

Finiteness of Symbolic Derivatives in Lean

Ekaterina Zhuchko¹, Hendrik Maarand¹,
Margus Veanes², Gabriel Ebner²

¹Tallinn University of Technology, Estonia

²Microsoft Research, USA

ITP, October 2025

Introduction

- *Brzozowski derivatives* of regular expressions

$$\begin{aligned}\mathcal{L}(\text{der}(c, R)) &:= \{w \in \Sigma^* \mid c \cdot w \in \mathcal{L}(R)\} \\ \text{null}(R) &:= \epsilon \in \mathcal{L}(R)\end{aligned}$$

Introduction

- *Brzozowski derivatives* of regular expressions

$$\begin{aligned}\mathcal{L}(\text{der}(c, R)) &:= \{w \in \Sigma^* \mid c \cdot w \in \mathcal{L}(R)\} \\ \text{null}(R) &:= \epsilon \in \mathcal{L}(R)\end{aligned}$$

- (Brzozowski, 64): *finiteness* of all iterated derivatives

$$\mathcal{D}\text{er}(R) = \{\text{der}_w^*(R) \mid w \in \sigma^*\} / \cong$$

quotiented by a relation called *ACI-similarity*:

$$\begin{array}{ll}(L \uplus R) \uplus S \cong L \uplus (R \uplus S) & \text{Associativity} \\ L \uplus (R \uplus L) \cong L \uplus R & \text{Commutativity} \\ R \uplus R \cong R & \text{Idempotence}\end{array}$$

This work

- ▶ We prove finiteness of *symbolic derivatives*:

This work

- ▶ We prove finiteness of *symbolic derivatives*:
 1. Derivatives do not take concrete characters

This work

- ▶ We prove finiteness of *symbolic derivatives*:
 1. Derivatives do not take concrete characters
 2. Derivatives return a *transition term* (instead of a regex)

This work

- ▶ We prove finiteness of *symbolic derivatives*:
 1. Derivatives do not take concrete characters
 2. Derivatives return a *transition term* (instead of a regex)
 3. We build an *overapproximation* for the set of all iterated derivatives

This work

- ▶ We prove finiteness of *symbolic derivatives*:
 1. Derivatives do not take concrete characters
 2. Derivatives return a *transition term* (instead of a regex)
 3. We build an *overapproximation* for the set of all iterated derivatives
- ▶ We consider *symbolic regular expressions* with *lookarounds*

This work

- ▶ We prove finiteness of *symbolic derivatives*:
 1. Derivatives do not take concrete characters
 2. Derivatives return a *transition term* (instead of a regex)
 3. We build an *overapproximation* for the set of all iterated derivatives
- ▶ We consider *symbolic regular expressions with lookarounds*
 - ▶ The alphabet is symbolic and represented by an *Effective Boolean Algebra* $\mathcal{A} = (\Sigma, \alpha, \models, \perp, \top, \sqcup, \sqcap, ^c)$

This work

- ▶ We prove finiteness of *symbolic derivatives*:
 1. Derivatives do not take concrete characters
 2. Derivatives return a *transition term* (instead of a regex)
 3. We build an *overapproximation* for the set of all iterated derivatives
- ▶ We consider *symbolic regular expressions with lookarounds*
 - ▶ The alphabet is symbolic and represented by an *Effective Boolean Algebra* $\mathcal{A} = (\Sigma, \alpha, \models, \perp, \top, \sqcup, \sqcap, ^c)$
- ▶ We do not assume commutativity of union

This work

- ▶ We prove finiteness of *symbolic derivatives*:
 1. Derivatives do not take concrete characters
 2. Derivatives return a *transition term* (instead of a regex)
 3. We build an *overapproximation* for the set of all iterated derivatives
- ▶ We consider *symbolic regular expressions with lookarounds*
 - ▶ The alphabet is symbolic and represented by an *Effective Boolean Algebra* $\mathcal{A} = (\Sigma, \alpha, \models, \perp, \top, \sqcup, \sqcap, ^c)$
- ▶ We do not assume commutativity of union
 - ▶ PCRE (leftmost-greedy) vs POSIX (leftmost-longest)

This work

- ▶ We prove finiteness of *symbolic derivatives*:
 1. Derivatives do not take concrete characters
 2. Derivatives return a *transition term* (instead of a regex)
 3. We build an *overapproximation* for the set of all iterated derivatives
- ▶ We consider *symbolic regular expressions* with *lookarounds*
 - ▶ The alphabet is symbolic and represented by an *Effective Boolean Algebra* $\mathcal{A} = (\Sigma, \alpha, \models, \perp, \top, \sqcup, \sqcap, ^c)$
- ▶ We do not assume commutativity of union
 - ▶ PCRE (leftmost-greedy) vs POSIX (leftmost-longest)
 - ▶ Let $R = (a \sqcup ab)^*$ and $s = "abab"$

This work

- ▶ We prove finiteness of *symbolic derivatives*:
 1. Derivatives do not take concrete characters
 2. Derivatives return a *transition term* (instead of a regex)
 3. We build an *overapproximation* for the set of all iterated derivatives
- ▶ We consider *symbolic regular expressions with lookarounds*
 - ▶ The alphabet is symbolic and represented by an *Effective Boolean Algebra* $\mathcal{A} = (\Sigma, \alpha, \models, \perp, \top, \sqcup, \sqcap, ^c)$
- ▶ We do not assume commutativity of union
 - ▶ PCRE (leftmost-greedy) vs POSIX (leftmost-longest)
 - ▶ Let $R = (a \sqcup ab)^*$ and $s = "abab"$
 - ▶ "**a**bab" vs "**a**bab"

Regular Expressions with Lookarounds

$$R, S ::= \psi \in \alpha \mid \varepsilon \mid R \uplus S \mid R \sqcap S \mid R \cdot S \mid R^* \mid \sim R \\ \mid (?=R) \mid (?<=R) \mid (?!R) \mid (?<!R)$$

- We work modulo an alphabet theory

$$\mathcal{A} = (\Sigma, \alpha, \models, \perp, \top, \sqcup, \sqcap, ^c)$$

For example, $\psi_{upper} \in \alpha$ and $\llbracket \psi_{upper} \rrbracket = [A - Z]$

Regular Expressions with Lookarounds

$$R, S ::= \psi \in \alpha \mid \varepsilon \mid R \uplus S \mid R \sqcap S \mid R \cdot S \mid R^* \mid \sim R \\ \mid (?=R) \mid (?<=R) \mid (?!R) \mid (?<!R)$$

- ▶ We work modulo an alphabet theory
 $\mathcal{A} = (\Sigma, \alpha, \models, \perp, \top, \sqcup, \sqcap, ^c)$
For example, $\psi_{upper} \in \alpha$ and $\llbracket \psi_{upper} \rrbracket = [A - Z]$
- ▶ Positive lookahead $(?=R)$ and lookbehind $(?<=R)$
Negative lookahead $(?!R)$ and lookbehind $(?<!R)$

Regular Expressions with Lookarounds

- ▶ Lookaround conditions do not consume any characters

Regular Expressions with Lookarounds

- ▶ Lookaround conditions do not consume any characters
- ▶ They describe a context in which a match should appear

Regular Expressions with Lookarounds

- ▶ Lookaround conditions do not consume any characters
- ▶ They describe a context in which a match should appear
- ▶ Given a word "aAbc"

Regular Expressions with Lookarounds

- ▶ Lookaround conditions do not consume any characters
- ▶ They describe a context in which a match should appear
- ▶ Given a word "aAbc"
 - ▶ a *location* is of the form ("aAb", "c")

Regular Expressions with Lookarounds

- ▶ Lookaround conditions do not consume any characters
- ▶ They describe a context in which a match should appear
- ▶ Given a word "aAbc"
 - ▶ a *location* is of the form ("aAb", "c")
 - ▶ In our setting, both the derivative and nullability functions take a location rather than a character

Regular Expressions with Lookarounds

- ▶ Lookaround conditions do not consume any characters
- ▶ They describe a context in which a match should appear
- ▶ Given a word "aAbc"
 - ▶ a *location* is of the form ("aAb", "c")
 - ▶ In our setting, both the derivative and nullability functions take a location rather than a character
 - ▶ a *span* is of the form ("a", "Ab", "c")

Regular Expressions with Lookarounds

- ▶ Lookaround conditions do not consume any characters
- ▶ They describe a context in which a match should appear
- ▶ Given a word "aAbc"
 - ▶ a *location* is of the form ("aAb", "c")
 - ▶ In our setting, both the derivative and nullability functions take a location rather than a character
 - ▶ a *span* is of the form ("a", "Ab", "c")

Regular Expressions with Lookarounds

- ▶ Lookaround conditions do not consume any characters
- ▶ They describe a context in which a match should appear
- ▶ Given a word "aAbc"
 - ▶ a *location* is of the form ("aAb", "c")
 - ▶ In our setting, both the derivative and nullability functions take a location rather than a character
 - ▶ a *span* is of the form ("a", "Ab", "c")

Semantics

$$(xs, ys, zs) \models (?=R) \iff ys = \epsilon \wedge (xs, zs, \epsilon) \models R \cdot \top^*$$

Example: $R = (?=\psi_{upper})$ and $s = \text{"aAbc"}$

Then $(\text{"a"}, \epsilon, \text{"Abc"}) \models (?=\psi_{upper})$

since $(\text{"a"}, \text{"Abc"}, \epsilon) \models \psi_{upper} \cdot \top^*$ is a valid *future match*

Transition terms and symbolic derivatives

A symbolic derivative is a transition term of the following type:

```
inductive TTerm ( $\alpha$  : Type) : Type where  
| Leaf : RE  $\alpha \rightarrow$  TTerm  $\alpha$   
| Node : RE  $\alpha \rightarrow$  TTerm  $\alpha \rightarrow$  TTerm  $\alpha \rightarrow$  TTerm  $\alpha$ 
```


Transition terms and symbolic derivatives

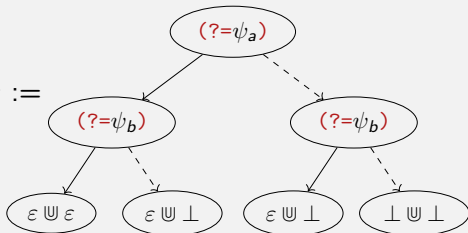
A symbolic derivative is a transition term of the following type:

```
inductive TTerm ( $\alpha$  : Type) : Type where  
| Leaf : RE  $\alpha \rightarrow$  TTerm  $\alpha$   
| Node : RE  $\alpha \rightarrow$  TTerm  $\alpha \rightarrow$  TTerm  $\alpha \rightarrow$  TTerm  $\alpha$ 
```

Example

Let ψ_a and ψ_b be atomic predicates.

$\delta(\psi_a \sqcup \psi_b) : \text{TTerm } \alpha :=$



Transition terms and symbolic derivatives

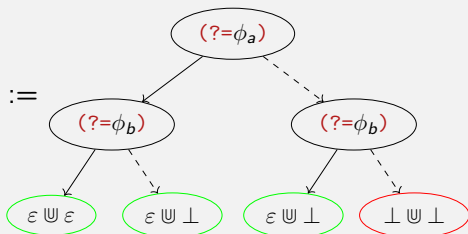
A symbolic derivative is a transition term i.e. trees of regexes

```
inductive TTerm ( $\alpha$  : Type) : Type where  
| Leaf : RE  $\alpha$   $\rightarrow$  TTerm  $\alpha$   
| Node : RE  $\alpha$   $\rightarrow$  TTerm  $\alpha$   $\rightarrow$  TTerm  $\alpha$   $\rightarrow$  TTerm  $\alpha$ 
```

Example

Let ϕ_a and ϕ_b be atomic predicates.

$\delta(\phi_a \sqcup \phi_b) : \text{TTerm } \alpha :=$



Semantics of transition terms

- ▶ Transition term \approx a function from locations to regexes
- ▶ ... but it postpones all the nullability tests
- ▶ We define the evaluation function of type

$\text{Loc } \sigma \rightarrow \text{TTerm } \alpha \rightarrow \text{RE } \alpha$

$$L[x] := L \quad (R, f, g)[x] := \begin{cases} f[x], & \text{if } \text{null } R \ x; \\ g[x], & \text{otherwise.} \end{cases}$$

Symbolic derivatives

Let $\ell \in \text{LA}$, $\psi \in \alpha$ then $\delta : \text{RE } \alpha \rightarrow \text{TTerm } \alpha$

$$\delta \varepsilon := \perp$$

$$\delta \psi := ((\textcolor{red}{?}=\psi), \varepsilon, \perp)$$

$$\delta (L \sqcap R) := \delta L \sqcap \delta R$$

$$\delta (L \sqcup R) := \delta L \sqcup \delta R$$

$$\delta \ell := \perp$$

$$\delta (L \cdot R) := (L, \delta L \cdot R \sqcup \delta R, \delta L \cdot R)$$

$$\delta (\sim R) := \sim(\delta R)$$

$$\delta (R^*) := \delta R \cdot R^*$$

Symbolic derivatives

Let $\ell \in \text{LA}$, $\psi \in \alpha$ then $\delta : \text{RE } \alpha \rightarrow \text{TTerm } \alpha$

$$\delta \varepsilon := \perp$$

$$\delta \ell := \perp$$

$$\delta \psi := ((\text{?}=\psi), \varepsilon, \perp)$$

$$\delta (L \cdot R) := (L, \delta L \cdot R \uplus \delta R, \delta L \cdot R)$$

$$\delta (L \sqcap R) := \delta L \sqcap \delta R$$

$$\delta (\sim R) := \sim(\delta R)$$

$$\delta (L \uplus R) := \delta L \uplus \delta R$$

$$\delta (R^*) := \delta R \cdot R^*$$

We show the equivalence of symbolic and location-based derivatives:

Theorem 1. $\forall x \in \text{Loc}, R \in \text{RE } \alpha : (\delta R)[x] = \text{der } R \ x$

Symbolic derivatives

Let $\ell \in \text{LA}$, $\psi \in \alpha$ then $\delta : \text{RE } \alpha \rightarrow \text{TTerm } \alpha$

$$\delta \varepsilon := \perp$$

$$\delta \ell := \perp$$

$$\delta \psi := ((\text{?}=\psi), \varepsilon, \perp)$$

$$\delta (L \cdot R) := (L, \delta L \cdot R \uplus \delta R, \delta L \cdot R)$$

$$\delta (L \sqcap R) := \delta L \sqcap \delta R$$

$$\delta (\sim R) := \sim(\delta R)$$

$$\delta (L \uplus R) := \delta L \uplus \delta R$$

$$\delta (R^*) := \delta R \cdot R^*$$

We show the equivalence of symbolic and location-based derivatives:

Theorem 1. $\forall x \in \text{Loc}, R \in \text{RE } \alpha : (\delta R)[x] = \text{der } R \ x$

→ from now on, we can just work with the symbolic definition.

Iterated derivatives

- We can now compute the immediate derivatives of R :

`lvs : TTerm α \rightarrow RE α`

`step (R : RE α) : List (RE α) := lvs (δ R)`

Iterated derivatives

- We can now compute the immediate derivatives of R :

$\text{lvs} : \text{TTerm } \alpha \rightarrow \text{RE } \alpha$

$\text{step } (R : \text{RE } \alpha) : \text{List } (\text{RE } \alpha) := \text{lvs } (\delta R)$

- The step function is well-behaved wrt operations on $\text{RE } \alpha$:

$$\text{step } (L \sqcap R) = \text{step } L \sqcap \text{step } R$$

$$\text{step } (L \sqcup R) = \text{step } L \sqcup \text{step } R$$

$$\text{step } (\sim R) = \sim(\text{step } R)$$

$$\text{step } (L \cdot R) = \text{step } L \cdot R \sqcup \text{step } R \text{ ++ } \text{step } L \cdot R$$

$$\text{step } (R^*) = \text{step } R \cdot R^*$$

Iterated derivatives

- ▶ We can now compute the immediate derivatives of R :

$\text{lvs} : \text{TTerm } \alpha \rightarrow \text{RE } \alpha$

$\text{step } (R : \text{RE } \alpha) : \text{List } (\text{RE } \alpha) := \text{lvs } (\delta R)$

- ▶ The step function is well-behaved wrt operations on $\text{RE } \alpha$:

$$\text{step } (L \sqcap R) = \text{step } L \sqcap \text{step } R$$

$$\text{step } (L \sqcup R) = \text{step } L \sqcup \text{step } R$$

$$\text{step } (\sim R) = \sim(\text{step } R)$$

$$\text{step } (L \cdot R) = \text{step } L \cdot R \sqcup \text{step } R ++ \text{step } L \cdot R$$

$$\text{step } (R^*) = \text{step } R \cdot R^*$$

- ▶ We compute the n -th derivatives (words of length n):

$\text{steps} : \text{RE } \alpha \rightarrow \text{Nat} \rightarrow \text{List } (\text{RE } \alpha)$

Finiteness of the state space

- ▶ **Classical approach** (Brzozowski/DFA construction)
 - ▶ $\mathcal{D}er(R) = \{der_w(R) \mid w \in \Sigma^*\} / \cong$
 - ▶ \cong is the equivalence induced by ACI for union \cup

Finiteness of the state space

- ▶ **Classical approach** (Brzozowski/DFA construction)
 - ▶ $\mathcal{D}er(R) = \{der_w(R) \mid w \in \Sigma^*\} / \cong$
 - ▶ \cong is the equivalence induced by ACI for union \cup
- ▶ **Antimirov (partial) derivatives** (NFA construction)

Finiteness of the state space

- ▶ **Classical approach** (Brzozowski/DFA construction)
 - ▶ $\mathcal{D}er(R) = \{der_w(R) \mid w \in \Sigma^*\} / \cong$
 - ▶ \cong is the equivalence induced by ACI for union \cup
- ▶ **Antimirov (partial) derivatives** (NFA construction)
 - ▶ Derivative function returns a set of expressions

Finiteness of the state space

- ▶ **Classical approach** (Brzozowski/DFA construction)
 - ▶ $\mathcal{D}er(R) = \{der_w(R) \mid w \in \Sigma^*\} / \cong$
 - ▶ \cong is the equivalence induced by ACI for union \cup
- ▶ **Antimirov (partial) derivatives** (NFA construction)
 - ▶ Derivative function returns a set of expressions
 - ▶ Proving finiteness is more straightforward but hard to deal with intersection and complement

Finiteness of the state space

- ▶ **Classical approach** (Brzozowski/DFA construction)
 - ▶ $\mathcal{D}er(R) = \{der_w(R) \mid w \in \Sigma^*\} / \cong$
 - ▶ \cong is the equivalence induced by ACI for union \cup
- ▶ **Antimirov (partial) derivatives** (NFA construction)
 - ▶ Derivative function returns a set of expressions
 - ▶ Proving finiteness is more straightforward but hard to deal with intersection and complement
- ▶ **Our approach**

Finiteness of the state space

- ▶ **Classical approach** (Brzozowski/DFA construction)
 - ▶ $\mathcal{D}er(R) = \{der_w(R) \mid w \in \Sigma^*\} / \cong$
 - ▶ \cong is the equivalence induced by ACI for union \cup
- ▶ **Antimirov (partial) derivatives** (NFA construction)
 - ▶ Derivative function returns a set of expressions
 - ▶ Proving finiteness is more straightforward but hard to deal with intersection and complement
- ▶ **Our approach**
 - ▶ Follow Antimirov's strategy for finiteness

Finiteness of the state space

- ▶ **Classical approach** (Brzozowski/DFA construction)
 - ▶ $\mathcal{D}er(R) = \{der_w(R) \mid w \in \Sigma^*\} / \cong$
 - ▶ \cong is the equivalence induced by ACI for union \cup
- ▶ **Antimirov (partial) derivatives** (NFA construction)
 - ▶ Derivative function returns a set of expressions
 - ▶ Proving finiteness is more straightforward but hard to deal with intersection and complement
- ▶ **Our approach**
 - ▶ Follow Antimirov's strategy for finiteness
 - ▶ While dealing with the extended class of expressions

Similarity

We define helpers to reason *up-to* a relation R

- ▶ List membership

$$x \in [R] \text{ } ys := \exists y, R \text{ } x \text{ } y \wedge y \in ys$$

- ▶ List inclusion

$$xs \subseteq [R] \text{ } ys := \forall x \in xs, x \in [R] \text{ } ys$$

- ▶ List equality

$$xs = [R] \text{ } ys := xs \subseteq [R] \text{ } ys \wedge ys \subseteq [R] \text{ } xs$$

Similarity

We define helpers to reason *up-to* a relation R

► List membership

$$x \in [R] \text{ } ys := \exists y, R \text{ } x \text{ } y \wedge y \in ys$$

► List inclusion

$$xs \subseteq [R] \text{ } ys := \forall x \in xs, x \in [R] \text{ } ys$$

► List equality

$$xs = [R] \text{ } ys := xs \subseteq [R] \text{ } ys \wedge ys \subseteq [R] \text{ } xs$$

Our relation: the ADI-similarity relation used for quotienting:

$$(L \uplus R) \uplus S \cong L \uplus (R \uplus S)$$

Associativity

$$L \uplus (R \uplus L) \cong L \uplus R$$

right Deduplication

$$R \uplus R \cong R$$

Idempotence

Finiteness of Antimirov derivatives

- Why is proving finiteness easy for Antimirov derivatives?

$$\text{support}(\perp) := \emptyset$$

$$\text{support}(\varepsilon) := \emptyset$$

$$\text{support}(c) := \{\varepsilon\} \text{ with } c \in \Sigma$$

$$\text{support}(L \uplus R) := \text{support}(L) \cup \text{support}(R)$$

$$\text{support}(L \cdot R) := \text{support}(L) \cdot R \cup \text{support}(R)$$

$$\text{support}(R^*) := \text{support}(R) \cdot R^*$$

Finiteness of Antimirov derivatives

- Why is proving finiteness easy for Antimirov derivatives?

$$\text{support}(\perp) := \emptyset$$

$$\text{support}(\varepsilon) := \emptyset$$

$$\text{support}(c) := \{\varepsilon\} \text{ with } c \in \Sigma$$

$$\text{support}(L \uplus R) := \text{support}(L) \cup \text{support}(R)$$

$$\text{support}(L \cdot R) := \text{support}(L) \cdot R \cup \text{support}(R)$$

$$\text{support}(R^*) := \text{support}(R) \cdot R^*$$

- All Antimirov derivatives are contained in the set:

$$\{R\} \cup \text{support}(R)$$

Finiteness of Antimirov derivatives

- Why is proving finiteness easy for Antimirov derivatives?

$$\text{support}(\perp) := \emptyset$$

$$\text{support}(\varepsilon) := \emptyset$$

$$\text{support}(c) := \{\varepsilon\} \text{ with } c \in \Sigma$$

$$\text{support}(L \uplus R) := \text{support}(L) \cup \text{support}(R)$$

$$\text{support}(L \cdot R) := \text{support}(L) \cdot R \cup \text{support}(R)$$

$$\text{support}(R^*) := \text{support}(R) \cdot R^*$$

- All Antimirov derivatives are contained in the set:

$$\{R\} \cup \text{support}(R)$$

- ACI is built into the set representation

Finiteness of Antimirov derivatives

- ▶ Why is proving finiteness easy for Antimirov derivatives?

$$\text{support}(\perp) := \emptyset$$

$$\text{support}(\varepsilon) := \emptyset$$

$$\text{support}(c) := \{\varepsilon\} \text{ with } c \in \Sigma$$

$$\text{support}(L \uplus R) := \text{support}(L) \cup \text{support}(R)$$

$$\text{support}(L \cdot R) := \text{support}(L) \cdot R \cup \text{support}(R)$$

$$\text{support}(R^*) := \text{support}(R) \cdot R^*$$

- ▶ All Antimirov derivatives are contained in the set:

$$\{R\} \cup \text{support}(R)$$

- ▶ ACI is built into the set representation
- ▶ Can we use a similar strategy for Brzozowski-style derivatives?

Constructing the overapproximation

- ▶ What we have: a way to reason about derivatives and their iterated forms to describe all states reachable from R

Constructing the overapproximation

- ▶ What we have: a way to reason about derivatives and their iterated forms to describe all states reachable from R
- ▶ We have to show finiteness of this set

Constructing the overapproximation

- ▶ What we have: a way to reason about derivatives and their iterated forms to describe all states reachable from R
- ▶ We have to show finiteness of this set
- ▶ Solution: **finite overapproximation** (modulo ADI)

Constructing the overapproximation

- ▶ What we have: a way to reason about derivatives and their iterated forms to describe all states reachable from R
- ▶ We have to show finiteness of this set
- ▶ Solution: **finite overapproximation** (modulo ADI)

$$\begin{array}{ccc} a \uplus b \cdot c & \xrightarrow{\text{der}_a} & \varepsilon \uplus \perp \cdot c \\ \downarrow \text{pieces} & \swarrow \ni & \\ [\perp, \varepsilon, a] \uplus [\perp, \varepsilon, c, \perp \cdot c, \varepsilon \cdot c, b \cdot c] & & \end{array}$$

Constructing the overapproximation

- ▶ What we have: a way to reason about derivatives and their iterated forms to describe all states reachable from R
- ▶ We have to show finiteness of this set
- ▶ Solution: **finite overapproximation** (modulo ADI)

$$\begin{array}{ccc} a \uplus b \cdot c & \xrightarrow{\text{der}_a} & \varepsilon \uplus \perp \cdot c \\ \text{pieces} \downarrow & \nearrow \ni & \\ [\perp, \varepsilon, a] \uplus [\perp, \varepsilon, c, \perp \cdot c, \varepsilon \cdot c, b \cdot c] & & \end{array}$$

- ▶ One step: ε and $\perp \cdot c$ are contained in *pieces* ($a \uplus b \cdot c$)

Constructing the overapproximation

- ▶ What we have: a way to reason about derivatives and their iterated forms to describe all states reachable from R
- ▶ We have to show finiteness of this set
- ▶ Solution: **finite overapproximation** (modulo ADI)

$$\begin{array}{ccc} a \uplus b \cdot c & \xrightarrow{\text{der}_a} & \varepsilon \uplus \perp \cdot c \\ \text{pieces} \downarrow & \nearrow \supseteq & \\ [\perp, \varepsilon, a] \uplus [\perp, \varepsilon, c, \perp \cdot c, \varepsilon \cdot c, b \cdot c] & & \end{array}$$

- ▶ One step: ε and $\perp \cdot c$ are contained in *pieces* ($a \uplus b \cdot c$)
- ▶ **Key idea:** all derivatives can be given as union of pieces

Pieces

- ▶ We don't have commutativity of union so we have to consider all permutations of a list:

$$\oplus[a, b] = [a, a \cup b, b \cup a, b]$$

Pieces

- ▶ We don't have commutativity of union so we have to consider all permutations of a list:

$$\oplus[a, b] = [a, a \cup b, b \cup a, b]$$

- ▶ For intersection we use the Cartesian product:

$$\text{productWith } (\cdot + \cdot) [1,2] [3,4,5] = [4,5,6,5,6,7]$$

Pieces

- ▶ We don't have commutativity of union so we have to consider all permutations of a list:

$$\oplus[a, b] = [a, a \cup b, b \cup a, b]$$

- ▶ For intersection we use the Cartesian product:

$$\text{productWith } (\cdot + \cdot) [1,2] [3,4,5] = [4,5,6,5,6,7]$$

Pieces

- ▶ We don't have commutativity of union so we have to consider all permutations of a list:

$$\oplus[a, b] = [a, a \uplus b, b \uplus a, b]$$

- ▶ For intersection we use the Cartesian product:

$$\text{productWith } (\cdot + \cdot) [1,2] [3,4,5] = [4,5,6,5,6,7]$$

```
def pieces : RE  $\alpha$   $\rightarrow$  List (RE  $\alpha$ )
|  $\varepsilon$       => [ $\varepsilon$ , Pred  $\perp$ ]
| Pred  $\varphi$  => [Pred  $\varphi$ ,  $\varepsilon$ , Pred  $\perp$ ]
|  $?= r$     => [ $?= r$ ,  $\varepsilon$ , Pred  $\perp$ ] | ...
|  $l \uplus r$  => pieces l ++ pieces r
|  $l \cap r$  => productWith ( $\cdot \cap \cdot$ )  $\oplus$ (pieces l)  $\oplus$ (pieces r)
|  $\sim r$     => map ( $\sim \cdot$ )  $\oplus$ (pieces r)
|  $l \cdot r$   => map ( $\cdot \cdot r$ )  $\oplus$ (pieces l) ++ pieces r
|  $r^*$      =>  $r^* ::$  map ( $\cdot \cdot r^*$ )  $\oplus$ (pieces r)
```


Main theorem

1. Reflexivity:

$\forall r,$
 $\exists xs, \text{toSum } xs \cong r \wedge xs \in \text{neSublists } (\text{pieces } r)$

2. Transitivity:

$e \in \text{pieces } f$
 $\rightarrow f \in \text{pieces } g$
 $\rightarrow e \in [(\cdot \cong \cdot)] \text{pieces } g$

3. One-step reconstruction:

$\forall r d, d \in \text{step } r$
 $\exists xs, \text{toSum } xs \cong d \wedge xs \in \text{neSubsets } (\text{pieces } r)$

Main theorem

1. Reflexivity:

$\forall r,$
 $\exists xs, \text{toSum } xs \cong r \wedge xs \in \text{neSublists } (\text{pieces } r)$

2. Transitivity:

$e \in \text{pieces } f$
 $\rightarrow f \in \text{pieces } g$
 $\rightarrow e \in [(\cdot \cong \cdot)] \text{pieces } g$

3. One-step reconstruction:

$\forall r d, d \in \text{step } r$
 $\exists xs, \text{toSum } xs \cong d \wedge xs \in \text{neSubsets } (\text{pieces } r)$

- **Main result:** every iterated derivative of R can be reconstructed as a sum of regexes from pieces R

Main theorem

1. Reflexivity:

$\forall r,$
 $\exists xs, \text{toSum } xs \cong r \wedge xs \in \text{neSublists } (\text{pieces } r)$

2. Transitivity:

$e \in \text{pieces } f$
 $\rightarrow f \in \text{pieces } g$
 $\rightarrow e \in [(\cdot \cong \cdot)] \text{pieces } g$

3. One-step reconstruction:

$\forall r d, d \in \text{step } r$
 $\exists xs, \text{toSum } xs \cong d \wedge xs \in \text{neSubsets } (\text{pieces } r)$

- **Main result:** every iterated derivative of R can be reconstructed as a sum of regexes from pieces R

theorem finiteness [DecidableEq α] {r : RE α } :
 $\exists (xs : \text{List } (\text{RE } \alpha)),$
 $\forall \{n : \mathbb{N}\}, \text{steps } r \ n \subseteq [(\cdot \cong \cdot)] \ xs$

Main theorem

1. Reflexivity:

$\forall r,$
 $\exists xs, \text{toSum } xs \cong r \wedge xs \in \text{neSublists } (\text{pieces } r)$

2. Transitivity:

$e \in \text{pieces } f$
 $\rightarrow f \in \text{pieces } g$
 $\rightarrow e \in [(\cdot \cong \cdot)] \text{pieces } g$

3. One-step reconstruction:

$\forall r d, d \in \text{step } r$
 $\exists xs, \text{toSum } xs \cong d \wedge xs \in \text{neSubsets } (\text{pieces } r)$

- **Main result:** every iterated derivative of R can be reconstructed as a sum of regexes from pieces R

theorem finiteness [DecidableEq α] {r : RE α } :
 $\exists (xs : \text{List } (\text{RE } \alpha)),$
 $\forall \{n : \mathbb{N}\}, \text{steps } r \ n \subseteq [(\cdot \cong \cdot)] xs$

- The witness is $xs := \oplus(\text{pieces } R)$

Previous work

- ▶ Coquand & Siles (2011), Nipkow & Traytel (2014) use canonical/normal forms

Previous work

- ▶ Coquand & Siles (2011), Nipkow & Traytel (2014) use canonical/normal forms
 - ▶ Two versions of the normalisation function: one that only implements ACI and one which implements more aggressive simplifications

Previous work

- ▶ Coquand & Siles (2011), Nipkow & Traytel (2014) use canonical/normal forms
 - ▶ Two versions of the normalisation function: one that only implements ACI and one which implements more aggressive simplifications
- ▶ We instead compute a finite overapproximation of all reachable derivatives

Previous work

- ▶ Coquand & Siles (2011), Nipkow & Traytel (2014) use canonical/normal forms
 - ▶ Two versions of the normalisation function: one that only implements ACI and one which implements more aggressive simplifications
- ▶ We instead compute a finite overapproximation of all reachable derivatives
- ▶ (Moreira et al., 2012) avoid the need for normalisation modulo ACI by using Antimirov derivatives

Previous work

- ▶ Coquand & Siles (2011), Nipkow & Traytel (2014) use canonical/normal forms
 - ▶ Two versions of the normalisation function: one that only implements ACI and one which implements more aggressive simplifications
- ▶ We instead compute a finite overapproximation of all reachable derivatives
- ▶ (Moreira et al., 2012) avoid the need for normalisation modulo ACI by using Antimirov derivatives
- ▶ We take inspiration from the Antimirov finiteness proof, but adapt it to handle intersection and negation

Simplifications

Which simplifications preserve the finiteness result?

Simplifications

Which simplifications preserve the finiteness result?

```
def NonIncreasing (f : RE  $\alpha$   $\rightarrow$  RE  $\alpha$ ) : Prop :=  
   $\forall$  r, pieces (f r)  $\subseteq$  pieces r
```

Simplifications

Which simplifications preserve the finiteness result?

```
def NonIncreasing (f : RE  $\alpha$   $\rightarrow$  RE  $\alpha$ ) : Prop :=  
   $\forall$  r, pieces (f r)  $\subseteq$  pieces r
```

```
 $\forall$  (f : RE  $\alpha$   $\rightarrow$  RE  $\alpha$ ) r,  
  NonIncreasing f  
   $\rightarrow$  map f (step r)  $\subseteq$  [ ( $\cdot \cong \cdot$ ) ]  $\oplus$  (pieces r)
```

Simplifications

Which simplifications preserve the finiteness result?

```
def NonIncreasing (f : RE  $\alpha$   $\rightarrow$  RE  $\alpha$ ) : Prop :=  
   $\forall$  r, pieces (f r)  $\subseteq$  pieces r
```

```
 $\forall$  (f : RE  $\alpha$   $\rightarrow$  RE  $\alpha$ ) r,  
  NonIncreasing f  
   $\rightarrow$  map f (step r)  $\subseteq$  [ ( $\cdot \cong \cdot$ ) ]  $\oplus$  (pieces r)
```

Allowed simplifications

$$\begin{array}{lll} \perp \sqcup s \rightsquigarrow s & \sim \perp \sqcup s \rightsquigarrow \sim \perp & \varepsilon \cdot s \rightsquigarrow s \\ r \sqcup \perp \rightsquigarrow r & r \sqcup \sim \perp \rightsquigarrow \sim \perp & r \cdot \perp \rightsquigarrow \perp \\ (r \sqcup s) \cdot t \rightsquigarrow r \cdot t \sqcup s \cdot t & & \end{array}$$

► $r \cdot s \rightsquigarrow s$ and $r \sqcup s \rightsquigarrow r$ are allowed

Simplifications

Which simplifications preserve the finiteness result?

```
def NonIncreasing (f : RE  $\alpha$   $\rightarrow$  RE  $\alpha$ ) : Prop :=  
   $\forall$  r, pieces (f r)  $\subseteq$  pieces r
```

```
 $\forall$  (f : RE  $\alpha$   $\rightarrow$  RE  $\alpha$ ) r,  
  NonIncreasing f  
   $\rightarrow$  map f (step r)  $\subseteq$  [ ( $\cdot \cong \cdot$ ) ]  $\oplus$  (pieces r)
```

Allowed simplifications

$$\begin{array}{lll} \perp \sqcup s \rightsquigarrow s & \sim \perp \sqcup s \rightsquigarrow \sim \perp & \varepsilon \cdot s \rightsquigarrow s \\ r \sqcup \perp \rightsquigarrow r & r \sqcup \sim \perp \rightsquigarrow \sim \perp & r \cdot \perp \rightsquigarrow \perp \\ (r \sqcup s) \cdot t \rightsquigarrow r \cdot t \sqcup s \cdot t & & \end{array}$$

- ▶ $r \cdot s \rightsquigarrow s$ and $r \sqcup s \rightsquigarrow r$ are allowed
- ▶ $r \cdot s \rightsquigarrow r$ is **not** allowed

Conclusion

We formally prove in Lean that the set of symbolic derivatives of regexes with lookarounds is finite modulo ADI

- ▶ Almost 2000 loc of Lean, modularly reusing previous work

Conclusion

We formally prove in Lean that the set of symbolic derivatives of regexes with lookarounds is finite modulo ADI

- ▶ Almost 2000 loc of Lean, modularly reusing previous work
- ▶ We show which simplifications can be applied to derivatives, preserving finiteness

Conclusion

We formally prove in Lean that the set of symbolic derivatives of regexes with lookarounds is finite modulo ADI

- ▶ Almost 2000 loc of Lean, modularly reusing previous work
- ▶ We show which simplifications can be applied to derivatives, preserving finiteness
- ▶ No assumption that the alphabet is finite; the alphabet algebra can even be undecidable or semidecidable

Conclusion

We formally prove in Lean that the set of symbolic derivatives of regexes with lookarounds is finite modulo ADI

- ▶ Almost 2000 loc of Lean, modularly reusing previous work
- ▶ We show which simplifications can be applied to derivatives, preserving finiteness
- ▶ No assumption that the alphabet is finite; the alphabet algebra can even be undecidable or semidecidable
- ▶ How to make this into a reusable framework for finiteness? (e.g. for other regex classes/logics)

Conclusion

We formally prove in Lean that the set of symbolic derivatives of regexes with lookarounds is finite modulo ADI

- ▶ Almost 2000 loc of Lean, modularly reusing previous work
- ▶ We show which simplifications can be applied to derivatives, preserving finiteness
- ▶ No assumption that the alphabet is finite; the alphabet algebra can even be undecidable or semidecidable
- ▶ How to make this into a reusable framework for finiteness? (e.g. for other regex classes/logics)

Conclusion

We formally prove in Lean that the set of symbolic derivatives of regexes with lookarounds is finite modulo ADI

- ▶ Almost 2000 loc of Lean, modularly reusing previous work
- ▶ We show which simplifications can be applied to derivatives, preserving finiteness
- ▶ No assumption that the alphabet is finite; the alphabet algebra can even be undecidable or semidecidable
- ▶ How to make this into a reusable framework for finiteness? (e.g. for other regex classes/logics)

Thank you!

`github.com/ezhuchko/finiteness-derivatives`