

Regular Expression Parsing, Greedily and Stingily

Bjørn Bugge Grathwohl Fritz Henglein
Lasse Nielsen Ulrik Terp Rasmussen

Department of Computer Science,
University of Copenhagen

WG 2.8 Meeting, Aussois
October 18, 2013

Regular Expressions

- ▶ Regular Expressions (RE):

$$E ::= 0 \mid 1 \mid a \mid E_1 \times E_2 \mid E_1 + E_2 \mid E_1^* \quad (a \in \Sigma)$$

- ▶ \times binds tighter than $+$.
- ▶ Assume **non-problematic** REs: No REs containing sub-REs of the form E^* where E nullable.
 - ▶ All results extend to problematic REs, but are more complicated to state and prove.

What is Regular Expression “Matching”?

Given $s \in \Sigma^*$.

1. **Acceptance testing**: Is $s \in \mathcal{L}[E]$?
 - ▶ **String searching**: Find some substring s' of s such that $s' \in \mathcal{L}[E]$. (Variation: Find *all* substrings.)
2. **Pattern matching**: Given $s \in \Sigma^*$, find substrings of s such that each matches a *sub-RE* in E . (Variation: Return multiple matches for each sub-RE.)
3. **Parsing**: Return complete parse tree of s under E , if it exists

Note:

- ▶ Increasing information content.
- ▶ Classical automata theory (NFA \rightarrow DFA, DFA minimization, etc.) applies only to acceptance testing.
- ▶ Pattern matching returns only one element match under $*$.

Example

RE = $((a + b) \times (c + d))^*$. Input string = $acbd$.

1. Acceptance testing: Yes!
2. Pattern matching: $(0, 4), (2, 4), (2, 3), (3, 4)$
3. Parsing: $[(inl\ a, inl\ c), (inr\ b, inr\ d)]$

Regular Expressions as Types

- ▶ **Type interpretation** $\mathcal{T}[E]$:

$$\mathcal{T}[0] = \emptyset$$

$$\mathcal{T}[1] = \{()\}$$

$$\mathcal{T}[a] = \{a\}$$

$$\mathcal{T}[E_1 \times E_2] = \{(V_1, V_2) \mid V_1 \in \mathcal{T}[E_1], V_2 \in \mathcal{T}[E_2]\}$$

$$\mathcal{T}[E_1 + E_2] = \{\text{inl } V_1 \mid V_1 \in \mathcal{T}[E_1]\} \\ \cup \{\text{inr } V_2 \mid V_2 \in \mathcal{T}[E_2]\}$$

$$\mathcal{T}[E^*] = \{(V_1, \dots, V_n) \mid n \geq 0 \wedge \\ \forall 1 \leq i \leq n. V_i \in \mathcal{T}[E]\}$$

- ▶ Parse tree = value

Unparsing (“Flattening”)

- ▶ Flattening yields underlying string:

$$\begin{aligned}\text{flat}(\epsilon) &= \epsilon \\ \text{flat}(a) &= a \\ \text{flat}((V_1, V_2)) &= \text{flat}(V_1) \text{flat}(V_2) \\ \text{flat}(\text{inl } V_1) &= \text{flat}(V_1) \\ \text{flat}(\text{inr } V_2) &= \text{flat}(V_2) \\ \text{flat}([V_1, \dots, V_n]) &= \text{flat}(V_1) \cdots \text{flat}(V_n)\end{aligned}$$

- ▶ The parse trees for a given string s :

$$\mathcal{T}_s[[E]] = \{V \in \mathcal{T}[[E]] \mid \text{flat}(V) = s\}.$$

Proposition

$$\mathcal{L}[[E]] = \{\text{flat}(V) \mid V \in \mathcal{T}[[E]]\}.$$

Challenges

- ▶ Grammatical ambiguity: Which parse tree to return?
- ▶ How to represent parse trees compactly?
- ▶ Time: Straightforward backtracking algorithm, but impractical: $\Theta(m2^n)$ time, where $m = |E|$, $n = |s|$.
- ▶ Space: How to minimize RAM consumption?

Disambiguation

- ▶ RE E **ambiguous** iff $|\mathcal{T}_s[E]| > 1$ for some s .
- ▶ How to **deterministically** choose one $V \in \mathcal{T}_s[E]$ among several possible candidates?
- ▶ **Greedy** matching: Intuitively, choose what a backtracking parser returns:
 1. Try left alternative first,
 2. If it fails, backtrack and try the right alternative.
 3. Treat E^* as $E \times E^* + 1$.

Greedy Order \prec_V

$\text{inl } V$	\prec_V	$\text{inr } V$	
$[V_1, \dots]$	\prec_V	\square	
(V_1, V_2)	\prec_V	(V'_1, V'_2)	if $V_1 \prec_V V_2 \vee$ $(V_1 = V'_1 \wedge V_2 \prec_V V'_2)$
$\text{inl } V$	\prec_V	$\text{inl } V'$	if $V \prec_V V'$
$\text{inr } V$	\prec_V	$\text{inr } V'$	if $V \prec_V V'$
$[V_1, \dots]$	\prec_V	$[V'_1, \dots]$	if $V_1 \prec_V V'_1$
$[V_1, V_2, \dots]$	\prec_V	$[V_1, V'_2, \dots]$	if $[V_2, \dots] \prec_V [V'_2, \dots]$

Proposition (Frisch/Cardelli)

For any RE E , string s , \prec_V is a strict well-founded total order on $\mathcal{T}_s[E]$.

Definition

Greedy parse for $s \in \mathcal{L}[E]$: $\min_{\prec_V} \mathcal{T}_s[E]$.

Bit-Coding

- ▶ Compact representation of parse trees where the RE is known.
- ▶ Encoding $\lceil \cdot \rceil : \mathcal{V} \rightarrow \{1, 0\}^*$,

$$\lceil () \rceil = \epsilon$$

$$\lceil a \rceil = \epsilon$$

$$\lceil (V_1, V_2) \rceil = \lceil V_1 \rceil \lceil V_2 \rceil$$

$$\lceil \text{inl } (V_1) \rceil = 0 \lceil V_1 \rceil$$

$$\lceil \text{inr } (V_2) \rceil = 1 \lceil V_2 \rceil$$

$$\lceil [V_1, \dots, V_n] \rceil = 0 \lceil V_1 \rceil \dots 0 \lceil V_n \rceil 1$$

- ▶ **Type-indexed** decoding $\lfloor \cdot \rfloor_E : \{1, 0\}^* \rightarrow \mathcal{T}[E]$: Interpret RE as nondeterministic algorithm to construct parse tree, with bit-code as oracle. (“Every bit counts.”)
- ▶ $\mathcal{B}[E] = \{\lceil V \rceil \mid V \in \mathcal{T}[E]\}$
 - ▶ $\mathcal{B}_s[E] = \{\lceil V \rceil \mid V \in \mathcal{T}_s[E]\}$.

Example

RE = $((a + b) \times (c + d))^*$. Input string = $acbd$.



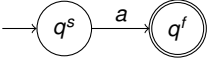
1. Acceptance testing: Yes!
2. Pattern matching: $(0, 4), (2, 4), (2, 3), (3, 4)$
3. Parsing: $[(inl\ a, inl\ c), (inr\ b, inr\ d)]$
 - ▶ Bit-code: 0000111.

Augmented NFAs

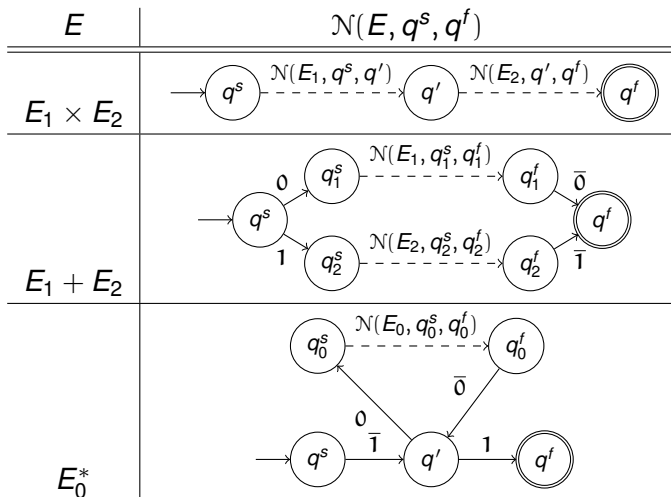
- ▶ Augmented NFA (aNFA) is a 5-tuple $M \in (Q, \Sigma, \Delta, q^s, q^f)$.
- ▶ States Q ; $q^s, q^f \in Q$ start and finishing states.
- ▶ Input alphabet Σ .
- ▶ Labeled transition relation $\Delta \subseteq Q \times (\Sigma \cup \{1, 0\} \cup \{\bar{1}, \bar{0}\}) \times Q$.
 - ▶ Σ input labels; $\{1, 0\}$ output labels; $\{\bar{1}, \bar{0}\}$ log labels.
- ▶ Write $q \overset{p}{\rightsquigarrow} q'$ if there is a walk from q to q' ; p sequence of labels.
 - ▶ $\text{in}(p)$ = input label subsequence;
 - ▶ $\text{out}(p)$ = output labels;
 - ▶ $\text{log}(p)$ = log labels.

aNFA Construction (1/2)

- Define $\mathcal{N}(E, q^s, q^f)$ as set of aNFAs for E , with start and finishing states q^s, q^f :

E	$\mathcal{N}(E, q^s, q^f)$
0	
1	 (implies $q^s = q^f$)
a	

aNFA Construction (2/2)



Representation Theorem

Theorem

Let $M = \mathcal{N}(E, q^s, q^f)$. The paths of M are in one-to-one correspondence with the parse trees of E :

$$\mathcal{B}_s[E] = \{\text{out}(p) \mid q^s \overset{p}{\rightsquigarrow} q^f, \text{in}(p) = s\}$$

Greedy parse = Lexicographically least bitcode

Proposition

For all $E, V, V' \in \mathcal{T}[E]$:

$$V \prec_{\mathcal{V}} V' \iff \lceil V \rceil \prec_{\mathcal{B}} \lceil V' \rceil$$

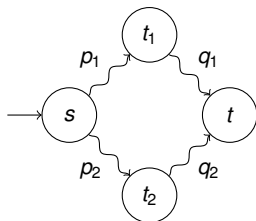
where $\prec_{\mathcal{B}}$ is lexicographic ordering on $\{0, 1\}^*$.

Corollary

Let $M = \mathcal{N}(E, q^s, q^f)$. For all $s \in \mathcal{L}[E]$:

$$\min_{\prec_{\mathcal{V}}} \mathcal{T}_s[E] = \lfloor \min_{\prec_{\mathcal{B}}} \{\text{out}(p) \mid q^s \xrightarrow{p} q^f, \text{in}(p) = s\} \rfloor_E.$$

Proposition



If p_1 not prefix of p_2 , then

$$\text{out}(p_1) \prec_{\mathcal{B}} \text{out}(p_2) \Rightarrow \text{out}(p_1 q_1) \prec_{\mathcal{B}} \text{out}(p_2 q_2)$$

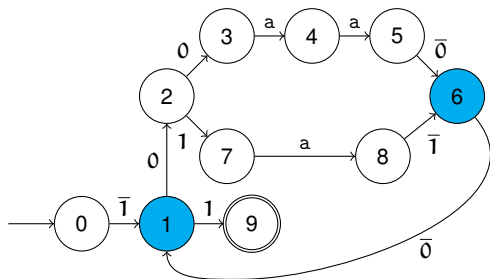
Lean-log algorithm

- ▶ Simulate aNFA for input s , using **ordered** state sets.
 - ▶ Each state represents lexicographically least path from initial state to it.
 - ▶ States are ordered according to the lexicographic ordering on the paths they represent.
- ▶ Perform state-ordered ϵ -closure: Log 1 bit per **join state** for each input character.
- ▶ Use reverse aNFA and log bits to construct bit-code.
- ▶ (Construct parse tree from bit-code, if desired.)

Example: Parse aaa with RE $(aa + a)^*$

► Input: aaa

Log	ϵ	a	a	a
1:				
6:				



Example: Parse aaa with RE $(aa + a)^*$

► Input: aaa

ε

3

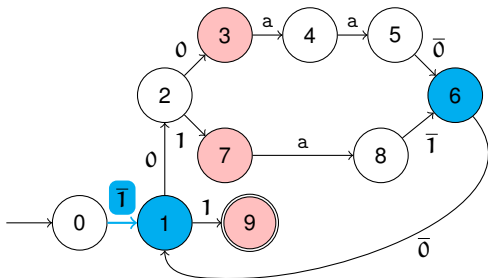
7

9

Log ε a a a

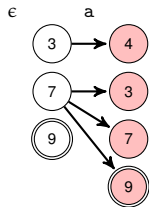
1: 1

6: -

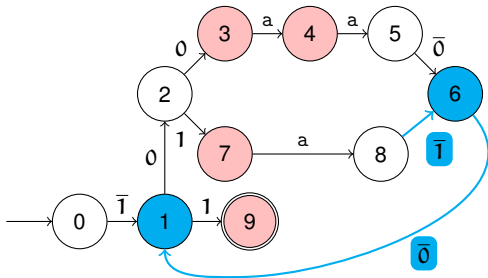


Example: Parse aaa with RE $(aa + a)^*$

► Input: a aa

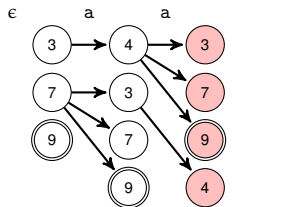


Log	ϵ	a	a	a
1:	$\bar{1}$	$\bar{0}$		
6:	-	$\bar{1}$		

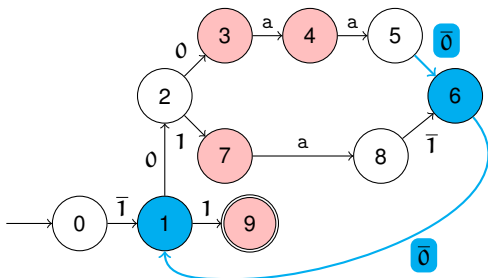


Example: Parse aaa with RE $(aa + a)^*$

► Input: aa a

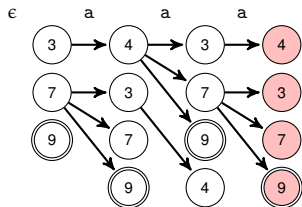


Log	ϵ	a	a	a
1:	$\bar{1}$	$\bar{0}$	$\bar{0}$	
6:	-	$\bar{1}$	$\bar{0}$	

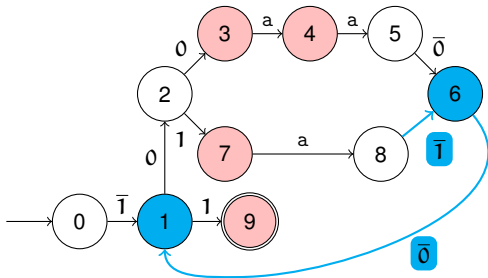


Example: Parse aaa with RE $(aa + a)^*$

► Input: **aaa**



Log	ϵ	a	a	a
1 :	$\bar{1}$	$\bar{0}$	$\bar{0}$	$\bar{0}$
6 :	-	$\bar{1}$	$\bar{0}$	$\bar{1}$



Key properties of lean-log algorithm

- ▶ Semi-streaming: Forward streaming pass over input, logging join-state bits; backward pass for constructing bit-code.
 - ▶ Two passes required because of disambiguation requiring unbounded look-ahead.
- ▶ Input string read in streaming fashion, using $O(m)$ working memory and kn bits of LIFO memory for the log.
- ▶ Input string need not be stored. (Consider input coming from a generator.)
- ▶ Runs in time $O(mn)$.

Implementation

- ▶ Implementations of lean-log algorithm
 - ▶ Straightforward Haskell version
 - ▶ Optimized Haskell version, based on Conduit (10 times faster and)
 - ▶ Straightforward C version (10 times faster than fast Haskell version)
- ▶ No NFA-minimization, no DFA generation, no word-level parallelism, no special RE-processing, no special handling of bounded iteration.

Performance

- ▶ Better performance than Play
- ▶ Competitive with RE2 when RE2 does not employ static optimizations, or when subjected to REs that are not “tuned” to Perl (made deterministic)
- ▶ Otherwise competitive with Grep and other tools, but not with RE2.
- ▶ Note: These tools perform only acceptance testing or RE pattern matching, not full parsing; and they don't always do it correctly.
- ▶ Best amongst all tested full RE parsers (both greedy and other).

See paper at CIAA 2013.

Questions?

Questions?

What does this have to do with FP?

- ▶ REs are a **declarative DSL**
 - ▶ Widely used, but still underused (notably REs with nested *, ambiguous REs)
- ▶ REs as **types**
 - ▶ Already in FP languages: unit, singleton, Cartesian product, direct sum, tail-recursive types.
 - ▶ RE containment as type *coercions*: order-preserving linear functions.
 - ▶ Types capture programming intension of REs, are elegant theoretical framework (e.g. definition of ambiguity)
- ▶ Bit-coding as **efficient** unboxed data type representation
 - ▶ for strings: bit-code of string as element of $\Sigma^* \cong$ the string itself; as element of $E \subseteq \Sigma^*$, fewer bits.
 - ▶ for simple and recursive types: **unboxed** data representation, with certain type isomorphisms as identities (noop-coercions); e.g. $A \times (B \times C) = (A \times B) \times C$.
- ▶ ...