



On database query languages for \mathcal{K} -relations

Floris Geerts^{a,*}, Antonella Poggi^b

^a University of Edinburgh, United Kingdom

^b Sapienza Università di Roma, Italy

ARTICLE INFO

Article history:

Available online 22 September 2009

Keywords:

Relational model

Query language

Annotations

Provenance

Language completeness

ABSTRACT

The relational model has recently been extended to so-called \mathcal{K} -relations in which tuples are assigned a unique value in a semiring \mathcal{K} . A query language, denoted by $\mathcal{RA}_{\mathcal{K}}^+$, similar to the classical positive relational algebra, allows for the querying of \mathcal{K} -relations. In this paper, we define more expressive query languages for \mathcal{K} -relations that extend $\mathcal{RA}_{\mathcal{K}}^+$ with the *difference* and *constant annotations* operations on annotated tuples. The latter are natural extensions of the duplicate elimination operator of the relational algebra on bags. We investigate conditions on semirings under which these operations can be added to $\mathcal{RA}_{\mathcal{K}}^+$ in a natural way, and establish basic properties of the resulting query languages. Moreover, we show how the provenance semiring of Green et al. can be extended to record provenance of data in the presence of difference and constant annotations. Finally, we investigate the *completeness* of $\mathcal{RA}_{\mathcal{K}}^+$ and extensions thereof in the sense of Bancilhon and Paredaens.

© 2009 Elsevier B.V. All rights reserved.

1. Introduction

Annotated relations appear in various contexts in the database literature. The querying of such relations involves the generalization of the relational algebra to perform corresponding operations on the annotations. Recently, a general data model (referred to as \mathcal{K} -relations) has been proposed for annotated relations in which tuples in a relation are assigned a unique value coming from a semiring \mathcal{K} [12]. By varying the semiring \mathcal{K} , \mathcal{K} -relations can model the standard relational model with both set [1] and bag semantics [16], incomplete databases (positive Boolean c -tables to be more precise) [13, 15] and probabilistic databases [10,19]. Moreover, operations that queries in the relational algebra perform on tuples can be naturally extended to operations on annotated tuples. More specifically, operations on tuples naturally translate into the algebraic operations (sum and product) in semirings. This leads to the definition of the positive relational algebra on \mathcal{K} -relations, or $\mathcal{RA}_{\mathcal{K}}^+$ for short [12].

The generality of semirings further allows for the definition of new data models which are of particular interest for the study of provenance of data [6,12]. A notable example is the *provenance semiring* that allows to record provenance information of data obtained as result of positive relational algebra queries. A crucial property of this semiring, named *factorization property*, is that it is the most general semiring. That is, for any semiring \mathcal{K} , to evaluate queries in $\mathcal{RA}_{\mathcal{K}}^+$ on \mathcal{K} -relations it is sufficient to know how to evaluate these queries on the provenance semiring.

In this paper, we study query languages for \mathcal{K} -relations. Indeed, while some basic properties of $\mathcal{RA}_{\mathcal{K}}^+$ are already established in [12], less is known about its expressive power. Furthermore, it was left open in [12] how to incorporate difference in $\mathcal{RA}_{\mathcal{K}}^+$ to get a full relational algebra on \mathcal{K} -relations. Hence, our goal is twofold. On one hand, we define more expressive query languages for \mathcal{K} -relations that extend $\mathcal{RA}_{\mathcal{K}}^+$ with operations on annotated tuples that are natural extensions of the

* Corresponding author.

E-mail address: fgeerts@inf.ed.ac.uk (F. Geerts).

difference and duplicate elimination operations of the standard relational algebra. On the other hand, we investigate the expressive power of $\mathcal{RA}_{\mathcal{K}}^+$ and extensions thereof. In particular, we investigate the *completeness* of these query languages. Recall that Codd qualified a query language on relational databases as *complete* if its expressive power is at least that of the relational calculus [8]. Bancilhon [4] and Paredaens [18] independently provided a *language-independent* characterization of completeness. This characterization, known as BP-completeness, can be stated as follows: a relation R_2 is the result of a relational algebra query applied to a database R_1 if and only if (i) the active domain of R_2 is included in the active domain of R_1 ; and (ii) every automorphism of R_1 is also an automorphism of R_2 .

The contributions of the paper can be summarized as follows:

- First, we define the query languages $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$, $\mathcal{RA}_{\mathcal{K}}^+(\delta)$ and $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$, obtained by extending $\mathcal{RA}_{\mathcal{K}}^+$ with *difference*, *constant annotations*, and with both *difference* and *constant annotations*, respectively. Here, constant annotations correspond to a family of operators that assign annotations to tuples among a finite set of elements of the semiring, that are the *semiring generators*. Note, in particular, that extending $\mathcal{RA}_{\mathcal{K}}^+$ with these operators forces to restrict the class of semirings under consideration. Specifically, on one hand, adding difference requires the definition of a *monus* operator on the underlying semiring, which might not always be possible. We call *m-semirings* the class of semirings admitting a monus operator. On the other hand, constant annotations require the underlying semiring to be *finitely generated*, i.e., to have a finite set of semiring generators. Interestingly, we observe that most semirings encountered in the literature are indeed finitely generated *m-semirings*.
- Second, we show how to extend the provenance semiring of [12], so that it can be used to record the provenance of data obtained as result of queries in $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$, $\mathcal{RA}_{\mathcal{K}}^+(\delta)$ and $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$. We show that, similarly to $\mathcal{RA}_{\mathcal{K}}^+$, the extended provenance semirings also satisfy the factorization property.
- Finally, we naturally extend the notion of BP-completeness to the setting of \mathcal{K} -relations and investigate whether query languages on \mathcal{K} -relations proposed so far are BP-complete. In particular, we show that none of the languages $\mathcal{RA}_{\mathcal{K}}^+$, $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$ and $\mathcal{RA}_{\mathcal{K}}^+(\delta)$ is BP-complete on \mathcal{K} -relations for arbitrary semirings, *m-semirings*, and finitely generated semirings, respectively. In contrast, $\mathcal{RA}_{\mathcal{K}}^+$ was shown to be BP-complete in the standard relational case [4,18]. We show, however, that $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$ is BP-complete on \mathcal{K} -relations for arbitrary finitely generated *m-semirings* \mathcal{K} .

Organization. The paper is organized as follows. After recalling in Section 2 the basic notions of \mathcal{K} -relations and the positive query language $\mathcal{RA}_{\mathcal{K}}^+$, we present in Section 3, the query languages $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$, $\mathcal{RA}_{\mathcal{K}}^+(\delta)$ and $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$, obtained by extending $\mathcal{RA}_{\mathcal{K}}^+$ with *difference* and *constant annotations*. Then, in Section 4, we discuss the relationship between provenance and \mathcal{K} -relations, and show how the provenance semiring can be extended to record provenance for $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$, $\mathcal{RA}_{\mathcal{K}}^+(\delta)$ and $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$. Section 5 discusses BP-completeness of $\mathcal{RA}_{\mathcal{K}}^+$ and extensions thereof. We conclude the paper in Section 6.

2. Preliminaries

In this section we recall the notions of \mathcal{K} -relation and the query language $\mathcal{RA}_{\mathcal{K}}^+$ that were introduced by Green et al. [12]. Then, we conclude the section by discussing an important property of $\mathcal{RA}_{\mathcal{K}}^+$, named *homomorphism property*.

2.1. \mathcal{K} -relations

A (commutative) semiring $\mathcal{K} = (\mathbb{K}, \oplus, \otimes, 0, 1)$ is an algebraic structure consisting of a set \mathbb{K} equipped with two binary operations, i.e., sum (\oplus) and product (\otimes), such that $(\mathbb{K}, \oplus, 0)$ is a commutative monoid with identity element 0; $(\mathbb{K}, \otimes, 1)$ is a commutative monoid with identity element 1; the operation \otimes distributes over \oplus ; and finally 0 is an annihilating element. Recall that a monoid consists of a set equipped with a binary operation that is associative and that has an identity element. Furthermore, the set is closed under the binary operation, i.e., the result of the operation on any two elements in the set belongs to the set as well.

Example 1. It is easily verified that the following structures are semirings: (1) the Boolean semiring $\mathcal{K}_{\mathbb{B}} = (\mathbb{B}, \vee, \wedge, \text{false}, \text{true})$ with $\mathbb{B} = \{\text{true}, \text{false}\}$; (2) the natural numbers semiring $\mathcal{K}_{\mathbb{N}} = (\mathbb{N}, +, \times, 0, 1)$; (3) the positive Boolean expressions semiring $\mathcal{K}_{\text{c-table}^+} = (\text{PosBool}(X), \vee, \wedge, \text{false}, \text{true})$, where $\text{PosBool}(X)$ is the set of all Boolean expressions (over a finite set of variables X) that involve only disjunction, conjunction, and constants for *true* and *false* and in which any two equivalent expressions are identified; and (4) the probabilistic semiring $\mathcal{K}_{\text{prob}} = (\mathcal{P}(\Omega), \cup, \cap, \emptyset, \Omega)$, where Ω is a finite set of events and $\mathcal{P}(\Omega)$ stands for the powerset of Ω .

To formally introduce semirings into the relational data model, we next recall the definition of \mathcal{K} -relations (see [12] for more details). Let \mathbb{D} be an (infinite) domain of data values and let U be a finite set of attributes. We define an *U-tuple* \bar{t} to be a mapping from $U \rightarrow \mathbb{D}$. The set of *U-tuples* is denoted by $U\text{-Tup}$. Let $\mathcal{K} = (\mathbb{K}, \oplus, \otimes, 0, 1)$ be a semiring. A *\mathcal{K} -relation* R over U is then a function $R: U\text{-Tup} \rightarrow \mathbb{K}$. The *support* of a \mathcal{K} -relation R , denoted by $\text{supp}(R)$, is defined as $\text{supp}(R) = \{\bar{t} \mid R(\bar{t}) \neq 0\}$; it is the standard relational database underlying R . The *active domain* of a \mathcal{K} -relation R , denoted by $\text{adom}(R)$, is defined as the set of data values (in \mathbb{D}) occurring in $\text{supp}(R)$.

$R_1 =$	<table border="1"><thead><tr><th>drink</th><th>kind</th><th>origin</th><th></th></tr></thead><tbody><tr><td>Montefalco</td><td>wine</td><td>Italy</td><td>true</td></tr><tr><td>Pinot</td><td>grappa</td><td>Italy</td><td>true</td></tr></tbody></table>	drink	kind	origin		Montefalco	wine	Italy	true	Pinot	grappa	Italy	true
drink	kind	origin											
Montefalco	wine	Italy	true										
Pinot	grappa	Italy	true										

$R_2 =$	<table border="1"><thead><tr><th>drink</th><th>kind</th><th>origin</th><th></th></tr></thead><tbody><tr><td>Stella</td><td>beer</td><td>Belgium</td><td>2</td></tr><tr><td>Montefalco</td><td>wine</td><td>Italy</td><td>1</td></tr><tr><td>Pinot</td><td>grappa</td><td>Italy</td><td>1</td></tr></tbody></table>	drink	kind	origin		Stella	beer	Belgium	2	Montefalco	wine	Italy	1	Pinot	grappa	Italy	1
drink	kind	origin															
Stella	beer	Belgium	2														
Montefalco	wine	Italy	1														
Pinot	grappa	Italy	1														

$R_3 =$	<table border="1"><thead><tr><th>drink</th><th>kind</th><th>origin</th><th></th></tr></thead><tbody><tr><td>Stella</td><td>beer</td><td>Belgium</td><td>party</td></tr><tr><td>Montefalco</td><td>wine</td><td>Italy</td><td>tasting</td></tr><tr><td>Pinot</td><td>grappa</td><td>Italy</td><td>party \vee tasting</td></tr></tbody></table>	drink	kind	origin		Stella	beer	Belgium	party	Montefalco	wine	Italy	tasting	Pinot	grappa	Italy	party \vee tasting
drink	kind	origin															
Stella	beer	Belgium	party														
Montefalco	wine	Italy	tasting														
Pinot	grappa	Italy	party \vee tasting														

$R_4 =$	<table border="1"><thead><tr><th>drink</th><th>kind</th><th>origin</th><th></th></tr></thead><tbody><tr><td>Stella</td><td>beer</td><td>Belgium</td><td>P</td></tr><tr><td>Montefalco</td><td>wine</td><td>Italy</td><td>T</td></tr><tr><td>Pinot</td><td>grappa</td><td>Italy</td><td>P \cup T</td></tr></tbody></table>	drink	kind	origin		Stella	beer	Belgium	P	Montefalco	wine	Italy	T	Pinot	grappa	Italy	P \cup T
drink	kind	origin															
Stella	beer	Belgium	P														
Montefalco	wine	Italy	T														
Pinot	grappa	Italy	P \cup T														

Fig. 1. Examples of \mathcal{K} -relations.

As already mentioned in the introduction, \mathcal{K} -relations have recently been used to unify a variety of data models, including the standard relational model with both set and bag semantics, incomplete databases (positive Boolean c -tables to be more precise) and probabilistic databases [12].

Example 2. Consider the set of attributes $U = \{\text{drink, kind, origin}\}$. Fig. 1 shows \mathcal{K} -relations over U , for the four different semirings described in Example 1. Strictly speaking, a \mathcal{K} -relation assigns a semiring value to every possible tuple. In Fig. 1 we only show the support of the \mathcal{K} -relations. The semiring value associated with each tuple is shown in the last column. (1) R_1 is a $\mathcal{K}_{\mathbb{B}}$ -relation and corresponds to a standard relational table with set semantics; specifically, the standard relational table corresponding to R_1 contains the tuples $\bar{t}_m = (\text{Montefalco, wine, Italy})$ and $\bar{t}_p = (\text{Pinot, grappa, Italy})$; (2) R_2 is a $\mathcal{K}_{\mathbb{N}}$ -relation and corresponds to a relational table with bag semantics; the bag corresponding to R_2 contains two tuples $\bar{t}_s = (\text{Stella, beer, Belgium})$, one tuple \bar{t}_m and one tuple \bar{t}_p ; (3) R_3 is a $\mathcal{K}_{c\text{-table}^+}$ and corresponds to a positive Boolean c -table [13]; Boolean c -tables are a restricted form of c -tables [15] in which tuples are annotated with conditions that can be any Boolean expression and variables can only take Boolean values and appear in conditions (not in the attributes); positive Boolean c -tables are Boolean c -tables in which annotation are positive Boolean expressions; hence, the c -table corresponding to R_3 represents a set of possible worlds, according to the closed-world semantics as defined in [15]; finally, (4) R_4 is a $\mathcal{K}_{\text{prob}}$ -relation and corresponds to a probabilistic event table introduced in [10,19]; assuming that both P and T denote probabilistic events, then R_4 corresponds to a probabilistic event table stating that the tuple \bar{t}_s occurs with the probability of event P , the tuple \bar{t}_m with probability of event T and the tuple \bar{t}_p with probability of the event $P \cup T$.

The real strength of \mathcal{K} -relations becomes apparent, however, when considering provenance information. Indeed, the flexibility of semirings allows for the definition of new provenance models at different levels of granularity. We will illustrate this in more detail in Section 4 after we describe query languages on \mathcal{K} -relations.

2.2. The query language $\mathcal{RA}_{\mathcal{K}}^+$

The introduction of semirings in the relational model requires the redefinition of the semantics of the standard relational algebra operators. Recall that the relational algebra consists of projection, selection, union, renaming and difference [1]. When difference is omitted, one obtains the so-called positive fragment of the relational algebra or positive algebra for short. In [12], the semantics of the positive algebra on \mathcal{K} -relations has been introduced. We next recall the definition of the positive relational algebra on \mathcal{K} -relations, denoted by $\mathcal{RA}_{\mathcal{K}}^+$. As before, $\mathcal{K} = (\mathbb{K}, \oplus, \otimes, 0, 1)$ denotes a semiring. Then $\mathcal{RA}_{\mathcal{K}}^+$ includes the following operators:

empty relation For any set of attributes U , we have $\emptyset : U\text{-Tup} \rightarrow \mathbb{K}$ such that $\emptyset(\bar{t}) = 0$ for any \bar{t} .

union If $R_1, R_2 : U\text{-Tup} \rightarrow \mathbb{K}$ then $R_1 \cup R_2 : U\text{-Tup} \rightarrow \mathbb{K}$ is defined by

$$(R_1 \cup R_2)(\bar{t}) = R_1(\bar{t}) \oplus R_2(\bar{t}).$$

projection If $R : U\text{-Tup} \rightarrow \mathbb{K}$ and $V \subseteq U$ then $\pi_V(R) : V\text{-Tup} \rightarrow \mathbb{K}$ is defined by

$$(\pi_V R)(\bar{t}) = \bigoplus_{\substack{\bar{t}' = \bar{t} \text{ on } V \\ R(\bar{t}') \neq 0}} R(\bar{t}').$$

selection If $R : U\text{-Tup} \rightarrow \mathbb{K}$ and the selection predicate \mathbf{P} maps each U -tuple to either 0 or 1 depending on the (in-)equality of attribute values, then $\sigma_{\mathbf{P}}(R) : U\text{-Tup} \rightarrow \mathbb{K}$ is defined by

$$(\sigma_{\mathbf{P}}(R))(\bar{t}) = R(\bar{t}) \otimes \mathbf{P}(\bar{t}).$$

natural join If $R_i : U_i\text{-Tup} \rightarrow \mathbb{K}$, for $i = 1, 2$, then $R_1 \bowtie R_2$ is the \mathcal{K} -relation over $U_1 \cup U_2$ defined by

$$(R_1 \bowtie R_2)(\bar{t}) = R_1(\bar{t}) \otimes R_2(\bar{t}).$$

renaming If $R : U\text{-Tup} \rightarrow \mathbb{K}$ and $\beta : U \rightarrow U'$ is a bijection then $\rho_{\beta}(R)$ is the \mathcal{K} -relation over U' defined by

$$(\rho_{\beta} R)(\bar{t}) = R(\bar{t} \circ \beta^{-1}).$$

It is observed in [12] that the semantics of $\mathcal{RA}_{\mathcal{K}}^+$ coincides with standard positive relational algebras for various semirings encountered in the database literature, i.e., for $\mathcal{K}_{\mathbb{B}}$ (set semantics) [1], $\mathcal{K}_{\mathbb{N}}$ (bag semantics) [16], $\mathcal{K}_{c\text{-tables}^+}$ (positive Boolean c-tables under closed world semantics) [13,15] and $\mathcal{K}_{\text{prob}}$ (probabilistic event tables) [10,19].

2.3. The homomorphism property of $\mathcal{RA}_{\mathcal{K}}^+$

A desirable property of query languages is that they provide the user with a conceptual interface of the underlying data, independent of how exactly that data is stored and without interpreting the exact data objects [2]. In this spirit, intuitively, the *homomorphism property* ensures that the $\mathcal{RA}_{\mathcal{K}}^+$ operations do not interpret the values of the underlying semiring. Formally, let $\mathcal{K} = (\mathbb{K}, \oplus_{\mathbb{K}}, \otimes_{\mathbb{K}}, 0_{\mathbb{K}}, 1_{\mathbb{K}})$ and $\mathcal{K}' = (\mathbb{K}', \oplus_{\mathbb{K}'}, \otimes_{\mathbb{K}'}, 0_{\mathbb{K}'}, 1_{\mathbb{K}'})$ be two semirings and let $h: \mathbb{K} \rightarrow \mathbb{K}'$ be a mapping. It is shown in [12] that the transformation from \mathcal{K} -relations to \mathcal{K}' -relations induced by h , which we also denote by h , satisfies the property that $Q(h(R)) = h(Q(R))$ for any $Q \in \mathcal{RA}_{\mathcal{K}}^+$ iff h is a semiring homomorphism [12]. That is, h satisfies the following properties: $h(0_{\mathbb{K}}) = 0_{\mathbb{K}'}$, $h(1_{\mathbb{K}}) = 1_{\mathbb{K}'}$, and for any $x, y \in \mathbb{K}$, $h(x \oplus_{\mathbb{K}} y) = h(x) \oplus_{\mathbb{K}'} h(y)$ and $h(x \otimes_{\mathbb{K}} y) = h(x) \otimes_{\mathbb{K}'} h(y)$.

3. The query languages $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$, $\mathcal{RA}_{\mathcal{K}}^+(\delta)$ and $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$

In this section we provide three extensions of $\mathcal{RA}_{\mathcal{K}}^+$: First, we extend $\mathcal{RA}_{\mathcal{K}}^+$ with a *difference operator* (\setminus) resulting in the algebra $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$ over \mathcal{K} -relations. Second, we extend $\mathcal{RA}_{\mathcal{K}}^+$ with (a family of) operators called *constant annotations* (δ). These can be thought of as a generalization of the duplicate elimination operator, an operator that is normally included in query languages over bags. The resulting query language is denoted by $\mathcal{RA}_{\mathcal{K}}^+(\delta)$. Finally, we extend $\mathcal{RA}_{\mathcal{K}}^+$ with both the difference and constant annotations, resulting in $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$.

3.1. The query language $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$

We first extend $\mathcal{RA}_{\mathcal{K}}^+$ with a difference operator. More specifically, we identify a large class of semirings that can be equipped with a so-called *monus* operator \ominus . The addition of the monus operator on semirings will then allow to extend $\mathcal{RA}_{\mathcal{K}}^+$ with a difference operator (\setminus). Finally, we show that $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$ satisfies a homomorphism property similar to $\mathcal{RA}_{\mathcal{K}}^+$.

3.1.1. Semirings with monus

We follow the standard approach for introducing a *monus* operator, denoted by \ominus , into additive commutative monoids [3]. As we will see shortly, when introducing \ominus one has to pose some restrictions on the class of semirings. More specifically, we first assume that \mathcal{K} is *naturally ordered*. That is, the quasi-order $x \preccurlyeq y$ on \mathbb{K} defined as $x \preccurlyeq y$ iff there exists a $z \in \mathbb{K}$ such that $x \oplus z = y$, must define a *partial order* on \mathbb{K} . This means that apart from being reflexive and transitive, \preccurlyeq should also be antisymmetric.

It is easily verified that all examples of semirings described in this paper are naturally ordered. We additionally require the following property (\dagger): for each pair of elements $x, y \in \mathbb{K}$, the set $\{z \in \mathbb{K} \mid x \preccurlyeq y \oplus z\}$ has a smallest element. Note that the assumption that \preccurlyeq defines a partial order guarantees that $\{z \in \mathbb{K} \mid x \preccurlyeq y \oplus z\}$ has a *unique* smallest element, provided that it exists.

Definition 1. Let \mathcal{K} be a naturally ordered semiring that satisfies property (\dagger). For any $x, y \in \mathbb{K}$, we define the *monus* $x \ominus y$ to be the smallest element z such that $x \preccurlyeq y \oplus z$. A semiring \mathcal{K} which can be equipped with a monus operator \ominus is called a *semiring with monus* or *m-semiring* for short.

A classical result in theory of additive commutative monoids with monus, or CMM for short, identifies two “natural” classes of CMMs [3]. Indeed, Amer shows that there are only two equationally complete classes of CMMs in the variety of CMMs. These are respectively Boolean algebras (or prime ideals thereof), for which the monus behaves like set difference, and so-called positive cones of lattice-ordered commutative groups, for which the monus behaves like the truncated minus of the natural numbers. Translated to the setting of *m-semirings*, this dichotomy translates to *m-semirings* that are Boolean algebras on the one hand, and *m-semirings* that are the positive cone of a lattice-ordered commutative ring on the other hand [14,17]. In the following example, we revisit the semirings described in Example 1 and discuss their extension to *m-semirings*.

Example 3. One can easily verify that the semirings described in Example 1 in Section 2 all satisfy property (\dagger). Hence, they can *all* be extended to *m-semirings*. Moreover, it is easily verified that they all fall in one of the two natural classes of *m-semirings* described above, except for $\mathcal{K}_{c\text{-table}^+}$. More specifically, $\mathcal{K}_{\mathbb{B}}$ and $\mathcal{K}_{\text{prob}}$ are both Boolean algebras and the monus behaves like set difference. On the other hand, $\mathcal{K}_{\mathbb{N}}$ is the positive cone of the ring \mathbb{Z} , i.e., $\mathbb{N} = \{n \mid n \in \mathbb{Z}, 0 \preccurlyeq n\}$. Consequently, the monus on $\mathcal{K}_{\mathbb{N}}$ corresponds to the truncated minus, i.e., $m \ominus n = m \dot{-} n$ which is defined as $m - n$ if $m > n$ and 0 otherwise. Finally, the case of $\mathcal{K}_{c\text{-table}^+}$ is more subtle since the corresponding *m-semiring* is neither a Boolean algebra nor the positive cone of a lattice-ordered ring. In fact, the semiring $\mathcal{K}_{c\text{-table}^+} = (\text{PosBool}(X), \vee, \wedge, \text{false}, \text{true})$ was defined

in [12] for positive queries only and therefore only positive Boolean expressions over X were allowed. The original definition of Boolean c -tables, however, does allow for *arbitrary* Boolean expressions [13]. Similar to general c -tables [15], the inclusion of difference only makes sense under the closed-world semantics. Recall, however, that \mathcal{K} -relations fully specify a relation and hence correspond to the closed-world semantics. We therefore define the semiring $\mathcal{K}_{c\text{-table}}$ as $(\text{Bool}(X), \vee, \wedge, \text{false}, \text{true})$, where $\text{Bool}(X)$ is the set of Boolean expressions over X in which any two equivalent expressions are identified. Then, each $\mathcal{K}_{c\text{-table}}$ corresponds to the Boolean c -table representing a set of *possible worlds* under the *closed-world semantics*. Clearly, $\mathcal{K}_{c\text{-table}}$ is a Boolean algebra. Furthermore, for any two expressions ϕ_1, ϕ_2 in $\text{Bool}(X)$, we have that $\phi_1 \ominus \phi_2$ is a Boolean expression that is equivalent to $\phi_1 \wedge \neg\phi_2$, as expected.

It is not surprising that not every semiring can be extended to an m -semiring.

Example 4. From the definition of m -semiring it follows that a semiring cannot be extended to an m -semiring if the semiring is not naturally ordered or it is naturally ordered but property (\dagger) fails to hold. For instance, consider the semiring $\mathcal{K}_{\mathbb{R}} = (\mathbb{R}, +, \times, 0, 1)$. Clearly, $r \leq s$ for any two elements $r, s \in \mathbb{R}$ and hence \leq is not antisymmetric. Therefore, $r \ominus s$ cannot be defined in $\mathcal{K}_{\mathbb{R}}$. Consider next the semiring $\mathcal{K}_{\mathbb{R}_{\min}} = (\mathbb{R} \cup \{+\infty\}, \min, +, +\infty, 0)$ where $\min\{x, y\}$ returns the minimum of x and y according to the usual ordering on $\mathbb{R} \cup \{+\infty\}$. It is easily verified $\mathcal{K}_{\mathbb{R}_{\min}}$ is naturally ordered. Indeed, if there exists a z such that $\min\{x, z\} = y$ and if in addition there exists a z' such that $\min\{y, z'\} = x$, then it follows that $x = y$. However, for any $x, y \in \mathbb{R} \cup \{+\infty\}$, the set $\{z \in \mathbb{R} \cup \{+\infty\} \mid x \leq \min\{y, z\}\}$ is equal to $\{z \in \mathbb{R} \cup \{+\infty\} \mid \exists z' \min\{x, z'\} = \min\{y, z\}\}$. Clearly, this is not bounded below since one can take arbitrary small values for z . Hence, although $\mathcal{K}_{\mathbb{R}_{\min}}$ is naturally ordered, it does not satisfy property (\dagger) and the monus operator cannot be defined in this semiring.

3.1.2. The difference operator

We are now ready to extend $\mathcal{RA}_{\mathcal{K}}^+$ with the difference operator. Let \mathcal{K} be an arbitrary m -semiring. Then, we obtain $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$ by extending $\mathcal{RA}_{\mathcal{K}}^+$ with the operator

difference If $R_1, R_2 : U\text{-Tup} \rightarrow \mathbb{K}$ then $R_1 \ominus R_2 : U\text{-Tup} \rightarrow \mathbb{K}$ is defined by

$$(R_1 \setminus R_2)(t) = R_1(t) \ominus R_2(t).$$

As a sanity check, from Example 3, it immediately follows that $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$ coincides with the (full) relational algebra on relational databases for $\mathcal{K}_{\mathbb{B}}$ (set semantics), and the bag algebra with the monus operator for $\mathcal{K}_{\mathbb{N}}$ [16]. Furthermore, in the case of $\mathcal{K}_{c\text{-table}}$ it coincides with the semantics of the relational algebra on Boolean c -tables under closed world semantics [15] and for $\mathcal{K}_{\text{prob}}$ it coincides with the semantics of the relational algebra provided on probabilistic event tables [10,19].

3.1.3. The homomorphism property for $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$

When looking at m -semirings the notion of semiring homomorphism needs to be revisited. Specifically, let $\mathcal{K} = (\mathbb{K}, \oplus_{\mathbb{K}}, \otimes_{\mathbb{K}}, \ominus_{\mathbb{K}}, \mathbf{0}_{\mathbb{K}}, \mathbf{1}_{\mathbb{K}})$ and $\mathcal{K}' = (\mathbb{K}', \oplus_{\mathbb{K}'}, \otimes_{\mathbb{K}'}, \ominus_{\mathbb{K}'}, \mathbf{0}_{\mathbb{K}'}, \mathbf{1}_{\mathbb{K}'})$ be two m -semirings. A mapping $h : \mathbb{K} \rightarrow \mathbb{K}'$ is an *m -semiring homomorphism* if it is a semiring homomorphism and, furthermore, h preserves \ominus , i.e., for any two elements $x, y \in \mathbb{K}$ we have that $h(x \ominus_{\mathbb{K}} y) = h(x) \ominus_{\mathbb{K}'} h(y)$. The following is easily verified:

Proposition 1. Let \mathcal{K} and \mathcal{K}' be two m -semirings. Let $h : \mathbb{K} \rightarrow \mathbb{K}'$ be a mapping. Then, for every query Q in $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$ and for every R , the transformation induced by h from \mathcal{K} -relations to \mathcal{K}' -relations commutes, i.e., $Q(h(R)) = h(Q(R))$, if and only if h is an m -homomorphism.

Proof. We first prove that if h is an m -semiring homomorphism, then for every Q in $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$ and for every R , $Q(h(R)) = h(Q(R))$. We proceed by induction on the structure of queries in $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$. Since $\mathcal{RA}_{\mathcal{K}}^+$ is embedded in $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$ and since every m -semiring homomorphism is a semiring homomorphism, by the homomorphism property for $\mathcal{RA}_{\mathcal{K}}^+$, we only need to treat the case of Q having the form $Q = Q_1 \setminus Q_2$ and can refer to [12] for the other cases. By the induction hypothesis, we have that $Q(h(R)) = Q_1(h(R)) \setminus Q_2(h(R)) = h(Q_1(R)) \setminus h(Q_2(R))$. Furthermore, since h is an m -homomorphism and by the definition of \setminus we have that $h(Q_1(R)(\bar{t})) \ominus_{\mathbb{K}'} h(Q_2(R)(\bar{t})) = h(Q_1(R)(\bar{t}) \ominus_{\mathbb{K}} Q_2(R)(\bar{t}))$ for every \bar{t} . Hence, $Q(h(R)) = h(Q(R))$.

Conversely, let h be a mapping from \mathbb{K} to \mathbb{K}' . We next show that if for every Q in $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$ and for every R , $Q(h(R)) = h(Q(R))$, then it follows that h is an m -semiring homomorphism. Since $\mathcal{RA}_{\mathcal{K}}^+$ is embedded in $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$, by the result for $\mathcal{RA}_{\mathcal{K}}^+$, h is a semiring homomorphism. Now, suppose by contradiction that h is not an m -semiring homomorphism. Let \bar{Q} and \bar{R} be such that $\bar{Q} = (\pi_A(\sigma_{A=B}(\bar{R})) \setminus \pi_A(\sigma_{A \neq B}(\bar{R})))$ and $\bar{R} = \{(a, a) \mapsto x, (a, b) \mapsto y\}$ for $a \neq b$ and arbitrary $x, y \in \mathbb{K}$. Then, on one hand, $\bar{Q}(h(\bar{R}))$ contains one tuple (a) associated with $h(x) \ominus_{\mathbb{K}'} h(y)$. On the other hand, $h(\bar{Q}(\bar{R}))$ contains one tuple (a) associated with $h(x \ominus_{\mathbb{K}} y)$. Hence, from $\bar{Q}(h(\bar{R})) = h(\bar{Q}(\bar{R}))$, it follows that for every $x, y \in \mathbb{K}$, $h(x) \ominus_{\mathbb{K}'} h(y) = h(x \ominus_{\mathbb{K}} y)$. Clearly, this contradicts the fact that h is not an m -semiring homomorphism. \square

3.2. The query language $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\delta)$

We next extend the positive algebra $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+$ on \mathcal{K} -relations with a family of operators called *constant annotations*. These operators are a generalization of the duplicate elimination operator present in most algebras over bags [16]. The intuition behind these operators is that they are “forgetful”, i.e., they allow to replace all values of tuples in \mathcal{K} -relations by some constant value. Similar to $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+$ and $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\lambda)$, we show that $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\delta)$ satisfies a homomorphism property.

3.2.1. Constant annotations

When considering $\mathcal{K}_{\mathbb{N}}$ -relations it is common to include the duplicate elimination operator δ in the query language. Intuitively, when δ is applied on a bag-relation, the result is a relation with the same support but in which each tuple is counted only once. In the language of \mathcal{K} -relations, $\delta(R)(\bar{t}) = 1$ for all \bar{t} in $\text{supp}(R)$ and $\delta(R)(\bar{t}) = 0$ otherwise.

To introduce duplicate elimination in $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+$ on general \mathcal{K} -relations, we restrict our attention to semirings $\mathcal{K} = (\mathbb{K}, \oplus, \otimes, 0, 1)$ that are *finitely generated*, i.e., every element in \mathbb{K} can be written as a finite sequence of sums and products of a finite set of elements $\mathbf{k}_1, \dots, \mathbf{k}_m$ in \mathbb{K} , called *generators* of \mathbb{K} . We denote a set of generators of \mathbb{K} by $\text{Gen}(\mathbb{K})$ and, for convenience, assume it is minimal.

Example 5. The semirings considered so far are all finitely generated. Indeed, it is easily verified that $\text{Gen}(\mathbb{B}) = \{\text{true}\}$, $\text{Gen}(\mathbb{N}) = \{1\}$, $\text{Gen}(\text{Bool}(X)) = X$, and $\text{Gen}(\mathcal{P}(\Omega)) = \Omega$. The two semirings $\mathcal{K}_{\mathbb{R}}$ and $\mathcal{K}_{\mathbb{R}_{\min}}$ given in Example 4 are not finitely generated since they consist of uncountably many elements.

We now formally define the notion of constant annotations. Given a finitely generated semiring $\mathcal{K} = (\mathbb{K}, \oplus, \otimes, 0, 1)$ with generators $\text{Gen}(\mathbb{K}) = \{\mathbf{k}_1, \dots, \mathbf{k}_m\}$, we define the following set of *constant annotation* operators:

constant annotation If $R : U\text{-Tup} \rightarrow \mathbb{K}$ and \mathbf{k}_i is a generator of \mathbb{K} then $\delta_{\mathbf{k}_i} : U\text{-Tup} \rightarrow \mathbb{K}$ is defined by

$$(\delta_{\mathbf{k}_i}(R))(\bar{t}) = \mathbf{k}_i \quad \text{for each } \bar{t} \in \text{supp}(R) \quad \text{and} \quad (\delta_{\mathbf{k}_i}(R))(\bar{t}) = 0 \quad \text{otherwise.}$$

We denote by $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\delta)$ the query language obtained by extending $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+$ with the constant annotation operators for the semiring \mathcal{K} and set of generators of \mathbb{K} under consideration. Note that for some semirings, e.g., the Boolean semiring, constant annotations do not add expressive power.

3.2.2. The homomorphism property for $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\delta)$

When considering the homomorphism property of queries in $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\delta)$ one has to make the choice of generators in \mathcal{K} and \mathcal{K}' explicit. Let $\text{Gen}(\mathbb{K}) = \{\mathbf{k}_1, \dots, \mathbf{k}_n\}$ and $\text{Gen}(\mathbb{K}') = \{\mathbf{l}_1, \dots, \mathbf{l}_m\}$. We say that a mapping $h : \mathbb{K} \rightarrow \mathbb{K}'$ is a *generator preserving* semiring homomorphism from \mathcal{K} to \mathcal{K}' if h is a semiring homomorphism and furthermore, $h(\text{Gen}(\mathbb{K})) = \text{Gen}(\mathbb{K}')$. Given a query $Q \in \mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\delta)$, let $h(Q)$ be the query in $\mathcal{R}\mathcal{A}_{\mathcal{K}'}^+(\delta)$ obtained by replacing each occurrence of $\delta_{\mathbf{k}_i}$ by $\delta_{h(\mathbf{k}_i)}$. Observe that for generator preserving homomorphisms h , each $\delta_{h(\mathbf{k}_i)}$ is of the form $\delta_{\mathbf{l}_j}$ for some $j = 1, \dots, m$. In other words, $h(Q)$ is well-defined. The following is now easily verified:

Proposition 2. Let \mathcal{K} and \mathcal{K}' be two semirings with generators $\text{Gen}(\mathbb{K})$ and $\text{Gen}(\mathbb{K}')$, respectively. Let $h : \mathbb{K} \rightarrow \mathbb{K}'$ be a mapping. Then, for every query Q in $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\delta)$ and for every R , $h(Q)(h(R)) = h(Q(R))$, if and only if h is a generator-preserving homomorphism from \mathcal{K} to \mathcal{K}' .

3.3. The query language $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\lambda, \delta)$

Finally, we introduce the query language obtained by extending $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+$ with both the difference and constant annotations operators. The resulting language is denoted by $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\lambda, \delta)$. It is easily verified that $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\lambda, \delta)$ satisfies the following homomorphism property:

Proposition 3. Let \mathcal{K} and \mathcal{K}' be two m -semirings with generators $\text{Gen}(\mathbb{K})$ and $\text{Gen}(\mathbb{K}')$, respectively. Let $h : \mathbb{K} \rightarrow \mathbb{K}'$ be a mapping. Then, for every query Q in $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\lambda, \delta)$ and for every R , $h(Q)(h(R)) = h(Q(R))$ if and only if h is a generator-preserving m -semiring homomorphism from \mathcal{K} to \mathcal{K}' .

4. \mathcal{K} -relations and provenance

Besides providing a general framework capturing many data models encountered in the literature, \mathcal{K} -relations are particularly useful for tracking various kinds of provenance information [6,12]. We illustrate this with two examples: the *lineage semiring* and the *provenance semiring*. We refer again to Green et al. [12,11] for more details concerning these and other provenance models. In particular, in this section we recall how to compute the why- and how-provenance for positive

drink	kind	origin	
Stella	beer	Belgium	{x}
Montefalco	wine	Italy	{y}
Pinot	grappa	Italy	{z}

 $R_5 =$

drink	kind	origin	
Pinot	wine	France	{v}
Ardbeg	whiskey	Scotland	{w}

 $R_6 =$

drink	kind		
Stella	beer		{x}
Montefalco	wine		{y}
Montefalco	grappa		{y, z}
Pinot	wine		{y, z, v}
Pinot	grappa		{z}
Ardbeg	whiskey		{w}

 $R_7 =$

Fig. 2. The lineage semiring.

drink	kind	origin	
Stella	beer	Belgium	x
Montefalco	wine	Italy	y
Pinot	grappa	Italy	z

 $\bar{R}_5 =$

drink	kind	origin	
Pinot	wine	France	v
Ardbeg	whiskey	Scotland	w

 $\bar{R}_6 =$

drink	kind		
Stella	beer		x ²
Montefalco	wine		y ²
Montefalco	grappa		yz
Pinot	wine		yz + v
Pinot	grappa		z ²
Ardbeg	whiskey		w

 $R_8 =$

Fig. 3. The provenance semiring.

queries and present m -semirings that allow for computing provenance information in the presence of difference in the relational algebra queries. We conclude this section by describing how to compute provenance in the presence of constant annotations.

4.1. The lineage semiring

Lineage/why-provenance was defined in [5,9] as a way of relating the tuples in a query output to the tuples in the source relations that contribute to them. Let X be a finite set representing the ids of the tuples in the source relations. Then, the lineage semiring $\mathcal{K}_{\text{lin}} = (\mathcal{P}(X), \cup, \cap, \emptyset, \emptyset)$ can be used to represent and compute the why-provenance, as we illustrate in the following example.

Example 6. Consider the \mathcal{K}_{lin} -relations R_5, R_6 shown in Fig. 2, where the set of source tuples ids is $X = \{x, y, z, v, w\}$. In both R_5 and R_6 tuples are annotated with the singleton containing their respective id. Next, let $Q(R', R'')$ be the following query over the relations R' and R'' of schema $U = \{\text{drink}, \text{kind}, \text{origin}\}$:

$$Q(R', R'') = \pi_{\text{drink, kind}}(\pi_{\text{drink, origin}} R' \bowtie \pi_{\text{kind, origin}} R') \cup \pi_{\text{drink, kind}} R''.$$

It is easily verified that R_7 (see Fig. 2) is the query result $Q(R_5, R_6)$. The \mathcal{K}_{lin} -values associated with the tuples in R_7 now provide their why-provenance. For example, they state that the tuple $\bar{s}_p = (\text{Pinot}, \text{wine})$ was obtained from the contribution of the tuples in R_5 and R_6 identified by y, z and v . Note, however, that why-provenance does not provide any information on the how-provenance, e.g., on the way the tuple \bar{s}_p was obtained. In particular, it is not possible to infer from the why-provenance information that \bar{s}_p can be obtained either from joining the tuples identified by y and z together or from the tuple identified by v alone.

4.2. The provenance semiring

In order to overcome the limitations of why-provenance a more powerful provenance semiring was proposed in [12]. This semiring allows to represent and compute the how-provenance of tuples in the query result. More precisely, the (positive algebra) provenance semiring is defined as $\mathcal{K}_{\text{prov}} = (\mathbb{N}[X], +, \times, 0, 1)$, where X is a set of source tuple ids and $\mathbb{N}[X]$ consists of all polynomials with variables taken from X and with coefficients in \mathbb{N} . Hence, $\mathcal{K}_{\text{prov}}$ -relations consist of tuples that are annotated with polynomials. These polynomials are to be interpreted as symbolic expressions over the source tuples ids that describe how the tuples were obtained from the source. This is illustrated in the following example:

Example 7. Consider the $\mathcal{K}_{\text{prov}}$ -relations \bar{R}_5, \bar{R}_6 and R_8 shown in Fig. 3. It can be easily checked that R_8 is the query result $Q(\bar{R}_5, \bar{R}_6)$ for the query Q given in Example 6. Consider again the tuple $\bar{s}_p = (\text{Pinot}, \text{wine})$. The $\mathcal{K}_{\text{prov}}$ -value of \bar{s}_p is the polynomial $R_8(\bar{s}_p) = yz + v$ and states that \bar{s}_p can be obtained either by joining together the tuples in \bar{R}_5 and \bar{R}_6 identified by y and z or by simply using the tuple in \bar{R}_6 identified by v . On the contrary, the tuple $\bar{s}_m = (\text{Montefalco}, \text{grappa})$ can only be obtained by joining together the tuples identified by y and z . Clearly, $\mathcal{K}_{\text{prov}}$ -relations provide more information about the provenance of tuples than \mathcal{K}_{lin} -relations.

drink	kind	origin	
Pinot	wine	France	2
Ardbeq	whiskey	Scotland	1

 $R_9 =$

drink	kind	
Stella	beer	4
Montefalco	wine	1
Montefalco	grappa	1
Pinot	wine	3
Pinot	grappa	1
Ardbeq	whiskey	1

 $R_{10} =$

Fig. 4. The factorization property for $\mathcal{RA}_{\mathcal{K}}^+$.

A nice property of the provenance semiring is that for any semiring \mathcal{K} , to evaluate queries in $\mathcal{RA}_{\mathcal{K}}^+$ on \mathcal{K} -relations it is sufficient to know how to evaluate these queries over $\mathcal{K}_{\text{prov}}$ -relations [12]. This property, called the *factorization property* for $\mathcal{RA}_{\mathcal{K}}^+$, crucially relies on the existence of a universal object in the class of semirings which in this case is precisely the provenance semiring $\mathcal{K}_{\text{prov}} = (\mathbb{N}[X], +, \times, 0, 1)$. More formally, let \mathcal{K} be a semiring, R a \mathcal{K} -relation and $Q \in \mathcal{RA}_{\mathcal{K}}^+$. Suppose that $\text{supp}(R) = \{\bar{t}_1, \dots, \bar{t}_k\}$ and let $X = \{x_1, \dots, x_k\}$ be a set of tuple ids for the tuples in $\text{supp}(R)$. That is, x_i is the tuple id for tuple \bar{t}_i for $i = 1, \dots, k$. Let \bar{R} be the *abstractly tagged version* of R , obtained by letting $\bar{R}(\bar{t}_i) = x_i$ for $\bar{t}_i \in \text{supp}(R)$ and $\bar{R}(\bar{t}) = 0$ otherwise. Let $\nu : X \rightarrow \mathbb{K}$ be the *valuation* that maps x_i to $\bar{R}(\bar{t}_i)$.

Because $\mathcal{K}_{\text{prov}} = (\mathbb{N}[X], +, \times, 0, 1)$ is the free semiring generated by X , we have the property that there exists a unique semiring homomorphism $\text{Eval}_{\nu} : \mathbb{N}[X] \rightarrow \mathbb{K}$ such that for one-variable monomials we have that $\text{Eval}_{\nu}(x) = \nu(x)$. Combined with the homomorphism property for $\mathcal{RA}_{\mathcal{K}}^+$ (see Section 2.3) and observing that $\text{Eval}_{\nu}(\bar{R}) = R$, we recall from [12] that

$$Q(R) = \text{Eval}_{\nu} \circ Q(\bar{R}).$$

In other words, the semantics of queries in $\mathcal{RA}_{\mathcal{K}}^+$ over arbitrary semirings *factors through* its semantics in the provenance semiring.

Example 8. Consider the \mathcal{K}_{lin} -relations R_5 and R_6 shown in Fig. 2. Their respective abstractly tagged versions \bar{R}_5 and \bar{R}_6 are shown in Fig. 3. Consider again the query Q of Example 6. Then, the $\mathcal{K}_{\text{prov}}$ -relation R_8 is the query result $Q(\bar{R}_5, \bar{R}_6)$. Let ν be the valuation that maps η to $\{\eta\}$, for $\eta \in \{x, y, z, v, w\}$. The factorization property then tells us that the \mathcal{K}_{lin} -relation R_7 , shown in Fig. 2, is equal to $\text{Eval}_{\nu}(R_8)$. Indeed, consider the tuple $\bar{s}_p = (\text{Pinot}, \text{grappa})$ annotated with $yz + v$. Then, $\text{Eval}_{\nu}(yz + v) = (\nu(y) \cup \nu(z)) \cup \nu(v) = \{y, z, v\}$, as desired. Similarly, consider the $\mathcal{K}_{\mathbb{N}}$ -relations R_2 shown in Fig. 1 and R_9 shown in Fig. 4. Their abstractly tagged versions \bar{R}_2 and \bar{R}_9 are identical to \bar{R}_5 and \bar{R}_6 , respectively. Let ν be the valuation that maps x and v to 2 and y, z and w to 1. Then the factorization property tells that $Q(R_2, R_9) = R_{10}$, shown in Fig. 4, is equal to $\text{Eval}_{\nu}(R_8)$. Indeed, consider again the tuple \bar{s}_p associated with $yz + v$. In this case we have that $\text{Eval}_{\nu}(yz + v) = (\nu(y) \times \nu(z)) + \nu(v) = 1 + 2 = 3$, as desired.

4.3. The provenance semiring with monus

We next describe how to represent and compute why and how provenance in the presence of difference. It is easily verified that both \mathcal{K}_{lin} and $\mathcal{K}_{\text{prov}}$ can be extended to m -semirings:

Example 9. In the case of \mathcal{K}_{lin} the monus operator simply coincides with set difference. For the provenance semiring, let $X = \{x_1, \dots, x_n\}$ be the set of variables and for $\alpha \in \mathbb{N}^n$, denote by x^α the monomial $x_1^{\alpha_1} x_2^{\alpha_2} \dots x_n^{\alpha_n}$, where by definition $x_i^0 = 1$. Let I be a finite subset of \mathbb{N}^n and let $f[X] = \sum_{\alpha \in I} f_\alpha x^\alpha$ and $g[X] = \sum_{\alpha \in I} g_\alpha x^\alpha$ be two polynomials in $\mathbb{N}[X]$. Then it is easily verified that $f[X] \ominus g[X] = \sum_{\alpha \in I} (f_\alpha \dot{-} g_\alpha) x^\alpha$, where $\dot{-}$ denotes the truncated minus on \mathbb{N} .

Unfortunately, the m -semiring $\mathcal{K}_{\text{prov}} = (\mathbb{N}[X], +, \times, \ominus, 0, 1)$ is not the universal object in the variety of all m -semirings and as a consequence it does not satisfy the factorization property for $\mathcal{RA}_{\mathcal{K}}^+(\wedge)$:

Example 10. Let R_2 be the $\mathcal{K}_{\mathbb{N}}$ -relation shown in Fig. 1 and consider the query

$$Q'(R) = (R \bowtie R) - R.$$

It is easily verified that $Q'(R_2)$ is the $\mathcal{K}_{\mathbb{N}}$ -relation R_{11} shown in Fig. 5. The straightforward generalization of the factorization property to $\mathcal{RA}_{\mathcal{K}}^+(\wedge)$ and using $\mathcal{K}_{\text{prov}}$ as factoring m -semiring would imply that $Q'(R_2)$ can be obtained from the query evaluation $Q'(\bar{R}_2)$ on the abstractly tagged version of R_2 (now interpreted as a $\mathcal{K}_{\text{prov}}$ -relation) and from the valuation ν that maps x to 2, and y, z to 1. The $\mathcal{K}_{\text{prov}}$ -relation $Q'(\bar{R}_2)$ is shown as relation R_{12} in Fig. 5. Here, each tuple is associated with $\eta^2 \ominus \eta = (0 \cdot \eta + 1 \cdot \eta^2) \ominus (1 \cdot \eta + 0 \cdot \eta^2) = (0 \dot{-} 1) \cdot \eta + (1 \dot{-} 0) \cdot \eta^2 = \eta^2$, for some id $\eta \in \{x, y, z\}$. Then, $Q'(R_2) = R_{11} \neq \text{Eval}_{\nu}(R_{12}) = R_{13}$. It is easily verified that a similar counterexample works when we consider the $\mathcal{K}_{\mathbb{B}}$ -relation R_1 shown in Fig. 1 and query Q' . Indeed, in this case $Q'(R_1)$ returns the empty relation, i.e., all tuples are associated with false. On the contrary, if we consider the valuation ν maps x and y to true, then we have that $\text{Eval}_{\nu}(Q'(\bar{R}_1))$ contains two tuples associated with $\nu(x^2) = \nu(x) \wedge \nu(x) = \text{true}$ and $\nu(y^2) = \nu(y) \wedge \nu(y) = \text{true}$, respectively.

$R_{11} =$	<table border="1" style="border-collapse: collapse; text-align: center;"> <thead> <tr><th>drink</th><th>kind</th><th>origin</th><th></th></tr> </thead> <tbody> <tr><td>Stella</td><td>beer</td><td>Belgium</td><td>2</td></tr> <tr><td>Montefalco</td><td>wine</td><td>Italy</td><td>0</td></tr> <tr><td>Pinot</td><td>grappa</td><td>Italy</td><td>0</td></tr> </tbody> </table>	drink	kind	origin		Stella	beer	Belgium	2	Montefalco	wine	Italy	0	Pinot	grappa	Italy	0
drink	kind	origin															
Stella	beer	Belgium	2														
Montefalco	wine	Italy	0														
Pinot	grappa	Italy	0														

$R_{12} =$	<table border="1" style="border-collapse: collapse; text-align: center;"> <thead> <tr><th>drink</th><th>kind</th><th>origin</th><th></th></tr> </thead> <tbody> <tr><td>Stella</td><td>beer</td><td>Belgium</td><td>x^2</td></tr> <tr><td>Montefalco</td><td>wine</td><td>Italy</td><td>y^2</td></tr> <tr><td>Pinot</td><td>grappa</td><td>Italy</td><td>z^2</td></tr> </tbody> </table>	drink	kind	origin		Stella	beer	Belgium	x^2	Montefalco	wine	Italy	y^2	Pinot	grappa	Italy	z^2
drink	kind	origin															
Stella	beer	Belgium	x^2														
Montefalco	wine	Italy	y^2														
Pinot	grappa	Italy	z^2														

$R_{13} =$	<table border="1" style="border-collapse: collapse; text-align: center;"> <thead> <tr><th>drink</th><th>kind</th><th>origin</th><th></th></tr> </thead> <tbody> <tr><td>Stella</td><td>beer</td><td>Belgium</td><td>4</td></tr> <tr><td>Montefalco</td><td>wine</td><td>Italy</td><td>1</td></tr> <tr><td>Pinot</td><td>grappa</td><td>Italy</td><td>1</td></tr> </tbody> </table>	drink	kind	origin		Stella	beer	Belgium	4	Montefalco	wine	Italy	1	Pinot	grappa	Italy	1
drink	kind	origin															
Stella	beer	Belgium	4														
Montefalco	wine	Italy	1														
Pinot	grappa	Italy	1														

Fig. 5. The failure of the factorization property for $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$ and $\mathcal{K}_{\text{dprov}}$.

We next show how a factorization property for $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$ can be obtained. Indeed, from universal algebra it follows that there exists a unique free m -semiring. We next describe the construction of this semiring and then show how it can be used to represent and compute provenance for $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$.

First, we observe that the class of m -semirings is an equational variety. Indeed, an algebraic structure $(\mathbb{K}, \oplus, \otimes, \ominus, 0, 1)$ is an m -semiring iff it satisfies (i) the defining equations of $(\mathbb{K}, \oplus, \otimes, 0, 1)$ being a semiring; and (ii) the defining equations of $(\mathbb{K}, \oplus, \ominus, 0)$ being a commutative monoid with monus [3]. Hence, by Birkhoff's Theorem, the class of m -semirings is indeed a variety and furthermore admits free objects [7].

We recall the standard universal algebra construction for the unique free object $T[X]$ generated by $X = \{x_1, \dots, x_n\}$ in the equational variety of m -semirings [7]. In a nutshell, elements of $T[X]$ consist of terms constructed inductively as follows: $x_i, 1$ and 0 are terms; and moreover, if t and s are terms then so are $(t \oplus s)$, $(t \ominus s)$ and $(t \otimes s)$; and finally, nothing else is a term.

We next need the notion of congruence relation. A congruence relation C over $T[X]$ is an equivalence relation over $T[X]$ that is compatible with \oplus , \otimes and \ominus , i.e., if $C(s_1, t_1)$ and $C(s_2, t_2)$ then also $C(s_1 \text{ op } s_2, t_1 \text{ op } t_2)$ for $\text{op} \in \{\oplus, \otimes, \ominus\}$. We next specialize C to correspond to the congruence relation that identifies terms based on the equations of m -semirings. It is then easily verified that the quotient structure $T[X]/C$ that consists of expressions in $T[X]$ in which any two equivalent expressions are identified (as specified by C), is indeed an m -semiring. Furthermore, it follows that $T[X]/C$ is the free m -semiring generated by X [7]. Hence, for any m -semiring \mathcal{K} and any valuation $\nu : X \rightarrow \mathbb{K}$, we have that ν can be lifted to an m -semiring homomorphism $\text{Eval}_{\nu} : T[X]/C \rightarrow \mathbb{K}$ that coincides with ν on X . We denote by $\mathcal{K}_{\text{dprov}}$ the free m -semiring $(T[X]/C, \oplus, \otimes, \ominus, 0, 1)$ obtained in this way.

The following example illustrates $\mathcal{K}_{\text{dprov}}$ and its corresponding factorization property.

Example 11. Consider again the relation \bar{R}_2 (which is equal to \bar{R}_5 shown in Fig. 3). This can obviously be seen as a $\mathcal{K}_{\text{dprov}}$ relation. Let Q' be the query of Example 10. It is easily verified that the $\mathcal{K}_{\text{dprov}}$ -relation $Q'(\bar{R}_2)$ is similar to the relation R_{12} shown in Fig. 5, except that each tuple is now associated with $(\eta \otimes \eta) \ominus \eta$ for $\eta \in \{x, y, z\}$. If we consider the valuation ν that maps x to 2 and y, z to 1 and extend ν to an m -homomorphism $\text{Eval}_{\nu} : T[X]/C \rightarrow \mathbb{N}$ in the natural way, then $Q'(\bar{R}_2) = R_{11} = \text{Eval}_{\nu}(Q'(\bar{R}_2))$. Indeed, this follows from the fact that $\text{Eval}_{\nu}((\eta \otimes \eta) \ominus \eta) = (\nu(\eta) \times \nu(\eta)) \dot{-} \nu(\eta)$. Similarly, if we consider the valuation ν that maps x and y to true and let $\text{Eval}_{\nu} : T[X]/C \rightarrow \mathbb{B}$, then $Q'(\bar{R}_1) = \text{Eval}_{\nu}(Q'(\bar{R}_1))$. This follows again from the fact that $\text{Eval}_{\nu}((\eta \otimes \eta) \ominus \eta) = (\nu(\eta) \wedge \nu(\eta)) \ominus \nu(\eta) = \nu(\eta) \wedge \bar{\nu}(\eta) = \text{false}$, for $\eta \in \{x, y\}$.

The following proposition is an immediate consequence of Proposition 1 and the fact that $\mathcal{K}_{\text{dprov}}$ is a free m -semiring over X :

Proposition 4. Let \mathcal{K} be an m -semiring. For any query $Q \in \mathcal{RA}_{\mathcal{K}}^+(\setminus)$ and any \mathcal{K} -relation R with tuple id set X , $Q(R) = \text{Eval}_{\nu} \circ Q(\bar{R})$, where \bar{R} denotes the $\mathcal{K}_{\text{dprov}}$ -relation obtained by tagging each tuple in R with its own tuple id.

4.4. The provenance semiring with monus and constant annotations

We can easily extend the construction of the provenance m -semiring $\mathcal{K}_{\text{dprov}}$ to obtain an extended provenance m -semiring for $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$ for which a factorization property holds. We first note that the provenance semirings discussed in this and other papers [12,11] are all finitely generated. Similarly for the extended provenance m -semiring described next.

In a nutshell, this m -semiring is constructed in the same way as $\mathcal{K}_{\text{dprov}}$, with the proviso that if t is a term of the m -semiring, then so are $\delta_{y_i}(t)$ for $y_i \in Y$. Here, Y is a set of variables disjoint from X . Intuitively, the factorization property holds also for $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$, after extending the valuation also to variables in Y . Formally, let \mathcal{K} be a finitely generated m -semiring with $\text{Gen}(\mathbb{K}) = \{\mathbf{k}_1, \dots, \mathbf{k}_n\}$. Let R be \mathcal{K} -relation and Q be a query in $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$. Let Y be a set of n fresh variables y_i , one for each generator in \mathcal{K} , and let ν be the valuation of $X \cup Y$ that maps, as before, x_i to $R(\bar{t}_i)$ and y_i to \mathbf{k}_i . Furthermore, we define Q' to be Q in which each occurrence of $\delta_{\mathbf{k}_i}$ is replaced by δ_{y_i} . Then, $Q(R) = \text{Eval}_{\nu} \circ Q'(\bar{R})$ where \bar{R} is viewed as an extended provenance m -semiring relation.

$$S_1 = \begin{array}{|c|c|c|} \hline A & B & \\ \hline a & a & 2 \\ b & b & 2 \\ \hline \end{array} \quad S_2 = \begin{array}{|c|c|c|} \hline A & B & \\ \hline a & a & 1 \\ b & b & 2 \\ \hline \end{array} \quad S_3 = \begin{array}{|c|c|c|} \hline A & B & \\ \hline a & a & 2 \\ b & b & 2 \\ \hline \end{array} \quad S_4 = \begin{array}{|c|c|c|} \hline A & B & \\ \hline a & a & 1 \\ b & b & 1 \\ \hline \end{array} \quad S_5 = \begin{array}{|c|c|c|} \hline A & B & \\ \hline a & a & 2 \\ b & b & 1 \\ \hline \end{array}$$
Fig. 6. Example $\mathcal{K}_{\mathbb{N}}$ -relations.

5. BP-completeness for \mathcal{K} -relations

In this section, we initiate our study of the completeness of query languages over \mathcal{K} -relations in the sense of Bancilhon and Paredaens [4,18]. First, recall that Codd qualified a query language on standard relational databases as *complete* if its expressive power is at least that of the relational calculus [8]. Bancilhon [4] and Paredaens [18] independently provided a *language-independent* characterization of completeness. This characterization, now known as *BP-completeness*, can be stated as follows: a relation T is the result of a *generic* relational algebra query applied to a database S if and only if (i) the active domain of T is included in the domain of S ; and (ii) every automorphism of S is also an automorphism of T . In fact, Paredaens [18] observed that once inequality conditions are allowed in the selection predicate, one does not require difference in the relational algebra for it to be BP-complete.

Recall that a generic query is one which is oblivious to the constants appearing in the relation, *i.e.*, for any permutation τ of the domain \mathbb{D} , we have that $Q(\tau(R)) = \tau(Q(R))$. Furthermore, an automorphism of a relation R is a permutation τ of \mathbb{D} that leaves R invariant, *i.e.*, for any $\bar{t} \in R$, $\tau(\bar{t}) \in R$. Hence, intuitively, the set of automorphisms of a relation R , denoted by $\text{Aut}(R)$, allows to identify values that are “indistinguishable” for the relation, *i.e.* values that can be switched without changing the relation itself.

In order to study BP-completeness in the setting of \mathcal{K} -relations, we first need to define the notion of *automorphism of a \mathcal{K} -relation*. Given that \mathcal{K} -relations are annotated relations, by analogy to the case of standard relations, \mathcal{K} -relations should allow to identify values in the support that can be switched without changing neither the tuples, nor the respective tuples annotations. That is, apart from being an automorphism of the underlying relational database, an automorphism of a \mathcal{K} -relation should additionally preserve the semiring values associated with the tuples. Hence, formally, the set of *automorphisms* of R , denoted by $\text{Aut}_{\mathcal{K}}(R)$, is defined as

$$\text{Aut}_{\mathcal{K}}(R) = \{ \tau \mid \tau \in \text{Aut}(\text{supp}(R)) \text{ and } R(\tau(\bar{t})) = R(\bar{t}), \forall \bar{t} \in \mathbb{D}^n \}.$$

Example 12. Consider the relations given in Fig. 6 and assume that $\mathbb{D} = \{a, b\}$. When considering the underlying standard relations, *i.e.*, ignoring the annotations, we have that $\text{Aut}(S_1) = \text{Aut}(S_2) = \text{Aut}(S_4) = \text{Aut}(S_5) = \{(a \mapsto a, b \mapsto b), (a \mapsto b, b \mapsto a)\}$ and $\text{Aut}(S_3) = \{(a \mapsto a, b \mapsto b)\}$. When viewed as $\mathcal{K}_{\mathbb{N}}$ -relations, however, with the multiplicities of each tuple shown in the last column, we have that $\text{Aut}_{\mathcal{K}}(S_1) = \text{Aut}_{\mathcal{K}}(S_4) = \{(a \mapsto a, b \mapsto b), (a \mapsto b, b \mapsto a)\}$ and $\text{Aut}_{\mathcal{K}}(S_2) = \text{Aut}_{\mathcal{K}}(S_5) = \text{Aut}_{\mathcal{K}}(S_3) = \{(a \mapsto a, b \mapsto b)\}$.

The set of \mathcal{K} -relations that are *preserved* by $\text{Aut}_{\mathcal{K}}(R)$, denoted by $\text{Inv}_{\mathbb{D}}(R)$, is defined as:

$$\text{Inv}_{\mathbb{D}}(R) = \{ S \mid \text{adom}(S) \subseteq \text{adom}(R), \text{Aut}_{\mathcal{K}}(R) \subseteq \text{Aut}_{\mathcal{K}}(S) \}.$$

Example 13. Consider again the relations given in Fig. 6. From the definition above, it follows that $\text{Inv}_{\mathbb{D}}(S_1) = \text{Inv}_{\mathbb{D}}(S_4) \subseteq \text{Inv}_{\mathbb{D}}(S_2) = \text{Inv}_{\mathbb{D}}(S_5)$ and moreover, $\text{Inv}_{\mathbb{D}}(S_3) \subseteq \text{Inv}_{\mathbb{D}}(S_i)$ for $i \in \{2, 5\}$. In particular, $S_3 \in \text{Inv}_{\mathbb{D}}(S_i)$ for $i \in \{2, 5\}$.

Finally, the expressiveness of a query language can be described in terms of the “information” that can be deduced from a \mathcal{K} -relation using queries in that query language. Following Paredaens [18] we define: Let \mathcal{Q} be a query language and R a \mathcal{K} -relation, then the *basic information* of R with respect to \mathcal{Q} is the set of \mathcal{K} -relations:

$$\text{BI}(R, \mathcal{Q}) = \{ S \mid Q(R) = S \text{ for some generic query } Q \in \mathcal{Q} \}.$$

Finally, BP-completeness links the notions of basic information and invariant relations together:

Definition 2. A query language \mathcal{Q} is *BP-complete* if $\text{BI}(R, \mathcal{Q}) = \text{Inv}_{\mathbb{D}}(R)$ for all \mathcal{K} -relations R .

It is worth noting that the above definitions coincide with the standard notions in the relational setting under the set semantics, *i.e.*, when considering $\mathcal{K} = \mathcal{K}_{\mathbb{B}}$.

We first study BP-completeness for $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+$. A straightforward induction on the structure of queries in $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+$ shows that the inclusion of $\text{BI}(R, \mathcal{R}\mathcal{A}_{\mathcal{K}}^+) \subseteq \text{Inv}_{\mathbb{D}}(R)$ holds for any semiring \mathcal{K} and \mathcal{K} -relation R :

Lemma 1. For any semiring \mathcal{K} , any (generic) $Q \in \mathcal{R}\mathcal{A}_{\mathcal{K}}^+$ and any \mathcal{K} -relation R , we have that

- (i) $\text{adom}(Q(R)) \subseteq \text{adom}(R)$ and
- (ii) $\text{Aut}_{\mathcal{K}}(R) \subseteq \text{Aut}_{\mathcal{K}}(Q(R))$.

The other direction, *i.e.*, whether $\text{Inv}_{\mathbb{D}}(R) \subseteq \text{BI}(R, \mathcal{R}\mathcal{A}_{\mathcal{K}}^+)$ holds for any semiring \mathcal{K} and \mathcal{K} -relation R , is not true. Indeed, a counterexample can be found for the semiring $\mathcal{K}_{\mathbb{N}}$.

Proposition 5. *There exists a semiring \mathcal{K} such that $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+$ is not BP-complete on \mathcal{K} -relations.*

Proof. Let \mathcal{K} be the semiring $\mathcal{K}_{\mathbb{N}}$ and consider the relations S_1 and S_4 in Fig. 6. From Example 13 we know that $S_4 \in \text{Inv}_{\mathbb{D}}(S_1)$. However, it is easily verified, by induction on the structure of queries, that for every generic query $Q \in \mathcal{R}\mathcal{A}_{\mathcal{K}}^+$, $Q(S_1)$ is either: (i) empty, or (ii) the empty tuple, or (iii) such that it contains only tuples having *even multiplicity*. In other words, S_4 cannot be the result of a generic query in $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+$, *i.e.* $S_4 \notin \text{BI}(S_1, \mathcal{R}\mathcal{A}_{\mathcal{K}}^+)$. \square

In fact, the counterexample in the previous proof shows that more expressive power is needed to make $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+$ BP-complete for the semiring $\mathcal{K}_{\mathbb{N}}$. It is easy to see that considering the language $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\delta)$ obtained by adding constant annotation operators to $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+$, resolves the previous counterexample for $\mathcal{K}_{\mathbb{N}}$. Indeed, $S_4 = \delta_1(S_1)$ and therefore $S_4 \in \text{BI}(S_1, \mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\delta))$ for $\mathcal{K}_{\mathbb{N}}$. It turns out, however, that the query language $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\delta)$ is still not BP-complete for arbitrary finitely generated semirings.

Proposition 6. *There exists a finitely generated semiring \mathcal{K} such that $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\delta)$ is not BP-complete on \mathcal{K} -relations.*

Proof. Let \mathcal{K} be the semiring $\mathcal{K}_{\mathbb{N}}$ and consider the relations S_2 and S_5 in Fig. 6. From Example 13 we know that $S_5 \in \text{Inv}_{\mathbb{D}}(S_2)$. It is easily verified, however, that for any generic query $Q \in \mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\delta)$, the query result $Q(S_2)$ satisfies the property that for any two tuples \bar{t}_1 and \bar{t}_2 in $Q(S_2)$, \bar{t}_1 occurs with less or equal multiplicity than \bar{t}_2 if and only if \bar{t}_1 contains a less or equal number of b 's than \bar{t}_2 . Hence, $S_5 \notin \text{BI}(S_2, \mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\delta))$. \square

The counterexample in the previous proof can, however, be resolved when considering $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\setminus, \delta)$ instead of $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\delta)$. Indeed, it is easily verified that for the m -semiring $\mathcal{K}_{\mathbb{N}} = (\mathbb{N}, +, \times, \dot{-}, 0, 1)$,

$$S_5 = ((\delta_1(S_2) \cup \delta_1(S_2)) \setminus S_2) \cup \delta_1(S_2).$$

In other words, in this case, we have that $S_5 \in \text{BI}(S_2, \mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\setminus, \delta))$.

At this point, one may wonder whether the extension of $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+$ with difference alone, *i.e.*, $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\setminus)$ results in a BP-complete language over arbitrary m -semirings. The proof of Proposition 5, however, carries through for $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\setminus)$. Hence, $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\setminus)$ is not BP-complete for arbitrary m -semirings.

We next show that the fact $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\setminus, \delta)$ resolves both counterexamples given in the proofs of Propositions 5 and 6 is not a coincidence.

Theorem 1. *The query language $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\setminus, \delta)$ is BP-complete on \mathcal{K} -relations for all finitely generated m -semirings \mathcal{K} .*

Proof. We first observe that Lemma 1 extends to $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\setminus, \delta)$ for any finitely generated m -semiring \mathcal{K} and any \mathcal{K} -relation R . Indeed, a straightforward induction on the queries in $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\setminus, \delta)$ shows that $\text{BI}(R, \mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\setminus, \delta)) \subseteq \text{Inv}_{\mathbb{D}}(R)$ for any \mathcal{K} -relation R .

For the opposite direction, *i.e.*, given a \mathcal{K} -relation R , whether $\text{Inv}_{\mathbb{D}}(R) \subseteq \text{BI}(R, \mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\setminus, \delta))$ holds, we show that for any \mathcal{K} -relation $S \in \text{Inv}_{\mathbb{D}}(R)$, there exists a generic query $Q \in \mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\setminus, \delta)$ such that $Q(R) = S$. In other words, we show that $S \in \text{BI}(R, \mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\setminus, \delta))$.

Let R and S be \mathcal{K} -relations and assume that $S \in \text{Inv}_{\mathbb{D}}(R)$. The desired query $Q \in \mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\setminus, \delta)$ such that $Q(R) = S$ is constructed in a number of steps:

First, we define a query $Q_{\text{Aut}} \in \mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\setminus, \delta)$ such that $Q_{\text{Aut}}(R) = \text{Aut}_{\mathcal{K}}(R)$. Strictly speaking, $Q_{\text{Aut}}(R)$ returns permutations of $\text{adom}(R)$ instead of permutation of \mathbb{D} . This, however, is sufficient for our purpose since $\text{adom}(S) \subseteq \text{adom}(R)$. More specifically, let (a_1, \dots, a_n) be a tuple that represents all values in $\text{adom}(R)$. Then $Q_{\text{Aut}}(R)$ consists of all tuples $(b_1, \dots, b_n) \in \text{adom}(R)^n$ such that the mapping $\tau(a_i) = b_i$, for $i = 1, \dots, n$, is an automorphism of R , *i.e.*, $\tau \in \text{Aut}_{\mathcal{K}}(R)$. Observe that (a_1, \dots, a_n) is always present in $Q_{\text{Aut}}(R)$ since it corresponds to the trivial automorphism of R . The query $Q_{\text{Aut}}(R)$ is constructed as follows. Assume that $\text{supp}(R) = \{\bar{t}_1, \dots, \bar{t}_p\}$ and denote the corresponding \mathcal{K} -values by $\ell_i = R(\bar{t}_i)$, for $i = 1, \dots, p$. For $i = 1, \dots, p$, we construct the following queries:

- Q_{ℓ_i} : A query such that $Q_{\ell_i}(R)(\bar{t}) = \ell_i$ for all $\bar{t} \in \text{supp}(R)$ and $Q_{\ell_i}(R)(\bar{t}) = 0$ otherwise. This query can be expressed in $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\setminus, \delta)$ using the constant annotation operators; indeed, these operators allow to generate arbitrary \mathcal{K} -values and assign them to tuples in a relation. In particular, one can assign each tuple in R the constant value ℓ_i .
- $Q_{=\ell_i}$: A query such that $Q_{=\ell_i}(R)(\bar{t}) = \ell_i$ if $R(\bar{t}) = \ell_i$ and $Q_{=\ell_i}(R)(\bar{t}) = 0$ otherwise. That is, this query extracts all tuples \bar{t} from R that satisfy $R(\bar{t}) = \ell_i$. This query is expressible in $\mathcal{R}\mathcal{A}_{\mathcal{K}}^+(\setminus, \delta)$. Indeed, we claim that $Q_{=\ell_i} = (R \setminus Q_{\ell_i \neq}) \setminus Q_{\neq \ell_i}$, where $Q_{\ell_i \neq} = R \bowtie Q_1(Q_{\ell_i}(R) \setminus R)$ and $Q_{\neq \ell_i} = R \bowtie Q_1(R \setminus Q_{\ell_i}(R))$. Here, Q_{ℓ_i} is the query previously constructed and

Q_1 is a similar query except that it replaces all semiring values of the tuples with the unit value $\mathbf{1}$ in \mathcal{K} instead of with ℓ_i . To show the correctness of this query we first observe that $R(\bar{t}) = \ell_i$ iff both $R(\bar{t}) \preceq \ell_i$ and $\ell_i \preceq R(\bar{t})$. This follows from the fact we consider m -semirings for which the natural order \preceq is, by assumption, a partial order. Consider first the query $Q_{\ell_i \not\preceq}$. It is easily verified that $Q_{\ell_i \not\preceq}(R)(\bar{t}) = R(\bar{t})$ if $\ell_i \not\preceq R(\bar{t})$ and $Q_{\ell_i \not\preceq}(R)(\bar{t}) = 0$ otherwise. Similarly, $Q_{\not\preceq \ell_i}(R)(\bar{t}) = R(\bar{t})$ if $R(\bar{t}) \not\preceq \ell_i$ and $Q_{\not\preceq \ell_i}(R)(\bar{t}) = 0$ otherwise. As a consequence, $(R \setminus Q_{\ell_i \not\preceq})(\bar{t}) = R(\bar{t})$ if $\ell_i \preceq R(\bar{t})$ and $(R \setminus Q_{\not\preceq \ell_i})(\bar{t}) = 0$ otherwise. From this and the definition of $Q_{\not\preceq \ell_i}$ the correctness of $Q_{=\ell_i}$ then follows.

- Q_{Aut}^S : A query that computes the automorphisms of domain values in $\text{adom}(R)$ of a \mathcal{K} -relation in which each tuple is assigned the same \mathcal{K} -value. Observe that for such \mathcal{K} -relations T , $\text{Aut}_{\mathcal{K}}(T) = \text{Aut}(\text{supp}(T))$. As a consequence, the query Q_{Aut}^S can be expressed in the same way as for the classical relational case [18].

Finally, we obtain $Q_{\text{Aut}}(R)$ by taking the intersection of $Q_{\text{Aut}}^S(Q_{=\ell_i}(R))$ for $i = 1, \dots, p$. That is,

$$Q_{\text{Aut}}(R) = Q_{\text{Aut}}^S(Q_{=\ell_1}(R)) \bowtie \dots \bowtie Q_{\text{Aut}}^S(Q_{=\ell_p}(R)).$$

It is easily verified that $Q_{\text{Aut}}(R) = \text{Aut}_{\mathcal{K}}(R)$, as desired.

We next proceed as follows. First we define a query Q^S in $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$ such that $\text{supp}(Q^S(R)) = \text{supp}(S)$, i.e., $Q^S(R)$ and S agree as standard relations. Second, we show how Q^S can be modified into a query Q in $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$ such that $Q(R) = S$, i.e., $Q(R)$ and S also agree as \mathcal{K} -relations.

- By assumption we have that $\text{adom}(S) \subseteq \text{adom}(R)$. Furthermore, recall that $\text{Aut}_{\mathcal{K}}(R)$ contains a tuple (a_1, \dots, a_n) that represents all values in $\text{adom}(R)$. Let $\bar{s} \in \text{supp}(S)$. It is clear that $\bar{s} \in \text{supp}(\tilde{\pi}_{\bar{s}}(Q_{\text{Aut}}(R)))$, where $\tilde{\pi}_{\bar{s}}$ stands for an appropriate generalized projection. Recall that a generalized projection is a projection in which the same attribute can be repeated several times. This operator can be simulated using the standard projection and join operator and therefore does not add to the expressive power of $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$. For instance, suppose that $\text{adom}(R) = \{a, b, c\}$ is represented by the tuple (a, b, c) in $\text{Aut}_{\mathcal{K}}(R)$. Furthermore, assume that $\bar{s} = (a, b, b)$. Then $\bar{s} \in \text{supp}(\tilde{\pi}_{1,2,2}(Q_{\text{Aut}}(R)))$. By assumption we also have that $\text{Aut}_{\mathcal{K}}(R) \subseteq \text{Aut}_{\mathcal{K}}(S)$. As a consequence, for each $\bar{s} \in \text{supp}(S)$ we have that $\text{supp}(\tilde{\pi}_{\bar{s}}(Q_{\text{Aut}}(R))) \subseteq \text{supp}(S)$. In other words, $\text{supp}(S) = \bigcup_{\bar{s} \in S} \text{supp}(\tilde{\pi}_{\bar{s}}(Q_{\text{Aut}}(R)))$. Finally, we observe that for any two tuples $\bar{s}, \bar{t} \in S$, if $\bar{s} \in \text{supp}(\tilde{\pi}_{\bar{t}}(Q_{\text{Aut}}(R)))$ then $\text{supp}(\tilde{\pi}_{\bar{t}}(Q_{\text{Aut}}(R))) = \text{supp}(\tilde{\pi}_{\bar{s}}(Q_{\text{Aut}}(R)))$. As a consequence, $\text{supp}(S)$ can be partitioned as

$$\text{supp}(S) = \text{supp}(\tilde{\pi}_{\bar{s}_1}(Q_{\text{Aut}}(R))) \uplus \dots \uplus \text{supp}(\tilde{\pi}_{\bar{s}_r}(Q_{\text{Aut}}(R))),$$

for some tuples $\bar{s}_i \in \text{supp}(S)$, $i = 1, \dots, r$. We define $Q^S = \bigcup_{i=1}^r \tilde{\pi}_{\bar{s}_i}(Q_{\text{Aut}}(R))$. Clearly this query satisfies $\text{supp}(S) = \text{supp}(Q^S(R))$.

- We next show how to modify Q^S into Q such that $S = Q(R)$. Observe that it only remains to correctly set the \mathcal{K} -values of the tuples in $Q^S(R)$. For this, we observe that since $\text{Aut}_{\mathcal{K}}(R) \subseteq \text{Aut}_{\mathcal{K}}(S)$, we have that for each $i = 1, \dots, r$, all tuples in $\text{supp}(\tilde{\pi}_{\bar{s}_i}(Q_{\text{Aut}}(R)))$ have the same \mathcal{K} -value in S , say μ_i . We therefore use the constant annotation operators in $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$ to set, for each $i = 1, \dots, r$, the \mathcal{K} -value of the tuples in $\tilde{\pi}_{\bar{s}_i}(Q_{\text{Aut}}(R))$ to μ_i . It is now easily verified that the query $Q = Q_{\mu_1}(\tilde{\pi}_{\bar{s}_1}(Q_{\text{Aut}}(R))) \cup \dots \cup Q_{\mu_r}(\tilde{\pi}_{\bar{s}_r}(Q_{\text{Aut}}(R)))$ satisfies $Q(R) = S$, i.e., Q is the desired query. \square

It is interesting to observe that in case of $\mathcal{K}_{\mathbb{B}} = (\mathbb{B}, \vee, \wedge, \ominus, \text{false}, \text{true})$, i.e., when considering the standard relational algebra with the set semantics, the construction of Q in the previous proof reduces to the construction given by Paredaens [18]. More specifically, neither difference nor duplicate elimination are needed in this case to obtain BP-completeness, in accordance with the results in [18].

Example 14. Consider the relations S_3 and S_5 given in Fig. 6. When viewed as $\mathcal{K}_{\mathbb{N}}$ -relations, Theorem 1 guarantees the existence of a query Q in $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$ such that $Q(S_5) = S_3$. Although the query constructed in the proof of Theorem 1 is such a query, this query is by no means the unique (and most elegant) query with this property. Indeed, it is easily verified that $S_3 = ((\delta_1(S_5) \cup \delta_1(S_5)) \setminus S_5) \cup ((\delta_1(S_5) \cup \delta_1(S_5)) \setminus S_5)$.

6. Conclusion

In view of the lack of expressive power of $\mathcal{RA}_{\mathcal{K}}^+$, we extended $\mathcal{RA}_{\mathcal{K}}^+$ with a difference operator, resulting in the query language $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$, constant annotation operators δ , resulting in the query language $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$, and both operators resulting in $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$. We proposed extended provenance semirings for $\mathcal{RA}_{\mathcal{K}}^+(\setminus)$ and $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$ and established crucial properties of the newly defined query languages, in particular the factorization property. This naturally extends previous work on the positive relational algebra. Finally, we initiated the study of BP-completeness of query languages on \mathcal{K} -relations. In particular, we showed that for some semirings \mathcal{K} , $\mathcal{RA}_{\mathcal{K}}^+$ is not BP-complete. Our main result is that $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$ is BP-complete on \mathcal{K} -relations for a general class of semirings \mathcal{K} . More specifically, $\mathcal{RA}_{\mathcal{K}}^+(\setminus, \delta)$ is BP-complete for semirings that can be extended with a monus operator and that are finitely generated. This class of semirings covers most of the semirings considered in the database literature so far. We also showed that neither the difference nor duplicate elimination can be omitted while still retaining BP-completeness.

In future work, we plan to find an exact characterization of when two \mathcal{K} -relations are related by means of a query in $\mathcal{RA}_{\mathcal{K}}^+$ and establish the complexity of deciding this problem. Also, it is interesting to study the semantics of $\mathcal{RA}_{\mathcal{K}}^+(\cdot, \delta)$ for provenance models different than the why- and how-provenance. Finally, in the spirit of Codd, it is challenging to find a characterization of the completeness of $\mathcal{RA}_{\mathcal{K}}^+$ and extensions thereof in terms of first-order logic.

Acknowledgements

We would like to thank Leonid Libkin for helpful discussions, Jan Van den Bussche, Todd J. Green and Val Tannen for comments on a preliminary version of this paper.

References

- [1] S. Abiteboul, R. Hull, V. Vianu, Foundations of Databases, Addison–Wesley, 1995.
- [2] A.V. Aho, J.D. Ullman, Universality of data retrieval languages, in: POPL '79, ACM, 1979, pp. 110–119.
- [3] K. Amer, Equationally complete classes of commutative monoids with monus, Algebra Universalis 18 (1) (1984) 129–131.
- [4] F. Bancilhon, On the completeness of query languages for relational data bases, in: MFCS '79, in: Lecture Notes in Computer Science, vol. 64, Springer, 1978, pp. 112–123.
- [5] P. Buneman, S. Khanna, W.C. Tan, Why and where: A characterization of data provenance, in: ICDT '01, in: Lecture Notes in Computer Science, vol. 1973, Springer, 2001, pp. 316–330.
- [6] P. Buneman, W.C. Tan, Provenance in databases, in: SIGMOD '07, ACM, 2007, pp. 1171–1173.
- [7] S. Burris, H. Sankappanavar, A Course in Universal Algebra, Springer-Verlag, 1981.
- [8] E.F. Codd, Relational completeness of data base sublanguages, IBM Research Report RJ 987, San Jose, California.
- [9] Y. Cui, J. Widom, J.L. Wiener, Tracing the lineage of view data in a warehousing environment, ACM TODS 25 (2) (2000) 179–227.
- [10] N. Fuhr, T. Rölleke, A probabilistic relational algebra for the integration of information retrieval and database systems, ACM Trans. Inf. Syst. 15 (1) (1997) 32–66.
- [11] T.J. Green, Containment of conjunctive queries on annotated relations, in: ICDT, 2009, pp. 296–309.
- [12] T.J. Green, G. Karvounarakis, V. Tannen, Provenance semirings, in: PODS '07, ACM, 2007, pp. 31–40.
- [13] T.J. Green, V. Tannen, Models for incomplete and probabilistic information, IEEE Data Eng. Bull. 29 (1) (2006) 17–24.
- [14] M. Henriksen, J.R. Isbell, Lattice-ordered rings and function rings, Pacific J. Math. 12 (1962) 533–565.
- [15] T. Imieliński, J.W. Lipski, Incomplete information in relational databases, J. ACM 31 (4) (1984) 761–791.
- [16] L. Libkin, L. Wong, Query languages for bags and aggregate functions, J. Comput. Syst. Sci. 55 (2) (1997) 241–272.
- [17] F. Montagna, V. Sebastiani, Equational fragments of systems for arithmetic, Algebra Universalis 46 (3) (2001) 417–441.
- [18] J. Paredaens, On the expressive power of the relational algebra, Inf. Process. Lett. 7 (2) (1978) 107–111.
- [19] E. Zimányi, Query evaluation in probabilistic relational databases, in: Selected Papers from the International Workshop on Uncertainty in Databases and Deductive Systems, Elsevier, 1997, pp. 179–219.