# Automata and Formal Languages (5)
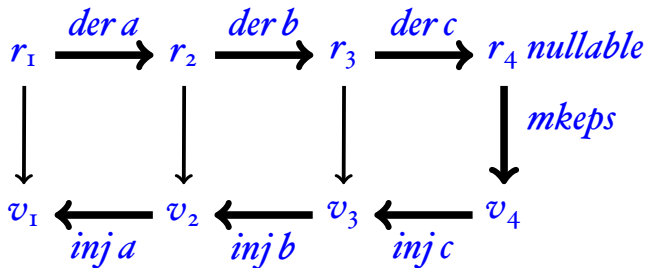
Email:    christian.urban at kcl.ac.uk
Office:   S1.27 (1st floor Strand Building)
Slides:   KEATS (also home work is there)
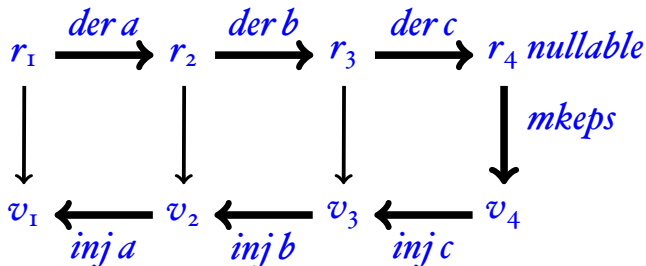
# Last Week
# Regexes and Values

Regular expressions and their corresponding values:

$$
\begin{array}{llll}
r & ::= & \varnothing & \\
  &  | & \epsilon & \\
  &  | & c & \\
  &  | & r_1 \cdot r_2 & \\
  &  | & r_1 + r_2 & \\
  &  | & r^* &
\end{array}
\qquad
\begin{array}{lll}
v & ::= & \\
  &  & \textit{Empty} \\
  & | & \textit{Char}(c) \\
  & | & \textit{Seq}(v_1, v_2) \\
  & | & \textit{Left}(v) \\
  & | & \textit{Right}(v) \\
  & | & [v_1, \ldots v_n]
\end{array}
$$

$r_1$:  $a \cdot (b \cdot c)$
$r_2$:  $\epsilon \cdot (b \cdot c)$
$r_3$:  $(\varnothing \cdot (b \cdot c)) + (\epsilon \cdot c)$
$r_4$:  $(\varnothing \cdot (b \cdot c)) + ((\varnothing \cdot c) + \epsilon)$

$$r_1 \xrightarrow{der\ a} r_2 \xrightarrow{der\ b} r_3 \xrightarrow{der\ c} r_4\ nullable$$

*mkeps*

$$v_1 \xleftarrow{inj\ a} v_2 \xleftarrow{inj\ b} v_3 \xleftarrow{inj\ c} v_4$$

$v_1$:  $Seq(Char(a), Seq(Char(b), Char(c)))$
$v_2$:  $Seq(Empty, Seq(Char(b), Char(c)))$
$v_3$:  $Right(Seq(Empty, Char(c)))$
$v_4$:  $Right(Right(Empty))$

# Mkeps

Finding a (posix) value for recognising the empty string:

$$mkeps\ \epsilon \quad \overset{\text{def}}{=} \quad Empty$$

$$mkeps\ r_1 + r_2 \quad \overset{\text{def}}{=} \quad \text{if } nullable(r_1)$$
$$\text{then } Left(mkeps(r_1))$$
$$\text{else } Right(mkeps(r_2))$$

$$mkeps\ r_1 \cdot r_2 \quad \overset{\text{def}}{=} \quad Seq(mkeps(r_1), mkeps(r_2))$$

$$mkeps\ r^* \quad \overset{\text{def}}{=} \quad []$$

# Inject

Injecting ("Adding") a character to a value

$$inj\,(c)\,c\,Empty \stackrel{\text{def}}{=} Char\,c$$

$$inj\,(r_1 + r_2)\,c\,Left(v) \stackrel{\text{def}}{=} Left(inj\,r_1\,c\,v)$$

$$inj\,(r_1 + r_2)\,c\,Right(v) \stackrel{\text{def}}{=} Right(inj\,r_2\,c\,v)$$

$$inj\,(r_1 \cdot r_2)\,c\,Seq(v_1, v_2) \stackrel{\text{def}}{=} Seq(inj\,r_1\,c\,v_1, v_2)$$

$$inj\,(r_1 \cdot r_2)\,c\,Left(Seq(v_1, v_2)) \stackrel{\text{def}}{=} Seq(inj\,r_1\,c\,v_1, v_2)$$

$$inj\,(r_1 \cdot r_2)\,c\,Right(v) \stackrel{\text{def}}{=} Seq(mkeps(r_1), inj\,r_2\,c\,v)$$

$$inj\,(r^*)\,c\,Seq(v, vs) \stackrel{\text{def}}{=} inj\,r\,c\,v \,::\, vs$$

*inj*: 1st arg $\mapsto$ a rexp; 2nd arg $\mapsto$ a character; 3rd arg $\mapsto$ a value

# Flatten

Obtaining the string underlying a value:

$$
\begin{aligned}
|Empty| &\overset{\text{def}}{=} [] \\
|Char(c)| &\overset{\text{def}}{=} [c] \\
|Left(v)| &\overset{\text{def}}{=} |v| \\
|Right(v)| &\overset{\text{def}}{=} |v| \\
|Seq(v_1, v_2)| &\overset{\text{def}}{=} |v_1| @ |v_2| \\
|[v_1, \ldots, v_n]| &\overset{\text{def}}{=} |v_1| @ \ldots @ |v_n|
\end{aligned}
$$

$r_1$:   $a \cdot (b \cdot c)$
$r_2$:   $\epsilon \cdot (b \cdot c)$
$r_3$:   $(\varnothing \cdot (b \cdot c)) + (\epsilon \cdot c)$
$r_4$:   $(\varnothing \cdot (b \cdot c)) + ((\varnothing \cdot c) + \epsilon)$

$r_1 \xrightarrow{der\,a} r_2 \xrightarrow{der\,b} r_3 \xrightarrow{der\,c} r_4$ *nullable*

*mkeps*

$v_1 \xleftarrow{inj\,a} v_2 \xleftarrow{inj\,b} v_3 \xleftarrow{inj\,c} v_4$

$v_1$:   $Seq(Char(a), Seq(Char(b), Char(c)))$
$v_2$:   $Seq(Empty, Seq(Char(b), Char(c)))$
$v_3$:   $Right(Seq(Empty, Char(c)))$
$v_4$:   $Right(Right(Empty))$

$|v_1|$:   $abc$
$|v_2|$:   $bc$
$|v_3|$:   $c$
$|v_4|$:   $[]$

# Simplification

- If we simplify after the derivative, then we are builing the value for the simplified regular expression, but *not* for the original regular expression.



$$(\varnothing \cdot (b \cdot c)) + ((\varnothing \cdot c) + \epsilon) \mapsto \epsilon$$

# Rectification

$simp(r)$:
   case $r = r_1 + r_2$
     let $(r_{1s}, f_{1s}) = simp(r_1)$
         $(r_{2s}, f_{2s}) = simp(r_2)$
    case $r_{1s} = \varnothing$: return $(r_{2s}, \lambda v.\, Right(f_{2s}(v)))$
    case $r_{2s} = \varnothing$: return $(r_{1s}, \lambda v.\, Left(f_{1s}(v)))$
    case $r_{1s} = r_{2s}$: return $(r_{1s}, \lambda v.\, Left(f_{1s}(v)))$
    otherwise: return $(r_{1s} + r_{2s}, f_{alt}(f_{1s}, f_{2s}))$

$f_{alt}(f_1, f_2) \stackrel{\text{def}}{=}$
    $\lambda v.$ case $v = Left(v')$:   return $Left(f_1(v'))$
         case $v = Right(v')$: return $Right(f_2(v'))$

# Rectification

$simp(r)$:...

    case $r = r_1 \cdot r_2$

        let $(r_{1s}, f_{1s}) = simp(r_1)$

             $(r_{2s}, f_{2s}) = simp(r_2)$

      case $r_{1s} = \varnothing$: return $(\varnothing, f_{error})$

      case $r_{2s} = \varnothing$: return $(\varnothing, f_{error})$

      case $r_{1s} = \epsilon$: return $(r_{2s}, \lambda v.\, Seq(f_{1s}(Empty), f_{2s}(v)))$

      case $r_{2s} = \epsilon$: return $(r_{1s}, \lambda v.\, Seq(f_{1s}(v), f_{2s}(Empty)))$

      otherwise: return $(r_{1s} \cdot r_{2s}, f_{seq}(f_{1s}, f_{2s}))$
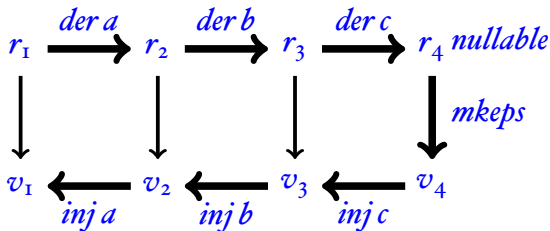
$f_{seq}(f_1, f_2) \overset{\text{def}}{=}$

    $\lambda v.$ case $v = Seq(v_1, v_2)$: return $Seq(f_1(v_1), f_2(v_2))$

# Lexing with Simplification

$$lex\ r\ [] \overset{\text{def}}{=} \text{if } nullable(r) \text{ then } mkeps(r) \text{ else } error$$

$$lex\ r\ c :: s \overset{\text{def}}{=} \text{let } (r', frect) = simp(der(c, r))$$
$$inj\ r\ c\ (frect(lex(r', s)))$$

# Records

- new regex: $(x : r)$     new value: $Rec(x, v)$

# Records

- new regex: $(x : r)$     new value: $Rec(x, v)$

- $nullable(x : r) \overset{\text{def}}{=} nullable(r)$

- $der\, c\, (x : r) \overset{\text{def}}{=} (x : der\, c\, r)$

- $mkeps(x : r) \overset{\text{def}}{=} Rec(x, mkeps(r))$

- $inj\, (x : r)\, c\, v \overset{\text{def}}{=} Rec(x, inj\, r\, c\, v)$

# Records

- new regex: $(x : r)$     new value: $Rec(x, v)$

- $nullable(x : r) \overset{\text{def}}{=} nullable(r)$

- $der\, c\, (x : r) \overset{\text{def}}{=} (x : der\, c\, r)$

- $mkeps(x : r) \overset{\text{def}}{=} Rec(x, mkeps(r))$

- $inj\, (x : r)\, c\, v \overset{\text{def}}{=} Rec(x, inj\, r\, c\, v)$

for extracting subpatterns $(z : ((x : ab) + (y : ba)))$

# While Tokens

$$\text{WHILE\_REGS} \stackrel{\text{def}}{=} ((\text{"k"} : \text{KEYWORD}) +$$
$$(\text{"i"} : \text{ID}) +$$
$$(\text{"o"} : \text{OP}) +$$
$$(\text{"n"} : \text{NUM}) +$$
$$(\text{"s"} : \text{SEMI}) +$$
$$(\text{"p"} : (\text{LPAREN} + \text{RPAREN})) +$$
$$(\text{"b"} : (\text{BEGIN} + \text{END})) +$$
$$(\text{"w"} : \text{WHITESPACE}))^*$$

"if true then then 42 else +"

```
KEYWORD(if),
WHITESPACE,
IDENT(true),
WHITESPACE,
KEYWORD(then),
WHITESPACE,
KEYWORD(then),
WHITESPACE,
NUM(42),
WHITESPACE,
KEYWORD(else),
WHITESPACE,
OP(+)
```

"if true then then 42 else +"

```
KEYWORD(if),
IDENT(true),
KEYWORD(then),
KEYWORD(then),
NUM(42),
KEYWORD(else),
OP(+)
```

There is one small problem with the tokenizer.
How should we tokenize:

$$"x - 3"$$

OP:
  "+", "-"
NUM:
  (NONZERODIGIT · DIGIT$^*$) + "0"
NUMBER:
  NUM + ("-" · NUM)

# Two Rules

- Longest match rule ("maximal munch rule"): The longest initial substring matched by any regular expression is taken as next token.

- Rule priority: For a particular longest initial substring, the first regular expression that can match determines the token.

# Environment

Obtaining the "recorded" parts of a regular expression:

$$env(Empty) \overset{\text{def}}{=} []$$
$$env(Char(c)) \overset{\text{def}}{=} []$$
$$env(Left(v)) \overset{\text{def}}{=} env(v)$$
$$env(Right(v)) \overset{\text{def}}{=} env(v)$$
$$env(Seq(v_1, v_2)) \overset{\text{def}}{=} env(v_1) @ env(v_2)$$
$$env([v_1, \ldots, v_n]) \overset{\text{def}}{=} env(v_1) @ \ldots @ env(v_n)$$
$$env(Rec(x : v)) \overset{\text{def}}{=} (x : |v|) :: env(v)$$

- Regular expression for email addresses

$$(\text{name: } [a\text{-}z0\text{-}9\_.-]^+) \cdot @ \cdot$$
$$(\text{domain: } [a\text{-}z0\text{-}9\,.-]^+) \cdot . \cdot$$
$$(\text{top\_level: } [a\text{-}z\,.]^{\{2,6\}})$$

`christian.urban@kcl.ac.uk`

- result environment:

$$[(\textit{name} : \texttt{christian.urban}),$$
$$(\textit{domain} : \texttt{kcl}),$$
$$(\textit{top\_level} : \texttt{ac.uk})]$$

# Coursework

$$nullable\left(\left[c_1 c_2 \ldots c_n\right]\right) \overset{def}{=} \;?$$

$$nullable(r^+) \overset{def}{=} \;?$$

$$nullable(r^?) \overset{def}{=} \;?$$

$$nullable(r^{\{n,m\}}) \overset{def}{=} \;?$$

$$nullable(\sim r) \overset{def}{=} \;?$$

$$der\, c\,\left(\left[c_1 c_2 \ldots c_n\right]\right) \overset{def}{=} \;?$$

$$der\, c\,(r^+) \overset{def}{=} \;?$$

$$der\, c\,(r^?) \overset{def}{=} \;?$$

$$der\, c\,(r^{\{n,m\}}) \overset{def}{=} \;?$$

$$der\, c\,(\sim r) \overset{def}{=} \;?$$