# Symbolic Solving of Extended Regular Expression Inequalities
Technical Report

**Matthias Keil and Peter Thiemann**

**Institute for Computer Science**
**University of Freiburg**
**Freiburg, Germany**
`{keilr,thiemann}@informatik.uni-freiburg.de`

─── **Abstract** ───

This paper presents a new solution to the containment problem for *extended regular expressions* that extends *basic regular expressions* with intersection and complement operators and consider regular expressions on *infinite alphabets* based on potentially *infinite character sets*. Standard approaches deciding the containment do not take extended operators or character sets into account. The algorithm avoids the translation to an expression-equivalent automaton and provides a purely symbolic term rewriting systems for solving regular expressions inequalities.

We give a new symbolic decision procedure for the containment problem based on Brzozowski's regular expression derivatives and Antimirov's rewriting approach to check containment. We generalize Brzozowski's syntactic derivative operator to two derivative operators that work with respect to (potentially infinite) representable character sets.

**1998 ACM Subject Classification** F.4.3 Formal Languages

**Keywords and phrases** Extended Regular Expressions, Containment, Infinite Alphabtes, Infinite Character Sets

## 1 Introduction

Regular expressions have many applications in the context of software development and information technology: text processing, program analysis, compiler construction, query processing, and so on. Modern programming languages either come with standard libraries for regular expression processing or they provide built-in facilities (e.g., Perl, Ruby, and JavaScript). Many of these implementations augment the basic regular operations $+$, $\cdot$, and $*$ (union, concatenation, and Kleene star) with enhancements like character classes and wildcard literals, cardinalities, sub-matching, intersection, or complement.

Regular expressions (RE) are advantageous in these domains because they provide a concise means to encode many interesting problems. REs are well suited for verification applications, because there are decision procedures for many problems involving them: the word problem ($w \in [\![r]\!]$), emptiness ($[\![r]\!] = \emptyset$), finiteness, containment ($[\![r]\!] \subseteq [\![s]\!]$), and equivalence ($[\![r]\!] = [\![s]\!]$). Here we let $r$ and $s$ range over RE and write $[\![\cdot]\!]$ for the function that maps a regular expression to the regular language that it denotes. There are also effective constructions for operations like union, intersection, complement, prefixes, suffixes, etc on regular languages.

Recent applications impose new demands on operations involving regular expressions. The Unicode character set with its more than 1.1 million code points requires the ability to deal effectively with very large character sets and hence character classes. Similarly,

formalizing access contracts for objects in scripting languages even requires regular expressions over an infinite alphabet: in this application, the alphabet itself is an infinite formal language (the language of field names) and a "character class" (i.e., a set of field names) is described by a regular expression [12, 8]. Hence, a "character class" may also have infinitely many elements.

We study the containment problem for regular expressions with two enhancements. First, we consider *extended regular expressions* (ERE) that contain intersection and complement operators beyond the standard regular operators of union, concatenation, and Kleene star. An ERE also denotes a regular language but it can be much more concise than a standard RE. Second, we consider EREs on any alphabet that is presented as an effective boolean algebra. This extension encompasses some infinite alphabets like the set of all field names in a scripting language.

The first enhancement is known to be decidable, but we give a new symbolic decision procedure based on Brzozowski's regular expression derivatives [4] and Antimirov's rewriting approach to check containment [1]. The second enhancement has been studied previously [20, 18, 19], but in the context of automata and finite state transducers. It has not been investigated on the level of regular expressions and in particular not in the context of Brzozowski's and Antimirov's work. We give sufficient conditions to ensure applicability of our modification of Brzozowski's and Antimirov's approach to the containment problem while retaining decidability.

## 1.1   Related Work

The practical motivation for considering this extension is drawn from the authors' work on checking access contracts for objects in a scripting language at run time [12]. In that work, an access contract specifies a set of access paths that start from a specific anchor object. An access path is a word over the field names of the objects traversed by the path and we specify a set of such paths by a regular expression on the field names. We claim that such a regular expression draws from an infinite alphabet because a field name in a scripting language is an arbitrary string (of characters). For succinctness, we specify sets of field names using a second level of regular expressions on characters.

In our implementation, checking containment is required to reduce memory consumption. If the same object is restricted by more than one contract, then we apply containment checking to remove redundant contracts. In that previous work, contracts were limited to basic regular expressions and the field-level expressions were limited to disjunctions of literals. Applying the results of the present paper enables us to lift both restrictions.

The standard approach to checking regular expression containment is via translation to finite automata, which may involve an exponential blowup, and then construction of a simulation (or a bisimulation for equivalence) [9]. A related approach based on non-deterministic automata is given by Bonchi and Pous [3].

The exponential blowup is due to the construction of a deterministic automaton from the regular expression. Thompson's construction [17], creates a non-deterministic finite automaton with $\epsilon$-transitions where the number of states and transitions is linear to the length of the (standard) regular expression. Glushkov's [7] and McNaughton and Yamada's [13] position automaton computes a $n + 1$-state non-deterministic automaton with up to $n^2$ transitions from a $n$-symbol expression. They are the first to use the notion of a first symbol. Brzozowski's regular expression derivatives [4] directly calculate a deterministic automaton from an ERE. Antimirov's partial derivative approach [2] computes a $n+1$-state non-deterministic automation, but again without intersection and complement. We are not

aware of an extension of Glushkov's algorithm to extended regular expressions.

Owens and other have implemented an extension of Brzozowski's approach with character classes and wildcards [15].

Antimirov [1] also proposes a symbolic method for solving regular expression inequalities, based on partial derivatives, with exponential worst-case run time. His *containment calculus* is closely related to the simulation technique used by Hopcroft and Karp [9] for proving equivalence of automata. In fact, a decision procedure for containment of regular expressions leads to one for equivalence and vice versa. Ginzburg [6] gives an equivalence procedure based on Brzozowski derivatives. Antimirov's original work does not consider intersection and complement. Caron and coworkers [5] extend Antimirov's work to ERE using antichains, but the resulting procedure is very complex compared to ours.

A shortcoming of all existing approaches is their restriction to finite alphabets. Supporting both makes a significant difference in practice: an iteration over the alphabet $\Sigma$ is feasible for small alphabets, but it is impractical for very large alphabets (e.g., Unicode) or infinite ones (e.g., another level of regular languages as for our contracts). Furthermore, most regular expressions used in practice contain character sets. We apply techniques developed for symbolic finite automata to address these issues [19].

## 1.2 Overview

This paper is organized as follows. In Section 2, we recall notations and concepts used in this paper. Section 3 introduces the notion of an effective boolean algebra for representing sets of symbols abstractly. Section 4 explains Antimirov's algorithm for checking containment, which is the starting point of our work. Next, Section 5 defines two notions of derivatives on regular expressions with respect to symbol sets. It continues to introduce the key notion of *next literals*, which ensures finiteness of our extension to Antimirov's algorithm. Section 6 contains the heart of our extended algorithm, a deduction system that determines containment of extended regular expressions along with a soundness proof.

This paper concludes with an appendix with further technical details, examples, and proofs of theorems.

## 2   Regular Expressions

An *alphabet* $\Sigma$ is a denumerable, potentially infinite set of symbols. $\Sigma^*$ is the set of all finite words over symbols from $\Sigma$ with $\epsilon$ denoting the empty word. Let $a, b, c \in \Sigma$ range over symbols; $u, v, w \in \Sigma^*$ over words; and $A, B, C \subseteq \Sigma$ over sets of symbols.

Let $\mathcal{L}, \mathcal{L}' \subseteq \Sigma^*$ be languages. The *left quotient* of $\mathcal{L}$ by a word $u$, written $u^{-1}\mathcal{L}$, is the language $\{v \mid uv \in \mathcal{L}\}$. It is immediate from the definition that $(au)^{-1}\mathcal{L} = u^{-1}(a^{-1}\mathcal{L})$ and that $u \in \mathcal{L}$ iff $\epsilon \in u^{-1}\mathcal{L}$. Furthermore, $\mathcal{L} \subseteq \mathcal{L}'$ iff $u^{-1}\mathcal{L} \subseteq u^{-1}\mathcal{L}'$ for all words $u \in \Sigma^*$. The left quotient of one language by another is defined by $\mathcal{L}^{-1}\mathcal{L}' = \{v \mid uv \in \mathcal{L}', u \in \mathcal{L}\}$. We abbreviate the concatenation of languages $\{uv \mid u \in \mathcal{L}, v \in \mathcal{L}'\}$ to $\mathcal{L} \cdot \mathcal{L}'$ and we write $\mathcal{L}^*$ for the iteration $\mathcal{L} \cdot \mathcal{L}^*$. We sometimes write $\overline{\mathcal{L}}$ for the complement $\Sigma^* \setminus \mathcal{L}$ and $\overline{A}$ for $\Sigma \setminus A$.

An *extended regular expression* (ERE) on an alphabet $\Sigma$ is a syntactic phrase derivable from non-terminals $r, s, t$. It comprises the the empty word, literals, union, concatenation, Kleene star, as well as negation and intersection operators.

$$r, s, t \ := \ \epsilon \mid A \mid r+s \mid r \cdot s \mid r^* \mid r \& s \mid !r$$

Compared to standard definitions, a *literal* is a set $A$ of symbols, which stands for an abstract, possibly empty, character class. We write $a$ instead of $\{a\}$ for the frequent case

of a single letter literal. We consider regular expressions up to similarity [4], that is, up to associativity and commutativity of the union operator with the empty set as identity.

The language $[\![r]\!] \subseteq \Sigma^*$ of a regular expression $r$ is defined inductively by:

$$
\begin{aligned}
[\![\epsilon]\!] &= \{\epsilon\} & [\![r{+}s]\!] &= [\![r]\!] \cup [\![s]\!] & [\![r\&s]\!] &= [\![r]\!] \cap [\![s]\!] \\
[\![A]\!] &= \{a \mid a \in A\} & [\![r{\cdot}s]\!] &= [\![r]\!]{\cdot}[\![s]\!] & [\![!r]\!] &= \overline{[\![r]\!]} \\
& & [\![r^*]\!] &= [\![r]\!]^*
\end{aligned}
$$

For finite alphabets, $[\![r]\!]$ is a regular language. For arbitrary alphabets, we *define* a language to be regular, if it is equal to $[\![r]\!]$, for some ERE $r$.

We write $r \sqsubseteq s$ ($r$ is *contained* in $s$) to express that $[\![r]\!] \subseteq [\![s]\!]$.

The *nullable* predicate $\nu(r)$ indicates whether $[\![r]\!]$ contains the empty word, that is, $\nu(r)$ iff $\epsilon \in [\![r]\!]$. It is defined inductively by:

$$
\begin{aligned}
\nu(\epsilon) &= \textit{true} & \nu(r{+}s) &= \nu(r) \vee \nu(s) & \nu(r\&s) &= \nu(r) \wedge \nu(s) \\
\nu(A) &= \textit{false} & \nu(r{\cdot}s) &= \nu(r) \wedge \nu(s) & \nu(!r) &= \neg\nu(r) \\
& & \nu(r^*) &= \textit{true}
\end{aligned}
$$

The *Brzozowski derivative* $\partial_a(r)$ of an expression $r$ w.r.t. a symbol $a$ computes a regular expression for the left quotient $a^{-1}[\![r]\!]$ (see [4]). It is defined inductively as follows:

$$
\begin{aligned}
\partial_a(\epsilon) &= \emptyset \\
\partial_a(A) &= \begin{cases} \epsilon, & a \in A \\ \emptyset, & a \notin A \end{cases} & \partial_a(r{\cdot}s) &= \begin{cases} \partial_a(r){\cdot}s{+}\partial_a(s), & \nu(r) \\ \partial_a(r){\cdot}s, & \neg\nu(r) \end{cases} \\
& & \partial_a(r^*) &= \partial_a(r){\cdot}r^* \\
\partial_a(r{+}s) &= \partial_a(r){+}\partial_a(s) & \partial_a(r\&s) &= \partial_a(r)\&\partial_a(s) \\
& & \partial_a(!r) &= !\partial_a(r)
\end{aligned}
$$

The case for the set literal $A$ generalizes Brzozowski's definition. The definition is extended to words by $\partial_{au}(r) = \partial_u(\partial_a(r))$ and $\partial_\epsilon(r) = r$. Hence, $u \in [\![r]\!]$ iff $\epsilon \in [\![\partial_u(r)]\!]$.

## 3    Representing Sets of Symbols

The definition of an ERE in Section 2 just states that a literal is a set of symbols $A \subseteq \Sigma$. However, to define tractable algorithms, we require that $A$ is an element of an effective boolean algebra [19] $(U, \sqcup, \sqcap, \bar{\cdot}, \bot, \top)$ where $U \subseteq \wp(\Sigma)$ is closed under the boolean operations. Here $\sqcup$ and $\sqcap$ denote union and intersection of symbol sets, $\bar{\cdot}$ the complement, and $\bot$ and $\top$ the empty set and the full set $\Sigma$, respectively. In this algebra, we need to be able to decide equality of sets (hence the term *effective*) and to represent singleton symbols.

- For finite (small) alphabets, we may just take $U = \wp(\Sigma)$. A set of symbols may be enumerated and ranges of symbols may be represented by character classes, as customarily supported in regular expression implementations. Alternatively, a bitvector representation may be used.
- If the alphabet is infinite (or just too large), then the boolean algebra of finite and cofinite sets of symbols is the basis for a suitable representation. That is, the set $U = \{A \in \wp(\Sigma) \mid A \text{ finite} \vee \overline{A} \text{ finite}\}$ is effectively closed under the boolean operations.
- In our application to checking access contracts in scripting languages [12], the alphabet itself is a set of words (the field names of objects) composed from another set $\Gamma$ of symbols: $\Sigma \subseteq \wp(\Gamma^*)$. To obtain an effective boolean algebra, we choose the set $U = \{A \subseteq \wp(\Gamma^*) \mid A \text{ is regular}\}$, which is effectively closed under the boolean operations.

- Sets of symbols may also be represented by formulas drawn from a decidable first-order theory over a (finite or infinite) alphabet. For example, the character range `[a-z]` would be represented by the formula $x \geq \text{'a'} \wedge x \leq \text{'z'}$. In this case, the boolean operations get mapped to the disjunction, conjunction, or negation of predicates; bottom and top are false and true, respectively. An SMT solver can decide equality and subset constraints. This approach has been demonstrated to be effective for very large character sets in the work on symbolic finite automata [19].

The rest of this paper is generic with respect to the choice of an effective boolean algebra.

## 4 Antimirov's algorithm for checking containment

Given two regular expressions $r$, $s$, the *containment problem* asks whether $r \sqsubseteq s$. This problem is decidable using standard techniques from automata theory: construct a deterministic finite automaton for $r\&!s$ and check it for emptiness. The drawback of this approach is the expensive construction of the automaton. In general, this expense cannot be avoided because problem is PSPACE-complete [10, 11, 14].

Antimirov [1] proposed an algorithm for deciding containment of standard regular expressions (without intersection and negation) that is based on rewriting of inequalities. His algorithm has the same asymptotic complexity as the automata construction, but it can fail early and is therefore better behaved in practice. We phrase the algorithm in terms of Brzozowski derivatives to avoid introducing Antimirov's notion of partial derivatives.

▶ **Theorem 1** (Containment [1, Proposition 7(2)]). *For regular expressions $r$ and $s$,*

$$r \sqsubseteq s \Leftrightarrow (\forall u \in \Sigma^*) \ \partial_u(r) \sqsubseteq \partial_u(s).$$

Antimirov's algorithm applies this theorem exhaustively to an inequality $r \mathrel{\dot{\sqsubseteq}} s$ (i.e., a proposed containment) to generate all pairs $\partial_u(r) \mathrel{\dot{\sqsubseteq}} \partial_u(s)$ of iterated derivatives until it finds a contradiction or saturation. More precisely, Antimirov defines a *containment calculus* $\mathcal{CC}$ which works on sets $S$ of atoms, where an atom is either an inequality $r \mathrel{\dot{\sqsubseteq}} s$ or a boolean constant *true* or *false*. It consists of the rule CC-DISPROVE which infers *false* from a trivially inconsistent inequality and the rule CC-UNFOLD that applies Theorem 1 to generate new inequalities.

$$
\begin{array}{cc}
\text{CC-DISPROVE} & \text{CC-UNFOLD} \\[4pt]
\dfrac{\nu(r) \wedge \neg\nu(s)}{r \mathrel{\dot{\sqsubseteq}} s \vdash_{\mathcal{CC}} \textit{false}} & \dfrac{\nu(r) \Rightarrow \nu(s)}{r \mathrel{\dot{\sqsubseteq}} s \vdash_{\mathcal{CC}} \{\partial_a(r) \mathrel{\dot{\sqsubseteq}} \partial_a(s) \mid a \in \Sigma\}}
\end{array}
$$

An inference in the calculus for checking whether $r_0 \sqsubseteq s_0$ is a sequence $S_0 \vdash_{\mathcal{CC}} S_1 \vdash_{\mathcal{CC}} S_2 \vdash_{\mathcal{CC}} \ldots$ where $S_0 = \{r_0 \mathrel{\dot{\sqsubseteq}} s_0\}$ and $S_{i+1}$ is an extension of $S_i$ by selecting an inequality in $S_i$ and adding the consequences of applying one of the $\mathcal{CC}$ rules to it. That is, if $r \mathrel{\dot{\sqsubseteq}} s \in S_i$ and $r \mathrel{\dot{\sqsubseteq}} s \vdash_{\mathcal{CC}} S$, then $S_{i+1} = S_i \cup S$.

Antimirov argues [1, Theorem 8] that this algorithm is sound and complete by proving (using Theorem 1) that $r \sqsubseteq s$ does not hold if and only if a set of atoms containing *false* is derivable from $r \mathrel{\dot{\sqsubseteq}} s$. The algorithm terminates because there are only finitely many different inequalities derivable from $r \mathrel{\dot{\sqsubseteq}} s$ using rule CC-UNFOLD.

The containment calculus $\mathcal{CC}$ has two drawbacks. First, the choice of an inequality for the next inference step is nondeterministic. Second, an adaptation to a setting with an infinite alphabet seems doomed because rule CC-UNFOLD requires us to compute the derivative for infinitely many $a \in \Sigma$ at each application. We address the second drawback next.

## 5   Derivatives on Literals

In this section, we develop a variant of Theorem 1 that enables us to define an CC-UNFOLD rule that is guaranteed to add finitely many atoms, even if the alphabet is infinite. First, we observe that we may restrict the symbols considered in rule CC-UNFOLD to initial symbols of the left hand side of an inequality.

▶ **Definition 2** (First). Let $\mathsf{first}(r) := \{a \mid aw \in [\![r]\!]\}$ be the set of initial symbols derivable from regular expression $r$.

Clearly, $(\forall a \in \Sigma)\ \partial_a(r) \sqsubseteq \partial_a(s)$ iff $(\forall b \in \mathsf{first}(r))\ \partial_b(r) \sqsubseteq \partial_b(s)$ because $\partial_b(r) = \emptyset$ for all $b \notin \mathsf{first}(r)$. Thus, CC-UNFOLD does not have to consider the entire alphabet, but unfortunately $\mathsf{first}(r)$ may still be an infinite set of symbols. For that reason, we propose to compute derivatives with respect to *literals* (i.e., non-empty sets of symbols) instead of single symbols. However, generalizing derivatives to literals has some subtle problems.

To illustrate these problems, let us recall the specification of the Brzozowski derivative:

$$[\![\partial_a(r)]\!] = a^{-1}[\![r]\!]$$

Now we might be tempted to consider the following naive extension of the derivative to a set of symbols $A$.

$$[\![\partial_A(r)]\!] \;=\; A^{-1}[\![r]\!] \;=\; \bigcup_{a \in A} a^{-1}[\![r]\!] = \bigcup_{a \in A} [\![\partial_a(r)]\!] \tag{wrong}$$

However, this attempt at a specification yields inconsistent results. To see why, consider the case where $r = \,!s$. Generalizing from $\partial_a(!s) = !\partial_a(s)$, we might try to define $\partial_A(!s) := !\partial_A(s)$. If this definition was sensible, then (1) and (2) should yield the same results:

$$[\![\partial_A(!s)]\!] \;\overset{\text{(wrong)}}{=}\; \bigcup_{a \in A} [\![\partial_a(!s)]\!] \;\overset{\text{def } \partial_a}{=}\; \bigcup_{a \in A} \overline{[\![\partial_a(s)]\!]} \tag{1}$$

$$[\![!\partial_A(s)]\!] \;\overset{\text{def } \partial_a}{=}\; \overline{[\![\partial_A(s)]\!]} \;\overset{\text{(wrong)}}{=}\; \overline{\bigcup_{a \in A} [\![\partial_a(s)]\!]} \;\overset{\text{de Morgan}}{=}\; \bigcap_{a \in A} \overline{[\![\partial_a(s)]\!]} \tag{2}$$

However, we obtain a contradiction: with $A = \{a, b\}$ and $s = a{\cdot}a + b{\cdot}b$, (1) yields $\Sigma^*$ whereas (2) yields $\overline{\{a, b\}}$, which is clearly different.

### 5.1   Positive and Negative Derivatives

To address this problem, we introduce two types of derivative operators with respect to symbol sets. The *positive derivative* $\Delta_A(r)$ computes an expression that contains the union of all $\partial_a(r)$ with $a \in A$, whereas the *negative derivative* $\nabla_A(r)$ computes an expression contained in the intersection of all $\partial_a(r)$ with $a \in A$.

The positive and negative derivative operators are defined by mutual induction and flip at the complement operator. Most cases of their definition are identical to the Brzozowski derivative (cf. Section 2), thus we only show the cases that are different[1]. For all literals $A$ with $[\![A]\!] \neq \emptyset$:

$$\begin{aligned}
\Delta_B(A) &:= \begin{cases} \epsilon, & A \sqcap B \neq \bot \\ \emptyset, & otherwise \end{cases} & \nabla_B(A) &:= \begin{cases} \epsilon, & \overline{A} \sqcap B = \bot \\ \emptyset, & otherwise \end{cases} \\
\Delta_B(!r) &:= \;!\nabla_B(r) & \nabla_B(!r) &:= \;!\Delta_B(r)
\end{aligned}$$

---

[1] See also Appendix A.

For single symbol literals of the form $B = \{a\}$, it holds that $\Delta_a(r) = \nabla_a(r) = \partial_a(r)$. Derivatives with respect to the empty set are defined as $\Delta_\emptyset(r) = \emptyset$ and $\nabla_\emptyset(r) = \Sigma^*$.

The following lemma states the connection between the derivative by a literal and the derivative by a symbol.

▶ **Lemma 3** (Positive and negative derivatives). *For any $r$ and $B$, it holds that:*

$$[\![\Delta_B(r)]\!] \supseteq \bigcup_{a \in B} [\![\partial_a(r)]\!] \qquad\qquad [\![\nabla_B(r)]\!] \subseteq \bigcap_{a \in B} [\![\partial_a(r)]\!]$$

**Proof of Lemma 3.** Both inclusions are proved simultaneously by induction on $r$. See Appendix C. ◀

The following examples illustrate the properties of the derivatives.

▶ **Example 4** (Positive derivative). Let $r$ be $(a \cdot c)\&(b \cdot c)$ and let the literal $A = \{a, b\}$.

$$\Delta_A(r) = \Delta_A(a \cdot c)\&\Delta_A(b \cdot c) = c\&c \sqsupseteq \partial_a(r) + \partial_b(r) = \emptyset + \emptyset$$

▶ **Example 5** (Negative derivative). Let $r$ be $(a \cdot c) + (b \cdot c)$ and let the literal $A = \{a, b\}$.

$$\nabla_A(r) = \nabla_A(a \cdot c) + \nabla_A(b \cdot c) = \emptyset + \emptyset \sqsubseteq \partial_a(r)\&\partial_b(r) = c\&c$$

Positive (negative) derivatives yield an upper (lower) approximation to the information expected from a derivative. This approximation arises because we tried to define the derivative with respect to an *arbitrary* literal $A$. To obtain the precise information, we need to restrict these literals suitably to *next literals*.

## 5.2 Next Literals

An occurrence of a literal $A$ in a regular expression $r$ is *initial* if there is some $a \in \Sigma$ such that $\partial_a(r)$ reduces this occurrence. That is, the computation of $\partial_a(r)$ involves $\partial_a(A)$. Intuitively, $A$ helps determine the first symbol of an element of $[\![r]\!]$.

▶ **Example 6** (Initial Literals).
1. Let $r_1 = \{a, b\}.a^*$. Then $\{a, b\}$ is an initial literal.
2. Let $r_2 = \{a, b\}.a^* + \{b, c\}.c^*$. Then $\{a, b\}$ and $\{b, c\}$ are initial.

Generalizing from the first example, we might be tempted to conjecture that if $A$ is initial in $r$, then $(\forall a, b \in A)\ \partial_a(r) = \partial_b(r)$. However, the second example shows that this conjecture is wrong: $\{a, b\}$ is initial in $r_2$, but $\partial_a(r_2) = a^*$ and $\partial_b(r_2) = a^* + c^*$.

The problem with the second example is that $\{a, b\} \cap \{b, c\} \neq \emptyset$. Hence, instead of identifying initial literals of an ERE $r$, we define a set $\mathsf{next}(r)$ of next literals which are mutually disjoint, whose union contains $\mathsf{first}(r)$, and where the symbols in each literal yield the same derivative. In the second example, it must be that $\mathsf{next}(r_2) = \{\{a\}, \{b\}, \{c\}\}$.

It turns out that this problem arises in a number of cases when defining $\mathsf{next}(r)$ inductively. Hence, we define an operation ⋈ that builds a set of mutually disjoint literals that cover the union of two sets of mutually disjoint literals.

▶ **Definition 7** (Join). Let $\mathfrak{L}_1$ and $\mathfrak{L}_2$ be two sets of mutually disjoint literals.

$$\mathfrak{L}_1 \bowtie \mathfrak{L}_2 := \{(A_1 \sqcap A_2), (A_1 \sqcap \overline{\bigsqcup \mathfrak{L}_2}), (\overline{\bigsqcup \mathfrak{L}_1} \sqcap A_2) \mid A_1 \in \mathfrak{L}_1, A_2 \in \mathfrak{L}_2\}$$

The following lemma states the properties of the join operation.

$$
\begin{aligned}
\mathsf{next}(\epsilon) &= \{\emptyset\} & \mathsf{next}(r{+}s) &= \mathsf{next}(r) \bowtie \mathsf{next}(s) \\
\mathsf{next}(A) &= \{A\} & \mathsf{next}(r{\cdot}s) &= \begin{cases} \mathsf{next}(r) \bowtie \mathsf{next}(s), & \nu(r) \\ \mathsf{next}(r), & \neg\nu(r) \end{cases} \\
& & \mathsf{next}(r^*) &= \mathsf{next}(r) \\
& & \mathsf{next}(r\&s) &= \mathsf{next}(r) \sqcap \mathsf{next}(s) \\
& & \mathsf{next}(!r) &= \mathsf{next}(r) \cup \{\sqcap\{\overline{A} \mid A \in \mathsf{next}(r)\}\}
\end{aligned}
$$

■ **Figure 1** Computing next literals.

▶ **Lemma 8** (Properties of Join). *Let $\mathfrak{L}_1$ and $\mathfrak{L}_2$ be non-empty sets of mutually disjoint literals.*
1.  $\bigcup(\mathfrak{L}_1 \bowtie \mathfrak{L}_2) = \bigcup \mathfrak{L}_1 \cup \bigcup \mathfrak{L}_2$.
2.  $(\forall A \neq A' \in \mathfrak{L}_1 \bowtie \mathfrak{L}_2)\ A \sqcap A' = \emptyset$.
3.  $(\forall A \in \mathfrak{L}_1 \bowtie \mathfrak{L}_2)\ (\forall A_i \in \mathfrak{L}_i)\ A \sqcap A_i \neq \emptyset \Rightarrow A \sqsubseteq A_i$.

**Proof of Lemma 8.** See Appendix D.                                                                    ◀

Figure 1 contains the definition of $\mathsf{next}(r)$. For $\epsilon$ the set of next literals consists of the empty set. The next literal of a literal $A$ is $A$. The next literals of a union $r{+}s$ are computed as the join of the next literals of $r$ and $s$ as explained in Example 6. The next literals of a concatenation $r{\cdot}s$ are the next literals of $r$ if $r$ is not nullable. Otherwise, they are the join of the next literals of both operands. The next literals of a Kleene star expression $r^*$ are the next literals of $r$. For an intersection $r\&s$, the set of next literals is the set of all intersections $A \sqcap A'$ of the next literals of both operands. In this case, the join operation $\bowtie$ is not needed because symbols that only appear in literals from one operand can be elided. To see this, consider $\mathsf{next}(a\&b) = \{\{a\} \sqcap \{b\}\} = \{\emptyset\}$ whereas $\{\{a\}\} \bowtie \{\{b\}\} = \{\emptyset, \{a\}, \{b\}\}$.

The set of next literals of $!r$ comprises the next literals of $r$ and a new literal, which is the intersection of the complements of all literals in $\mathsf{next}(r)$. We might contemplate to exclude literals that contain symbols $a$ such that $\partial_a(r)$ is equivalent to $\Sigma^*$, but we refrain from doing so because this equivalence cannot be decided with a finite set of rewrite rules [16].

The function $\mathsf{next}(r) \setminus \{\emptyset\}$ computes the equivalence classes of a partial equivalence relation $\sim$ on $\Sigma$ such that equivalent symbols yield the same derivative on $r$. The relation is defined by $a \sim b$ if there exists $A \in \mathsf{next}(r)$ such that $a \in A$ and $b \in A$. Furthermore, the derivative by a symbol that is not part of the relation yields the empty set.

▶ **Lemma 9** (Partial Equivalence). *Let $\mathfrak{L} = \mathsf{next}(r)$.*
1.  $(\forall A \in \mathfrak{L})\ (\forall a, b \in A)\ \partial_a(r) = \partial_b(r)$
2.  $(\forall a \notin \bigcup \mathfrak{L})\ \partial_a(r) = \emptyset$

**Proof of Lemma 9.** See Appendix E.                                                                    ◀

It remains to show that $\mathsf{next}(r)$ covers all symbols in $\mathsf{first}(r)$.

▶ **Lemma 10** (First). *For all $r$, $\bigcup \mathsf{next}(r) \supseteq \mathsf{first}(r)$.*

**Proof of Lemma 10.** See Appendix F.                                                                   ◀

Moreover, there are only finitely many different next literals for each regular expression.

▶ **Lemma 11** (Finiteness). *For all $r$, $|\mathsf{next}(r)|$ is finite.*

**Proof of Lemma 11.** By induction on $r$. The base cases construct finite sets and the inductive cases build a finite number of combinations of the results from the subexpressions. ◄

Now, we put next literals to work. If we only take positive or negative derivatives with respect to next literals, then the inclusions in Lemma 3 turn into equalities. The result is that both the positive and the negative derivative, when applied to a next literal $A$, calculate a regular expression for the left quotient $A^{-1}[\![r]\!]$.

▶ **Theorem 12** (Left Quotient). *For all $r$, $A \in \mathsf{next}(r) \setminus \{\emptyset\}$, and $a \in [\![A]\!]$:*

$$[\![\Delta_A(r)]\!] = [\![\nabla_A(r)]\!] = [\![\partial_a(r)]\!]$$

**Proof of Lemma 12.** By induction on $r$. See Appendix G. ◄

Motivated by this result, we define the Brzozowski derivative for a non-empty subset $A$ of a literal in $\mathsf{next}(r)$. This definition involves an arbitrary choice of $a \in A$, but this choice does not influence the calculated derivative according to Lemma 9, part 1.

▶ **Definition 13.** Let $A' \in \mathsf{next}(r)$. For each $\emptyset \neq A \subseteq A'$ define $\partial_A(r) := \partial_a(r)$, where $a \in A$.

▶ **Lemma 14** (Coverage). *For all $a$, $u$, and $r$ it holds that:*

$$u \in [\![\partial_a(r)]\!] \iff \exists A \in \mathsf{next}(r) : a \in A \land u \in [\![\Delta_A(r)]\!] \land u \in [\![\nabla_A(r)]\!]$$

**Proof of Lemma 14.** This result follows from Theorem 12 and Lemma 10. ◄

We conclude that to determine a finite set of representatives for all derivatives of a regular expression $r$ it is sufficient to select one symbol $a$ from each equivalence class $A \in \mathsf{next}(r) \setminus \{\emptyset\}$ and calculate $\partial_a(r)$. Alternatively, we may calculate $\Delta_A(r)$ or $\nabla_A(r)$ according to Theorem 12. It remains to lift this result to solving inequalities.

## 6 Solving Inequalities

Theorem 1 is the foundation of Antimirov's algorithm. It turns out that we can prove a stronger version of this theorem, which makes the rules CC-DISPROVE and CC-UNFOLD sound and complete and which also encompasses the soundness of the restriction to first sets.

▶ **Theorem 15** (Containment).

$$r \sqsubseteq s \iff (\nu(r) \Rightarrow \nu(s)) \land (\forall a \in \mathsf{first}(r)) \ \partial_a(r) \sqsubseteq \partial_a(s)$$

**Proof of Theorem 15.** See Appendix H. ◄

As we remarked before, it may be very expensive (or even impossible) to construct all derivatives with respect to the first symbols, particularly for negated expressions and for large or infinite alphabets. To obtain a decision procedure for containment, we need a finite set of derivatives. Therefore, we use next literals as representatives of the first symbols and use Brzozowski derivatives on literals (Definition 13) on both sides.

To define the next literals of an inequality $r \mathrel{\dot{\sqsubseteq}} s$, it would be sound to use the join of the next literals of both sides: $\mathsf{next}(r) \bowtie \mathsf{next}(s)$. However, we can do slightly better. Theorem 15 proves that the first symbols of $r$ are sufficient to prove containment. Using the full join operation, however, would cover $\mathsf{first}(r) \cup \mathsf{first}(s)$ (by Lemma 10). Hence, we define a left-biased version of the join operator that only covers the symbols of its left operand.

▶ **Definition 16** (Left Join). Let $\mathfrak{L}_1$ and $\mathfrak{L}_2$ be two sets of mutually disjoint literals.

$$\mathfrak{L}_1 \ltimes \mathfrak{L}_2 := \{(A_1 \sqcap A_2), (A_1 \sqcap \overline{\bigsqcup \mathfrak{L}_2}) \mid A_1 \in \mathfrak{L}_1, A_2 \in \mathfrak{L}_2\}$$

The following lemma states the properties of the left join operation.

▶ **Lemma 17** (Properties of Left Join). *Let $\mathfrak{L}_1$ and $\mathfrak{L}_2$ be non-empty sets of mutually disjoint literals.*
*1.* $\bigcup(\mathfrak{L}_1 \ltimes \mathfrak{L}_2) = \bigcup \mathfrak{L}_1$.
*2.* $(\forall A \neq A' \in \mathfrak{L}_1 \ltimes \mathfrak{L}_2)\ A \sqcap A' = \emptyset$.
*3.* $(\forall A \in \mathfrak{L}_1 \ltimes \mathfrak{L}_2)\ (\forall A_i \in \mathfrak{L}_i)\ A \sqcap A_i \neq \emptyset \Rightarrow A \sqsubseteq A_i$.

**Proof of Lemma 17.** Analogous to the proof of Lemma 8 in Appendix D.          ◀

▶ **Definition 18** (Next Literals of an Inequality). Let $r \mathrel{\dot{\sqsubseteq}} s$ be an inequality.

$$\mathsf{next}(r \mathrel{\dot{\sqsubseteq}} s) := \mathsf{next}(r) \ltimes \mathsf{next}(s)$$

Finally, we can state a generalization of Antimirov's containment theorem for EREs, where each unfolding step generates only finitely many derivatives.

▶ **Theorem 19** (Containment). *For all regular expressions $r$ and $s$,*

$$r \sqsubseteq s \ \Leftrightarrow\ (\nu(r) \Rightarrow \nu(s)) \ \wedge\ (\forall A \in \mathbf{\textit{next}}(r \mathrel{\dot{\sqsubseteq}} s))\ \partial_A(r) \sqsubseteq \partial_A(s).$$

**Proof of Theorem 19.** The proof is by contraposition. If $r \not\sqsubseteq s$ then $\exists A \in \mathsf{next}(r \mathrel{\dot{\sqsubseteq}} s)$ : $\partial_A(r) \not\sqsubseteq \partial_A(s)$ or $\neg(\nu(r) \Rightarrow \nu(s))$. See also Appendix I.          ◀

For $A \in \mathsf{next}(r \mathrel{\dot{\sqsubseteq}} s)$ define $\nabla_A(r \mathrel{\dot{\sqsubseteq}} s) := (\nabla_A(r) \mathrel{\dot{\sqsubseteq}} \Delta_A(s)) = (\partial_A(r) \mathrel{\dot{\sqsubseteq}} \partial_A(s))$.

▶ **Theorem 20** (Finiteness). *Let $R$ be a finite set of regular inequalities. Define*

$$F(R) = R \cup \{\nabla_A(r \mathrel{\dot{\sqsubseteq}} s) \mid r \mathrel{\dot{\sqsubseteq}} s \in R, A \in \mathbf{\textit{next}}(r \mathrel{\dot{\sqsubseteq}} s)\}$$

*For each $r$ and $s$, the set $\bigcup_{i \in \mathbb{N}} F^{(i)}(\{r \sqsubseteq s\})$ is finite.*

**Proof of Theorem 20.** As we consider regular expressions up to similarity (as defined by Brzozowski [4]) and $\nabla_A(r \mathrel{\dot{\sqsubseteq}} s) = \partial_A(r) \mathrel{\dot{\sqsubseteq}} \partial_A(s)$ is essentially applying the Brzozowski derivative to a pair of (extended) regular expressions, we know that the set of these pairs is finite (because there are only finitely many dissimilar iterated Brzozowski derivatives for a regular expression [4]).          ◀

These results are the basis for a complete decision procedure for solving inequalities on extended regular expressions where literals are defined via an effective boolean algebra. Figure 2 defines this procedure as a judgment of the form $\Gamma \vdash r \mathrel{\dot{\sqsubseteq}} s \ :\ b$, where $\Gamma$ is a set of previous visited inequalities $r \mathrel{\dot{\sqsubseteq}} s$ with $\nu(r) \Rightarrow \nu(s)$ that are assumed to be true and $b \in \{true, false\}$. The effective boolean algebra comes into play in the computation of the next literals and in the computation of the derivatives.

Rule (DISPROVE) detects contradictory inequalities in the same way as Antimirov's rule CC-DISPROVE. Rule (CYCLE) detects circular reasoning: Under the assumption that $r \mathrel{\dot{\sqsubseteq}} s$ holds we were not (yet) able to derive a contradiction and thus conclude that $r \mathrel{\dot{\sqsubseteq}} s$ holds. This rule guarantees termination because of the finiteness result (Theorem 20). The rules (UNFOLD-TRUE) and (UNFOLD-FALSE) apply only if $r \mathrel{\dot{\sqsubseteq}} s$ is neither contradictory nor in the context. A deterministic implementation would generate the literals $A \in \mathsf{next}(r \mathrel{\dot{\sqsubseteq}} s)$ and recursively check $\nabla_A(r \mathrel{\dot{\sqsubseteq}} s)$. If any of these checks returns false, then (UNFOLD-FALSE) fires. Otherwise (UNFOLD-TRUE) signals a successful containment proof. Theorem 19 is the basis for soundness and completeness of the unfolding rules.

$$
\begin{array}{cc}
(\text{Disprove}) & (\text{Cycle}) \\[4pt]
\dfrac{\nu(r) \qquad \neg\nu(s)}{\Gamma \;\vdash\; r \mathrel{\dot{\sqsubseteq}} s \;:\; \mathit{false}} & \dfrac{r \mathrel{\dot{\sqsubseteq}} s \in \Gamma}{\Gamma \;\vdash\; r \mathrel{\dot{\sqsubseteq}} s \;:\; \mathit{true}}
\end{array}
$$

$(\text{Unfold-True})$
$$
\dfrac{r \mathrel{\dot{\sqsubseteq}} s \notin \Gamma \qquad \nu(r) \Rightarrow \nu(s) \qquad \forall A \in \mathsf{next}(r \mathrel{\dot{\sqsubseteq}} s):\ \Gamma \cup \{r \mathrel{\dot{\sqsubseteq}} s\} \;\vdash\; \partial_A(r) \mathrel{\dot{\sqsubseteq}} \partial_A(s) \;:\; \mathit{true}}{\Gamma \;\vdash\; r \mathrel{\dot{\sqsubseteq}} s \;:\; \mathit{true}}
$$

$(\text{Unfold-False})$
$$
\dfrac{r \mathrel{\dot{\sqsubseteq}} s \notin \Gamma \qquad \nu(r) \Rightarrow \nu(s) \qquad \exists A \in \mathsf{next}(r \mathrel{\dot{\sqsubseteq}} s):\ \Gamma \cup \{r \mathrel{\dot{\sqsubseteq}} s\} \;\vdash\; \partial_A(r) \mathrel{\dot{\sqsubseteq}} \partial_A(s) \;:\; \mathit{false}}{\Gamma \;\vdash\; r \mathrel{\dot{\sqsubseteq}} s \;:\; \mathit{false}}
$$

**Figure 2** Decision procedure for containment.

$$
\begin{array}{cccc}
(\text{Prove-Identity}) & (\text{Prove-Empty}) & (\text{Prove-Nullable}) & (\text{Disprove-Empty}) \\[4pt]
 & & \dfrac{\nu(s)}{\Gamma \;\vdash\; \epsilon \sqsubseteq s \;:\; \mathit{true}} & \dfrac{\exists A \in \mathsf{next}(r):\ A \neq \emptyset}{\Gamma \;\vdash\; r \sqsubseteq \emptyset \;:\; \mathit{false}} \\[4pt]
\Gamma \;\vdash\; r \sqsubseteq r \;:\; \mathit{true} & \Gamma \;\vdash\; \emptyset \sqsubseteq s \;:\; \mathit{true} & &
\end{array}
$$

**Figure 3** Prove and disprove axioms.

▶ **Theorem 21** (Soundness). *For all regular expression $r$ and $s$:*

$$
\emptyset \;\vdash\; r \mathrel{\dot{\sqsubseteq}} s \;:\; \top \;\Leftrightarrow\; r \sqsubseteq s
$$

**Proof of Theorem 21.** We prove that $\Gamma \vdash r \mathrel{\dot{\sqsubseteq}} s \;:\; \mathit{false}$ iff $r \not\sqsubseteq s$, for all contexts $\Gamma$ where $r \mathrel{\dot{\sqsubseteq}} s \notin \Gamma$. This is sufficient because each regular inequality gives rise to a finite derivation by Theorem 20. See Appendix J for details.    ◀

In addition to the rules from Figure 2, we may add auxiliary rules to detect trivially consistent or inconsistent inequalities early (Figure 3 contains some examples). Such rules may be used to improve efficiency. They decide containment directly instead of unfolding repeatedly.

## 7    Conclusion

We extended Antimirov's algorithm for proving containment of regular expressions to extended regular expressions on potentially infinite alphabets. To work effectively with such alphabets, we require that literals in regular expressions are drawn from an effective boolean algebra. As a slight difference, we work with Brzozowski derivatives instead of Antimirov's notion of partial derivative.

The main effort in lifting Antimirov's algorithm is to identify, for each regular inequality $r \mathrel{\dot{\sqsubseteq}} s$, a finite set of symbols such that calculating the derivation with respect to these symbols covers all possible derivations with all symbols. We regard the construction of the set of suitable representatives, embodied in the notion of next literals $\mathsf{next}(r \mathrel{\dot{\sqsubseteq}} s)$, as a key contribution of this work.

────── **References** ──────

**1** Valentin M. Antimirov. Rewriting regular inequalities. In Horst Reichel, editor, *FCT*, volume 965 of *LNCS*, pages 116–125. Springer, 1995.

**2** Valentin M. Antimirov. Partial derivatives of regular expressions and finite automaton constructions. *Theoretical Computer Science*, 155(2):291–319, 1996.

**3** Filippo Bonchi and Damien Pous. Checking NFA equivalence with bisimulations up to congruence. In Roberto Giacobazzi and Radhia Cousot, editors, *POPL*, pages 457–468, Rome, Italy, January 2013. ACM.

**4** Janusz A. Brzozowski. Derivatives of regular expressions. *J. ACM*, 11(4):481–494, 1964.

**5** Pascal Caron, Jean-Marc Champarnaud, and Ludovic Mignot. Partial derivatives of an extended regular expression. In Adrian Horia Dediu, Shunsuke Inenaga, and Carlos Martín-Vide, editors, *LATA*, volume 6638 of *LNCS*, pages 179–191. Springer, 2011.

**6** A. Ginzburg. A procedure for checking equality of regular expressions. *J. ACM*, 14(2):355–362, April 1967.

**7** Victor M. Glushkov. The abstract theory of automata. *Russian Mathematical Surveys*, 16(5):1–53, 1961.

**8** Phillip Heidegger, Annette Bieniusa, and Peter Thiemann. Access permission contracts for scripting languages. In John Field and Michael Hicks, editors, *Proc. 39th ACM Symp. POPL*, pages 111–122, Philadelphia, USA, January 2012. ACM Press.

**9** John Edward Hopcroft and Richard Manning Karp. A linear algorithm for testing equivalence of finite automata. Technical report, Cornell University, 1971.

**10** Harry B. Hunt III, Daniel J. Rosenkrantz, and Thomas G. Szymanski. On the equivalence, containment, and covering problems for the regular and context-free languages. *J. Comput. Syst. Sci.*, 12(2):222–268, 1976.

**11** Tao Jiang and Bala Ravikumar. Minimal NFA problems are hard. *SIAM J. Comput.*, 22(6):1117–1141, 1993.

**12** Matthias Keil and Peter Thiemann. Efficient dynamic access analysis using JavaScript proxies. In *Proceedings of the 9th Symposium on Dynamic Languages*, DLS '13, pages 49–60, New York, NY, USA, 2013. ACM.

**13** Robert McNaughton and Hisao Yamada. Regular expressions and state graphs for automata. *Electronic Computers, IRE Transactions on*, EC-9(1):39–47, 1960.

**14** Albert R. Meyer and Larry J. Stockmeyer. The equivalence problem for regular expressions with squaring requires exponential space. In *SWAT (FOCS)*, pages 125–129. IEEE Computer Society, 1972.

**15** Scott Owens, John H. Reppy, and Aaron Turon. Regular-expression derivatives re-examined. *J. Funct. Program.*, 19(2):173–190, 2009.

**16** Valentin N. Redko. On defining relations for the algebra of regular events. *Ukrain. Mat.*, 16:120–126, 1964.

**17** Ken Thompson. Regular expression search algorithm. *Commun. ACM*, 11(6):419–422, 1968.

**18** Gertjan van Noord and Dale Gerdemann. Finite state transducers with predicates and identities. *Grammars*, 4(3):263–286, 2001.

**19** Margus Veanes. Applications of symbolic finite automata. In Stavros Konstantinidis, editor, *CIAA*, volume 7982 of *Lecture Notes in Computer Science*, pages 16–23, Halifax, NS, Canada, 2013. Springer.

**20** Bruce W. Watson. Implementing and using finite automata toolkits. *Nat. Lang. Eng.*, 2(4):295–302, December 1996.

## A Positive and Negative Derivatives

This sections shows the full definition of the *positive* and *negative* derivative operator. The operators are defined by induction and flip on the complement operator.

### A.1 Positive Derivatives

For all literals $A \neq \emptyset$:

$$
\begin{aligned}
\Delta_A(\epsilon) \quad &:= \quad \emptyset \\
\Delta_A(B) \quad &:= \quad \begin{cases} \epsilon, & A \sqcap B \neq \bot \\ \emptyset, & \text{otherwise} \end{cases} \\
\Delta_A(r^*) \quad &:= \quad \Delta_A(r) \cdot r^* \\
\Delta_A(r+s) \quad &:= \quad \Delta_A(r) + \Delta_A(s) \\
\Delta_A(r\&s) \quad &:= \quad \Delta_A(r) \& \Delta_A(s) \\
\Delta_A(!r) \quad &:= \quad !\nabla_A(r) \\
\Delta_A(r \cdot s) \quad &:= \quad \begin{cases} \Delta_A(r) \cdot s + \Delta_A(s), & \nu(r) \\ \Delta_A(r) \cdot s, & \text{otherwise} \end{cases}
\end{aligned}
$$

### A.2 Negative Derivatives

For all literals $A \neq \emptyset$:

$$
\begin{aligned}
\nabla_A(\epsilon) \quad &:= \quad \emptyset \\
\nabla_A(B) \quad &:= \quad \begin{cases} \epsilon, & A \sqcap \overline{B} = \bot \\ \emptyset, & \text{otherwise} \end{cases} \\
\nabla_A(r^*) \quad &:= \quad \nabla_A(r) \cdot r^* \\
\nabla_A(r+s) \quad &:= \quad \nabla_A(r) + \nabla_A(s) \\
\nabla_A(r\&s) \quad &:= \quad \nabla_A(r) \& \nabla_A(s) \\
\nabla_A(!r) \quad &:= \quad !\Delta_A(r) \\
\nabla_A(r \cdot s) \quad &:= \quad \begin{cases} \nabla_A(r) \cdot s + \nabla_A(s), & \nu(r) \\ \nabla_A(r) \cdot s, & \text{otherwise} \end{cases}
\end{aligned}
$$

## B     Complexity

This section comprises the complexity of the decision procedure. The complexity has two sources: building the next literals and computing the derivatives.

We express the complexity in terms of the size of a regular expression. The size is directly related to the number of derivation steps and to the number of operations if gathering the next literals.

▶ **Definition 22** (Size). The size $\mathcal{S}(r)$ of a regular expression $r$ is the number of expression constructors and literals.

$$
\begin{array}{llllll}
\mathcal{S}(\epsilon) & = & 1 & \mathcal{S}(r^*) & = & \mathcal{S}(r) + 1 \\
\mathcal{S}(A) & = & 1 & \mathcal{S}(r{+}s) & = & \mathcal{S}(r) + \mathcal{S}(s) + 1 \\
& & & \mathcal{S}(r{\cdot}s) & = & \mathcal{S}(r) + \mathcal{S}(s) + 1
\end{array}
\qquad
\begin{array}{lll}
\mathcal{S}(r\&s) & = & \mathcal{S}(r) + \mathcal{S}(s) + 1 \\
\mathcal{S}(!r) & = & \mathcal{S}(r) + 1
\end{array}
$$

The number of literals in a regular expression is another useful measure.

▶ **Definition 23** (Literal Width). The literal width $\|r\|$ of a regular expression $r$ denotes the total number of literals $A$ in $r$.

Calculating the next literals for $r$ may require a number of operations on the symbol set representation which is exponential in the literal width $\|r\|$, because there are regular expressions where the number of next literals is already exponential. For example, consider $r = (A_1{+}B_1)\&(A_2{+}B_2)\&\ldots(A_n{+}B_n)$ with $\|r\| = 2n$. With a sufficiently large alphabet, we may choose the sets $A_i$ and $B_i$ such that $|\mathsf{next}(r)| = 2^n$.

The number of different derivatives of a regular expression is bounded by $2^{\mathcal{S}(r)}$ analogously to Brzozowski's result. Hence, the number of different derivatives of a regular expression inequality $r \sqsubseteq s$ is bounded by $2^{\mathcal{S}(r)+\mathcal{S}(s)}$. Taken together, our decision procedure requires the computation of an exponential number of derivative operations and, for the result of each of these operations, a new set of next literals has to be determined, in the worst case.

The derivative itself runs in constant time in most cases. However, in the case where the argument expression is a symbol-set literal, a calculation on the representation of symbol sets is required.

## C Lemma 3: Positive and Negative Derivatives

**Proof of Lemma 3.** For any ERE $r$ and for any literal $A$, the following equation holds:

$$[\![\nabla_A(r)]\!] \subseteq \bigcap_{a \in A} [\![\partial_a(r)]\!] \tag{3}$$

$$[\![\Delta_A(r)]\!] \supseteq \bigcup_{a \in A} [\![\partial_a(r)]\!] \tag{4}$$

Proof by induction on $r$.

**Case** $r = \epsilon$: Claim holds because $[\![\nabla_A(\epsilon)]\!] = [\![\Delta_A(\epsilon)]\!] = \emptyset$.

**Case** $r = B$: Claim holds because

$$[\![\nabla_A(B)]\!] = \bigcap_{a \in A} [\![\partial_a(B)]\!] = \begin{cases} \{\epsilon\}, & A \subseteq B \\ \emptyset, & otherwise \end{cases} \tag{5}$$

and

$$[\![\Delta_A(B)]\!] = \bigcup_{a \in A} [\![\partial_a(B)]\!] = \begin{cases} \{\epsilon\}, & A \cap B \neq \emptyset \\ \emptyset, & otherwise \end{cases} \tag{6}$$

**Case** $r = s^*$: By induction

$$[\![\nabla_A(s)]\!] \overset{\mathsf{IH}}{\subseteq} \bigcap_{a \in A} [\![\partial_a(s)]\!] \tag{7}$$

$$[\![\Delta_A(s)]\!] \overset{\mathsf{IH}}{\supseteq} \bigcup_{a \in A} [\![\partial_a(s)]\!] \tag{8}$$

holds. We obtain that

$$\forall a : \ [\![\partial_a(s^*)]\!] = [\![\partial_a(s) \cdot s^*]\!] \tag{9}$$

$$\forall A : \ [\![\nabla_A(s^*)]\!] = [\![\nabla_A(s) \cdot s^*]\!] \tag{10}$$

$$\forall A : \ [\![\Delta_A(s^*)]\!] = [\![\Delta_A(s) \cdot s^*]\!] \tag{11}$$

holds. Claim holds because

$$\forall A : \ [\![\nabla_A(s^*)]\!] = [\![\nabla_A(s) \cdot s^*]\!] \tag{12}$$

$$\overset{\mathsf{IH}}{\subseteq} \{uv \mid u \in \bigcap_{a \in A} [\![\partial_a(s)]\!], v \in [\![s^*]\!]\} \tag{13}$$

$$= \bigcap_{a \in A} \{uv \mid u \in [\![\partial_a(s)]\!], v \in [\![s^*]\!]\} \tag{14}$$

$$= \bigcap_{a \in A} [\![\partial_a(s^*)]\!] \tag{15}$$

and

$$\forall A : \ [\![\Delta_A(s^*)]\!] = [\![\Delta_A(s) \cdot s^*]\!] \tag{16}$$

$$\overset{\mathsf{IH}}{\supseteq} \{uv \mid u \in \bigcup_{a \in A} [\![\partial_a(s)]\!], v \in [\![s^*]\!]\} \tag{17}$$

$$= \bigcup_{a \in A} \{uv \mid u \in [\![\partial_a(s)]\!], v \in [\![s^*]\!]\} \tag{18}$$

$$= \bigcup_{a \in A} [\![\partial_a(s^*)]\!] \tag{19}$$

**Case** $r = s+t$: By induction

$$\llbracket \nabla_A(s) \rrbracket \overset{\text{IH}}{\subseteq} \bigcap_{a \in A} \llbracket \partial_a(s) \rrbracket \tag{20}$$

$$\llbracket \nabla_A(t) \rrbracket \overset{\text{IH}}{\subseteq} \bigcap_{a \in A} \llbracket \partial_a(t) \rrbracket \tag{21}$$

$$\llbracket \Delta_A(s) \rrbracket \overset{\text{IH}}{\supseteq} \bigcup_{a \in A} \llbracket \partial_a(s) \rrbracket \tag{22}$$

$$\llbracket \Delta_A(t) \rrbracket \overset{\text{IH}}{\supseteq} \bigcup_{a \in A} \llbracket \partial_a(t) \rrbracket \tag{23}$$

holds. We obtain that

$$\llbracket \partial_a(s+t) \rrbracket = \llbracket \partial_a(s) \rrbracket \cup \llbracket \partial_a(t) \rrbracket \tag{24}$$
$$\llbracket \nabla_A(s+t) \rrbracket = \llbracket \nabla_A(s) \rrbracket \cup \llbracket \nabla_A(t) \rrbracket \tag{25}$$
$$\llbracket \Delta_A(s+t) \rrbracket = \llbracket \Delta_A(s) \rrbracket \cup \llbracket \Delta_A(t) \rrbracket \tag{26}$$

holds. Claim holds because

$$\llbracket \nabla_A(s+t) \rrbracket = \llbracket \nabla_A(s) \rrbracket \cup \llbracket \nabla_A(t) \rrbracket \tag{27}$$

$$\overset{\text{IH}}{\subseteq} \bigcap_{a \in A} \llbracket \partial_a(s) \rrbracket \cup \bigcap_{a \in A} \llbracket \partial_a(t) \rrbracket \tag{28}$$

$$\subseteq \bigcap_{a \in A} \llbracket \partial_a(s) \rrbracket \cup \llbracket \partial_a(t) \rrbracket \tag{29}$$

$$= \bigcap_{a \in A} \llbracket \partial_a(s+t) \rrbracket \tag{30}$$

and

$$\llbracket \Delta_A(s+t) \rrbracket = \llbracket \Delta_A(s) \rrbracket \cup \llbracket \Delta_A(t) \rrbracket \tag{31}$$

$$\overset{\text{IH}}{\supseteq} \bigcup_{a \in A} \llbracket \partial_a(s) \rrbracket \cup \bigcup_{a \in A} \llbracket \partial_a(t) \rrbracket \tag{32}$$

$$= \bigcup_{a \in A} \llbracket \partial_a(s) \rrbracket \cup \llbracket \partial_a(t) \rrbracket \tag{33}$$

$$= \bigcup_{a \in A} \llbracket \partial_a(s+t) \rrbracket \tag{34}$$

**Case** $r = s\&t$: By induction

$$\llbracket \nabla_A(s) \rrbracket \overset{\text{IH}}{\subseteq} \bigcap_{a \in A} \llbracket \partial_a(s) \rrbracket \tag{35}$$

$$\llbracket \nabla_A(t) \rrbracket \overset{\text{IH}}{\subseteq} \bigcap_{a \in A} \llbracket \partial_a(t) \rrbracket \tag{36}$$

$$\llbracket \Delta_A(s) \rrbracket \overset{\text{IH}}{\supseteq} \bigcup_{a \in A} \llbracket \partial_a(s) \rrbracket \tag{37}$$

$$\llbracket \Delta_A(t) \rrbracket \overset{\text{IH}}{\supseteq} \bigcup_{a \in A} \llbracket \partial_a(t) \rrbracket \tag{38}$$

holds. We obtain that

$$\llbracket \partial_a(s\&t) \rrbracket = \llbracket \partial_a(s) \rrbracket \cap \llbracket \partial_a(t) \rrbracket \tag{39}$$

$$\llbracket \nabla_A(s\&t) \rrbracket = \llbracket \nabla_A(s) \rrbracket \cap \llbracket \nabla_A(t) \rrbracket \tag{40}$$

$$\llbracket \Delta_A(s\&t) \rrbracket = \llbracket \Delta_A(s) \rrbracket \cap \llbracket \Delta_A(t) \rrbracket \tag{41}$$

holds. Claim holds because

$$\llbracket \nabla_A(s\&t) \rrbracket = \llbracket \nabla_A(s) \rrbracket \cap \llbracket \nabla_A(t) \rrbracket \tag{42}$$

$$\overset{\text{IH}}{\subseteq} \bigcap_{a \in A} \llbracket \partial_a(s) \rrbracket \cap \bigcap_{a \in A} \llbracket \partial_a(t) \rrbracket \tag{43}$$

$$= \bigcap_{a \in A} \llbracket \partial_a(s) \rrbracket \cap \llbracket \partial_a(t) \rrbracket \tag{44}$$

$$= \bigcap_{a \in A} \llbracket \partial_a(s\&t) \rrbracket \tag{45}$$

and

$$\llbracket \Delta_A(s\&t) \rrbracket = \llbracket \Delta_A(s) \rrbracket \cap \llbracket \Delta_A(t) \rrbracket \tag{46}$$

$$\overset{\text{IH}}{\supseteq} \bigcup_{a \in A} \llbracket \partial_a(s) \rrbracket \cap \bigcup_{a \in A} \llbracket \partial_a(t) \rrbracket \tag{47}$$

$$\supseteq \bigcup_{a \in A} \llbracket \partial_a(s) \rrbracket \cap \llbracket \partial_a(t) \rrbracket \tag{48}$$

$$= \bigcup_{a \in A} \llbracket \partial_a(s\&t) \rrbracket \tag{49}$$

**Case** $r = {!s}$: By induction

$$\llbracket \nabla_A(s) \rrbracket \overset{\text{IH}}{\subseteq} \bigcap_{a \in A} \llbracket \partial_a(s) \rrbracket \tag{50}$$

$$\llbracket \Delta_A(s) \rrbracket \overset{\text{IH}}{\supseteq} \bigcup_{a \in A} \llbracket \partial_a(s) \rrbracket \tag{51}$$

holds. We obtain that

$$\forall a : \quad \llbracket \partial_a({!s}) \rrbracket = \Sigma^* \backslash \llbracket \partial_a(s) \rrbracket \tag{52}$$

$$\forall A : \quad \llbracket \nabla_A({!s}) \rrbracket = \Sigma^* \backslash \llbracket \Delta_A(s) \rrbracket \tag{53}$$

$$\forall A : \quad \llbracket \Delta_A({!s}) \rrbracket = \Sigma^* \backslash \llbracket \nabla_A(s) \rrbracket \tag{54}$$

holds. Claim holds because

$$\forall A : \quad \llbracket \nabla_A({!s}) \rrbracket = \Sigma^* \backslash \llbracket \Delta_A(s) \rrbracket \tag{55}$$

$$\overset{\text{IH}}{\subseteq} \Sigma^* \backslash \bigcup_{a \in A} \llbracket \partial_a(s) \rrbracket \tag{56}$$

$$= \bigcap_{a \in A} \Sigma^* \backslash \llbracket \partial_a(s) \rrbracket \tag{57}$$

$$= \bigcap_{a \in A} \llbracket \partial_a({!s}) \rrbracket \tag{58}$$

and

$$\forall A : \ [\![\Delta_A(!s)]\!] \ = \Sigma^*\backslash[\![\nabla_A(s)]\!] \tag{59}$$

$$\overset{\text{IH}}{\supseteq} \Sigma^*\backslash \bigcap_{a\in A} [\![\partial_a(s)]\!] \tag{60}$$

$$= \bigcup_{a\in A} \Sigma^*\backslash[\![\partial_a(s)]\!] \tag{61}$$

$$= \bigcup_{a\in A} [\![\partial_a(!s)]\!] \tag{62}$$

**Case** $r = s{\cdot}t$:  By induction

$$[\![\nabla_A(s)]\!] \ \overset{\text{IH}}{\subseteq} \ \bigcap_{a\in A} [\![\partial_a(s)]\!] \tag{63}$$

$$[\![\nabla_A(t)]\!] \ \overset{\text{IH}}{\subseteq} \ \bigcap_{a\in A} [\![\partial_a(t)]\!] \tag{64}$$

$$[\![\Delta_A(s)]\!] \ \overset{\text{IH}}{\supseteq} \ \bigcup_{a\in A} [\![\partial_a(s)]\!] \tag{65}$$

$$[\![\Delta_A(t)]\!] \ \overset{\text{IH}}{\supseteq} \ \bigcup_{a\in A} [\![\partial_a(t)]\!] \tag{66}$$

holds. We obtain that

$$\forall a : \ [\![\partial_a(s{\cdot}t)]\!] \ = \begin{cases} [\![\partial_a(s)\cdot t]\!] \cup [\![\partial_a(t)]\!], & \nu(s) \\ [\![\partial_a(s)\cdot t]\!], & otherwise \end{cases} \tag{67}$$

$$\forall A : \ [\![\nabla_A(s{\cdot}t)]\!] \ = \begin{cases} [\![\nabla_A(s)\cdot t]\!] \cup [\![\nabla_A(t)]\!], & \nu(s) \\ [\![\nabla_A(s)\cdot t]\!], & otherwise \end{cases} \tag{68}$$

$$\forall A : \ [\![\Delta_A(s{\cdot}t)]\!] \ = \begin{cases} [\![\Delta_A(s)\cdot t]\!] \cup [\![\Delta_A(t)]\!], & \nu(s) \\ [\![\Delta_A(s)\cdot t]\!], & otherwise \end{cases} \tag{69}$$

holds.

**Subcase** $\nu(s)$:  Claim holds because

$$[\![\nabla_A(s{\cdot}t)]\!] \ = [\![\nabla_A(s)\cdot t]\!] \cup [\![\nabla_A(t)]\!] \tag{70}$$

$$\overset{\text{IH}}{\subseteq} \{uv \mid u \in \bigcap_{a\in A} [\![\partial_a(s)]\!], v \in [\![t]\!]\} \cup \bigcap_{a\in A} [\![\partial_a(t)]\!] \tag{71}$$

$$\subseteq \bigcap_{a\in A} \{uv \mid [\![\partial_a(s)]\!], v \in [\![t]\!]\} \cup [\![\partial_a(t)]\!] \tag{72}$$

$$= \bigcap_{a\in A} [\![\partial_a(s)\cdot t]\!] \cup [\![\partial_a(t)]\!] \tag{73}$$

$$= \bigcap_{a\in A} \partial_a(s \cdot t) \tag{74}$$

and

$$\llbracket \Delta_A(s{\cdot}t) \rrbracket = \llbracket \Delta_A(s) \cdot t \rrbracket \cup \llbracket \Delta_A(t) \rrbracket \tag{75}$$

$$\overset{\text{IH}}{\supseteq} \{uv \mid u \in \bigcup_{a \in A} \llbracket \partial_a(s) \rrbracket, v \in \llbracket t \rrbracket\} \cup \bigcup_{a \in A} \llbracket \partial_a(t) \rrbracket \tag{76}$$

$$= \bigcup_{a \in A} \{uv \mid \llbracket \partial_a(s) \rrbracket, v \in \llbracket t \rrbracket\} \cup \llbracket \partial_a(t) \rrbracket \tag{77}$$

$$= \bigcup_{a \in A} \llbracket \partial_a(s) \cdot t \rrbracket \cup \llbracket \partial_a(t) \rrbracket \tag{78}$$

$$= \bigcup_{a \in A} \partial_a(s \cdot t) \tag{79}$$

**Subcase** $\neg\nu(s)$: Claim holds because

$$\llbracket \nabla_A(s{\cdot}t) \rrbracket = \llbracket \nabla_A(s) \cdot t \rrbracket \tag{80}$$

$$\overset{\text{IH}}{\subseteq} \{uv \mid u \in \bigcap_{a \in A} \llbracket \partial_a(s) \rrbracket, v \in \llbracket t \rrbracket\} \tag{81}$$

$$= \bigcap_{a \in A} \{uv \mid \llbracket \partial_a(s) \rrbracket, v \in \llbracket t \rrbracket\} \tag{82}$$

$$= \bigcap_{a \in A} \llbracket \partial_a(s) \cdot t \rrbracket \tag{83}$$

$$= \bigcap_{a \in A} \partial_a(s \cdot t) \tag{84}$$

and

$$\llbracket \Delta_A(s{\cdot}t) \rrbracket = \llbracket \Delta_A(s) \cdot t \rrbracket \tag{85}$$

$$\overset{\text{IH}}{\supseteq} \{uv \mid u \in \bigcup_{a \in A} \llbracket \partial_a(s) \rrbracket, v \in \llbracket t \rrbracket\} \tag{86}$$

$$= \bigcup_{a \in A} \{uv \mid \llbracket \partial_a(s) \rrbracket, v \in \llbracket t \rrbracket\} \tag{87}$$

$$= \bigcup_{a \in A} \llbracket \partial_a(s) \cdot t \rrbracket \tag{88}$$

$$= \bigcup_{a \in A} \partial_a(s \cdot t) \tag{89}$$

◀

## D Lemma 8: Properties of Join

Let $\mathfrak{L}_1$ and $\mathfrak{L}_2$ be non-empty sets of mutually disjoint literals.
1. $\bigcup(\mathfrak{L}_1 \bowtie \mathfrak{L}_2) = \bigcup \mathfrak{L}_1 \cup \bigcup \mathfrak{L}_2$.
2. $(\forall A \neq A' \in \mathfrak{L}_1 \bowtie \mathfrak{L}_2)\ A \sqcap A' = \emptyset$.
3. $(\forall A \in \mathfrak{L}_1 \bowtie \mathfrak{L}_2)\ (\forall A_i \in \mathfrak{L}_i)\ A \sqcap A_i \neq \emptyset \Rightarrow A \sqsubseteq A_i$.

**Proof of Lemma 8.**
1. Inclusion from left to right "$\subseteq$": Suppose that $a \in \bigcup(\mathfrak{L}_1 \bowtie \mathfrak{L}_2)$. Then there exists some $A_1 \in \mathfrak{L}_1$ and $A_2 \in \mathfrak{L}_2$ such that
   - $a \in A_1 \sqcap A_2$, but then $a \in A_1 \subseteq \bigcup \mathfrak{L}_1 \cup \bigcup \mathfrak{L}_2$;
   - $a \in A_1 \sqcap \overline{\bigsqcup \mathfrak{L}_2}$, but then $a \in A_1 \subseteq \bigcup \mathfrak{L}_1 \cup \bigcup \mathfrak{L}_2$; or
   - $a \in \overline{\bigsqcup \mathfrak{L}_1} \sqcap A_2$, but then $a \in A_2 \subseteq \bigcup \mathfrak{L}_1 \cup \bigcup \mathfrak{L}_2$.

   Inclusion from right to left "$\supseteq$": Suppose that $a \in \bigcup \mathfrak{L}_1 \cup \bigcup \mathfrak{L}_2$. There are three cases.
   - If there are $A_1 \in \mathfrak{L}_1$ such that $a \in A_1$ and $A_2 \in \mathfrak{L}_2$ such that $a \in A_2$, then $a \in A_1 \sqcap A_2 \in \mathfrak{L}_1 \bowtie \mathfrak{L}_2$.
   - If there is some $A_1 \in \mathfrak{L}_1$ such that $a \in A_1$ but there is no $A_2 \in \mathfrak{L}_2$ such that $a \in A_2$, then $a \in A_1 \sqcap \overline{\bigsqcup \mathfrak{L}_2} \in \mathfrak{L}_1 \bowtie \mathfrak{L}_2$.
   - Symmetric to previous case (exchange indices 1 and 2).
2. Suppose that $A_{01} = \overline{\bigsqcup \mathfrak{L}_1}$ and $A_{02} = \overline{\bigsqcup \mathfrak{L}_2}$. Clearly, $A_{01}$ is disjoint to any element of $\mathfrak{L}_1$ and $A_{02}$ is disjoint to any element of $\mathfrak{L}_2$. There are nine possible cases for $A$ and $A'$. To construct arbitrary elements of $\mathfrak{L}_1 \bowtie \mathfrak{L}_2$, we pick some $A_1, A_1' \in \mathfrak{L}_1$ and $A_2, A_2' \in \mathfrak{L}_2$.
   - $A = A_1 \sqcap A_2$ and $A' = A_1' \sqcap A_2'$. If $A \neq A'$, then $(A_1, A_2) \neq (A_1', A_2')$ and the claim follows from disjointness of $\mathfrak{L}_1$ and $\mathfrak{L}_2$.
   - $A = A_1 \sqcap A_2$ and $A' = A_1' \sqcap A_{02}$. The claim follows from $A_2 \sqcap A_{02} = \emptyset$.
   - $A = A_1 \sqcap A_2$ and $A' = A_{01} \sqcap A_2'$. The claim follows from $A_1 \sqcap A_{01} = \emptyset$.
   - $A = A_1 \sqcap A_{02}$ and $A' = A_1' \sqcap A_2'$. The claim follows from $A_2' \sqcap A_{02} = \emptyset$.
   - $A = A_1 \sqcap A_{02}$ and $A' = A_1' \sqcap A_{02}$. If $A \neq A'$, then $A_1 \neq A_1'$ and the claim follows from disjointness of $\mathfrak{L}_1$.
   - $A = A_1 \sqcap A_{02}$ and $A' = A_{01} \sqcap A_2'$. The claim follows from $A_2' \sqcap A_{02} = \emptyset$.
   - $A = A_{01} \sqcap A_2$ and $A' = A_1' \sqcap A_2'$. The claim follows from $A_{01} \sqcap A_1' = \emptyset$.
   - $A = A_{01} \sqcap A_2$ and $A' = A_1' \sqcap A_{02}$. The claim follows from $A_{01} \sqcap A_1' = \emptyset$.
   - $A = A_{01} \sqcap A_2$ and $A' = A_{01} \sqcap A_2'$. If $A \neq A'$, then $A_2 \neq A_2'$ and the claim follows from disjointness of $\mathfrak{L}_2$.
3. Immediate from the definition.

◀

## E   Lemma 9: Partial Equivalence

Let $\mathfrak{L} = \mathsf{next}(r)$.

**1.** $(\forall A \in \mathfrak{L})\ (\forall a, b \in A)\ \partial_a(r) = \partial_b(r)$

**2.** $(\forall a \notin \bigcup \mathfrak{L})\ \partial_a(r) \sqsubseteq \emptyset$

**Proof of Lemma 9.** We write $a \sim_{\mathfrak{L}} b$ if there exists some $A \in \mathfrak{L}$ such that $\{a, b\} \subseteq A$. The proof is by induction on $r$. The equality in item 1 has to be read as semantic equality. It is not necessarily syntactic.

**Cases $\epsilon$:** trivial.

**Case $A$:** In this case, $\mathfrak{L} = \{A\}$. By definition of the derivative: For each $a \in A$, $\partial_a(A) = \epsilon$. For each $b \notin A$, $\partial_b(A) = \emptyset$.

**Case $r+s$:** Let $\mathfrak{L}_r = \mathsf{next}(r)$, $\mathfrak{L}_s = \mathsf{next}(s)$, $\mathfrak{L} = \mathfrak{L}_r \bowtie \mathfrak{L}_s$, and $A \in \mathfrak{L}$.

There are three cases. If there exist $A_r \in \mathfrak{L}_r$ and $A_s \in \mathfrak{L}_s$ such that $A \sqsubseteq A_r$ and $A \sqsubseteq A_s$, then for all $a, b \in A$ it holds that $a \sim_{\mathfrak{L}_r} b$ and $a \sim_{\mathfrak{L}_s} b$ such that, by induction, $\partial_a(r) = \partial_b(r)$ and $\partial_a(s) = \partial_b(s)$. Hence, $\partial_a(r+s) = \partial_b(r+s)$ by definition of the derivative.

If there exist $A_r \in \mathfrak{L}_r$ such that $A \sqsubseteq A_r$, but for all $A_s \in \mathfrak{L}_s$ it is the case that $A \not\sqsubseteq A_s$, then for all $a, b \in A$ it holds that $a \sim_{\mathfrak{L}_r} b$ such that, by induction, $\partial_a(r) = \partial_b(r)$ and $\partial_a(s) \sqsubseteq \emptyset$ and $\partial_b(s) \sqsubseteq \emptyset$. Hence, $\partial_a(r+s) = \partial_b(r+s)$ by definition of the derivative.

If there exist $A_s \in \mathfrak{L}_s$ such that $A \sqsubseteq A_s$, but for all $A_r \in \mathfrak{L}_r$ it is the case that $A \not\sqsubseteq A_r$, then for all $a, b \in A$ it holds that $a \sim_{\mathfrak{L}_s} b$ such that, by induction, $\partial_a(s) = \partial_b(s)$ and $\partial_a(r) \sqsubseteq \emptyset$ and $\partial_b(r) \sqsubseteq \emptyset$. Hence, $\partial_a(r+s) = \partial_b(r+s)$ by definition of the derivative.

**Case $r \cdot s$:** Similar.

**Case $r^*$:** Let $\mathfrak{L} = \mathsf{next}(r)$ and $a \sim_{\mathfrak{L}} b$. Now $\partial_a(r^*) = \partial_a(r) \cdot r^* = \partial_b(r) \cdot r^* = \partial_b(r^*)$ where the middle equality holds by induction.

If $a \notin \bigcup \mathfrak{L}$, then $\partial_a(r) \sqsubseteq \emptyset$. Hence $\partial_a(r^*) = \partial_a(r) \cdot r^* \sqsubseteq \emptyset \cdot r^* \sqsubseteq \emptyset$.

**Case $r \& s$:** Let $\mathfrak{L}_r = \mathsf{next}(r)$, $\mathfrak{L}_s = \mathsf{next}(s)$, $\mathfrak{L} = \mathfrak{L}_r \sqcap \mathfrak{L}_s$, and $A \in \mathfrak{L}$.

By construction of $\mathfrak{L}$, there exist $A_r \in \mathfrak{L}_r$ and $A_s \in \mathfrak{L}_s$ such that $A \sqsubseteq A_r$ and $A \sqsubseteq A_s$. Thus, for all $a, b \in A$ it holds that $a \sim_{\mathfrak{L}_r} b$ and $a \sim_{\mathfrak{L}_s} b$ such that, by induction, $\partial_a(r) = \partial_b(r)$ and $\partial_a(s) = \partial_b(s)$. Hence, $\partial_a(r \& s) = \partial_b(r \& s)$ by definition of the derivative.

If $a \notin \bigcup \mathfrak{L}$, then assume that $a \notin \bigcup \mathfrak{L}_r$ (the case for $s$ is symmetric). By induction, $\partial_a(r) \sqsubseteq \emptyset$ so that $\partial_a(r \& s) = \partial_a(r) \& \partial_a(s) \sqsubseteq \emptyset \& \partial_a(s) \sqsubseteq \emptyset$.

**Case $!r$:** Let $\mathfrak{L}_r = \mathsf{next}(r)$ so that $\mathfrak{L} = \mathsf{next}(!r) = \mathfrak{L}_r \cup \{\bigsqcap \{\overline{A} \mid A \in \mathfrak{L}_r\}\}$. Clearly, $\bigcup \mathfrak{L} = \Sigma$.

If $a \sim_{\mathfrak{L}} b$, then there are two cases. If $a \sim_{\mathfrak{L}_r} b$, then $\partial_a(!r) = !\partial_a(r) = !\partial_b(r) = \partial_b(!r)$ by induction.

If $\{a, b\} \subseteq \bigsqcap \{\overline{A} \mid A \in \mathfrak{L}_r\}$, then $\{a, b\} \in \overline{\bigcup \mathfrak{L}_r}$ so that, by induction, $\partial_a(r) \sqsubseteq \emptyset$ and $\partial_b(r) \sqsubseteq \emptyset$. Hence, $!\partial_a(r) = !\partial_b(r)$. ◀

<span style="background-color:gold">**F**</span>    **Lemma 10: First and Next**

For all $r$, $\bigcup \mathsf{next}(r) \supseteq \mathsf{first}(r)$.

**Proof of Lemma 10.** The proof is by induction on $r$.

    **Cases** $\epsilon$, $A$: trivial.

    **Case** $r{+}s$: Let $\mathfrak{L}_r = \mathsf{next}(r)$, $\mathfrak{L}_s = \mathsf{next}(s)$, and $\mathfrak{L} = \mathfrak{L}_r \bowtie \mathfrak{L}_s$. By induction, $\bigcup \mathfrak{L}_r \supseteq$ $\mathsf{first}(r)$ and $\bigcup \mathfrak{L}_s \supseteq \mathsf{first}(s)$. By Lemma 8, $\bigcup \mathfrak{L} = \bigcup \mathfrak{L}_r \cup \bigcup \mathfrak{L}_s \supseteq \mathsf{first}(r) \cup \mathsf{first}(s) = \mathsf{first}(r{+}s)$.

    **Case** $r{\cdot}s$: Let $\mathfrak{L}_r = \mathsf{next}(r)$, $\mathfrak{L}_s = \mathsf{next}(s)$, and $\mathfrak{L} = \mathfrak{L}_r \bowtie \mathfrak{L}_s$.

    If $\neg \nu(r)$, then $\bigcup \mathsf{next}(r{\cdot}s) = \bigcup \mathsf{next}(r) \supseteq \mathsf{first}(r) = \mathsf{first}(r{\cdot}s)$.

    If $\nu(r)$, then $\bigcup \mathsf{next}(r{\cdot}s) = \bigcup (\mathsf{next}(r) \bowtie \mathsf{next}(s)) \supseteq (\mathsf{first}(r) \cup \mathsf{first}(s)) = \mathsf{first}(r{\cdot}s)$ by induction and using Lemma 8.

    **Case** $r^*$: $\bigcup \mathsf{next}(r^*) = \bigcup \mathsf{next}(r) \supseteq \mathsf{first}(r) = \mathsf{first}(r^*)$ by induction

    **Case** $r\&s$: $\bigcup \mathsf{next}(r\&s) = \bigcup (\mathsf{next}(r) \sqcap \mathsf{next}(s)) = \bigcup (\mathsf{next}(r)) \cap \bigcup (\mathsf{next}(s)) \supseteq \mathsf{first}(r) \cap$ $\mathsf{first}(s) \supseteq \mathsf{first}(r\&s)$.

    **Case** $!r$: $\bigcup \mathsf{next}(!r) = \Sigma \supseteq \mathsf{first}(!r)$. ◀

## G  Theorem 12: Left Quotient

▶ **Definition 24** (Next2). Let $\mathsf{next}^*(r) = \mathsf{next}(r) \setminus \{\emptyset\}$ be the set of first literals of ERE $r$ exlcuding the ehe empty set $\{\emptyset\}$.

**Proof of Theorem 12.** For any ERE $r$, for any literal $A \in \mathsf{next}^*(r)$, and for any symbol $a \in A$, the following equation holds:

$$[\![\nabla_A(r)]\!] = [\![\partial_a(r)]\!] \tag{90}$$

$$[\![\Delta_A(r)]\!] = [\![\partial_a(r)]\!] \tag{91}$$

Proof by induction on $r$.

**Case** $r = \epsilon$: Claim holds because $[\![\nabla_A(\epsilon)]\!] = [\![\Delta_A(\epsilon)]\!] = [\![\partial_a(\emptyset)]\!] = \emptyset$.
**Case** $r = B$: Claim holds because

$$[\![\nabla_A(B)]\!] = [\![\partial_a(B)]\!] = \begin{cases} \{\epsilon\}, & A \subseteq B \\ \emptyset, & otherwise \end{cases} \tag{92}$$

and

$$[\![\Delta_B(A)]\!] = [\![\partial_a(A)]\!] = \begin{cases} \{\epsilon\}, & A \subseteq B \\ \emptyset, & otherwise \end{cases} \tag{93}$$

**Case** $r = s^*$:  By induction

$$[\![\nabla_A(s)]\!] \overset{\mathsf{IH}}{=} [\![\partial_a(s)]\!] \tag{94}$$

$$[\![\Delta_A(s)]\!] \overset{\mathsf{IH}}{=} [\![\partial_a(s)]\!] \tag{95}$$

holds. We obtain that

$$\forall a : [\![\partial_a(s^*)]\!] = [\![\partial_a(s) \cdot s^*]\!] \tag{96}$$

$$\forall A : [\![\nabla_A(s^*)]\!] = [\![\nabla_A(s) \cdot s^*]\!] \tag{97}$$

$$\forall A : [\![\Delta_A(s^*)]\!] = [\![\Delta_A(s) \cdot s^*]\!] \tag{98}$$

holds. Claim holds because

$$\forall A : [\![\nabla_A(s^*)]\!] = [\![\nabla_A(s) \cdot s^*]\!] \tag{99}$$

$$\overset{\mathsf{IH}}{=} \{uv \mid u \in [\![\partial_a(s)]\!], v \in [\![s^*]\!]\} \tag{100}$$

$$= [\![\partial_a(s^*)]\!] \tag{101}$$

and

$$\forall A : [\![\Delta_A(s^*)]\!] = [\![\Delta_A(s) \cdot s^*]\!] \tag{102}$$

$$\overset{\mathsf{IH}}{=} \{uv \mid u \in [\![\partial_a(s)]\!], v \in [\![s^*]\!]\} \tag{103}$$

$$= [\![\partial_a(s^*)]\!] \tag{104}$$

**Case** $r = s+t$:  By induction

$$\llbracket \nabla_A(s) \rrbracket \overset{\text{IH}}{=} \llbracket \partial_a(s) \rrbracket \tag{105}$$

$$\llbracket \nabla_A(t) \rrbracket \overset{\text{IH}}{=} \llbracket \partial_a(t) \rrbracket \tag{106}$$

$$\llbracket \Delta_A(s) \rrbracket \overset{\text{IH}}{=} \llbracket \partial_a(s) \rrbracket \tag{107}$$

$$\llbracket \Delta_A(t) \rrbracket \overset{\text{IH}}{=} \llbracket \partial_a(t) \rrbracket \tag{108}$$

holds. We obtain that

$$\llbracket \partial_a(s+t) \rrbracket \;=\; \llbracket \partial_a(s) \rrbracket \cup \llbracket \partial_a(t) \rrbracket \tag{109}$$

$$\llbracket \nabla_A(s+t) \rrbracket \;=\; \llbracket \nabla_A(s) \rrbracket \cup \llbracket \nabla_A(t) \rrbracket \tag{110}$$

$$\llbracket \Delta_A(s+t) \rrbracket \;=\; \llbracket \Delta_A(s) \rrbracket \cup \llbracket \Delta_A(t) \rrbracket \tag{111}$$

holds. Claim holds because

$$\llbracket \nabla_A(s+t) \rrbracket \;=\; \llbracket \nabla_A(s) \rrbracket \cup \llbracket \nabla_A(t) \rrbracket \tag{112}$$

$$\overset{\text{IH}}{=} \llbracket \partial_a(s) \rrbracket \cup \llbracket \partial_a(t) \rrbracket \tag{113}$$

$$= \llbracket \partial_a(s+t) \rrbracket \tag{114}$$

and

$$\llbracket \Delta_A(s+t) \rrbracket \;=\; \llbracket \Delta_A(s) \rrbracket \cup \llbracket \Delta_A(t) \rrbracket \tag{115}$$

$$\overset{\text{IH}}{=} \llbracket \partial_a(s) \rrbracket \cup \llbracket \partial_a(t) \rrbracket \tag{116}$$

$$= \llbracket \partial_a(s+t) \rrbracket \tag{117}$$

**Case** $r = s\&t$:  By induction

$$\llbracket \nabla_A(s) \rrbracket \overset{\text{IH}}{=} \llbracket \partial_a(s) \rrbracket \tag{118}$$

$$\llbracket \nabla_A(t) \rrbracket \overset{\text{IH}}{=} \llbracket \partial_a(t) \rrbracket \tag{119}$$

$$\llbracket \Delta_A(s) \rrbracket \overset{\text{IH}}{=} \llbracket \partial_a(s) \rrbracket \tag{120}$$

$$\llbracket \Delta_A(t) \rrbracket \overset{\text{IH}}{=} \llbracket \partial_a(t) \rrbracket \tag{121}$$

holds. We obtain that

$$\llbracket \partial_a(s\&t) \rrbracket \;=\; \llbracket \partial_a(s) \rrbracket \cap \llbracket \partial_a(t) \rrbracket \tag{122}$$

$$\llbracket \nabla_A(s\&t) \rrbracket \;=\; \llbracket \nabla_A(s) \rrbracket \cap \llbracket \nabla_A(t) \rrbracket \tag{123}$$

$$\llbracket \Delta_A(s\&t) \rrbracket \;=\; \llbracket \Delta_A(s) \rrbracket \cap \llbracket \Delta_A(t) \rrbracket \tag{124}$$

holds. Claim holds because

$$\llbracket \nabla_A(s\&t) \rrbracket \;=\; \llbracket \nabla_A(s) \rrbracket \cap \llbracket \nabla_A(t) \rrbracket \tag{125}$$

$$\overset{\text{IH}}{=} \llbracket \partial_a(s) \rrbracket \cap \llbracket \partial_a(t) \rrbracket \tag{126}$$

$$= \llbracket \partial_a(s\&t) \rrbracket \tag{127}$$

and

$$\llbracket \Delta_A(s\&t) \rrbracket \;=\; \llbracket \Delta_A(s) \rrbracket \cap \llbracket \Delta_A(t) \rrbracket \tag{128}$$

$$\overset{\text{IH}}{=} \llbracket \partial_a(s) \rrbracket \cap \llbracket \partial_a(t) \rrbracket \tag{129}$$

$$= \llbracket \partial_a(s\&t) \rrbracket \tag{130}$$

**Case** $r = !s$: By induction

$$\llbracket \nabla_A(s) \rrbracket \overset{\text{IH}}{=} \llbracket \partial_a(s) \rrbracket \tag{131}$$

$$\llbracket \Delta_A(s) \rrbracket \overset{\text{IH}}{=} \llbracket \partial_a(s) \rrbracket \tag{132}$$

holds. We obtain that

$$\forall a: \ \llbracket \partial_a(!s) \rrbracket \ = \ \Sigma^* \backslash \llbracket \partial_a(s) \rrbracket \tag{133}$$

$$\forall A: \ \llbracket \nabla_A(!s) \rrbracket \ = \ \Sigma^* \backslash \llbracket \Delta_A(s) \rrbracket \tag{134}$$

$$\forall A: \ \llbracket \Delta_A(!s) \rrbracket \ = \ \Sigma^* \backslash \llbracket \nabla_A(s) \rrbracket \tag{135}$$

holds. Claim holds because

$$\forall A: \ \llbracket \nabla_A(!s) \rrbracket \ = \Sigma^* \backslash \llbracket \Delta_A(s) \rrbracket \tag{136}$$

$$\overset{\text{IH}}{=} \Sigma^* \backslash \llbracket \partial_a(s) \rrbracket \tag{137}$$

$$= \llbracket \partial_a(!s) \rrbracket \tag{138}$$

and

$$\forall A: \ \llbracket \Delta_A(!s) \rrbracket \ = \Sigma^* \backslash \llbracket \nabla_A(s) \rrbracket \tag{139}$$

$$\overset{\text{IH}}{=} \Sigma^* \backslash \llbracket \partial_a(s) \rrbracket \tag{140}$$

$$= \llbracket \partial_a(!s) \rrbracket \tag{141}$$

**Case** $r = s{\cdot}t$: By induction

$$\llbracket \nabla_A(s) \rrbracket \overset{\text{IH}}{=} \llbracket \partial_a(s) \rrbracket \tag{142}$$

$$\llbracket \nabla_A(t) \rrbracket \overset{\text{IH}}{=} \llbracket \partial_a(t) \rrbracket \tag{143}$$

$$\llbracket \Delta_A(s) \rrbracket \overset{\text{IH}}{=} \llbracket \partial_a(s) \rrbracket \tag{144}$$

$$\llbracket \Delta_A(t) \rrbracket \overset{\text{IH}}{=} \llbracket \partial_a(t) \rrbracket \tag{145}$$

holds. We obtain that

$$\forall a: \ \llbracket \partial_a(s{\cdot}t) \rrbracket \ = \ \begin{cases} \llbracket \partial_a(s) \cdot t \rrbracket \cup \llbracket \partial_a(t) \rrbracket, & \nu(s) \\ \llbracket \partial_a(s) \cdot t \rrbracket, & otherwise \end{cases} \tag{146}$$

$$\forall A: \ \llbracket \nabla_A(s{\cdot}t) \rrbracket \ = \ \begin{cases} \llbracket \nabla_A(s) \cdot t \rrbracket \cup \llbracket \nabla_A(t) \rrbracket, & \nu(s) \\ \llbracket \nabla_A(s) \cdot t \rrbracket, & otherwise \end{cases} \tag{147}$$

$$\forall A: \ \llbracket \Delta_A(s{\cdot}t) \rrbracket \ = \ \begin{cases} \llbracket \Delta_A(s) \cdot t \rrbracket \cup \llbracket \Delta_A(a)t \rrbracket, & \nu(s) \\ \llbracket \Delta_A(s) \cdot t \rrbracket, & otherwise \end{cases} \tag{148}$$

holds.

**Subcase** $\nu(s)$: Claim holds because

$$\llbracket \nabla_A(s{\cdot}t) \rrbracket \ = \llbracket \nabla_A(s) \cdot t \rrbracket \cup \llbracket \nabla_A(t) \rrbracket \tag{149}$$

$$\overset{\text{IH}}{=} \{uv \mid u \in \llbracket \partial_a(s) \rrbracket, v \in \llbracket t \rrbracket\} \cup \llbracket \partial_a(t) \rrbracket \tag{150}$$

$$= \llbracket \partial_a(s) \cdot t \rrbracket \cup \llbracket \partial_a(t) \rrbracket \tag{151}$$

$$= \partial_a(s \cdot t) \tag{152}$$

and

$$\llbracket \Delta_A(s{\cdot}t) \rrbracket \;=\; \llbracket \Delta_A(s) \cdot t \rrbracket \cup \llbracket \Delta_A(t) \rrbracket \tag{153}$$

$$\overset{\text{IH}}{=} \{uv \mid u \in \llbracket \partial_a(s) \rrbracket, v \in \llbracket t \rrbracket\} \cup \llbracket \partial_a(t) \rrbracket \tag{154}$$

$$= \llbracket \partial_a(s) \cdot t \rrbracket \cup \llbracket \partial_a(t) \rrbracket \tag{155}$$

$$= \partial_a(s \cdot t) \tag{156}$$

**Subcase** $\neg\nu(s)$:  Claim holds because

$$\llbracket \nabla_A(s{\cdot}t) \rrbracket \;=\; \llbracket \nabla_A(s) \cdot t \rrbracket \tag{157}$$

$$\overset{\text{IH}}{=} \{uv \mid u \in \llbracket \partial_a(s) \rrbracket, v \in \llbracket t \rrbracket\} \tag{158}$$

$$= \llbracket \partial_a(s) \cdot t \rrbracket \tag{159}$$

$$= \partial_a(s \cdot t) \tag{160}$$

and

$$\llbracket \Delta_A(s{\cdot}t) \rrbracket \;=\; \llbracket \Delta_A(s) \cdot t \rrbracket \tag{161}$$

$$\overset{\text{IH}}{=} \{uv \mid u \in \llbracket \partial_a(s) \rrbracket, v \in \llbracket t \rrbracket\} \tag{162}$$

$$= \llbracket \partial_a(s) \cdot t \rrbracket \tag{163}$$

$$= \partial_a(s \cdot t) \tag{164}$$

◀

## H     Theorem 15: Semantic Containment

▶ **Lemma 25** (Word Inclusion). *For all ERE $r$ and words $w$ in $\Sigma^*$,*

$$w \in [\![r]\!] \Leftrightarrow \nu(\partial_w(r))$$

**Proof of Lemma 25.** Proof by the definition of $\delta$ and $\nu$.     ◀

▶ **Lemma 26** (Word Containment). *For all ERE $r$ and $s$,*

$$r \sqsubseteq s \Leftrightarrow \nu(\partial_w(s)) \text{ for all } w \in [\![r]\!]$$

**Proof of Lemma 26.** An ERE $r$ is subset of another ERE $s$ iff for all words $w \in [\![r]\!]$ the derivation of $s$ w.r.t. word $w$ is nullable. For all $w \in \Sigma^*$ it holds that $w \in [\![s]\!]$ iff $\nu(\partial_w(s))$. It is trivial to see that

$$r \sqsubseteq s \tag{165}$$
$$\Leftrightarrow [\![r]\!] \subseteq [\![s]\!] \tag{166}$$
$$\Leftrightarrow \forall w \in [\![r]\!] : w \in [\![s]\!] \tag{167}$$
$$\Leftrightarrow \forall w \in [\![r]\!] : \nu(\partial_w(s)) \tag{168}$$

holds.     ◀

**Proof of Theorem 15.** For all regular expressions $r$ and $s$,

$$r \sqsubseteq s \Leftrightarrow (\nu(r) \Rightarrow \nu(s)) \ \wedge \ (\forall a \in \mathsf{first}(r)) \ \partial_a(r) \sqsubseteq \partial_a(s)$$

An ERE $r$ is subset of another ERE $s$ iff for all symbols $a$ in $\mathsf{first}(r)$ the derivation of $r$ w.r.t. symbol $a$ is subset of the derivation of $s$ w.r.t. $a$. We obtain that

$$[\![\partial_a(r)]\!] = a^{-1}[\![r]\!] \tag{169}$$

and this leads to

$$\{\epsilon \mid \nu(r)\} \cup \{au \mid a \in [\![\mathsf{first}(r)]\!], u \in a^{-1}[\![r]\!]\} = [\![r]\!] \tag{170}$$

Claim holds because

$$r \sqsubseteq s \tag{171}$$
$$\Leftrightarrow [\![r]\!] \subseteq [\![s]\!] \tag{172}$$
$$\Leftrightarrow \forall u \in [\![r]\!] : \ u \in [\![s]\!] \tag{173}$$
$$\Leftrightarrow \epsilon \in [\![r]\!] \Rightarrow \epsilon \in [\![s]\!] \ \wedge \ \forall a, u : \ au \in [\![r]\!] \Rightarrow \nu(\partial_{au}(s)) \tag{174}$$
$$\Leftrightarrow \nu(r) \Rightarrow \nu(s) \ \wedge \ \forall a \in \mathsf{first}(r), \forall u : \ au \in [\![r]\!] \Rightarrow \nu(\partial_u(\partial_a(s))) \tag{175}$$
$$\Leftrightarrow \nu(r) \Rightarrow \nu(s) \ \wedge \ \forall a \in \mathsf{first}(r), \forall u \in [\![\partial_a(r)]\!] : \ \nu(\partial_u(\partial_a(s))) \tag{176}$$
$$\Leftrightarrow \nu(r) \Rightarrow \nu(s) \ \wedge \ \forall a \in \mathsf{first}(r) : \ [\![\partial_a(r)]\!] \subseteq [\![\partial_a(s)]\!] \tag{177}$$
$$\Leftrightarrow \nu(r) \Rightarrow \nu(s) \ \wedge \ \forall a \in \mathsf{first}(r) : \ \partial_a(r) \sqsubseteq \partial_a(s) \tag{178}$$

    ◀

## I    Theorem 19: Symbolic Containment

**Proof of Theorem 19.** The proof is by contraposition. If $r \not\sqsubseteq s$ then $\exists A \in \mathsf{next}(r \sqsubseteq s) : \nabla_A(r) \not\sqsubseteq \nabla_A(s)$ or $\neg(\nu(r) \Rightarrow \nu(s))$.

We obtain that:

$$r \not\sqsubseteq s \;\Leftrightarrow\; [\![r]\!] \not\subseteq [\![s]\!] \tag{179}$$

$$\Leftrightarrow\; \exists u \in [\![r]\!]\backslash[\![s]\!] \tag{180}$$

**Case** $u = \epsilon$:

Claim holds because $\neg(\nu(r) \Rightarrow \nu(s))$.

**Case** $u \neq \epsilon$:

It must be that $u = av$ with $a \in \mathsf{first}(r) = \mathsf{next}(r)$. Therefore $\exists A \in \mathsf{next}(r) : a \in A$.

**Subcase** $a \notin \mathsf{first}(s)$:

Claim holds by Lemma 12 and Lemma 27 because $\exists A \in \mathsf{next}(r) : \nabla_A(r) \neq \emptyset$ and $\nabla_A(s) = \emptyset$ implies that $\nabla_A(r) \not\sqsubseteq \nabla_A(s)$.

**Subcase** $a \in \mathsf{first}(s)$:

By Lemma 12 and Lemma 27 claim holds because $v \in [\![\partial_a(r)]\!]\backslash[\![\partial_a(s)]\!]$ implies that $v \in [\![\nabla_A(r)]\!]\backslash[\![\nabla_A(s)]\!]$

◀

## J Theorem 21: Soundness

For all regular expression $r$ and $s$:

$$\emptyset \vdash r \mathbin{\dot{\sqsubseteq}} s \,:\, \top \;\Leftrightarrow\; r \sqsubseteq s$$

**Proof.** We prove that $\Gamma \vdash r \mathbin{\dot{\sqsubseteq}} s \,:\, \bot$ iff $r \not\sqsubseteq s$, for all contexts $\Gamma$ where $r \mathbin{\dot{\sqsubseteq}} s \notin \Gamma$. This is sufficient because each regular inequality gives rise to a finite derivation by Theorem 20.

The "only-if" direction is by rule induction on the derivation of $\Gamma \vdash r \mathbin{\dot{\sqsubseteq}} s \,:\, \bot$.

- Suppose the last rule is (DISPROVE). By inversion, $\nu(r)$ and $\neg\nu(s)$ so that $r \not\sqsubseteq s$.
- Suppose the last rule is (UNFOLD-FALSE). By inversion,

$$r \mathbin{\dot{\sqsubseteq}} s \notin \Gamma \tag{181}$$

$$\nu(r) \Rightarrow \nu(s) \tag{182}$$

$$\exists A \in \mathsf{next}(r \mathbin{\dot{\sqsubseteq}} s) : \; \Gamma \cup \{r \mathbin{\dot{\sqsubseteq}} s\} \;\vdash\; \partial_A(r) \mathbin{\dot{\sqsubseteq}} \partial_A(s) \,:\, \bot \tag{183}$$

By induction, we obtain that $\partial_A(r) \not\sqsubseteq \partial_A(s)$, for some $A \in \mathsf{next}(r \mathbin{\dot{\sqsubseteq}} s)$. By Theorem 19, we obtain $r \not\sqsubseteq s$.

For the "if" direction, the assumption that $r \not\sqsubseteq s$ implies that $[\![r]\!] \backslash [\![s]\!] \neq \emptyset$. Let $u \in [\![r]\!] \backslash [\![s]\!]$ a word of shortest length. We continue by induction on $u$. If $u = \epsilon$, then $\nu(r)$ but not $\nu(s)$ must hold. By our assumption on $\Gamma$, it cannot be that $r \mathbin{\dot{\sqsubseteq}} s \notin \Gamma$. By rule (DISPROVE), $\Gamma \vdash r \mathbin{\dot{\sqsubseteq}} s \,:\, \bot$.

If $u = au'$, then there exists $A_a \in \mathsf{next}(r \mathbin{\dot{\sqsubseteq}} s)$ such that $a \in A_a$ (by Lemma 10). It must be that $\nu(r) \Rightarrow \nu(s)$: otherwise, we get a contradiction against the minimality of $u$'s length. By Theorem 19 it must be that $\partial_{A_a}(r) \not\sqsubseteq \partial_{A_a}(s)$ so that induction yields $\Gamma \cup \{r \mathbin{\dot{\sqsubseteq}} s\} \vdash \partial_{A_a}(r) \mathbin{\dot{\sqsubseteq}} \partial_{A_a}(s) \,:\, \bot$. Applying rule (UNFOLD-FALSE) yields $\Gamma \vdash r \mathbin{\dot{\sqsubseteq}} s \,:\, \bot$. ◀

### K     Lemma 27: Coverage

▶ **Lemma 27** (Coverage)**.** *For all symbols $a \in \Sigma$, words $u \in \Sigma^*$, and EREs on $\Sigma$ it holds that:*

$$u \in [\![\partial_a(r)]\!] \;\Leftrightarrow\; \exists A \in \mathsf{next}^*(r) : \; u \in [\![\Delta_A(r)]\!]$$
$$u \in [\![\partial_a(r)]\!] \;\Leftrightarrow\; \exists A \in \mathsf{next}^*(r) : \; u \in [\![\nabla_A(r)]\!]$$

**Proof of Lemma 27.** Suppose $[\![\partial_a(r)]\!] \neq \emptyset$. Because $\nabla_A(r) = \Delta_A(r)$ for all $A \in \mathsf{next}^*(r)$ show $\exists A \in \mathsf{next}^*(r) : \; w \in [\![\nabla_A(r)]\!]$. Proof by induction on $r$.

**Case** $r = \epsilon$, $\mathsf{next}^*(r) = \emptyset$: Contradicts assumption.

**Case** $r = A$, $\mathsf{next}^*(r) = \{A\}$:
 We obtain that $a \in A \;\Rightarrow\; \partial_a(A) = \epsilon$. Claim holds because $\mathsf{next}^*(r) = \{A\}$, $\nabla_A(A) = \epsilon$, and thus $w = \epsilon$ and $\epsilon \in [\![\epsilon]\!]$.

**Case** $r = s^*$, $\mathsf{next}^*(r) = \mathsf{next}^*(s)$:
 We obtain that $w \in [\![\partial_a(s^*)]\!] = [\![\partial_a(s) \cdot s^*]\!] \neq \emptyset$. By induction $\exists A' \in \mathsf{next}^*(s) : \; u \in [\![\nabla_{A'}(s)]\!]$. The chain holds because $\mathsf{next}^*(s^*) = \mathsf{next}^*(s)$ and $\nabla_A(s^*) = \nabla_A(s) \cdot s*$ and $u \in [\![\nabla_A(s)]\!], v \in [\![\nabla_A(s^*)]\!]$ implies $w = u \cdot v \in [\![\nabla_A(s^*)]\!]$.

**Case** $r = (s + t)$, $\mathsf{next}^*(r) = \mathsf{next}^*(s) \bowtie \mathsf{next}^*(t)$:
 We obtain that $w \in [\![\partial_a(s + t)]\!] = [\![\partial_a(s)]\!] \cup [\![\partial_a(t)]\!] \neq \emptyset$. By induction $\exists A' \in \mathsf{next}^*(s) : \; u \in [\![\nabla_{A'}(s)]\!]$ and $\exists A'' \in \mathsf{next}^*(t) : \; v \in [\![\nabla_{A''}(t)]\!]$. The chain holds because $\mathsf{next}^*(s + t) = \mathsf{next}^*(s) \bowtie \mathsf{next}^*(t)$ and $\nabla_A(s + t) = \nabla_A(s) + \nabla_A(t)$ and $w \in [\![\nabla_A(s)]\!]$ or $w \in [\![\nabla_A(t)]\!]$ implies $w \in [\![\nabla_A(s + t)]\!]$.

**Case** $r = (s \& t)$, $\mathsf{next}^*(r) = \{A' \sqcap A'' \mid A' \in \mathsf{next}^*(s), A'' \in \mathsf{next}^*(t)\}$:
 We obtain that $w \in [\![\partial_a(s \& t)]\!] = [\![\partial_a(s)]\!] \cap [\![\partial_a(t)]\!]$ implies $w \in [\![\partial_a(s)]\!]$ and $w \in [\![\partial_a(t)]\!]$. By induction $\exists A' \in \mathsf{next}^*(s) : \; w \in [\![\nabla_A(s)]\!]$ and $\exists A'' \in \mathsf{next}^*(t) : \; w \in [\![\nabla_{A''}(t)]\!]$. Let $A = A' \sqcap A'' \in \mathsf{next}^*(s \& t)$. If $a \in A'$ and $a \in A'$ then $a \in A$. The chain holds because $\mathsf{next}^*(s \& t) = \{A' \sqcap A'' \mid A' \in \mathsf{next}^*(s), A'' \in \mathsf{next}^*(t)\}$ and $\nabla_A(s \& t) = \nabla_A(s) \& \nabla_A(t)$, and $w \in [\![\nabla_A(s)]\!]$ and $w \in [\![\nabla_A(t)]\!]$ implies $w \in [\![\nabla_A(s \& t)]\!]$.

**Case** $r = (! s)$, $\mathsf{next}^*(r) = \bigsqcap\{\overline{A} \mid A \in \mathsf{next}^*(s)\} \;\bowtie\; \{A \in \mathsf{next}^*(s) \mid \nabla_A(s) \neq \Sigma^*\}$:
 We obtain that $w \in [\![\partial_a(! s)]\!] = \Sigma^* \setminus [\![\partial_a(s)]\!]$ implies $w \notin [\![\partial_a(s)]\!]$. By induction $\exists A' \in \mathsf{next}^*(s) : \; w \in [\![\nabla_A(s)]\!]$. Let $A = \bigsqcap\{\overline{A} \mid A \in \mathsf{next}^*(s)\} \;\bowtie\; \{A \in \mathsf{next}^*(s) \mid \nabla_A(s) \neq \Sigma^*\}$. If $[\![\partial_a(! s)]\!] \neq \emptyset$ implies $[\![\partial_a(s)]\!] \neq \Sigma^*$. The chain holds because $\mathsf{next}^*(! s) = \bigsqcap\{\overline{A} \mid A \in \mathsf{next}^*(s)\} \;\bowtie\; \{A \in \mathsf{next}^*(s) \mid \nabla_A(s) \neq \Sigma^*\}$ and $\nabla_A(! s) = ! \nabla_A(s)$, and $w \notin [\![\nabla_A(s)]\!]$ implies $w \in [\![\nabla_A(! s)]\!]$.

**Case** $r = (s \cdot t)$:
 **Subcase** $\nu(s)$, $\mathsf{next}^*(r) = \mathsf{next}^*(s) \bowtie \mathsf{next}^*(t)$:
  We obtain that $w \in [\![\partial_a(s \cdot t)]\!] = [\![\partial_a(s) \cdot t]\!] \cup [\![\partial_a(t)]\!]$ implies $w \in [\![\partial_a(s) \cdot t]\!]$ or $w \in [\![\partial_a(t)]\!]$. By induction $\exists A' \in \mathsf{next}^*(s) : \; u \in [\![\nabla_A(s)]\!]$ and $\exists A'' \in \mathsf{next}^*(t) : \; v \in [\![\nabla_{A''}(t)]\!]$. The chain holds because $\mathsf{next}^*(s \cdot t) = \mathsf{next}^*(s) \bowtie \mathsf{next}^*(t)$ and $\nabla_A(s \cdot t) = (\nabla_A(s) \cdot t) + \nabla_A(t)$, and $u \in [\![\nabla_A(s)]\!]$ and $v \in [\![t]\!]$ implies $w = u \cdot v \in [\![\nabla_A(s \cdot t)]\!]$ or $w = \epsilon \cdot v \in [\![\nabla_A(s \cdot t)]\!]$.
 **Subcase** $\neg\nu(s)$, $\mathsf{next}^*(r) = \mathsf{next}^*(s)$:
  We obtain that $w \in [\![\partial_a(s \cdot t)]\!] = [\![\partial_a(s) \cdot t]\!]$ implies $w \in [\![\partial_a(s) \cdot t]\!]$. By induction $\exists A' \in \mathsf{next}^*(s) : \; u \in [\![\nabla_A(s)]\!]$. The chain holds because $\mathsf{next}^*(s \cdot t) = \mathsf{next}^*(s)$ and $\nabla_A(s \cdot t) = \nabla_A(s) \cdot t$, and $u \in [\![\nabla_A(s)]\!]$ and $v \in [\![\nabla_A(t)]\!]$ implies $w = u \cdot v \in [\![\nabla_A(s \cdot t)]\!]$.

◀

## L Lemma 28: Equivalence

▶ **Lemma 28** (Equivalence)**.** *For all ERE $r$, literals $A \in \mathsf{next}(r)$, and literals $A'$ with $A' \subseteq A \mid A' \neq \emptyset$ holds:*

$$\llbracket \Delta_A(r) \rrbracket \iff \llbracket \Delta_{A'}(r) \rrbracket$$
$$\llbracket \nabla_A(r) \rrbracket \iff \llbracket \nabla_{A'}(r) \rrbracket$$

**Proof of Lemma 27.** Suppose $\mathsf{next}(r) \neq \{\emptyset\}$. Because $\nabla_A(r) = \Delta_A(r)$ for all $A \in \mathsf{next}(r)$ show $\llbracket \nabla_A(r) \rrbracket = \llbracket \nabla_{A'}(r) \rrbracket$. Proof by induction on $r$.

**Case** $r = \epsilon$, $\mathsf{next}(r) = \{\emptyset\}$: Contradicts assumption.

**Case** $r = A$, $\mathsf{next}(r) = \{A\}$:

Claim holds because for all $A' \subseteq A \Rightarrow \nabla_{A'}(r) = \nabla_A(r) = \epsilon$ and thus $\llbracket \nabla_A(r) \rrbracket = \llbracket \nabla_{A'}(r) \rrbracket$.

**Case** $r = s^*$, $\mathsf{next}(r) = \mathsf{next}(s)$:

We obtain that $\llbracket \nabla_A(s^*) \rrbracket = \llbracket \nabla_A(s) \cdot s^* \rrbracket \neq \emptyset$. By induction $\forall A_s \in \mathsf{next}(s), A_s' \subset A_s$ : $\llbracket \nabla_{A_s}(s) \rrbracket = \llbracket \nabla_{A_s'}(s) \rrbracket$. The chain holds because $\mathsf{next}(s^*) = \mathsf{next}(s)$ and $\llbracket \nabla_{A'}(s^*) \rrbracket = \llbracket \nabla_{A'}(s) \cdot s* \rrbracket$.

**Case** $r = (s + t)$, $\mathsf{next}(r) = \mathsf{next}(s) \bowtie \mathsf{next}(t)$:

We obtain that $\llbracket \nabla_A(s+t) \rrbracket = \llbracket \nabla_A(s) \rrbracket \cup \llbracket \nabla_A(t) \rrbracket \neq \emptyset$. By induction $\forall A_s \in \mathsf{next}(s), A_s' \subset A_s : \llbracket \nabla_{A_s}(s) \rrbracket = \llbracket \nabla_{A_s'}(s) \rrbracket$ and $\forall A_t \in \mathsf{next}(t), A_t' \subset A_t : \llbracket \nabla_{A_t}(t) \rrbracket = \llbracket \nabla_{A_t'}(t) \rrbracket$. The chain holds because $\mathsf{next}(s + t) = \mathsf{next}(s) \bowtie \mathsf{next}(t)$ and $\forall A'' \in \mathsf{next}(r) \cup \mathsf{next}(s) : \exists A''' \in \mathsf{next}(s + t) : A''' \subseteq A''$ and $\llbracket \nabla_{A'}(s+t) \rrbracket = \llbracket \nabla_{A'}(s) \rrbracket \cup \llbracket \nabla_{A'}(t) \rrbracket$.

**Case** $r = (s \& t)$, $\mathsf{next}(r) = \{A' \sqcap A'' \mid A' \in \mathsf{next}(s), A'' \in \mathsf{next}(t)\}$:

We obtain that $\llbracket \nabla_A(s\&t) \rrbracket = \llbracket \nabla_A(s) \rrbracket \cap \llbracket \nabla_A(t) \rrbracket \neq \emptyset$. By induction $\forall A_s \in \mathsf{next}(s), A_s' \subset A_s : \llbracket \nabla_{A_s}(s) \rrbracket = \llbracket \nabla_{A_s'}(s) \rrbracket$ and $\forall A_t \in \mathsf{next}(t), A_t' \subset A_t : \llbracket \nabla_{A_t}(t) \rrbracket = \llbracket \nabla_{A_t'}(t) \rrbracket$. Let $A = A_s \sqcap A_t \in \mathsf{next}(s\&t)$. If $A' \subseteq A$ then $A' \subseteq A_s$ and $A' \subseteq A_t$. The chain holds because $\mathsf{next}(s\&t) = \mathsf{next}(s) \sqcap \mathsf{next}(t)$ and $\llbracket \nabla_{A'}(s\&t) \rrbracket = \llbracket \nabla_{A'}(s) \rrbracket \cap \llbracket \nabla_{A'}(t) \rrbracket$.

**Case** $r = (!s)$, $\mathsf{next}(r) = \bigsqcap\{\overline{A} \mid A \in \mathsf{next}(s)\} \cup \{A \in \mathsf{next}(s)\}$:

We obtain that $\llbracket \nabla_A(!s) \rrbracket = \Sigma^* \setminus \llbracket \nabla_A(s) \rrbracket$. By induction $\forall A_s \in \mathsf{next}(s), A_s' \subset A_s : \llbracket \nabla_{A_s}(s) \rrbracket = \llbracket \nabla_{A_s'}(s) \rrbracket$. The chain holds because $\mathsf{next}(!s) = \bigsqcap\{\overline{A} \mid A \in \mathsf{next}(s)\} \cup \{A \in \mathsf{next}(s)\}$ and $\llbracket \nabla_{A'}(!s) \rrbracket = \Sigma^* \setminus \llbracket \nabla_{A'}(s) \rrbracket$ and for all $A'' \notin \mathsf{next}(s)$ $\nabla_{A''}(s) = \emptyset$.

**Case** $r = (s \cdot t)$:

**Subcase** $\nu(s)$, $\mathsf{next}(r) = \mathsf{next}(s) \bowtie \mathsf{next}(t)$:

We obtain that $\llbracket \nabla_A(s \cdot t) \rrbracket = \llbracket \nabla_A(s) \cdot t \rrbracket \cup \llbracket \nabla_A(t) \rrbracket$. By induction $\forall A_s \in \mathsf{next}(s), A_s' \subset A_s : \llbracket \nabla_{A_s}(s) \rrbracket = \llbracket \nabla_{A_s'}(s) \rrbracket$ and $\forall A_t \in \mathsf{next}(t), A_t' \subset A_t : \llbracket \nabla_{A_t}(t) \rrbracket = \llbracket \nabla_{A_t'}(t) \rrbracket$. The chain holds because $\mathsf{next}(s + t) = \mathsf{next}(s) \bowtie \mathsf{next}(t)$ and $\forall A'' \in \mathsf{next}(r) \cup \mathsf{next}(s) : \exists A''' \in \mathsf{next}(s + t) : A''' \subseteq A''$ and $\llbracket \nabla_{A'}(s \cdot t) \rrbracket = \llbracket \nabla_{A'}(s) \cdot t \rrbracket \cup \llbracket \nabla_{A'}(t) \rrbracket$.

**Subcase** $\neg\nu(s)$, $\mathsf{next}(r) = \mathsf{next}(s)$:

We obtain that $\llbracket \nabla_A(s \cdot t) \rrbracket = \llbracket \nabla_A(s) \cdot t \rrbracket$. By induction $\forall A_s \in \mathsf{next}(s), A_s' \subset A_s : \llbracket \nabla_{A_s}(s) \rrbracket = \llbracket \nabla_{A_s'}(s) \rrbracket$. The chain holds because $\mathsf{next}(s \cdot t) = \mathsf{next}(s)$ and $\llbracket \nabla_{A'}(s \cdot t) \rrbracket = \llbracket \nabla_{A'}(s) \cdot t \rrbracket$.

◀

## M   Containment Example

▶ **Example 29** (Containment). Consider the regular expressions $r = ((a+b)+c)$ and $s = (a+b)$ and the inequality $r \sqsubseteq s$ which is obviously invalid. The computation of one derivation step is as follows:

$$\nabla_A(r) \sqsubseteq \nabla_A(s) \;\Leftrightarrow\; \nabla_A(((a+b)+c)) \sqsubseteq \nabla_A((a+b)) \tag{184}$$
$$\Leftrightarrow\; (\nabla_A((a+b))+\nabla_A(c)) \sqsubseteq (\nabla_A(a)+\nabla_A(b)) \tag{185}$$
$$\Leftrightarrow\; ((\nabla_A(a)+\nabla_A(b))+\nabla_A(c)) \sqsubseteq (\nabla_A(a)+\nabla_A(b)) \tag{186}$$

To solve the inequality $r \sqsubseteq s$ the inequality gets derived in respect to the next literals of $r \sqsubseteq s$. The calculation of $\mathsf{next}(r \sqsubseteq s)$ is split into several sub-calculation concerning to the calculation of $\mathsf{next}$.

$$\mathsf{next}(r) \;\Leftrightarrow\; \mathsf{next}((a+b)+c) \tag{187}$$
$$\Leftrightarrow\; \mathsf{next}(a+b) \bowtie \mathsf{next}(c) \tag{188}$$
$$\Leftrightarrow\; (\mathsf{next}(a) \bowtie \mathsf{next}(b)) \bowtie \mathsf{next}(c) \tag{189}$$
$$\Leftrightarrow\; (\{a\} \bowtie \{b\}) \bowtie \{c\} \tag{190}$$
$$\Leftrightarrow\; \{a,b,c\} \tag{191}$$

$$\mathsf{next}(s) \;\Leftrightarrow\; \mathsf{next}(a+b) \tag{192}$$
$$\Leftrightarrow\; \mathsf{next}(a) \bowtie \mathsf{next}(b) \tag{193}$$
$$\Leftrightarrow\; \{a\} \bowtie \{b\} \tag{194}$$
$$\Leftrightarrow\; \{a,b\} \tag{195}$$

$$\mathsf{next}(!s) \;\Leftrightarrow\; \textstyle\prod\{\overline{A} \mid A \in \mathsf{next}(s)\} \;\bowtie\; \{A \in \mathsf{next}(r)\} \tag{196}$$
$$\Leftrightarrow\; \textstyle\prod\{\overline{a},\overline{b}\} \;\bowtie\; \{a,b\} \tag{197}$$
$$\Leftrightarrow\; \{\overline{\{a,b\}}\} \;\bowtie\; \{a,b\} \tag{198}$$
$$\Leftrightarrow\; \{\overline{\{a,b\}},a,b\} \tag{199}$$

$$\mathsf{next}(r \sqsubseteq s) \;\Leftrightarrow\; \mathsf{next}(r \& !s) \tag{200}$$
$$\Leftrightarrow\; \{A \sqcap A' \mid A \in \mathsf{next}(r), A' \in \mathsf{next}(!s)\} \tag{201}$$
$$\Leftrightarrow\; \{A \sqcap A' \mid A \in \{a,b,c\}, A' \in \{\overline{\{a,b\}},a,b\}\} \tag{202}$$
$$\Leftrightarrow\; \{a,b,c\} \tag{203}$$

Finally, the inequality gets derived in respect to the next literals.

$$\forall A \in \mathsf{next}(r):\; \nabla_A(r) \sqsubseteq \nabla_A(s) \mid \mathsf{next}(r) = \{a,b,c\} \tag{204}$$

This results in three iterations:

$$\nabla_a(r) \sqsubseteq \nabla_a(s) \;\Leftrightarrow\; \nabla_a(((a{+}b){+}c)) \sqsubseteq \nabla_a((a{+}b)) \tag{205}$$

$$\Leftrightarrow\; (\nabla_a((a{+}b)){+}\nabla_a(c)) \sqsubseteq (\nabla_a(a){+}\nabla_a(b)) \tag{206}$$

$$\Leftrightarrow\; ((\nabla_a(a){+}\nabla_a(b)){+}\nabla_a(c)) \sqsubseteq (\nabla_a(a){+}\nabla_a(b)) \tag{207}$$

$$\Leftrightarrow\; ((\epsilon{+}\emptyset){+}\emptyset) \sqsubseteq (\epsilon{+}\emptyset) \tag{208}$$

$$\Leftrightarrow\; (\epsilon{+}\emptyset) \sqsubseteq \epsilon \tag{209}$$

$$\Leftrightarrow\; \epsilon \sqsubseteq \epsilon \tag{210}$$

$$\nabla_b(r) \sqsubseteq \nabla_b(s) \;\Leftrightarrow\; \nabla_b(((a{+}b){+}c)) \sqsubseteq \nabla_b((a{+}b)) \tag{211}$$

$$\Leftrightarrow\; (\nabla_b((a{+}b)){+}\nabla_b(c)) \sqsubseteq (\nabla_b(a){+}\nabla_b(b)) \tag{212}$$

$$\Leftrightarrow\; ((\nabla_b(a){+}\nabla_b(b)){+}\nabla_b(c)) \sqsubseteq (\nabla_b(a){+}\nabla_b(b)) \tag{213}$$

$$\Leftrightarrow\; ((\emptyset{+}\epsilon){+}\emptyset) \sqsubseteq (\emptyset{+}\epsilon) \tag{214}$$

$$\Leftrightarrow\; (\epsilon{+}\emptyset) \sqsubseteq \epsilon \tag{215}$$

$$\Leftrightarrow\; \epsilon \sqsubseteq \epsilon \tag{216}$$

$$\nabla_c(r) \sqsubseteq \nabla_c(s) \;\Leftrightarrow\; \nabla_c(((a{+}b){+}c)) \sqsubseteq \nabla_c((a{+}b)) \tag{217}$$

$$\Leftrightarrow\; (\nabla_c((a{+}b)){+}\nabla_c(c)) \sqsubseteq (\nabla_c(a){+}\nabla_c(b)) \tag{218}$$

$$\Leftrightarrow\; ((\nabla_c(a){+}\nabla_c(b)){+}\nabla_c(c)) \sqsubseteq (\nabla_c(a){+}\nabla_c(b)) \tag{219}$$

$$\Leftrightarrow\; ((\nabla_c(a){+}\nabla_c(b)){+}\nabla_c(c)) \sqsubseteq (\nabla_c(a){+}\nabla_c(b)) \tag{220}$$

$$\Leftrightarrow\; ((\emptyset{+}\emptyset){+}\epsilon) \sqsubseteq (\emptyset{+}\emptyset) \tag{221}$$

$$\Leftrightarrow\; (\emptyset{+}\epsilon) \sqsubseteq \emptyset \tag{222}$$

$$\Leftrightarrow\; \epsilon \sqsubseteq \emptyset \tag{223}$$