

# Derivatives of Quantitative Regular Expressions

Rajeev Alur<sup>1</sup>, Konstantinos Mamouras<sup>1(✉)</sup>, and Dogan Ulus<sup>2</sup>

<sup>1</sup> University of Pennsylvania, Philadelphia, PA, USA

mamouras@seas.upenn.edu

<sup>2</sup> Verimag, Université Grenoble-Alpes, Grenoble, France

**Abstract.** Quantitative regular expressions (QREs) have been recently proposed as a high-level declarative language for specifying complex numerical queries over data streams in a modular way. QREs have appealing theoretical properties, and each QRE can be compiled into an efficient streaming algorithm for its evaluation. In this paper, we generalize the notion of Brzozowski derivatives for classical regular expressions to QREs. Such derivatives immediately lead to an algorithm for incremental evaluation of QREs. While this algorithm does not have better time or space complexity than the previously known evaluation technique, it has the benefit of being simpler to explain and easier to prove correct.

## 1 Introduction

There are numerous applications that require the real-time processing of data generated at high rates such as: analyzing stock market data, monitoring production and manufacturing using sensors, network traffic monitoring, and click-stream analysis on the web. A core computational problem that is relevant to all such applications is the incremental aggregation of a stream of data items into numerical values that are useful for real-time decision making. Due to the enormous volume of data, these applications have hard constraints regarding space usage and the time required to process each newly arriving element.

There is a large body of prior research on stream processing which focuses on algorithmic techniques, often involving approximation and randomization, for computing efficiently specific numerical quantities such as the median [18], the number of distinct elements [17], the frequency moments [2], and aggregates over sliding windows [15]. There have also been several proposals for languages and systems that integrate stream processing with the data processing capabilities of traditional relational database systems [1, 8, 9].

The formalism of Quantitative Regular Expressions (QREs) was recently introduced in [5] with the orthogonal goal of providing convenient high-level programming abstractions for specifying complex queries over data streams in a modular way. QREs extend regular expressions, a well-established formalism for imparting hierarchical structure to sequences of symbols, with numerical operations such as sum, difference, min, max, average, and median. A QRE thus

describes both the regular parsing of the stream prefix seen so far and the hierarchical calculation of a quantitative aggregate that reflects the structure of the parse tree. This combination gives rise to a powerful declarative language, which can express conveniently many useful queries and is amenable to space- and time-efficient evaluation in the streaming model of computation. An implementation of QREs extended with extra features for processing realistic workloads is reported in [21]. The expressiveness of QREs coincides with the class of regular functions, which can be characterized with the model of Cost Register Automata (CRAs) [4] or, equivalently, with the MSO-definable graph transformations [14].

The main computational problem for QREs is their evaluation in the streaming model of computation. An efficient algorithm has already been described in [5] and implemented in [21], but it is not based on automata-theoretic techniques, and the question remains of whether there exists a simple model of automata for the streaming evaluation of QREs. In the simpler setting of classical regular expressions, the translation into a Nondeterministic Finite Automaton (NFA) gives rise to a very efficient streaming evaluation algorithm: the state of the algorithm consists of the active states of the NFA, and upon the arrival of a new symbol the state is updated by performing all possible transitions. Another approach for the evaluation problem is based on a technique proposed by Brzozowski in 1964 [11], where he introduced the notion of *derivation* for regular expressions extended with arbitrary Boolean operations. The *derivative* of an expression  $e$  with respect to a symbol  $a$ , typically denoted as  $D_a(e)$ , is an expression given by a simple recursive definition on the structure of  $e$ . The crucial property of these derivatives is that a string of the form  $aw$  (starting with the symbol  $a$ ) matches an expression  $e$  iff the suffix  $w$  matches the derivative  $D_a(e)$ . This suggests a streaming evaluation algorithm for regular expressions: the state is an expression, and upon arrival of a new symbol  $a$  the state is replaced by its derivative with respect to  $a$ . A refinement of Brzozowski's ideas was proposed by Antimirov [7] under the name of *partial derivatives*. He described a representation of derivatives as sets of partial derivatives, which corresponds closely to the construction of a NFA from an expression.

Given the success of automata-based techniques for the evaluation of plain regular expressions, it is worthwhile investigating whether similar ideas can be used for QRE evaluation. The well-studied model of weighted automata over semirings (see the monograph [16] for a broad survey) seems relevant, but unfortunately it is not expressive enough to handle the complex nesting of several different quantitative operations found in QREs. In particular, by the Kleene-Schützenberger theorem [22], weighted regular expressions can be easily translated into equivalent weighted automata and evaluated efficiently. On the other hand, QREs can be translated into the model of deterministic CRAs [4], but this translation incurs a doubly exponential blowup and is therefore not conducive to efficient evaluation. A hierarchical automaton model for the streaming computation of an interesting class of quantitative queries is introduced in [6], but its precise relationship to QREs and other automata formalisms remains to be clarified. A conclusively appropriate notion of automaton for the efficient evaluation of general quantitative queries has not been proposed yet, therefore a meaningful

**Table 1.** Complexity results under syntactic and semantic restrictions

Query language	Time-per-element and space complexity
Unrestricted, multiset semantics	Exponential in stream/query
Unrestricted, unambiguous semantics	Constant in stream, exponential in query
Strongly typed	Constant in stream, polynomial in query

investigation is the development of a notion of derivative. Indeed, in the present paper, we extend the notion of Brzozowski derivatives to QREs, and we show that there is a representation of QRE derivatives that gives rise to an efficient evaluation algorithm. This result offers a simple and clean alternative proof of why QREs can be efficiently evaluated, and it strongly suggests the possibility of an automata-based formulation of the evaluation algorithm. We should note here that derivatives have already been studied in the weighted setting [20], but the case of QREs is substantially different. See also [10] for an investigation of how Brzozowski derivatives can be extended to various algebraic structures.

**Outline of Paper.** In Sect. 2 we present the syntax and meaning of QREs. We consider two different natural semantics: (1) The *multiset semantics* allows for several output values for a given stream prefix, each of which corresponds to a different parse tree. (2) The *unambiguous semantics*, on the other hand, specifies the output to be undefined when the input stream can be parsed in more than one way. Thus, the unambiguous semantics ensures a single output value for each input sequence. In Sect. 3 we define derivatives of QREs by generalizing the classical notion of Brzozowski derivatives of regular expressions, and we propose an incremental evaluation algorithm based on derivatives. We also consider in Sect. 4 a representation of QRE derivatives that is analogous to Antimirov’s partial derivatives [7]. In the presence of intersection the number of distinct Antimirov derivatives (for plain regular expressions) is exponentially bounded by the size of the expression. We show how to obtain a similar bound for QRE evaluation using the unambiguous semantics. Finally, we consider in Sect. 5 a syntactic restriction for QREs [5] that guarantees unambiguous parsing. Our complexity results for the streaming evaluation problem are summarized in Table 1. Although the proposed derivative-based algorithm has the same time and space complexity as the previously known method [5], our approach here is cleaner and the analysis of the algorithm much simpler. We conclude in Sect. 6 with a summary of our results and directions for future work.

## 2 Quantitative Regular Expressions

The formalism of Quantitative Regular Expressions offers a declarative language for describing complex hierarchical computations on streams of data values. As an illustrative example, consider the application of patient monitoring, where the data stream is a time series of timestamped measurements produced by a

sensor attached to a patient. We want to analyze the data stream by first identifying regions of interest which we call “episodes”. These are maximal intervals where the measurements are above a fixed threshold value. Every episode is summarized by recording the average measurement, the maximum measurement, and the duration of the episode. The top-level query is an aggregation over the last 30 days (e.g., by calculating the average) of the episode statistics. This query imparts a hierarchical structure on the data stream by splitting it into episodes, bucketing episodes into days, and considering the last 30 days every time the aggregate statistics are computed. To describe this computation we need a language that supports regular constructs such as iteration (to express, e.g., that an episode is a sequence of measurements exceeding the threshold) and concatenation (to express, e.g., that an episode is followed by a sequence of measurements below the threshold), extended with quantitative operations for computing numerical aggregates (e.g., the maximum measurement of the sequence of measurements that constitute an episode).

In this section we present the syntax and semantics of Quantitative Regular Expressions. A QRE is interpreted over a stream of data values, and specifies for each finite prefix of the stream an output value. We consider two different semantics: the multiset semantics which records several different output possibilities, and the unambiguous semantics which only allows the output to be defined when it is uniquely determined. Since we are interested in computing well-defined *functions* on data streams, the multiset semantics is not satisfactory. We consider it here, however, because it is a natural semantics from a mathematical perspective, and it is useful for formulating and proving our results regarding efficient evaluation. The unambiguous semantics is simply a projection of the multiset semantics, and therefore several results w.r.t. to the multiset semantics transfer essentially unchanged to the unambiguous semantics.

To define QREs, we first choose a typed signature which describes the basic data types and operations for manipulating them. We fix a collection of *basic types*, and we write  $A, B, \dots$  to range over them, as well as a collection of *basic operations* on them, e.g.  $op : A_1 \times \dots \times A_k \rightarrow B$ . The identity function on  $D$  is written  $id_D : D \rightarrow D$ . For every basic type  $D$ , assume that we have fixed a collection of *atomic predicates*, so that the satisfiability of their Boolean combinations is decidable. We write  $\phi : D \rightarrow \mathbb{B}$  to indicate that  $\phi$  is a predicate on  $D$ , and  $true_D : D \rightarrow \mathbb{B}$  for the predicate that is always true. The *unit type*, with unique inhabitant  $\mathbf{def}$ , is  $\mathbb{U}$  and  $!_D : D \rightarrow \mathbb{U}$  is the unique function from  $D$  to  $\mathbb{U}$ . We also write  $\pi_1 : A \times B \rightarrow A$  and  $\pi_2 : A \times B \rightarrow B$  for the left and right projection respectively. We assume that the collection of basic operations contains all identities and projections, and is closed under pairing and function composition. For example, if  $op : A \times B \rightarrow C$  and  $a \in A$  are basic operations, then so is  $(op a) = \lambda b. op(a, b) : B \rightarrow C$ .

Every QRE is defined on a regular subset of stream prefixes, so we first introduce a variant of regular expressions with unary predicates to describe the domains of definition of QREs. For a basic type  $D$ , we define (*Symbolic*) *Regular Expressions (REs)* over  $D$  with the grammar  $r ::= \perp \mid \varepsilon \mid \phi \mid r \sqcup r \mid r \cdot r \mid r^* \mid r \sqcap r$ , where  $\phi : D \rightarrow \mathbb{B}$ . The expression  $r^n$  is abbreviation for  $r \cdot r \cdots r$  with  $r$  repeated

$$\begin{array}{c}
 \frac{}{\perp : \text{QRE}\langle D, C \rangle} \text{ (bottom)} \qquad \frac{c \in C}{\text{eps}(c) : \text{QRE}\langle D, C \rangle} \text{ (empty)} \\
 \frac{\text{satisfiable } \phi : D \rightarrow \mathbb{B} \quad op : D \rightarrow C}{\text{atom}(\phi, op) : \text{QRE}\langle D, C \rangle} \text{ (single item)} \qquad \frac{\mathbf{f}, \mathbf{g} : \text{QRE}\langle D, C \rangle}{\mathbf{f} \sqcup \mathbf{g} : \text{QRE}\langle D, C \rangle} \text{ (choice)} \\
 \frac{\mathbf{f} : \text{QRE}\langle D, A \rangle \quad \mathbf{g} : \text{QRE}\langle D, B \rangle \quad op : A \times B \rightarrow C}{\text{split}(\mathbf{f}, \mathbf{g}, op) : \text{QRE}\langle D, C \rangle} \text{ (split)} \\
 \frac{\text{init} : \text{QRE}\langle D, B \rangle \quad \text{body} : \text{QRE}\langle D, A \rangle \quad \mathcal{E}(\text{body}) = \emptyset \quad op : B \times A \rightarrow B}{\text{iter}(\text{init}, \text{body}, op) : \text{QRE}\langle D, B \rangle} \text{ (iteration)} \\
 \frac{\mathbf{f} : \text{QRE}\langle D, A \rangle \quad op : A \rightarrow B}{op(\mathbf{f}) : \text{QRE}\langle D, B \rangle} \text{ (application)} \\
 \frac{\mathbf{f} : \text{QRE}\langle D, A \rangle \quad \mathbf{g} : \text{QRE}\langle D, B \rangle \quad op : A \times B \rightarrow C}{op(\mathbf{f}, \mathbf{g}) : \text{QRE}\langle D, C \rangle} \text{ (combination)} \\
 \\
 \mathcal{E}(\perp) = \emptyset \qquad \mathcal{E}(\mathbf{f} \sqcup \mathbf{g}) = \mathcal{E}(\mathbf{f}) \uplus \mathcal{E}(\mathbf{g}) \\
 \mathcal{E}(\text{eps}(c)) = \{c\} \qquad \mathcal{E}(\text{split}(\mathbf{f}, \mathbf{g}, op)) = \text{M}(op)(\mathcal{E}(\mathbf{f}), \mathcal{E}(\mathbf{g})) \\
 \mathcal{E}(\text{atom}(\phi, op)) = \emptyset \qquad \mathcal{E}(\text{iter}(\mathbf{f}, \mathbf{g}, op)) = \mathcal{E}(\mathbf{f}) \\
 \mathcal{E}(op(\mathbf{f})) = \text{M}(op)(\mathcal{E}(\mathbf{f})) \qquad \mathcal{E}(op(\mathbf{f}, \mathbf{g})) = \text{M}(op)(\mathcal{E}(\mathbf{f}), \mathcal{E}(\mathbf{g}))
 \end{array}$$

**Fig. 1.** Syntax of Quantitative Regular Expressions (QREs) without  $\varepsilon$ -cycles.

$n$  times. We write  $r : \text{RE}\langle D \rangle$  to indicate that  $r$  is a regular expression over  $D$ . We define  $\llbracket r \rrbracket : D^* \rightarrow \mathbb{N}$  to be the weighted semantics of regular expressions without  $\varepsilon$ -cycles (i.e., no  $\varepsilon$ -cycles in the corresponding  $\varepsilon$ -NFA [19]) that counts the number of different parse trees.

$$\begin{array}{lll}
 \llbracket \perp \rrbracket w = 0 & \llbracket \varepsilon \rrbracket w = 0 \ (w \neq \varepsilon) & \llbracket r_1 \sqcup r_2 \rrbracket w = (\llbracket r_1 \rrbracket w) + (\llbracket r_2 \rrbracket w) \\
 \llbracket \varepsilon \rrbracket \varepsilon = 1 & \llbracket \phi \rrbracket d = 1 \ (d \models \phi) & \llbracket r_1 \cdot r_2 \rrbracket w = \sum_{w=uv} (\llbracket r_1 \rrbracket u) \cdot (\llbracket r_2 \rrbracket v) \\
 & \llbracket \phi \rrbracket d = 0 \ (d \not\models \phi) & \llbracket r^* \rrbracket w = \sum_{w=u_1 \dots u_n} (\llbracket r \rrbracket u_1) \cdots (\llbracket r \rrbracket u_n) \\
 & \llbracket \phi \rrbracket w = 0 \ (w \notin D) & \llbracket r_1 \sqcap r_2 \rrbracket w = (\llbracket r_1 \rrbracket w) \cdot (\llbracket r_2 \rrbracket w)
 \end{array}$$

This semantics corresponds to standard operations for formal power series [16].

In Fig. 1 we define *Quantitative Regular Expressions* (QREs), which we also call *queries*. The queries are typed, and we write  $\mathbf{f} : \text{QRE}\langle D, C \rangle$  to indicate that the query  $\mathbf{f}$  has *input type*  $D$  and *output type*  $C$ . The original definition of QREs in [5] was more general in that it involved an extra sort of typed parameters, and the outputs were essentially algebraic terms built from the parameters and the operations of the signature. The definition of [5] was motivated by expressiveness considerations, i.e. so that QREs capture exactly the class of regular functions over the same signature [3, 4, 14]. We consider here a simpler language, where the outputs are just values. This simplification makes QREs significantly more usable for practical queries and obviates the need for a *term simplification* procedure (see Sect. 3.2 of [5]) that depends on the nature of the data types and operations of the signature.

The language of Fig. 1 has eight core constructs: (1)  $\perp$  is undefined for every input sequence; (2)  $\mathbf{eps}(c)$  maps the empty stream to the output value  $c$ ; (3)  $\mathbf{atom}(\phi, op)$  maps an input stream consisting of a single item satisfying the predicate  $\phi$  to an output computed by applying the operation  $op$ ; (4)  $\mathbf{f} \sqcup \mathbf{g}$  nondeterministically chooses either  $\mathbf{f}$  or  $\mathbf{g}$  to apply; (5)  $\mathbf{split}(\mathbf{f}, \mathbf{g}, op)$  splits the input stream into two parts, applies the queries  $\mathbf{f}$  and  $\mathbf{g}$  to the left and right parts respectively, and combines their results using the operation  $op$ ; (6)  $\mathbf{iter}(\mathbf{init}, \mathbf{body}, op)$  splits the input into multiple parts  $uv_1v_2\dots v_n$ , applies  $\mathbf{init}$  to  $u$  (which gives  $b$ ) and  $\mathbf{body}$  to each  $v_i$  (which gives  $a_i$ ), and combines the sequence of values  $a_1a_2\dots a_n$  using the initial value  $b$  and the binary aggregation operation  $op$  in the style of the *fold combinator* used in functional programming; (7)  $op(\mathbf{f})$  applies the query  $\mathbf{f}$  and transforms its output using the operation  $op$ ; (8)  $op(\mathbf{f}, \mathbf{g})$  applies both  $\mathbf{f}$ ,  $\mathbf{g}$  and combines their results using the operation  $op$ .

The definition of queries is by mutual induction with the function  $\mathcal{E}$ , which sends a query of type  $\mathbf{QRE}\langle D, C \rangle$  to a finite multiset over its output type  $C$ . The  $\mathcal{E}$  function is meant to give the output of a query on the empty sequence, and it is used for the assumption  $\mathcal{E}(\mathbf{body}) = \emptyset$  for the iteration query  $\mathbf{iter}(\mathbf{init}, \mathbf{body}, op)$ . It is necessary to define the syntax of queries simultaneously with  $\mathcal{E}$  in order to eliminate queries with  $\varepsilon$ -cycles, that is, queries where the body  $\mathbf{g}$  of the iteration query  $\mathbf{iter}(\mathbf{f}, \mathbf{g}, op)$  matches the empty sequence  $\varepsilon$ . When the semantics is defined later, we will see that  $\varepsilon$ -cycles result in having an infinity of output values. This complication of  $\varepsilon$ -cycles appears also in the context of weighted automata and expressions [16]. The operation  $\uplus$  in Fig. 1 is multiset union. For a set  $X$ , we write  $\mathbf{M}(X)$  to denote the set of all finite multisets over  $X$ . For a unary operation  $op : A \rightarrow B$ , we define the lifting  $\mathbf{M}(op) : \mathbf{M}(A) \rightarrow \mathbf{M}(B)$  by  $\mathbf{M}(op)(X) = \{op(a) \mid a \in X\}$ . Similarly, for an operation  $op : A \times B \rightarrow C$ , we define the lifting  $\mathbf{M}(op) : \mathbf{M}(A) \times \mathbf{M}(B) \rightarrow \mathbf{M}(C)$  by  $\mathbf{M}(op)(X, Y) = \{op(a, b) \mid a \in X \text{ and } b \in Y\}$ .

**Multiset Semantics.** We give a denotational semantics of queries in terms of functions of type  $D^* \rightarrow \mathbf{M}(C)$ . We call this the *multiset semantics* of queries. The *domain* of  $f$  is the set of sequences for which the value of  $f$  is nonempty, i.e.  $\text{dom}(f) = \{w \in D^* \mid f(w) \neq \emptyset\}$ . The denotation of a query  $\mathbf{f} : \mathbf{QRE}\langle D, C \rangle$  is the function  $\llbracket \mathbf{f} \rrbracket : D^* \rightarrow \mathbf{M}(C)$ , where  $\llbracket \cdot \rrbracket$  is called the interpretation function and is defined by induction on the structure of queries as shown in Fig. 2. To reduce the notational clutter we sometimes write  $\llbracket \mathbf{f} \rrbracket w$  instead of  $\llbracket \mathbf{f} \rrbracket(w)$ . The semantics of iteration involves the *multiset fold* combinator  $\mathbf{mfold}$ , which is a generalization of the familiar fold combinator to multisets of values. The definitions of  $\mathbf{mfold}$  and  $\mathbf{fold}$  are by recursion on the length of the sequence. For example,  $\mathbf{fold}(s, op, a_1a_2) = op(op(s, a_1), a_2)$  and

$$\mathbf{mfold}(\{b_1, b_2\}, op, \{a_1\} \{a_2, a_3\}) = \\ \{op(op(b_1, a_1), a_2), op(op(b_1, a_1), a_3), op(op(b_2, a_1), a_2), op(op(b_2, a_1), a_3))\}.$$

For an iteration query  $\mathbf{h} = \mathbf{iter}(\mathbf{f}, \mathbf{g}, op)$ , the typing restriction  $\mathcal{E}(\mathbf{g}) = \emptyset$  implies that  $\llbracket \mathbf{g} \rrbracket \varepsilon = \emptyset$  (formally proved later in Theorem 5). So, to calculate the value

$$\begin{array}{ll}
 \llbracket \perp \rrbracket w = \emptyset & \llbracket \mathbf{f} \sqcup \mathbf{g} \rrbracket w = \llbracket \mathbf{f} \rrbracket w \uplus \llbracket \mathbf{g} \rrbracket w \\
 \llbracket \mathbf{eps}(c) \rrbracket \varepsilon = \{c\} & \llbracket \mathbf{split}(\mathbf{f}, \mathbf{g}, op) \rrbracket w = \biguplus_{w=uv} M(op)(\llbracket \mathbf{f} \rrbracket u, \llbracket \mathbf{g} \rrbracket v) \\
 \llbracket \mathbf{atom}(\phi, op) \rrbracket d = \{op(d)\}, \text{ if } d \models \phi & \llbracket op(\mathbf{f}) \rrbracket w = M(op)(\llbracket \mathbf{f} \rrbracket w) \\
 \llbracket \mathbf{atom}(\phi, op) \rrbracket d = \emptyset, \text{ if } d \not\models \phi & \llbracket op(\mathbf{f}, \mathbf{g}) \rrbracket w = M(op)(\llbracket \mathbf{f} \rrbracket w, \llbracket \mathbf{g} \rrbracket w) \\
 \\ 
 \llbracket \mathbf{iter}(\mathbf{f}, \mathbf{g}, op) \rrbracket w = \biguplus_{w=uv_1 \dots v_n} \mathbf{mfold}(\llbracket \mathbf{f} \rrbracket u, op, \llbracket \mathbf{g} \rrbracket v_1 \dots \llbracket \mathbf{g} \rrbracket v_n) \\
 \\ 
 \mathbf{fold} : B \times (B \times A \rightarrow B) \times A^* \rightarrow B & \mathbf{mfold} : M(B) \times (B \times A \rightarrow B) \times M(A)^* \rightarrow M(B) \\
 \mathbf{fold}(s, op, \varepsilon) = s & \mathbf{mfold}(S, op, \varepsilon) = S \\
 \mathbf{fold}(s, op, wa) = op(\mathbf{fold}(s, op, w), a) & \mathbf{mfold}(S, op, WX) = M(op)(\mathbf{mfold}(S, op, W), X)
 \end{array}$$

**Fig. 2.** Finite multiset semantics of Quantitative Regular Expressions without  $\varepsilon$ -cycles.

$\llbracket \mathbf{h} \rrbracket w$  (see Fig. 2) we only need to consider the splittings  $w = uv_1 \dots v_n$  of  $w$  where  $n \geq 0$  and  $v_i \neq \varepsilon$  for every  $i = 1, \dots, n$ .

For every input sequence  $w$ , the value  $\llbracket \mathbf{f} \rrbracket w$  is a finite multiset whose size is at most exponential in the size of  $w$ . More precisely,  $(\text{size of } \llbracket \mathbf{f} \rrbracket w) \leq 2^{|\mathbf{f}| \cdot |w|}$  for every query  $\mathbf{f} : \text{QRE}(D, C)$  and every sequence  $w \in D^*$ .

The multiset semantics of queries induces an equivalence relation on them, written as  $\equiv$ . Two queries are equivalent if their denotations are equal. We can then write equations such as  $\mathbf{split}(\mathbf{f}, \mathbf{eps}(b), op) \equiv op_b(\mathbf{f})$ , where  $op_b$  is the unary operation given by  $op_b(x) = op(x, b)$  for all  $x$ .

**Example 1.** The query  $\mathbf{f} = \mathbf{atom}(true_{\mathbb{N}}, id_{\mathbb{N}}) : \text{QRE}(\mathbb{N}, \mathbb{N})$  matches a single number and returns it. Using it as the body of an iteration, we write the query  $\mathbf{g} = \mathbf{iter}(\mathbf{f}, \mathbf{f}, +)$  which processes a nonempty sequence of numbers and returns their sum. Now, the query  $\mathbf{g}' = \mathbf{iter}(\mathbf{f}', \mathbf{f}', +)$ , where  $\mathbf{f}' = \mathbf{atom}(true_{\mathbb{N}}, \lambda x.1)$ , processes a nonempty sequence and returns its length. If  $\mathbf{div} : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{Q}$  and  $\mathbf{div}(x, y)$  is the result of dividing  $x$  by  $y$ , the query  $\mathbf{h} = \mathbf{div}(\mathbf{g}, \mathbf{g}') : \text{QRE}(\mathbb{N}, \mathbb{Q})$  calculates the average of a nonempty sequence of natural numbers.

**Rates.** The *rate* of a query is a symbolic regular expression that denotes its domain. It is defined by induction:

$$\begin{array}{ll}
 \mathbf{rate}(\perp) = \perp & \mathbf{rate}(\mathbf{f} \sqcup \mathbf{g}) = \mathbf{rate}(\mathbf{f}) \sqcup \mathbf{rate}(\mathbf{g}) \\
 \mathbf{rate}(\mathbf{eps}(c)) = \varepsilon & \mathbf{rate}(\mathbf{split}(\mathbf{f}, \mathbf{g}, op)) = \mathbf{rate}(\mathbf{f}) \cdot \mathbf{rate}(\mathbf{g}) \\
 \mathbf{rate}(\mathbf{atom}(\phi, op)) = \phi & \mathbf{rate}(\mathbf{iter}(\mathbf{f}, \mathbf{g}, op)) = \mathbf{rate}(\mathbf{f}) \cdot \mathbf{rate}(\mathbf{g})^* \\
 \mathbf{rate}(op(\mathbf{f})) = \mathbf{rate}(\mathbf{f}) & \mathbf{rate}(op(\mathbf{f}, \mathbf{g})) = \mathbf{rate}(\mathbf{f}) \sqcap \mathbf{rate}(\mathbf{g})
 \end{array}$$

Notice that the value of  $\mathbf{rate}$  is always an expression without  $\varepsilon$ -cycles.

**Unambiguous Semantics.** We defined previously the multiset semantics of queries, which allows a query to have several (finitely many) outputs for a given input sequence. Now, we take the viewpoint that a query should specify a unique output value (or be undefined) for a given input sequence. This means that we want to ignore output multisets of cardinality greater than one, which we do

by setting the output to be undefined. This is the *unambiguous semantics* for queries. For a query  $\mathbf{f} : \text{QRE}\langle D, C \rangle$ , this is given formally as follows:

$$\langle\langle \mathbf{f} \rangle\rangle w = \theta(\llbracket \mathbf{f} \rrbracket w), \text{ where } \theta(X) = X \text{ if } |X| = 1 \text{ and } \theta(X) = \emptyset \text{ otherwise}$$

So,  $\langle\langle \mathbf{f} \rangle\rangle$  is a function  $D^* \rightarrow \mathbf{M}(C)$  so that each output multiset is of cardinality at most one, which we call an *unambiguous function*. This means that  $\langle\langle \mathbf{f} \rangle\rangle$  can be also represented as a partial function  $D^* \rightarrow C$ . We say that a query  $\mathbf{f}$  is *unambiguous* when the multiset meaning  $\llbracket \mathbf{f} \rrbracket$  is an unambiguous function. This is equivalent to  $\llbracket \mathbf{f} \rrbracket = \langle\langle \mathbf{f} \rangle\rangle$ , which says that the multiset and unambiguous semantics coincide. See the papers [12, 13] for recent surveys of unambiguity in traditional automata theory.

It is necessary that the  $\llbracket - \rrbracket$  semantics records the multiplicity of each output value, otherwise  $\langle\langle - \rangle\rangle$  cannot be defined as a projection or  $\llbracket - \rrbracket$ . Indeed, we see in the example below that we can write queries that return exactly one output value of multiplicity greater than one.

**Example 2.** The query  $\mathbf{f} = \text{iter}(\text{eps}(0), \text{atom}(\text{true}_{\mathbb{N}}, \text{id}_{\mathbb{N}}), +) : \text{QRE}\langle \mathbb{N}, \mathbb{N} \rangle$ , which processes a sequence of natural numbers and returns their sum, is unambiguous. The query  $\mathbf{g} = \text{iter}(\text{eps}(\text{def}), \text{atom}(\text{true}_{\mathbb{N}}, !_{\mathbb{N}}), !_{\mathbb{N} \times \mathbb{N}}) : \text{QRE}\langle \mathbb{N}, \mathbb{U} \rangle$ , which matches a sequence of natural numbers and returns nothing, is also unambiguous. The query  $\text{split}(\mathbf{g}, \mathbf{f}, \pi_2) : \text{QRE}\langle \mathbb{N}, \mathbb{N} \rangle$ , which matches a sequence of natural numbers and returns the sum of every suffix of the sequence, is ambiguous because every sequence of length  $\ell$  can be parsed in  $\ell + 1$  different ways. The query that matches sequences of length at least two and returns the sum of the last two elements is  $\text{split}(\mathbf{g}, \text{split}(\text{atom}(\text{true}_{\mathbb{N}}, \text{id}_{\mathbb{N}}), \text{atom}(\text{true}_{\mathbb{N}}, \text{id}_{\mathbb{N}}), +), \pi_2) : \text{QRE}\langle \mathbb{N}, \mathbb{N} \rangle$  and is unambiguous. The query  $\text{split}(\mathbf{f}, \mathbf{f}, +) : \text{QRE}\langle \mathbb{N}, \mathbb{N} \rangle$  is ambiguous but single-valued: for a sequence of length  $\ell$ , it returns the sum of its elements with multiplicity equal to  $\ell + 1$ .

The unambiguous semantics of queries induces an equivalence relation on them, written as  $\approx$ . Two queries  $\mathbf{f}$  and  $\mathbf{g}$  are  $\approx$ -equivalent if their denotations  $\langle\langle \mathbf{f} \rangle\rangle$  and  $\langle\langle \mathbf{g} \rangle\rangle$  are equal. We observe that the equivalence relation  $\approx$  is strictly coarser than  $\equiv$ , that is,  $\mathbf{f} \equiv \mathbf{g}$  implies  $\mathbf{f} \approx \mathbf{g}$  and there exist queries that are  $\approx$ -equivalent but  $\equiv$ -inequivalent.

**Observation 3.** A finite multiset  $Q$  of queries of the same type can also be thought of as a query, namely the finite choice over the queries of  $Q$ . We now want to find some sufficient conditions for reducing the cardinality of  $Q$ , while preserving its meaning under the unambiguous semantics. This will turn out to be useful later in Sect. 4, where an evaluation algorithm for QREs using Antimirov derivatives is presented (Fig. 3 and Theorem 9).

Suppose  $Q$  contains the queries  $\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_k$  with  $k \geq 2$  that have the same rate, that is  $r = \text{rate}(\mathbf{f}_i)$  for every  $i$ , and there is no other query in  $Q$  that has this rate. The condition on the rates implies that all functions  $\llbracket \mathbf{f}_i \rrbracket$  have the same domain  $D_1 = \text{dom}(\llbracket \mathbf{f}_i \rrbracket)$ , and moreover all functions  $\langle\langle \mathbf{f}_i \rangle\rangle$  have the same domain  $D_2 = \text{dom}(\langle\langle \mathbf{f}_i \rangle\rangle) \subseteq D_1$ .

- (1) Remove all of  $\mathbf{f}_1, \dots, \mathbf{f}_k$  (*wrong*): We claim that  $Q' = Q \setminus \{\mathbf{f}_1, \dots, \mathbf{f}_k\}$  is not necessarily equivalent to  $Q$ . Suppose that  $Q = \{\mathbf{f}_1, \dots, \mathbf{f}_k, \mathbf{g}\}$ , the queries of  $Q$  are all unambiguous, and the domains  $\text{dom}(\llbracket \mathbf{g} \rrbracket) = \text{dom}(\llbracket \mathbf{g}' \rrbracket)$  and  $D_1 = D_2$  intersect. So, all queries of  $Q$  are defined on some sequence  $w$ . Then,  $\llbracket Q \rrbracket$  is not defined on  $w$  because the cardinality of  $\llbracket Q \rrbracket w$  is greater than one, but  $\llbracket Q' \rrbracket$  is defined on  $w$  and equal to  $\llbracket \mathbf{g}' \rrbracket w$ . So,  $Q \not\approx Q'$ .
- (2) If  $k \geq 3$  then remove  $\mathbf{f}_3, \dots, \mathbf{f}_k$  and keep  $\mathbf{f}_1, \mathbf{f}_2$  (*correct*): Define the multiset  $Q' = Q \setminus \{\mathbf{f}_3, \dots, \mathbf{f}_k\}$ . We claim that  $Q \approx Q'$ . Let  $w$  be an arbitrary sequence. If  $w$  belongs to the domain of the functions  $\llbracket \mathbf{f}_i \rrbracket$ , then both  $\llbracket Q \rrbracket$  and  $\llbracket Q' \rrbracket$  are undefined on  $w$ , because the cardinalities of  $\llbracket Q \rrbracket w$  and  $\llbracket Q' \rrbracket w$  are greater than one. Suppose now that  $w$  does not belong to the domain of the functions  $\llbracket \mathbf{f}_i \rrbracket$ . Then,  $\llbracket Q \rrbracket w = \llbracket Q \setminus \{\mathbf{f}_i \mid i\} \rrbracket w = \llbracket Q' \rrbracket w$  and hence  $\llbracket Q \rrbracket w = \llbracket Q' \rrbracket w$ .

This means that we can always reduce  $Q$  so that it has at most two queries with the same rate, while preserving its unambiguous semantics.

### 3 Brzowski Derivative

We introduce in this section the *Brzowski derivative* [11] of Quantitative Regular Expressions, which is a straightforward adaption of derivatives for classical regular expressions. The main property of these derivatives is that their semantic counterpart agrees with the syntactic counterpart. This agreement property gives as a corollary the existence of an incremental algorithm for evaluating QREs on streams of data items.

**Example 4.** The Brzowski derivative  $\mathcal{D}_a^{\mathbb{B}}(r)$  of a regular expression  $r$  w.r.t. the letter  $a$  denotes the language that results from the language of  $r$  by removing the starting  $a$  letter from those words that start with  $a$ . For example,  $r = (a \sqcup b)^*bb$  denotes the language of all strings over  $\Sigma = \{a, b\}$  that end in  $bb$ ,  $\mathcal{D}_a^{\mathbb{B}}(r) = (a \sqcup b)^*bb = r$ ,  $\mathcal{D}_{ab}^{\mathbb{B}}(r) = \mathcal{D}_b^{\mathbb{B}}(\mathcal{D}_a^{\mathbb{B}}(r)) = \mathcal{D}_b^{\mathbb{B}}(r) = r \sqcup b$  and  $\mathcal{D}_{abb}^{\mathbb{B}}(r) = \mathcal{D}_b^{\mathbb{B}}(\mathcal{D}_{ab}^{\mathbb{B}}(r)) = \mathcal{D}_b^{\mathbb{B}}(r \sqcup b) = r \sqcup b \sqcup \varepsilon$ . The string  $abb$  matches  $r$  because the empty string  $\varepsilon$  matches the derivative  $\mathcal{D}_{abb}^{\mathbb{B}}(r)$ .

A query of type  $\text{QRE}(\mathbb{N}, \mathbb{N})$  that is similar in form to the regex  $r$  is  $\mathbf{k} = \text{split}(\mathbf{g}, \mathbf{h}, \text{max})$ , where  $\mathbf{g} = \text{iter}(\text{eps}(0), \mathbf{f}_e \sqcup \mathbf{f}_o, +)$ ,  $\mathbf{f}_e = \text{atom}(\text{even}_{\mathbb{N}}, \text{id}_{\mathbb{N}})$ ,  $\mathbf{f}_o = \text{atom}(\text{odd}_{\mathbb{N}}, \text{id}_{\mathbb{N}})$ , and  $\mathbf{h} = \text{split}(\mathbf{f}_o, \mathbf{f}_o, +)$ . The unambiguous query  $\mathbf{k}$  matches sequences that end with two odd numbers and returns the maximum of  $x, y$  where  $y$  is the sum of the last two numbers and  $x$  is the sum of the rest of the numbers. The extension of derivatives to QREs should give rise to the following calculations:

$$\begin{aligned}
 \mathcal{D}_4^{\mathbb{B}}(\mathbf{k}) &= \text{split}(\text{iter}((\lambda x.0 + x)(\text{eps}(4)), \mathbf{f}_e \sqcup \mathbf{f}_o, +), \mathbf{h}, \text{max}) \\
 \mathcal{D}_{43}^{\mathbb{B}}(\mathbf{k}) &= \text{split}(\text{iter}((\lambda x.4 + x)(\text{eps}(3)), \mathbf{f}_e \sqcup \mathbf{f}_o, +), \mathbf{h}, \text{max}) \sqcup \\
 &\quad (\lambda x. \text{max}(4, x))((\lambda x.3 + x)(\mathbf{f}_o)) \\
 \mathcal{D}_{435}^{\mathbb{B}}(\mathbf{k}) &= \text{split}(\text{iter}((\lambda x.7 + x)(\text{eps}(5)), \mathbf{f}_e \sqcup \mathbf{f}_o, +), \mathbf{h}, \text{max}) \sqcup \\
 &\quad (\lambda x. \text{max}(7, x))((\lambda x.5 + x)(\mathbf{f}_o)) \sqcup \\
 &\quad (\lambda x. \text{max}(4, x))((\lambda x.3 + x)(\text{eps}(5)))
 \end{aligned}$$

$$\begin{aligned} \mathcal{D}_{4351}^{\mathbb{B}}(\mathbf{k}) &= \text{split}(\text{iter}((\lambda x.12 + x)(\mathbf{eps}(1)), \mathbf{f}_e \sqcup \mathbf{f}_o, +), \mathbf{h}, \text{max}) \sqcup \\ &\quad (\lambda x. \text{max}(12, x))((\lambda x.1 + x)(\mathbf{f}_o)) \sqcup \\ &\quad (\lambda x. \text{max}(7, x))((\lambda x.5 + x)(\mathbf{eps}(1))) \end{aligned}$$

From the above we notice that  $\varepsilon$  matches  $\mathcal{D}_{435}^{\mathbb{B}}(\mathbf{k})$  with value  $\text{max}(4, 3 + 5) = 8$ , and it also matches  $\mathcal{D}_{4351}^{\mathbb{B}}(\mathbf{k})$  with value  $\text{max}(7, 5 + 1) = 7$ .

A simple streaming algorithm that implements the computation described by  $\mathbf{h}$  can be given by maintaining the following state: the sum of all elements so far except for the last two, and the two most recent elements. The reader can observe that the derivatives calculated earlier record this information, and it can be found inside the queries that are constructed. For example, in the derivative  $\mathcal{D}_{4351}^{\mathbb{B}}(\mathbf{k})$  we see (last line) the sum 7 and the elements 5, 1. The structure of the derivative also encodes how these three numbers should be combined to produce the output  $\text{max}(7, 5 + 1) = 7$ .  $\square$

For a function  $f : D^* \rightarrow \mathbb{M}(C)$ , we define the *semantic derivative*  $\mathcal{D}_u(f) : D^* \rightarrow \mathbb{M}(C)$  with respect to the sequence  $u \in D^*$  of data items as

$$\mathcal{D}_u(f) w = f(uw) \text{ for all } w \in D^*.$$

An immediate consequence is that  $\mathcal{D}_v(\mathcal{D}_u(f)) = \mathcal{D}_{uv}(f)$  for all sequences  $u$  and  $v$  in  $D^*$ . Moreover,  $f(u) = \mathcal{D}_u(f)(\varepsilon)$  for every sequence  $u \in D^*$ . For a query  $\mathbf{f}$  of type  $\text{QRE}\langle D, C \rangle$  and a data item  $d \in D$ , the (*syntactic*) *Brzozowski derivative*  $\mathcal{D}_d^{\mathbb{B}}(\mathbf{f})$  of  $\mathbf{f}$  w.r.t.  $d$  is also a query of type  $\text{QRE}\langle D, C \rangle$ .

$$\begin{aligned} \mathcal{D}_d^{\mathbb{B}}(\mathbf{eps}(c)) &= \mathcal{D}_d^{\mathbb{B}}(\perp) = \perp & \mathcal{D}_d^{\mathbb{B}}(\mathbf{f} \sqcup \mathbf{g}) &= \mathcal{D}_d^{\mathbb{B}}(\mathbf{f}) \sqcup \mathcal{D}_d^{\mathbb{B}}(\mathbf{g}) \\ \mathcal{D}_d^{\mathbb{B}}(\mathbf{atom}(\phi, op)) &= \perp, \text{ if } d \not\models \phi & \mathcal{D}_d^{\mathbb{B}}(op(\mathbf{f})) &= op(\mathcal{D}_d^{\mathbb{B}}(\mathbf{f})) \\ \mathcal{D}_d^{\mathbb{B}}(\mathbf{atom}(\phi, op)) &= \mathbf{eps}(op(d)), \text{ if } d \models \phi & \mathcal{D}_d^{\mathbb{B}}(op(\mathbf{f}, \mathbf{g})) &= op(\mathcal{D}_d^{\mathbb{B}}(\mathbf{f}), \mathcal{D}_d^{\mathbb{B}}(\mathbf{g})) \\ \mathcal{D}_d^{\mathbb{B}}(\mathbf{split}(\mathbf{f}, \mathbf{g}, op)) &= \mathbf{split}(\mathcal{D}_d^{\mathbb{B}}(\mathbf{f}), \mathbf{g}, op) \sqcup \bigsqcup_{a \in \mathcal{E}(\mathbf{f})} (op a)(\mathcal{D}_d^{\mathbb{B}}(\mathbf{g})) \\ \mathcal{D}_d^{\mathbb{B}}(\mathbf{iter}(\mathbf{f}, \mathbf{g}, op)) &= \mathbf{iter}(\mathcal{D}_d^{\mathbb{B}}(\mathbf{f}), \mathbf{g}, op) \sqcup \bigsqcup_{b \in \mathcal{E}(\mathbf{f})} \mathbf{iter}((op b)(\mathcal{D}_d^{\mathbb{B}}(\mathbf{g})), \mathbf{g}, op) \end{aligned}$$

The derivative  $\mathcal{D}_w^{\mathbb{B}}(\mathbf{f})$  w.r.t. a sequence  $w \in D^*$  is defined by induction on  $w$ :  $\mathcal{D}_\varepsilon^{\mathbb{B}}(\mathbf{f}) = \mathbf{f}$  and  $\mathcal{D}_{dw}^{\mathbb{B}}(\mathbf{f}) = \mathcal{D}_w^{\mathbb{B}}(\mathcal{D}_d^{\mathbb{B}}(\mathbf{f}))$ . A crucial result is the correspondence between semantic and syntactic derivatives:

**Theorem 5 (Derivative Agreement).** For every query  $\mathbf{f}$  of type  $\text{QRE}\langle D, C \rangle$  and every data item  $d \in D$ , we have that  $\mathcal{E}(\mathbf{f}) = \llbracket \mathbf{f} \rrbracket \varepsilon$  and  $\mathcal{D}_d(\llbracket \mathbf{f} \rrbracket) = \llbracket \mathcal{D}_d^{\mathbb{B}}(\mathbf{f}) \rrbracket$ .

Theorem 5 suggests immediately an algorithm for evaluating queries. Given a sequence  $w = d_1 d_2 \dots d_n$  and a query  $\mathbf{f}$ , we notice that  $\llbracket \mathbf{f} \rrbracket w$  is equal to

$$\mathcal{D}_w(\llbracket \mathbf{f} \rrbracket) \varepsilon = \mathcal{D}_{d_n} \dots \mathcal{D}_{d_1}(\llbracket \mathbf{f} \rrbracket) \varepsilon = \llbracket \mathcal{D}_{d_n}^{\mathbb{B}} \dots \mathcal{D}_{d_1}^{\mathbb{B}}(\mathbf{f}) \rrbracket \varepsilon = \mathcal{E}(\mathcal{D}_{d_n}^{\mathbb{B}} \dots \mathcal{D}_{d_1}^{\mathbb{B}}(\mathbf{f})).$$

So, to compute the value of a query on a given input, we iteratively calculate the derivative of the query w.r.t. each input item, and finally apply the  $\mathcal{E}$  function. We suggest an optimization of this evaluation procedure by incorporating

query rewriting to eliminate subqueries that cannot contribute to the result. The equations  $op(\perp) \equiv \perp$  and

$$\begin{array}{llll} \perp \sqcup \mathbf{f} \equiv \mathbf{f} & \text{split}(\perp, \mathbf{f}, op) \equiv \perp & op(\perp, \mathbf{f}) \equiv \perp & \text{iter}(\perp, \mathbf{f}, op) \equiv \perp \\ \mathbf{f} \sqcup \perp \equiv \mathbf{f} & \text{split}(\mathbf{f}, \perp, op) \equiv \perp & op(\mathbf{f}, \perp) \equiv \perp & \text{iter}(\mathbf{f}, \perp, op) \equiv \mathbf{f} \end{array}$$

are all valid, that is, they are true for every query  $\mathbf{f}$ . If we orient these equations from left to right, then we get a rewrite system for simplifying queries. We will assume that our Brzozowski derivative-based evaluation procedure simplifies all intermediate queries as much as possible (according to this rewrite system).

**Example 6.** Consider the always-true predicate  $true_{\mathbb{N}} : \mathbb{N} \rightarrow \mathbb{B}$ , the identity function  $id_{\mathbb{N}} : \mathbb{N} \rightarrow \mathbb{N}$ , the constant zero function  $\lambda x.0 : \mathbb{N} \rightarrow \mathbb{N}$ , and the binary sum  $+: \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$ . Using these operations we write the queries

$$\mathbf{f} = \text{atom}(true_{\mathbb{N}}, id_{\mathbb{N}}) \quad \mathbf{g} = \text{atom}(true_{\mathbb{N}}, \lambda x.0) \quad \mathbf{h} = \text{iter}(\text{eps}(0), \mathbf{f} \sqcup \mathbf{g}, +)$$

of type  $\text{QRE}(\mathbb{N}, \mathbb{N})$ . The top-level query is  $\mathbf{h}$ , which maps an input sequence of natural numbers into the multiset of all possible partial sums over those numbers. First, notice that the derivatives of  $\mathbf{f}$  and  $\mathbf{g}$  w.r.t. the number  $d \in \mathbb{N}$  are  $\mathcal{D}_d^{\mathbb{B}}(\mathbf{f}) = \text{eps}(d)$  and  $\mathcal{D}_d^{\mathbb{B}}(\mathbf{g}) = \text{eps}(0)$  respectively. Moreover,  $\mathcal{E}(\text{eps}(0)) = \{0\}$ . We will use the derivative-based evaluation algorithm to find the value of the query  $\mathbf{h}$  on the input sequence  $abc \in \mathbb{N}^*$ . This amounts to calculating  $\mathcal{E}(\mathcal{D}_c^{\mathbb{B}}(\mathcal{D}_b^{\mathbb{B}}(\mathcal{D}_a^{\mathbb{B}}(\mathbf{h}))))$ , since Theorem 5 implies that it equals  $\llbracket \mathbf{h} \rrbracket abc$ . The steps of this calculation are:

$$\begin{aligned} \mathcal{D}_a^{\mathbb{B}}(\mathbf{h}) &= \text{iter}(\text{eps}(a) \sqcup \text{eps}(0), \mathbf{f} \sqcup \mathbf{g}, +) \\ \mathcal{E}(\mathcal{D}_a^{\mathbb{B}}(\mathbf{h})) &= \{a, 0\} \\ \mathcal{D}_b^{\mathbb{B}}(\mathcal{D}_a^{\mathbb{B}}(\mathbf{h})) &= \text{iter}((\lambda x.a + x)(\text{eps}(b) \sqcup \text{eps}(0)), \mathbf{f} \sqcup \mathbf{g}, +) \sqcup \\ &\quad \text{iter}((\lambda x.0 + x)(\text{eps}(b) \sqcup \text{eps}(0)), \mathbf{f} \sqcup \mathbf{g}, +) \\ \mathcal{E}(\mathcal{D}_b^{\mathbb{B}}(\mathcal{D}_a^{\mathbb{B}}(\mathbf{h}))) &= \{a + b, a, b, 0\} \\ \mathcal{D}_c^{\mathbb{B}}(\mathcal{D}_b^{\mathbb{B}}(\mathcal{D}_a^{\mathbb{B}}(\mathbf{h}))) &= \text{iter}((\lambda x.(a + b) + x)(\text{eps}(c) \sqcup \text{eps}(0)), \mathbf{f} \sqcup \mathbf{g}, +) \sqcup \\ &\quad \text{iter}((\lambda x.a + x)(\text{eps}(c) \sqcup \text{eps}(0)), \mathbf{f} \sqcup \mathbf{g}, +) \sqcup \\ &\quad \text{iter}((\lambda x.b + x)(\text{eps}(c) \sqcup \text{eps}(0)), \mathbf{f} \sqcup \mathbf{g}, +) \sqcup \\ &\quad \text{iter}((\lambda x.0 + x)(\text{eps}(c) \sqcup \text{eps}(0)), \mathbf{f} \sqcup \mathbf{g}, +) \\ \mathcal{E}(\mathcal{D}_c^{\mathbb{B}}(\mathcal{D}_b^{\mathbb{B}}(\mathcal{D}_a^{\mathbb{B}}(\mathbf{h})))) &= \{a + b + c, a + b, a + c, a, b + c, b, c, 0\} \end{aligned}$$

We have used implicitly the rewrite rules for eliminating  $\perp$  as much as possible.

The previous example shows that there are queries whose evaluation requires an enormous amount of computational resources. Given that the size of the output can be exponential in the size of the input sequence and the size of the query, we have exponential time and space requirements for every evaluation algorithm.

## 4 Antimirov Derivative

The evaluation of the general QREs of Fig. 1 w.r.t. the multiset semantics is inherently expensive, since the output itself can be of size exponential in the size of the query and the input sequence. So, we focus here on the evaluation of QREs w.r.t. the unambiguous semantics, and our goal is to describe a streaming algorithm that uses resources that are independent of the size of the input stream. Using a variant of the partial derivatives of Antimirov [7], we show that this is indeed possible. We obtain a streaming algorithm for evaluation that uses space and time-per-element exponential in the size of the query and independent of the stream length. The crucial idea for this algorithm is that we can prune the Antimirov derivative to contain only a couple of QREs with the same rate without changing its unambiguous meaning. Antimirov-style derivatives are preferable for the results of this section, because the representation itself encodes many valid equations on queries that are useful for proving the result. This is similar to classical regular expressions where the Antimirov derivatives encode the ACI rules (associativity, commutativity and idempotence) by virtue of the set-based representation.

**Example 7.** We continue with Example 4 to consider Antimirov derivatives. Recall that  $r = (a \sqcup b)^*bb$ , and that Antimirov derivatives are a set-based representation of Brzozowski derivatives. That is,  $\mathcal{D}_a^\Delta(r) = \{r\}$ ,  $\mathcal{D}_{ab}^\Delta(r) = \mathcal{D}_b^\Delta(\{r\}) = \{r, b\}$  and  $\mathcal{D}_{abb}^\Delta(r) = \mathcal{D}_b^\Delta(\{r, b\}) = \{r, b, \varepsilon\}$ . For the query  $k$  of Example 4 we calculate using Antimirov-style derivatives:

$$\begin{aligned} \mathcal{D}_4^\Delta(k) &= \{\text{split}(\text{iter}((\lambda x.0 + x)(\text{eps}(4)), \mathbf{f}_e \sqcup \mathbf{f}_o, +), \mathbf{h}, \text{max})\} \\ \mathcal{D}_{43}^\Delta(k) &= \{\text{split}(\text{iter}((\lambda x.4 + x)(\text{eps}(3)), \mathbf{f}_e \sqcup \mathbf{f}_o, +), \mathbf{h}, \text{max}), \\ &\quad (\lambda x. \text{max}(4, x))((\lambda x.3 + x)(\mathbf{f}_o))\} \\ \mathcal{D}_{435}^\Delta(k) &= \{\text{split}(\text{iter}((\lambda x.7 + x)(\text{eps}(5)), \mathbf{f}_e \sqcup \mathbf{f}_o, +), \mathbf{h}, \text{max}), \\ &\quad (\lambda x. \text{max}(7, x))((\lambda x.5 + x)(\mathbf{f}_o)), \\ &\quad (\lambda x. \text{max}(4, x))((\lambda x.3 + x)(\text{eps}(5)))\} \end{aligned}$$

Since the choice operation  $\sqcup$  for queries is not idempotent, we need multisets for the representation of QRE Antimirov derivatives.  $\square$

For a regular expression  $r : \text{RE}\langle D \rangle$  we define  $\mathcal{E}(r) \in \{0, 1\}$  (the two-element Boolean algebra with join operation  $+$  and meet operation  $\cdot$ ) and the *Antimirov derivative*  $\mathcal{D}_d^\Delta(r)$  w.r.t.  $d \in D$ , which is a set of regular expressions of type  $\text{RE}\langle D \rangle$ :

$$\begin{aligned} \mathcal{E}(\perp) &= 0 & \mathcal{D}_d^\Delta(\varepsilon) &= \mathcal{D}_d^\Delta(\perp) = \emptyset \\ \mathcal{E}(\varepsilon) &= 1 & \mathcal{D}_d^\Delta(\phi) &= \emptyset, \text{ if } d \not\equiv \phi \\ \mathcal{E}(\phi) &= 0 & \mathcal{D}_d^\Delta(\phi) &= \{\varepsilon\}, \text{ if } d \models \phi \\ \mathcal{E}(r_1 \sqcup r_2) &= \mathcal{E}(r_1) + \mathcal{E}(r_2) & \mathcal{D}_d^\Delta(r_1 \sqcup r_2) &= \mathcal{D}_d^\Delta(r_1) \cup \mathcal{D}_d^\Delta(r_2) \\ \mathcal{E}(r_1 \sqcap r_2) &= \mathcal{E}(r_1) \cdot \mathcal{E}(r_2) & \mathcal{D}_d^\Delta(r_1 \sqcap r_2) &= \mathcal{D}_d^\Delta(r_1) \sqcap \mathcal{D}_d^\Delta(r_2) \\ \mathcal{E}(r_1 \cdot r_2) &= \mathcal{E}(r_1) \cdot \mathcal{E}(r_2) & \mathcal{D}_d^\Delta(r_1 \cdot r_2) &= \mathcal{D}_d^\Delta(r_1) \cdot r_2 \cup \mathcal{E}(r_1) \cdot \mathcal{D}_d^\Delta(r_2) \\ \mathcal{E}(r^*) &= 1 & \mathcal{D}_d^\Delta(r^*) &= \mathcal{D}_d^\Delta(r) \cdot r^* \end{aligned}$$

**State** : Finite multiset  $S$  of queries of type  $\text{QRE}\langle D, C \rangle$ .

**Initialization**() : Set the state  $S$  to the singleton multiset  $\{\mathbf{f}\}$ .

**Update**( $d$ ) : Compute the Antimirov derivative  $S' = \mathcal{D}_d^\Delta(S) = \uplus\{\mathcal{D}_d^\Delta(\mathbf{g}) \mid \mathbf{g} \in S\}$ . Now, iterate the following conditional modification until no further changes can be made: if  $S'$  contains at least three queries that have the same rate, keep only two of them and remove the rest. Finally, set  $S$  to be equal to  $S'$ .

**Output**() : If  $\mathcal{E}(S) = \uplus\{\mathcal{E}(\mathbf{g}) \mid \mathbf{g} \in S\}$  is the singleton multiset  $\{c\}$ , then return the value  $c$ . Otherwise, the output is undefined.

**Fig. 3.** Streaming evaluation algorithm for an arbitrary query  $\mathbf{f} : \text{QRE}\langle D, C \rangle$ , with respect to the unambiguous semantics of QREs.

For subsets  $X, Y$  of expressions we have used in the definition the abbreviations:  $X \sqcap Y = \{r \sqcap s \mid r \in X, s \in Y\}$ ,  $X \cdot s = \{r \cdot s \mid r \in X\}$ ,  $0 \cdot X = \emptyset$  and  $1 \cdot X = X$ .

**Lemma 8.** For every regular expression  $r$ , the set  $\bigcup_{w \in D^*} \mathcal{D}_w^\Delta(r)$  of all derivatives of  $r$  is of size exponential in  $r$ .

For a query  $\mathbf{f}$  of type  $\text{QRE}\langle D, C \rangle$  and a data item  $d \in D$ , the *Antimirov derivative*  $\mathcal{D}_d^\Delta(\mathbf{f})$  of  $\mathbf{f}$  w.r.t.  $d$  is a finite *multiset* of queries of type  $\text{QRE}\langle D, C \rangle$ .

$$\begin{aligned}
 \mathcal{D}_d^\Delta(\mathbf{eps}(c)) &= \mathcal{D}_d^\Delta(\perp) = \emptyset & \mathcal{D}_d^\Delta(\mathbf{f} \sqcup \mathbf{g}) &= \mathcal{D}_d^\Delta(\mathbf{f}) \uplus \mathcal{D}_d^\Delta(\mathbf{g}) \\
 \mathcal{D}_d^\Delta(\mathbf{atom}(\phi, op)) &= \emptyset, \text{ if } d \not\models \phi & \mathcal{D}_d^\Delta(op(\mathbf{f})) &= op(\mathcal{D}_d^\Delta(\mathbf{f})) \\
 \mathcal{D}_d^\Delta(\mathbf{atom}(\phi, op)) &= \{\mathbf{eps}(op(d))\}, \text{ if } d \models \phi & \mathcal{D}_d^\Delta(op(\mathbf{f}, \mathbf{g})) &= op(\mathcal{D}_d^\Delta(\mathbf{f}), \mathcal{D}_d^\Delta(\mathbf{g})) \\
 \mathcal{D}_d^\Delta(\mathbf{split}(\mathbf{f}, \mathbf{g}, op)) &= \mathbf{split}(\mathcal{D}_d^\Delta(\mathbf{f}), \mathbf{g}, op) \uplus \uplus_{a \in \mathcal{E}(\mathbf{f})} (op a)(\mathcal{D}_d^\Delta(\mathbf{g})) \\
 \mathcal{D}_d^\Delta(\mathbf{iter}(\mathbf{f}, \mathbf{g}, op)) &= \mathbf{iter}(\mathcal{D}_d^\Delta(\mathbf{f}), \mathbf{g}, op) \uplus \uplus_{b \in \mathcal{E}(\mathbf{f})} \mathbf{iter}((op b)(\mathcal{D}_d^\Delta(\mathbf{g})), \mathbf{g}, op)
 \end{aligned}$$

We have used above convenient abbreviations like  $op(X) = \{op(\mathbf{f}) \mid \mathbf{f} \in X\}$ . The derivative  $\mathcal{D}_w^\Delta(\mathbf{f})$  w.r.t. a sequence  $w \in D^*$  is defined by induction on  $w$ :  $\mathcal{D}_\varepsilon^\Delta(\mathbf{f}) = \{\mathbf{f}\}$  and  $\mathcal{D}_{dw}^\Delta(\mathbf{f}) = \bigcup\{\mathcal{D}_w^\Delta(\mathbf{g}) \mid \mathbf{g} \in \mathcal{D}_d^\Delta(\mathbf{f})\}$ .

**Theorem 9.** The algorithm of Fig. 3 solves the evaluation problem for QREs w.r.t. the unambiguous semantics. It requires space and time-per-element that are constant in the length of the stream, and exponential in the size of the query.

*Proof.* To prove that the algorithm of Fig. 3 is correct, we observe that it satisfies the following crucial invariant: after consuming the input  $w \in D^*$ , the multiset  $S$  is  $\approx$ -equivalent to the derivative  $\mathcal{D}_w^\Delta(\mathbf{f})$ . This is because the removal step of the **Update** procedure preserves the unambiguous meaning of  $S$  (recall Observation 3). It remains to see that the Antimirov derivative is a streamlined representation of the Brzozowski derivative. That is, for every query  $\mathbf{f} : \text{QRE}\langle D, C \rangle$  and every data item  $d \in D$ , it holds that  $\mathcal{D}_d^\Delta(\mathbf{f}) \equiv \bigsqcup \mathcal{D}_d^\Delta(\mathbf{f})$ . The equivalences

$$\begin{aligned}
\mathbf{f} \sqcup \mathbf{g} &\equiv \mathbf{g} \sqcup \mathbf{f} & op(\mathbf{f} \sqcup \mathbf{g}, \mathbf{h}) &\equiv op(\mathbf{f}, \mathbf{h}) \sqcup op(\mathbf{g}, \mathbf{h}) \\
op(\mathbf{f} \sqcup \mathbf{g}) &\equiv op(\mathbf{f}) \sqcup op(\mathbf{g}) & op(\mathbf{f}, \mathbf{g} \sqcup \mathbf{h}) &\equiv op(\mathbf{f}, \mathbf{g}) \sqcup op(\mathbf{f}, \mathbf{h}) \\
&& split(\mathbf{f} \sqcup \mathbf{g}, \mathbf{h}) &\equiv split(\mathbf{f}, \mathbf{h}) \sqcup split(\mathbf{g}, \mathbf{h}) \\
&& iter(\mathbf{f} \sqcup \mathbf{g}, \mathbf{h}, op) &\equiv iter(\mathbf{f}, \mathbf{h}, op) \sqcup iter(\mathbf{g}, \mathbf{h}, op)
\end{aligned}$$

are used in the proof of this claim. Essentially, the above equations are encoded in the derivatives by way of their representation as multisets of queries. The output procedure returns  $\langle\langle S \rangle\rangle \varepsilon = \langle\langle \mathcal{D}_w^A(\mathbf{f}) \rangle\rangle \varepsilon = \theta(\llbracket \mathcal{D}_w^E(\mathbf{f}) \rrbracket \varepsilon) = \theta(\llbracket \mathbf{f} \rrbracket w) = \langle\langle \mathbf{f} \rangle\rangle w$ .

The Antimirov derivatives on QREs correspond closely to the Antimirov derivatives on regular expressions. For every query  $\mathbf{g} \in \mathcal{D}_d^A(\mathbf{f})$ , it holds that  $rate(\mathbf{g}) \in \mathcal{D}_d^A(rate(\mathbf{f}))$ . Because of the pruning step, the state  $S$  contains at most two queries for every possible rate of the derivatives. Since there are at most exponentially many derivatives of  $rate(\mathbf{f})$  by Lemma 8, the cardinality of  $S$  is also at most exponential in  $\mathbf{f}$  and independent of the size of the input sequence. The time to process each element is also exponential in  $|\mathbf{f}|$ .  $\square$

**Example 10.** The query  $\mathbf{f} = \mathbf{atom}(true_{\mathbb{N}}, id_{\mathbb{N}})$  is of type  $\text{QRE}(\mathbb{N}, \mathbb{N})$  and it maps a single natural number to itself. The query  $\mathbf{g} = \mathbf{iter}(\mathbf{eps}(0), \mathbf{f}, +)$  maps any sequence of natural numbers to their sum. We calculate the derivative of  $\mathbf{g}$ :

$$\begin{aligned}
\mathcal{E}(\mathbf{eps}(0)) &= \{0\} & \mathcal{D}_d^A(\mathbf{f}) &= \{\mathbf{eps}(d)\} & \mathcal{D}_d^A(\mathbf{eps}(0)) &= \emptyset \\
\mathcal{D}_d^A(\mathbf{g}) &= \mathbf{iter}(\mathcal{D}_d^A(\mathbf{eps}(0)), \mathbf{f}, +) \uplus \mathbf{iter}((\lambda x.0 + x)(\mathcal{D}_d^A(\mathbf{f})), \mathbf{f}, +) \\
&= \{\mathbf{iter}((\lambda x.0 + x)(\mathbf{eps}(d)), \mathbf{f}, +)\}
\end{aligned}$$

For  $x, y \in \mathbb{N}$  define  $\mathbf{g}_{xy} = \mathbf{iter}((\lambda z.x + z)(\mathbf{eps}(y)), \mathbf{f}, +)$ . The derivative of  $\mathbf{g}_{xy}$  is:

$$\begin{aligned}
\mathcal{D}_d^A(\mathbf{g}_{xy}) &= \mathbf{iter}(\mathcal{D}_d^A((\lambda z.x + z)(\mathbf{eps}(y))), \mathbf{f}, +) \uplus \\
&\quad \mathbf{iter}((\lambda z.(x + y) + z)(\mathcal{D}_d^A(\mathbf{f})), \mathbf{f}, +) \\
&= \{\mathbf{iter}((\lambda z.(x + y) + z)(\mathbf{eps}(d)), \mathbf{f}, +)\} = \mathbf{g}_{x+y,d}
\end{aligned}$$

because  $\mathcal{E}((\lambda z.x + z)(\mathbf{eps}(y))) = \{x + y\}$  and  $\mathcal{D}_d^A((\lambda z.x + z)(\mathbf{eps}(y))) = \emptyset$ .  $\square$

## 5 Strongly Typed Queries and Hierarchical Derivatives

Following [5], we consider a syntactic restriction of QREs that ensures unambiguity of parsing. This means that the multiset and unambiguous semantics coincide, thus this subclass of QREs inherently describes well-defined functions on data streams. Together with an additional restriction that demands in expressions of the form  $op(\mathbf{f}, \mathbf{g})$  the subqueries  $\mathbf{f}$  and  $\mathbf{g}$  to have the same domain, it can be ensured that the evaluation of QREs can be performed efficiently, that is, using space and time-per-element that is polynomial in the size of the query and independent of the stream. Such an algorithm is proposed in [5] and implemented in [21], but it is hard to describe and even harder to analyze. We consider

here an alternative approach based on derivatives, which gives rise to an evaluation algorithm of the same complexity that is much easier to describe and prove correct. The main technical tool is a novel representation of derivatives, which we call *hierarchical derivatives*. This is a very space-efficient representation and is crucial for obtaining our complexity bounds.

Before defining the subclass of strongly typed queries, we introduce some definitions that will be used for formalizing the idea of uniqueness of parsing. The languages  $L_1, L_2$  are said to be *unambiguously concatenable* if for every word  $w \in L_1 \cdot L_2$  there are unique  $w_1 \in L_1, w_2 \in L_2$  with  $w = w_1 w_2$ . The language  $L$  is said to be *unambiguously iterable* if for every word  $w \in L^*$  there is a unique integer  $n \geq 0$  and unique  $w_i \in L$  with  $w = w_1 \cdots w_n$ . These definitions extend to regular expressions in the obvious way.

We say that a query is *strongly typed* if the following hold: (1) for every subquery  $\mathbf{f} \sqcup \mathbf{g}$  the rates  $\text{rate}(\mathbf{f})$  and  $\text{rate}(\mathbf{g})$  are disjoint, (2) for every subquery  $\text{split}(\mathbf{f}, \mathbf{g}, \text{op})$  the rates  $\text{rate}(\mathbf{f})$  and  $\text{rate}(\mathbf{g})$  are unambiguously concatenable. (3) for every subquery  $\text{iter}(\mathbf{f}, \mathbf{g}, \text{op})$  the rate  $\text{rate}(\mathbf{g})$  is unambiguously iterable and  $\text{rate}(\mathbf{f}), \text{rate}(\mathbf{g})^*$  are unambiguously concatenable, and (4) for every subquery  $\text{op}(\mathbf{f}, \mathbf{g})$  the rates  $\text{rate}(\mathbf{f})$  and  $\text{rate}(\mathbf{g})$  are equivalent. It is shown in [5] that checking whether a query is strongly typed can be done in polynomial-time.

**Lemma 11.** If the query  $\mathbf{f}$  is strongly typed, then  $\llbracket \mathbf{f} \rrbracket$  is an unambiguous function. So, the multiset and unambiguous semantics coincide, i.e.  $\llbracket \mathbf{f} \rrbracket = \llbracket \mathbf{f} \rrbracket$ .

The main problem with the Antimirov derivative of Sect. 4 is the treatment of  $\text{op}(\mathbf{f}, \mathbf{g})$ , where  $\mathcal{D}_d^\wedge(\text{op}(\mathbf{f}, \mathbf{g})) = \{\text{op}(\mathbf{f}', \mathbf{g}') \mid \mathbf{f}' \in \mathcal{D}_d^\wedge(\mathbf{f}), \mathbf{g}' \in \mathcal{D}_d^\wedge(\mathbf{g})\}$ . This definition corresponds to a cartesian product of derivatives, thus causing a quadratic blowup in the number of possible derivatives. As we will see later in Example 14, this blowup can materialize even in the context of strongly typed queries. Since the output combination operation can nest with other regular constructs, the Antimirov representation can result in an exponential blowup, which is avoidable in the strongly typed case. In order to prove this, we need a new representation that avoids this “cartesian product” problem of the Antimirov derivative by allowing sets of queries to be used as subexpressions.

We thus generalize the syntax of strongly typed queries to allow finite sets of queries as subexpressions. Intuitively, these finite sets of queries extend the choice operation  $\sqcup$  to a finite number of arguments. Since the subqueries of  $\mathbf{f} \sqcup \mathbf{g}$  must have disjoint domains, the  $\sqcup$  constructor is associative and commutative, which means that we can represent  $(\mathbf{f} \sqcup \mathbf{g}) \sqcup \mathbf{h}$  as  $\{\mathbf{f}, \mathbf{g}, \mathbf{h}\}$ . This finite choice constructor can be nested arbitrarily with the other query constructors. In Fig. 4 we see the formal definition, where queries and *q-sets* (finite sets of queries) are defined by mutual induction. To reduce the notational clutter, we sometimes write the query  $\mathbf{f}$  instead of the q-set  $\{\mathbf{f}\}$ , for example  $\text{iter}(\mathbf{f}, \mathbf{g}, \text{op})$  instead of  $\text{iter}(\{\mathbf{f}\}, \{\mathbf{g}\}, \text{op})$ . The rate of a generalized query is defined in the usual way, the only difference being  $\text{rate}(\{\mathbf{f}_1, \dots, \mathbf{f}_k\}) = \text{rate}(\mathbf{f}_1) \sqcup \dots \sqcup \text{rate}(\mathbf{f}_k)$ .

For the expressions of Fig. 4 we define in Fig. 5 a new kind of derivative, called the (*syntactic*) *hierarchical derivative*  $\mathcal{D}_d^h(\cdot)$ , which maps a query of type

$$\begin{array}{c}
\frac{c \in C}{\text{eps}(c) : \text{QRE}\langle D, C \rangle} \quad \frac{\text{satisfiable } \phi : D \rightarrow \mathbb{B} \quad \text{op} : D \rightarrow C}{\text{atom}(\phi, \text{op}) : \text{QRE}\langle D, C \rangle} \\
\frac{\text{F} : \text{QSET}\langle D, A \rangle \quad \text{G} : \text{QSET}\langle D, B \rangle \quad \text{op} : A \times B \rightarrow C}{\text{F}, \text{G} \neq \emptyset \quad \text{rate}(\text{F}), \text{rate}(\text{G}) \text{ unambiguously concatenable}} \\
\frac{\text{split}\langle \text{F}, \text{G}, \text{op} \rangle : \text{QRE}\langle D, C \rangle}{} \\
\frac{\text{op} : B \times A \rightarrow B \quad \text{F} : \text{QSET}\langle D, B \rangle \quad \text{rate}(\text{F}), \text{rate}(\text{G})^* \text{ unambiguously concatenable}}{\text{F}, \text{G} \neq \emptyset \quad \text{G} : \text{QSET}\langle D, A \rangle \quad \text{rate}(\text{G}) \text{ unambiguously iterable}} \\
\text{iter}\langle \text{F}, \text{G}, \text{op} \rangle : \text{QRE}\langle D, B \rangle \\
\frac{\text{F} \neq \emptyset : \text{QSET}\langle D, A \rangle \quad \text{op} : A \rightarrow B}{\text{op}\langle \text{F} \rangle : \text{QRE}\langle D, C \rangle} \quad \frac{\text{F} : \text{QSET}\langle D, A \rangle \quad \text{G} : \text{QSET}\langle D, B \rangle \quad \text{op} : A \times B \rightarrow C}{\text{F}, \text{G} \neq \emptyset \quad \text{rate}(\text{F}) \text{ and } \text{rate}(\text{G}) \text{ are equivalent}} \\
\text{op}\langle \text{F}, \text{G} \rangle : \text{QRE}\langle D, C \rangle \\
\frac{\text{f}_1, \dots, \text{f}_k : \text{QRE}\langle D, C \rangle \quad \text{rate}(\text{f}_1), \dots, \text{rate}(\text{f}_k) \text{ pairwise disjoint}}{\{\text{f}_1, \dots, \text{f}_k\} : \text{QSET}\langle D, C \rangle}
\end{array}$$

**Fig. 4.** Generalized syntax for strongly typed QREs.

$$\begin{array}{l}
\mathcal{D}_d^{\text{H}}(\{\mathbf{f}_i \mid i \in I\}) = \bigcup_{i \in I} \mathcal{D}_d^{\text{H}}(\mathbf{f}_i) \quad \mathcal{D}_d^{\text{H}}(\text{atom}(\phi, \text{op})) = \emptyset, \text{ if } d \not\models \phi \\
\mathcal{D}_d^{\text{H}}(\text{eps}(c)) = \emptyset \quad \mathcal{D}_d^{\text{H}}(\text{atom}(\phi, \text{op})) = \{\text{eps}(\text{op}(d))\}, \text{ if } d \models \phi \\
\mathcal{D}_d^{\text{H}}(\text{op}\langle \text{F} \rangle) = \{\text{op}\langle \mathcal{D}_d^{\text{H}}(\text{F}) \rangle\} \quad \mathcal{D}_d^{\text{H}}(\text{op}\langle \text{F}, \text{G} \rangle) = \{\text{op}\langle \mathcal{D}_d^{\text{H}}(\text{F}), \mathcal{D}_d^{\text{H}}(\text{G}) \rangle\} \\
\mathcal{D}_d^{\text{H}}(\text{split}\langle \text{F}, \text{G}, \text{op} \rangle) = \{\text{split}\langle \mathcal{D}_d^{\text{H}}(\text{F}), \text{G}, \text{op} \rangle\} \cup \{(\text{op } a)\langle \mathcal{D}_d^{\text{H}}(\text{G}) \rangle \mid a \in \mathcal{E}(\mathbf{f})\} \\
\mathcal{D}_d^{\text{H}}(\text{iter}\langle \text{F}, \text{G}, \text{op} \rangle) = \{\text{iter}\langle \mathcal{D}_d^{\text{H}}(\text{F}), \text{G}, \text{op} \rangle\} \cup \{\text{iter}\langle (\text{op } b)\langle \mathcal{D}_d^{\text{H}}(\text{G}) \rangle, \text{G}, \text{op} \rangle \mid b \in \mathcal{E}(\text{F})\} \\
\mathcal{E}(\text{eps}(c)) = \{c\} \quad \mathcal{E}(\text{op}\langle \text{F} \rangle) = \text{M}(\text{op})(\mathcal{E}(\text{F})) \\
\mathcal{E}(\text{atom}(\phi, \text{op})) = \emptyset \quad \mathcal{E}(\text{op}\langle \text{F}, \text{G} \rangle) = \text{M}(\text{op})(\mathcal{E}(\text{F}), \mathcal{E}(\text{G})) \\
\mathcal{E}(\{\mathbf{f}_i\}_i) = \bigcup_i \mathcal{E}(\mathbf{f}_i) \quad \mathcal{E}(\text{split}\langle \text{F}, \text{G}, \text{op} \rangle) = \text{M}(\text{op})(\mathcal{E}(\text{F}), \mathcal{E}(\text{G})) \\
\mathcal{E}(\text{iter}\langle \text{F}, \text{G}, \text{op} \rangle) = \mathcal{E}(\text{F})
\end{array}$$

**Fig. 5.** Hierarchical derivatives for generalized QREs.

$\text{QRE}\langle D, C \rangle$  or a q-set of type  $\text{QSET}\langle D, C \rangle$  to a q-set of type  $\text{QSET}\langle D, C \rangle$ . The hierarchical derivative  $\mathcal{D}_w^{\text{H}}(\text{F})$  w.r.t. a sequence  $w \in D^*$  is defined by induction on  $w$ :  $\mathcal{D}_\varepsilon^{\text{H}}(\text{F}) = \text{F}$  and  $\mathcal{D}_{dw}^{\text{H}}(\text{F}) = \mathcal{D}_w^{\text{H}}(\mathcal{D}_d^{\text{H}}(\text{F}))$ . For every sequence  $u \in D^*$ , the derivative  $\mathcal{D}_u^{\text{H}}(\text{split}\langle \text{F}, \text{G}, \text{op} \rangle)$  is of the form:

$$\{\text{split}\langle \mathcal{D}_u^{\text{H}}(\text{F}), \text{G}, \text{op} \rangle, \text{op}\langle a_1, \mathcal{D}_{v_1}^{\text{H}}(\text{G}) \rangle, \dots, \text{op}\langle a_n, \mathcal{D}_{v_n}^{\text{H}}(\text{G}) \rangle\}$$

for some  $a_i \in A$  and sequences  $v_i \in D^*$ . Because of unambiguity, the rates of the derivatives  $\mathcal{D}_{v_1}^{\text{H}}(\text{G}), \dots, \mathcal{D}_{v_n}^{\text{H}}(\text{G})$  are pairwise disjoint. Similarly, for every sequence  $u \in D^*$ , the derivative  $\mathcal{D}_u^{\text{H}}(\text{iter}\langle \text{F}, \text{G}, \text{op} \rangle)$  is of the form:

$$\{\text{iter}\langle \mathcal{D}_u^{\text{H}}(\text{F}), \text{G}, \text{op} \rangle, \text{iter}\langle \text{op}\langle b_1, \mathcal{D}_{v_1}^{\text{H}}(\text{G}) \rangle, \text{G}, \text{op} \rangle, \dots, \text{iter}\langle \text{op}\langle b_n, \mathcal{D}_{v_n}^{\text{H}}(\text{G}) \rangle, \text{G}, \text{op} \rangle\}$$

for some  $b_i \in B$  and sequences  $v_i \in D^*$ . Again, because of unambiguity, the rates of the derivatives  $\mathcal{D}_u^h(\mathbf{F})$ ,  $\mathcal{D}_{v_1}^h(\mathbf{G})$ ,  $\dots$ ,  $\mathcal{D}_{v_n}^h(\mathbf{G})$  are pairwise disjoint. A key technical lemma to obtain an efficient evaluation algorithm is that the derivatives of strongly typed queries are of polynomial size. To show this, first we define a reasonable notion of size of queries:

$$\begin{aligned}
 \text{size}(\{\mathbf{f}_i\}_i) &= \sum_i \text{size}(\mathbf{f}_i) & \text{size}(\text{op}\langle \mathbf{F} \rangle) &= 1 + \text{size}(\mathbf{F}) \\
 \text{size}(\text{eps}(c)) &= 1 & \text{size}(\text{op}\langle \mathbf{F}, \mathbf{G} \rangle) &= 1 + \text{size}(\mathbf{F}) + \text{size}(\mathbf{G}) \\
 \text{size}(\text{atom}(\phi, \text{op})) &= 2 & \text{size}(\text{split}\langle \mathbf{F}, \mathbf{G}, \text{op} \rangle) &= 1 + \text{size}(\mathbf{F}) + \text{size}(\mathbf{G}) \\
 & & \text{size}(\text{iter}\langle \mathbf{F}, \mathbf{G}, \text{op} \rangle) &= 1 + \text{size}(\mathbf{F}) + \text{size}(\mathbf{G})
 \end{aligned}$$

In order to obtain the desired lemma, the claim has to be strengthened:

**Lemma 12.** Let  $\mathbf{F}$  be a q-set in  $\text{QSET}\langle D, C \rangle$  and  $u_1, \dots, u_n$  be sequences over  $D$  such that the rates of  $\mathcal{D}_{u_1}^h(\mathbf{F})$ ,  $\dots$ ,  $\mathcal{D}_{u_n}^h(\mathbf{F})$  are pairwise disjoint. Then:

- (1) At most  $\text{size}(\mathbf{F})$  of the sets  $\mathcal{D}_{u_1}^h(\mathbf{F})$ ,  $\dots$ ,  $\mathcal{D}_{u_n}^h(\mathbf{F})$  are nonempty.
- (2) The space needed to represent  $\mathcal{D}_{u_1}^h(\mathbf{F}) \cup \dots \cup \mathcal{D}_{u_n}^h(\mathbf{F})$  is bounded by  $\text{size}(\mathbf{F})^2$ .

*Proof.* To prove the lemma, we must extend the same claim to queries as well. The proof then proceeds by induction on the structure of queries and q-sets. The base cases  $\text{eps}(c)$  and  $\text{atom}(\phi, \text{op})$  and the step case  $\text{op}\langle \mathbf{F} \rangle$  are easy.

For the case of a q-set  $\mathbf{F} = \{\mathbf{f}_1, \dots, \mathbf{f}_m\}$  we first notice for an arbitrary  $j$  in  $\{1, \dots, m\}$  that:  $\text{rate}(\mathbf{f}_j) \subseteq \text{rate}(\mathbf{F})$  and therefore  $\text{rate}(\mathcal{D}_v^h(\mathbf{f}_j)) \subseteq \text{rate}(\mathcal{D}_v^h(\mathbf{F}))$  for every  $v \in D^*$ . It follows that the rates of the derivatives  $\mathcal{D}_{u_1}^h(\mathbf{f}_j)$ ,  $\dots$ ,  $\mathcal{D}_{u_n}^h(\mathbf{f}_j)$  are pairwise disjoint. To show part (1), we observe that

$$\{i \mid \mathcal{D}_{u_i}^h(\mathbf{F}) \neq \emptyset\} = \{i \mid \mathcal{D}_{u_i}^h(\mathbf{f}_1) \cup \dots \cup \mathcal{D}_{u_i}^h(\mathbf{f}_m) \neq \emptyset\} = \bigcup_{j=1}^m \{i \mid \mathcal{D}_{u_i}^h(\mathbf{f}_j) \neq \emptyset\}.$$

By the induction hypothesis, the size of this set is bounded by  $\sum_{j=1}^m \text{size}(\mathbf{f}_j) = \text{size}(\mathbf{F})$ . Now,

$$\begin{aligned}
 \text{space}(\bigcup_{i=1}^n \mathcal{D}_{u_i}^h(\mathbf{F})) &= \text{space}(\bigcup_{i=1}^n \bigcup_{j=1}^m \mathcal{D}_{u_i}^h(\mathbf{f}_j)) \\
 &= \sum_{j=1}^m \text{space}(\bigcup_{i=1}^n \mathcal{D}_{u_i}^h(\mathbf{f}_j)) \\
 &\leq \sum_{j=1}^m \text{size}(\mathbf{f}_j)^2,
 \end{aligned}$$

which is less than  $(\text{size}(\mathbf{f}_1) + \dots + \text{size}(\mathbf{f}_m))^2 = \text{size}(\mathbf{F})^2$ .

For the case of the query  $\mathbf{h} = \text{op}\langle \mathbf{F}, \mathbf{G} \rangle$ , we first recall that  $\text{rate}(\mathbf{h}) = \text{rate}(\mathbf{F}) \equiv \text{rate}(\mathbf{G})$ . So, the hypotheses of the lemma hold for both  $\mathbf{F}$  and  $\mathbf{G}$ . This means that at most  $\min(\text{size}(\mathbf{F}), \text{size}(\mathbf{G})) \leq \text{size}(\mathbf{h})$  of the sets  $\mathcal{D}_{u_1}^h(\mathbf{h})$ ,  $\dots$ ,  $\mathcal{D}_{u_n}^h(\mathbf{h})$  are nonempty. For part (2), we have that

$$\begin{aligned}
 \text{space}(\bigcup_{i=1}^n \mathcal{D}_{u_i}^h(\mathbf{h})) &= \text{space}(\{\text{op}\langle \mathcal{D}_{u_i}^h(\mathbf{F}), \mathcal{D}_{u_i}^h(\mathbf{G}) \rangle \mid i = 1, \dots, n\}) \\
 &\leq \text{size}(\mathbf{F}) + \left(\sum_{i=1}^n \text{space}(\mathcal{D}_{u_i}^h(\mathbf{F}))\right) + \left(\sum_{i=1}^n \text{space}(\mathcal{D}_{u_i}^h(\mathbf{G}))\right) \\
 &\leq \text{size}(\mathbf{F}) + \text{size}(\mathbf{F})^2 + \text{size}(\mathbf{G})^2,
 \end{aligned}$$

which is less than  $(1 + \text{size}(\mathbf{F}) + \text{size}(\mathbf{G}))^2 = \text{size}(\text{op}\langle \mathbf{F}, \mathbf{G} \rangle)^2$ .

**State** : Q-set  $\mathbf{H}$  of type  $\text{QSET}\langle D, C \rangle$ .

**Initialization**() : Set the state  $\mathbf{H}$  to the singleton q-set  $\{\mathbf{f}\}$ .

**Update**( $d$ ) : Replace  $\mathbf{H}$  by its hierarchical derivative  $\mathcal{D}_d^{\mathbf{H}}(\mathbf{H})$ .

**Output**() : If  $\mathcal{E}(\mathbf{H})$  is the singleton multiset  $\{c\}$ , then return the value  $c$ . Otherwise, the output is undefined.

**Fig. 6.** Streaming evaluation algorithm for a strongly typed query  $\mathbf{f} : \text{QRE}\langle D, C \rangle$ .

We handle now the case  $\mathbf{h} = \text{split}\langle \mathbf{F}, \mathbf{G}, \text{op} \rangle$ . As discussed previously, the union  $\mathcal{D}_{u_1}^{\mathbf{h}}(\mathbf{h}) \cup \dots \cup \mathcal{D}_{u_n}^{\mathbf{h}}(\mathbf{h})$  of the derivatives is of the form

$$\mathbf{H} = \{\text{split}\langle \mathcal{D}_{u_1}^{\mathbf{h}}(\mathbf{F}), \mathbf{G}, \text{op} \rangle, \dots, \text{split}\langle \mathcal{D}_{u_n}^{\mathbf{h}}(\mathbf{F}), \mathbf{G}, \text{op} \rangle, \\ \text{op}\langle a_1, \mathcal{D}_{v_1}^{\mathbf{h}}(\mathbf{G}) \rangle, \dots, \text{op}\langle a_n, \mathcal{D}_{v_m}^{\mathbf{h}}(\mathbf{G}) \rangle\}$$

for some  $a_1, \dots, a_m \in A$  and  $v_1, \dots, v_m \in D^*$ . Now,  $\mathbf{H}$  is unambiguous and therefore the q-sets  $\mathcal{D}_{u_1}^{\mathbf{h}}(\mathbf{F}), \dots, \mathcal{D}_{u_n}^{\mathbf{h}}(\mathbf{F})$  are pairwise disjoint, and similarly the q-sets  $\mathcal{D}_{v_1}^{\mathbf{h}}(\mathbf{G}), \dots, \mathcal{D}_{v_m}^{\mathbf{h}}(\mathbf{G})$  are pairwise disjoint. From the induction hypothesis (part 1) for  $\mathbf{F}$  and  $\mathbf{G}$ , we get that  $\mathbf{H}$  is of size  $\leq \text{size}(\mathbf{F}) + \text{size}(\mathbf{G})$ , which implies part (1) for  $\mathbf{h}$ . To measure the space needed to represent  $\mathbf{H}$ , we first observe that the subquery  $\mathbf{G}$  is shared by several queries in  $\mathbf{H}$  and therefore we can replace every  $\mathbf{G}$  occurrence with a pointer to a representation of  $\mathbf{G}$ . Total space:

$$\text{space}(\mathbf{H}) \leq \text{size}(\mathbf{F})^2 + 2 \cdot \text{size}(\mathbf{F}) + \text{size}(\mathbf{G})^2 + \text{size}(\mathbf{G}),$$

which is less than  $(1 + \text{size}(\mathbf{F}) + \text{size}(\mathbf{G}))^2 = \text{size}(\mathbf{h})^2$ .

Finally, we consider the case  $\mathbf{h} = \text{iter}\langle \mathbf{F}, \mathbf{G}, \text{op} \rangle$  of iteration. As discussed previously, the union  $\mathcal{D}_{u_1}^{\mathbf{h}}(\mathbf{h}) \cup \dots \cup \mathcal{D}_{u_n}^{\mathbf{h}}(\mathbf{h})$  of the derivatives is of the form

$$\mathbf{H} = \{\text{iter}\langle \mathcal{D}_{u_1}^{\mathbf{h}}(\mathbf{F}), \mathbf{G}, \text{op} \rangle, \dots, \text{iter}\langle \mathcal{D}_{u_n}^{\mathbf{h}}(\mathbf{F}), \mathbf{G}, \text{op} \rangle, \\ \text{iter}\langle \text{op}\langle b_1, \mathcal{D}_{v_1}^{\mathbf{h}}(\mathbf{G}) \rangle, \mathbf{G}, \text{op} \rangle, \dots, \text{iter}\langle \text{op}\langle b_n, \mathcal{D}_{v_m}^{\mathbf{h}}(\mathbf{G}) \rangle, \mathbf{G}, \text{op} \rangle\}$$

for some  $b_1, \dots, b_m \in B$  and  $v_1, \dots, v_m \in D^*$ . The total space requirements are:

$$\text{space}(\mathbf{H}) \leq \text{size}(\mathbf{F})^2 + 2 \cdot \text{size}(\mathbf{F}) + \text{size}(\mathbf{G})^2 + 3 \cdot \text{size}(\mathbf{G}),$$

which is less than  $(1 + \text{size}(\mathbf{F}) + \text{size}(\mathbf{G}))^2 = \text{size}(\mathbf{h})^2$ .  $\square$

Lemma 12 establishes, in particular, that the hierarchical derivative of any strongly typed q-set  $\mathbf{F}$  w.r.t. any sequence is of size at most quadratic in the size of  $\mathbf{F}$ . Using this fact, we can prove the main theorem of this paper:

**Theorem 13.** The algorithm of Fig. 6 solves the evaluation problem for strongly typed QREs. It requires space and time-per-element that are constant in the length of the stream, and polynomial in the size of the query.

*Proof.* First, we observe that the hierarchical derivative is simply a different representation of the Brzozowski derivative, which is streamlined for space efficiency. This implies that  $\mathcal{D}_d^{\mathbf{B}}(\mathbf{f})$ ,  $\mathcal{D}_d^{\mathbf{A}}(\mathbf{f})$  and  $\mathcal{D}_d^{\mathbf{H}}(\mathbf{f})$  are all  $\equiv$ -equivalent when

$f : \text{QRE}\langle D, C \rangle$  is strongly typed. The algorithm of Fig. 6 satisfies the invariant: after consuming input  $w \in D^*$ , the q-set  $\mathbb{H}$  is equal to the hierarchical derivative  $\mathcal{D}_w^{\mathbb{H}}(f)$ . The correctness of the algorithm then follows immediately from the semantic agreement of the hierarchical derivative with the Brzozowski derivative. At every step of the computation the state  $\mathbb{H}$  is of the form  $\mathcal{D}_w^{\mathbb{H}}(f)$ , where  $w$  is the input sequence seen so far. Lemma 12 give us immediately that  $\mathbb{H}$  can be represented using space that is quadratic in the size of the input query. It follows that the time to process each element is also quadratic in the query.  $\square$

**Example 14.** The queries  $\mathbf{h}_1, \mathbf{h}_2$  below calculate the maximum and minimum respectively of two consecutive natural numbers, where  $\mathbf{f} = \text{atom}(true_{\mathbb{N}}, id_{\mathbb{N}})$ . The top-level query  $\mathbf{k}$  shown below is strongly typed, and it calculates the average of the maximum and minimum of the last two elements of the stream.

$$\begin{aligned} \mathbf{h}_1 &= \text{split}(\mathbf{f}, \mathbf{f}, \text{max}) & \mathbf{k}_1 &= \text{split}(\mathbf{g}, \mathbf{h}_1, \pi_2) \\ \mathbf{h}_2 &= \text{split}(\mathbf{f}, \mathbf{f}, \text{min}) & \mathbf{k}_2 &= \text{split}(\mathbf{g}, \mathbf{h}_2, \pi_2) \\ \mathbf{g} &= \text{iter}(\text{eps}(0), \text{atom}(true_{\mathbb{N}}, \lambda x.0), \pi_2) & \mathbf{k} &= \text{avg}(\mathbf{k}_1, \mathbf{k}_2) \end{aligned}$$

The Antimirov and hierarchical derivatives of  $\mathbf{g}$  are the following:

$$\begin{aligned} \{\mathbf{g}'\} &= \mathcal{D}_w^{\mathbb{A}}(\mathbf{g}) = \{\text{iter}((\lambda x.\pi_2(0, x))(\text{eps}(0)), \text{atom}(true_{\mathbb{N}}, \lambda x.0), \pi_2)\} \\ \mathbf{g}'' &= \mathcal{D}_w^{\mathbb{H}}(\mathbf{g}) = \{\text{iter}((\lambda x.\pi_2(0, x))(\text{eps}(0)), \text{atom}(true_{\mathbb{N}}, \lambda x.0), \pi_2)\} \end{aligned}$$

for every  $w \neq \varepsilon$ . For the Antimirov derivative of  $\mathbf{k}$  w.r.t. 3 we calculate:

$$\begin{aligned} \mathcal{D}_3^{\mathbb{A}}(\mathbf{k}_1) &= \{\text{split}(\mathbf{g}', \mathbf{h}_1, \pi_2), (\lambda x.\pi_2(0, x))((\lambda x.\text{max}(3, x))(\mathbf{f}))\} \\ \mathcal{D}_3^{\mathbb{A}}(\mathbf{k}_2) &= \{\text{split}(\mathbf{g}', \mathbf{h}_2, \pi_2), (\lambda x.\pi_2(0, x))((\lambda x.\text{min}(3, x))(\mathbf{f}))\} \\ \mathcal{D}_3^{\mathbb{A}}(\mathbf{k}) &= \{\text{avg}(\text{split}(\mathbf{g}', \mathbf{h}_1, \pi_2), \text{split}(\mathbf{g}', \mathbf{h}_2, \pi_2)), \\ &\quad \text{avg}(\text{split}(\mathbf{g}', \mathbf{h}_1, \pi_2), (\lambda x.\pi_2(0, x))((\lambda x.\text{min}(3, x))(\mathbf{f}))), \\ &\quad \text{avg}((\lambda x.\pi_2(0, x))((\lambda x.\text{max}(3, x))(\mathbf{f})), \text{split}(\mathbf{g}', \mathbf{h}_2, \pi_2)), \\ &\quad \text{avg}((\lambda x.\pi_2(0, x))((\lambda x.\text{max}(3, x))(\mathbf{f})), \\ &\quad (\lambda x.\pi_2(0, x))((\lambda x.\text{min}(3, x))(\mathbf{f})))\} \end{aligned}$$

and then the Antimirov derivative  $\mathcal{D}_{35}^{\mathbb{A}}(\mathbf{k})$  contains  $3 \cdot 3 = 9$  queries. For the hierarchical derivative of  $\mathbf{k}$  w.r.t. 3 we calculate:

$$\begin{aligned} \mathcal{D}_3^{\mathbb{H}}(\mathbf{k}_1) &= \{\text{split}(\mathbf{g}'', \mathbf{h}_1, \pi_2), (\lambda x.\pi_2(0, x))((\lambda x.\text{max}(3, x))(\mathbf{f}))\} \\ \mathcal{D}_3^{\mathbb{H}}(\mathbf{k}_2) &= \{\text{split}(\mathbf{g}'', \mathbf{h}_2, \pi_2), (\lambda x.\pi_2(0, x))((\lambda x.\text{min}(3, x))(\mathbf{f}))\} \\ \mathcal{D}_3^{\mathbb{A}}(\mathbf{k}) &= \{\text{avg}(\mathcal{D}_3^{\mathbb{H}}(\mathbf{k}_1), \mathcal{D}_3^{\mathbb{H}}(\mathbf{k}_2))\} \end{aligned}$$

In the hierarchical derivative  $\mathcal{D}_{35}^{\mathbb{H}}(\mathbf{k})$ , on the other hand, the subexpressions  $\mathcal{D}_{35}^{\mathbb{H}}(\mathbf{k}_1)$  and  $\mathcal{D}_{35}^{\mathbb{H}}(\mathbf{k}_2)$  contain a total of  $3 + 3 = 6$  queries. This example illustrates the quadratic blowup for Antimirov derivatives, which is avoided using hierarchical derivatives.

## 6 Conclusion

This paper introduces syntactic derivatives for the Quantitative Regular Expressions (QREs) of [5]. The most natural generalization of the classical Brzozowski derivative to QREs is appropriate for the so-called multiset semantics of QREs, which records the possibility of several (finitely many) outputs. Since QREs are meant to describe well-defined functions on streams, we consider a projection of the multiset semantics into the so-called *unambiguous semantics*. Using a representation of derivatives that is inspired from Antimirov’s variant of classical derivatives, we obtain an evaluation algorithm for QREs with streaming space and time complexity that is constant in the stream and exponential in the query. We then restrict attention to the strongly-typed QREs, also considered in [5], which admit more efficient streaming evaluation. We devise a novel representation of derivatives on QREs, which we call *hierarchical derivatives*, and we obtain an evaluation algorithm that streaming space and time complexity that is polynomial in the query. This matches the bounds of [5, 21], but the algorithm presented here is much easier to describe, prove correct, and analyze.

The treatment of QRE evaluation using derivatives is a significant step towards developing automata models for QREs that play the same role as NFAs do for plain regular expressions. The definition of the space-efficient hierarchical derivatives of Sect. 5 suggests that the *parallel evaluation* of  $\mathbf{f}$ ,  $\mathbf{g}$  in subqueries of the form  $op(\mathbf{f}, \mathbf{g})$  and some form of *hierarchical nesting* are essential features of a model that can support efficient evaluation of QREs. A hierarchical automaton model for the streaming computation of quantitative queries is described in [6], but its precise relationship to the QREs of [5] remains to be clarified. Finding the appropriate model of automata for QREs is an important direction for future work, since it would also open the door for query optimization by applying equivalence preserving transformations on the automata.

## References

1. Abadi, D.J., Carney, D., Cetintemel, U., Cherniack, M., Convey, C., Lee, S., Stonebraker, M., Tatbul, N., Zdonik, S.: Aurora: a new model and architecture for data stream management. *VLDB J.* **12**(2), 120–139 (2003)
2. Alon, N., Matias, Y., Szegedy, M.: The space complexity of approximating the frequency moments. *J. Comput. Syst. Sci.* **58**(1), 137–147 (1999)
3. Alur, R., D’Antoni, L.: Streaming tree transducers. In: Czumaj, A., Mehlhorn, K., Pitts, A., Wattenhofer, R. (eds.) *ICALP 2012*. LNCS, vol. 7392, pp. 42–53. Springer, Heidelberg (2012). doi:[10.1007/978-3-642-31585-5\\_8](https://doi.org/10.1007/978-3-642-31585-5_8)
4. Alur, R., D’Antoni, L., Deshmukh, J., Raghothaman, M., Yuan, Y.: Regular functions and cost register automata. In: *Proceedings of the 28th Annual ACM/IEEE Symposium on Logic in Computer Science (LICS 2013)*, pp. 13–22 (2013)
5. Alur, R., Fisman, D., Raghothaman, M.: Regular programming for quantitative properties of data streams. In: Thiemann, P. (ed.) *ESOP 2016*. LNCS, vol. 9632, pp. 15–40. Springer, Heidelberg (2016). doi:[10.1007/978-3-662-49498-1\\_2](https://doi.org/10.1007/978-3-662-49498-1_2)
6. Alur, R., Mamouras, K., Stanford, C.: Automata-based stream processing. In: *ICALP 2017* (2017, to appear)

7. Antimirov, V.: Partial derivatives of regular expressions and finite automaton constructions. *Theoret. Comput. Sci.* **155**(2), 291–319 (1996)
8. Arasu, A., Babu, S., Widom, J.: The CQL continuous query language: Semantic foundations and query execution. *VLDB J.* **15**(2), 121–142 (2006)
9. Babu, S., Widom, J.: Continuous queries over data streams. *ACM Sigmod Rec.* **30**(3), 109–120 (2001)
10. Bonchi, F., Bonsangue, M.M., Hansen, H.H., Panangaden, P., Rutten, J., Silva, A.: Algebra-coalgebra duality in Brzozowski’s minimization algorithm. *ACM Trans. Comput. Log. (TOCL)* **15**(1), 3:1–3:29 (2014)
11. Brzozowski, J.A.: Derivatives of regular expressions. *J. ACM* **11**(4), 481–494 (1964)
12. Colcombet, T.: Forms of determinism for automata (invited talk). In: *Proceedings of the 29th International Symposium on Theoretical Aspects of Computer Science (STACS 2012)*. Leibniz International Proceedings in Informatics (LIPIcs), vol. 14, pp. 1–23 (2012)
13. Colcombet, T.: Unambiguity in automata theory. In: Shallit, J., Okhotin, A. (eds.) *DCFS 2015*. LNCS, vol. 9118, pp. 3–18. Springer, Cham (2015). doi:[10.1007/978-3-319-19225-3\\_1](https://doi.org/10.1007/978-3-319-19225-3_1)
14. Courcelle, B.: Monadic second-order definable graph transductions: a survey. *Theoret. Comput. Sci.* **126**(1), 53–75 (1994)
15. Datar, M., Gionis, A., Indyk, P., Motwani, R.: Maintaining stream statistics over sliding windows. *SIAM J. Comput.* **31**(6), 1794–1813 (2002)
16. Droste, M., Kuich, W., Vogler, H. (eds.): *Handbook of Weighted Automata*. Springer, Heidelberg (2009)
17. Flajolet, P., Martin, G.N.: Probabilistic counting algorithms for data base applications. *J. Comput. Syst. Sci.* **31**(2), 182–209 (1985)
18. Greenwald, M., Khanna, S.: Quantiles and equidepth histograms over streams. In: Garofalakis, M., Gehrke, J., Rastogi, R. (eds.) *Data Stream Management: Processing High-Speed Data Streams*. Data-Centric Systems and Applications, pp. 45–86. Springer, Heidelberg (2016). doi:[10.1007/978-3-540-28608-0\\_3](https://doi.org/10.1007/978-3-540-28608-0_3)
19. Hopcroft, J.E., Motwani, R., Ullman, J.D.: *Introduction to Automata Theory, Languages, and Computation*, 3rd edn. Pearson Education, Upper Saddle River (2006)
20. Lombardy, S., Sakarovitch, J.: Derivatives of rational expressions with multiplicity. *Theoret. Comput. Sci.* **332**(1), 141–177 (2005)
21. Mamouras, K., Raghathan, M., Alur, R., Ives, Z.G., Khanna, S.: StreamQRE: modular specification and efficient evaluation of quantitative queries over streaming data. In: *PLDI 2017* (2017, to appear)
22. Schützenberger, M.P.: On the definition of a family of automata. *Inf. Control* **4**(2), 245–270 (1961)