

Extension of Tuple Calculus to Multisets

Iryna Lysenko^{1*}, Olha Moroz²

¹Nizhyn Mykola Gogol State University, Graftska str., 2, Nizhyn, 16600, Ukraine

²International Research and Training Centre for Information Technologies and Systems of the NASU and MESU, Academician Hlushkov Ave., 40, Kyiv, 03680, Ukraine

Abstract

This paper is a logical continuation of research devoted to the actual problem of developing the theoretical foundations of table (relational) databases. The issue of using multisets in table databases is important and relevant. Many database-oriented languages require a relational model with multiset semantics. There are many applied problems, the feature of which is multiplicity and repeatability of data. For example, these are sociological polls of different population groups, calculations on DNA, and others. In this context, the question of constructing a tuple calculus for a multiset table algebra is considered, in which the concept of a table is refined using the concept of a multiset. In the article, the formalization of tuple calculus for multiset table algebra is carried out. The alphabet, and the syntax of terms, atoms, and formulas are defined. A set of legal formulas is introduced through the concept of the free and bound variable. The concept of a scheme and set of attributes with which a tuple variable occurs in a formula are also introduced. The definition of tuple calculus expression for multiset table algebra is given, according to which it is a multiset of tuples that satisfy the condition defined by the legal formula. The article provides rules for determining the number of tuple duplicates in the resulting multiset. Another important result consists in proving that constructed tuple calculus is as expressive as multiset table algebra. This research opens up new possibilities for database theory development and may be useful for information technology and database professionals. It contributes to a deeper understanding of construction query principles, an important aspect of modern computer science and industry.

Keywords

Relation Databases, multiset, multiset table algebra, tuple calculus

1. Introduction

Relational calculus underlies most relational query languages, specifying only the expected result, while relational algebra involves constructing a relational expression and performing operations. There are two main approaches to relational calculus:

- Tuple calculus (E. Codd) which operates on table tuples [1];
- Domain calculus (M. Lacroix, A. Pirotte) which focuses on table domains [2].

The clarification of the concept of relation in terms of naminal sets was carried out by V.N. Redko, Yu.Y. Bron, D.B. Buy, and S.A. Polyakov [3]. The monograph [4] introduces the consideration of table algebra for infinite (finite) tables, which significantly generalizes and extends classical Codd's relational algebra. The generalization is that a relation is understood as an arbitrary set of single-scheme tuples, in particular, an infinite one, while each table is assigned a certain scheme. Tuple (domain) calculi have been also constructed for these algebras. Tuple calculus and domain calculus are supplemented with functional and predicate signatures on the universal domain. Two main results are presented. The first result demonstrates the equivalence between the

14th International Scientific and Practical Conference from Programming UkrPROG'2024, May 14-15, 2024, Kyiv, Ukraine

* Corresponding author.

✉ iryna.glushko@ndu.edu.ua (I. Lysenko); mog_91@ukr.net (O. Moroz)

📄 0000-0003-2549-5356 (I. Lysenko); 0000-0002-0356-8780 (O. Moroz)



© 2024 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

table algebra for infinite tables and the corresponding generalized relational calculi. The second result focuses on the equivalence between the subalgebra for finite tables of this algebra and the corresponding restricted generalized relational calculi. For the second result, the consideration is limited to only "safe" expressions. The methodological basis of these formalisms is the compositional approach to programming.

Additionally, in the work [4] multiset table algebra is constructed, which extends database capabilities through multisets. The signature of multiset table algebra is supplemented with new operations: inner and outer join operations, semijoin operations, and aggregate operations.

In [5, 6] the theorem-plural and logical-algebraic methods found that table algebra of infinite tables does not form subalgebra of multiset table algebra since it is not closed with respect to the union, projection and active complement. Thus, multiset table algebra is not a wider formalism than table algebra of infinite tables

Many database-oriented languages, such as SQL, require a relational model with multiset semantics. However, the traditional tuple (domain) calculi do not support queries that produce tables with duplicates. Therefore, there is a need to develop a calculus for multiset table algebra.

The research aims to construct a tuple calculus for multiset table algebra and to show that it is no less expressive than multiset table algebra.

2. Basics of Multisets

Many languages oriented towards working with databases require a relational data model with multiset semantics because, firstly, relations that allow duplicates are useful in many applied fields where duplicate objects may exist; secondly, in the relational data model, removing duplicates after performing projection and union operations implies merging identical elements or performing other labor-intensive actions. Various researchers have focused on the utilization of multisets in databases, including Paul Grefen and Rolf A. de By [7, 8], G. Lamperti, M. Melchiori, and M. Zanella [9], G. Garcia-Molina, J. Ullman, and J. Widom [10, 11], A. Silbeschatz, H. Korth, and S. Sudarshan [12], D.B. Buy, S.A. Polyakov, Yu.Y. Brona, and V.N. Redko [3]. However, this issue still requires clarification and expansion.

Let's start by considering the key terms of multisets based on sources [3, 13].

Let U be an arbitrary set. A multiset α with basis U is a function

$$\alpha: U \rightarrow N,$$

where $N = \{1, 2, \dots\}$ is the set of natural numbers.

Let \mathbf{D} be the universe of elements of multiset bases. A characteristic function of multiset α is a function $\chi_\alpha: \mathbf{D} \rightarrow Z_+$, the values of which are specified by the following piecewise schema:

$$\chi_\alpha(d) = \begin{cases} \alpha(d) & \text{if } d \in \text{dom } \alpha \\ 0, & \text{else} \end{cases}$$

for all $d \in \mathbf{D}$. Here $\text{dom } \alpha$ is the range of definition of multiset as a function ($\text{dom } \alpha = U_\alpha$ – basis of multiset α).

The empty multiset \emptyset_m is defined as a multiset which basis is an empty set.

The 1-multisets are multisets whose range of values is the empty set or a single-element set $\{1\}$. These multisets are the analogues of ordinary sets.

The operations over multisets are defined in terms of characteristic functions in monograph [3]. Authors define operations of multiset \cup_1 , intersection \cap_1 , difference \setminus_1 , which build 1-multisets, and operations of multiset union \cup_{AU}^R , intersection \cap_{AU}^R , difference \setminus_{AU}^R , which build multisets of

general view. The Cartesian product of multiset \otimes , the operation Dist α , which build 1-multiset, and analog of a full image for multisets are defined too.

3. Multiset Table Algebra

Let's examine the main concepts and statements of multiset table algebra, based on the monograph [4]. In this case, under relation, understand a multiset, in particular, infinite.

Let's consider two sets: \mathbf{A} is the set of attributes, and \mathbf{D} is the universal domain.

A scheme will be defined as any arbitrary finite set of attributes $R \subseteq \mathbf{A}$. A tuple of the scheme R is a nominal set on pair R, \mathbf{D} . The projection of this nominal set for the first component is equal to R . We will use the following notations: $S(R)$ is the set of all tuples of the schema R , and S is the set of all tuples.

A table of the scheme R ($R \subseteq \mathbf{A}$) is a pair

$$\langle \psi, R \rangle,$$

where the first component ψ is an arbitrary multiset, the basis of which $\Theta(\psi)$ is the set of tuples of the same scheme and the second component R is a scheme of the table.

The set of all table on scheme R is designated as $\Psi(R)$ and the set of all table is designated as $\Psi = \bigcup_R \Psi(R)$.

The notation $Occ(s, \psi)$ represents the number of duplicate tuple s in the multiset ψ . A multiset ψ can also be written as $\{s_1^{n_1}, \dots, s_k^{n_k}\}$, where $n_i = Occ(s_i, \psi)$, $i = 1, \dots, k$, and $\Theta(\psi) = \{s_1, \dots, s_k\}$ is a of the multiset ψ .

Under multiset table algebra is understood an algebra

$$\langle \Psi, \Omega_{P, \Xi} \rangle,$$

where Ψ is the set of all tables, $\Omega_{P, \Xi} = \{\bigcup_{All}^R, \bigcap_{All}^R, \setminus_{All}^R, \sigma_{p,R}, \pi_{X,R}, \otimes_{R_1, R_2}, Rt_{\xi,R}, \sim_R\}_{X, R, R_1, R_2 \subseteq \mathbf{A}}^{p \in P, \xi \in \Xi}$ is the signature, P, Ξ are the sets of parameters. The signature $\Omega_{P, \Xi}$ contains operations of multiset (union \bigcup_{All}^R , intersection \bigcap_{All}^R , difference \setminus_{All}^R) and special operations (selection $\sigma_{p,R}$, projection $\pi_{X,R}$, join \otimes_{R_1, R_2} , renaming $Rt_{\xi,R}$, and active complement \sim_R).

Multiset table algebra is also filled up with additional operations such as inner join (Cartesian join, natural join, join using attributes and join on predicate), outer join (outer left join, outer right join, outer full join and union join), semijoin and aggregate operations (Sum, Avg, Max, Min, Count). The special element NULL is inserted in the universal domain for to define of outer join and outer set operations. The operations of signature $\Omega_{P, \Xi}$ and a formal mathematical semantics of additional operations are defined in [4].

The following statement holds.

Theorem 1. Any expression of the multiset table algebra can be replaced with an equivalent expression that uses only the operations of selection, join, projection, union, difference, and renaming.

Proof. To prove the first statement, we will show that the operations of intersection and active complement can be expressed through other operations of multiset table algebra. The operation of multiset intersection can be replaced by the difference [13]:

$$\langle \psi_1, R \rangle \cap_{All}^R \langle \psi_2, R \rangle = \langle \psi_1, R \rangle \setminus_{All}^R (\langle \psi_1, R \rangle \setminus_{All}^R \langle \psi_2, R \rangle).$$

The operation of active complement is expressed through the operations of projection, difference, and join (by definition):

$$\sim_R (\langle \psi, R \rangle) = C(\langle \psi, R \rangle) \setminus_{All}^R \langle \psi, R \rangle,$$

where $C(\langle\psi, R\rangle) = \pi_{\{A_1\}, R}(\langle\psi, R\rangle) \otimes_{\{A_1\}, \{A_2\}} \dots \otimes_{\{A_1, \dots, A_{n-1}\}, \{A_n\}} \pi_{\{A_n\}, R}(\langle\psi, R\rangle)$, and $R = \{A_1, \dots, A_n\}$ is a scheme of the table $\langle\psi, R\rangle$.

According to the proof, the operations of intersection and active complement are derived from other operations in the signature of multiset table algebra and can henceforth be excluded from consideration.

4. Tuple Calculus for Multiset Table Algebra

The basis of relational calculus is the calculus of first-order predicates. We will start the construction of tuple calculus for multiset table algebra with the definition of the alphabet.

The alphabet of tuple calculus for multiset table algebra consists of:

- a set of attributes \mathbf{A} and universal domain \mathbf{D} ;
- a set of object variables (tuple variables) x_1, x_2, \dots ;
- a set of object constants d_1, d_2, \dots ;
- a set of function symbols $f_1^{n_1}, f_2^{n_2}, \dots, n_i \geq 1$;
- a set of predicate symbols $p_1^{m_1}, p_2^{m_2}, \dots, m_i \geq 1$;
- a set of symbols of constant tables $\langle\alpha, R\rangle$;
- a set of symbols of variable tables $\langle X, R\rangle$;
- the signs of logical operations \neg, \wedge, \vee , and quantifiers \forall, \exists ;
- punctuation marks – parentheses $()$ and commas.

The universal domain \mathbf{D} is the domain of interpretation of object constants, and the set of all tuples over \mathbf{D} is the domain of interpretation of object variables.

We will use

- \mathbf{x} as syntactic variable, the domain of change of which is the set of variables;
- \mathbf{f} as syntactic variable, the domain of change of which is the set of function symbols;
- \mathbf{p} as syntactic variable, the domain of change of which is the set of predicate symbols;
- \mathbf{d} as syntactic variable;
- \mathcal{A} as syntactic variable, the domain of change of which is the set of attributes.

Terms and formulas are distinguished among the words written using alphabet symbols. Let us define these syntactic objects by induction on their length.

The following expressions are terms:

- a) \mathbf{d} is a term;
- b) $\mathbf{x}(\mathcal{A})$ is a term;
- c) if $\mathbf{u}_1, \dots, \mathbf{u}_n$ are terms and \mathbf{f} is a function symbol of arity n then $\mathbf{f}(\mathbf{u}_1, \dots, \mathbf{u}_n)$ is a term;
- d) there are no terms other than those specified in points a)-c).

Further in the text, \mathbf{u} is a syntactic variable, the domain of change of which is the set of terms.

Let's define the formulas, starting with atomic formulas (atoms), which come in three types:

- a1. For any constant table $\langle\alpha, R\rangle$ and for any tuple variable \mathbf{x} , $\alpha_R(\mathbf{x})$ is an atom. $\alpha_R(\mathbf{x})$ stands for $\mathbf{x} \in \langle\alpha, R\rangle$.
- a2. For any variable table $\langle X, R\rangle$ and for any tuple variable \mathbf{x} , $X_R(\mathbf{x})$ is an atom. $X_R(\mathbf{x})$ stands for $\mathbf{x} \in \langle X, R\rangle$.
- a3. For any terms $\mathbf{u}_1, \dots, \mathbf{u}_m$, and for any predicate \mathbf{p} of arity m on the universal domain \mathbf{D} , $\mathbf{p}(\mathbf{u}_1, \dots, \mathbf{u}_m)$ is an atom.
 - f1. Let's construct formulas from atoms using the logical connectives \neg, \wedge, \vee , quantifiers \forall, \exists , and parentheses.
 - f2. Every atom is a formula.
 - f3. If \mathbf{P} and \mathbf{Q} are formulas, then $\neg\mathbf{P}, \mathbf{P}\wedge\mathbf{Q}, \mathbf{P}\vee\mathbf{Q}$ are formulas.

- f4. If \mathbf{x} is a tuple variable, \mathbf{P} is a formula, $R \subseteq \mathbf{A}$ is a scheme, then $\forall \mathbf{x}(R)\mathbf{P}, \exists \mathbf{x}(R)\mathbf{P}$ are formulas.
- f5. If \mathbf{P} is a formula, then (\mathbf{P}) is a formula.
- f6. There are no other formulas besides those specified in points f1-f4.

We will use \mathbf{P}, \mathbf{Q} and \mathbf{G} as syntactic variables, the domain of change of which is the set of formulas.

Let's define the class of legal formulas, using the concepts of free and bound variables, the scheme $scheme(\mathbf{x}, \mathbf{P})$ for tuple variable \mathbf{x} , and the set of attributes with which tuple variable \mathbf{x} occurs in a formula \mathbf{P} , $attr(\mathbf{x}, \mathbf{P})$. The expressions $scheme(\mathbf{x}, \mathbf{P})$ and $attr(\mathbf{x}, \mathbf{P})$ are defined if the tuple \mathbf{x} has at least one free occurrence in the formula \mathbf{P} . Additionally, the including $attr(\mathbf{x}, \mathbf{P}) \subseteq scheme(\mathbf{x}, \mathbf{P})$ holds if these expressions are defined.

We will define expression $attr$ for terms first:

- 1. if $\mathbf{u} = \mathbf{d}$, then $attr(\mathbf{x}, \mathbf{u}) = \emptyset$;
- 2. if $\mathbf{u} = \mathbf{x}(\mathcal{A})$, then $attr(\mathbf{x}, \mathbf{u}) = \{\mathcal{A}\}$, and $attr(\mathbf{x}, \mathbf{y}(\mathcal{A})) = \emptyset$, where $\mathbf{x} \neq \mathbf{y}$;
- 3. if $\mathbf{u} = \mathbf{f}(\mathbf{u}_1, \dots, \mathbf{u}_n)$, where \mathbf{u}_i are terms then $attr(\mathbf{x}, \mathbf{u}) = \bigcup_{i=1}^n attr(\mathbf{x}, \mathbf{u}_i)$.

In other words, $attr(\mathbf{x}, \mathbf{u})$ is the set of attributes that the scheme of the tuple \mathbf{x} must have.

Consider the cases where \mathbf{P} is an atomic formula, then

- a1. if $\mathbf{P} = \alpha_R(\mathbf{x})$, then \mathbf{x} is free in \mathbf{P} and $scheme(\mathbf{x}, \mathbf{P}) = attr(\mathbf{x}, \mathbf{P}) = R$;
- a2. similarly if $\mathbf{P} = X_R(\mathbf{x})$, then \mathbf{x} is free in \mathbf{P} and $scheme(\mathbf{x}, \mathbf{P}) = attr(\mathbf{x}, \mathbf{P}) = R$;
- a3. if $\mathbf{P} = \mathbf{p}(\mathbf{u}_1, \dots, \mathbf{u}_m)$, where \mathbf{u}_j are terms, and $\mathbf{x}_1, \dots, \mathbf{x}_k$ are all variables of these terms, then this tuple variables are free in formula \mathbf{P} , $scheme(\mathbf{x}_i, \mathbf{P})$ is undefined, and $attr(\mathbf{x}_i, \mathbf{P}) = \bigcup_{j=1}^m attr(\mathbf{x}_i, \mathbf{u}_j)$, $i = 1, \dots, k$.

Atomic formulas are all legal. The construction of all legal formulas proceeds by induction on the length of formulas. Assume \mathbf{Q} and \mathbf{G} are both legal formulas.

- f1. If $\mathbf{P} = \neg \mathbf{G}$, then \mathbf{P} is legal, and all occurrences of variables in \mathbf{P} free or bound as they are in \mathbf{G} . For every tuple \mathbf{x} that occurs free in \mathbf{G} , $scheme(\mathbf{x}, \mathbf{P}) \simeq scheme(\mathbf{x}, \mathbf{G})$ and $attr(\mathbf{x}, \mathbf{P}) = attr(\mathbf{x}, \mathbf{G})$, where \simeq is a generalized equality, meaning that both sides of the equality are either undefined or defined and have equal values [14].
- f2. If $\mathbf{P} = \mathbf{G} \wedge$ or $\mathbf{P} = \mathbf{G} \vee \mathbf{Q}$, then all occurrences of variables in \mathbf{P} are free or bound as their corresponding occurrences are in \mathbf{G} and \mathbf{Q} . Assume variable \mathbf{x} occurs free in subformulas \mathbf{G} and/or \mathbf{Q} . Define the scheme, and the set of attributes with which tuple variable \mathbf{x} occurs in a formula for formula \mathbf{P} . Next cases take place.
 - a. The schemes of formulas $scheme(\mathbf{x}, \mathbf{G})$ and $scheme(\mathbf{x}, \mathbf{Q})$ are both defined. Formula \mathbf{P} is legal if equality $scheme(\mathbf{x}, \mathbf{G}) = scheme(\mathbf{x}, \mathbf{Q})$ holds. According to the definition $scheme(\mathbf{x}, \mathbf{P}) = scheme(\mathbf{x}, \mathbf{G})$.
 - b. The scheme is defined for only one subformula. Assume $scheme(\mathbf{x}, \mathbf{G})$ is defined, and $scheme(\mathbf{x}, \mathbf{Q})$ is undefined. Formula \mathbf{P} is legal if including $attr(\mathbf{x}, \mathbf{Q}) \subseteq scheme(\mathbf{x}, \mathbf{G})$ holds. According to the definition $scheme(\mathbf{x}, \mathbf{P}) = scheme(\mathbf{x}, \mathbf{G})$.
 - c. The scheme is undefined for both subformulas, then formula \mathbf{P} is legal but $scheme(\mathbf{x}, \mathbf{P})$ is undefined.

In all these three cases $attr(\mathbf{x}, \mathbf{P}) = attr(\mathbf{x}, \mathbf{G}) \cup attr(\mathbf{x}, \mathbf{Q})$.

- f3. If $\mathbf{P} = \exists \mathbf{x}(R)\mathbf{G}$ then \mathbf{x} must occur free in \mathbf{G} for \mathbf{P} to be legal. Furthermore, if $scheme(\mathbf{x}, \mathbf{G})$ is defined, then equality $scheme(\mathbf{x}, \mathbf{G}) = R$ must hold if including $attr(\mathbf{x}, \mathbf{G}) \subseteq R$ is held. $scheme(\mathbf{x}, \mathbf{P})$ and $attr(\mathbf{x}, \mathbf{P})$ are not defined, since \mathbf{x} does not occur free in \mathbf{P} . Any occurrence of a variable $\mathbf{y} \neq \mathbf{x}$ is free or bound in \mathbf{P} as it was in \mathbf{G} . If \mathbf{y} occur free in \mathbf{P} , then $scheme(\mathbf{y}, \mathbf{P}) \simeq scheme(\mathbf{y}, \mathbf{G})$ and $attr(\mathbf{y}, \mathbf{P}) = attr(\mathbf{y}, \mathbf{G})$.
- f4. If $\mathbf{P} = \forall \mathbf{x}(R)\mathbf{G}$, then all restrictions and definitions are the same as in f3.
- f5. If $\mathbf{P} = (\mathbf{G})$, then \mathbf{P} is legal, and free and bound variables, $scheme$ and $attr$ are the same as for \mathbf{G} .

In other words, the equality $attr(\mathbf{x}, \mathbf{P}) = R$ means that for a specific interpretation of the formula \mathbf{P} , when the variable \mathbf{x} takes on a value in the form of a tuple s of the scheme R' , the inclusion $R \subseteq R'$ must hold.

Expressions of tuple calculus for multiset table algebra have the form

$$\{\mathbf{x}^n(R)|\mathbf{P}(\mathbf{x})\},$$

where:

1. The formula \mathbf{P} is legal.
2. The variable \mathbf{x} is the only variable that occurs free in the formula \mathbf{P} .
3. If $scheme(\mathbf{x}, \mathbf{P})$ is defined, then $scheme(\mathbf{x}, \mathbf{P}) = R$, otherwise $attr(\mathbf{x}, \mathbf{P}) \subseteq R$.
4. n is the number of duplicates of the tuple \mathbf{x} .

It is worth emphasizing that the result of executing the query defined by the expression $\{\mathbf{x}^n(R)|\mathbf{P}(\mathbf{x})\}$ is a multiset of tuple described by the expression $\mathbf{P}(\mathbf{x})$.

Let $\mathbf{P}(\mathbf{x})$ be a legal formula, $R \subseteq \mathbf{A}$. Then s substituted for \mathbf{x} in formula \mathbf{P} is the formula $\mathbf{P}(s/\mathbf{x})$. The value of the formula $\mathbf{P}(s/\mathbf{x})$ is determined by modifying each atom from \mathbf{P} according to the following rules:

- a1. Let the tuple variable \mathbf{x} in $\alpha_{R'}(\mathbf{x})$ be free in formula \mathbf{P} . By the definition of an legal formula, we have the inclusion $R' \subseteq R$. Atom $\alpha_{R'}(\mathbf{x})$ acquires the value of true at s substituted for \mathbf{x} , if $s|R' \in \theta(\alpha)$, otherwise, atom $\alpha_{R'}(\mathbf{x})$ acquires the value of false.
- a2. Let the tuple variable \mathbf{x} in $X_{R'}(\mathbf{x})$ be free in formula \mathbf{P} . By the definition of an legal formula, we have the inclusion $R' \subseteq R$. Atom $X_{R'}(\mathbf{x})$ acquires the value of true at s substituted for \mathbf{x} , if $s|R' \in \theta(X)$, otherwise, atom $X_{R'}(\mathbf{x})$ acquires the value of false.
- a3. Let the tuple variable \mathbf{x} in $\mathbf{P} = \mathbf{p}(\mathbf{u}_1, \dots, \mathbf{u}_m)$ be free in formula \mathbf{P} , then at s substituted for \mathbf{x} , replace $\mathbf{x}(A_i)$ by $d_i \in \mathbf{D}$, where $\langle A_i, d_i \rangle \in s$ (d_i is value of attribute A_i in tuple s). Atom $\mathbf{p}(\mathbf{u}_1, \dots, \mathbf{u}_m)$ acquires the value of true, if predicate \mathbf{p} is true on the proper values, otherwise atom acquires the value of false.

The set of values of truth of all atoms of formula is called interpretation. Let \mathbf{P} be a legal formula with no free tuple variables. The interpretation of formula \mathbf{P} is defined as follows.

- f1. If $\mathbf{P} = \neg \mathbf{G}$, then \mathbf{G} must have no free variables. The formula \mathbf{P} is true, if \mathbf{G} is false, and it is false, if formula \mathbf{G} is true.
- f2. If $\mathbf{P} = \mathbf{G} \wedge \mathbf{Q}$ or $\mathbf{P} = \mathbf{G} \vee \mathbf{Q}$, then neither \mathbf{G} or \mathbf{Q} have free variables. If $\mathbf{P} = \mathbf{G} \wedge \mathbf{Q}$, then \mathbf{P} is true exactly when \mathbf{G} and \mathbf{Q} are both true, otherwise, \mathbf{P} is false. If $\mathbf{P} = \mathbf{G} \vee \mathbf{Q}$, then \mathbf{P} is false exactly when \mathbf{G} and \mathbf{Q} are both false, otherwise, \mathbf{P} is true.
- f3. If $\mathbf{P} = \exists \mathbf{x}(R)\mathbf{G}$, then \mathbf{x} is the only variable that occurs free in \mathbf{G} . The formula \mathbf{P} is true, if there is at least one tuple $s \in S(R)$ such that formula $\mathbf{G}(s/\mathbf{x})$ is true, otherwise, formula \mathbf{P} is false.
- f4. If $\mathbf{P} = \forall \mathbf{x}(R)\mathbf{G}$, then \mathbf{x} is the only variable that occurs free in \mathbf{G} . **The formula \mathbf{P} is true**, if for every tuple $s \in S(R)$ formula $\mathbf{G}(s/\mathbf{x})$ is true, otherwise, formula \mathbf{P} is false.

If $\mathbf{P} = (\mathbf{G})$, then \mathbf{P} is true, if formula \mathbf{G} is true and \mathbf{P} is false, if \mathbf{G} is false.

Let $E = \{\mathbf{x}^n(R)|\mathbf{P}(\mathbf{x})\}$ be a tuple calculus expression. The value of expression E is the table $\langle \varphi, R \rangle$ of multiset table algebra containing every tuple $s \in S(R)$ such that formula $\mathbf{P}(s/\mathbf{x})$ is true. The number of duplicates of the tuple s in the table $\langle \varphi, R \rangle$ is determined as follows:

- 1) if $\mathbf{P} = \alpha_R(\mathbf{x})$, then $n = Occ(s, \varphi) = Occ(s, \alpha)$;
- 2) if $\mathbf{P} = X_R(\mathbf{x})$, then $n = Occ(s, \varphi) = Occ(s, X)$;
- 3) if $\mathbf{P} = \mathbf{p}(\mathbf{u}_1, \dots, \mathbf{u}_m)$, where \mathbf{u}_j are terms, $j = \overline{1, m}$, and $\mathbf{x}_1, \dots, \mathbf{x}_k$ are all variables of these terms, then

$$n = Occ(s, \varphi) = \sum_{\substack{s' \in \theta(\psi), s'|R'=s \\ \mathbf{p}(\mathbf{u}_1, \dots, \mathbf{u}_m)=true}} Occ(s', \psi),$$

where $\langle \psi, R \rangle$ is the table to which the query is constructed, $R' = \bigcup_{j=1}^m attr(\mathbf{x}_i, \mathbf{u}_j), i = \overline{1, k}$;

- 4) if $\mathbf{P} = \neg \mathbf{G}$ and the formula \mathbf{G} generates m duplicates of the tuple s , then
- 5)

$$n = Occ(s, \varphi) = Occ(s, C(\psi)) \dot{-} m,$$

where $C(\psi)$ is a multiset of the table $C(\langle\psi, R\rangle)$ which is the saturation of the table $\langle\psi, R\rangle$, which is the value of the expression $\{\mathbf{y}^n(R)|\mathbf{P}(\mathbf{y})\}$, and $x \dot{-} y = \begin{cases} x - y & \text{if } x \geq y \\ 0 & \text{if } x < y \end{cases}$;

- 6) if $\mathbf{P} = \mathbf{G} \wedge \mathbf{Q}$, the formula \mathbf{G} generates k duplicates of the tuple s , and the formula \mathbf{Q} generates m duplicates of the tuple s , then $n = Occ(s, \varphi) = \min(k, m)$;
- 7) if $\mathbf{P} = \mathbf{G} \vee \mathbf{Q}$, the formula \mathbf{G} generates k duplicates of the tuple s , and the formula \mathbf{Q} generates m duplicates of the tuple s , then $n = Occ(s, \varphi) = k + m$;
- 8) if $\mathbf{P} = \exists \mathbf{x}(R)\mathbf{G}$ and the formula \mathbf{G} generates k duplicates of the tuple s , then $n = Occ(s, \varphi) = k$.
- 9) if $\mathbf{P} = \forall \mathbf{x}(R)\mathbf{G}$ and the formula \mathbf{G} generates k duplicates of the tuple s , then $n = Occ(s, \varphi) = k$.
- 10) if $\mathbf{P} = (\mathbf{G})$ and the formula \mathbf{G} generates k duplicates of the tuple s , then $n = Occ(s, \varphi) = k$.

Theorem 2. If F is a multiset table algebra expression, then it is possible effectively to build equivalent to it expression E of tuple calculation.

Proof. According to Theorem 1, in the proof, it is sufficient to consider expressions of multiset table algebra that contain only the operations of union, difference, selection, projection, join, and renaming. We will prove this theorem by mathematical induction on the number of operations in the expression F .

Basis. In this case, the expression F does not contain any operations. There are two cases.

First case. If $F = \langle\alpha, R\rangle$ is constant table, where α is a multiset of the scheme R , then $E = \{\mathbf{x}^n(R)|\alpha_R(\mathbf{x})\}$.

Second case. If $F = \langle X, R\rangle$ is variable table, then $E = \{\mathbf{x}^n(R)|X_R(\mathbf{x})\}$.

Induction. Assume the theorem holds for any multiset table algebra expression with fewer than i operators. Let F have i operators. There are six cases.

Case 1 (union). $F = F_1 \cup_{All}^R F_2$,

where expressions F_1 and F_2 each have less than i operators. By the inductive hypothesis we can find tuple calculus expressions $\{\mathbf{x}^n(R)|\mathbf{P}(\mathbf{x})\}$ and $\{\mathbf{x}^m(R)|\mathbf{Q}(\mathbf{x})\}$ equivalent to F_1 and F_2 respectively. The values of these expressions are tables in which the tuple \mathbf{x} appears n and m times, respectively. Then E is $\{\mathbf{x}^k(R)|\mathbf{P}(\mathbf{x}) \vee \mathbf{Q}(\mathbf{x})\}$, where $k = n + m$.

Case 2 (difference). $F = F_1 \setminus_{All}^R F_2$,

where expressions F_1 and F_2 each have less than i operators. Tuple calculus expressions $\{\mathbf{x}^n(R)|\mathbf{P}(\mathbf{x})\}$ and $\{\mathbf{x}^m(R)|\mathbf{Q}(\mathbf{x})\}$ are equivalent to F_1 and F_2 respectively as in Case 1. Then $E = \{\mathbf{x}^k(R)|\mathbf{P}(\mathbf{x}) \wedge \neg \mathbf{Q}(\mathbf{x})\}$, where $k = n - m$.

Case 3 (selection). $F = \sigma_{\tilde{p}, R}(F_1)$,

where expression F_1 has less than i operators. Let $\{\mathbf{x}^n(R)|\mathbf{P}(\mathbf{x})\}$ be tuple calculus expression equivalent to F_1 . Then $E = \{\mathbf{x}^k(R)|\mathbf{P}(\mathbf{x}) \wedge \mathbf{p}(\mathbf{x}(A_1), \dots, \mathbf{x}(A_m))\}$, where $R = \{A_1, \dots, A_m\}$ is the scheme of table that is the value of expression F_1 . The number of duplicates of the tuple \mathbf{x} in the output table does not change, therefore $k = n$. Assume that predicate-parameter of select is defined as $\tilde{p}(s) = true \Leftrightarrow \mathbf{p}(s(A_1), \dots, s(A_m)) = true$, $s \in S(R)$, where \mathbf{p} is a predicate symbol of arity m . \square

Case 4 (projection). $F = \pi_{X, R}(F_1)$.

Let $\{\mathbf{x}^n(R)|\mathbf{P}(\mathbf{x})\}$ be tuple calculus expression equivalent to F_1 . Then $E = \{\mathbf{y}^k(X \cap R)|\exists \mathbf{x}(R)(\mathbf{P}(\mathbf{x}) \wedge \bigwedge_{A \in X \cap R} \mathbf{y}(A) = \mathbf{x}(A))\}$, where $k = \sum_{\bigwedge_{A \in X \cap R} \mathbf{y}(A) = \mathbf{x}(A)} n$.

Case 5 (join). $F = F_1 \otimes_{R_1, R_2} F_2$.

Tuple calculus expressions $\{\mathbf{x}^n(R)|\mathbf{P}(\mathbf{x})\}$ and $\{\mathbf{x}^m(R)|\mathbf{Q}(\mathbf{x})\}$ are equivalent to F_1 and F_2 respectively. Then E is

$\{\mathbf{z}^k(R_1 \cup R_2)|\exists \mathbf{x}(R_1)\exists \mathbf{y}(R_2)(\mathbf{P}(\mathbf{x}) \wedge \mathbf{Q}(\mathbf{y}) \wedge \bigwedge_{A \in R_1} \mathbf{z}(A) = \mathbf{x}(A) \wedge \bigwedge_{A \in R_2} \mathbf{z}(A) = \mathbf{y}(A))\}$, where $k = n \times m$. \square

Case 6 (renaming). $F = Rt_{\xi, R}(F_1)$,

where $\xi: A \rightarrow A$ is injective function that renames attributes. Let $\{\mathbf{x}^n(R)|P(\mathbf{x})\}$ be tuple calculus expression equivalent to F_1 . Then

$$E = \{\mathbf{y}^k(R_2)|\exists \mathbf{x}(R)(P(\mathbf{x}) \wedge \bigwedge_{C \in R \setminus \text{dom} \xi} \mathbf{y}(C) = \mathbf{x}(C) \wedge \bigwedge_{A \in R \cap \text{dom} \xi} \mathbf{x}(A) = \mathbf{y}(\xi(A)))\}, \text{ where } R_2 = R \setminus \text{dom} \xi \cup \xi[R], \text{ and } k = n.$$

5. Relation to the SQL query language

The fundamental data object in the SQL language is not the classical relation of E. Codd but rather a table. Moreover, SQL tables do not contain sets of tuples but multisets of tuples, meaning duplicates are allowed. The basic SQL operators are not relational operators in the strict sense but are analogs of relational operators designed to work with multisets. When creating a new table using a query, an SQL system typically does not remove duplicate tuples but returns a result in which the same tuple can appear multiple times. To exclude duplicates, the keyword DISTINCT must be placed after the SELECT operator.

Let's demonstrate the appropriateness of constructing tuple calculus for multiset table algebra with the following example.

Example 1. Consider the table $\langle \text{Scores}, R \rangle$, where the scheme $R = \{\mathbb{N}^0, \text{Name}, \text{Topic 1}, \text{Topic 2}, \text{Topic 3}, \text{Quiz}\}$ (see Table 1).

Table 1

Table $\langle \text{Scores}, R \rangle$

\mathbb{N}^0	Name	Topic 1	Topic 2	Topic 3	Quiz
1.	Student 1	5	15	14	16
2.	Student 2	6	14	15	16
3.	Student 3	7	17	20	20
4.	Student 4	9	20	19	20
5.	Student 5	5	18	15	18
6.	Student 6	6	19	13	20
7.	Student 7	4	8	9	16

The answer to the question "What scores did the first five students get for the quiz?" in SQL would look like this:

```
SELECT Quiz
FROM Scores
LIMIT 5;
```

The result is a table $\langle \text{Quiz}, R_1 \rangle$, where the scheme $R_1 = \{\text{Quiz}\}$, which contains duplicate values (see Table 2).

Table 2

Table $\langle \text{Quiz}, R_1 \rangle$

Quiz
16
16
20
20
18

It is impossible to implement this query in terms of classical tuple calculus, since the result will be a set of tuples, not a multiset (as expected), meaning duplicates will not be considered in the result.

In the tuple calculus for multiset table algebra, the expression equivalent to this query will have the form:

$$\{x^n(\text{Quiz}) \mid \exists y(R)(\text{Scores}(y) \wedge (\text{No} = 1 \vee (\text{No} = 2 \vee (\text{No} = 3 \vee (\text{No} = 4 \vee (\text{No} = 5) \wedge \wedge x(\text{Quiz}) = y(\text{Quiz}))))))\}.$$

The result described by this tuple calculus expression is similar to the result obtained when executing the corresponding query in SQL.

Example 2. Consider the table $\langle \text{Results}, R' \rangle$, where the scheme $R' = \{ \text{No}, \text{Name}, \text{Total}, \text{ECTS} \}$ (see Table 3).

Let's consider the query "How many points for the quiz have students who scored more than 83 points in total?".

We will write the corresponding expression of multiset table algebra, assuming that we only need to know the quiz scores:

$$F = \pi_{\text{Quiz}} \left(\langle \text{Scores}, R \rangle \otimes_{R, R'} \sigma_{\text{Total} > 83} \langle \text{Results}, R' \rangle \right).$$

The result is a table $\langle \text{Query}, \{ \text{Quiz} \} \rangle$, which contains duplicate values (see Table 4).

Table 3

Table $\langle \text{Results}, R' \rangle$

No	Name	Total	ECTS
1.	Student 1	77	C
2.	Student 2	80	C
3.	Student 3	93	A
4.	Student 4	90	A
5.	Student 5	81	C
6.	Student 6	79	C
7.	Student 7	51	Fx

From Table 4, it is clear that only two students in total have more than 83 points. Let's write the tuple calculus expression equivalent to this algebraic expression.

Table 4

Table $\langle \text{Query}, \{ \text{Quiz} \} \rangle$

Quiz
20
20

Tuple calculus expressions $\{x^n(R) \mid \text{Scores}(x)\}$ and $\{y^m(R') \mid \text{Results}(y)\}$ are equivalent to $\langle \text{Scores}, R \rangle$ and $\langle \text{Results}, R' \rangle$ respectively, where $R = \{ \text{No}, \text{Name}, \text{Topic 1}, \text{Topic 2}, \text{Topic 3}, \text{Quiz} \}$ and $R' = \{ \text{No}, \text{Name}, \text{Total}, \text{ECTS} \}$.

We have tuple calculus expression

$$\{y^m(R') \mid \text{Results}(y) \wedge y(\text{Total}) > 83\}$$

for algebraic expression

$$\sigma_{\text{Total} > 83} \langle \text{Results}, R' \rangle.$$

Its value is a table without duplicates because table $\langle \text{Results}, R' \rangle$ does not have duplicates. Tuple calculus expressions which equivalent to algebraic expression

$$\langle \text{Scores}, R \rangle \otimes_{R, R'}^{\sigma_{\text{Total} > 83}} \langle \text{Results}, R' \rangle$$

is

$$\{\mathbf{z}^{k=n \times m}(R \cup R') | \exists \mathbf{x}(R) \exists \mathbf{y}(R') (\text{Scores}(\mathbf{x}) \wedge \text{Results}(\mathbf{y}) \wedge \mathbf{y}(\text{Total}) > 83 \wedge \mathbf{z}(\text{No}) = \mathbf{x}(\text{No}) \wedge \mathbf{z}(\text{Name}) = \mathbf{x}(\text{Name}) \wedge \mathbf{z}(\text{No}) = \mathbf{y}(\text{No}) \wedge \mathbf{z}(\text{Name}) = \mathbf{y}(\text{Name}))\}.$$

Since there are no duplicate tuples in the tables $\langle \text{Scores}, R \rangle$ and $\langle \text{Results}, R' \rangle$, there will also be no duplicate tuples in the resulting table after the join.

Finally, the tuple calculus expression equivalent to the algebraic expression F has the form:

$$\{\mathbf{w}^q(\{\text{Quiz}\} \cap (R \cup R')) | \exists \mathbf{z}(R \cup R') (\exists \mathbf{x}(R) \exists \mathbf{y}(R') (\text{Scores}(\mathbf{x}) \wedge \text{Results}(\mathbf{y}) \wedge \mathbf{y}(\text{Total}) > 83 \wedge \mathbf{z}(\text{No}) = \mathbf{x}(\text{No}) \wedge \mathbf{z}(\text{Name}) = \mathbf{x}(\text{Name}) \wedge \mathbf{z}(\text{No}) = \mathbf{y}(\text{No}) \wedge \mathbf{z}(\text{Name}) = \mathbf{y}(\text{Name})))\}.$$

In the table obtained after the join, the tuple \mathbf{z} appears once. Therefore, in the output table, the number of duplicates of the tuple \mathbf{w} is

$$\mathbf{q} = \sum_{A \in \{\text{Quiz}\} \cap (R \cup R')} \mathbf{w}(A) = \mathbf{z}(A) \mathbf{k}.$$

The result is a table $\langle \text{Query}, \{\text{Quiz}\} \rangle$, where $\text{Occ}(s, \text{Query}) = 1 + 1 = 2$, $s = \{\{\text{Quiz}\}, 20\}$. The corresponding SQL query has the form:

```
SELECT Quiz
FROM scores
INNER JOIN result ON scores.`No`=result.`No` AND scores.Name=result.Name
WHERE result.Total>83
```

Thus, as can be seen from the examples, the constructed tuple calculus for multiset table algebra allows for the adequate formalization of query languages, particularly SQL, considering the multiset semantics embedded in them.

Conclusions

The article proposes a tuple calculus for multiset table algebra. The alphabet, syntax of terms, atoms, and formulas of the tuple calculus are defined. Using the concept of free and bound tuple, the notion of scheme, and the set of attributes with which a tuple appears in formulas, a class of legal formulas is introduced. It is shown that the proposed tuple calculus is no less expressive than multiset table algebra. An example demonstrates the feasibility of constructing tuple calculus for multiset table algebra, considering that query languages oriented towards database work imply the repeatability of elements in a table.

The next research challenge is to establish the corresponding dual result.

References

- [1] E.F. Codd, Relational Completeness of Data Base Sublanguages, in: Data Base Systems (1972) 65-98.
- [2] M. Lacroix, A. Pirotte, Domain-oriented Relational Languages, in: Proceedings of 3rd Int. Conf. on Very Large Data Bases., 1977, pp. 370-378.
- [3] V.N. Redko, et al., Relational Databases: Table Algebras and SQL-like Language. Kyiv: Publishing house Academperiodica, 2001. [in Ukrainian]

- [4] D.B Buy, I.M. Glushko, Calculi and extensions of table algebras signature. Nizhyn: NDU im. M. Gogol, 2016. [in Ukrainian]
- [5] I.Glushko, About relationship between table algebra of infinite tables and multiset table algebra, in CEUR Workshop Proceedings, 2139, 2018, pp. 159–163.
- [6] I. Lysenko, Extended table algebra, extended multiset table algebra, and their relationship, in: Proceedings of International Conference on Software Engineering “SoftEngine 2020”, 2020, pp 28-32.
- [7] Paul W.P.J. Grefen, Rolf A. de By, A Multi-Set Extended Relational Algebra. A Formal Approach to a Practical Issue, in: Proceedings of 10th International Conference on Data Engineering, ICDE, 1994, pp. 80-88.
- [8] Paul W.P.J. Grefen, J. Flokstra, Extending a Multi-Set Relational Algebra to a Parallel Environment, in: Distributed and Parallel Databases, Vol 4. (1996) 81-99.
- [9] G. Lamperti, M. Melchiori, M. Zanella, On Multisets in Database Systems, in Multiset Processing: Mathematical, Computer Science, and Molecular Computing Points of View, number 2235 in Lecture Notes in Computing Since (2001)147-215.
- [10] H. Garcia-Molina, J.D. Ullman, J. Widom, Database Systems: The Complete Book, 2nd. ed., Prentice Hall, 2008.
- [11] J.D. Ullman, J. Widom. Ullman, A First Course in Database Systems, 3rd ed., Prentice Hall, 2007.
- [12] A. Silbeschatz, H. Korth, S. Sudarshan, Database System Concepts, 6th ed. McGraw-Hill, 2011.
- [13] J.A. Bogatyreva, Multisets theory and its applications. Ph.D. thesis, Kyiv National Taras Shevchenko University, 2011. [in Ukrainian]
- [14] N.Cutland, Computability. An introduction to recursive function theory, Cambridge University Press, 1980.