may come up with ideas for D, e.g. new algorithms (paper brick 14), and eventually for an HLSI (paper brick 15). In contrast to the existing peer review system the student can easily publish the PE-part

(assuming an existing implementation) and then later send D and HLSI.

Example 10 (I+PS+HLSI+D+PE): The backward compatibility mode. These type of publications will seize to be the default, as they tend to come with all the problems discussed in Section 1.

4. ADVANTAGES

Shorter publications. There should be page limits on paper bricks, e.g. two pages for I, one for PS and so forth. Writing a shorter publication requires usually less work (unless you prove $P \neq NP$ on a single page). This also allows you to **stream a publication** to a conference: first send I. While I is being reviewed, you work on PS and HLSI. Once you receive **early feedback** on I, you adjust PS and HLSI, send a revised I plus PS and HLSI, and so forth. This will help you to **avoid wrong directions**. It will lower your time spent investing on material that will never make it to a publication. Shorter publications also means less reading, which in turn means **faster review times**. It is also likely that a shorter publication will be reviewed entirely; rather than leaving some ‘technical detail in the middle’ unreviewed. An additional advantage is that **reviewers may specialize** for certain types of paper bricks, e.g. some may focus on HLSI while others prefer to dive down into D or PE.

Accelerated innovation process. I-, PS-, and HLSI-paper bricks may be reviewed in days and can then be immediately published online. Thus review times will be considerably shorter as (1) paper bricks are shorter than “full-story”-papers, and (2) the review load is spread out over the entire

year (no fixed deadlines whatsoever — not even one deadline per month as with PVLDB). Overall, paper bricks creates a big **market of problems, ideas, and details**. That market is easily accessible by both people from industry as well as people from academia at different stages of their career (see Examples 6&9).

Less risk. Assume you have a high-level idea, but are unsure about the details. With the current system there is a high risk for you: if you wait too long to assemble all the pieces for a “full-story”-publication, another researcher may publish your idea before you. Then you will not receive *any* credit. With paper bricks this risk does not exist: rather than not publishing your idea at all, you could at least submit an HLSI-paper brick early on. If your idea is considered a valuable idea and accepted, you will **receive the credit**. Others may then receive the credit for some follow-up D-paper brick, but you will receive the credit for HLSI.

Independent experimentation as a principle. Many fields suffer from publications that present a new algorithm and its performance evaluation at the same time. As the publication pressure is hard and conferences are selective, the experimental results of those publications are often foreseeable: whatever is published is typically considerably better than *some* previous algorithm on *some* dataset for *some* queries. In addition, we have **no culture for negative results** — which often would be scientifically extremely valuable to have. This is also problematic for other researchers who would like to have a good overview on the relative performance of different algorithms. With paper bricks it is easy to write a performance evaluation for any other existing algorithm. There is no special need to justify this or criticize original publications. Repetition of experimental results from a different group is no “strange special case of a publication”. It is just fine to send such type of paper brick; in other communities this is good scientific practice anyway. Still, the quality of this PE-paper brick must meet certain standards (see Section 5).

No early crushing of high-level ideas. Fancy ideas may be proposed as an HLSI-paper brick. Yet, neither details nor algorithms are required. It is fine to have a big market of fancy ideas. People will decide in retrospect which one was the best idea.

Less investment on selling. Currently a considerable portion of the paper writing process goes into selling, i.e. justifying the work in the Introduction, contrasting it with other related work, and making sure it is different or has some other twist that was not investigated before. This goes away with paper

bricks: it is perfectly fine to have a dozen different D-paper bricks for the same HLSI-paper brick. The different D-paper bricks may have overlap in the details. However as long as they provide some improvement over existing D-paper bricks attacking the same HLSI, they may be considered. Like that a community of researchers will iterate over the D-paper bricks until they get it right. None of these works may be rejected with the argument “X considered this high-level idea before” or “a majority of the techniques have already been used in X.”

Less niche topics. Currently we often see arguments like “Although paper X provided a general solution for problem Y, it did not consider the case where <whatever>. This paper Z fills the gap.”. This defense is not required anymore with paper bricks: If there is an HLSI- or D-paper brick X solving problem Y, it is still fine to write another paper brick Z solving problem Y, if Z differs from X without inventing another niche. Hence, there will be less *forced* niches, i.e. artificially splitting Y into subproblems Y_1, \dots, Y_{42} to defend against X. Yet there will be a big market of solutions.

Better competition for best solutions: competition happens at the level of paper bricks rather than “full-story”-publications. In addition, a “full-story”-paper does not *block* (or *kill*) a topic anymore (see Section 5).

Exciting conferences. Rather than presenting each and every paper brick that gets accepted, **conferences should select** certain paper bricks for presentation. Presentation slots should be assigned based on relevance to the entire subfield. Paper bricks may be grouped into sessions, e.g. D-paper bricks solving the same PS should be in the same session. The same holds for PE-paper bricks. Notice that Q&A-sessions after a publication presentation are a form of lightweight post-reviewing anyway. As major conferences tend to attract hundreds of domain experts, the collected wisdom of these people should be explored more systematically. Therefore a publication session should be ended with a mini-panel discussing the pros and cons of the different proposals. Panelists should be authors, reviewers, as well as additional domain experts.

Furthermore, conferences may pick up the market idea by having **high-level idea sessions** where the audience votes for the best idea. The Computing Community Consortium [2] already started an initiative in this direction last year [3]. Although this initiative is a step in the right direction, we believe that it is important to clearly define what is part of such type of “vision-paper” and what is not (see Example 7). Again, other than audience voting there

Paper brick	Benchmark
Introduction	interesting area? non-existence of similar Introduction in related work? understandable by a non-domain expert? etc.
Problem Statement	link to some I or PS? correctness? non-existence of similar Problem Statement in related work? etc.
High-Level Solution Idea	link to some PS or HLSI? non-existence of similar High-Level Solution Idea in related work? etc. Notice: reviewing could use audience voting (see HLSI-post-reviewing).
Details	link to some HLSI or D? completeness of the description? pseudo-code? runtime complexity? space complexity? etc.
Performance Evaluation	link to some D? clear specification whether simulation, experiments, or analytical modeling? meaningful scenario? meaningful reference hardware? meaningful datasets and queries? appropriate benchmark? scaling experiments? appropriate baselines and competitors? etc.

Table 2: Possible paper brick reviewing guidelines

should be **panel-style discussions** allowing for systematic, immediate feedback on the presented ideas — not just three questions or a question by the sessions chair to break the silence. An extension of this could be **HLSI-post-reviewing**. An HLSI-publication may be submitted as a 4-page paper. If the contribution gets ‘accepted’, this means the authors will be allowed to give a short presentation at the conference. The attendees will vote for the best contributions. Only if the vote by the audience is above a certain threshold, the HLSI-publication will obtain a larger slot in the proceedings, i.e. up to the originally submitted 4 pages. All other papers are only allowed to publish a 1-page extended abstract. Technically, this may be implemented by asking authors of accepted HLSI-publications to provide two camera-ready versions: a 1-page and a 4-page version.

Connecting the dots. Paper bricks form the nodes in a huge graph of research (see Figure 1). Related work are the edges. Therefore, another idea for conferences could be sessions where a panel plus the auditorium identifying edges, e.g. “this PS is also relevant for I X!”, “this HLSI may also solve PS Y!”, etc.

Best-of journal publications. There may still be “full-story” research publications. However, they will look different: In paper bricks a journal publication is created by picking the best paper bricks from possibly different authors and republish them as one big publication at a journal. This best-of journal publication identifies the best paper bricks and shows the complete story. For instance, in Figure 1 the path 1, 9, 18, 21, 26 might define a best-of publication.

Domain interfacing. Each paper brick has links

to paper bricks on different levels. These links serve as **interfaces to other domains**. For instance, assume you write a PS-paper brick in the domain of databases. A domain expert from algorithmic optimization may pick up that PS and give hints for a solution in an HLSI-paper brick. Details may then be worked out by the database researcher again in a D-paper brick. Thus, paper bricks allows people from different domains of computer science to collaborate more easily on the same problem.

Better assessment of researchers. Paper bricks helps students to find a job better matching their qualification after graduating. It is also helpful for job committees in both academia and industry.

Which type of paper bricks did she publish? Is she creative in finding high-level ideas? Is she good in details and/or algorithms? Or does she do well in performance evaluations? With the current system peer review system it is not always clear who of the authors contributed to which part of a long publication. With paper bricks, this becomes more clear. The specific skills of a researcher become more explicit and thus it becomes easier to assess a student.

Clear contributions. Most “full-story” research publications make several contributions. With paper bricks, the individual contributions are submitted and reviewed independently. Hence, publications are **easier to judge** by reviewers and readers.

5. SOME ISSUES AND POSSIBLE SOLUTIONS

In this section we briefly sketch some issues and possible solutions with paper bricks.

Title, Abstract, and Conclusion belong to every publication including single paper brick publi-

cations.

Where does **related work** go? Paper bricks form the nodes of a large graph; related work defines the edges. Each paper brick must provide appropriate edges (see Figure 1 for an example). It is up to the authors to decide whether related work is discussed in a separate section or integrated into the main sections.

How to review paper brick publications? Overall: the quality of the write-up matters, i.e. English, brevity, up to the point yet readable. In addition, there should be **different reviewing guidelines** for the different types of paper bricks. Table 2 shows a proposal for these reviewing guidelines. That table is just an initial proposal and can definitely be extended in many ways.

For instance, for I-paper bricks reviewers should evaluate whether the publication describes an interesting area. In addition, reviewers should make sure that that area was not published before. Furthermore, reviewers should enforce that an I-paper brick is understandable by a non-domain expert to facilitate interaction with other domains. For D-paper bricks reviewers should make sure that there is at least one link to an HLSI- or a D-paper brick. Reviewers should check whether the description provided by the authors is complete. Would someone else be able to implement the proposed solution just reading that description? Do the authors provide meaningful pseudo-code? Do the authors discuss runtime and space complexities of their approach?

Notice that the overall goal of having a reviewing guideline is to avoid that reviewers use different standards for reviewing — which would lead to high variance in paper assessments. There should be a small catalogue of reviewing rules. That catalogue should be accessible to both reviewers and authors.

6. CONCLUSIONS

We have proposed *paper bricks*. *Paper bricks* allows for a fine-granular yet peer-reviewed way of doing research. We believe that paper bricks has strong advantages over the existing peer review systems. As future work we are planning to explore paper bricks in a D-paper brick. We would also like to see major conferences trying out our approach. A possible way to bootstrap paper bricks is to add a special **paper bricks track** to an existing conference. If the experience with such a track is positive, it could be extended to other tracks.

This is not “us vs. them”.

There is only us.

Don’t wait for “them” to change things.

We are “them”. [6]

Acknowledgements. I would like to thank the attendees from CIDR 2011, numerous colleagues from Saarland University, as well as peers watching my [youtube-video](#) for their valuable feedback on the initial idea.

7. REFERENCES

- [1] Challenges and Visions Track at VLDB 2011, <http://www.vldb.org/2011/?q=node/17>.
- [2] Computing Community Consortium, <http://www.cra.org/ccc/>.
- [3] Crazy-Idea Sessions at PLDI 2010, <http://www.cccblog.org/2010/07/26/pldisfun-ideas-thoughts-stimulating-new-research-visions/>.
- [4] Dan S. Wallach, Rebooting the CS Publication Process, <http://www.cs.rice.edu/~dwallach/pub/reboot-2010-06-14.pdf>.
- [5] Experiments&Analyses Track at VLDB 2009, <http://vldb2009.org/?q=node/4#ea>.
- [6] Jeff Naughton, ICDE 2010 Keynote, <http://pages.cs.wisc.edu/~naughton/naughtonicde.pptx>.
- [7] Nature Peer Review Debate, <http://www.nature.com/nature/peerreview/debate/index.html>.
- [8] Opening up BMJ, BMJ 1999; 318 : 4, <http://www.bmj.com/content/318/7175/4.full>.
- [9] Outrageous Ideas and Vision Track at CIDR 2011, <http://www.cccblog.org/2011/01/18/outrageous-ideas-at-cidr-seeking-to-stimulate-innovative-research-directions/>.
- [10] SIGMOD 2001 Double-Blind Reviews, <http://www.dbnet.ece.ntua.gr/~timos/submitinfo-2001.html>.
- [11] The Scientist: Is Peer Review Broken?, <http://www.the-scientist.com/article/display/23061>.
- [12] R. Jain. *The Art of Computer Systems Performance Analysis: Techniques for Experimental Design, Measurement, Simulation, and Modeling*. Wiley-Interscience, 1991.
- [13] J. C. Mogul and T. Anderson. Open Issues in Organizing Computer Systems Conferences. *SIGCOMM Comput. Commun. Rev.*, 38:93–102, July 2008.



ADVANCE CALL FOR PARTICIPATION

14th International Workshop on the Web and Databases (WebDB 2011)

Collocated with ACM SIGMOD 2011

Athens, Greece - June 12, 2011

WebDB Goals

The WebDB workshop provides a forum where researchers, theoreticians, and practitioners can share their insights and their knowledge on problems and solutions at the intersection of data management and the Web. WebDB has high impact and has been a forum in which a number of seminal papers have been presented.

Website: <http://webdb2011.rutgers.edu/>

Invited talk by Nick Koudas

We are happy to announce [Nick Koudas](#) (University of Toronto, Sysomos) as our keynote speaker. Nick is co-founder of Sysomos, a Marketwire company that is one of the pioneers of social media analytics.

Topics

WebDB 2011 will include papers on some of the following topics:

- Cloud computing and distributed computing over the Web
- Collaborative data management on the Web
- Data integration over the Web
- Data models and query languages for Web information
- Data-intensive applications on the Web
- Database support for social network and Web 2.0 applications
- Human computing in Web databases
- Information retrieval in semistructured data and the Web
- Filtering and recommendation systems
- Location-aware Web applications
- Modeling, mining and querying user generated content
- Personal information management systems
- Quality of user generated content and other Web data
- Semantic search on the Web
- Semi-structured data management

- Social and tagged data management
- The Semantic Web and reasoning on Web data
- Visualization of social network and related data
- Web community data management systems
- Web information extraction
- Web privacy and security
- Web source discovery, analysis, and retrieval
- Web-based distributed data management

The list of accepted papers will be announced in the first week of May 2011.

Registration

<http://webdb2011.rutgers.edu/register.html>

Fees, in US\$:

ACM Members: \$80

Non members: \$100

Students (ACM members and non-members): \$35

Workshop registration fees are fixed; there is no early registration deadline for WebDB.

We are looking forward to seeing you in Athens!

Amelie Marian, Rutgers University, USA

Vasilis Vassalos, Athens University of Economics and Business, Greece