# Derivatives for Enhanced Regular Expressions

Peter Thiemann

University of Freiburg

**Abstract.** Regular languages are closed under a wealth of formal language operators. Incorporating such operators in regular expressions leads to concise language specifications, but the transformation of such enhanced regular expressions to finite automata becomes more involved. We present an approach that enables the direct construction of finite automata from regular expressions enhanced with further operators that preserve regularity. Our construction is based on an extension of the theory of derivatives for regular expressions. To retain the standard results about derivatives, we develop a derivability criterion for the compatibility of the extra operators with derivatives.

Some derivable operators do not preserve regularity. Derivatives provide a decision procedure for the word problem of regular expressions enhanced with such operators.

**Keywords:** automata and logic, regular languages, derivatives

## 1 Introduction

Brzozowski derivatives [4] and Antimirov's partial derivatives [2] are well-known tools to transform regular expressions to automata and to define algorithms for equivalence and containment on them [1,10]. Brzozowski's automaton construction relies on the finiteness of the set of iterated derivatives when considered up to similarity (commutativity, associativity, and idempotence for union). Derivatives had quite some impact on the study of algorithms for regular languages on finite words and trees [13,5].

While derivative-based algorithms have been deprecated for performance reasons [17], there has been renewed interest in the study of derivatives and partial derivatives. On the practical side, Owens and coworkers [12] report a functional implementation that revives many features. Might and coworkers [11] implement parsing for context-free languages using derivatives.

A common theme on the theory side is the study of derivative structures for enhancements of regular expressions. While Brzozowski's original work covered extended regular expressions, partial derivatives were originally limited to simple expressions without intersection and complement. It is a significant effort to define partial derivatives for extended regular expressions [5].

Derivatives have also been used to study various shuffle operators for applications in modeling concurrent programs [14]. Later extensions consider forkable expressions with a new operator that abstracts process creation [15].

Caron and coworkers [6] study derivatives for multi-tilde-bar expressions. The tilde (bar) operator adds (removes) $\varepsilon$ from a language. Multi-tilde-bar applies to a list of languages and (roughly) defines a selective concatenation operation that can be configured to include or exclude certain languages of the list.

Champarnaud and coworkers [8] consider derivatives of approximate regular expressions (ARE). AREs extend regular expressions with a family of unary operators $\mathbb{F}_k$, for $k \in \mathbb{N}$, which enhance their argument language $L$ with all words $u$ such that $d(u, w) \leq k$, for some word $w \in L$. Here, $d$ is a suitable distance function, for example, Hamming distance or Levenshtein distance.

Traytel and Nipkow [16] obtain decision procedures for MSO using a suitably defined derivative operation on regular expressions with a projection operation.

The general framework of Caron and coworkers [7] generalizes the syntactic structure underlying a derivative construction to a *support*. A support generalizes expressions (for constructing Brzozowski derivatives), sets of expressions (for Antimirov's partial derivatives), and sets of clausal forms over sets of regular expressions, and thus yields an encompassing framework in which different kinds of derivative constructions can be formalized and compared. The authors give a sufficient criterion for a support to generate a finite number of iterated derivatives from a given expression along with automata constructions for deterministic, nondeterministic, and alternating finite automata. Their work applies to extended regular expressions with arbitrary boolean functions.

Champernaud and coworkers [9] propose constrained regular expressions with a notion of comprehension (filtering by a predicate) and matching. The resulting languages are in general not regular, but the authors propose expression derivation for the membership test and study its decidability.

## Contributions

In this work, we identify a pattern in the definition of (standard) derivatives for enhancements of regular expressions that go beyond boolean functions. Concretely, we consider regular expressions enhanced with further operators on languages (e.g., shuffle, homomorphism, approximation). Then we propose *left derivability* and *$\varepsilon$-testability* as a sufficient condition for the set of operators such that a syntactic derivative operation is definable for enhanced expressions. This condition gives rise to a decision procedure for the membership test for enhanced expressions via expression derivation.

A refinement, *linear left derivability*, is a sufficient condition to guarantee finiteness of the set of dissimilar derivatives of an enhanced expression. The finiteness condition enables the direct construction of a deterministic finite automaton. We show that every linear left derivable operator can be defined by a rational finite state transducer and thus preserves regularity.

Proofs and further examples may be found in the appendix.

## 2 Preliminaries

We write $\mathbb{N}$ for the set of natural numbers, $\mathbb{B} = \{\mathbf{0}, \mathbf{1}\}$ for the set of booleans, and $X \uplus Y$ for the disjoint union of sets $X$ and $Y$. We sometimes write $(\overline{E_k}^{k=1,\dots,n})$ for the tuple $(E_1, \dots, E_n)$ where $E_k$ is some entity depending on $k$.

An alphabet $\Sigma$ is a finite set of symbols. The set $\Sigma^*$ denotes the set of finite words over $\Sigma$, $\varepsilon \in \Sigma^*$ stands for the empty word, and $\Sigma^+ = \Sigma^* \setminus \{\varepsilon\}$. For $u, v, w \in \Sigma^*$, we write $|u| \in \mathbb{N}$ for the length of $u$, $u \cdot v$ (or just $uv$) for the concatenation of words, and $w \succ v$ if $v$ is a proper suffix of $w$, that is, $\exists u \in \Sigma^+$ such that $w = u \cdot v$.

Given languages $U, V, W \subseteq \Sigma^*$, concatenation extends to languages as usual: $U \cdot V = \{u \cdot v \mid u \in U, v \in V\}$. The Kleene closure is defined as the smallest set $U^* \subseteq \Sigma^*$ such that $U^* = \{\varepsilon\} \cup U \cdot U^*$. We write the *left quotient* as $U \backslash W = \{v \mid v \in \Sigma^*, \exists u \in U : uv \in W\}$ and the *right quotient* as $W/U = \{v \mid v \in \Sigma^*, \exists u \in U : vu \in W\}$. For a singleton language $U = \{u\}$, we write $u \backslash W$ ($W/u$) for the left (right) quotient.

A *ranked alphabet* $\mathcal{F}$ is a finite set of *operator symbols* with a function $\# : \mathcal{F} \to \mathbb{N}$ that determines the *arity* of each symbol. We write $\mathcal{F}^{(n)} = \{F \in \mathcal{F} \mid \#(F) = n\}$ for the symbols of arity $n$. The set $T_{\mathcal{F}}(X)$ *of $\mathcal{F}$-terms over a set $X$* is defined inductively. If $x \in X$, then $x \in T_{\mathcal{F}}(X)$. If $n \in \mathbb{N}$, $F \in \mathcal{F}^{(n)}$, and $t_1, \dots, t_n \in T_{\mathcal{F}}(X)$, then $F(t_1, \dots, t_n) \in T_{\mathcal{F}}(X)$.

An *$\mathcal{F}$-algebra* consists of a carrier set $M$ and an interpretation function $\mathcal{I} : (n : \mathbb{N}) \to \mathcal{F}^{(n)} \to M^n \to M$. Given a function $\mathcal{I}_0 : X \to M$, the *term interpretation* $\hat{\mathcal{I}}(t)$, for $t \in T_{\mathcal{F}}(X)$, is defined inductively as follows. If $x \in X$, then $\hat{\mathcal{I}}(x) = \mathcal{I}_0(x)$. If $F \in \mathcal{F}^{(n)}$ and $t_1, \dots, t_n \in T_{\mathcal{F}}(X)$, then $\hat{\mathcal{I}}(F(t_1, \dots, t_n)) = \mathcal{I}(n)(F)(\hat{\mathcal{I}}(t_1), \dots, \hat{\mathcal{I}}(t_n))$. We often write $T_{\mathcal{F}}$ in place of $T_{\mathcal{F}}(\emptyset)$.

To avoid notational clutter, we fix an arbitrary alphabet $\Sigma$.

**Definition 1.** *The* regular alphabet *is defined by* $\mathcal{R} = \Sigma \uplus \{\mathbf{0}, \mathbf{1}, \cdot, +, *\}$ *with arities* $\#(x) = 0$, *for* $x \in \Sigma$, $\#(\mathbf{1}) = \#(\mathbf{0}) = 0$, $\#(*) = 1$, *and* $\#(\cdot) = \#(+) = 2$.

*Similarity is defined as the smallest equivalence relation* $\equiv \,\subseteq T_{\mathcal{R}} \times T_{\mathcal{R}}$ *that enforces left and right unit, idempotence, commutativity, and associativity for the $+$ operator. For all $r, s, t \in T_{\mathcal{R}}$, the relation $\equiv$ contains the pairs:*

$$r + \mathbf{0} \equiv r \quad \mathbf{0} + s \equiv s \quad r + r \equiv r \quad r + s \equiv s + r \quad (r+s) + t \equiv r + (s+t)$$

*The set $R$ of regular expressions over $\Sigma$ is defined as the quotient term algebra $R = T_{\mathcal{R}}/(\equiv)$.*

*The language of $r \in R$ is defined by $\mathcal{L}(r) = \hat{\mathcal{I}}(r)$, that is, the interpretation of the term in the $\mathcal{R}$-algebra with carrier set $\wp(\Sigma^*)$ and interpretation function*

$$
\begin{aligned}
\mathcal{I}(0)(\mathbf{0}) &= \{\} & \mathcal{I}(1)(*) &= U \mapsto U^* \\
\mathcal{I}(0)(\mathbf{1}) &= \{\varepsilon\} & \mathcal{I}(2)(\cdot) &= (U, V) \mapsto U \cdot V \\
\mathcal{I}(0)(x) &= \{x\} & \mathcal{I}(2)(+) &= (U, V) \mapsto U \cup V
\end{aligned}
$$

The interpretation function $\mathcal{I}$ is compatible with the definition of $R$ as a quotient term algebra because the interpretation of $+$ maps equivalent expressions

to the same language. We usually work with a unique representative for each equivalence class computed by a function $\mathsf{nf}$ (see [10]). We use parenthesized infix notation for the binary operators $\cdot$ and $+$ and postfix superscript for the unary $*$. We adopt the convention that $\cdot$ binds stronger than $+$ to omit parentheses. The overloading of $\mathbf{0}$ and $\mathbf{1}$ as regular expressions and boolean values is deliberate.

**Definition 2.** *The operations* $\odot, \oplus : R \times R \to R$ *are* smart concatenation *and* union *constructors for regular expressions. Operator* $\odot$ *binds stronger than* $\oplus$.

$$r \odot s = \begin{cases} \mathbf{0} & r = \mathbf{0} \vee s = \mathbf{0} \\ r & s = \mathbf{1} \\ s & r = \mathbf{1} \\ (r \cdot s) & \text{otherwise} \end{cases} \qquad r \oplus s = \mathsf{nf}(r + s)$$

**Lemma 3.** *For all* $r$, $s$: $\mathcal{L}(r \odot s) = \mathcal{L}(r \cdot s)$; $\mathcal{L}(r \oplus s) = \mathcal{L}(r + s)$.

**Definition 4.** *A regular expression* $r$ *is* nullable *if* $\varepsilon \in \mathcal{L}(r)$. *The function* $N : R \to \{\mathbf{0}, \mathbf{1}\}$ *detects nullable expressions:* $N(\mathbf{1}) = \mathbf{1}$. $N(\mathbf{0}) = \mathbf{0}$. $N(x) = \mathbf{1}$. $N(r \cdot s) = N(r) \odot N(s)$. $N(r + s) = N(r) \oplus N(s)$. $N(r^*) = \mathbf{1}$.

**Lemma 5.** *For all* $r \in R$. $N(r) = \mathbf{1}$ *iff* $\varepsilon \in \mathcal{L}(r)$.

**Definition 6.** *The* Brzozowski derivative *[4] is a function* $D : \Sigma \times T_{\mathcal{R}} \to T_{\mathcal{R}}$ *defined inductively for* $a \neq b \in \Sigma$ *and* $r, s \in T_{\mathcal{R}}$.

$$\begin{aligned} D(a, \mathbf{0}) &= \mathbf{0} & D(a, r + s) &= D(a, r) \oplus D(a, s) \\ D(a, \mathbf{1}) &= \mathbf{0} & D(a, r \cdot s) &= D(a, r) \odot s \oplus N(r) \odot D(a, s) \\ D(a, a) &= \mathbf{1} & D(a, r^*) &= D(a, r) \odot r^* \\ D(a, b) &= \mathbf{0} \end{aligned}$$

*It extends to a function on words and languages* $D : \Sigma^* \times T_{\mathcal{R}} \to T_{\mathcal{R}}$ *and* $D : \wp(\Sigma^*) \times T_{\mathcal{R}} \to \wp(T_{\mathcal{R}})$ *as usual* $(a \in \Sigma, w \in \Sigma^*, U \subseteq \Sigma^*)$:

$$D(\varepsilon, r) = r \qquad D(a \cdot w, r) = D(w, D(a, r)) \qquad D(U, r) = \{D(w, r) \mid w \in U\}$$

**Theorem 7 ([4]).** *For all* $w \in \Sigma^*$, $r \in T_{\mathcal{R}}$, $\mathcal{L}(D(w, r)) = w \setminus \mathcal{L}(r)$.

**Theorem 8 ([4]).** *For all* $r \in T_{\mathcal{R}}$, $\mathcal{L}(r) = \mathcal{L}\left(N(r) + \sum_{a \in \Sigma} D(a, r)\right)$.

**Definition 9.** *A (nondeterministic) finite automaton (NFA) is a tuple* $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$ *where* $Q$ *is a finite set of states,* $\Sigma$ *an alphabet,* $\delta : Q \times \Sigma \to \wp(Q)$ *the transition function,* $q_0 \in Q$ *the initial state, and* $F \subseteq Q$ *the set of final states.*

*Let* $n \in \mathbb{N}$. *A* run *of* $\mathcal{A}$ *on* $w = a_0 \ldots a_{n-1} \in \Sigma^*$ *is a sequence* $q_0 \ldots q_n \in Q^*$ *such that, for all* $0 \leq i < n$, $q_{i+1} \in \delta(q_i, a_i)$. *The run is* accepting *if* $q_n \in F$. *The language* $\mathcal{L}(\mathcal{A}) = \{w \in \Sigma^* \mid \exists \text{ accepting run of } \mathcal{A} \text{ on } w\}$ *is* recognized *by* $\mathcal{A}$.

*The automaton* $\mathcal{A}$ *is* total deterministic *if* $|\delta(q, a)| = 1$, *for all* $q \in Q$, $a \in \Sigma$.

## 3 Enhanced Derivatives

An operation on languages takes one or more languages as arguments and yields another language. In this section, we enhance the syntax and semantics of regular expressions with extra operations and consider conditions for the existence of a syntactic derivative for such enhanced expressions. Many examples can be drawn from the closure properties of regular languages.

**Definition 10.** *A function $f : (\Sigma^*)^n \to \Sigma^*$ is* regularity-preserving *if for all regular languages $R_1, \ldots, R_n$ the image $f(R_1, \ldots, R_n)$ is a regular language.*

*Example 11.* We give a range of examples for operators on languages. All operators, except shuffle closure, are regularity-preserving. Proofs may be found in textbooks on formal languages unless otherwise indicated. We let $U, V, L \subseteq \Sigma^*$ range over regular languages; $a, b \in \Sigma$ range over symbols.

1. The intersection $U \cap V$ and the complement $\neg U$ of regular languages are regular.
2. The shuffle of two regular languages is defined by $U \| V = \bigcup \{u \| v \mid u \in U, v \in V\}$ where $\varepsilon \| v = \{v\}$, $u \| \varepsilon = \{u\}$, and $au \| bv = \{a\} \cdot (u \| bv) \cup \{b\} \cdot (au \| v)$, is regular.
   The shuffle closure operation $L^{\|} = \{\varepsilon\} \cup L \cup (L \| L) \cup (L \| L \| L) \cup \ldots$ does **not** preserve regularity.
3. The inverse homomorphism, i.e., $h^{-1}(U) = \{w \in \Sigma^* \mid h(w) \in U\}$ is regular for a function $h : \Sigma \to \Sigma^*$ that is extended homomorphically to a function $\Sigma^* \to \Sigma^*$ (for simplicity, we do not consider homomorphisms between different alphabets, which can be simulated by using the disjoint union of the alphabets).
   The non-erasing homomorphism $h(L) = \{h(w) \mid w \in L\}$ is regular where $h : \Sigma \to \Sigma^+$.
4. The language of every $k$-th symbol starting from position $i$ from words in a regular language $L$ is regular: for $k > 0$ and $0 < i \le k$
   $f_{i,k}(L) = \{a_i a_{i+k} a_{i+2k} \cdots a_{i+k\lfloor (n-i)/k \rfloor} \mid n \in \mathbb{N}, a_1 \ldots a_n \in L\}$.
5. The left quotient $\backslash$ and the right quotient $/$ of regular languages are regular.
6. Functions $suffixes(L) = \Sigma^* \backslash L$ and $prefixes(L) = L / \Sigma^*$ preserve regularity.
7. The function $reverse(L) = \{a_n \cdots a_1 \mid n, i \in \mathbb{N}, 1 \le i \le n, a_i \in \Sigma, a_1 \ldots a_n \in L\}$ preserves regularity.
8. For each $k \in \mathbb{N}$, the function $\mathbb{H}_k(L) = \{v \mid v \in \Sigma^*, \exists u \in L . d(u, v) \le k\}$ is regularity preserving where the *Hamming distance* of words $a_1 \cdots a_n$ and $b_1 \cdots b_m$ is defined by $h = d(a_1 \cdots a_n, b_1 \cdots b_m)$. If $m = n$, then $h = |\{i \mid 1 \le i \le n, a_i \neq b_i\}|$. Otherwise $h = \infty$.
   Analogously, $\mathbb{L}_k(L)$ is a regularity preserving approximation that uses the Levenshtein distance (see [8]).
9. The tilde and bar operators defined by $\tilde{L} = L \cup \{\varepsilon\}$ and $\bar{L} = L \backslash \{\varepsilon\}$ preserve regularity (they are the primitive building blocks of multi-tilde-bar expressions [6], which we do not consider to save space).

The notion of a nullable expression is an important ingredient in the definition of the derivative (Definition 6). Nullability can be computed by induction on a regular expression because each regular operator corresponds to a boolean function on the nullability of the operator's arguments. The following definition imposes exactly this condition on the extra operators in regular expressions.

**Definition 12.** *A function $f : (\Sigma^*)^n \to \Sigma^*$ is $\varepsilon$-testable, if there is a boolean function $B_f : \mathbb{B}^n \to \mathbb{B}$ such that $\varepsilon \in f(L_1, \ldots, L_n)$ iff $B_f((\varepsilon \in L_1), \ldots, (\varepsilon \in L_n))$.*

*Example 13.* Some of the functions from Example 11 are $\varepsilon$-testable.

1. intersection, complement: $B_\cap = \wedge$, $B_\neg = \neg$;
2. shuffle: $B_\| = \wedge$; the shuffle closure operation $L^\|$ is $\varepsilon$-testable using $B^\|(b) = \mathbf{1}$;
3. inverse homomorphism: $B_{h^{-1}}(b) = b$, for $b \in \mathbb{B}$; homomorphism $h$: if $h$ is non-erasing, then $B_h(b) = b$; erasing homomorphism is not $\varepsilon$-testable: consider $L_1 = \{a\}$, $L_2 = \{b\}$, and an erasing homomorphism $h$ defined by $h(a) = \varepsilon$ and $h(b) = b$. Thus, $h(L_1) = \{\varepsilon\}$ and $h(L_2) = \{b\}$. If there was a boolean function $f_h$ to vouch for $\varepsilon$-testability of $h$, then $L_1$ shows that $f_h(\mathbf{0}) = \mathbf{1}$ and $L_2$ yields $f_h(\mathbf{0}) = \mathbf{0}$, a contradiction.
4. $k$-th letter extraction: $\varepsilon \in f_{i,k}(L)$ if $\exists w \in L$ such that $|w| < i$, so $f_{i,k}$ is not $\varepsilon$-testable. To see this let $i = k = 2$, $L_1 = \{a\}$, and $L_2 = \{aa\}$ and assume that $B_f$ is the boolean function required for $\varepsilon$-testability. Now $f_{2,2}(L_1) = \{\varepsilon\}$ and $f_{2,2}(L_2) = \{a\}$, so that $B_f(\mathbf{0}) = \mathbf{1}$ (by $L_1$) and $B_f(\mathbf{0}) = \mathbf{0}$ (by $L_2$), a contradiction.
5. The left quotient is not $\varepsilon$-testable because $\varepsilon \in U \backslash W$ iff $U \cap W \neq \emptyset$: consider $U = \Sigma^*$ with $a \in \Sigma$, $W_1 = \emptyset$, and $W_2 = \{a\}$ so that $U \backslash W_1 = \emptyset$ and $U \backslash W_2 = \{\varepsilon, a\}$. A binary boolean function $B_\backslash$ for $\varepsilon$-testability would have to satisfy $B_\backslash(\mathbf{0}, \mathbf{0}) = \mathbf{0}$ (for $W_1$) and $B_\backslash(\mathbf{0}, \mathbf{0}) = \mathbf{1}$ (for $W_2$), a contradiction. The same reasoning applies, mutatis mutandis, to the right quotient.
6. The *suffixes* function is not $\varepsilon$-testable by the proof for the left quotient. The proof for *prefixes* is analogous to the one for the right quotient.
7. The *reverse* function is $\varepsilon$-testable: $B_{reverse}(b) = b$.
8. The approximation for Hamming distance is $\varepsilon$-testable by $B_{\mathbb{H}_k}(b) = b$ . The approximation for Levenshtein distance $\mathbb{L}_k$ is not $\varepsilon$-testable for $k > 0$. The argument here is similar as for erasing homomorphism because a word at distance $k$ from a given word $w$ may be up to $k$ symbols shorter than $w$.
9. The tilde and bar operators are obviously $\varepsilon$-testable with the constants $\mathbf{1}$ and $\mathbf{0}$, respectively.

**Definition 14 (Enhanced regular expression).** *Let $\mathcal{F} \supseteq \mathcal{R}$ be a ranked alphabet, an* enhanced regular alphabet. *Let further $\mathcal{J}$ be an interpretation function for $\mathcal{F}$ on the carrier $\wp(\Sigma^*)$ extending the regular interpretation $\mathcal{I}$ from Definition 1. The set of $\mathcal{F}$-regular expressions over a set $X$ is the set of terms $T_\mathcal{F}(X)$. For $t \in T_\mathcal{F}(X)$ we define its language $L(t) = \hat{\mathcal{J}}(t)$. The resulting $\mathcal{F}$-algebra $(\wp(\Sigma^*), \mathcal{J})$ is a* regular enhancement *if every symbol $F \in \mathcal{F}^{(n)}$ is interpreted by a regularity-preserving function $\mathcal{J}(n)(F)$.*

*Example 15.* To extend regular expressions with a shuffle operator, consider $\mathcal{F}^{\|} = \mathcal{R} \cup \{\|\}$ with $\#\| = 2$.

To extend expressions with $k$th-letter extraction, we consider $\mathcal{F}^{x-k} = \mathcal{R} \cup \{f_{i,k} \mid 0 < i \leq k\}$ with $\#f_{i,k} = 1$.

**Lemma 16.** *If $\mathcal{J}(F)$ is $\varepsilon$-testable, for each $F \in \mathcal{F}$, then the nullability function $N$ can be extended to $\mathcal{F}$.*

To obtain syntactic derivability for an enhanced regular expression, it must be possible to express the derivative of an operator in terms of a regular expression that applies the derivative to the arguments of the operator. We first define a suitable property semantically as an algebraic property of a regular enhancement.

**Definition 17.** *Let $\mathcal{F}$ be an enhanced regular alphabet and $\mathcal{J}$ an extension of the regular interpretation $\mathcal{I}$. The $\mathcal{F}$-algebra $(\wp(\Sigma^*), \mathcal{J})$ is left derivable if, for each $F \in \mathcal{F}^{(k)}$ and $a \in \Sigma$, there exists a finite subset $X \subset \{x_{v,j} \mid v \in \Sigma^*, 1 \leq j \leq k\}$ and an $\mathcal{F}$-regular expression $r \in T_{\mathcal{F}}(X)$ such that, for all $L_1, \dots, L_k \subseteq \Sigma^*$ the left quotient $a\backslash(\mathcal{J}(F)(L_1, \dots, L_k))$ can be expressed as $\hat{\mathcal{J}}(r)$ using the interpretation $\mathcal{J}_0(x_{v,j}) = v\backslash L_j$.*

*Example 18.* We revisit the previous examples of functions on languages and examine them for being left derivable.

1. Intersection is left derivable: $a\backslash(L_1 \cap L_2) = \hat{\mathcal{J}}(x_{a,1} \cap x_{a,2}) = a\backslash L_1 \cap a\backslash L_2$.
   For negation $\neg$, the pattern is the same.
2. Shuffle is left derivable:
   $a\backslash(L_1\|L_2) = (a\backslash L_1)\|L_2 \cup L_1\|(a\backslash L_2) = \hat{\mathcal{J}}(x_{a,1}\|x_{\epsilon,2} + x_{\epsilon,1}\|x_{a,2})$;
   shuffle closure is also left derivable:

$$a\backslash L^{\|} = (a\backslash L)\|L^{\|} = \hat{\mathcal{J}}(x_{a,1}\|x_{\epsilon,2} + x_{\epsilon,1}^{\|})$$

3. Inverse homomorphism is left derivable:
   $a\backslash(h^{-1}(L)) = h^{-1}(h(a)\backslash L) = \hat{\mathcal{J}}(h^{-1}(x_{h(a),1}))$.
   Non-erasing homomorphism is left derivable:
   $a\backslash(h(L)) = \bigcup_{b \in \Sigma, h(b)=av} v \cdot h(b\backslash L) = \hat{\mathcal{J}}(\sum_{b \in \Sigma, h(b)=av} v \cdot h(x_{b,1}))$.
4. For $k > 1$, the set $\{f_{i,k} \mid 0 < i \leq k\}$ is left derivable.
   $a\backslash(f_{i,k}(L)) = f_k(\bigcup_{|w|=i-1} wa\backslash L) = \hat{\mathcal{J}}(\sum_{|w|=i-1} f_k(x_{wa,1}))$.
5. The left and right quotients are left derivable.
   $a\backslash(L_1\backslash L_2) = (L_1 \cdot a)\backslash L_2 = \hat{\mathcal{J}}((x_{\varepsilon,1} \cdot a)\backslash x_{\varepsilon,2})$.
   $a\backslash(L_1/L_2) = (a\backslash L_1)/L_2 = \hat{\mathcal{J}}(x_{a,1}/x_{\varepsilon,2})$.
6. The function *suffixes* is not left derivable because $a\backslash suffixes(L) = \{w \mid \exists u.uaw \in L\} = (\Sigma^* \cdot a)\backslash L$ cannot be finitely expressed using just derivatives, the *suffixes* function, and the regular operators.
   To see this, consider the family of languages $L_n = w_n^*$ where $w_n = (abab^2 \cdot ab^n)^*$, for all $n \in \mathbb{N}$, and find that

$$L_n' = a\backslash suffixes(L_n) = bab^2 \cdot ab^n w_n^* + b^2 \cdot ab^n w_n^* + \cdots + b^n w_n^*$$

Suppose there is a *suffixes*-enhanced regular expression $r$ for $a\backslash L$ that only depends on $a$ and $\Sigma$ and that refers to finitely many derivatives, say, $v_1\backslash L, \ldots, v_m\backslash L$. Considering $r$ for $L'_n$, we find that $r$ cannot contain the *suffixes* function because that would introduce words starting with $a$, which cannot be in $L'_n$ and which cannot be amended by prepending a fixed word without breaking the $a$-$b$ pattern. There must exist some $v \in w_n^*$ such that each $v_j$ is either a prefix of $v$ that ends with an $a$ or it is not a prefix of $v$.

Now, if we consider $L_k$ where $k = \max(n, |v_1|, \ldots, |v_m|) + 1$ then none of the $v_j\backslash L_k$ can contain $b^k w_k^*$. Note that if $v_j$ is not a prefix of $w_n*$, then it is not a prefix of $w_k^*$, for any $k \geq n$, either. Hence, $r$ cannot describe $L'_k$.

If we assume that $\mathcal{F}$ contains *suffixes* and the left quotient operator, then we could consider *suffixes*$(L)$ as an abbreviation for $\Sigma^*\backslash L$ and we would regain left derivability. Furthermore, with a suitable variation of Definition 17, *suffixes* is *right derivable*:

*suffixes*$(L)/a = \{v \mid \exists u.uv \in L\}/a = \{v \mid \exists u.uva \in L\} = $ *suffixes*$(L/a)$.

The function *prefixes* is left derivable:

$a\backslash$*prefixes*$(L) = a\backslash\{v \mid \exists u.vu \in L\} = $ *prefixes*$(a\backslash L) = \hat{\mathcal{J}}($*prefixes*$(x_{a,1}))$.

7. The function *reverse* is neither left derivable nor right derivable, but swaps between left and right quotients:

   $a\backslash$*reverse*$(L) = $ *reverse*$(L/a)$.

   To see that *reverse* is no left derivable, consider the language $L = b^*a$. Clearly, *reverse*$(L) = ab^*$ and $a\backslash$*reverse*$(L) = b^*$. Now suppose we can obtain $b^*$ by a regular expression with *reverse* on arbitrary derivatives of $L$. There are only two distinct derivatives: $a\backslash(b^*a) = \{\varepsilon\}$ and $b\backslash(b^*a) = b^*a$. Hence, for any $w \in \{a, b\}^*$, $w\backslash(b^*a)$ will be either empty, $\{\varepsilon\}$, or $b^*a$. Now consider a language $U$ constructed from these derivatives by application of regular operators or *reverse*. It can be shown that any word in $U$ is either $\varepsilon$ or it contains the symbol $a$. Thus, $U$ cannot be equal to $b^*$.

8. The enhancement with the approximation operators $\mathbb{H}_k, \mathbb{H}_{k-1}, \ldots, \mathbb{H}_1, \mathbb{H}_0$ operators (for Hamming distance) is left derivable because

$$a\backslash\mathbb{H}_k(L) = \mathbb{H}_k(a\backslash L) + \sum_{\substack{k>0 \\ x\neq a}} \mathbb{H}_{k-1}(x\backslash L)$$

For approximation with operators $\mathbb{L}_k, \ldots, \mathbb{L}_0$ that rely on Levenshtein distance, we also obtain left closure (assuming that $\mathbb{L}_{-1}(L) = \emptyset$):

$$a\backslash\mathbb{L}_k(L) = \sum_{\substack{w\in\Sigma^* \\ |w|\leq k \\ k>0}} \left( \mathbb{L}_{k-|w|}(wa\backslash L) + \sum_{x\neq a} \mathbb{L}_{k-|w|-1}(wx\backslash L) + \mathbb{L}_{k-|w|-1}(w\backslash L) \right)$$

The terms correspond to the actions "delete $w$, then match $a$", "delete $w$, then substitute $a$ by some $x$", and "delete $w$, then insert $a$".

9. Tilde and bar are trivially left derivable: $a\backslash\tilde{L} = a\backslash L$ and $a\backslash\bar{L} = a\backslash L$.

$$D^+(\mathbf{0}) = \{\mathbf{0}\} \qquad D^+(r + s) = D^+(r) \oplus D^+(s)$$
$$D^+(\mathbf{1}) = \{\mathbf{0}\} \qquad D^+(r \cdot s) \ = D^+(r) \odot s \oplus \bigoplus D^+(s)$$
$$D^+(a) = \{\mathbf{0}, \mathbf{1}\} \qquad D^+(r^*) \quad = \bigoplus D^+(r) \odot r^*$$

**Fig. 1.** Iterated Brzozowski derivatives for $T_{\mathcal{R}}$

## 4 Word Problem

To obtain a decision procedure for the word problem of left derivable enhanced regular expressions, we first define the corresponding syntactic derivative and then extend Brzozowski's result that $w \in \mathrm{L}(r)$ iff $\varepsilon \in \mathrm{L}(D(w, r))$ (which follows from Theorem 7). It is interesting to remark that, for example, we obtain a decision procedure for the word problem for the language of regular expressions enhanced with the shuffle-closure operator is no longer regular.

**Theorem 19.** *If $(\wp(\Sigma^*), \mathcal{J})$ is a left derivable $\mathcal{F}$-algebra which is $\varepsilon$-testable, then there is a syntactic derivative function $D : \Sigma \times T_{\mathcal{F}} \to T_{\mathcal{F}}$ such that $\hat{\mathcal{J}}(D(a, t)) = a \backslash \hat{\mathcal{J}}(t)$, for all $a \in \Sigma$ and $t \in T_{\mathcal{F}}$.*

*Proof.* Define $D$ inductively as an extension of Definition 6 for $F \in \mathcal{F} \setminus \mathcal{R}$:

$$D(a, F(r_1, \ldots, r_n)) = R(F, a)[x_{v,j} \mapsto D(v, r_j) \mid x_{v,j} \in X(F, a)]$$

where $N$ extends to $T_{\mathcal{F}}$ by Lemma 16 and where $D$ extends to words as before.

$$D(\varepsilon, r) = r \qquad\qquad D(aw, r) = D(w, D(a, r))$$

The statement about the semantics follows by induction on the augmented term using the definition of left derivability. $\qquad\square$

**Theorem 20.** *If $(\wp(\Sigma^*), \mathcal{J})$ is a left derivable $\mathcal{F}$-algebra which is $\varepsilon$-testable, then the word problem for $\hat{\mathcal{J}}(t)$ is decidable, for any $t \in T_{\mathcal{F}}$.*

*Proof.* By Theorem 19, there is a nullability function $N$ and a derivative $D$ for $T_{\mathcal{F}}$. By induction on the length of $w \in \Sigma^*$, we obtain that $w \in \hat{\mathcal{J}}(t)$ iff $\varepsilon \in \hat{\mathcal{J}}(D(w, t))$ iff $N(D(w, t))$. $\qquad\square$

## 5 Finiteness

For classical derivatives on $T_{\mathcal{R}}$ (cf. Definition 6), Brzozowski showed that the set of iterated derivatives $D(\Sigma^*, r)$ of a given regular expression $r$ is finite, when considered modulo similarity (i.e., associativity, commutativity, and idempotence of union). Hence, we now look for conditions such that the set of dissimilar iterated derivatives is finite for enhanced regular expressions. First, we set up a framework for reasoning about finiteness.

Recent work on determining the number of iterated *partial* derivatives starts with an inductive definition for the set of iterated partial derivatives [3]. We

transfer that definition to the classical case and define an upper approximation $D^+(r)$ of the set of iterated derivatives of expression $r$ in Figure 1 by induction on $r$. In the definition, we lift $\odot$ and $\oplus$ to sets of expressions (i.e., if $R, S \subseteq T_{\mathcal{R}}$, then $R \odot S = \{r \odot s \mid r \in R, s \in S\}$ and $R \oplus S = \{r \oplus s \mid r \in R, s \in S\}$). We further write $\bigoplus S$ for the set $\{s_1 \oplus \cdots \oplus s_n \mid n \in \mathbb{N}, s_i \in S\}$ of finite sums of elements drawn from $S$ where the nullary sum stands for $\mathbf{0}$ and where we assume sums to be identified modulo associativity, commutativity, and idempotence to obtain the following results.[1]

**Theorem 21.** *The set $D^+(r)$ is finite, for all $r \in T_{\mathcal{R}}$.*

Clearly, the set $D^*(r) = \{r\} \cup D^+(r)$ is also finite for all $r$.

**Theorem 22 (Closure under derivation).**

1. *For all $r$ and $a$, $D(a, r) \in D^+(r)$.*
2. *For all $r$ and $a$, if $t \in D^+(r)$, then $D(a, t) \in D^+(r)$.*

**Corollary 23.** *The set $\{D(w, r) \mid w \in \Sigma^+\} \subseteq D^+(r)$, for all $r$.*

To obtain finiteness for enhanced regular expressions, we strengthen the notion of left derivability. Essentially, we restrict the form of a derivative to a linear combination of enhancement functions applied to derivatives of the arguments.

**Definition 24.** *Let $\mathcal{F} = \{F_1, \ldots, F_m\}$ be a ranked alphabet. The $\mathcal{F}$-algebra $(\wp(\Sigma^*), \mathcal{J})$ is* linear left derivable *if, for each $n \in \mathbb{N}$, $F \in \mathcal{F}^{(n)}$, and $a \in \Sigma$, there exists a finite index set $J$ such that, for each $j \in J$, there is a word $v_j \in \Sigma^*$, an index $i_j \in \{1, \ldots, m\}$ of an element of $\mathcal{F}$ with arity $\#(F_{i_j}) = n_j$, and, for $1 \leq k \leq n_j$, words $w_k^j \in \Sigma^*$ and indexes $\alpha_k^j \in \{1, \ldots, n\}$ of left-hand-side languages, such that for all $L_1, \ldots L_n \subseteq \Sigma^*$, the left quotient can be expressed by:*

$$a \backslash (\mathcal{J}(F)(L_1, \ldots, L_n)) = \bigcup_{j \in J} v_j \cdot \mathcal{J}(F_{i_j})(\overline{w_k^j \backslash (L_{\alpha_k^j})}^{k=1,\ldots,n_j}) \tag{1}$$

Of the standard regular operators, only union (and in fact all boolean functions) is linear left derivable. Concatenation $U \cdot V$ does not fit the pattern because it has a summand which is conditional on $\varepsilon \in U$. The Kleene star does not fit, either, because it concatenates the derivative of the argument with the original term (Definition 6). But many useful operators are linear left derivable (Example 28).

**Theorem 25.** *Suppose that $\mathcal{F} = \{F_1, \ldots, F_m\} \cup \mathcal{R}$ is an enhanced regular alphabet with interpretation $\mathcal{J}$ such that $(\wp(\Sigma^*), \mathcal{J}_{|\{F_1,\ldots,F_m\}})$ is linear left derivable. Then, for all $n \in \mathbb{N}$, $F \in \mathcal{F}^{(n)}$, and $a \in \Sigma$ there exists a finite index set $J$, for each $j \in J$, there is a word $v_j \in \Sigma^*$, an index $i_j \in \{1, \ldots, m\}$ of an element of $\mathcal{F} \setminus \mathcal{R}$ with arity $\#(F_{i_j}) = n_j$, for each $1 \leq k \leq n_j$, a word $w_k^j \in \Sigma^*$, and an*

---

[1] See the technical report for auxiliary lemmas and proofs.

index $\alpha_k^j \in \{1, \ldots, n\}$ *that selects one of the left-hand-side regular expressions as an argument. Then, for each $r_1, \ldots, r_n \in T_{\mathcal{F}}$, the syntactic derivative of $F(r_1, \ldots, r_n)$ by $a$ is given in the form*

$$D(a, F(r_1, \ldots, r_n)) = \sum_{j \in J(F,a)} v_j \cdot F_{i_j}(\overline{D(w_k^j, r_{\alpha_k^j})}^{k=1,\ldots,n_j}) \tag{2}$$

*In this setting, the set of iterated derivatives of any $\mathcal{F}$-regular expression $r$ is finite. Specifically, in extension of the definition in Figure 1, we claim that for each $F \in \mathcal{F} \setminus \mathcal{R}$, the set of iterated derivatives*

$$D^+(F(r_1, \ldots, r_n)) = \bigoplus \{v \cdot G(\overline{r_i'}) \mid V \succeq v, G \in \mathcal{F} \setminus \mathcal{R}, r_i' \in \bigcup_j D^*(r_j)\} \tag{3}$$

*is finite. Here $V = \{v_j \mid j \in J, F \in \mathcal{F}, a \in \Sigma\}$, and we write $V \succeq v$ for $\exists v' \in V. v' \succeq v$.*

**Corollary 26.** *Let $\mathcal{F}$ be an enhanced regular alphabet and $(\wp(\Sigma^*), \mathcal{J})$ be an $\varepsilon$-testable, linear left derivable $\mathcal{F}$-algebra. Then any $\mathcal{F}$-regular expression defines a regular language.*

*Proof.* Let $r \in T_{\mathcal{F}}$ and let $Q_r$ be the set of dissimilar derivatives of $r$. As $Q_r \subseteq D^*(r)$, $Q_r$ is finite. Hence $M = (Q_r, \Sigma, D, r, F)$ with $F = \{q \in Q_r \mid N(q)\}$ is a total deterministic finite automaton that recognizes $\mathrm{L}(r)$, which is thus regular. $\square$

**Corollary 27.** *Let $\mathcal{F}$ be an enhanced regular alphabet and $(\wp(\Sigma^*), \mathcal{J})$ be an $\varepsilon$-testable, linear left derivable $\mathcal{F}$-algebra. Then, for each $F \in \mathcal{F}$, the operation $\mathcal{J}(F)$ preserves regularity.*

*Proof.* Let $F \in \mathcal{F}^{(n)}$, for some $n \in \mathbb{N}$. Let $R_1, \ldots, R_n$ be regular languages defined by regular expressions $r_1, \ldots, r_n \in T_{\mathcal{R}} \subseteq T_{\mathcal{F}}$. By Corollary 26, $\hat{\mathcal{J}}(F(r_1, \ldots, r_n))$ is regular. Hence $\mathcal{J}(F)$ preserves regularity. $\square$

*Example 28.* Many operators are in fact linear left derivable.

1. Intersection and complement are linear left derivable.
2. The shuffle operation is linear left derivable, but the derivative of the shuffle closure contains a nested application of shuffle closure.
3. Inverse and non-erasing homomorphism are linear left derivable.
4. For $k > 0$, the set $\{f_{i,k} \mid 0 < i \leq k\}$ is linear left derivable.
5. The left quotient is not linear left derivable, but the right quotient is linear left derivable.
6. The function *suffixes* is not left derivable; the function *prefixes* is linear left derivable.
7. The function *reverse* is not left derivable.
8. Both, $\mathbb{H}_k$ and $\mathbb{L}_k$ are linear left derivable.
9. Tilde and bar are linear left derivable.

By Corollary 26, regular languages are closed under $\varepsilon$-testable operators that are linear left derivable: $\cap$, $\neg$, $\|$, $h^{-1}$, non-erasing $h$, $\mathbb{H}_k$, tilde, bar.

For a set of unary operators, linear left derivability amounts to definability by a rational finite state transducer.

**Theorem 29.** *Let $\mathcal{F} = \mathcal{F}^{(1)} = \{F_1, \ldots, F_m\}$ be a ranked alphabet of unary operators and $(\wp(\Sigma^*), \mathcal{J})$ be a linear left derivable $\mathcal{F}$-algebra which is $\varepsilon$-testable using the identity function. Then, for each $1 \leq l \leq m$ and $L \subseteq \Sigma^*$, $\mathcal{J}(F_l)(L)$ is equal to $T(L)$ where $T$ is a rational finite state transducer.*

The reverse implication does not hold because transducers may, in general, consume an unbounded amount of input before producing an output. The transducers resulting from Theorem 29 only consume bounded input before producing at least one output symbol.

## 6 Conclusion

We introduce a framework for constructing derivatives for regular expressions enhanced with new operators. If these operators are left derivable, we obtain an algorithm for the word problem; if they are linear left derivable, we can construct a DFA from an enhanced expression. In fact, unary operators with this property are rational transductions.

Some of the operators considered in this paper are known to be regularity preserving, yet, they fail to be linear left derivable or to be $\varepsilon$-testable. In future work, we plan to address these restrictions by generalizing linear derivability as well as the nullability test.

## References

1. Antimirov, V.M.: Rewriting regular inequalities. In: Proc. of FCT'95. LNCS, vol. 965, pp. 116–125. Springer (1995)
2. Antimirov, V.M.: Partial derivatives of regular expressions and finite automaton constructions. Theoretical Computer Science 155(2), 291–319 (1996)
3. Broda, S., Machiavelo, A., Moreira, N., Reis, R.: Study of the average size of Glushkov and partial derivative automata (Oct 2011)
4. Brzozowski, J.A.: Derivatives of regular expressions. J. ACM 11(4), 481–494 (1964)
5. Caron, P., Champarnaud, J.M., Mignot, L.: Partial derivatives of an extended regular expression. In: LATA. LNCS, vol. 6638, pp. 179–191. Springer (2011)
6. Caron, P., Champarnaud, J., Mignot, L.: Multi-tilde-bar derivatives. In: CIAA. LNCS, vol. 7381, pp. 321–328. Springer (2012)
7. Caron, P., Champarnaud, J., Mignot, L.: A general framework for the derivation of regular expressions. RAIRO - Theor. Inf. and Applic. 48(3), 281–305 (2014)
8. Champarnaud, J., Jeanne, H., Mignot, L.: Approximate regular expressions and their derivatives. In: LATA. LNCS, vol. 7183, pp. 179–191. Springer (2012)
9. Champarnaud, J., Mignot, L., Nicart, F.: Constrained expressions and their derivatives. Tech. Rep. 1406.6144v2, arXiv (Oct 2015), http://arxiv.org/abs/1406.6144v2

10. Grabmayer, C.: Using proofs by coinduction to find "traditional" proofs. In: Proc. of CALCO'05. pp. 175–193. Springer (2005)
11. Might, M., Darais, D., Spiewak, D.: Parsing with derivatives: a functional pearl. In: Proc. of ICFP'11. pp. 189–195. ACM (2011)
12. Owens, S., Reppy, J., Turon, A.: Regular-expression derivatives reexamined. Journal of Functional Programming 19(2), 173–190 (2009)
13. Rosu, G., Viswanathan, M.: Testing extended regular language membership incrementally by rewriting. In: RTA'03. LNCS, vol. 2706, pp. 499–514. Springer (2003)
14. Sulzmann, M., Thiemann, P.: Derivatives for regular shuffle expressions. In: LATA'15. LCNS, vol. 8977, pp. 275–286. Springer, Nice, France (Mar 2015)
15. Sulzmann, M., Thiemann, P.: Forkable regular expressions. In: Proc. of LATA'16. Prague, Czech Republic (2016)
16. Traytel, D., Nipkow, T.: Verified decision procedures for MSO on words based on derivatives of regular expressions. J. Funct. Program. 25 (2015)
17. Watson, B.W.: FIRE lite: FAs and REs in C++. In: Workshop on Implementing Automata. LNCS, vol. 1260, pp. 167–188. Springer (1996)

# A   Proofs and auxiliary lemmas

*Proof (of Lemma 16).* If $f = \mathcal{J}(F)$ for $F \in \mathcal{F}^{(n)}$ is $\varepsilon$-testable, then there exists a boolean function $B_f$ such that $\varepsilon \in f(L_1, \ldots, L_n)$ iff $B_f(\varepsilon \in L_1, \ldots, \varepsilon \in L_n)$. Now set

$$N(F(t_1, \ldots, t_n)) := B_f(N(t_1), \ldots, N(t_n))$$

for the desired extension of $N$. □

*Proof (of Theorem 25).* In equation (3), the set $V$ is finite as it is a union of finite sets. Hence, the set $\{v \mid V \succeq v\}$ is also finite. As $\mathcal{F}$ and $D^+(r_j)$ are also assumed finite, the set $D^+(F(r_1, \ldots, r_n))$ is finite.

The set $D^+(F(r_1, \ldots, r_n))$ is also closed under formation of derivatives. Consider the derivative of a summand of the form $D(a, v \cdot F(r'_1, \ldots, r'_n))$:

$$D(a, v \cdot F(\overline{r'_i})) = \begin{cases} \mathbf{0} & v = bv', a \neq b \\ v' \cdot F(\overline{r'_i}) & v = av' \\ \sum_{j \in J(F,a)} v_j \cdot F_{i_j}(\overline{D(w_k^j, r'_{\alpha_k^j})}^{k=1,\ldots,n_j}) & v = \varepsilon \end{cases}$$

In each case, the outcome is covered by the right hand side of Equation (3). □

We need a few auxiliary lemmas before we can prove that $D^+(r)$ is closed under the derivative operation.

**Lemma 30.** *For all $r \in T_{\mathcal{R}}$, $\mathbf{0} \in D^+(r)$.*

*Proof (Lemma 30).* Induction on $r$.
 **Case $\mathbf{0}$, $\mathbf{1}$, $a$:** Immediate.
 **Case $r + s$:** $\mathbf{0} = \mathbf{0} \oplus \mathbf{0} \in D^+(r) \oplus D^+(s)$, by induction.
 **Case $r \cdot s$:** $\mathbf{0} = \mathbf{0} \odot s \oplus \mathbf{0} \in D^+(r) \odot s \oplus D^+(s)$, by induction.
 **Case $r^*$:** $\mathbf{0} = \mathbf{0} \odot r^* \in D^+(r) \odot r^*$.

**Lemma 31.** *For all $r$ and $s$, $D(a, r \oplus s) = D(a, r) \oplus D(a, s)$.*

*Proof (Lemma 31).* By cases in the definition of $\oplus$.
    **Case $s = \mathbf{0}$.** $D(a, r \oplus \mathbf{0}) = D(a, r) = D(a, r) \oplus \mathbf{0} = D(a, r) \oplus D(a, \mathbf{0})$.
    **Case $r = \mathbf{0}$.** Analogously.
    **Case $r = s$.** $D(a, r \oplus r) = D(a, r) = D(a, r) \oplus D(a, r)$.
    **Case $r \oplus s = r + s$.** $D(a, r \oplus s) = D(a, r + s) = D(a, r) \oplus D(a, s)$.

**Lemma 32.** *For all $r$ and $s$, $D(a, r \odot s) = D(a, r \cdot s)$.*

*Proof (Lemma 32).* By cases in the definition of $\odot$.
    **Case $r = \mathbf{0}$.** $D(a, r \odot s) = D(a, \mathbf{0}) = \mathbf{0}$. $D(a, r \cdot s) = D(a, \mathbf{0} \cdot s) = D(a, \mathbf{0}) \odot s \oplus N(\mathbf{0}) \odot D(a, s) = \mathbf{0} \oplus \mathbf{0} = \mathbf{0}$.
    **Case $s = \mathbf{0}$.** Similar.
    **Case $r = \mathbf{1}$.** $D(a, r \odot s) = D(a, s)$. $D(a, r \cdot s) = D(a, \mathbf{1} \cdot s) = D(a, \mathbf{1}) \odot s \oplus N(\mathbf{1}) \odot D(a, s) = \mathbf{0} \oplus D(a, s) = D(a, s)$.
    **Case $s = \mathbf{1}$.** Similar.
    **Case $r \odot s = r \cdot s$.** Immediate.

*Proof (of Theorem 22).* **Part 1.** By induction on $r$.
    **Case 0**: $D(a, \mathbf{0}) = \mathbf{0} \in D^+(\mathbf{0})$.
    **Case 1**: $D(a, \mathbf{1}) = \mathbf{0} \in D^+(\mathbf{1})$.
    **Case $b$**: $D(a, b) \in \{\mathbf{0}, \mathbf{1}\} = D^+(b)$.
    **Case $r + s$**: immediate.
    **Case $r \cdot s$**: $D(a, r \cdot s) = D(a, r) \odot s \oplus N(r) \odot D(a, s)$. By induction, $D(a, r) \in D^+(r)$ and $D(a, s) \in D^+(s)$. We distinguish two cases on the outcome of $N(r)$.
    **Subcase $N(r) = \mathbf{0}$**: In this case $D(a, r \cdot s) = D(a, r) \odot s \oplus \mathbf{0} \in D^+(r) \odot s \oplus \mathbf{0}$, by induction, and by Lemma 30, $\mathbf{0} \in D^+(s)$, so $D^+(r) \odot s \oplus \mathbf{0} \subseteq D^+(r) \odot s \oplus D^+(s)$.
    **Subcase $N(r) = \mathbf{1}$**: In this case $D(a, r \cdot s) = D(a, r) \odot s \oplus D(a, s) \in D^+(r) \odot s \oplus D^+(s)$, by induction.
    **Case $r^*$**: $D(a, r^*) = D(a, r) \odot r^* \in D^+(r) \odot r^*$ by induction.
    **Part 2.** By induction on $r$.
    **Case 0**, $\mathbf{1}$, $a$: immediate.
    **Case $r + s$**: if $t \in D^+(r + s)$, then $t = r' \oplus s'$ for some $r' \in D^+(r)$ and $s' \in D^+(s)$. By Lemma 31, $D(a, r' \oplus s') = D(a, r') \oplus D(a, s') \in D^+(r) \oplus D^+(s)$, by induction.
    **Case $r \cdot s$**: if $t \in D^+(r \cdot s)$, then $t = r' \odot s \oplus s'$ for some $r' \in D^+(r)$ and $s' \in D^+(s)$. By Lemmas 31 and 32, $D(a, r' \odot s \oplus s') = D(a, r' \odot s) \oplus D(a, s') = D(a, r' \cdot s) \oplus D(a, s') = D(a, r') \odot s \oplus N(r') \odot D(a, s) \oplus D(a, s')$.
    By induction, $D(a, r') \odot s \in D^+(r) \odot s$.
    By item 1, $D(a, s) \in D^+(s)$.
    By induction, $D(a, s') \in D^+(s)$.
    Hence, $N(r') \odot D(a, s) \oplus D(a, s') \in \bigoplus D^+(s)$, which proves the claim.
    **Case $r^*$**: if $t \in D^+(r^*)$, then $t = r' \odot r^*$ for some $r' \in D^+(r)$. Now, $D(a, r' \odot r^*) = D(a, r' \cdot r^*) = D(a, r') \odot r^* \oplus N(r') \odot D(a, r) \odot r^*$.
    By induction, $D(a, r') \odot r^* \in D^+(r) \odot r^*$.
    By item 1, $D(a, r) \odot r^* \in D^+(r) \odot r^*$.
    Hence, $D(a, r') \odot r^* \oplus N(r') \odot D(a, r) \odot r^* \in \bigoplus D^+(r) \odot r^*$.

*Proof (Theorem 29).* For $F_l$, define the finite state transducer $T = (Q, \Sigma, \Sigma, I, A, \delta)$ as follows: set of states $Q = \{F_1, \ldots, F_m\}$; initial states $I = \{F_l\}$; accepting states $A = Q$; the transition relation is the smallest relation $\delta \subseteq Q \times \Sigma^* \times \Sigma^* \times Q$ (relating a state and an input word to an output word and a next state) such that: For $F \in \mathcal{F}$ and $a \in \Sigma$, if, by linear left derivability,

$$a \backslash (\mathcal{J}(F)(L)) = \bigcup_{j \in J(F,a)} v_j \cdot \mathcal{J}(F_{i_j})(w_j \backslash L)$$

where $v_j, w_j \in \Sigma^*$, then $\delta \supseteq \{(F, w_j, a \cdot v_j, F_{i_j}) \mid j \in J(F,a)\}$.

Now recall the definition of $T_l(L)$ (where $l = 1, \ldots, m$), the language translated from state $F_l$, for the transducer $T = (Q, \Sigma, \Sigma, I, F, \delta)$:

$$v \in T_i(L) \Leftrightarrow \exists w \in L : \exists k \in \{1, \ldots, m\} : (F_i, v, w, F_k) \in \delta^*$$

where $\delta^* \subseteq Q \times \Sigma^* \times \Sigma^* \times Q$ is the iterated transition relation defined as the smallest relation such that:

1. for all $F \in Q$, $(F, \varepsilon, \varepsilon, F) \in \delta^*$;
2. for all $F_i, F_j, F_k \in Q$, $v', v, w', w \in \Sigma$,
   $(F_i, v'v, w'w, F_k) \in \delta^*$ iff $(F_i, v', w', F_j) \in \delta$ and $(F_j, v, w, F_k) \in \delta^*$.

Now suppose that $v \in \mathcal{J}(F)(L)$ and show that $v \in T(L)$ by induction on the length of $v$.

If $v = \varepsilon$, then $\varepsilon$-testability with the identity function implies that $\varepsilon \in L$. Hence $\varepsilon \in T(L)$ by the definition of $\delta^*$ (Item 1).

If $v \neq \varepsilon$, then consider $a \backslash v \in a \backslash (\mathcal{J}(F)(L))$. By assumption of linear left derivability, we find that

$$a \backslash (\mathcal{J}(F)(L)) = \bigcup_{j \in J(F,a)} v_j \cdot \mathcal{J}(F_{i_j})(w_j \backslash L)$$

for some $v_j, w_j \in \Sigma^*$. Thus, there exists some $j \in J(F, a)$ such that

$$a \backslash v \in v_j \cdot \mathcal{J}(F_{i_j})(w_j \backslash L)$$

Hence, there is some $v' \in \Sigma^*$ such that $v = a \cdot v_j \cdot v'$ with $v' \in \mathcal{J}(F_{i_j})(w_j \backslash L)$. By induction, $v' \in T(w_j \backslash L)$ which means that there exists some $w' \in w_j \backslash L$ and $F_k \in Q$ such that $(F_{i_j}, w', v', F_k) \in \delta^*$. By construction, $(F, w_j, a \cdot v_j, F_{i_j}) \in \delta$, so by definition of $\delta^*$ (Item 2), $(F, w_j \cdot w', a \cdot v_j \cdot v', F_k) \in \delta^*$ and thus $v \in T(L)$.

The reasoning for the reverse inclusion $T(L) \subseteq \mathcal{J}(F)(L)$ is analogous. $\square$

# B  On upward and downward closures

The upward and downward closure of any language is regular \cite{HAINES196994}:
$\uparrow L = \{y \mid \exists x \in L.x \sqsubseteq y\}$ and $\downarrow L\{x \mid \exists y \in L.x \sqsubseteq y\}$. Here, $\sqsubseteq\, \subseteq \Sigma^* \times \Sigma^*$ is the
*subword ordering* where $x \sqsubseteq y$ iff $x = a_1 \cdots a_n$ and $y = u_0 a_1 u_1 \cdots u_{n-1} a_n u_n$ for
some $n \in \mathbb{N}$, $a_i \in \Sigma$, for $1 \le i \le n$, and $u_i \in \Sigma^*$, for $0 \le i \le n$.

The upward closure is $\varepsilon$-testable using the identity function, but the downward closure is not, in fact, $\varepsilon \in \downarrow L$ iff $L \ne \emptyset$. To see this, consider $L_1 = \emptyset$ and
$L_2 = \{a\}$ with downward closures $\downarrow L_1 = \emptyset$ and $\downarrow L_2 = \{\varepsilon, a\}$. Thus, a hypothetical boolean function $B_\downarrow$ would have to satisfy $B_\downarrow(\mathbf{0}) = \mathbf{0}$ (according to $L_1$) and
$B_\downarrow(\mathbf{0}) = \mathbf{1}$ (according to $L_2$), a contradiction.

The upward closure operation is left derivable as $a\backslash(\uparrow L) = \uparrow L \cup \uparrow(a\backslash L)$.
The downward closure operation does not appear to be left derivable as the left quotient involves skipping ahead arbitrarily many characters: $a\backslash(\downarrow L) = \sum_{w \in \Sigma^*} \downarrow(wa\backslash L)$, but we have neither proof nor disproof for this conjecture.

Upward closure is linear left derivable, but downward closure is not.