



Proceedings of the
11th International Workshop on Graph Transformation and
Visual Modeling Techniques
(GTVMT 2012)

Incremental Pattern Matching for Regular Expressions

Arash Jalali, Amir Hossein Ghamarian, Arend Rensink

12 pages

Incremental Pattern Matching for Regular Expressions

Arash Jalali², Amir Hossein Ghamarian¹, Arend Rensink³

¹ a.h.ghamarian@cs.utwente.nl

³ rensink@cs.utwente.nl

Department of Computer Science,
University of Twente, The Netherlands

² arash@netstairs.com

NetStairs.com, Inc.

Abstract: Graph pattern matching lies at the heart of any graph transformation-based system. Incremental pattern matching is one approach proposed for reducing the overall cost of pattern matching over successive transformations by preserving the matches that stay relevant after a rule application. An important issue in any matching scheme, is the ability to properly and consistently deal with various facilities that add to the expressiveness of a GT-tool's rule language. One such feature is the support for regular path expressions, which would let two nodes to be considered as a “match”, if a certain path of edges exists between them. In this paper, the incorporation of regular expression support into incremental pattern matching is discussed within the context of the GROOVE tool set. This includes laying down a formal foundation for incremental pattern matching for regular expressions which is then used to justify the extension proposed to add regular expression support to a well-known pattern matching algorithm.

Keywords: Incremental Matching, Path Matching, Regular Expressions, RETE Algorithm, State Space Exploration

1 Introduction

Tool performance is essential to the success of any specification formalism. This is certainly also true for graph transformation, which has an innate disadvantage that it crucially relies on the computationally expensive problem of pattern matching.

Previous work of Bergmann et al [BÖR⁺08] has shown that the principle of *incremental matching*, in the form of the RETE algorithm [BGT91, For82], can outperform plan-based searching (as advocated in, e.g., [GBG⁺06, VVF06]). The fundamental observation is that graph transformation is all about the gradual evolution of graphs, and therefore search results found for a given graph mostly continue to be valid for any graph obtained by the application of a single rule. We have ourselves reported that the same is true for the exhaustive state space exploration in GROOVE; see [GJR10].

In this paper we extend the scope of the RETE algorithm to matching regular expressions. These were already included in PROGRES [SWZ99] (where they are called path expressions) and are also offered by GROOVE. A regular expression can be used on a rule edge to indicate that the

host graph is expected to have a path between the image of the rule edge's source to the image of the rule edge's target, such that the path labels, when concatenated, form a word in the language of the regular expression. Typically, one is only interested in the *existence* (or *absence*) of such a path, and not in the precise path itself. A limited version of the problem of dynamic transitive closure, where one is merely interested in answering queries of the form “*is node y reachable from node x?*” in the face of updates to the graph, has already been dealt with in the past (see [DI06] for an exhaustive survey).

The existing RETE algorithm for incremental matching is not geared towards checking regular expressions. The main hurdle is the fact that the match of a regular expression does not have a fixed size; in particular, transitively closed expressions allow paths of arbitrary length. This does not sit well with the basic principle of RETE, which relies on the propagation of a match through a static, non-cyclic network of subgraph checker nodes. For regular expressions this network needs to be cyclic; this could potentially make the propagation non-terminating.

Research question and approach. The questions we set out to answer in this paper are: (i) can RETE be extended to cope with regular expression matching, and (ii) does the resulting algorithm continue to outperform search plan-based matching?

In order to answer the first question, we will extend the concept of a RETE network with path checker nodes; in contrast to the existing subgraph checkers, path checkers may be connected in a cycle. We will then give a cut-off criterion to stop the propagation and guarantee termination. The main proof obligation is to show that the set of paths found at the cut-off moment is large enough to include at least one representative path between every pair of connected nodes.

However, it is not enough to have a terminating algorithm: it should also perform well in practice. This means that the cut-off criterion should be as sharp as possible, so that the propagation is fast enough to beat ordinary regular expression matching. In order to evaluate this, we have added regular expression matching to the RETE implementation in GROOVE and compared the performance with the pre-existing, search plan-based matching algorithm.

The results show that incremental matching continues to pay off in the presence of regular expressions, but that the actual gain depends on the dynamics of the rules. In particular, if the structure that is tested by regular expression matching changes frequently, then incremental matching loses out. This is indeed not surprising, as match propagation for path checker nodes is a potentially complex operation which, when repeated often, does not scale as well as the time-honoured word checking algorithm for regular automata.

2 Definitions

We first define some basic concepts: graphs, paths and match morphisms.

Graphs consist of nodes and labelled edges. We assume the existence of global universes *Vertex* (ranged over by v), *Edge* (ranged over by e), and *Label* of vertex (or node), edge and label entities. *Label* is further partitioned into the universe *Atom* of atomic labels (ranged over by a) and *RegExp* of *regular expressions* (ranged over by R). There are also global functions $src, tgt : Edge \rightarrow Vertex$ and $lab : Edge \rightarrow Label$ associating with every edge a source and target node and label.

Definition 1 (graph) A graph G is a tuple $\langle V_G, E_G \rangle$ in which $V_G \subseteq \text{Vertex}$ is the set of vertices (or nodes) and $E_G \subseteq \text{Edge}$ is the set of edges, such that $\text{src}(E_G) \subseteq V_G$ and $\text{tgt}(E_G) \subseteq V_G$.

We distinguish *host graphs*, in which only Atom is used as label universe, and *rule graphs*, in which also RegExp labels may occur. Regular expressions are defined by the following grammar:

$$R ::= \lambda \mid a \mid \neg R \mid R_1 | R_2 \mid R_1 \cdot R_2 \mid R^*$$

where λ stands for the empty word and $a \in \text{Atom}$. The only non-standard operation is $\neg R$. In graphs, this will be matched by an inversed path; to capture this formally, we introduce the *complementary actions* $\overline{\text{Atom}}$, disjoint from Label . For arbitrary $a \in \text{Atom} \cup \overline{\text{Atom}}$ we use \bar{a} to denote the complement of a ; this operator is its own inverse, hence $\bar{\bar{a}} = a$. This is extended to words $w \in \text{Word}$ where $\text{Word} = (\text{Atom} \cup \overline{\text{Atom}})^*$ by also reversing them: $\bar{\varepsilon} = \varepsilon$ and $\overline{a\bar{w}} = \bar{w}a$.

The language of a regular expression R is then inductively defined by $L_R \subseteq \text{Word}$ such that:

$$\begin{aligned} L_\lambda &= \{\varepsilon\} \\ L_a &= \{a\} \\ L_{\neg R} &= \{\bar{w} \mid w \in L_R\} \\ L_{R_1 \cdot R_2} &= \{w_1 w_2 \mid w_1 \in L_{R_1}, w_2 \in L_{R_2}\} \\ L_{R_1 | R_2} &= L_{R_1} \cup L_{R_2} \\ L_{R^*} &= \{w_1 \cdots w_n \mid n \geq 0, w_1, \dots, w_n \in L_R\} \end{aligned}$$

To define the meaning of regular expressions over graphs, we need the concept of *paths* in a host graph. For this purpose, we first introduce *inverse edges*: \bar{E}_G for a graph G will denote the set of inverse edges, such that $\text{src}(\bar{e}) = \text{tgt}(e)$, $\text{tgt}(\bar{e}) = \text{src}(e)$ and $\text{lab}(\bar{e}) = \overline{\text{lab}(e)}$.

Definition 2 (path) A *path* is an alternating sequence $p = v_1 e_2 v_2 \cdots e_n v_n$ of vertices and edges, starting and ending on a vertex, such that $v_{i-1} = \text{src}(e_i)$ and $\text{tgt}(e_i) = v_i$ for $1 < i \leq n$.

We use Π to denote the set of all paths and $\Pi(G)$ for the subset of all paths in G , and extend the source and target functions to paths by $\text{src}(p) = v_1$ and $\text{tgt}(p) = v_n$. Path concatenation of $p = v_1 e_2 v_2 \cdots e_n v_n$ and $p' = v'_1 e'_2 v'_2 \cdots e'_n v'_n$ is defined by

$$p \cdot p' = v_1 e_2 v_2 \cdots e_n v_n e'_2 v'_2 \cdots e'_n v'_n \quad \text{if } v_n = v'_1$$

Two further relevant functions on paths are *word* : $\Pi \rightarrow \text{Word}$ that maps each path in G to the corresponding sequence of labels, and path reversal $\text{rev} : \Pi \rightarrow \Pi$, inductively defined by:

$$\begin{aligned} \text{word}(v) &= \varepsilon & \text{rev}(v) &= v \\ \text{word}(p \cdot e) &= \text{word}(p) \text{lab}(e) & \text{rev}(p \cdot e) &= v \bar{e} \text{rev}(p) \end{aligned}$$

We can now formally define the relation between regular expressions and graphs. For an arbitrary regular expression R and graph G , we define $\Pi^R(G)$ as the set of paths in G satisfying R , thus:

$$\Pi^R(G) = \{p \in \Pi(G) \mid \text{word}(p) \in L_R\} . \quad (1)$$

Note that $\Pi^R(G)$ is in general infinite, due to the possible presence of loops in G .

Finally, a *match morphism* is a function from a rule graph L to a host graph G that provides images in G for the atom-labelled edges of L and guarantees the existence of paths in G for the regular expression-labelled edges of L :

Definition 3 (match morphism) Given a rule graph L and a host graph G , a match morphism of L into G is a partial function $f : (V_L \cup E_L) \rightarrow (V_G \cup E_G)$ that is total on V_L such that $f(V_L) \subseteq V_G$ and $f(E_L) \subseteq E_G$, satisfying for all $e \in E_L$:

- if $lab(e) = a \in \text{Atom}$, then $f(e) = e' \in E_G$ such that $src(e') = f(src(e))$, $tgt(e') = f(tgt(e))$ and $lab(e') = a$;
- if $lab(e) = R \in \text{RegExp}$, then there is a $p \in \Pi^R(G)$ such that $src(p) = f(src(e))$ and $tgt(p) = f(tgt(e))$.

For the purpose of our discussions, we shall only be concerned with the second condition in the above definition. We use F^L to denote the set of all match morphisms of L , and $F^L(G)$ for the subset of match morphisms into G . It is important to note that the paths themselves are not part of the match morphism images. Among other things, this implies that $F^L(G)$ is always finite.

It should also be noted that, in this paper, we deal with positive conditions only. Negative application conditions can be supported using the same techniques as described in [GJR10].

3 Incremental Matching

The RETE algorithm for incremental matching uses the notion of *subgraph checkers* that keep a record of their match morphisms into a host graph, and are *updated* upon every graph change. Graph changes are additions and deletions of nodes and edges, using the following (partially defined) operations:

$$\begin{aligned}
 G + v &= \langle V \cup \{v\}, E \rangle && \text{for all } v \notin V_G \\
 G + e &= \langle V, E \cup \{e\} \rangle && \text{for all } e \notin E_G \text{ with } src(e), tgt(e) \in V_G \\
 G - v &= \langle V \setminus \{v\}, E \rangle && \text{for all } v \in V_G \text{ such that } \nexists e \in E_G : v \in \{src(e), tgt(e)\} \\
 G - e &= \langle V, E \setminus \{e\} \rangle && \text{for all } e \in E_G
 \end{aligned}$$

The essence of incremental matching is the existence of “cheap” updates $\delta_{+x}^L(G), \delta_{-x}^L(G) \subseteq F^L$ for every rule graph L host graph G and graph element $x \in \text{Vertex} \cup \text{Edge}$, which capture the *change* to $F^L(G)$ when a vertex or edge is added or removed:

$$\begin{aligned}
 F^L(G + v) &= F^L(G) \cup \delta_{+v}^L(G) && F^L(G - v) = F^L(G) \setminus \delta_{-v}^L(G) \\
 F^L(G + e) &= F^L(G) \cup \delta_{+e}^L(G) && F^L(G - e) = F^L(G) \setminus \delta_{-e}^L(G)
 \end{aligned}$$

whenever the graph operations are defined. In RETE, this is achieved by creating a directed acyclic graph $\langle N, \sqsubset \rangle$ (called a network) of *checker nodes* $n \in N$, each of which is associated with a subgraph $S_n \subseteq L$ such that $n_1 \sqsubset n_2$ implies $S_{n_1} \subset S_{n_2}$. For the \sqsubset -minimal nodes, the S_n are single nodes and edges of L , whereas $S_n = L$ when n is the \sqsubset -maximum.

Each checker node n continually maintains $F^{S_n}(G)$, and computes its own δ 's upon every graph update, propagating from small to large in the \sqsubset -ordering. For instance, for a rule graph $L = \langle \{v'_0, v'_1\}, \{e'\} \rangle$ consisting only of a single edge e' with $src(e') = v'_0$ and $tgt(e') = v'_1$, the δ 's are given by

$$\delta_{+v}^L(G) = \delta_{-v}^L(G) = \emptyset$$

$$\delta_{+e}^L(G) = \delta_{-e}^L(G) = \begin{cases} \{(v'_0 \mapsto v_0, v'_1 \mapsto v_1, e' \mapsto e)\} & \text{if } lab(e) = lab(e') \\ \emptyset & \text{otherwise} \end{cases}$$

Only non-empty δ are propagated to the \sqsubset -successors in the network, reducing the amount of work required to update the entire structure.

3.1 Path matches

Since this paper deals with the extension of RETE to regular expressions, we will not go further into the details of the standard algorithm. The extension involves the definition of new checker nodes for paths and extending the RETE network with those nodes. There are two important novelties:

- The path checker nodes do not have to maintain all paths, but only have to test for the existence of a path (see Def. 3).
- The \sqsubset -ordering over path checkers is no longer acyclic; instead, we will have $n \sqsubset n$ whenever R_n (the regular expression associated with n) is of the form R^* .

The main issue is to decide what kind of information every path checker node n should maintain, given that it is only required to know whether a satisfying path exists between given (arbitrary) $v_1, v_2 \in V_G$. Important criteria for this decision are: (i) the amount of information should be kept as small as possible and (ii) incremental updates should be cheap. To meet these criteria, we introduce an *algebra of path matches*:

$$m ::= \langle v_0 \rangle \mid \langle e \rangle \mid inv(m) \mid dot(m_1, m_2) \mid close(m_1, m_2)$$

Here, $v_0 \in \text{Vertex}$ and $e \in \text{Edge}$; the other operators are syntactic constructors that give rise to a tree-like structure of path matches. The close constructor, for instance, represents a match for a closure regular expression. The universe of all path matches will be denoted M . The following table defines functions $start, end : M \rightarrow \text{Vertex}$ that retrieve the start and end nodes of a path match, and a function $int : M \rightarrow 2^{\text{Vertex}}$ that retrieves the set of internal nodes:

m	$start(m)$	$end(m)$	$int(m)$
$\langle v_0 \rangle$	v_0	v_0	\emptyset
$\langle e \rangle$	$src(e)$	$tgt(e)$	\emptyset
$inv(m_1)$	$end(m_1)$	$start(m_1)$	\emptyset
$dot(m_1, m_2)$	$start(m_1)$	$end(m_2)$	\emptyset
$close(m_1, m_2)$	$start(m_1)$	$end(m_2)$	$\{start(m_1)\} \cup int(m_2)$

m is said to be a path match *in* G if it only uses nodes and edges from G . We call a path match m *consistent* if

- $m = \text{inv}(m_1)$ implies m_1 is consistent;
- $m = \text{dot}(m_1, m_2)$ implies (i) m_1 and m_2 are consistent and (ii) $\text{end}(m_1) = \text{start}(m_2)$;
- $m = \text{close}(m_1, m_2)$ implies (i) m_1 and m_2 are consistent, (ii) $\text{end}(m_1) = \text{start}(m_2)$, (iii) $\text{start}(m_1) \notin \text{int}(m_2)$ and (iv) $m_2 = \langle v_0 \rangle$ or $m_2 = \text{close}(m_{21}, m_{22})$.

Note in particular condition (iv) of `close`, which states that the second operand is always either an empty path match or itself of the form `close`.

Obviously, every path match is essentially a path with additional structure. Indeed, we can easily retrieve the path from the path match:

$$\begin{aligned} \text{path}(\langle v_0 \rangle) &= v_0 \\ \text{path}(\langle e \rangle) &= \text{src}(e) \text{ } e \text{ } \text{tgt}(e) \\ \text{path}(\text{inv}(m)) &= \text{rev}(\text{path}(m)) \\ \text{path}(\text{dot}(m_1, m_2)) &= \text{path}(m_1) \cdot \text{path}(m_2) \\ \text{path}(\text{close}(m_1, m_2)) &= \text{path}(m_1) \cdot \text{path}(m_2) \end{aligned}$$

We now recursively define the path matches of arbitrary regular expressions in a given graph G , as the smallest sets satisfying the following equations:

$$\begin{aligned} M^\lambda(G) &= \{\langle v \rangle \mid v \in V_G\} \\ M^a(G) &= \{\langle e \rangle \mid e \in E_G, \text{lab}(e) = a\} \\ M^{-R}(G) &= \{\text{inv}(m) \mid m \in M^R(G)\} \\ M^{R_1|R_2}(G) &= M^{R_1}(G) \cup M^{R_2}(G) \\ M^{R_1 \cdot R_2}(G) &= \{m = \text{dot}(m_1, m_2) \mid m_1 \in M^{R_1}(G), m_2 \in M^{R_2}(G), m \text{ is consistent}\} \\ M^{R^*}(G) &= M^\lambda(G) \cup \{m = \text{close}(m_1, m_2) \mid m_1 \in M^R(G), m_2 \in M^{R^*}(G), m \text{ is consistent}\} \end{aligned}$$

In particular, the path matches of the Kleene star are defined recursively. The following property is the core of the correctness proof for our incremental match algorithm.

Theorem 1 *Let G be a host graph and R a regular expression.*

1. $M^R(G)$ is finite;
2. For every $m \in M^R(G)$, $\text{path}(m) \in \Pi^R(G)$;
3. For every $p \in \Pi^R(G)$, there is an $m \in M^R(G)$ such that $\text{start}(m) = \text{src}(p)$ and $\text{end}(m) = \text{tgt}(p)$.

Clause 1 ensures that $M^R(G)$ can be effectively computed. It follows from the fact that, due to the consistency conditions on path matches, $\text{int}(\text{close}(m_1, m_2)) \supset \text{int}(m_2)$, combined with the fact that $\text{int}(m) \subseteq V_G$ for all $m \in M^R(G)$. Clause 2 guarantees that the path matches are sound: they only generate paths that satisfy the regular expression. Finally, Clause 3, called the *sufficient coverage* clause, states that the path matches are complete: they contain at least one representative element for every pair of host nodes between which there exists a path.

This provides us with an answer to the question above: the RETE path checker node for R will maintain the set $M^R(G)$. In the remainder of this section we will show that this can be updated incrementally.

3.2 Incremental Computation of Path Matches

Given the definition of the path matches, defining the δ -functions, which now range over the terms of our path match algebra, is relatively straightforward. For the base regular expressions they are given by

$$\begin{aligned} \delta_{+v}^\lambda(G) &= \delta_{-v}^\lambda(G) = \{\langle v \rangle\} & \delta_{+e}^\lambda(G) &= \delta_{-e}^\lambda(G) = \emptyset \\ \delta_{+v}^a(G) &= \delta_{-v}^a(G) = \emptyset & \delta_{+e}^a(G) &= \delta_{-e}^a(G) = \begin{cases} \{\langle e \rangle\} & \text{if } \text{lab}(e) = a \\ \emptyset & \text{otherwise} \end{cases} \end{aligned}$$

For the composed regular expressions we can ignore the distinction between node and edge updates: for $x \in \text{Vertex} \cup \text{Edge}$

$$\begin{aligned} \delta_{+x}^{-R}(G) &= \{\text{inv}(m) \mid m \in \delta_{+x}^R(G)\} \\ \delta_{+x}^{R_1|R_2}(G) &= \delta_{+x}^{R_1}(G) \cup \delta_{+x}^{R_2}(G) \\ \delta_{+x}^{R_1 \cdot R_2}(G) &= \{m = \text{dot}(m_1, m_2) \mid m_1 \in \delta_{+x}^{R_1}(G), m_2 \in M^{R_2}(G+x), m \text{ is consistent}\} \\ &\quad \cup \{m = \text{dot}(m_1, m_2) \mid m_1 \in M^{R_1}(G+x), m_2 \in \delta_{+x}^{R_2}(G), m \text{ is consistent}\} \\ \delta_{+x}^{R^*,0}(G) &= \delta_{+x}^\lambda(G) \cup \{m = \text{close}(m_1, m_2) \mid m_1 \in \delta_{+x}^R(G), m_2 \in M^{R^*}(G), m \text{ is consistent}\} \\ \delta_{+x}^{R^*,i+1}(G) &= \{m = \text{close}(m_1, m_2) \mid m_1 \in M^R(G+x), m_2 \in \delta_{+x}^{R^*,i}(G), m \text{ is consistent}\} \\ \delta_{+x}^{R^*}(G) &= \bigcup_{i=0}^\omega \delta_{+x}^{R^*,i}(G) \end{aligned}$$

The most interesting is of course the case for R^* . This is approximated by a sequence of steps, which is defined as unbounded but in practice will always result in $\delta_{+x}^{R^*,i}(G) = \emptyset$ from a certain value of i onward — the argument is the same as for the finiteness of $M^{R^*}(G)$.

The removal updates for the composed operators only have to look at existing matches; in other words, we do not have to compute them:

$$\begin{aligned} \delta_{-x}^{-R}(G) &= \{\text{inv}(m) \in M^R(G) \mid m \in \delta_{-x}^R(G)\} \\ \delta_{-x}^{R_1|R_2}(G) &= \delta_{-x}^{R_1}(G) \cup \delta_{-x}^{R_2}(G) \\ \delta_{-x}^{R_1 \cdot R_2}(G) &= \{\text{dot}(m_1, m_2) \in M^R(G) \mid m_1 \in \delta_{-x}^{R_1}(G) \text{ or } m_2 \in \delta_{-x}^{R_2}(G)\} \\ \delta_{-x}^{R^*,0}(G) &= \delta_{-x}^\lambda(G) \cup \{\text{close}(m_1, m_2) \in M^R(G) \mid m_1 \in \delta_{-x}^R(G)\} \\ \delta_{-x}^{R^*,i+1}(G) &= \{m = \text{close}(m_1, m_2) \mid m_2 \in \delta_{-x}^{R^*,i}(G)\} \\ \delta_{-x}^{R^*}(G) &= \bigcup_{i=0}^\omega \delta_{-x}^{R^*,i}(G) \end{aligned}$$

In an actual implementation, removal can be implemented efficiently by keeping a mapping from each match to the set of matches derived from it; i.e., from m to all existing matches of the form $\text{inv}(m)$, $\text{dot}(m, m')$, $\text{dot}(m', m)$, $\text{close}(m', m)$ and $\text{close}(m, m')$. This obviates all further lookups.

The correctness criterion of incremental updating is captured by the following theorem:

Theorem 2 For every graph G , regular expression R and element $x \in \text{Vertex} \cup \text{Edge}$,

$$\begin{aligned} M^R(G+x) &= M^R(G) \cup \delta_{+x}^R(G) \\ M^R(G-x) &= M^R(G) \setminus \delta_{-x}^R(G) . \end{aligned}$$

In order to show this theorem for the choice operator, we need an alternative characterisation of the set δ_{-x}^R , namely that it always removes precisely those match objects that contain x . For every $x \in \text{Vertex} \cup \text{Edge}$ define

$$M_x = \{m \in M \mid x \text{ occurs in } m\} .$$

Now we can easily show the following property, by induction on the structure of R :

Proposition 1 For every graph G , regular expression R and element $x \in \text{Vertex} \cup \text{Edge}$,

$$\delta_{-x}^R(G) = M^R(G) \cap M_x .$$

With the use of this property we can prove the main theorem above.

Proof of Theorem 2. By induction on the structure of R . The only really interesting cases are those of $M^{R_1|R_2}(G-x)$ and $M^{R^*}(G+x)$.

Case $M^{R_1|R_2}(G-x) = M^{R_1|R_2}(G) \setminus \delta_{-x}^{R_1|R_2}(G)$. This follows from

$$\begin{aligned} M^{R_1|R_2}(G-x) &= M^{R_1}(G-x) \cup M^{R_2}(G-x) && \text{(by definition)} \\ &= (M^{R_1}(G) \setminus \delta_{-x}^{R_1}) \cup (M^{R_2}(G) \setminus \delta_{-x}^{R_2}) && \text{(by the induction hypothesis)} \\ &= (M^{R_1}(G) \setminus M_x) \cup (M^{R_2}(G) \setminus M_x) && \text{(by Prop. 1)} \\ &= (M^{R_1}(G) \cup M^{R_2}(G)) \setminus M_x && \text{(by set algebra)} \\ &= M^{R_1|R_2}(G) \setminus \delta_{-x}^{R_1|R_2}(G) && \text{(by definition and Prop. 1).} \end{aligned}$$

Case $M^{R^*}(G+x) = M^{R^*}(G) \cup \delta_{+x}^{R^*}(G)$. The inclusion from right to left is immediate. Now assume $m \in M^{R^*}(G+x)$. We proceed by induction on $|m|$, where $|m| = 0$ if m is not of the form close , and $|m| = |m_2| + 1$ if $m = \text{close}(m_1, m_2)$. The hypothesis is that either $m \in M^{R^*}(G)$ or $m \in \delta_{+x}^{R^*,i}(G)$ for some $i \leq |m|$.

Base case: $|m| = 0$ It follows that $m \in M^\lambda(G+x)$. Then $m \in M^\lambda(G)$ or $m \in \delta_{+x}^\lambda(G)$ by the outer induction hypothesis; in the first case $m \in M^{R^*}(G)$ and in the second $m \in \delta_{+x}^{R^*,0}(G)$. Since $0 \leq |m|$ we are done.

Induction step: $|m| = n + 1$. It follows that $m = \text{close}(m_1, m_2)$ with $|m_2| = n$; hence $m_1 \in M^R(G+x)$ and $m_2 \in M^{R^*}(G+x)$. By the outer induction hypothesis, either (i) $m_1 \in M^R(G)$ or (ii) $m_1 \in \delta_{+x}^R(G)$, and by the inner induction hypothesis, either (iii) $m_2 \in M^{R^*}(G)$ or (iv) $m_2 \in \delta_{+x}^{R^*,i}(G)$ for $i \leq |m_2|$. This gives rise to the following combinations of cases:

- (i) and (iii): Then $m \in M^{R^*}(G)$, hence we are done.
- (ii) and (iii): Then $m \in \delta_{+x}^{R^*,0}(G)$; since $0 \leq |m|$ we are done.
- (iv): Then $m \in \delta_{+x}^{R^*,i+1}(G)$; since $i + 1 \leq |m|$ we are done.

□

4 Implementation

In this section we briefly discuss how the extended RETE algorithm has been implemented in GROOVE. Our discussion will be exclusive to regular expressions. For a more general explanation of RETE and its implementation in GROOVE, the reader is referred to [GJR10]. Our explanation will mainly focus on the static build of the RETE network, as the dynamic behavior of the network follows directly from the formal semantics of the δ operations.

4.1 Path-checkers

The RETE algorithm relies heavily on the network of checker nodes that are responsible for finding and passing along partial matches for one or more rules. The network in the classic RETE algorithm comes equipped with several types of checker nodes, including *edge checkers*, *node checkers*, *subgraph checkers*, as well as *condition/production rule checkers*. In order to add the functionality of finding matches for paths, as opposed to partial graphs, it is natural to think that one needs a new type of checker node that finds and passes along matches for a specific kind of path - hence the addition of *path checkers*.

A path checker node, in general, receives one or two partial path matches, and then tries to combine and/or transform them based on the semantics of the regular expression operator they represent. For instance, a *sequence path checker*, receives two path matches, and tries to see if they can be combined into a larger path by sequentially concatenating them.

4.2 The RETE Static Build

Given a rule edge with a regular expression label R, we use the syntax of R as a guide to build the RETE network, such that it would be able to incrementally find path matches. Figure 1 shows

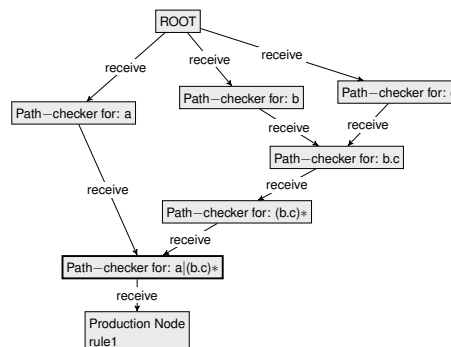


Figure 1: RETE network for $a|(b.c)^*$

#	Rules	Len	Pool	States / Transitions	RETE(ms)	SP(ms)	Prf. Gain(%)
1	move	27	-	2941 / 8042	1781	1984	10.2
2	move	54	-	19063 / 54439	18734	40437	53.7
3	move,mark	12	-	2048 / 8192	1016	2531	59.9
4	move,mark	15	-	16384 / 77824	7297	21758	66.5
5	move,add,del	2	5	23248 / 42615	6062	5750	-5.2
6	move,add,del	4	5	357605 / 634772	89719	87000	-3.0
7	all	2	5	179480 / 764760	90125	125375	28.1

Table 1: Experimental results for exhaustive state space exploration

chain (rows 1 to 4), the path matches found for next+ by RETE remain valid throughout the exploration, while SP is forced to follow and match the chain repeatedly after each rule application. When the rules ‘del’ and ‘add’ are present, RETE starts losing a lot of matches, or ends up creating more matches every time a next edge is deleted or created by those rules. This leads to RETE losing the race to SP (rows 5 and 6). This performance loss is easily compensated when the rule ‘mark’, which also does not modify the structure of the chain, is used along with the others (row 7), where a 5.2% loss is turned into a 28% performance gain. While the other experiments are meant to reveal the underlying mechanisms affecting performance, row 7 represents the typical and naturally occurring scenario in which RETE can be used to achieve gains. There is however another characteristic that should be noted about the above scenarios: RETE’s loss occurs at a more or less stable rate (rows 5 and 6) even when doubling the length of the chain has resulted in a 11 fold increase in the size of the state space. Rows 1-4 show that this is not the case when SP is losing, as increasing the length of the chain causes further declines in its performance compared to RETE.

6 Conclusion and Future Work

In this paper, we provided an extension to the RETE algorithm to support matching for regular path expressions. Finding paths in a graph that satisfy a given regular expression is particularly problematic as the matches for a regular expression with transitive closures are neither of a fixed size, nor is there always a finite number of them due to the presence of cycles. This gave rise to the question of what is actually meant by finding a match for a regular expression, which we addressed by formally defining paths and their matches. We also provided formal specifications of incremental updates to the finite set of matches found for each regular expression and showed that our formulation of the set of matches for a given regular expression gives rise to a finite and complete set, which makes it suitable for implementation.

Our experimental results showed that RETE can in fact provide performance gains in regular expression matching when the dynamics of the situation are such that RETE’s ‘find once, use many times’ behavior makes it possible to avoid repetitive finding of long, mostly unaltered paths.

As future work, one possible direction for improvement would be to boost RETE’s regular expression matching by enhancing its path checkers, especially the closure path checkers, with an *on-demand* capability, such that they would only produce matches when needed and only



those that are anchored at a certain node in the host graph. In our example grammar discussed in Section 5, the closure path checker does in fact find paths that lie in the middle of the chain, which are basically of little use, since all our operations happen at the end of the chain. This could lead to serious performance declines in finding matches for the expression a^* in cases where the host graph is dense with a -labelled edges and there are frequent updates in the form of addition and removal of such edges.

Bibliography

- [BGT91] H. Bunke, T. Glauser, T.-H. Tran. An Efficient Implementation of Graph Grammars Based on the RETE Matching Algorithm. In *Proceedings of the 4th International Workshop on Graph-Grammars and Their Application to Computer Science*. Pp. 174–189. Springer-Verlag, London, UK, 1991.
- [BÖR⁺08] G. Bergmann, A. Ökrös, I. Ráth, D. Varró, G. Varró. Incremental pattern matching in the VIATRA model transformation system. In *GRaMoT '08: Proceedings of the third international workshop on Graph and model transformations*. Pp. 25–32. ACM, NY, USA, 2008.
- [DI06] C. Demetrescu, G. F. Italiano. Dynamic shortest paths and transitive closure: Algorithmic techniques and data structures. *Journal of Discrete Algorithms* 4(3):353–383, 2006.
- [For82] C. Forgy. RETE, a fast algorithm for the many pattern / many object pattern match problem. *Artificial Intelligence* 19:17–37, 1982.
- [GBG⁺06] R. Geiß, G. V. Batz, D. Grund, S. Hack, A. Szalkowski. GRGEN: A Fast SPO-Based Graph Rewriting Tool. In Corradini et al. (eds.), *International Conference on Graph Transformations (ICGT)*. LNCS 4178, pp. 383–397. Springer, 2006.
- [GJR10] A. H. Ghamarian, A. Jalali, A. Rensink. Incremental Pattern Matching in Graph-Based State Space Exploration. *ECEASST* 32, 2010.
- [SWZ99] A. Schürr, A. J. Winter, A. Zündorf. The PROGRES approach: language and environment. In Ehrig et al. (eds.), *Handbook of graph grammars and computing by graph transformation: applications, languages, and tools*. Volume 2, pp. 487–550. World Scientific Publishing Co., Inc., 1999.
- [VVF06] G. Varró, D. Varró, K. Friedl. Adaptive Graph Pattern Matching for Model Transformations using Model-sensitive Search Plans. In Karsai and Taentzer (eds.), *GraMot 2005, International Workshop on Graph and Model Transformations*. ENTCS 152, p. 191–205. Elsevier, 2006.
http://www.inf.mit.bme.hu/FTSRG/Publications/varro/2005/gramot05_vvf.pdf