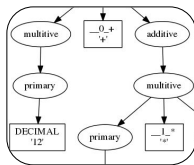


Recapitulation from Previous Lectures

Recap: Compilers

```
def show: String = {  
  val vars =  
    vars.distinct  
  val ctx =  
    vars.zipWithIndex.map {  
      case (tv, idx) =>  
        def nme = {  
          assert(idx <= 'z'  
            - 'a', "TODO handle case  
              of not enough chars")  
        }  
    }  
}
```

Program source



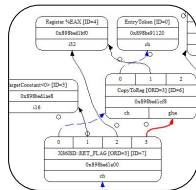
Abstract syntax tree (AST)



type checking,
semantic analysis

```
.globl "_add_forty_two<Int32>:Int32"  
.align 4, 0x90  
"_add_forty_two<Int32>:Int32":  
.cfi_startproc  
pushq %rbp  
Ltmp1992:  
.cfi_def_cfa_offset 16  
Ltmp1993:  
.cfi_offset %rbp, -16  
movq %rsp, %rbp  
Ltmp1994:  
.cfi_def_cfa_register %rbp  
addl $42, %edi  
movl %edi, %eax
```

Assembly



Internal representation



optimization

Recap: Compiler Phases

characters

↓ lexical analyzer

words

↓ parser

trees

↓ name analyzer

graphs

↓ type checker

graphs

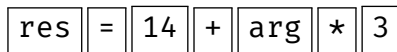
↓ intermediate code generator

intermediate code e.g. LLVM bitcode, JVM bytecode, Web Assembly

↓ JIT compiler or platform-specific back end

machine code e.g. x86, ARM, RISC-V

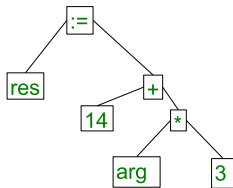
res = 14 + arg * 3



Assign(res, Plus(C(14), Times(V(arg),C(3))))

(variables mapped to declarations)

Assign(res:Int, Plus(C(14), Times(V(arg):Int,C(3)))):Unit



Recap: Transforming Text Into a Tree

characters `res = 14 + arg * 3`

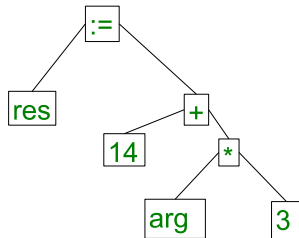
↓ lexical analyzer

words

res	=	14	+	arg	*	3
-----	---	----	---	-----	---	---

↓ parser

trees `Assign(res, Plus(C(14), Times(V(arg), C(3))))`



First two phases:

1. lexical analyzer (lexer): sequence of characters \rightarrow sequence of words
2. syntax analyzer (parser): sequence of words \rightarrow tree

We will study *linear-time algorithms* for these problems.

We start with the underlying *theory of formal languages*.

Recap: Words

Let A be an alphabet $\{a, b, c, \dots\}$

We define words of length n , denoted A^n , as follows:

$A^0 = \{\varepsilon\}$ (only one word of length zero, always denoted ε)

For $n > 0$, $A^n = \{aw \mid w \in A^{n-1}\}$

Notation ' $a\varepsilon$ ' abbreviated to ' a '.

Recap: Words

Let A be an alphabet $\{a, b, c, \dots\}$

We define words of length n , denoted A^n , as follows:

$A^0 = \{\varepsilon\}$ (only one word of length zero, always denoted ε)

For $n > 0$, $A^n = \{aw \mid w \in A^{n-1}\}$

Notation ' $a\varepsilon$ ' abbreviated to ' a '.

Concatenation: $u \cdot v$, also written uv , is associative.

We can decompose words arbitrarily, as in $w = ua$ (i.e., $w = u \cdot a\varepsilon$) if $|w| > 0$

Referring to letters by index: $w_{(0)} = \mathbf{1}$ $w_{(1)} = \mathbf{0}$ $w_{(2)} = \mathbf{1}$ $w_{(3)} = \mathbf{1}$

Recap: Words

Let A be an alphabet $\{a, b, c, \dots\}$

We define words of length n , denoted A^n , as follows:

$A^0 = \{\varepsilon\}$ (only one word of length zero, always denoted ε)

For $n > 0$, $A^n = \{aw \mid w \in A^{n-1}\}$

Notation ' $a\varepsilon$ ' abbreviated to ' a '.

Concatenation: $u \cdot v$, also written uv , is associative.

We can decompose words arbitrarily, as in $w = ua$ (i.e., $w = u \cdot a\varepsilon$) if $|w| > 0$

Referring to letters by index: $w_{(0)} = \mathbf{1}$ $w_{(1)} = \mathbf{0}$ $w_{(2)} = \mathbf{1}$ $w_{(3)} = \mathbf{1}$

Set of all words: $A^* = \bigcup_{n \geq 0} A^n$

which means: $w \in A^*$ if and only iff there exists n such that $w \in A^n$.

Recap: Languages

A *language* over alphabet A is a set $L \subseteq A^*$.

Recap: Languages

A *language* over alphabet A is a set $L \subseteq A^*$.

Examples for $A = \{0, 1\}$:

- ▶ empty language \emptyset ;
- ▶ finite languages like $L = \{1, 10, 1001\}$;
- ▶ language L described by a characteristic function f , i.e., $L = \{w \in A^* \mid f(w)\}$

Recap: Languages

A *language* over alphabet A is a set $L \subseteq A^*$.

Examples for $A = \{0, 1\}$:

- ▶ empty language \emptyset ;
- ▶ finite languages like $L = \{1, 10, 1001\}$;
- ▶ language L described by a characteristic function f , i.e., $L = \{w \in A^* \mid f(w)\}$

Empty language \emptyset and language of the empty word $\{\varepsilon\}$ are *very different*.

$\forall L$. we have $L\emptyset = \emptyset$ but $L\{\varepsilon\} = L$

Recap: Languages

A *language* over alphabet A is a set $L \subseteq A^*$.

Examples for $A = \{0, 1\}$:

- ▶ empty language \emptyset ;
- ▶ finite languages like $L = \{1, 10, 1001\}$;
- ▶ language L described by a characteristic function f , i.e., $L = \{w \in A^* \mid f(w)\}$

Empty language \emptyset and language of the empty word $\{\varepsilon\}$ are *very different*.

$\forall L$. we have $L\emptyset = \emptyset$ but $L\{\varepsilon\} = L$

Operations on languages:

- ▶ Set operations: union (\cup), intersection (\cap), difference (\setminus), etc.
- ▶ Language concatenation: $L_1 \cdot L_2 = \{w_1 \cdot w_2 \mid w_1 \in L_1, w_2 \in L_2\}$, also written $L_1 L_2$
- ▶ Language exponentiation: $L^0 = \{\varepsilon\}$ and $L^{n+1} = L \cdot L^n$

Recap: Regular Expressions

Syntax of regular expressions: $e ::= \emptyset \mid \varepsilon \mid c \mid (e_1 \mid e_2) \mid e_1 e_2 \mid e^*$

Recap: Regular Expressions

Syntax of regular expressions: $e ::= \emptyset \mid \varepsilon \mid c \mid (e_1 \mid e_2) \mid e_1 e_2 \mid e^*$

The *semantics* of regular expression e is the language $L(e)$ – also written L_e – where:

$$L(\emptyset) = \emptyset$$

$$L(\varepsilon) = \{\varepsilon\}$$

$$L(c) = \{c\} \quad (c \in A)$$

$$L(e_1 \mid e_2) = L(e_1) \cup L(e_2)$$

$$L(e_1 e_2) = L(e_1) \cdot L(e_2)$$

$$L(e^*) = L(e)^*$$

Example: $letter(letter \mid digit)^*$ where $letter = a \mid b \mid c \mid \dots$ and $digit = 0 \mid 1 \mid 2 \dots \mid 9$

Recap: Closed Operations on Regular Expressions

Regular languages/expressions are closed under these operations:

- ▶ Shorthands for finite languages, such as $[a..z] = a \mid b \mid \dots \mid z$
- ▶ Optionality $e^? = e \mid \varepsilon$
- ▶ Repeating at least once $e^+ = ee^*$
- ▶ Other repetitions $e^{k..*} = e^k e^*$ and $e^{p..q} = e^p (e^?)^{q-p}$
- ▶ Complementation: $!e$, denoting $A^* \setminus L(e)$ — Q: how to translate?
- ▶ Intersection $e_1 \& e_2 = !(!e \mid !e)$, denoting $L(e_1) \cap L(e_2)$

Recap: First and Nullable (Definitions)

Formal definitions of *first* and *nullable*:

Recap: First and Nullable (Definitions)

Formal definitions of *first* and *nullable*:

$$\text{first}(e) = \{ c \mid \exists w. cw \in L(e) \}$$

Recap: First and Nullable (Definitions)

Formal definitions of *first* and *nullable*:

$$\text{first}(e) = \{ c \mid \exists w. cw \in L(e) \}$$

$$\text{nullable}(e) = \varepsilon \in L(e)$$

Recap: First and Nullable (Definitions)

Formal definitions of *first* and *nullable*:

$$\text{first}(e) = \{ c \mid \exists w. cw \in L(e) \}$$

$$\text{nullable}(e) = \varepsilon \in L(e)$$

Examples:

$$\text{first}(ab^*) = \{a\}$$

$$\text{first}(ab^* \mid c) = \{a, c\}$$

$$\text{first}(a^*b^*c) = \{a, b, c\}$$

$$\text{first}((cb \mid a^*c^*)d^*e) = \{a, b, c, d\}$$

Recap: First and Nullable (Algorithms)

Algorithms for computing *first* and *nullable*:

$$\text{nullable}(\emptyset) = \text{false}$$

$$\text{nullable}(\varepsilon) = \text{true}$$

$$\text{nullable}(a) = \text{false}$$

$$\text{nullable}(e_1|e_2) = \text{nullable}(e_1) \vee \text{nullable}(e_2)$$

$$\text{nullable}(e^*) = \text{true}$$

$$\text{nullable}(e_1e_2) = \text{nullable}(e_1) \wedge \text{nullable}(e_2)$$

Recap: First and Nullable (Algorithms)

Algorithms for computing *first* and *nullable*:

$$\begin{aligned} nullable(\emptyset) &= false \\ nullable(\varepsilon) &= true \\ nullable(a) &= false \\ nullable(e_1|e_2) &= nullable(e_1) \vee nullable(e_2) \\ nullable(e^*) &= true \\ nullable(e_1e_2) &= nullable(e_1) \wedge nullable(e_2) \end{aligned}$$

$$\begin{aligned} first(\emptyset) &= \emptyset \\ first(\varepsilon) &= \emptyset \\ first(a) &= \{a\}, \text{ for } a \in A \\ first(e_1|e_2) &= first(e_1) \cup first(e_2) \\ first(e^*) &= first(e) \\ first(e_1e_2) &= \text{if } (nullable(e_1)) \text{ then } first(e_1) \cup first(e_2) \\ &\quad \text{else } first(e_1) \end{aligned}$$

Recap: Lexers

Definition of a *lexer* (aka *tokenizer*):
an ordered set of n labelled token definitions

$$\langle Token_1 := e_1; Token_2 := e_2; \dots; Token_n := e_n \rangle$$

Disambiguation rules:

Recap: Lexers

Definition of a *lexer* (aka *tokenizer*):
an ordered set of n labelled token definitions

$$\langle Token_1 := e_1; Token_2 := e_2; \dots; Token_n := e_n \rangle$$

Disambiguation rules:

- ▶ Longest-match

Recap: Lexers

Definition of a *lexer* (aka *tokenizer*):
an ordered set of n labelled token definitions

$$\langle Token_1 := e_1; Token_2 := e_2; \dots; Token_n := e_n \rangle$$

Disambiguation rules:

- ▶ Longest-match
- ▶ First-match

Trickiness of Lexing

Consider a lexer for these two classes of tokens:

$$\langle \text{Token}_1 := \text{aaa}; \text{Token}_2 := a^*b \rangle$$

Example run on aaaab:

Trickiness of Lexing

Consider a lexer for these two classes of tokens:

$$\langle \text{Token}_1 := \text{aaa}; \text{Token}_2 := a^*b \rangle$$

Example run on aaaab:

{aaaab Token_1 or Token_2 ?

Trickiness of Lexing

Consider a lexer for these two classes of tokens:

$$\langle \text{Token}_1 := \text{aaa}; \text{Token}_2 := a^*b \rangle$$

Example run on aaaab:

$\{\}$ aaaab Token_1 or Token_2 ?

$\{a\}$ aaab Token_1 or Token_2 ?

Trickiness of Lexing

Consider a lexer for these two classes of tokens:

$$\langle \text{Token}_1 := \text{aaa}; \text{Token}_2 := \text{a}^*\text{b} \rangle$$

Example run on aaaab:

{ }aaaab Token_1 or Token_2 ?

{a}aaab Token_1 or Token_2 ?

{aa}aab Token_1 or Token_2 ?

Trickiness of Lexing

Consider a lexer for these two classes of tokens:

$$\langle \text{Token}_1 := \text{aaa}; \text{Token}_2 := a^*b \rangle$$

Example run on aaaab:

{ }aaaab Token_1 or Token_2 ?

{a}aaab Token_1 or Token_2 ?

{aa}aab Token_1 or Token_2 ?

{aaa}ab **Token** $_1$ or Token_2 ?

Trickiness of Lexing

Consider a lexer for these two classes of tokens:

$$\langle \text{Token}_1 := \text{aaa}; \text{Token}_2 := a^*b \rangle$$

Example run on aaaab:

{ }aaaab Token_1 or Token_2 ?

{a}aaab Token_1 or Token_2 ?

{aa}aab Token_1 or Token_2 ?

{aaa}ab **Token** $_1$ or Token_2 ?

{aaaa}b Token_1 then (Token_1 or Token_2) or still Token_2 ?

Trickiness of Lexing

Consider a lexer for these two classes of tokens:

$$\langle \text{Token}_1 := \text{aaa}; \text{Token}_2 := \text{a}^*\text{b} \rangle$$

Example run on aaaab:

{ }aaaab	Token ₁ or Token ₂ ?
{a}aaab	Token ₁ or Token ₂ ?
{aa}aab	Token ₁ or Token ₂ ?
{aaa}ab	Token ₁ or Token ₂ ?
{aaaa}b	Token ₁ then (Token ₁ or Token ₂) or still Token ₂ ?
{aaaab}	Token₁ then (Token₁ or Token₂) or still Token ₂ !

Trickiness of Lexing

Consider a lexer for these two classes of tokens:

$$\langle \text{Token}_1 := \text{aaa}; \text{Token}_2 := \text{a}^*\text{b} \rangle$$

Example run on aaaab:

{ }aaaab	Token ₁ or Token ₂ ?
{a}aaab	Token ₁ or Token ₂ ?
{aa}aab	Token ₁ or Token ₂ ?
{aaa}ab	Token ₁ or Token ₂ ?
{aaaa}b	Token ₁ then (Token ₁ or Token ₂) or still Token ₂ ?
{aaaab}	Token₁ then (Token₁ or Token₂) or still Token ₂ !

How about a run on aaaaaa? And aaaaaaaaaa? And aaaaaaaaaab?

General Approach to Automatic Lexing

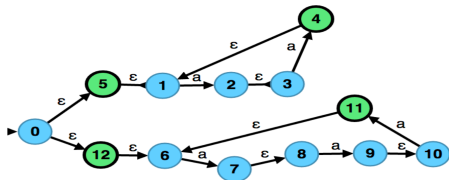
Compile regexps into some sort of *automata* with *transitions* between *states*.

General Approach to Automatic Lexing

Compile regexps into some sort of *automata* with *transitions* between *states*.

Traditional approach:

1. convert to nondeterministic finite-state automaton;
2. perform determinization (can be expensive);
3. run the resulting automaton on input (linear in the input size).

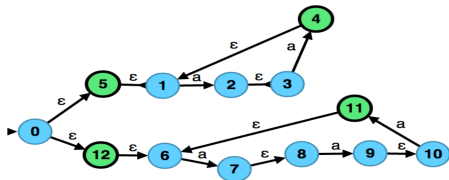


General Approach to Automatic Lexing

Compile regexps into some sort of *automata* with *transitions* between *states*.

Traditional approach:

1. convert to nondeterministic finite-state automaton;
2. perform determinization (can be expensive);
3. run the resulting automaton on input (linear in the input size).



In this course: we'll look at a more straightforward approach based on *derivatives*

Automatically Recognizing Regular Languages

Regular Expression Matcher: First Trial

accepts : $(e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

Regular Expression Matcher: First Trial

$accepts : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

$accepts(e, \varepsilon) =$

$accepts(e, cu) =$

Regular Expression Matcher: First Trial

$accepts : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

$accepts(e, \varepsilon) = nullable(e)$

$accepts(e, cu) =$

Regular Expression Matcher: First Trial

$accepts : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

$accepts(e, \varepsilon) = nullable(e)$

$accepts(e, cu) =$

if $e = \emptyset$,

Regular Expression Matcher: First Trial

$accepts : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

$accepts(e, \varepsilon) = nullable(e)$

$accepts(e, cu) =$

if $e = \emptyset$, **false**

Regular Expression Matcher: First Trial

$accepts : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

$accepts(e, \varepsilon) = nullable(e)$

$accepts(e, cu) =$

if $e = \emptyset$, **false**

if $e = \varepsilon$,

Regular Expression Matcher: First Trial

$accepts : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

$accepts(e, \varepsilon) = nullable(e)$

$accepts(e, cu) =$

if $e = \emptyset$, **false**

if $e = \varepsilon$, **false**

Regular Expression Matcher: First Trial

accepts : $(e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

accepts(e, ε) = *nullable*(e)

accepts(e, cu) =

if $e = \emptyset$, **false**

if $e = \varepsilon$, **false**

if $e = c$,

Regular Expression Matcher: First Trial

$accepts : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

$accepts(e, \varepsilon) = nullable(e)$

$accepts(e, cu) =$

if $e = \emptyset$, **false**

if $e = \varepsilon$, **false**

if $e = c$, **true** **if** $u = \varepsilon$, **false** **otherwise**

Regular Expression Matcher: First Trial

$accepts : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

$accepts(e, \varepsilon) = nullable(e)$

$accepts(e, cu) =$

if $e = \emptyset$, **false**

if $e = \varepsilon$, **false**

if $e = c$, **true** **if** $u = \varepsilon$, **false** **otherwise**

if $e = d$ ($d \neq c$),

Regular Expression Matcher: First Trial

$accepts : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

$accepts(e, \varepsilon) = nullable(e)$

$accepts(e, cu) =$

if $e = \emptyset$, **false**

if $e = \varepsilon$, **false**

if $e = c$, **true** **if** $u = \varepsilon$, **false** **otherwise**

if $e = d$ ($d \neq c$), **false**

Regular Expression Matcher: First Trial

$accepts : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

$accepts(e, \varepsilon) = nullable(e)$

$accepts(e, cu) =$

if $e = \emptyset$, **false**

if $e = \varepsilon$, **false**

if $e = c$, **true** **if** $u = \varepsilon$, **false** **otherwise**

if $e = d$ ($d \neq c$), **false**

if $e = e_1 \mid e_2$,

Regular Expression Matcher: First Trial

$accepts : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

$accepts(e, \varepsilon) = nullable(e)$

$accepts(e, cu) =$

if $e = \emptyset$, **false**

if $e = \varepsilon$, **false**

if $e = c$, **true** **if** $u = \varepsilon$, **false** **otherwise**

if $e = d$ ($d \neq c$), **false**

if $e = e_1 \mid e_2$, $accepts(e_1, w) \vee accepts(e_2, w)$

Regular Expression Matcher: First Trial

$accepts : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

$accepts(e, \varepsilon) = nullable(e)$

$accepts(e, cu) =$

if $e = \emptyset$, **false**

if $e = \varepsilon$, **false**

if $e = c$, **true** **if** $u = \varepsilon$, **false** **otherwise**

if $e = d$ ($d \neq c$), **false**

if $e = e_1 \mid e_2$, $accepts(e_1, w) \vee accepts(e_2, w)$

if $e = e_1 e_2$,

Regular Expression Matcher: First Trial

$accepts : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

$accepts(e, \varepsilon) = nullable(e)$

$accepts(e, cu) =$

if $e = \emptyset$, **false**

if $e = \varepsilon$, **false**

if $e = c$, **true** **if** $u = \varepsilon$, **false** **otherwise**

if $e = d$ ($d \neq c$), **false**

if $e = e_1 \mid e_2$, $accepts(e_1, w) \vee accepts(e_2, w)$

if $e = e_1 e_2$, ...what to do? try all splits of w ? (expensive)

Regular Expression Matcher: First Trial

$accepts : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

$accepts(e, \varepsilon) = nullable(e)$

$accepts(e, cu) =$

if $e = \emptyset$, **false**

if $e = \varepsilon$, **false**

if $e = c$, **true** **if** $u = \varepsilon$, **false** **otherwise**

if $e = d$ ($d \neq c$), **false**

if $e = e_1 \mid e_2$, $accepts(e_1, w) \vee accepts(e_2, w)$

if $e = e_1 e_2$, ...what to do? try all splits of w ? (expensive)

if $e = e_1^*$,

Regular Expression Matcher: First Trial

$accepts : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

$accepts(e, \varepsilon) = nullable(e)$

$accepts(e, cu) =$

if $e = \emptyset$, **false**

if $e = \varepsilon$, **false**

if $e = c$, **true** **if** $u = \varepsilon$, **false** **otherwise**

if $e = d$ ($d \neq c$), **false**

if $e = e_1 \mid e_2$, $accepts(e_1, w) \vee accepts(e_2, w)$

if $e = e_1 e_2$, ...what to do? try all splits of w ? (expensive)

if $e = e_1^*$, ...similar problem here

Regular Expression Matcher: First Trial

accepts : $(e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

accepts(e, ε) = *nullable*(e)

accepts(e, cu) =

if $e = \emptyset$, **false**

if $e = \varepsilon$, **false**

if $e = c$, **true** **if** $u = \varepsilon$, **false** **otherwise**

if $e = d$ ($d \neq c$), **false**

if $e = e_1 \mid e_2$, *accepts*(e_1, w) \vee *accepts*(e_2, w)

if $e = e_1 e_2$, ...what to do? try all splits of w ? (expensive)

if $e = e_1^*$, ...similar problem here

Regular Expression Matcher: First Trial

$accepts : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$

$accepts(e, \varepsilon) = nullable(e)$

$accepts(e, cu) =$

if $e = \emptyset$, **false**

if $e = \varepsilon$, **false**

if $e = c$, **true** **if** $u = \varepsilon$, **false** **otherwise**

if $e = d$ ($d \neq c$), **false**

if $e = e_1 \mid e_2$, $accepts(e_1, w) \vee accepts(e_2, w)$

if $e = e_1 e_2$, ...what to do? try all splits of w ? (expensive)

if $e = e_1^*$, ...similar problem here

We need a concept of “the rest of the regexp after consuming a given letter.”

Brzozowski's Derivatives, Definition

We need a concept of “the rest of the regexp after consuming a given letter.”

Brzowski's Derivatives, Definition

We need a concept of “the rest of the regexp after consuming a given letter.”

Definition (Brzowski Derivative)

The *derivative* of regexp e with respect to letter c , written $\delta^c(e)$, is defined as:

$$L(\delta^c(e)) = \{w \mid cw \in L(e)\}$$

Brzowski's Derivatives, Definition

We need a concept of “the rest of the regexp after consuming a given letter.”

Definition (Brzowski Derivative)

The *derivative* of regexp e with respect to letter c , written $\delta^c(e)$, is defined as:

$$L(\delta^c(e)) = \{w \mid cw \in L(e)\}$$

Theorem

The derivative of a regexp is a regexp.

Brzowski Derivatives, Definition

We need a concept of “the rest of the regexp after consuming a given letter.”

Definition (Brzowski Derivative)

The *derivative* of regexp e with respect to letter c , written $\delta^c(e)$, is defined as:

$$L(\delta^c(e)) = \{w \mid cw \in L(e)\}$$

Theorem

The derivative of a regexp is a regexp.

Proof.

We'll give an algorithm for computing derivatives, which only yields regexps. □

Brzozowski's Derivatives, Examples

Here are a few derivative examples:

$$\delta^a(aaa) =$$

Brzozowski's Derivatives, Examples

Here are a few derivative examples:

$$\delta^a(aaa) = aa$$

$$\delta^a(ab \mid ac \mid da) =$$

Brzozowski's Derivatives, Examples

Here are a few derivative examples:

$$\delta^a(aaa) = aa$$

$$\delta^a(ab|ac|da) = b|c$$

$$\delta^a((ab)^*) =$$

Brzozowski's Derivatives, Examples

Here are a few derivative examples:

$$\delta^a(aaa) = aa$$

$$\delta^a(ab|ac|da) = b|c$$

$$\delta^a((ab)^*) = b(ab)^*$$

$$\delta^a((ab|c)^*ad) =$$

Brzozowski's Derivatives, Examples

Here are a few derivative examples:

$$\delta^a(aaa) = aa$$

$$\delta^a(ab|ac|da) = b|c$$

$$\delta^a((ab)^*) = b(ab)^*$$

$$\delta^a((ab|c)^*ad) = b(ab|c)^*ad|d$$

Brzozowski's Derivatives, Algorithm

Derivative of a regexp e with respect to letter c , written $\delta^c(e)$, can be computed as:

$$\delta^c(\emptyset) = \emptyset$$

$$\delta^c(\epsilon) =$$

$$\delta^c(d) =$$

$$\delta^c(e_1 \mid e_2) =$$

$$\delta^c(e_1 e_2) =$$

$$\delta^c(e_1^*) =$$

Brzozowski's Derivatives, Algorithm

Derivative of a regexp e with respect to letter c , written $\delta^c(e)$, can be computed as:

$$\delta^c(\emptyset) = \emptyset$$

$$\delta^c(\epsilon) = \emptyset$$

$$\delta^c(d) =$$

$$\delta^c(e_1 \mid e_2) =$$

$$\delta^c(e_1 e_2) =$$

$$\delta^c(e_1^*) =$$

Brzozowski's Derivatives, Algorithm

Derivative of a regexp e with respect to letter c , written $\delta^c(e)$, can be computed as:

$$\delta^c(\emptyset) = \emptyset$$

$$\delta^c(\varepsilon) = \emptyset$$

$$\delta^c(d) = \begin{cases} \varepsilon & \text{if } d = c \\ \emptyset & \text{if } d \neq c \end{cases}$$

$$\delta^c(e_1 \mid e_2) =$$

$$\delta^c(e_1 e_2) =$$

$$\delta^c(e_1^*) =$$

Brzozowski's Derivatives, Algorithm

Derivative of a regexp e with respect to letter c , written $\delta^c(e)$, can be computed as:

$$\delta^c(\emptyset) = \emptyset$$

$$\delta^c(\varepsilon) = \emptyset$$

$$\delta^c(d) = \begin{cases} \varepsilon & \text{if } d = c \\ \emptyset & \text{if } d \neq c \end{cases}$$

$$\delta^c(e_1 \mid e_2) = \delta^c(e_1) \mid \delta^c(e_2)$$

$$\delta^c(e_1 e_2) =$$

$$\delta^c(e_1^*) =$$

Brzozowski's Derivatives, Algorithm

Derivative of a regexp e with respect to letter c , written $\delta^c(e)$, can be computed as:

$$\delta^c(\emptyset) = \emptyset$$

$$\delta^c(\varepsilon) = \emptyset$$

$$\delta^c(d) = \begin{cases} \varepsilon & \text{if } d = c \\ \emptyset & \text{if } d \neq c \end{cases}$$

$$\delta^c(e_1 \mid e_2) = \delta^c(e_1) \mid \delta^c(e_2)$$

$$\delta^c(e_1 e_2) = \begin{cases} \delta^c(e_1) e_2 \mid \delta^c(e_2) & \text{if } \text{nullable}(e_1) \\ \delta^c(e_1) e_2 & \text{otherwise} \end{cases}$$

$$\delta^c(e_1^*) =$$

Brzozowski's Derivatives, Algorithm

Derivative of a regexp e with respect to letter c , written $\delta^c(e)$, can be computed as:

$$\delta^c(\emptyset) = \emptyset$$

$$\delta^c(\varepsilon) = \emptyset$$

$$\delta^c(d) = \begin{cases} \varepsilon & \text{if } d = c \\ \emptyset & \text{if } d \neq c \end{cases}$$

$$\delta^c(e_1 \mid e_2) = \delta^c(e_1) \mid \delta^c(e_2)$$

$$\delta^c(e_1 e_2) = \begin{cases} \delta^c(e_1) e_2 \mid \delta^c(e_2) & \text{if } \text{nullable}(e_1) \\ \delta^c(e_1) e_2 & \text{otherwise} \end{cases}$$

$$\delta^c(e_1^*) = \delta^c(e_1) e_1^*$$

Matching Regular Expressions by Derivation

Idea: given a **regexp** e and a **word** w , to know whether e **accepts** w ,
simply ***derive** e **for all letters of** w , then check whether the **result is nullable!***

$$\text{accepts} : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$$

$$\text{accepts}(e, \varepsilon) = \text{nullable}(e)$$

$$\text{accepts}(e, cw) = \text{accepts}(\delta^c(e), w)$$

Matching Regular Expressions by Derivation

Idea: given a **regexp** e and a **word** w , to know whether e **accepts** w ,
simply ***derive** e **for all letters of** w , then check whether the **result is nullable!***

$$\text{accepts} : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$$

$$\text{accepts}(e, \varepsilon) = \text{nullable}(e)$$

$$\text{accepts}(e, cw) = \text{accepts}(\delta^c(e), w)$$

Matching Regular Expressions by Derivation

Idea: given a **regexp** e and a **word** w , to know whether e **accepts** w ,
simply *derive* e **for all letters of** w , then check whether the **result is nullable!**

$$\text{accepts} : (e, w) \rightarrow \{\mathbf{true}, \mathbf{false}\}$$

$$\text{accepts}(e, \varepsilon) = \text{nullable}(e)$$

$$\text{accepts}(e, cw) = \text{accepts}(\delta^c(e), w)$$

Important: need to *cache* each intermediate result to avoid duplicating work.

Regular Expressions in Scala

Regex Data Type in Scala

```
enum RegExp:
```

```
  // empty language  $\emptyset$ 
```

```
  case Failure
```

```
  // empty word  $\varepsilon$ 
```

```
  case EmptyStr
```

```
  // character  $a$  such that  $\text{predicate}(a)$ 
```

```
  case CharWhere(predicate: Character  $\Rightarrow$  Boolean)
```

```
  // union  $\text{left} \mid \text{right}$ 
```

```
  case Union(left: RegExp, right: RegExp)
```

```
  // concatenation  $\text{first} \mid \text{second}$ 
```

```
  case Concat(first: RegExp, second: RegExp)
```

```
  // Kleene star  $\text{underlying}^*$ 
```

```
  case Star(underlying: RegExp)
```

Regex Methods

```
enum RegExp:
```

```
  case ...
```

```
  // is this regexp nullable?
```

```
  def acceptsEmpty: Boolean = ...
```

```
  // can this regexp possibly accept some words?
```

```
  def isProductive: Boolean = this match
```

```
    case Failure ⇒ false
```

```
    case EmptyStr | CharWhere(_) | Star(_) ⇒ true // approx.
```

```
    case Union(l, r) ⇒ l.isProductive || r.isProductive
```

```
    case Concat(l, r) ⇒ l.isProductive && r.isProductive
```

Regex Combinators 1 (extension methods)

```
extension (expr: RegExpr):  
  def ~ (that: RegExpr): RegExpr = Concat(expr, that)  
  def | (that: RegExpr): RegExpr = Union(expr, that)  
  def * : RegExpr = Star(expr)  
  def ? : RegExpr = expr | EmptyStr  
  def + : RegExpr = expr ~ expr.*  
  def times(n: Int): RegExpr =  
    if (n <= 0) EmptyStr else expr ~ expr.times(n - 1) }
```

Examples:

```
e1 ~ e2 ~ e3 | e4  
e1.* | e2.+
```

Regex Combinators 2 (top-level definitions)

```
def elem(pred: Char  $\Rightarrow$  Boolean): RegExpr = CharWhere(pred)
```

```
def elem(char: Char): RegExpr = CharWhere(_ == char)
```

```
def elem(chars: Iterable[Char]): RegExpr =  
  chars.map(elem).foldLeft[RegExpr](Failure)(_ | _)
```

```
def word(chars: Iterable[Char]): RegExpr =  
  chars.map(elem).foldLeft[RegExpr](EmptyStr)(_ ~ _)
```

```
def inRange(low: Char, high: Char): RegExpr =  
  elem(c  $\Rightarrow$  c  $\geq$  low && c  $\leq$  high)
```

// Example:

```
elem(_.isLetter) ~ (elem(_.isLetter) | elem(_.isDigit)).*
```

Regex Derivation in Scala

```
def derive(char: Character): RegExp =  
  def work(expr: RegExp): RegExp = expr match  
    case Failure | EmptyStr  $\Rightarrow$  Failure  
    case CharWhere(pred)  $\Rightarrow$   
      if pred(char) then EmptyStr else Failure  
    case Union(left, right)  $\Rightarrow$  work(left) | work(right)  
    case Concat(left, right)  $\Rightarrow$   
      val w = work(left) ~ right  
      if left.acceptsEmpty then w | work(right) else w  
    case Star(inner)  $\Rightarrow$  work(inner) ~ expr  
  work(this)
```

Efficiently Recognizing Regular Languages

Problem of Naive Derivation: Inefficiency

Pathological cases easily arise.

Problem of Naive Derivation: Inefficiency

Pathological cases easily arise.

Example: consider a^*

Problem of Naive Derivation: Inefficiency

Pathological cases easily arise.

Example: consider a^*

$$\delta^a(a^*) = \varepsilon a^*$$

Problem of Naive Derivation: Inefficiency

Pathological cases easily arise.

Example: consider a^*

$$\delta^a(a^*) = \varepsilon a^*$$

$$\delta^{aa}(a^*) = \emptyset a^* \mid \varepsilon a^*$$

Problem of Naive Derivation: Inefficiency

Pathological cases easily arise.

Example: consider a^*

$$\delta^a(a^*) = \varepsilon a^*$$

$$\delta^{aa}(a^*) = \emptyset a^* \mid \varepsilon a^*$$

$$\delta^{aaa}(a^*) = \emptyset a^* \mid \emptyset a^* \mid \varepsilon a^*$$

...

Every step needs to go through the list of $\emptyset a^*$, which grows without limit:

\Rightarrow *unbounded number of “states” and wrong complexity*

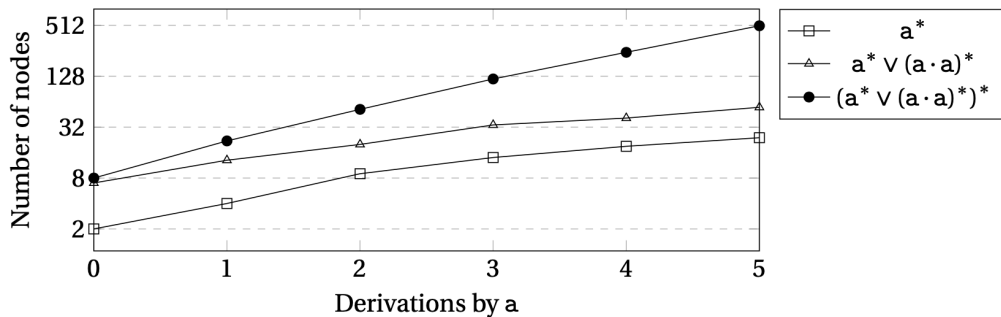
This is a simple case;

successive derivations are not always possible to “compact” on the fly...

– consider $(a^* \mid (aa)^*)^*$

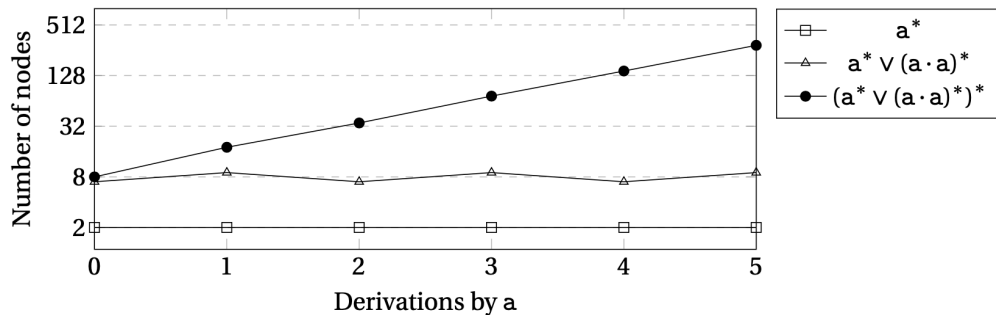
Problem of Naive Derivation: Inefficiency

Pathological cases easily arise.



Problem of Naive Derivation: Inefficiency

Not always possible to “compact” on the fly. Here it is with auto-compaction:



Idea: On-the-fly Normalization

To avoid getting into pathological cases,
it is sufficient to *normalize* the part of the expression we derive

Idea: On-the-fly Normalization

To avoid getting into pathological cases,

it is sufficient to *normalize* the part of the expression we derive

Normalization approach:

- ▶ associate all concatenations to the right: $(e_1 e_2) e_3 \Rightarrow e_1 (e_2 e_3)$
- ▶ avoid repetitions in unions: $(e_1 \mid e_2) \mid (e_3 \mid e_1) \Rightarrow e_1 \mid e_2 \mid e_3$

Normalizing Regexp Derivation in Scala

```
def deriveNorm(char: Character): RegExp =  
  val disjuncted = collection.mutable.SortedSet[RegExp]()  
  def work(expr: RegExp, rest: RegExp): Unit = expr match  
    case CharWhere(pred) => if pred(char) then disjuncted += rest  
    case Union(left, right) => work(left, rest); work(right, rest)  
    case Concat(left, right) =>  
      work(left, right ~ rest)  
      if left.acceptsEmpty then work(right, rest)  
    case Star(inner) => work(inner, expr ~ rest)  
    case Failure | EmptyStr => ()  
  work(this, EmptyStr) // register unions into `disjuncted`  
  disjuncted.foldLeft[RegExp](Failure)(_ | _) // rebuild regexp
```

The Pumping Lemma

What is the Essence of Regular Expressions/Languages?

Although these languages are usually infinite,
they contain, at their core, *simple repeating patterns*.

Pumping Lemma

Recall: $w^n = \underbrace{w w \dots w}_{n \text{ times}}$

Pumping Lemma

Recall: $w^n = \underbrace{w w \dots w}_{n \text{ times}}$

Lemma (Pumping)

For **all regexp** e , there **exists a constant** $p \geq 1$ such that
for all $w \in L(e)$ of length at least $|w| \geq p$,
we have $w = w_0 w_1 w_2$ such that:

Pumping Lemma

Recall: $w^n = \underbrace{w w \dots w}_{n \text{ times}}$

Lemma (Pumping)

For **all regexp** e , there **exists a constant** $p \geq 1$ such that
for all $w \in L(e)$ of length at least $|w| \geq p$,
we have $w = w_0 w_1 w_2$ such that:

1. $w_1 \neq \varepsilon$
2. $|w_0 w_1| \leq p$
3. $\forall n. w_0 (w_1)^n w_2 \in L(e)$

Pumping Lemma

Recall: $w^n = \underbrace{w w \dots w}_{n \text{ times}}$

Lemma (Pumping)

For **all regexp** e , there **exists a constant** $p \geq 1$ such that
for all $w \in L(e)$ of length at least $|w| \geq p$,
we have $w = w_0 w_1 w_2$ such that:

1. $w_1 \neq \varepsilon$
2. $|w_0 w_1| \leq p$
3. $\forall n. w_0 (w_1)^n w_2 \in L(e)$

Proof?

Proving the Pumping Lemma

Highly non-trivial.

By structural induction on e . Concatenation and Kleene star require careful case analysis on all possible length distributions.

Proving the Pumping Lemma

Highly non-trivial.

By structural induction on e . Concatenation and Kleene star require careful case analysis on all possible length distributions.

See §2.7 of Edelmann, Romain. Efficient Parsing with Derivatives and Zippers. PhD Thesis (EPFL). 2021.