

# Regular expression matching and operational semantics

Hayo Thielecke  
joint work with Asiri Rathnayake

University of Birmingham, UK

# Introduction: regular expression matching

- Regular expressions are everywhere
- In compiling: convert to Deterministic Finite Automaton (DFA) via tools like lex
- DFA matching is efficient, linear time
- Java, Perl etc do **not** build a DFA but use backtracking
- Naive backtracking:  
given  $(e_1 \mid e_2)$ , try to match  $e_1$  and then if it fails try  $e_2$
- functional programming classic: failure continuations
- But backtracking is inefficient: **exponential**
- In Java: matching  $\underbrace{a \dots a}_{1000} b$  to  $a^*$  takes forever

# Thomson's lockstep matcher

- Thomson's lockstep matcher from 1968 is more efficient than backtracking
- Interpret the regular expression as a non-deterministic program
- Cox: "Name the most widely used bytecode interpreter or virtual machine."
- Lockstep interpretation avoids exponential blowup
- Thomson's construction can be seen as constructing an NFA (Non-deterministic Finite Automaton)
- But looking at it as merely a naive NFA misses key ideas
- Topic of this talk: use programming language technology, not so much automata theory

# Preliminaries and background

- My background is in programming language theory (not automata theory)
- This talk is an extended version of our paper in Structural Operation Semantics 2011
- Operational semantics is much like writing an interpreter or `eval` in Lisp or Javascript
- Big-**step** semantics:  $((1 + 2) + 3) \Downarrow 6$
- Small **step** semantics  $((1 + 2) + 3) \rightarrow (3 + 3) \rightarrow 6$
- $\Downarrow$  and  $\rightarrow$  are relations, not always functions: non-determinism
- Compare derivations in grammars

$$S \Rightarrow E + E \Rightarrow 1 + E \Rightarrow 1 + 2$$

# Automata $\leftrightarrow$ operational semantics

## Automata theory

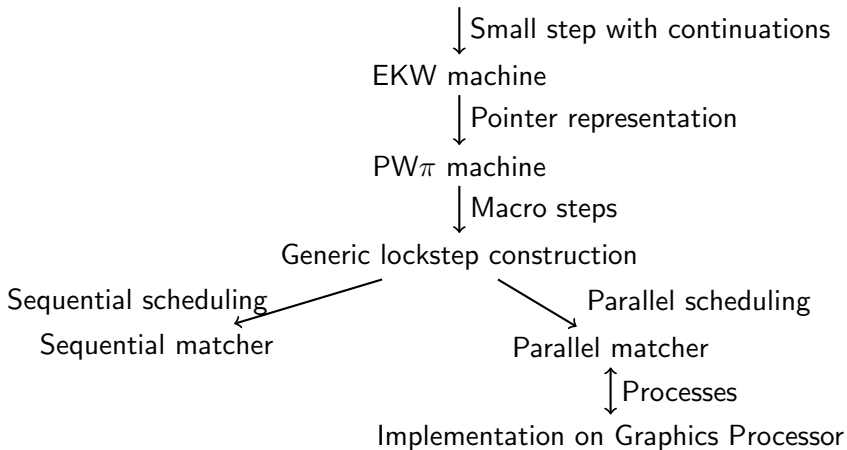
- $p \xrightarrow{a} q$ , where  $a \in \Sigma$ ,  $p, q \in Q$
- Typical questions: languages closed under complement
- $p$  and  $q$  are atomic states, elements of a *finite* set  $Q$  without further structure

## Programming language theory: operational semantics

- Structural: 
$$\frac{P \rightarrow Q}{PK \rightarrow QK}$$
 or  $\langle P, K \rangle \rightarrow \langle Q, K \rangle$
- $P$ ,  $Q$  and  $K$  are defined by induction
- $\rightarrow$  is defined by induction following their structure (SOS)

# Structure of the paper — not everything will be in the talk

Regular expression matching as big-step semantics



# Matching as big-step semantics

$$\frac{e_1 \downarrow w_1 \quad e_2 \downarrow w_2}{(e_1 e_2) \downarrow (w_1 w_2)} \text{ (SEQ)} \quad \frac{}{a \downarrow a} \text{ (MATCH)} \quad \frac{}{\varepsilon \downarrow \varepsilon} \text{ (EPSILON)}$$

$$\frac{e \downarrow w_1 \quad e^* \downarrow w_2}{e^* \downarrow (w_1 w_2)} \text{ (KLEENE1)} \quad \frac{}{e^* \downarrow \varepsilon} \text{ (KLEENE2)}$$

$$\frac{e_1 \downarrow w}{(e_1 \mid e_2) \downarrow w} \text{ (ALT1)} \quad \frac{e_2 \downarrow w}{(e_1 \mid e_2) \downarrow w} \text{ (ALT2)}$$

Like denotational semantics — but how do you match efficiently?

# The EKW machine

$$\langle e ; k ; w \rangle \rightarrow \langle e' ; k' ; w' \rangle$$

$$\langle e_1 \mid e_2 ; k ; w \rangle \rightarrow \langle e_1 ; k ; w \rangle$$

$$\langle e_1 \mid e_2 ; k ; w \rangle \rightarrow \langle e_2 ; k ; w \rangle$$

$$\langle e_1 e_2 ; k ; w \rangle \rightarrow \langle e_1 ; e_2 :: k ; w \rangle$$

$$\langle e^* ; k ; w \rangle \rightarrow \langle e ; e^* :: k ; w \rangle$$

$$\langle e^* ; k ; w \rangle \rightarrow \langle \varepsilon ; k ; w \rangle$$

$$\langle a ; k ; a w \rangle \rightarrow \langle \varepsilon ; k ; w \rangle$$

$$\langle \varepsilon ; e :: k ; w \rangle \rightarrow \langle e ; k ; w \rangle$$



# Regular expressions as trees/graphs in memory

## Expression as trees

We need to distinguish the two  $a$  positions in

$$(a b) | (a c)$$

Represent syntax tree as heap with pointers, using ideas from Separation Logic.

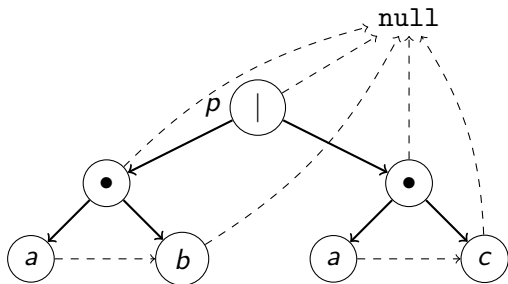
## Continuation pointers

The continuation ( $k$  in the EKW machine) is determined by the position in the tree.

Hardwire the continuation in the tree as a pointer.

# Syntax tree in memory: $(a b) | (a c)$

- child nodes in syntax tree
- > continuation pointers



Interleaved communication for evolving pointers  $p$  (CCS)

$$\boxed{M \longrightarrow M'}$$

$$\frac{M_1 \longrightarrow M_2}{M_1 \parallel M_3 \longrightarrow M_2 \parallel M_3} \quad \frac{}{((p.M) \parallel \bar{p}) \longrightarrow M}$$

Synchronous step for lockstep matching of symbol  $a$  (SCCS)

$$\boxed{M \xrightarrow{a} M'}$$

$$\frac{M' \not\equiv (\$a.M') \parallel M''' \quad M' \not\rightarrow}{(\$a.M_1 \parallel \dots \parallel \$a.M_n \parallel M') \xrightarrow{a} (M_1 \parallel \dots \parallel M_n)}$$

# Process translation

Each node  $p$  in the syntax tree become a process  $\llbracket p \rrbracket \pi$  for recognizing that expression.

$\llbracket p \rrbracket \pi = p.(\overline{q_1} \parallel \overline{q_2})$	if $\pi(p) = (q_1 \mid q_2)$
$\llbracket p \rrbracket \pi = p.\overline{q_1}$	if $\pi(p) = (q_1 \bullet q_2)$
$\llbracket p \rrbracket \pi = p.(\overline{q_1} \parallel \overline{q_2})$	if $\pi(p) = q_1^*$ and $\text{cont } p = q_2$
$\llbracket p \rrbracket \pi = p.\overline{q}$	if $\pi(p) = \varepsilon$ and $\text{cont } p = q$
$\llbracket p \rrbracket \pi = p.\$a.\overline{q}$	if $\pi(p) = a$ and $\text{cont } p = q$

Example:  $(a\ b) \mid (a\ c)$  with input  $a\ b$

$p.(\overline{p_1} \parallel \overline{p_2})$

$p_1.\$a.\overline{p_3}$

$p_2.\$a.\overline{p_4}$

$p_3.\$b.\overline{p_5}$

$p_4.\$c.\overline{p_5}$

Remaining input:  $a\ b$

$p$  is sent

Example:  $(a b) \mid (a c)$  with input  $a b$

$\overline{p_1} \parallel \overline{p_2}$

$p_1.\$a.\overline{p_3}$

$p_2.\$a.\overline{p_4}$

$p_3.\$b.\overline{p_5}$

$p_4.\$c.\overline{p_5}$

Remaining input:  $a b$

$p_2$  is sent

Example:  $(a b) \mid (a c)$  with input  $a b$

$\overline{p_1}$

$p_1.\$a.\overline{p_3}$

$\$a.\overline{p_4}$

$p_3.\$b.\overline{p_5}$

$p_4.\$c.\overline{p_5}$

Remaining input:  $a b$

$p_1$  is sent

Example:  $(a b) \mid (a c)$  with input  $a b$

$\$a.\overline{p_3}$

$\$a.\overline{p_4}$

$p_3.\$b.\overline{p_5}$

$p_4.\$c.\overline{p_5}$

Remaining input:  $a b$

$a$  is matched synchronously



Example:  $(a b) \mid (a c)$  with input  $a b$

$\overline{p_3}$

$\overline{p_4}$

$p_3.\$b.\overline{p_5}$

$p_4.\$c.\overline{p_5}$

Remaining input:  $b$

$p_3$  is sent

Example:  $(a b) \mid (a c)$  with input  $a b$

$\overline{p_4}$

$\$b.\overline{p_5}$

$p_4.\$c.\overline{p_5}$

Remaining input:  $b$

$p_4$  is sent

Example:  $(a b) \mid (a c)$  with input  $a b$

$\$b.\overline{p_5}$

$\$c.\overline{p_5}$

Remaining input:  $b$

$b$  is matched synchronously

Example:  $(a\ b) \mid (a\ c)$  with input  $a\ b$

$\overline{p_5}$

Remaining input: empty

Successful match if  $p_5$  is the final continuation

# Main correctness result

- Theorem: our concurrent matcher is correct.
- Like computing the  $\varepsilon$  closure in automata theory
- But concurrently.
- Proof relies on invariant using

$$\Box P = \{q \mid \exists p \in P. p \longrightarrow \dots \longrightarrow q \wedge q \not\rightarrow\}$$

- $\Box P$  not quite a closure operator, as some elements are removed.
- $\Box P$  is a kind of tree modality

# Summary

- We have formalized regular expression matching using abstract machines
- Thompson's technique from 1968 - more efficient than regex matchers in Perl or Java
- View expression as program and run it on an interpreter.
  - continuations:  $(a b) \mid (a c)$  leads to stack holding  $b$  or  $c$
  - pointers via separation logic: two copies of  $a$  at different addresses, intentionally  $\neq$
  - processes: alternation  $\mid$  becomes parallel composition, lockstep as in SCCS
- Bridge-building between programming language theory and automata theory

## Further work 1: **extended** reg exp matching

- Machines can be extended to non-regular constructs used in practice
- Examples: submatching, back references: no longer regular
- $a(b)\backslash 1$  and  $(ab)\backslash 1$  are not equivalent due to the backreference
- We try to adapt Thompson's technique to such constructs
- Compare: abstract machines for **impure** languages (e.g., SECD machine + assignment)
- Analogy: DFA  $\cong \lambda$ , non-regular  $\cong$  effects
- Some current and future work, perhaps Deep Packet Inspection for security

## Further work 2: GPGPU programming and op sem

- GPGPU = General Purpose Graphics Processing Unit
- A leading example of non-numeric GPGPU is state machines
- We made a toy reg exp implementation in CUDA based on message passing
- But existing NFA matchers like iNFAnt are faster: optimized,  $\epsilon$  transitions eliminated, data structure well suited to GPU memory
- We are working on abstract machine and operational semantics techniques for multi-core and GPGPU
- GPUs make both concurrent and interleaved transitions