

Towards a Theory of Search Queries

ANONYMOUS
Omitted for double-blind review

The need to manage diverse information sources has triggered the rise of very loosely structured data models, known as “dataspace models.” Such information management systems must allow querying in simple ways, mostly by a form of searching. Motivated by these developments, we propose a theory of search queries in a general model of dataspace. In this model, a dataspace is a collection of data objects, where each data object is a collection of data items. Basic search queries are expressed using filters on data items, following the basic model of boolean search in information retrieval. We characterize semantically the class of queries that can be expressed by searching. We apply our theory to classical relational databases, where we connect search queries to the known class of fully generic queries, and to dataspace where data items are formed by attribute–value pairs. We also extend our theory to a more powerful, associative form of searching where one can ask for objects that are similar to objects satisfying given search conditions. Such associative search queries are shown to correspond to a very limited kind of joins. Specifically, we show that the basic search language extended with associative search can define exactly the queries definable in a restricted fragment of the semijoin algebra working on an explicit relational representation of the dataspace.

Categories and Subject Descriptors: F.4.1 [Mathematical Logic and Formal Languages]: Mathematical Logic—*Modal logic*; H.2.1 [Database Management]: Logical Design—*Data Models*; H.2.3 [Database Management]: Languages—*Query languages*; H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval

General Terms: Design, Languages, Theory

Additional Key Words and Phrases: Dataspace, genericity, search

1. INTRODUCTION

One can say that querying by a form of searching has become the norm. Indeed, in most current information systems such as Web search engines, e-commerce sites, or desktop search systems, as well as in classical information retrieval systems such as library catalogs or document repositories, users query the database by means of a search interface. Searching is expressed by means of keywords which can be combined with boolean operators. Such a basic search facility is much weaker than the standard database query languages, where select–project–join queries are considered the minimum and full-fledged first-order logic is considered the norm (cf. Codd’s relational algebra and calculus). Contemporary languages such as SQL/PSM or XQuery are even computationally complete.

Database queries are a major theme of database theory [Abiteboul et al. 1995]. In general, queries are generic mappings from databases to relations, where ‘generic’ refers to invariance under isomorphisms [Chandra and Harel 1980; Aho and Ullman 1979]. Many classes of database queries have been identified and characterized in terms of semantic properties; expressibility in various query languages; or computability under various complexity limitations. Since searching is such a simple, natural, and important form of querying, it appears that search queries deserve to have their own chapter in the theory of database queries. Our goal in this paper

2 .

is to propose a first draft of such a chapter. Apart from the foundational motivation, our work is further motivated by two important trends in data management research: dataspace and usability.

Dataspace [Halevy et al. 2006; Dong and Halevy 2007; Dittrich and Vaz Salles 2006] are a new type of databases, characterized by a very loosely structured data model and geared towards the management of data coming from a diverse set of sources. Dataspace are queried by a form of searching. In essence, the data in a dataspace is modeled as a collection of objects, where each object is a collection of attribute-value pairs. Dataspace are queried by searching for conditions on attributes or values, or by following links between objects. Relationally complete querying of dataspace-like databases was studied earlier in the context of models for data integration [Jain et al. 1995]. Recent work on data processing for cloud computing has also focused on dataspace-like databases (e.g., [DeCandia et al. 2007; Olston et al. 2008]).

Improving the usability of database systems has been an issue pretty much since the beginning of database research, as witnessed, for example, by the past research on universal relation interfaces (surveyed by Ullman [1989, Chapter 17]). Recently, Jagadish et al. [2007] have revived our interest in this topic; he argues, among other things, that queries involving explicit joins or subqueries are too cumbersome to express, and stimulates us to ask how far we can get with simpler forms of querying such as searching, or with more implicit or automatic ways of joining information.

So, in a nutshell, in this paper, we try to understand formally the question of how much database querying can be done using a basic search language, possibly extended with a simple facility for following links between objects.

We should also clarify what we do *not* do in this paper. We do not investigate how searching can be implemented efficiently. Rather, we focus on expressive power. Also, because of this focus, we ignore other important issues that have been investigated in research on keyword search in relational, tree-structured (XML), and semistructured (graph) databases [Golenberg et al. 2008; Qin et al. 2009, we give just two recent references]. The two main such issues are automatically finding connections among objects in the database that contain the given keywords (which is a nice approach to the usability question mentioned above), and ranking the results of a keyword search.

Concretely, the contributions of this paper can be summarized as follows:

- (1) We define a general formal model of dataspace, where a dataspace is a collection of data objects, and where an object is a collection of data items. On these data items, we assume a number of abstract filter predicates to be defined. These filters naturally serve as atomic search conditions: searching a dataspace with a filter returns all objects containing an item satisfying the filter. We obtain a basic search language by combining atomic searches using the boolean set operators union, intersection, and difference.
- (2) Search queries are defined in general as functions on dataspace that map a dataspace to one of its subsets. The question then arises of exactly which such mappings are definable in the basic search language. This question turns out to be naturally answered by requiring the search queries to be invariant under a natural indistinguishability relation on objects. The concept of genericity

(invariance) has always been a central theme in the development of the theory of database queries [Abiteboul et al. 1995]; we have tried here to get at the right genericity concept for search queries.

- (3) We apply the above semantic characterization to the case of classical database relations. A relation can be viewed as a dataspace by viewing the tuples as sets of attribute–value pairs. When using as filters on such attribute–value pairs the classical selection conditions “attribute equals constant,” we obtain as basic search queries precisely the class of search queries that are “fully generic” in the sense of Beer, Milo and Ta-Shma [1996; 1997]. (We hasten to add that these authors looked at full genericity of queries in a complex-object setting much richer than mere search queries on flat relations.)
- (4) We extend the basic model to allow for so-called “associative search”, where one can search not just for all objects satisfying some search query, but also for all objects that are somehow linked to those objects. This amounts to adding a link operator to the basic search language, much in the spirit of modal logic [Blackburn et al. 2001; Blackburn et al. 2007]. Associative search is indeed a feature of most dataspace and keyword search query languages proposed in the literature [Halevy et al. 2006; Dong and Halevy 2007; Dittrich and Vaz Salles 2006; Golenberg et al. 2008]. Actually, in these proposals, links between objects are often assumed to be found automatically by the system. Since in this paper we are interested in a detailed analysis of expressive power, we focus on the case where the linking conditions are explicitly specified in the query.
- (5) A question we find interesting is how search query languages relate to standard database query languages. Those languages remain relevant in the search and dataspace setting. For example, SPARQL [spa 2008], the standard language to query RDF graphs [rdf 2004], is closely related to database queries on ternary relations [Fletcher 2008; Gutierrez et al. 2004; Pérez et al. 2009]; RDF graphs can be viewed as a dataspace model.
Under the natural assumption that linking between two objects is done in terms of a set of similarity relations among the items in these objects, we can indeed relate associative search to standard database query languages. Specifically, we observe that associative search queries are definable in the semijoin algebra: the version of the relational algebra where the join operator is replaced by the semijoin [Leinders et al. 2005; Leinders and Van den Bussche 2007; Ross and Janevski 2005]. Here, the semijoin algebra works on the natural representation of a dataspace as a binary relation; abstract filters are used as selection conditions, and abstract similarity predicates are used as join conditions.
Conversely, however, not every semijoin query is an associative search query. We actually identify three kinds of constructions that are definable in the semijoin algebra but not in our associative search language, and prove that the fragment of the semijoin algebra in which these three constructions are disallowed fully characterizes our associative search queries.
- (6) Finally, we show that our general abstract theory is workable by applying it to the most common concrete dataspace setting, where data items are attribute–value pairs. Unlike the application to classical relations mentioned earlier, here, there is no fixed schema and objects can have multiple values for the

same attribute. We consider a natural repertoire of filters that can test if the attribute, and the value, equals, or is different from, a finite number of possibilities. In this setting, the semantic characterization of basic search queries can be rephrased in a very intuitive manner, as we will show. We will also instantiate the semijoin algebra characterization of associative search queries to the attribute–value setting. Associative search queries over attribute–value dataspace are thus characterized as the queries expressible in a simple and attractive fragment of the semijoin algebra; this time, dataspace are represented by ternary relations with schema (Oid, Att, Val), selection conditions are constant equalities (on attributes as well as on values; remember that there is no dataspace schema), and all semijoins are simple equijoins.

The current paper extends earlier work on this topic, including non-trivial proofs of results, applications of the theory, and further discussion [Anonymous].

This paper is organized as follows. Section 2 develops the general theory of search queries on the model of abstract dataspace, both in terms of boolean search queries and associative search queries. Section 3 presents the application to classical database relations. Section 4 presents the application to attribute–value dataspace. We finish with a discussion in Section 5.

2. A THEORY OF SEARCH QUERIES

In this section we propose a theory of search queries in a general model of dataspace. In this model, a *dataspace* is a collection of data *objects*. Each data object is a collection of data element called *items*, drawn from a universe of items \mathcal{I} . Although we will keep \mathcal{I} abstract in this section in order to obtain a generally applicable theory, a rich set of well-known, practical data models can be obtained by picking concrete interpretations for \mathcal{I} . For instance:

- Items could be words over some alphabet. Objects then correspond to documents in the classical boolean model of information retrieval [Baeza-Yates and Ribeiro-Neto 1999, chapter 4.2.3], and dataspace correspond to collections of such documents. Figure 1(a) illustrates such a dataspace, consisting of two documents concerning events in St. Petersburg.
- Items could be pairs (a, v) of attributes a and values v . Objects, being finite sets of such pairs, then intuitively describe the attributes of a real-world entity. (Note that there can be multiple values for the same attribute.) A dataspace is just a collection of such descriptions. Figure 1(b), for example, depicts an attribute-value dataspace describing researchers, drinks, and papers.
- Staying with attribute–value pairs, we could consider those dataspace whose objects contain only attributes from a fixed relation schema, and contain a single value for each such attribute. Such dataspace are classical database relations.

Definition 1 (Dataspace). Formally, let \mathcal{I} be a set of data element called items. An *object over \mathcal{I}* (by short: \mathcal{I} -*object*) is a finite, non-empty set of items. A *dataspace over \mathcal{I}* (by short: \mathcal{I} -*dataspace*) is a finite set of \mathcal{I} -objects.

Throughout the remainder of this section we fix \mathcal{I} arbitrarily and simply talk about objects and dataspace instead of \mathcal{I} -objects and \mathcal{I} -dataspace, respectively.

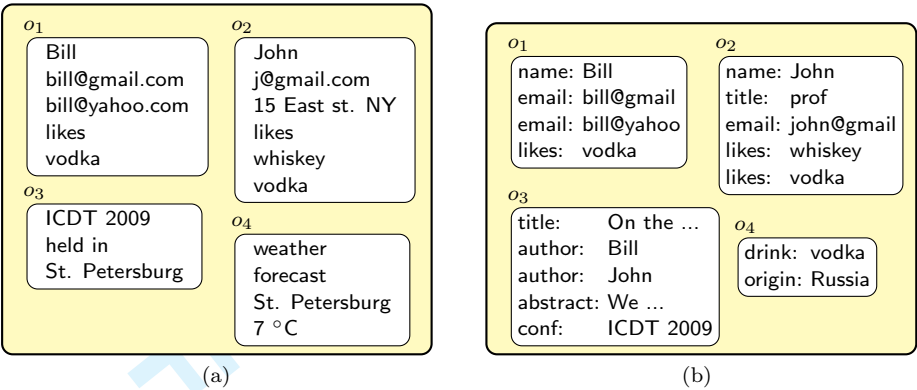


Fig. 1. (a) An example document space and (b) an example attribute-value space.

For ease of exposition, however, we will illustrate our definitions and results in the situation where \mathcal{I} is the set of all words, as in classical information retrieval. The other concrete situations above, where \mathcal{I} is a set of attribute-value pairs, are studied in greater depth in Sections 3 and 4.

2.1 Boolean Search

As a first basic language to express search queries on dataspaces, we consider the Boolean Search Language BSL. Similar to the boolean information retrieval setting, searches in BSL are formed by combining atomic search conditions (keywords) using the set operators union, intersection, and difference. First, let us make precise the notion of a keyword.

Definition 2 (Keywords). For our purposes, a *keyword* is defined semantically as a subset $k \subseteq \mathcal{I}$. An item i is said to *match* k if it is an element $i \in k$. An object o *satisfies* k , denoted $o \models k$, if it contains an item that matches k , i.e., if $o \cap k$ is non-empty.

Keywords naturally serve as atomic search conditions: searching a dataspace with a keyword returns all objects that satisfy the keyword. For example, let *contains* i denote the keyword that consists of all items that “contain” i (e.g., as a substring). Searching by *contains* ICDT then returns all objects in the dataspace that contain an item in which ICDT occurs as a substring.

In what follows, we write \mathcal{K} for the set of all keywords. Two kinds of keywords deserve special attention:

Definition 3 (Wildcard and literal keywords). The universe \mathcal{I} itself is the keyword that matches all items. We refer to it as the *wildcard* in what follows. To draw attention to the fact that it is used as a keyword rather than the universe of all items, we will denote it by \star instead of \mathcal{I} .

For every item i , the singleton keyword $\{i\}$ is called a *literal* keyword. We often abuse notation and omit the set braces from these keywords, writing i instead of $\{i\}$.

6 .

Equipped with the atomic search conditions, we are ready to formally define the Boolean Search Language BSL:

Definition 4. The expressions of BSL are given by the grammar

$$e ::= k \mid e \text{ and } e \mid e \text{ or } e \mid e \text{ except } e,$$

where k ranges over the keywords in \mathcal{K} . Semantically, an expression e can be applied to a dataspace D , resulting in a subset $e(D)$ of D defined as follows:

$$\begin{aligned} k(D) &:= \{o \in D \mid o \models k\} \\ (e_1 \text{ and } e_2)(D) &:= e_1(D) \cap e_2(D) \\ (e_1 \text{ or } e_2)(D) &:= e_1(D) \cup e_2(D) \\ (e_1 \text{ except } e_2)(D) &:= e_1(D) - e_2(D). \end{aligned}$$

Note that the language is a bit redundant as $e_1 \text{ and } e_2$ is equivalent to $e_1 \text{ except } (e_1 \text{ except } e_2)$.

It is important to appreciate the existential nature of the above-defined semantics. Thus, the expression $k_1 \text{ except } k_2$ does *not* return all objects that contain an item that matches k_1 but not k_2 ; rather, it returns all objects that contain an item that matches k_1 , but do not contain an item that matches k_2 .

Example 5. To illustrate the semantics, consider the setting where items are words. Then the expression

$$e_1 := \{\text{weather}\} \text{ and } \{\text{forecast}\} \text{ and } \{\text{St. Petersburg}\}$$

retrieves all documents in which the words *weather*, *forecast*, and *St. Petersburg* all occur. By our convention that we identify items i with the singleton keyword $\{i\}$ that matches i only, we will simply write this expression as

$$\text{weather and forecast and St. Petersburg}$$

in what follows. As another example, suppose that, for each item i , *similar to i* is a keyword that matches all items ‘similar’ to i according to some notion of similarity. Then the expression

$$e_2 := \text{vodka except (similar to whiskey)}$$

retrieves all documents in which the drink *vodka* occurs, but no drink similar to *whiskey* occurs. In particular, on the document space in Figure 1(a), e_1 and e_2 return $\{o_4\}$ and $\{o_1\}$, respectively.

Of course, in real-life situations, one rarely has *all* keywords available to search with. For that reason, we denote by $\text{BSL}(K)$ the fragment of BSL expressions that use only keywords in $K \subseteq \mathcal{K}$. Thus, e_1 in Example 5 is in $\text{BSL}(\text{weather, forecast, St. Petersburg})$, but e_2 is not.

As just defined, every BSL expression defines a *search query*:

Definition 6 (Search Query). A *search query* is a mapping q from dataspace to dataspace such that $q(D) \subseteq D$ for each dataspace D . We say that a search query is *definable in BSL* if there exists an expression e such that $e(D) = q(D)$ for every D .

Of course not all search queries are definable in BSL. Which ones are? We can answer this question by identifying three typical properties of BSL queries: *additivity*, *K-safety*, and *K-distinguishing*. We define these three properties next.

Definition 7 (Additivity, K-safety, K-distinguishing). Let q be a search query and let $K \subseteq \mathcal{K}$ be a set of keywords.

—We say that q is *additive* if

$$q(D) = \bigcup_{o \in D} q(\{o\})$$

for any dataspace D .

—We say that q is *K-safe* if for any dataspace D and any $o \in q(D)$, we have that $o \models k$ for at least one $k \in K$.

—Two objects o_1 and o_2 are called *K-equivalent* if for all $k \in K$ we have $o_1 \models k$ iff $o_2 \models k$. We denote this by $o_1 \simeq_K o_2$. (Clearly, \simeq_K is an equivalence relation.) We then say that q is *K-distinguishing* if for any two dataspaces D_1 and D_2 and objects $o_1 \in D_1$ and $o_2 \in D_2$ that are *K-equivalent*, we have $o_1 \in q(D_1)$ iff $o_2 \in q(D_2)$.

The above three properties represent three distinctive features of BSL queries. Additivity simply states that the query can be processed one object at a time. *K-safety* states that we cannot retrieve arbitrary objects from the dataspace, but only objects that satisfy at least one of the specified keywords. This is also the case in all real-life search engines and information retrieval systems. Finally, *K-distinguishing* naturally states that the query can only distinguish between objects based on their satisfaction of specified keywords.

As a matter of fact, *additivity already follows from K-distinguishing*, since the latter property implies $o \in q(D)$ iff $o \in q(\{o\})$ which readily implies additivity. We stated the property of additivity separately because we will need it later independently of *K-distinguishing*.

We establish the following semantic characterization:

THEOREM 8. *Let $K \subseteq \mathcal{K}$ be a set of keywords. The following are equivalent for any search query q :*

- (1) q is definable in BSL(K)
- (2) q is L -safe and L -distinguishing, for some finite set $L \subseteq K$.

PROOF. (1 \rightarrow 2) Suppose that $e \in \text{BSL}(K)$ defines q . Take L to be the finite set of keywords mentioned in e . The properties of L -safety and L -distinguishing are then readily verified by structural induction on e (details omitted).

(2 \rightarrow 1) Suppose that search query q is L -safe and L -distinguishing with $L \subseteq K$ finite. First observe that \simeq_L has only a finite number of equivalence classes. Indeed, since two objects are L -equivalent only if they satisfy the same keywords in L , we can find for each equivalence class C a subset of keywords $L_C \subseteq L$ such that

$$C = \{o \mid o \models k \text{ for all } k \in L_C \text{ and } o \not\models k \text{ for all } k \in L - L_C\}.$$

Since there are only a finite number of subsets $L_C \subseteq L$, there can only be a finite number of equivalence classes.

8 .

Then let C_1, \dots, C_s be those equivalence classes of \simeq_L that contain an object o with $o \in q(\{o\})$. We say that these equivalence classes *satisfy* q . Note that, since q is L -distinguishing, for each object o' we have $o' \in q(\{o'\})$ if and only if $o' \in C_1 \cup \dots \cup C_s$.

Then observe that each equivalence class $C \in \{C_1, \dots, C_s\}$ is definable in $\text{BSL}(L)$ in the sense that we have an expression e_C such that $o \in C$ iff $o \in e_C(\{o\})$. Indeed, pick $o \in C$ arbitrarily, let $\{k_1, \dots, k_m\} = L_C$ be the finite set of all keywords in L that o satisfies, and let $\{l_1, \dots, l_n\} = L - L_C$ be finite set of all keywords in L that o does not satisfy. Let

$$e_C := (k_1 \text{ and } \dots \text{ and } k_m) \text{ except } (l_1 \text{ or } \dots \text{ or } l_n).$$

Clearly, for any object o' we have $o' \in e_C(\{o'\})$ if and only if $o' \simeq_L o$ (and hence $o' \in C$). Note that e_C is a valid expression only if $m \geq 1$. However, since q is L -safe and C satisfies q , we know that the objects in C must satisfy at least one keyword in L . Thus, we are assured that $m \geq 1$.

Then q is defined in $\text{BSL}(L)$ by $e := e_{C_1} \text{ or } \dots \text{ or } e_{C_s}$. Indeed, by the only-if direction above we know that every expression in $\text{BSL}(L)$ is L -distinguishing, and thus additive. In particular, e is additive. Then observe that for any dataspace D and any $o \in D$ we have

$$\begin{aligned} o \in (e_{C_1} \text{ or } \dots \text{ or } e_{C_s})(D) \\ &\Leftrightarrow o \in (e_{C_1} \text{ or } \dots \text{ or } e_{C_s})(\{o\}) \quad (\text{since } e \text{ is additive}) \\ &\Leftrightarrow o \in e_{C_1}(\{o\}) \cup \dots \cup e_{C_s}(\{o\}) \\ &\Leftrightarrow o \in C_1 \cup \dots \cup C_s \quad (\text{by construction of } e_{C_i}) \\ &\Leftrightarrow o \in q(\{o\}) \\ &\Leftrightarrow o \in q(D) \quad (\text{since } q \text{ is } L\text{-distinguishing, and hence additive}) \end{aligned}$$

As such, e defines q . \square

We point out that in the presence of the wildcard keyword, the issue of L -safety becomes moot:

COROLLARY 9. *When $K \subseteq \mathcal{K}$ includes the wildcard keyword \star , a search query is definable in $\text{BSL}(K)$ if and only if it is L -distinguishing for some finite $L \subseteq K$.*

PROOF. The only-if direction follows directly from Theorem 8. For if-direction, first observe that $o \simeq_L o'$ iff $o \simeq_{L \cup \{\star\}} o'$. Hence, since q is L -distinguishing, it is also $(L \cup \{\star\})$ -distinguishing. Moreover, since all objects satisfy \star , every query is $(L \cup \{\star\})$ -safe. Therefore, q is both $(L \cup \{\star\})$ -safe and $(L \cup \{\star\})$ -distinguishing, and hence definable in $\text{BSL}(K)$ by Theorem 8. \square

The following example gives two illustrations of how one can show that certain queries are not definable in BSL with only a specific set of keywords.

Example 10. One may wonder if negations of keywords can be expressed: can we express the query “ $\neg k$ ”,

$$(\neg k)(D) = \{o \in D \mid \exists i \in o : i \notin k\},$$

that retrieves all objects containing an item that does not match some fixed keyword k ? The answer in general is no. In proof, suppose that k is a keyword with at least

one matching item i and at least one nonmatching item j . Then the query $\neg k$ is not L -distinguishing for any $L \subseteq \{k, \star\}$. Indeed, $\{i\} \simeq_L \{i, j\}$ for any such L , but $\{i, j\}$ satisfies the query whereas $\{i\}$ does not. Hence, $\neg k$ is not definable in $\text{BSL}(k, \star)$.

Of course, one can add negated keywords directly to the set of allowed keywords to express such queries. One may then wonder whether that is enough to allow already all boolean combinations of keywords to be expressed. For example, can we now retrieve all objects containing an item that matches neither k_1 nor k_2 ? The answer is still no. In proof, suppose that k_1 and k_2 are two keywords and i_1, i_2, j_0 , and j are items such that i_1 matches k_1 but not k_2 ; i_2 matches k_2 but not k_1 ; j_0 matches neither k_1 nor k_2 ; and j matches both k_1 and k_2 . Let $\neg k_1$ and $\neg k_2$ be the negated keywords of k_1 and k_2 , respectively. (I.e., $\neg k_1 := \mathcal{I} - k_1$, and similarly for $\neg k_2$.) Then the query “ $\neg k_1 \wedge \neg k_2$ ” is not L -distinguishing for any $L \subseteq \{k_1, k_2, \neg k_1, \neg k_2, \star\}$ and is thus not definable in $\text{BSL}(k_1, k_2, \neg k_1, \neg k_2, \star)$. Indeed, $\{i_1, i_2\} \simeq_L \{j_0, j\}$ for any such L , but $\{j_0, j\}$ satisfies the query whereas $\{i_1, i_2\}$ does not.

We conclude this section by pointing out that Theorem 8 continues to hold if one is only interested in definability of q with respect to a particular class of dataspace, in the following sense.

Definition 11. Let \mathcal{D} be a subclass of the class of all dataspace and let $K \subseteq \mathcal{K}$ be a set of keywords.

- A search query q is *definable on \mathcal{D} in BSL* if there exists an expression e such that $e(D) = q(D)$ for every $D \in \mathcal{D}$.
- A search query q is *K -safe on \mathcal{D}* if for any $D \in \mathcal{D}$ and any $o \in q(D)$, we have $o \models k$ for at least on $k \in K$.
- A search query q is *K -distinguishing on \mathcal{D}* if for any two dataspace $D_1 \in \mathcal{D}$ and $D_2 \in \mathcal{D}$ and any two objects $o_1 \in D_1$ and $o_2 \in D_2$ that are K -equivalent, we have $o_1 \in q(D_1)$ iff $o_2 \in q(D_2)$.

COROLLARY 12. *Let $K \subseteq \mathcal{K}$ be a set of keywords. The following are equivalent for any search query q and any class of dataspace \mathcal{D} :*

- (1) q is definable on \mathcal{D} in $\text{BSL}(K)$
- (2) q is L -safe and L -distinguishing on \mathcal{D} , for some finite set $L \subseteq K$.

As in Corollary 9, the condition that q be L -safe can be dropped when $\star \in K$.

PROOF. (1 \rightarrow 2) Suppose that $e \in \text{BSL}(K)$ defines q on \mathcal{D} . Since $e \in \text{BSL}(K)$, we know by Theorem 8 that the query defined by e is L -safe and L -distinguishing on the class of *all* dataspace, for some finite $L \subseteq K$. This implies in particular that e is L -safe and L -distinguishing on \mathcal{D} . Hence, since q coincides with e on \mathcal{D} , so is q .

(2 \rightarrow 1) Suppose that search query q is L -safe and L -distinguishing on \mathcal{D} , with $L \subseteq K$ finite. Define the search query p on the class of *all* dataspace as follows:

$$p(D) = \{o \in D \mid \text{there exists } D' \in \mathcal{D} \text{ and } o' \in q(D') \text{ with } o \simeq_L o'\}.$$

This search query is L -safe on the class of all dataspace. Indeed, if $o \in p(D)$ then there exists some $D' \in \mathcal{D}$ and some $o' \in q(D')$ with $o \simeq_L o'$. Since q is L -safe on \mathcal{D} we know that there exists $k \in L$ with $o' \models k$. Then, since $o \simeq_L o'$, also $o \models k$.

10 .

The query is also L -distinguishing on the class of all dataspace: let D_1 and D_2 be arbitrary dataspace, and suppose that $o_1 \in D_1$ and $o_2 \in D_2$ are L -equivalent. Then:

$$\begin{aligned} o_1 \in p(D_1) &\Leftrightarrow \text{there exists } D' \in \mathcal{D} \text{ and } o' \in q(D') \text{ with } o_1 \simeq_L o' \\ &\Leftrightarrow \text{there exists } D' \in \mathcal{D} \text{ and } o' \in q(D') \text{ with } o_2 \simeq_L o' \\ &\Leftrightarrow o_2 \in p(D_2) \end{aligned}$$

By Theorem 8, p is definable in $\text{BSL}(L)$ on the class of all dataspace. The corollary then follows since $p(D) = q(D)$ for all $D \in \mathcal{D}$. Indeed, if $D \in \mathcal{D}$ and $o \in q(D)$ then also $o \in p(D)$ by definition of p . Conversely, if $D \in \mathcal{D}$ and $o \in p(D)$ then there exists $D' \in \mathcal{D}$ and $o' \in q(D')$ with $o \simeq_L o'$. Since q is L -distinguishing on \mathcal{D} , also $o \in q(D)$. \square

2.2 Associative Search

Since BSL queries are additive, BSL cannot define queries that relate objects to other objects. In other words: BSL cannot do joins. For example, in the information retrieval setting (where items are words) the search query

“retrieve all objects that share a word with an object in which ICDT 2009 occurs”

is not additive, and therefore not definable in BSL. In dataspace systems, however, we often want to be able to retrieve not just all objects that satisfy some search query, but also all objects that are related to those objects [Dong and Halevy 2007; Dittrich and Vaz Salles 2006]. To this end, we extend our theory as follows.

Definition 13 (Link condition). Let \mathcal{O} denote the set of all objects. A *link condition* is a binary relationship $\lambda \subseteq \mathcal{O} \times \mathcal{O}$ on objects. We write $\lambda(o, o')$ to denote $(o, o') \in \lambda$, and write \mathcal{L} for the set of all link conditions.

We then define the *associative search language* ASL as an extension of the basic search language BSL with an operator for retrieving related objects:

Definition 14. The expressions of ASL are given by the grammar

$$e ::= k \mid e \text{ and } e \mid e \text{ or } e \mid e \text{ except } e \mid \text{link}\langle\lambda\rangle e,$$

where λ ranges over \mathcal{L} . The semantics of the new construct is given by

$$(\text{link}\langle\lambda\rangle e)(D) := \{o \in D \mid \exists o' \in e(D) : \lambda(o, o')\}.$$

Example 15. As a simple example on document spaces, let *shares* be the link condition that links o to o' if and only if o and o' share a common word. Then the expression $\text{link}\langle\text{shares}\rangle \text{ICDT 2009}$ returns all documents that share a word with a document in which ICDT 2009 occurs. In particular, on the dataspace in Figure 1(a), it returns $\{o_3, o_4\}$.

Of course, in real-life situations, we do not have *all* possible link conditions available to search with. For that reason, we are primarily interested in the queries definable in $\text{ASL}(K, L)$, where $\text{ASL}(K, L)$ denotes the fragment of ASL expressions that use only keywords in $K \subseteq \mathcal{K}$ and link conditions in $L \subseteq \mathcal{L}$.

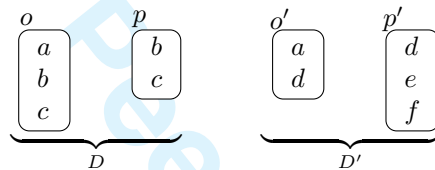
It turns out that there is a connection between ASL and modal logic; indeed, ASL is a modal logic. We can make this connection precise by defining a bisimulation notion appropriate for ASL.

Definition 16 (Bisimilarity). A *pointed dataspace* is a pair (D, o) with D a dataspace and o an object in D . Let $K \subseteq \mathcal{K}$ and $L \subseteq \mathcal{L}$ be sets of keywords and link conditions, respectively. Two pointed dataspaces (D, o) and (D', o') are *n-bisimilar*, or *n-bisimulation equivalent*, under K and L , denoted $(D, o) \rightleftharpoons_n^{K,L} (D', o')$, if $o \simeq_K o'$ and the following conditions hold for $n > 0$:

Forth. For any $\lambda \in L$, if $\lambda(o, p)$ for some $p \in D$, then there is some $p' \in D'$ such that $\lambda(o', p')$ and $(D, p) \rightleftharpoons_{n-1}^{K,L} (D', p')$.

Back. For any $\lambda \in L$, if $\lambda(o', p')$ for some $p' \in D'$, then there is some $p \in D$ such that $\lambda(o, p)$ and $(D, p) \rightleftharpoons_{n-1}^{K,L} (D', p')$.

Example 17. For example, if $K = \{\star, a\}$; $L = \{\text{shares}\}$ (with *shares* as in Example 15); and D and D' are as follows:



then $(D, o) \rightleftharpoons_n^{K,L} (D', o')$ and $(D, p) \rightleftharpoons_n^{K,L} (D', p')$ for any n . Indeed, clearly $o \simeq_K o'$ and $p \simeq_K p'$. To verify the forth property, *shares* links o in D only to o and p by which we can respond in D' by o' and p' respectively. Similarly, *shares* links p only to p and o , by which we can respond in D' by p' and o' respectively. The back property is verified similarly. In contrast, $(D, o) \not\rightleftharpoons_n^{K,L} (D', p')$ for any n since $o \not\simeq_K p'$.

Definition 18. Let \mathcal{D} be a class of dataspace. A query q is $\rightleftharpoons_n^{K,L}$ -invariant on \mathcal{D} if for any $(D, o) \rightleftharpoons_n^{K,L} (D', o')$ with $D, D' \in \mathcal{D}$ we have $o \in q(D) \Leftrightarrow o' \in q(D')$. A query is $\rightleftharpoons_n^{K,L}$ -invariant if it is $\rightleftharpoons_n^{K,L}$ -invariant on the entire class of all dataspace.

We then obtain the following lemma known from the model theory of modal logic [Goranko and Otto 2007, Theorem 32]. Hereto, first define the nesting depth of link operators in an expression e , denoted by $\text{depth}(e)$, as follows:

$$\begin{aligned} \text{depth}(k) &:= 0 \\ \text{depth}(e_1 \text{ and } e_2) &:= \max(\text{depth}(e_1), \text{depth}(e_2)) \\ \text{depth}(e_1 \text{ or } e_2) &:= \max(\text{depth}(e_1), \text{depth}(e_2)) \\ \text{depth}(e_1 \text{ except } e_2) &:= \max(\text{depth}(e_1), \text{depth}(e_2)) \\ \text{depth}(\text{link}(\lambda) e) &:= 1 + \text{depth}(e) \end{aligned}$$

LEMMA 19. Let K be a finite nonempty set of keywords and let L be a finite nonempty set of link conditions. Then the following are equivalent for any search query q and any class \mathcal{D} of dataspace:

- (1) q is definable on \mathcal{D} in $\text{ASL}(K \cup \{\star\}, L)$ by an expression e with nesting depth of link operators at most n ;

12 .

(2) q is $\rightleftharpoons_n^{K,L}$ -invariant on \mathcal{D} .

PROOF. (1 \rightarrow 2) Let e be an arbitrary expression in $\text{ASL}(K \cup \{\star\}, L)$ and let $n = \text{depth}(e)$ be the nesting depth of link constructs in e . It is readily verified by induction on e that the query defined by e is $\rightleftharpoons_n^{K,L}$ -invariant on the class of all dataspace, and therefore also $\rightleftharpoons_n^{K,L}$ -invariant on any subclass \mathcal{D} .

(2 \rightarrow 1) The crucial step in proving the converse direction from (2) to (1) is the construction, for every pointed dataspace (D, o) , of a so-called *characteristic expression* $e_{(D,o)}^n$ in $\text{ASL}(K \cup \{\star\}, L)$, with the property that for all pointed dataspace (D', o') ,

$$o' \in e_{(D,o)}^n(D') \Leftrightarrow (D', o') \rightleftharpoons_n^{K,L} (D, o). \quad (1)$$

This construction proceeds by induction on n , for all (D, o) simultaneously.

—When $n = 0$, $e_{(D,o)}^0$ is the $\text{BSL}(K \cup \{\star\})$ expression

$$e_{(D,o)}^0 := (k \text{ and } \dots \text{ and } k') \text{ except } (l \text{ or } \dots \text{ or } l')$$

where $\{k, \dots, k'\}$ are all keywords in K satisfied by o and $\{l, \dots, l'\}$ are all keywords in K not satisfied by o . If o does not satisfy any keyword in K then $e_{(D,o)}^0$ is $\star \text{ except } (l \text{ or } \dots \text{ or } l')$.

—Inductively, let $\text{link}[\lambda]e$ abbreviate the expression $\star \text{ except link}[\lambda](\star \text{ except } e)$ that returns all objects that are only linked by λ to objects returned by e . Then $e_{(D,o)}^{n+1}$ is the $\text{ASL}(K \cup \{\star\}, L)$ expression with link nesting depth $n + 1$ given by

$$e_{(D,o)}^{n+1} := e_{(D,o)}^0 \text{ and } \bigwedge_{\lambda \in L} \left(\underbrace{\bigwedge_{\substack{p \in D \\ \lambda(o,p)}} \text{link}[\lambda] e_{(D,p)}^n}_{\text{forth}} \wedge \underbrace{\text{link}[\lambda] \bigvee_{\substack{p \in D \\ \lambda(o,p)}} e_{(D,p)}^n}_{\text{back}} \right).$$

As such, if $\{(D_1, o_1), (D_2, o_2), \dots\}$ is the set of all pointed dataspace (D, o) with $D \in \mathcal{D}$ and $o \in q(D)$, q is defined on \mathcal{D} by the (possibly infinite) expression

$$e_{(D_1,o_1)}^n \text{ or } e_{(D_2,o_2)}^n \text{ or } \dots$$

Since K and L are finite, however, there are only a finite number of characteristic expressions $e_{(D,o)}^n$ up to logical equivalence. Therefore, the above expression can always be made finite, as desired. \square

Like Theorem 8 for BSL, Lemma 19 provides a tool to show certain queries are (un)definable in $\text{ASL}(K, L)$ for various K and L —a tool that we will repeatedly employ in the following sections.

2.3 On searching by similarity link conditions and the semijoin algebra

In their full generality, link conditions can be arbitrarily complex, and need not adhere to any reasonable notion of “associative search”. Consider, for instance, the link condition $\text{subset} := \{(o, o') \mid o \subseteq o'\}$. Then $\text{link}[\text{subset}] \text{ICDT 2009}$ returns all objects o for which there exists some o' that mentions ICDT 2009 and $o \subseteq o'$. This corresponds to a set join in the relational algebra, and has more to do with full-blown querying than searching.

Moreover, a flexible search query language should allow the link conditions to be expressed *within* the language itself.

For these reasons we introduce next a restricted set of link conditions, the *simlinks* that intuitively link two objects o and o' by (1) searching *within* o using a keyword k ; (2) searching *within* o' using a keyword l ; and (3) comparing the two search results and requiring that they contain a pair of similar items.

Definition 20. A *similarity relation* (or *simrel* for short) is a binary relation \sim on items, $\sim \subseteq \mathcal{I} \times \mathcal{I}$. Let k and l be keywords and let \sim be a simrel. Then the expression $k \sim l$ is called a *simlink* and can be used as a link condition with the following semantics:

$$k \sim l := \{(o, o') \in \mathcal{O} \times \mathcal{O} \mid \exists i \in o \cap k : \exists j \in o' \cap l : i \sim j\}.$$

Henceforth we write \mathcal{S} for the set of all simrels.

Example 21. For a simple example on document spaces consider the simrel *soviet* between words such that w_1 *soviet* w_2 if w_1 is a location name (city name, street name) from the Soviet era, and w_2 is the corresponding post-Soviet name. For example, Leningrad *soviet* St. Petersburg. Then the expression

$$\text{link}(\star \text{ soviet } \star)(\text{ICDT})$$

retrieves all documents containing Soviet versions of location names mentioned in documents about ICDT, and the expression

$$\text{link}(\star \text{ soviet Peterhof})(\text{ICDT})$$

retrieves all documents containing the Soviet names of the Russian town Peterhof, if mentioned in documents about ICDT. In Section 4.2 we will see more examples of simlinks.

For the remainder of this article, we will always use the language ASL with simlinks as link conditions. In particular, if $K \subseteq \mathcal{K}$ and $S \subseteq \mathcal{S}$ then we write $\text{ASL}(K, S)$ for the fragment of ASL expressions that use only keywords in K and link conditions in $\{k \sim l \mid k, l \in K \text{ and } \sim \in S\}$.

Relationship to the semijoin algebra. Search queries are a special kind of database queries. It is therefore natural to ask how the language ASL (with simlinks) compares to more standard query languages. Observing that the link operator is a kind of semijoin, a comparison with the semijoin algebra seems a good approach to this question. (Recall that semijoin algebra is the version of the relation algebra where the join operator is replaced by the semijoin operator [Leinders et al. 2005; Leinders and Van den Bussche 2007].)

Since the relational algebra works on relations instead of dataspace, we need a relational representation of a dataspace:

Definition 22. For each dataspace D let $\text{rep}(D)$ be the binary relation $\{(o, i) \mid o \in D \text{ and } i \in o\}$ over the relation schema $\{\text{id}, \text{item}\}$.

The objects in the *id*-column of this relation are regarded as object identifiers; the semijoin algebra will not be able to peek inside the objects in any other way than working with the *item* column. An example is shown in Figure 2.

14

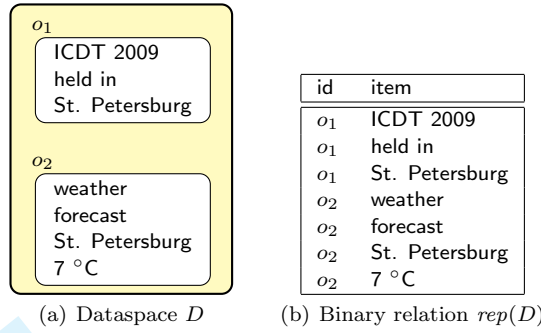


Fig. 2. A dataspace and its relational representation.

$$\begin{array}{c}
 \frac{}{T: \{id, item\}} \quad \frac{E: \Sigma \quad item \in \Sigma \quad k \in \mathcal{K}}{\sigma_k(E): \Sigma} \quad \frac{E: \Sigma \quad \Delta \subseteq \Sigma}{\pi_\Delta(E): \Delta} \quad \frac{E_1: \Sigma \quad E_2: \Sigma}{E_1 \cup E_2: \Sigma} \\
 \\
 \frac{E_1: \Sigma \quad E_2: \Sigma}{E_1 - E_2: \Sigma} \quad \frac{E_1: \Sigma_1 \quad E_2: \Sigma_2 \quad \Sigma_1 \cap \Sigma_2 = \{id\}}{E_1 \bowtie E_2: \Sigma_1} \\
 \\
 \frac{E_1: \Sigma_1 \quad E_2: \Sigma_2 \quad \Sigma_1 \cap \Sigma_2 = \{item\} \quad \sim \in \mathcal{S}}{E_1 \bowtie_{\sim} E_2: \Sigma_1} \quad \frac{E_1: \{id, item\} \quad E_2: \{id, item\} \quad \sim \in \mathcal{S}}{E_1 \bowtie_{=, \sim} E_2: \{id, item\}}
 \end{array}$$

Fig. 3. Syntax of SA.

In the version of the semijoin algebra we are using, keywords on items are used as selection conditions, and simrels on items and equality on objects (viewed as id's) are used as semijoin conditions. The relation symbol T (for 'Table') stands for the binary representation of the input dataspace. Every expression has an output schema that is either empty, the schema $\{id, item\}$ itself, or one of the unary schemas $\{id\}$ or $\{item\}$ (we do not have renaming). Expressions must be well-typed, e.g., we only allow the union of two relations over the same schema. The full syntax of semijoin algebra (abbreviated SA) expressions, together with the derivation of their output schemas, is given in Figure 3. There, the notation $E: \Sigma$ denotes that E is a legal expression with output schema Σ .

The semantics of projection π , union \cup and difference $-$ is well known; the semantics of the semijoin operators is as follows.

$$\begin{aligned}
 \sigma_k(R) &:= \{t \in R \mid t(item) \in k\} \\
 R_1 \bowtie R_2 &:= \{t_1 \in R_1 \mid \exists t_2 \in R_2 : t_2(id) = t_1(id)\} \\
 R_1 \bowtie_{\sim} R_2 &:= \{t_1 \in R_1 \mid \exists t_2 \in R_2 : t_1(item) \sim t_2(item)\} \\
 R_1 \bowtie_{=, \sim} R_2 &:= \{(o, i) \in R_1 \mid \exists j : (o, j) \in R_2 \text{ and } i \sim j\}
 \end{aligned}$$

So, \bowtie is a normal natural semijoin on the common id attribute; \bowtie_{\sim} is a \sim -semijoin on the common item attribute; and $\bowtie_{=, \sim}$ is a combination of the two.

Definition 23. A search query q is *definable in SA* if there exists an SA expression $E: \Sigma$ with $id \in \Sigma$ such that $q(D) = \pi_{id}(E)(rep(D))$, for all dataspace D . We also say that E *defines* q in this case.

The following is now expected:

PROPOSITION 24. *Let $K \subseteq \mathcal{K}$ and $S \subseteq \mathcal{S}$. Every search query definable in $\text{ASL}(K, S)$ is also definable in $\text{SA}(K, S)$.*

PROOF. Here is a straight syntactic translation:

$$\begin{aligned} \text{SA}[k] &:= \sigma_k(T) \\ \text{SA}[e_1 \text{ and } e_2] &:= \pi_{\text{id}}\text{SA}[e_1] - (\pi_{\text{id}}\text{SA}[e_1] - \pi_{\text{id}}\text{SA}[e_2]) \\ \text{SA}[e_1 \text{ or } e_2] &:= \pi_{\text{id}}\text{SA}[e_1] \cup \pi_{\text{id}}\text{SA}[e_2] \\ \text{SA}[e_1 \text{ except } e_2] &:= \pi_{\text{id}}\text{SA}[e_1] - \pi_{\text{id}}\text{SA}[e_2] \\ \text{SA}[\text{link}(k \sim l) e] &:= \pi_{\text{id}}(\text{SA}[k] \ltimes \pi_{\text{item}}(\text{SA}[l] \ltimes \pi_{\text{id}}\text{SA}[e])) \quad \square \end{aligned}$$

Is the converse true as well? The answer is no, and we will perform a thorough analysis of the situation, with the goal of arriving at a well-defined fragment of SA that is equivalent to ASL .

The first observation is that SA can express boolean combinations of keywords. For example, the SA expression $\pi_{\text{id}}(T - \sigma_k(T))$ defines the negated keyword $\neg k$, i.e., the query $q(D) = \{o \in D \mid \exists i \in o : i \notin k\}$. We have already seen in Example 10 that negated keywords, and more generally, boolean combinations of keywords, are not definable in ASL . This is easily repaired by closing the keywords under the boolean operators.

Definition 25. For $K \subseteq \mathcal{K}$, let K^* be the smallest set of keywords containing K that is closed under the operators \neg and \vee , where

$$\begin{aligned} \neg k &:= \{i \in \mathcal{I} \mid i \notin k\} \text{ and} \\ k \vee l &:= \{i \in \mathcal{I} \mid i \in k \text{ or } i \in l\}. \end{aligned}$$

Note that whenever K is non-empty, K^* necessarily includes the wildcard keyword \star that matches all items since $\star \equiv k \vee \neg k$.

We can now revise Proposition 24 as follows:

PROPOSITION 24 REVISED. *Let $K \subseteq \mathcal{K}$ and $S \subseteq \mathcal{S}$. Every search query definable in $\text{ASL}(K^*, S)$ is also definable in $\text{SA}(K, S)$.*

PROOF. The syntactic translation from the proof of Proposition 24 can be extended as follows: (φ and ψ stand for boolean combinations of keywords)

$$\begin{aligned} \text{SA}[k] &:= \sigma_k(T) \\ \text{SA}[\neg \varphi] &:= T - \text{SA}[\varphi] \\ \text{SA}[\varphi \vee \psi] &:= \text{SA}[\varphi] \cup \text{SA}[\psi] \quad \square \end{aligned}$$

The next proposition points at a number of distinct query constructions definable in SA but not in ASL :

PROPOSITION 26. *None of the following SA queries is definable in ASL , even*

16

over the boolean closure of keywords:

$$\begin{aligned}
 E_1 &:= \pi_{\text{id}}(\sigma_k(T) \bowtie_{=, \sim} \sigma_l(T)) \\
 E_2 &:= \pi_{\text{id}}(T - T \bowtie_{\sim} \pi_{\text{item}} \sigma_k(T)) \\
 E_3 &:= \pi_{\text{id}}(\sigma_k(T) \bowtie_{\sim} \pi_{\text{item}}(\sigma_l(T) \bowtie_{\sim} \pi_{\text{item}} \sigma_m(T))) \\
 E_4 &:= \pi_{\text{id}}(T \bowtie_{\sim} (\pi_{\text{item}}(T) - \pi_{\text{item}}(T \bowtie_{\sim} \pi_{\text{id}} \sigma_k(T))))
 \end{aligned}$$

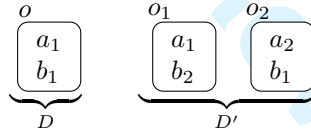
More precisely, for each of these queries we can pick the keywords k, l, m , and $\text{simrel} \sim$ such that the query is not definable in $\text{ASL}(K^*, S)$ with $K = \{k, l, m\}$ and $S = \{\sim\}$.

Note that, referring to the items $i \in k$ as k -items, the above four expressions define the following queries:

- (1) Retrieve all objects containing a k -item i and an l -item j such that $i \sim j$.
- (2) Retrieve all objects containing an item that is not \sim to any k -item in the dataspace.
- (3) Retrieve all objects containing a k -item that is \sim to some l -item j in the dataspace, and j itself is \sim to some m -item in the dataspace.
- (4) Retrieve all objects containing an item i that is \sim to an item that is not present in any object containing a k -item.

These queries are shown to be undefinable in $\text{ASL}(K^*, S)$ by showing that they are not bisimulation invariant, then invoking Lemma 19.

PROOF OF PROPOSITION 26. For the query q_1 expressed by E_1 , we pick $k = \{a_1, a_2\}$ and $l = \{b_1, b_2\}$, with $a_1 \sim b_1$ and $a_2 \sim b_2$. Now consider



Let $K = \{k, l\}$ and let L be the set of all link conditions $\{p \sim p' \mid p, p' \in K^*\}$. Then (D, o) , (D', o_1) , and (D', o_2) are all bisimilar (n -bisimilar for any n) under K^* and L . The verification of this claim proceeds by induction on n :

—Case $n = 0$. Here it suffices to show that $o \simeq_{K^*} o_1 \simeq_{K^*} o_2$, which readily follows from the observation that, writing $k \wedge l$ for $\neg(\neg k \vee \neg l)$,

$$K^* = \{k, l, \neg k, \neg l, k \vee l, k \vee \neg l, \neg k \vee l, k \wedge l, k \wedge \neg l, \neg k \wedge l, \emptyset, \star\}.$$

By construction, o, o_1 , and o_2 all satisfy the same keywords in this set.

—Case $n > 0$. Suppose that $(D, o) \rightleftharpoons_{n-1}^{K^*, L} (D', o_1) \rightleftharpoons_{n-1}^{K^*, L} (D', o_2)$. We only need to show that $(D, o) \rightleftharpoons_n^{K^*, L} (D', o_1)$ and $(D, o) \rightleftharpoons_n^{K^*, L} (D', o_2)$, since $(D', o_1) \rightleftharpoons_n^{K^*, L} (D', o_2)$ then follows by transitivity of $\rightleftharpoons_n^{K^*, L}$. We check the forth and back properties.

(Forth). Let $\lambda = p \sim p'$ be an arbitrary simlink in L , and let o' be an arbitrary object in D with $\lambda(o, o')$. We need to show that there exists o'_1 and o'_2 in D' such that (1) $\lambda(o_1, o'_1)$ and $(D, o) \rightleftharpoons_{n-1}^{K^*, L} (D', o'_1)$; and (2) $\lambda(o_2, o'_2)$ and $(D, o) \rightleftharpoons_{n-1}^{K^*, L} (D', o'_2)$.

(D', o'_2) . Hereto, first observe that, since D is a singleton, we must have $o' = o$. Moreover, since the only items in o that are related by \sim are a_1 and b_1 , $\lambda(o, o)$ implies $a_1 \in p$, $b_1 \in p'$, and $a_1 \sim b_1$.

- (1) Then, since $a_1 \in o_1$ and $a_2 \in o_2$, certainly also $\lambda(o_1, o_2)$. Furthermore, $(D, o) \rightleftharpoons_{n-1}^{K^*, L} (D', o_2)$ by the induction hypothesis. As such, the forth property holds for (D', o_1) .
- (2) To see that the forth property also holds for (D', o_2) , first observe that for any of the keywords in K^* above we have $a_1 \in p$ iff $a_2 \in p$ and $b_1 \in p'$ iff $b_2 \in p'$. Hence, since $a_1 \in p$ and $b_1 \in p'$, also $a_2 \in p$ and $b_2 \in p'$. Thus, $a_2 \in o_2 \cap p$, $b_2 \in o_1 \cap p$, and $a_2 \sim b_2$. Hence $\lambda(o_2, o_1)$. Furthermore, $(D, o) \rightleftharpoons_{n-1}^{K^*, L} (D', o_1)$ by the induction hypothesis. As such, the forth property holds for (D', o_2) .

(Back). Similar to Forth.

So, in particular (D, o) and (D', o_1) are bisimilar under K^* and L . Yet, $o \in q_1(D)$ whereas $o_1 \notin q_1(D')$. Hence q_1 is not definable in $\text{ASL}(K^*, S)$ by Lemma 19.

For the query q_2 expressed by E_2 , we use items a , a' , and b with $k = \{b\}$ and $a' \sim b$. Now consider $D = \{o_1, o_2\}$ and $D' = \{o'_1, o_2\}$ with $o_1 = \{a, a'\}$; $o_2 = \{b\}$; and $o'_1 = \{a'\}$. Let $K = \{k, l\}$ and let L be the set of all link conditions $\{p \sim p' \mid p, p' \in K^*\}$. Then (D, o_1) is bisimilar to (D', o'_1) under K^* and L (as can again be verified by induction), but $o_1 \in q_2(D)$ whereas $o'_1 \notin q_2(D')$.

For the query q_3 expressed by E_3 , we pick $k = \{a\}$; $l = \{b, b_1, b_2\}$; and $m = \{c\}$, with $a \sim b$, $b \sim c$, $a \sim b_1$, and $b_2 \sim c$. Now consider $D = \{o_1, o_2\}$ and $D' = \{o_1, o_2\}$ with $o_1 = \{a\}$; $o_2 = \{b, c\}$; and $o'_2 = \{b_1, b_2, c\}$. Let $K = \{k, l, m\}$ and let L be the set of all link conditions $\{p \sim p' \mid p, p' \in K^*\}$. Then (D, o_1) is bisimilar to (D', o_1) under K^* and L , but $o_1 \in q_3(D)$ whereas $o_1 \notin q_3(D')$.

For the query q_4 expressed by E_4 , we use items a , b and c , with $k = \{a\}$, and \sim interpreted as equality. Now consider $D = \{o_1, o_2\}$ and $D' = \{o'_1, o_2\}$ with $o_1 = \{b, c\}$; $o_2 = \{a, b\}$; and $o'_1 = \{b\}$. Let $K = \{k\}$ and let L be the set of all link conditions $\{p \sim p' \mid p, p' \in K^*\}$. Then (D, o_1) is bisimilar to (D', o'_1) under K^* and L , but $o_1 \in q_4(D)$ whereas $o'_1 \notin q_4(D')$. \square

Proposition 26 can inspire us to restrict SA so as to obtain a fragment equivalent to ASL . E_1 suggests that we should banish the combined semijoin $\bowtie_{=, \sim}$. E_2 suggests that we should not allow unrestricted use of the result of a \bowtie_{\sim} semijoin; we should project the result. But E_3 shows we should not project on $\{\text{item}\}$, so we conclude we must always project the result of \bowtie_{\sim} on $\{\text{id}\}$. Moreover, E_4 suggests that we should not allow projection on $\{\text{item}\}$ altogether, except of course to allow for a semijoin \bowtie_{\sim} to be applied. We thus arrive at the following fragment of SA , which we denote by $\text{SA}^{\text{search}}$:

Definition 27. The fragment $\text{SA}^{\text{search}}$ of SA is defined by the following rules:

- The operators $\bowtie_{=, \sim}$, $\pi_{\{\text{item}\}}$ and π_{\emptyset} are disallowed.
- The rule for \bowtie_{\sim} is changed as follows:

$$\frac{E_1 : \{\text{id}, \text{item}\} \quad E_2 : \{\text{id}, \text{item}\} \quad \sim \in \mathcal{S}}{\pi_{\{\text{id}\}}(E_1 \bowtie_{\sim} \pi_{\{\text{item}\}} E_2) : \{\text{id}\}}$$

18

With these restrictions, we are ready to establish the connection between associative search and the semijoin algebra:

THEOREM 28. *Let $K \subseteq \mathcal{K}$ and $S \subseteq \mathcal{S}$. A search query is definable in $\text{ASL}(K^*, S)$ if and only if it is definable in $\text{SA}^{\text{search}}(K, S)$.*

PROOF. The only-if direction follows from the observation that the translation from ASL to SA given in the proof of Proposition 24 stays within the fragment $\text{SA}^{\text{search}}$.

The if-direction is also proven by a translation, but a complication here is the translation of subexpressions with output schema $\{\text{id}, \text{item}\}$, since such intermediate results are not directly representable by the result of a search query (which can only return id's).

Formally, for each $\text{SA}^{\text{search}}$ expression E we construct a finite set χ_E of pairs (e, k) , with e an $\text{ASL}(K^*, S)$ expression and $k \in K^*$, such that for all D :

—if the output schema of E is $\{\text{id}, \text{item}\}$, then $E(\text{rep}(D))$ equals

$$\bigcup_{(e,k) \in \chi_E} \{(o, i) \mid o \in e(D), i \in o \cap k\};$$

—if the output schema of E is $\{\text{id}\}$, then $E(\text{rep}(D))$ equals

$$\bigcup_{(e,k) \in \chi_E} \{o \in e(D) \mid o \models k\}.$$

Of course the global SA expression E expressing the search query q has output schema Σ with $\{\text{id}\} \subseteq \Sigma \subseteq \{\text{id}, \text{item}\}$, and thus q is defined in $\text{ASL}(K^*, S)$ by

$$(e_1 \text{ and } k_1) \text{ or } \dots \text{ or } (e_n \text{ and } k_n),$$

where $\chi_E = \{(e_1, k_1), \dots, (e_n, k_n)\}$.

We construct χ_E by induction as follows:

$$\chi_T := \{(\star, \star)\}$$

$$\chi_{\sigma_l(E)} := \{(e, k \wedge l) \mid (e, k) \in \chi_E\}$$

$$\chi_{\pi_{\text{id}}(E)} := \{(e \text{ and } k, \star) \mid (e, k) \in \chi_E\}$$

$$\chi_{E_1 \cup E_2} := \chi_{E_1} \cup \chi_{E_2}$$

$$\chi_{\pi_{\text{id}}(E_1 \bowtie \sim \pi_{\text{item}} E_2)} := \{(e_1 \text{ and link}(k_1 \sim k_2) e_2, \star) \mid (e_1, k_1) \in \chi_{E_1}, (e_2, k_2) \in \chi_{E_2}\}$$

$$\chi_{E_1 \bowtie E_2} := \{(e_1 \text{ and } e_2 \text{ and } k_2, k_1) \mid (e_1, k_1) \in \chi_{E_1}, (e_2, k_2) \in \chi_{E_2}\}$$

For $E = E_1 - E_2$ we first need some terminology and notation. Let $\chi_{E_2} = \{(e_1, k_1), \dots, (e_n, k_n)\}$. A *splitter* of E_2 is a pair (S^+, S^-) of disjoint subsets of $\{1, \dots, n\}$ such that $S^+ \cup S^- = \{1, \dots, n\}$. Then

$$\chi_{E_1 - E_2} := \{(e \text{ except } (e_{i_1} \text{ or } \dots \text{ or } e_{i_r}), k \wedge \neg(k_{j_1} \vee \dots \vee k_{j_s})) \mid (e, k) \in \chi_{E_1}, \\ (\{i_1, \dots, i_r\}, \{j_1, \dots, j_s\}) \text{ a splitter of } E_2\}.$$

The correctness of χ_E is proven by induction on E . We only show the reasoning for $E = E_1 - E_2$; the other cases are similar. Suppose that E_1 and E_2 both have

schema $\{\text{id}, \text{item}\}$. Then, by induction hypothesis,

$$\begin{aligned}
 & (o, i) \in E(\text{rep}(D)) \\
 & \Leftrightarrow (o, i) \in E_1(\text{rep}(D)) \text{ and } (o, i) \notin E_2(\text{rep}(D)) \\
 & \Leftrightarrow \text{there exists } (e, k) \in \chi_{E_1} : o \in e(D), i \in o \cap k \\
 & \quad \text{and for all } 1 \leq m \leq n : o \notin e_m(D) \text{ or } i \notin o \cap k_m \\
 & \quad \text{where } \chi_{E_2} = \{(e_1, k_1), \dots, (e_n, k_n)\} \\
 & \Leftrightarrow \text{there exists } (e, k) \in \chi_{E_1} : o \in e(D), i \in o \cap k \\
 & \quad \text{and a splitter } (\{i_1, \dots, i_r\}, \{j_1, \dots, j_s\}) \text{ for } E_2 \\
 & \quad \text{such that } o \notin e_m(D) \text{ for } m \in \{i_1, \dots, i_r\} \\
 & \quad \text{and } i \notin o \cap k_m \text{ for } m \in \{j_1, \dots, j_s\} \\
 & \Leftrightarrow (o, i) \in \bigcup_{(e, k) \in \chi_{E_1 - E_2}} \{(o', i') \mid o' \in e(D), i' \in o' \cap k\}
 \end{aligned}$$

The reasoning is similar when E_1 and E_2 both have schema $\{\text{id}\}$. \square

We conclude this section with a remark concerning conjunctions in join conditions. We have defined simlinks as link conditions involving just a single condition $k \sim l$. But link conditions involving a conjunction of two or more such conditions are also very natural. For example, the query $\text{link}((k_1 \sim_1 l_1) \wedge (k_2 \sim_2 l_2))(e)$, applied to a dataspace D , would return those objects $o \in D$ for which there exists an object o' in $e(D)$ such that there exists a k_1 -item $i \in o$ and an l_1 -item $j \in o'$ such that $i \sim_1 j$, and there also exists a k_2 -item $i \in o$ (not necessarily the same i) and an l_2 -item $j \in o'$ (not necessarily the same j) such that $i \sim_2 j$.

Such conjunctive join conditions, however, can bring us outside the semijoin algebra. We can see this using the following example, which already gives a taste of what will be served in the next two sections.

Consider a chain of stores selling wine and cheese, with some stores located in France and others in Italy. We can model each store as a set of items: one item simply indicates the country France or Italy; other items are of the form (wine, v) listing all kinds of wine v sold by the store; and the remaining items are of the form (cheese, x) listing all kinds of cheese x sold by the store. Define the simrel **samewine** as the set of all identical pairs of the form $((\text{wine}, v), (\text{wine}, v))$ for arbitrary wines v ; define analogously the simrel **samecheese** as the set of all identical pairs of the form $((\text{cheese}, x), (\text{cheese}, x))$ for arbitrary cheeses x .

Now consider the query “list all stores in France for which there exists a store in Italy selling some same wine and some same cheese”, formally expressed as follows:

$$\text{france and link}((\star \text{ samewine } \star) \wedge (\star \text{ samecheese } \star))(\text{italy})$$

This query is not expressible in the semijoin algebra. This can be proven formally as follows. Consider relational databases over two unary relations France and Italy, and two binary relations Wine and Cheese. The relation France (Italy) holds the ids of stores located in France (Italy), and the relation Wine (Cheese) holds the pairs (o, u) such that wine (cheese) u is sold in store o . Consider the two specific database instances D_1 and D_2 , with the following content for D_1 :

20

France	Italy	Wine	Cheese
1	2	1 bordeaux 2 bordeaux	1 roquefort 2 roquefort

So, in D_1 , store 1 is french and store 2 is italian and they offer exactly the same products (one wine and one cheese). Obviously then, the query applied to D_1 will return store 1. The content for D_2 is the following:

France	Italy	Wine	Cheese
3	4 5 6	3 bordeaux 4 bordeaux 5 sangiovese 6 sangiovese	3 parmesan 4 roquefort 5 parmesan 6 roquefort

Observe that our query has the empty answer on D_2 . It can be verified, however, that $(D_1, 1)$ and $(D_2, 3)$ are *guarded bisimilar* and hence cannot be distinguished in the semijoin algebra [Leinders et al. 2005; Leinders and Van den Bussche 2007]. Consequently, any semijoin algebra query that returns 1 on D_1 will return 3 on D_2 . Since our query does not have this property, it is not expressible in the semijoin algebra.

It seems very well possible to define a fragment of the relational algebra appropriate for a generalization of Theorem 28 to conjunctions in link conditions. We leave this (involved) exercise to the interested reader.

3. APPLICATION 1: SEARCHING CLASSICAL RELATIONS

The theory presented in Section 2 holds for any concrete interpretation of \mathcal{I} . As an illustration of how this general theory can be applied, let us fix an infinite domain \mathcal{V} of values and a relation schema Σ , i.e., a finite set of attributes. We can use the set of attribute–value pairs $\Sigma \times \mathcal{V}$ as our set \mathcal{I} of items. Note that a tuple $t : \Sigma \rightarrow \mathcal{V}$ is then an object over \mathcal{I} , and a finite relation R over Σ is a dataspace over \mathcal{I} . Also note that conversely, not all objects over \mathcal{I} are proper tuples, because in a proper tuple every attribute occurs, with only one value. Hence, not all dataspaces in this setting are proper relations. Let us denote the subclass of dataspaces over $\Sigma \times \mathcal{V}$ that are finite relations over Σ by \mathcal{R} .

3.1 Boolean Search

To search \mathcal{R} we consider here the set of keywords

$$\text{LIT} := \{\star\} \cup \{(a, v) \mid a \in \Sigma, v \in \mathcal{V}\}$$

consisting of the wildcard \star and the literal keywords over $\mathcal{I} = \Sigma \times \mathcal{V}$. Note that a literal keyword (a, v) corresponds to the relational selection operator $\sigma_{a=v}$. Indeed, for all $R \in \mathcal{R}$ Definition 4 specializes to

$$(a, v)(R) = \{t \in R \mid t(a) = v\}$$

since R consists only of tuples over Σ . So, the language $\text{BSL}(\text{LIT})$ evaluated over relations in \mathcal{R} corresponds to the fragment of the relational algebra consisting of just the three operators of constant selection, union, and difference.

It turns out that our notion of K -distinguishing search queries corresponds in this setting to the so-called fully C -generic queries. Fully generic queries without

constants have been introduced and investigated (in a much richer setting than mere selection queries) by Beeri, Milo and Ta-Shma [Beeri et al. 1996; 1997]; here we consider the version with constants.

Definition 29 (C-epimorphism). Let $C \subset \mathcal{V}$ be a finite set of constants. A C -epimorphism is a mapping $f: \mathcal{V} \rightarrow \mathcal{V}$ such that both $f|_C$ and $f^{-1}|_C$ are the identity on C . We extend $f: \mathcal{V} \rightarrow \mathcal{V}$ to items, objects, and dataspace in the canonical, pointwise manner:

$$\begin{aligned} f(a, v) &:= (a, f(v)), \\ f(o) &:= \{f(a, v) \mid (a, v) \in o\}, \\ f(R) &:= \{f(o) \mid o \in R\}. \end{aligned}$$

Now a search query q is said to be *fully C-generic* (on the class \mathcal{R}) if for any relation R over Σ and any C -epimorphism f , we have $q(f(R)) = f(q(R))$.

We establish:

PROPOSITION 30. *Let C be a finite set of constants, and let $K = \{\{(a, v)\} \mid a \in \Sigma, v \in C\}$. The following are equivalent for any search query q :*

- (1) q is K -distinguishing on the class \mathcal{R}
- (2) q is additive and fully C -generic on \mathcal{R} .

PROOF. (1 \rightarrow 2) Since q is K -distinguishing on \mathcal{R} , we know that for every $R \in \mathcal{R}$ and every $t \in R$ we have $t \in q(R)$ iff $t \in q(\{t\})$. This readily implies additivity. It remains to show that q is fully C -generic on \mathcal{R} . Hereto, fix $R \in \mathcal{R}$ arbitrarily and let f be a C -epimorphism. Observe that $t \simeq_K f(t)$ for all $t \in R$. Indeed, let $k = \{(a, v)\}$ be an arbitrary keyword in K with $a \in \Sigma$ and $v \in C$. If $t \models k$, then $t(a) = v$. Since $f|_C$ is the identity on C , $f(v) = v$ and hence $f(t)(a) = v$. Thus $f(t) \models k$. Conversely, if $f(t) \models k$ then $f(t)(a) = v$. Since $f^{-1}|_C$ is the identity on C , also $t(a) = v$. Therefore, $t \models k$. As such, t satisfies the same keywords in K as $f(t)$. Since q is K -distinguishing, we hence have $t \in q(R) \Leftrightarrow f(t) \in q(f(R))$, which implies $q(f(R)) = f(q(R))$, as desired.

(2 \rightarrow 1) Let R and R' be arbitrary relations in \mathcal{R} and let $t \in R$ and $t' \in R'$ such that $t \simeq_K t'$. We need to show that $t \in q(R) \Leftrightarrow t' \in q(R')$. Hereto, first observe that for all $a \in \Sigma$, if $t(a) \in C$ then $t'(a) = t(a)$ (since $t \simeq_K t'$) and conversely if $t'(a) \in C$ then $t(a) = t'(a)$ (for the same reason). Then construct the tuple u such that for all $a \in \Sigma$

- if $t(a) \in C$ then $u(a) = t(a) = t'(a)$; and
- if $t(a) \notin C$ then $u(a) \notin C$ and $u(a) \neq u(b)$ for every $b \in \Sigma$ distinct from a .

It is readily verified that there exists C -epimorphism f and f' such that $f(u) = t$

22 .

and $f'(u) = t'$. Now reason as follows:

$$\begin{aligned}
 t \in q(R) &\Leftrightarrow t \in q(\{t\}) \\
 &\Leftrightarrow f(u) \in q(\{f(u)\}) \\
 &\Leftrightarrow u \in q(\{u\}) \\
 &\Leftrightarrow f'(u) \in q(\{f'(u)\}) \\
 &\Leftrightarrow t' \in q(\{t'\}) \\
 &\Leftrightarrow t' \in q(R').
 \end{aligned}$$

Hence the proposition. \square

From the above proposition; Corollary 12; and the fact that BSL(LIT) corresponds to fragment of the relation algebra consisting of just the three operators of constant selection, union, and difference we then obtain:

COROLLARY 31. *A search query q on relations over Σ is definable in the relational algebra using only the operators constant selection, union, and difference, if and only if q is additive and fully C -generic for some finite set C of constants.*

The main purpose of this modest theorem is to illustrate that our general theory can be effectively applied and connected to earlier work. Moving towards associative search over classical relations, current systems implement languages approaching the full relational algebra (e.g., [Golenberg et al. 2008; Qin et al. 2009]), the expressive power of which is already well understood [Abiteboul et al. 1995].

4. APPLICATION 2: SEARCHING ATTRIBUTE-VALUE DATASPACES

We next apply the abstract theory of Section 2 to the concrete setting of attribute-value dataspaces as they have been investigated in the literature [Halevy et al. 2006; Dong and Halevy 2007; Dittrich and Vaz Salles 2006]. In this setting, items are attribute-value pairs as in the relational setting of Section 3, so $\mathcal{I} = \Sigma \times \mathcal{V}$. The important difference, however, is that the universe of attributes Σ can now be infinite; there is no longer a fixed finite relation schema. Of course each dataspaces is finite, so in each dataspaces just a finite number of attributes will occur, but these attributes can vary from dataspaces to dataspaces and even from object to object. Moreover, attributes can appear multiple times in the same object, with different values. So, an object really is just a nonempty finite set of items, without any restriction.

The dataspaces models in the literature [Halevy et al. 2006; Dong and Halevy 2007; Dittrich and Vaz Salles 2006] give each object an id, and naturally represent a dataspaces D as a set of triples $\{(o, a, v) \mid o \in D \text{ and } (a, v) \in o\}$. We argue that our view of a dataspaces as just a set of objects is equivalent. Indeed, we just showed how to go from our representation to the set of triples; if one wants to go in the converse direction, the only difficulty one may encounter is that there might be two object id's in the triple set with exactly the same set of associated (a, v) pairs. In that case we can add an explicit id attribute to the objects, so that the sets become distinct. Having explicit id attributes is also necessary when we need to represent links between objects based on ids. We will see an example of such linking later (Example 37).

4.1 Boolean Search

The next question is what keywords to use for searching attribute–value dataspace. Surely we need all literal keywords (a, v) (as in the classical relational case), so that we can formulate basic queries like ‘retrieve all persons who live in Belgium, like the beer Duvel, but do not like the beer Heineken’:¹

$$((\text{country: Belgium}) \text{ and } (\text{likes: Duvel})) \text{ except } (\text{likes: Heineken})$$

We also want negation separately on attributes and values: for example,

$$(\text{likes: } \neg \text{Heineken})$$

retrieves objects containing a value for attribute `likes` that is different from `Heineken`, and $(\neg \text{likes: Heineken})$ retrieves objects containing `Heineken` as the value of an attribute different from `likes`. Note that, since the set of attributes is not fixed, we cannot express the last example by a disjunction using all possible attributes other than `likes`. Similarly, we need wildcards on values as well as on attributes, so $(\star: \text{Belgium})$ retrieves all objects with value `Belgium` for some attribute. Finally, we need disjunctions such as $(\text{likes: } \neg(\text{Heineken} \vee \text{Budweiser}))$. (We only need negated disjunctions; positive disjunctions can already be expressed in BSL using `or`.)

To sum up, we propose the following system of keywords for attribute–value pairs.

Definition 32. First define the following abbreviations for $a \in \Sigma$, $v \in \mathcal{V}$, $A \subseteq \Sigma$, and $V \subseteq \mathcal{V}$:

$$\begin{aligned} (a: v) &:= \{(a, v)\} \\ (a: \neg V) &:= \{(a, v) \mid v \in \mathcal{V} - V\} \\ (\neg A: v) &:= \{(a, v) \mid a \in \Sigma - A\} \\ (\neg A: \neg V) &:= \{(a, v) \mid a \in \Sigma - A, v \in \mathcal{V} - V\}. \end{aligned}$$

Note that the wildcard \star is readily obtained by $(\neg \emptyset, \neg \emptyset)$. The set **AV** of *attribute–value keywords* is the set consisting of all keywords of the form

$$(a: v) \mid (a: \neg V) \mid (\neg A: v) \mid (\neg A: \neg V),$$

where $a \in \Sigma$, $v \in \mathcal{V}$, and $A \subseteq \Sigma$, $V \subseteq \mathcal{V}$ are finite.

A first indication of the flexibility of this keyword system is that we do not need to add boolean combinations:

PROPOSITION 33. *BSL(AV) is equivalent to BSL(AV*), where AV* denotes boolean closure of AV keywords (cf. Definition 25).*

PROOF. It is readily verified that every boolean combination of **AV** keywords amounts to a disjunction of **AV** keywords. Such a disjunction can be expressed in

¹To improve readability, we will write attribute–value pairs (a, v) as $(a: v)$, conforming more to a programming language-like syntax.

BSL using `or`. To illustrate using a few concrete examples:

$$\begin{aligned}
 \neg(a: v) &\equiv (\neg a: \neg\emptyset) \text{ or } (\neg\emptyset: \neg\{v\}) \\
 \neg(\neg A: v) &\equiv (A: \neg\emptyset) \text{ or } (\neg\emptyset: \neg\{v\}) \\
 (\neg A: v) \vee (a: \neg V) &\equiv (\neg A: v) \text{ or } (a: \neg V) \\
 \neg((\neg A: v) \vee (a: \neg V)) &\equiv (A \setminus \{a\}: \neg\emptyset) \text{ or } (A \cap \{a\}: V) \text{ or } (\neg\{a\}: \neg\{v\}) \text{ or } (a: V) \quad \square
 \end{aligned}$$

Since we have the wildcard, Corollary 9 applies, so we know that $\text{BSL}(\text{AV})$ defines exactly all search queries that are K -distinguishing for some finite set $K \subseteq \text{AV}$. Recall that a search query q is K -distinguishing if it is invariant under the equivalence \simeq_K on objects (Definition 7). In the attribute-value setting, we can formulate a more intuitive alternative to this equivalence relation, directly in terms of attributes and values, rather than AV keywords. The idea, similar to full genericity, is that only a finite set of attributes and values can be distinguished.

Definition 34. Let W be a finite set of attributes and values. Let \diamond be a “blank value”, that is an arbitrary element not in W . For an attribute or value x , define

$$\text{blank}_W(x) := \begin{cases} x & \text{if } x \in W \\ \diamond & \text{otherwise.} \end{cases}$$

We extend blank_W to attribute-value pairs, objects, and dataspace in the canonical, pointwise manner:

$$\begin{aligned}
 \text{blank}_W(a, v) &:= (\text{blank}_W(a), \text{blank}_W(v)), \\
 \text{blank}_W(o) &:= \{\text{blank}_W(a, v) \mid (a, v) \in o\}, \\
 \text{blank}_W(D) &:= \{\text{blank}_W(o) \mid o \in D\}.
 \end{aligned}$$

Two objects o_1 and o_2 are called W -equivalent if $\text{blank}_W(o_1) = \text{blank}_W(o_2)$. We denote this by $o_1 \simeq_W o_2$. We now say that a search query q is W -distinguishing if for any two dataspace D_1 and D_2 and objects $o_1 \in D_1$ and $o_2 \in D_2$ that are W -equivalent, we have $o_1 \in q(D_1)$ iff $o_2 \in q(D_2)$.

PROPOSITION 35. *A search query is K -distinguishing for a finite set $K \subseteq \text{AV}$ if, and only if, it is W -distinguishing for a finite set of attributes and values W .*

PROOF. The crux of the matter is that if two objects are K -equivalent, for some finite set of AV keywords K , then they are W -equivalent, where W is the finite set of attributes and values explicitly mentioned in the keywords in K . Conversely, if two objects are W -equivalent, for some finite set of attributes and values W , then they are K -equivalent, where K is the finite set of all AV keywords that can be constructed using the elements in W . \square

4.2 Associative Search

Let us now turn to associative search in the attribute-value context. What are the simrels needed to join objects in an attribute-value dataspace? Focusing on equijoins, there are three natural possibilities: two pairs can be compared on their values, on their attributes, or on both together. So, for the set \mathcal{S} of simrels in the attribute-value setting, we will use the set $\{\text{eq}, \text{eq-attr}, \text{eq-val}\}$ defined as follows:

Definition 36. An *eqrel* (short for equality relation) is one of the three following simrels on attribute–value pairs:

$$\begin{aligned}(a, v) \text{ eq } (b, w) &\Leftrightarrow a = b \text{ and } v = w, \\ (a, v) \text{ eq-attr } (b, w) &\Leftrightarrow a = b, \\ (a, v) \text{ eq-val } (b, w) &\Leftrightarrow v = w.\end{aligned}$$

We denote by EQ the set {eq, eq-attr, eq-val} of all eqrels.

Example 37. For example, the following is an expression in the associative search language ASL(AV, EQ):

$$\text{link}((\text{name}: \star) \text{ eq-val } (\text{author}: \star))(\text{published in: ICDT 2009})$$

It defines the query that retrieves all authors of objects published in ICDT 2009. More precisely, it retrieves all objects with a value for attribute **name** that equals a value for attribute **author** in an object containing the item (published in, ICDT 2009).

In other dataspace models based on attribute–value pairs [Dong and Halevy 2007; Dittrich and Vaz Salles 2006], objects are not joined using eqrel simlinks, but through explicit named links (edges) between objects. So there, a dataspace is not merely a set of objects, but a graph of objects. Using eqrel simlinks as we do, an explicit graph model is redundant. Indeed, as just illustrated in Example 37, named edges (for instance linking papers to their authors) can easily be represented using id attributes for objects and “pointer” attributes (**name** and **author**) having these ids as values.

We next show that the abstract bisimulation Lemma 19 can well be applied in the present setting as an aid to understand the limits of expressive power of ASL(AV, EQ). For example, consider the generic equality selection query $q_{a=b}$, for two attributes a and b , defined as follows:

$$q_{a=b}(D) = \{o \in D \mid \exists v : (a, v) \in o \text{ and } (b, v) \in o\}$$

It is not immediately clear whether or not this query is definable in ASL(AV, EQ); it turns out it is not:

PROPOSITION 38. $q_{a=b}$ is not definable in ASL(AV, EQ).

PROOF. In a nutshell, the proof consist of showing that $q_{a=b}$ is not bisimulation invariant relative to any finite set of AV keywords and eqrel simlinks, and then invoking Lemma 19. Concretely, let e be an arbitrary expression in ASL(AV, EQ), let K be the finite set of keywords mentioned in e (be it as an expression or as a link condition), let

$$L = \{k \sim l \mid k, l \in K, \sim \in \text{EQ}\},$$

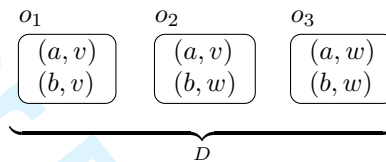
and let $n = \text{depth}(e)$ be the nesting depth of link expressions in e . To see that e cannot define $q_{a=b}$, first define, for every $M \subseteq K$, the set $\llbracket M \rrbracket$ of pairs matched by all keywords in M and none of the keywords in $K - M$:

$$\llbracket M \rrbracket := (\Sigma \times \mathcal{V}) \cap \bigcap_{k \in M} k - \bigcup_{k' \in K - M} k'.$$

Then the set $\{\llbracket M \rrbracket \mid M \subseteq K\}$ forms a partition of $\Sigma \times \mathcal{V}$. In particular, $\Sigma \times \mathcal{V} = \bigcup_{M \subseteq K} \llbracket M \rrbracket$. Since the subset $\{(a, u) \mid u \in \mathcal{V}\}$ of $\Sigma \times \mathcal{V}$ is infinite and since there are only a finite number of subsets of K , there must be at least one $M^a \subseteq K$ for which $\llbracket M^a \rrbracket \cap \{(a, u) \mid u \in \mathcal{V}\}$ is infinite. Then let $\{M_1^b, \dots, M_m^b\} = \{M \subseteq K \mid \exists u \in \mathcal{V} : (b, u) \in \llbracket M \rrbracket\}$. Since $\{u \mid (a, u) \in \llbracket M^a \rrbracket\}$ is infinite, since $\mathcal{V} = \bigcup_{m=1}^n \{u \mid (b, u) \in \llbracket M_i^b \rrbracket\}$, and since m is finite, there must exist some $1 \leq i \leq m$ such that

$$\{u \mid (a, u) \in \llbracket M^a \rrbracket\} \cap \{u \mid (b, u) \in \llbracket M_i^b \rrbracket\}$$

is infinite. Then fix v, w in this intersection with $v \neq w$. By construction, $(a, v), (a, w) \in \llbracket M^a \rrbracket$ and $(b, v), (b, w) \in \llbracket M_i^b \rrbracket$, i.e., (a, v) and (a, w) satisfy the same keywords over K and similarly for (b, v) and (b, w) . Now fix D as follows.



To prove the proposition, it suffices to show $(D, o_1) \rightleftharpoons_n^{K,L} (D, o_2)$. Indeed, then $o_1 \in e(D) \Leftrightarrow o_2 \in e(D)$ by Lemma 19, while $o_1 \in q_{a=b}(D)$ but $o_2 \notin q_{a=b}(D)$. As such, e cannot define $q_{a=b}$.

We actually show that $(D, o_1) \rightleftharpoons_n^{K,L} (D, o_2)$ and $(D, o_1) \rightleftharpoons_n^{K,L} (D, o_3)$ by simultaneous induction on n . When $n = 0$, $o_1 \simeq_K o_2$ and $o_1 \simeq_K o_3$ by construction ((a, v) matches the same keywords as (a, w) and (b, v) matches the same keywords as (b, w)), and hence $(D, o_1) \rightleftharpoons_0^{K,L} (D, o_2)$ and $(D, o_1) \rightleftharpoons_0^{K,L} (D, o_3)$. Inductively, it suffices to check the forth and back properties.

(Forth) Suppose that there is some $\lambda \in L$ and some $p_1 \in D$ such that $\lambda(o_1, p_1)$. We show that there exists $p_2, p_3 \in D$ such that (1) $\lambda(o_2, p_2)$ and $(D, p_1) \rightleftharpoons_{n-1}^{K,L} (D, p_2)$ and (2) $\lambda(o_3, p_3)$ and $(D, p_1) \rightleftharpoons_{n-1}^{K,L} (D, p_3)$.

—Case $\lambda = k \text{ eq } l$. Then $o_1 \cap k \cap p_1 \cap l$ is non-empty. There are two possibilities.

(1) When (a, v) is in the intersection, we can take $p_2 := p_1$ and $p_3 := o_3$. Indeed, to see that it suffices to take $p_2 := p_1$, observe that (a, v) is also in $o_2 \cap k \cap p_2 \cap l$ (which is hence non-empty). Moreover, since \rightleftharpoons is an equivalence relation and $p_1 = p_2$, also $(D, p_1) \rightleftharpoons_{n-1}^{K,L} (D, p_2)$. To see that it suffices to take $p_3 := o_3$, observe that by construction also $(a, w) \in o_3 \cap k \cap p_3 \cap l$ since (a, v) and (a, w) satisfy the same keywords in K . Moreover, since p_1 is either o_1, o_2 , or o_3 ; since $(D, o_1) \rightleftharpoons_{n-1}^{K,L} (D, o_2)$ and $(D, o_1) \rightleftharpoons_{n-1}^{K,L} (D, o_3)$ (by induction hypothesis); since hence also $(D, o_2) \rightleftharpoons_{n-1}^{K,L} (D, o_3)$ (by transitivity of \rightleftharpoons); and since $(D, o_3) \rightleftharpoons_{n-1}^{K,L} (D, o_3)$ (by reflexivity of \rightleftharpoons), certainly $(D, p_1) \rightleftharpoons_{n-1}^{K,L} (D, p_3)$. (2) When (b, v) is in the intersection, as symmetric argument shows that it suffices to take $p_2 := o_2$ and $p_3 := p_1$.

—Case $\lambda = k \text{ eq-attr } l$. Then $\pi_1(o_1 \cap k) \cap \pi_1(p_1 \cap l)$ is non-empty. As such, either a is in the intersection, or b is. In both event, it can be seen using a similar argument as in the case $\lambda = k \text{ eq } l$ that it suffices to take $p_2 := p_1$ and $p_3 := p_1$.

—Case $\lambda = k \text{ eq-val } l$. Then $\pi_2(o_1 \cap k) \cap \pi_2(p_1 \cap l)$ is non-empty. Since v is the only word in $\pi_2(o_1)$, $\pi_2(o_1 \cap k) \cap \pi_2(p_1 \cap l) = \{v\}$. There are four possibilities

why this is so:

- (1) $(a, v) \in o_1 \cap k$ and $(a, v) \in p_1 \cap l$;
- (2) $(a, v) \in o_1 \cap k$ and $(b, v) \in p_1 \cap l$;
- (3) $(b, v) \in o_1 \cap k$ and $(a, v) \in p_1 \cap l$;
- (4) $(b, v) \in o_1 \cap k$ and $(b, v) \in p_1 \cap l$.

Similar arguments as in the case $\lambda = k \text{ eq } l$ show that it suffices to take $p_2 := o_2$ and $p_3 := o_3$ in the first event; $p_2 := o_1$ and $p_3 := o_2$ in the second; $p_2 := o_3$ and $p_3 := o_3$ in the third; and $p_2 := o_2$ and $p_3 := o_3$ in the fourth.

(Back) Similar to Forth. \square

4.3 On the correspondence with the semijoin algebra

We conclude this section by returning to the equivalence between ASL and the semijoin algebra. Recall that in Section 2.3 we have defined the semijoin algebra to work on abstract dataspace represented as binary relations over the schema $\{\text{id}, \text{item}\}$. While the equivalence of ASL and $\text{SA}^{\text{search}}$ (Theorem 28) can be directly applied to the attribute–value setting, it is not so natural to store attribute–value pairs in a single item column. It is more natural to represent attribute–value dataspace as sets of triples, i.e., as *ternary* relations over the schema $\{\text{id}, \text{attr}, \text{val}\}$. Also the RDF query language SPARQL works over such ternary relations [spa 2008; Gutierrez et al. 2004; Pérez et al. 2009]; RDF graphs can also be viewed as a dataspace model [Dong and Halevy 2007].

So, it is worthwhile to define an alternative to $\text{SA}^{\text{search}}$ working on ternary relations. An added simplification is that, since we are working with eqrel simlinks rather than simlinks based on general abstract simrels, we will no longer need the \sim -semijoin operator and will have enough with the standard natural semijoin.

Definition 39. The fragment SA^{AV} of the semijoin algebra, defined on relations $T : \{\text{id}, \text{attr}, \text{val}\}$, is defined by the following grammar:

$$\begin{aligned} E ::= & T \mid \sigma_{\text{attr}=c}(E) \mid \sigma_{\text{val}=c}(E) \mid E \cup E \mid E - E \mid \\ & \mid \pi_{\alpha}(E) \mid E \bowtie \pi_{\text{id}}(E) \mid \pi_{\{\text{id}\}}(E \bowtie \pi_{\beta}(E)) \end{aligned}$$

where c ranges over attribute and value constants; α is either $\{\text{id}\}$, $\{\text{id}, \text{attr}\}$, or $\{\text{id}, \text{val}\}$; and β is either $\{\text{attr}\}$, $\{\text{val}\}$, or $\{\text{attr}, \text{val}\}$.

The semantics of \bowtie is the standard natural semijoin on equality of common attributes.

We have the following analog of Theorem 28:

THEOREM 40. *A search query is definable in $\text{ASL}(\text{AV}, \text{EQ})$, if and only if it is definable in SA^{AV} , where we regard a dataspace D as a ternary relation $\{(o, a, v) \mid o \in D \text{ and } (a, v) \in o\}$.*

PROOF. The only-if direction is similar to the proof of Theorem 28. AV keywords are expressed using combinations of constant selections using union and difference. Simlinks based on eq-attr, eq-val, or eq are expressed using semijoin with projection on the right (the set β in the syntax definition) equal to $\{\text{attr}\}$, $\{\text{val}\}$, or $\{\text{attr}, \text{val}\}$, respectively.

The if-direction is also similar, but we have the added complication that boolean combinations of keywords are not directly available in the AV setting. Yet, because

of Proposition 33, we can simulate them in the language. Proposition 33 is only formulated for BSL(AV), but the same argument holds for ASL(AV, EQ), because link distributes over a disjunction of keywords used in a simlink:

$$\text{link}(\langle(\varphi_1 \vee \varphi_2) \sim \psi\rangle(e) = \text{link}\langle\varphi_1 \sim \psi\rangle(e) \text{ or } \text{link}\langle\varphi_2 \sim \psi\rangle(e).$$

Semijoins are translated to simlinks using an eqrel depending on β : if $\beta = \{\text{attr}\}$ we use eq-attr, if $\beta = \{\text{val}\}$ we use eq-val, and if $\beta = \{\text{attr}, \text{val}\}$ we use eq. \square

5. DISCUSSION

Our goal has been to provide the beginnings of a theoretical foundation for search and associative search queries, motivated by the ubiquity of such queries in everyday information systems. Our approach has been to investigate search queries as restricted kinds of database queries, and to use the tools and the concepts already developed in the theory of database queries.

We have first presented a general abstract theory, then applied it to the concrete setting of attribute–value dataspace. It would be interesting to conduct a similar application to the XML data model, for example, with XPath playing the role of relational algebra.

Mainly inspired by dataspace [Dong and Halevy 2007], we have focused on selection queries, i.e., queries that always return a subset of the original objects. Current approaches to keyword search on structured databases [Golenberg et al. 2008; Qin et al. 2009, we give just two recent references] return tuples of objects that are related by patterns. In the semijoin algebra one can express such patterns as long as they are not cyclic, but the patterns themselves cannot be returned. It remains to be investigated if and how our theory should be extended so that patterns can be returned. Of course, one can simply move to the full relational algebra, but then there is less news to discover.

It would also be interesting to look at a similar theory of search queries on RDF graphs (i.e., ternary relations, discussed above in Section 4.3). Indeed, RDF graphs can be viewed as a particular instance of the classical relations, and hence our results of Section 3.1 for boolean search already apply here. The expressivity of SPARQL, which is a kind of specialized relational algebra for RDF, has been studied recently [Fletcher 2008; Gutierrez et al. 2004; Pérez et al. 2009]. Given their close similarity to AV dataspace, it might also be fruitful to investigate specializations of our semijoin algebra fragment SA^{AV} into interesting fragments of SPARQL, or of its navigational extensions (e.g., [Alkhateeb et al. 2009; Pérez et al. 2008]).

Received Month year; revised Month year; accepted Month year.

REFERENCES

- 2004. RDF primer. W3C Recommendation.
- 2008. SPARQL query language for RDF. W3C Recommendation.
- ABITEBOUL, S., HULL, R., AND VIANU, V. 1995. *Foundations of Databases*. Addison-Wesley.
- AHO, A. AND ULLMAN, J. 1979. Universality of data retrieval languages. In *Conference Record, 6th ACM Symposium on Principles of Programming Languages*. 110–120.
- ALKHATEEB, F., BAGET, J.-F., AND EUZENAT, J. 2009. Extending SPARQL with regular expression patterns (for querying RDF). *J. Web Sem.* 7, 2, 57–73.
- ACM Transactions on Database Systems, Vol. V, No. N, October 2009.

ANONYMOUS. details omitted due to double-blind reviewing.

BAEZA-YATES, R. AND RIBEIRO-NETO, B. 1999. *Modern Information Retrieval*. Addison-Wesley.

BEERI, C., MILO, T., AND TA-SHMA, P. 1996. On genericity and parametricity. In *Proceedings 15th ACM Symposium on Principles of Database Systems*. 104–116.

BEERI, C., MILO, T., AND TA-SHMA, P. 1997. Towards a languages for the fully generic queries. In *Database Programming Languages*, S. Cluet and R. Hull, Eds. Lecture Notes in Computer Science. Springer, 239–259.

BLACKBURN, P., DE RIJKE, M., AND VENEMA, Y. 2001. *Modal Logic*. Cambridge University Press.

BLACKBURN, P., VAN BENTHEM, J., AND WOLTER, F., Eds. 2007. *Handbook of Modal Logic*. Elsevier.

CHANDRA, A. AND HAREL, D. 1980. Computable queries for relational data bases. *J. Comput. Syst. Sci.* 21, 2, 156–178.

DECANDIA, G. ET AL. 2007. Dynamo: Amazon’s highly available key-value store. *SIGOPS Oper. Syst. Rev.* 41, 6, 205–220.

DITTRICH, J.-P. AND VAZ SALLES, M. 2006. iDM: A unified and versatile data model for personal dataspace management. In *Proceedings 32nd International Conference on Very Large Data Bases*. 367–378.

DONG, X. AND HALEVY, A. 2007. Indexing dataspaces. In *Proceedings ACM SIGMOD International Conference on Management of Data*. 43–54.

FLETCHER, G. 2008. An algebra for basic graph patterns. Presented at the workshop on Logic in Databases, Rome, Italy.

GOLENBERG, K., KIMELFELD, B., AND SAGIV, Y. 2008. Keyword proximity search in complex data graphs. In *Proceedings ACM SIGMOD International Conference on Management of Data*. 927–940.

GORANKO, V. AND OTTO, M. 2007. Model theory of modal logic. See Blackburn et al. [2007], Chapter 5.

GUTIERREZ, C., HURTADO, C., AND MENDELZON, A. 2004. Foundations of semantic web databases. In *Proceedings 23rd ACM Symposium on Principles of Database Systems*. 95–106.

HALEVY, A., FRANKLIN, M., AND MAIER, D. 2006. Principles of dataspace systems. In *Proceedings 25th ACM Symposium on Principles of Database Systems*. 1–9.

JAGADISH, H. ET AL. 2007. Making database systems usable. In *Proceedings ACM SIGMOD International Conference on Management of Data*. 13–24.

JAIN, M., MENDHEKAR, A., AND VAN GUCHT, D. 1995. A uniform data model for relational data and meta-data query processing. In *Advances in Data Management ’95*. Tata McGraw-Hill, 146–165.

LEINDERS, D., MARX, M., TYSZKIEWICZ, J., AND VAN DEN BUSSCHE, J. 2005. The semijoin algebra and the guarded fragment. *Journal of Logic, Language and Information* 14, 331–343.

LEINDERS, D. AND VAN DEN BUSSCHE, J. 2007. On the complexity of division and set joins in the relational algebra. *J. Comput. Syst. Sci.* 73, 4, 538–549.

OLSTON, C., REED, B., SRIVASTAVA, U., KUMAR, R., AND TOMKINS, A. 2008. Pig latin: a not-so-foreign language for data processing. In *Proceedings ACM SIGMOD International Conference on Management of Data*. 1099–1110.

PÉREZ, J., ARENAS, M., AND GUTIERREZ, C. 2008. nSPARQL: A navigational language for RDF. In *International Semantic Web Conference*. 66–81.

PÉREZ, J., ARENAS, M., AND GUTIERREZ, C. 2009. Semantics and complexity of SPARQL. *ACM Trans. Database Syst.* 34, 3.

QIN, L., YU, J. X., AND CHANG, L. 2009. Keyword search in databases: the power of RDBMS. In *Proceedings ACM SIGMOD International Conference on Management of Data*. 681–694.

ROSS, K. AND JANEVSKI, A. 2005. Querying faceted databases. In *Proceedings 2nd International Workshop on Semantic Web and Databases*. Lecture Notes in Computer Science, vol. 3372. Springer, 199–218.

ULLMAN, J. 1989. *Principles of Database and Knowledge-Base Systems*. Vol. II. Computer Science Press.