# Regular Expressions

## Language Hierarchy
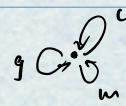
Unrestricted

Context - sensitive

Context - free

Regular { recognizer: finite Automata

Application: Tokenization



c: go to CMU. then go home

g: go to grocery. then go home

w: go for a walk, then go home

c: go to CMU once

cc: go to CMU twice

$c^*$: $\geq 0$ times go to CMU

g+w: got groceries OR a walk

$(g+w)^*$: $\geq 0$ times got groceries OR a walk

$\Sigma$ : an alphabet of characters

$\Sigma^*$ : the set of all finite-length strings over $\Sigma$ ( i.e. with chars in $\Sigma$ )

eng. "aabba" is in $\{a, b\}^*$

$\varepsilon$ : the empty string, containing no chars

A language over $\Sigma$ is a subset of $\Sigma^*$ ( may contain infinite strings, but via a finite representation )

regular expressions

A regular expression over $\Sigma$ is any of the following:

$a$    for every $a \in \Sigma$

$0$    a special symbol

$1$    a special symbol

$r_1 + r_2$   with $r_1, r_2$ regular expressions   alternation

$r_1 r_2$    with $r_1, r_2$ regular expressions   concatenation

$r^*$     with $r$ a regular expression   kleene star

# Regular Language

Given regular expression $r$, language: $L(r)$

$L(a) = \{a\}$   singleton set   $\forall a \in \Sigma$

$L(0) = \{ \}$   empty language, no strings

$L(1) = \{\varepsilon\}$   consist of empty string

$L(r_1 + r_2) = \{ s \mid s \in L(r_1) \text{ or } s \in L(r_2) \}$

$L(r_1 \cdot r_2) = \{ s_1 s_2 \mid s_1 \in L(r_1) \text{ and } s_2 \in L(r_2) \}$

$L(r^*) = \{ s \mid s = s_1 s_2 \ldots s_n, n \geq 0 \text{ with each } s_i \in L(r) \}$      $s = \varepsilon$ when $n = 0 \Rightarrow \varepsilon \in L(r^*) \ \forall \ r$

Let $L$ be a subset of $\Sigma^*$. $L$ is regular if $L = L(r)$ for some regular expression $r$.

minimal class closed under

Assume $\Sigma = \{a, b\}$

$L(a) = \{a\}$     $L((a+b)^*) = \Sigma^*$

$L(aa) = \{aa\}$   $L((a+b)^* aa (a+b)^*) = $ all strings in $\Sigma^*$ containing $\geq 2$ consecutive a's.

$L((a+1)(b+ba)^*) = $ all strings in $\Sigma^*$ not containing $\geq 2$ consecutive a   } negation!

$$L(r^*) = L(1 + rr^*)$$
<span>0 or more    0   1 or more</span>

# Acceptor

```
(* accept: regexp -> string -> bool
   REQ: true (?)
   ENS: (accept r s) => true if s ∈ L(r)
        (accept r s) => false otherwise
*)

fun accept r s = match r (String.explode s)
                       List.null
```
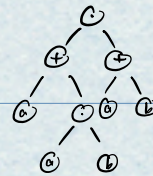    List.null              ↓
                      turn string into
   'a list -> bool        char list

Suppose  $r = (a + ab)(a + b)$
$L(r) = \{ aa, ab, aba, abb \}$



Split aba:
$(-, aba), (a, ba), (ab, a)$
×             ×          ✓

# Matcher

```
(* match : regexp -> char list -> (char list -> bool) -> bool
   REQ: k is total (?)                              matcher
   ENS:  match r cs k => true, if cs ≅ p @ s with p ∈ L(r) and k(s) ~> true
                      => false, otherwise
*)

datatype regexp =

fun match (Char a) cs k = (case cs of
                             [] => false    a ∉ []
                           | c::cs => (a = c) and also (k cs'))
| match Zero _ _ = false
| match One cs k = k cs
| match (Plus (r₁, r₂)) cs k = (match r₁ cs k) orelse (match r₂ cs k)
| match (Times (r₁, r₂)) cs k = (match r₁ cs (fn cs' => match r₂ cs' k)
| match (Star r) cs k = k cs orelse  match r cs (fn cs' => match (Star r) cs' k))   ?
```
                                          ^
                               not (cs = cs') andalso /  assume std form in REQ

$L(r^*) = L(1 + rr^*)$

```
match (Star One) [# "a"]  List.null  -> loop forever
since List.null [# "a"] ≅ false, match One cs k' passes all cs' to match
```
How to fix ? cs' must be a proper suffix of cs ! change spec / code !

# Proof

1. prove termination  (assume true)
2. prove soundness & completeness ⟵ if and only if

```
match r cs k => true, if cs ≅ p @ s with p ∈ L(r) and k(s) ~> true
             => false, otherwise
```

Base Case: Zero, One, Char (a)  ∀ a : char
IS: Plus, Times, Star

e.g. Plus



type matcher = char list -> (char list -> bool) -> bool

(true false)   orelse
               andalso

---

**Base Case**   | REJECT |   | CHECK_FOR |   Combinators  - - -
                | ACCEPT |
                                                    match

---

( char list → bool ) → bool

val REJECT : matcher = fn _ => fn _ => false   *always reject*

val ACCEPT : matcher = fn cs => fn k => k cs   *not always accept ( locally accept, pass into continuation to check the real outcome )*

fun CHECK - FOR (a : char) : matcher cs k = ( case cs of
  ( Version 1 )                                        [] => false
                                                     | c::cs' => a = c andalso k cs' )

fun CHECK - FOR (a : char) : matcher =
  ( Version 2 ) pass the work into continuation

    fn [] => REJECT []                        the input doesn't matter
     | c::cs' => if a = c then ACCEPT cs' otherwise REJECT (c::cs')

char list → ( char list → bool ) → bool

---

infix 8 ORELSE          ORELSE : matcher * matcher → matcher
infix 9 THEN            THEN   : matcher * matcher → matcher
                        REPEAT : matcher → matcher

fun (m₁ ORELSE m₂) cs k = m₁ cs k orelse m₂ cs k
fun (m₁ THEN m₂) cs k = m₁ cs (fn cs' => m₂ cs' k)
fun REPEAT m₁ cs k = (k cs) orelse m ( fn cs' => REPEAT m cs' k )

( Star Version 2 )
| match ( Star r ) cs k = let
                            fun mstar cs' = (k cs') orelse (match r cs' mstar)
                          in
                            mstar cs
                          end

( Similarly )

fun REPEAT m cs k = let
                      fun mstar cs' = (k cs') orelse ( m cs' mstar)
                    in
                      mstar cs
                    end

( Rewrite match )

fun match ( Char a : regexp ) : matcher = CHECK - FOR a

| match Zero _ _ = REJECT
| match One cs k = ACCEPT                                  *Seperate regular exp*
| match ( Plus (r₁,r₂)) cs k = match r₁ ORELSE match r₂    *from other implementations*
| match ( Times (r₁,r₂)) cs k = match r₁ THEN match r₂
| match ( Star r) cs k = REPEAT match r

---

    r = (a + b) c*
                                                    (THEN)

        (c)                    match           (ORELSE)    (REPEAT)
                              ⟹
     (⊕)    (R)                                /    \      /    \

```
(* accept: regexp -> string -> bool *)
fun accept (r: regexp): string -> bool =   fn s => match r (String.explode s) List.null

    ⇓ staging

fun accept (r: regexp): string -> bool = let
                                              val m = match r
                                          in
                                              fn s => m (String.explode s) List.null
                                          end
```

```
(* accept: regexp -> string -> bool *)
fun accept (r: regexp): string -> bool =   fn s => match r (String.explode s) List.null
```