# Strings and Regular Expressions

Stephen A. Edwards

Columbia University

Spring 2023



[^ ]*?@[^ ]*?\.[^ ]*

## Alphabets, Strings, and the Empty String

An *alphabet* $\Sigma$ is a finite set of symbols.
Strings over $\Sigma$ are members of $\Sigma^*$, defined by

$$\frac{}{\text{“”} \in \Sigma^*}\ \text{epsilon} \qquad \frac{c \in \Sigma \quad \text{“}s\text{”} \in \Sigma^*}{\text{“}cs\text{”} \in \Sigma^*}\ \text{char}$$

Judgments:   $c \in \Sigma$ "character in $\Sigma$"    $s \in \Sigma^*$ "sequence of zero or more characters"
Variables:   $c$   "character"    $s$   "sequence"
Symbols:    " "    "start and end of a string"

If $\Sigma = \{a, b, c, \ldots, z\}$,
  ""                           The empty string

$$\frac{}{\text{“”} \in \Sigma^*}\ \text{epsilon}$$

## Alphabets, Strings, and the Empty String

An *alphabet* $\Sigma$ is a finite set of symbols.
Strings over $\Sigma$ are members of $\Sigma^*$, defined by

$$\frac{}{\text{“”} \in \Sigma^*}\text{ epsilon} \qquad \frac{c \in \Sigma \quad \text{“}s\text{”} \in \Sigma^*}{\text{“}cs\text{”} \in \Sigma^*}\text{ char}$$

Judgments:   $c \in \Sigma$ "character in $\Sigma$"   $s \in \Sigma^*$ "sequence of zero or more characters"
Variables:   $c$ "character"   $s$ "sequence"
Symbols:   " " "start and end of a string"

If $\Sigma = \{a, b, c, \ldots, z\}$,
  ""                                  The empty string
  "a"                               The string consisting of just "a"

$$\frac{a \in \Sigma \quad \dfrac{}{\text{“”} \in \Sigma^*}\text{ epsilon}}{\text{“a”} \in \Sigma^*}\text{ char} \quad \leftarrow \text{Choose } s \text{ to be the empty sequence}$$

## Alphabets, Strings, and the Empty String

An *alphabet* $\Sigma$ is a finite set of symbols.

Strings over $\Sigma$ are members of $\Sigma^*$, defined by

$$\frac{}{\text{``''} \in \Sigma^*}\text{ epsilon} \qquad\qquad \frac{c \in \Sigma \quad \text{``}s\text{''} \in \Sigma^*}{\text{``}cs\text{''} \in \Sigma^*}\text{ char}$$

Judgments:   $c \in \Sigma$ "character in $\Sigma$"   $s \in \Sigma^*$ "sequence of zero or more characters"
Variables:   $c$  "character"   $s$  "sequence"
Symbols:    " "   "start and end of a string"

If $\Sigma = \{a, b, c, \ldots, z\}$,

| | |
|---|---|
| "" | The empty string |
| "a" | The string consisting of just "a" |
| "ba" | The string consisting of "b" followed by "a" |

$$\frac{b \in \Sigma \quad \dfrac{a \in \Sigma \quad \dfrac{}{\text{``''} \in \Sigma^*}\text{ epsilon}}{\text{``a''} \in \Sigma^*}\text{ char}}{\text{``ba''} \in \Sigma^*}\text{ char} \quad \leftarrow \text{Prepend characters from right to left}$$

## Alphabets, Strings, and the Empty String

An *alphabet* $\Sigma$ is a finite set of symbols.

Strings over $\Sigma$ are members of $\Sigma^*$, defined by

$$\frac{}{\text{""} \in \Sigma^*} \text{ epsilon} \qquad\qquad \frac{c \in \Sigma \quad \text{"}s\text{"} \in \Sigma^*}{\text{"}cs\text{"} \in \Sigma^*} \text{ char}$$

Judgments:    $c \in \Sigma$ "character in $\Sigma$"    $s \in \Sigma^*$ "sequence of zero or more characters"

Variables:    $c$ "character"    $s$ "sequence"

Symbols:    " "    "start and end of a string"

---

If $\Sigma = \{a, b, c, \ldots, z\}$,

| | |
|---|---|
| "" | The empty string |
| "a" | The string consisting of just "a" |
| "ba" | The string consisting of "b" followed by "a" |
| "aba" | The string "a" followed by "b" followed by "a" |
| "abcd" | The four-letter string "abcd" |
| "sphinxofblackquartzjudgemyvow" | A pangram with only a, o, and u repeated |

## Alphabets, Strings, and the Empty String

An *alphabet* $\Sigma$ is a finite set of symbols.

Strings over $\Sigma$ are members of $\Sigma^*$, defined by

$$\frac{}{\text{``''} \in \Sigma^*}\text{ epsilon} \qquad \frac{c \in \Sigma \quad \text{``}s\text{''} \in \Sigma^*}{\text{``}cs\text{''} \in \Sigma^*}\text{ char}$$

```haskell
infixr 5  :              -- Not legal Haskell    a : b : c = a : (b : c)
data [a] = []            -- [] is the empty list   "a" is a type variable
       | a : [a]         -- : is the list "cons" or prepend operator

type String = [Char]     -- In Haskell, strings are lists of characters

"Hello"                  -- shorthand for
'H' : 'e' : 'l' : 'l' : 'o' : []
```
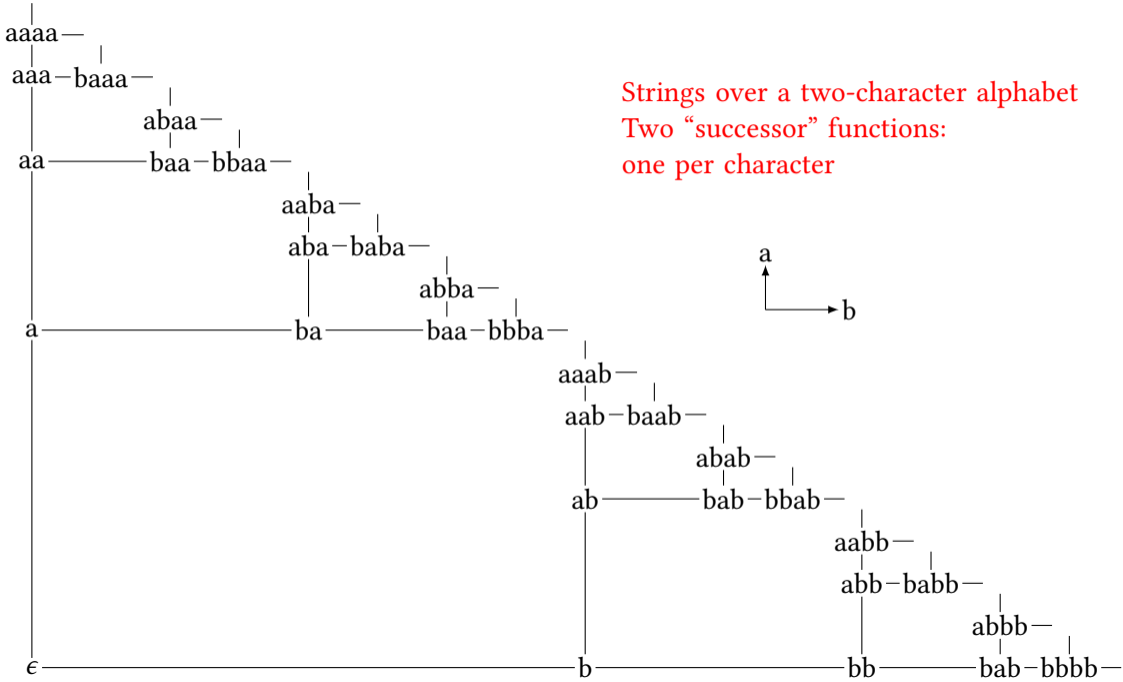
aaaa
aaa

aa

a

A visualization of the Peano encoding of natural numbers     $\epsilon$ = 0     a = successor

$\epsilon$

aaaa —

aaa — baaa —

abaa —

aa — baa — bbaa —

aaba —

aba — baba —

abba —

a — ba — baa — bbba —

aaab —

aab — baab —

abab —

ab — bab — bbab —

aabb —

abb — babb —

abbb —

ε — b — bb — bab — bbbb

Strings over a two-character alphabet
Two "successor" functions:
one per character

a
b

## Alphabets, Strings, and the Empty String

$$\overline{\text{"" } \in \Sigma^*} \text{ epsilon} \qquad \frac{c \in \Sigma \quad \text{"}s\text{" } \in \Sigma^*}{\text{"}cs\text{" } \in \Sigma^*} \text{ char}$$

## String Equality

$$\overline{\text{"" } = \text{""}} \text{ equal-epsilon} \qquad \frac{c \in \Sigma \quad \text{"}s_1\text{" } = \text{"}s_2\text{"}}{\text{"}cs_1\text{" } = \text{"}cs_2\text{"}} \text{ equal}$$

Judgments: $\quad$ "$s_1$" = "$s_2$" $\quad$ "Strings "$s_1$" and "$s_2$" are equal"

Variables: $\quad s_1, s_2 \quad$ "character sequence" $\quad c$ "character"

## Is "ab" = "ab"?

"ab" = "ab"

## Alphabets, Strings, and the Empty String

$$\frac{}{\text{“”} \in \Sigma^*} \text{ epsilon} \qquad\qquad \frac{c \in \Sigma \quad \text{“}s\text{”} \in \Sigma^*}{\text{“}cs\text{”} \in \Sigma^*} \text{ char}$$

## String Equality

$$\frac{}{\text{“”} = \text{“”}} \text{ equal-epsilon} \qquad\qquad \frac{c \in \Sigma \quad \text{“}s_1\text{”} = \text{“}s_2\text{”}}{\text{“}cs_1\text{”} = \text{“}cs_2\text{”}} \text{ equal}$$

Judgments:   $\text{“}s_1\text{”} = \text{“}s_2\text{”}$   "Strings $\text{“}s_1\text{”}$ and $\text{“}s_2\text{”}$ are equal"
Variables:   $s_1, s_2$   "character sequence"   $c$ "character"

### Is "ab" = "ab"?

$$\frac{\mathsf{a} \in \Sigma \quad \text{“b”} = \text{“b”}}{\text{“ab”} = \text{“ab”}} \text{ equal}$$

## Alphabets, Strings, and the Empty String

$$\frac{}{\text{""} \in \Sigma^*}\ \text{epsilon} \qquad\qquad \frac{c \in \Sigma \quad \text{"}s\text{"} \in \Sigma^*}{\text{"}cs\text{"} \in \Sigma^*}\ \text{char}$$

## String Equality

$$\frac{}{\text{""} = \text{""}}\ \text{equal-epsilon} \qquad\qquad \frac{c \in \Sigma \quad \text{"}s_1\text{"} = \text{"}s_2\text{"}}{\text{"}cs_1\text{"} = \text{"}cs_2\text{"}}\ \text{equal}$$

Judgments: $\text{"}s_1\text{"} = \text{"}s_2\text{"}$  "Strings $\text{"}s_1\text{"}$ and $\text{"}s_2\text{"}$ are equal"
Variables: $s_1, s_2$  "character sequence"  $c$ "character"

## Is "ab" = "ab"?

$$\frac{a \in \Sigma \quad \dfrac{b \in \Sigma \quad \dfrac{}{\text{""} = \text{""}}}{\text{"b"} = \text{"b"}}\ \text{equal}}{\text{"ab"} = \text{"ab"}}\ \text{equal}$$

## Alphabets, Strings, and the Empty String

$$\frac{}{\text{``''} \in \Sigma^*}\ \text{epsilon} \qquad\qquad \frac{c \in \Sigma \quad \text{``}s\text{''} \in \Sigma^*}{\text{``}cs\text{''} \in \Sigma^*}\ \text{char}$$

## String Equality

$$\frac{}{\text{``''} = \text{``''}}\ \text{equal-epsilon} \qquad\qquad \frac{c \in \Sigma \quad \text{``}s_1\text{''} = \text{``}s_2\text{''}}{\text{``}cs_1\text{''} = \text{``}cs_2\text{''}}\ \text{equal}$$

Judgments: $\text{``}s_1\text{''} = \text{``}s_2\text{''}$    "Strings $\text{``}s_1\text{''}$ and $\text{``}s_2\text{''}$ are equal"
Variables: $s_1, s_2$   "character sequence"   $c$ "character"

## Is "ab" = "ab"?

$$\frac{a \in \Sigma \quad \dfrac{b \in \Sigma \quad \dfrac{}{\text{``''} = \text{``''}}\ \text{equal-epsilon}}{\text{``}b\text{''} = \text{``}b\text{''}}\ \text{equal}}{\text{``}ab\text{''} = \text{``}ab\text{''}}\ \text{equal}$$

## Alphabets, Strings, and the Empty String

$$\frac{}{\text{""} \in \Sigma^*}\text{ epsilon} \qquad\qquad \frac{c \in \Sigma \quad \text{"}s\text{"} \in \Sigma^*}{\text{"}cs\text{"} \in \Sigma^*}\text{ char}$$

## String Equality

$$\frac{}{\text{""} = \text{""}}\text{ equal-epsilon} \qquad\qquad \frac{c \in \Sigma \quad \text{"}s_1\text{"} = \text{"}s_2\text{"}}{\text{"}cs_1\text{"} = \text{"}cs_2\text{"}}\text{ equal}$$

Judgments:   $\text{"}s_1\text{"} = \text{"}s_2\text{"}$    "Strings $\text{"}s_1\text{"}$ and $\text{"}s_2\text{"}$ are equal"
Variables:   $s_1, s_2$   "character sequence"   $c$ "character"

## Is "ab" = "ac"?

"ab" = "ac"

## Alphabets, Strings, and the Empty String

$$\frac{}{\text{""} \in \Sigma^*}\text{ epsilon} \qquad\qquad \frac{c \in \Sigma \quad \text{``}s\text{''} \in \Sigma^*}{\text{``}cs\text{''} \in \Sigma^*}\text{ char}$$

## String Equality

$$\frac{}{\text{""} = \text{""}}\text{ equal-epsilon} \qquad\qquad \frac{c \in \Sigma \quad \text{``}s_1\text{''} = \text{``}s_2\text{''}}{\text{``}cs_1\text{''} = \text{``}cs_2\text{''}}\text{ equal}$$

Judgments:  $\text{``}s_1\text{''} = \text{``}s_2\text{''}$  "Strings $\text{``}s_1\text{''}$ and $\text{``}s_2\text{''}$ are equal"
Variables:  $s_1, s_2$  "character sequence"  $c$ "character"

## Is "ab" = "ac"?

$$\frac{a \in \Sigma \quad \text{``}b\text{''} = \text{``}c\text{''}}{\text{``}ab\text{''} = \text{``}ac\text{''}}\text{ equal}$$

## Alphabets, Strings, and the Empty String

$$\frac{}{\text{“”} \in \Sigma^*} \text{ epsilon}$$

$$\frac{c \in \Sigma \quad \text{“}s\text{”} \in \Sigma^*}{\text{“}cs\text{”} \in \Sigma^*} \text{ char}$$

## String Equality

$$\frac{}{\text{“”} = \text{“”}} \text{ equal-epsilon}$$

$$\frac{c \in \Sigma \quad \text{“}s_1\text{”} = \text{“}s_2\text{”}}{\text{“}cs_1\text{”} = \text{“}cs_2\text{”}} \text{ equal}$$

Judgments:   $\text{“}s_1\text{”} = \text{“}s_2\text{”}$   "Strings $\text{“}s_1\text{”}$ and $\text{“}s_2\text{”}$ are equal"
Variables:   $s_1, s_2$   "character sequence"   $c$ "character"

## Is "ab" = "ac"?

$$\frac{a \in \Sigma \quad \dfrac{?}{\text{“b”} = \text{“c”}} ?}{\text{“ab”} = \text{“ac”}} \text{ equal}$$

← We are stuck: the *equal* rule requires identical initial characters

## Alphabets, Strings, and the Empty String

$$\frac{}{\text{``''} \in \Sigma^*} \text{ epsilon} \qquad\qquad \frac{c \in \Sigma \quad \text{``}s\text{''} \in \Sigma^*}{\text{``}cs\text{''} \in \Sigma^*} \text{ char}$$

## String Equality

$$\frac{}{\text{``''} = \text{``''}} \text{ equal-epsilon} \qquad\qquad \frac{c \in \Sigma \quad \text{``}s_1\text{''} = \text{``}s_2\text{''}}{\text{``}cs_1\text{''} = \text{``}cs_2\text{''}} \text{ equal}$$

Judgments:     "$s_1$" = "$s_2$"     "Strings "$s_1$" and "$s_2$" are equal"
Variables:     $s_1, s_2$     "character sequence"     $c$ "character"

## Additional Theorems

Reflexive: For any $s \in \Sigma^*$, $s = s$
Symmetric: For any $s_1, s_2 \in \Sigma^*$ with $s_1 = s_2$, $s_2 = s_1$.
Transitive: For any $s_1, s_2, s_3 \in \Sigma^*$ with $s_1 = s_2$ and $s_2 = s_3$, $s_1 = s_3$.

## Alphabets, Strings, and the Empty String

$$\frac{}{\text{``''} \in \Sigma^*} \text{ epsilon} \qquad \frac{c \in \Sigma \quad \text{``}s\text{''} \in \Sigma^*}{\text{``}cs\text{''} \in \Sigma^*} \text{ char}$$

## String Equality

$$\frac{}{\text{``''} = \text{``''}} \text{ equal-epsilon} \qquad \frac{c \in \Sigma \quad \text{``}s_1\text{''} = \text{``}s_2\text{''}}{\text{``}cs_1\text{''} = \text{``}cs_2\text{''}} \text{ equal}$$

Judgments: $\text{``}s_1\text{''} = \text{``}s_2\text{''}$    "Strings $\text{``}s_1\text{''}$ and $\text{``}s_2\text{''}$ are equal"

Variables: $s_1, s_2$   "character sequence"   $c$ "character"

```haskell
(==)                    :: [Char] -> [Char] -> Bool
     [] == []          = True              -- equal-epsilon
c1 : s1 == c2 : s2      = c1 == c2 && s1 == s2   -- equal
     _ == _            = False             -- default case

data [a] = [] | a : [a] deriving (Eq, Ord) -- Default implementation of Eq
```

## Alphabets, Strings, and the Empty String

$$\frac{}{\text{``''} \in \Sigma^*} \text{ epsilon} \qquad \frac{c \in \Sigma \quad \text{``}s\text{''} \in \Sigma^*}{\text{``}cs\text{''} \in \Sigma^*} \text{ char}$$

## String Concatenation

$$\frac{\text{``}s\text{''} \in \Sigma^*}{\text{``''} ++ \text{``}s\text{''} = \text{``}s\text{''}} \text{ concat-epsilon} \qquad \frac{c \in \Sigma \quad \text{``}s_1\text{''} ++ \text{``}s_2\text{''} = \text{``}s_3\text{''}}{\text{``}cs_1\text{''} ++ \text{``}s_2\text{''} = \text{``}cs_3\text{''}} \text{ concat}$$

Judgments:     "$s_1$" ++ "$s_2$" = "$s_3$"    "Concatenating strings "$s_1$" and "$s_2$" gives string "$s_3$""
Variables:     $s, s_1, s_2, s_3$     "character sequence"    $c$ "character"

## Alphabets, Strings, and the Empty String

$$\frac{}{\text{``''} \in \Sigma^*} \text{ epsilon} \qquad \frac{c \in \Sigma \quad \text{``}s\text{''} \in \Sigma^*}{\text{``}cs\text{''} \in \Sigma^*} \text{ char}$$

## String Concatenation

$$\frac{\text{``}s\text{''} \in \Sigma^*}{\text{``''} ++ \text{``}s\text{''} = \text{``}s\text{''}} \text{ concat-epsilon} \qquad \frac{c \in \Sigma \quad \text{``}s_1\text{''} ++ \text{``}s_2\text{''} = \text{``}s_3\text{''}}{\text{``}cs_1\text{''} ++ \text{``}s_2\text{''} = \text{``}cs_3\text{''}} \text{ concat}$$

## Is "ab" ++ "cde" = "abcde"?

"ab" ++ "cde" = "abcde"

## Alphabets, Strings, and the Empty String

$$\frac{}{\text{``''} \in \Sigma^*}\text{ epsilon} \qquad \frac{c \in \Sigma \quad \text{``}s\text{''} \in \Sigma^*}{\text{``}cs\text{''} \in \Sigma^*}\text{ char}$$

## String Concatenation

$$\frac{\text{``}s\text{''} \in \Sigma^*}{\text{``''} ++ \text{``}s\text{''} = \text{``}s\text{''}}\text{ concat-epsilon} \qquad \frac{c \in \Sigma \quad \text{``}s_1\text{''} ++ \text{``}s_2\text{''} = \text{``}s_3\text{''}}{\text{``}cs_1\text{''} ++ \text{``}s_2\text{''} = \text{``}cs_3\text{''}}\text{ concat}$$

## Is "ab" ++ "cde" = "abcde"?

$$\frac{a \in \Sigma \qquad \text{``b''} ++ \text{``cde''} = \text{``bcde''}}{\text{``ab''} ++ \text{``cde''} = \text{``abcde''}}\text{ concat}$$

## Alphabets, Strings, and the Empty String

$$\frac{}{\text{``''} \in \Sigma^*}\ \text{epsilon} \qquad\qquad \frac{c \in \Sigma \quad \text{``}s\text{''} \in \Sigma^*}{\text{``}cs\text{''} \in \Sigma^*}\ \text{char}$$

## String Concatenation

$$\frac{\text{``}s\text{''} \in \Sigma^*}{\text{``''} {+}{+}\ \text{``}s\text{''} = \text{``}s\text{''}}\ \text{concat-epsilon} \qquad \frac{c \in \Sigma \quad \text{``}s_1\text{''} {+}{+}\ \text{``}s_2\text{''} = \text{``}s_3\text{''}}{\text{``}cs_1\text{''} {+}{+}\ \text{``}s_2\text{''} = \text{``}cs_3\text{''}}\ \text{concat}$$

## Is "ab" ++ "cde" = "abcde"?

$$\frac{\mathsf{a} \in \Sigma \quad \dfrac{\mathsf{b} \in \Sigma \quad \dfrac{\text{``''} {+}{+}\ \text{``cde''} = \text{``cde''}}{\text{``b''} {+}{+}\ \text{``cde''} = \text{``bcde''}}\ \text{concat}}{\text{``ab''} {+}{+}\ \text{``cde''} = \text{``abcde''}}}{}\ \text{concat}$$

## Alphabets, Strings, and the Empty String

$$\frac{}{\text{``''} \in \Sigma^*}\text{ epsilon} \qquad \frac{c \in \Sigma \quad \text{``}s\text{''} \in \Sigma^*}{\text{``}cs\text{''} \in \Sigma^*}\text{ char}$$

## String Concatenation

$$\frac{\text{``}s\text{''} \in \Sigma^*}{\text{``''} \mathbin{++} \text{``}s\text{''} = \text{``}s\text{''}}\text{ concat-epsilon} \qquad \frac{c \in \Sigma \quad \text{``}s_1\text{''} \mathbin{++} \text{``}s_2\text{''} = \text{``}s_3\text{''}}{\text{``}cs_1\text{''} \mathbin{++} \text{``}s_2\text{''} = \text{``}cs_3\text{''}}\text{ concat}$$

## Is "ab" ++ "cde" = "abcde"?

$$\cfrac{a \in \Sigma \quad \cfrac{b \in \Sigma \quad \cfrac{\text{``cde''} \in \Sigma^*}{\text{``''} \mathbin{++} \text{``cde''} = \text{``cde''}}\text{ concat-epsilon}}{\text{``b''} \mathbin{++} \text{``cde''} = \text{``bcde''}}\text{ concat}}{\text{``ab''} \mathbin{++} \text{``cde''} = \text{``abcde''}}\text{ concat}$$

## Alphabets, Strings, and the Empty String

$$\frac{}{\text{“”} \in \Sigma^*}\ \text{epsilon} \qquad\qquad \frac{c \in \Sigma \quad \text{“}s\text{”} \in \Sigma^*}{\text{“}cs\text{”} \in \Sigma^*}\ \text{char}$$

## String Concatenation

$$\frac{\text{“}s\text{”} \in \Sigma^*}{\text{“”} ++ \text{“}s\text{”} = \text{“}s\text{”}}\ \text{concat-epsilon} \qquad \frac{c \in \Sigma \quad \text{“}s_1\text{”} ++ \text{“}s_2\text{”} = \text{“}s_3\text{”}}{\text{“}cs_1\text{”} ++ \text{“}s_2\text{”} = \text{“}cs_3\text{”}}\ \text{concat}$$

## Is "ab" ++ "cde" = "abcde"?

$$\cfrac{a \in \Sigma \quad \cfrac{b \in \Sigma \quad \cfrac{c \in \Sigma \quad \cfrac{d \in \Sigma \quad \cfrac{e \in \Sigma \quad \cfrac{}{\text{“”} \in \Sigma^*}\ \text{epsilon}}{\text{“}e\text{”} \in \Sigma^*}\ \text{char}}{\text{“}de\text{”} \in \Sigma^*}\ \text{char}}{\text{“}cde\text{”} \in \Sigma^*}\ \text{char}}{\text{“”} ++ \text{“cde”} = \text{“cde”}}\ \text{concat-epsilon}}{\text{“b”} ++ \text{“cde”} = \text{“bcde”}}\ \text{concat}}{\text{“ab”} ++ \text{“cde”} = \text{“abcde”}}\ \text{concat}$$

## Alphabets, Strings, and the Empty String

$$\frac{}{\text{``''} \in \Sigma^*}\text{ epsilon} \qquad\qquad \frac{c \in \Sigma \quad \text{``}s\text{''} \in \Sigma^*}{\text{``}cs\text{''} \in \Sigma^*}\text{ char}$$

## String Concatenation

$$\frac{\text{``}s\text{''} \in \Sigma^*}{\text{``''} ++ \text{``}s\text{''} = \text{``}s\text{''}}\text{ concat-epsilon} \qquad \frac{c \in \Sigma \quad \text{``}s_1\text{''} ++ \text{``}s_2\text{''} = \text{``}s_3\text{''}}{\text{``}cs_1\text{''} ++ \text{``}s_2\text{''} = \text{``}cs_3\text{''}}\text{ concat}$$

## Theorem: "" is also a right identity for concatenation

Assume "$s$" $\in \Sigma^*$. From the definition of $\Sigma^*$, $s$ must be of the form $c_1 c_2 \cdots c_n$ where $c_i \in \Sigma$.

$$\frac{c_1 \in \Sigma \quad \frac{\vdots}{\frac{c_n \in \Sigma \quad \dfrac{\dfrac{}{\text{``''} \in \Sigma^*}\text{ epsilon}}{\text{``''} ++ \text{``''} = \text{``''}}\text{ concat-epsilon}}{\text{``}c_n\text{''} ++ \text{``''} = \text{``}c_n\text{''}}\text{ concat}}{\text{``}c_2 \cdots c_n\text{''} ++ \text{``''} = \text{``}c_2 \cdots c_n\text{''}}}{\text{``}c_1 c_2 \cdots c_n\text{''} ++ \text{``''} = \text{``}c_1 c_2 \cdots c_n\text{''}}\text{ concat}$$

## Alphabets, Strings, and the Empty String

$$\frac{}{\text{``''} \in \Sigma^*}\ \text{epsilon} \qquad\qquad \frac{c \in \Sigma \quad \text{``}s\text{''} \in \Sigma^*}{\text{``}cs\text{''} \in \Sigma^*}\ \text{char}$$

## String Concatenation

$$\frac{\text{``}s\text{''} \in \Sigma^*}{\text{``''} ++ \text{``}s\text{''} = \text{``}s\text{''}}\ \text{concat-epsilon} \qquad \frac{c \in \Sigma \quad \text{``}s_1\text{''} ++ \text{``}s_2\text{''} = \text{``}s_3\text{''}}{\text{``}cs_1\text{''} ++ \text{``}s_2\text{''} = \text{``}cs_3\text{''}}\ \text{concat}$$

## Theorem: String concatenation is a function

If "$s_1$" ++ "$s_2$" = "$s_3$" and "$s_1$" ++ "$s_2$" = "$s_4$" then "$s_3$" = "$s_4$".

## Alphabets, Strings, and the Empty String

$$\frac{}{\text{``''} \in \Sigma^*} \text{ epsilon} \qquad \frac{c \in \Sigma \quad \text{``}s\text{''} \in \Sigma^*}{\text{``}cs\text{''} \in \Sigma^*} \text{ char}$$

## String Concatenation

$$\frac{\text{``}s\text{''} \in \Sigma^*}{\text{``''} \,{+}{+}\, \text{``}s\text{''} = \text{``}s\text{''}} \text{ concat-epsilon} \qquad \frac{c \in \Sigma \quad \text{``}s_1\text{''} \,{+}{+}\, \text{``}s_2\text{''} = \text{``}s_3\text{''}}{\text{``}cs_1\text{''} \,{+}{+}\, \text{``}s_2\text{''} = \text{``}cs_3\text{''}} \text{ concat}$$

```
infixr 5  :                  -- Not legal Haskell    a : b : c = a : (b : c)
data [a] = [] | a : [a]      -- Not legal Haskell    [] is empty list   : is cons
type String = [Char]


infixr 5  ++                 -- a ++ b ++ c = a ++ (b ++ c)
(++) :: [a] -> [a] -> [a]    -- Concatenate two lists
(++) []      s  = s          -- concat-epsilon
(++) (c:s1) s2 = c : s1 ++ s2 -- concat    c : (s1 ++ s2)
```

### Regular Expressions

| | |
|---|---|
| A character matches itself | "a" $\sim$ a    "x" $\sim$ x |
| Juxtaposition matches a sequence | "abc" $\sim$ abc |
| \| indicates a choice | "ab" $\sim$ ab\|bc    "bc" $\sim$ ab\|bc |
| $\star$ means "zero or more" | "", "a", "aa", "aaa", "aaaa", "aaaaa", ... $\sim$ a$\star$ |

## Regular Expressions

A character matches itself $\quad\quad$ "a" $\sim$ a $\quad$ "x" $\sim$ x

Juxtaposition matches a sequence $\quad$ "abc" $\sim$ abc

$$\frac{}{\text{""} \sim \epsilon}\,\text{epsilon} \quad\quad \frac{c \in \Sigma}{\text{"}c\text{"} \sim c}\,\text{char} \quad\quad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \mathbin{++} s_2 = s_3}{s_3 \sim r_1 r_2}\,\text{seq}$$

Judgments: $\quad s \sim r \quad$ "string $s$ matches regular expression $r$"

Variables: $\quad c$ character $\quad r$ regular expression $\quad s$ string

Symbols: $\quad \epsilon \quad$ " " $\quad | \quad * \quad$ a b c d...

## Regular Expressions

A character matches itself $\qquad$ "a" $\sim$ a $\quad$ "x" $\sim$ x

Juxtaposition matches a sequence $\quad$ "abc" $\sim$ abc

$$\frac{}{\text{""} \sim \epsilon}\text{epsilon} \qquad \frac{c \in \Sigma}{\text{"}c\text{"} \sim c}\text{char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \mathbin{++} s_2 = s_3}{s_3 \sim r_1 r_2}\text{seq}$$

## "abc" $\sim$ abc?

$$\text{"abc"} \sim \text{abc}$$

## Regular Expressions

A character matches itself $\qquad$ "a" $\sim$ a $\quad$ "x" $\sim$ x

Juxtaposition matches a sequence $\qquad$ "abc" $\sim$ abc

$$\frac{}{\text{""} \sim \epsilon}\,\text{epsilon} \qquad \frac{c \in \Sigma}{\text{"}c\text{"} \sim c}\,\text{char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \,\texttt{++}\, s_2 = s_3}{s_3 \sim r_1 r_2}\,\text{seq}$$

## "abc" $\sim$ abc?

$$\frac{\text{"a"} \sim \text{a} \qquad\qquad \text{"bc"} \sim \text{bc} \qquad\qquad \text{"a"} \,\texttt{++}\, \text{"bc"} = \text{"abc"}}{\text{"abc"} \sim \text{abc}}\,\text{seq}$$

## Regular Expressions

A character matches itself           "a" $\sim$ a     "x" $\sim$ x

Juxtaposition matches a sequence    "abc" $\sim$ abc

$$\frac{}{\text{""} \sim \epsilon}\text{ epsilon} \qquad \frac{c \in \Sigma}{\text{"}c\text{"} \sim c}\text{ char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \mathrel{++} s_2 = s_3}{s_3 \sim r_1 r_2}\text{ seq}$$

## "abc" $\sim$ abc?

$$\frac{\dfrac{a \in \Sigma}{\text{"a"} \sim a}\text{char} \quad \dfrac{\text{"b"} \sim b \quad \text{"c"} \sim c \quad \text{"b"} \mathrel{++} \text{"c"} = \text{"bc"}}{\text{"bc"} \sim bc}\text{seq} \quad \dfrac{\vdots}{\text{"a"} \mathrel{++} \text{"bc"} = \text{"abc"}}\text{concat}}{\text{"abc"} \sim abc}\text{seq}$$

## Regular Expressions

A character matches itself      "a" $\sim$ a    "x" $\sim$ x

Juxtaposition matches a sequence      "abc" $\sim$ abc

$$\frac{}{\text{""} \sim \epsilon}\text{epsilon} \qquad \frac{c \in \Sigma}{\text{"}c\text{"} \sim c}\text{char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \mathbin{++} s_2 = s_3}{s_3 \sim r_1 r_2}\text{seq}$$

## "abc" $\sim$ abc?

$$\frac{a \in \Sigma}{\text{"a"} \sim \text{a}}\text{char} \quad \frac{\dfrac{b \in \Sigma}{\text{"b"} \sim \text{b}}\text{char} \quad \dfrac{c \in \Sigma}{\text{"c"} \sim \text{c}}\text{char} \quad \dfrac{\vdots}{\text{"b"} \mathbin{++} \text{"c"} = \text{"bc"}}\text{concat}}{\text{"bc"} \sim \text{bc}}\text{seq} \quad \frac{\vdots}{\text{"a"} \mathbin{++} \text{"bc"} = \text{"abc"}}\text{concat}}{\text{"abc"} \sim \text{abc}}\text{seq}$$

## Regular Expressions

A character matches itself $\quad$ "a" $\sim$ a $\quad$ "x" $\sim$ x

Juxtaposition matches a sequence $\quad$ "abc" $\sim$ abc

$$\frac{}{\text{""} \sim \epsilon}\,\text{epsilon} \qquad \frac{c \in \Sigma}{\text{"c"} \sim c}\,\text{char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 +\!+ s_2 = s_3}{s_3 \sim r_1 r_2}\,\text{seq}$$

## "abc" $\sim$ abc? $\quad$ What should we do about ambiguity?

$$\frac{a \in \Sigma}{\text{"a"} \sim a}\,\text{char} \quad \frac{\dfrac{b \in \Sigma}{\text{"b"} \sim b}\,\text{char} \quad \dfrac{c \in \Sigma}{\text{"c"} \sim c}\,\text{char} \quad \dfrac{\vdots}{\text{"b"} +\!+ \text{"c"} = \text{"bc"}}\,\text{concat}}{\text{"bc"} \sim bc}\,\text{seq} \quad \frac{\vdots}{\text{"a"} +\!+ \text{"bc"} = \text{"abc"}}\,\text{concat}}{\text{"abc"} \sim abc}\,\text{seq}$$

$$\frac{\dfrac{a \in \Sigma}{\text{"a"} \sim a}\,\text{char} \quad \dfrac{b \in \Sigma}{\text{"b"} \sim b}\,\text{char} \quad \dfrac{\vdots}{\text{"a"} +\!+ \text{"b"} = \text{"ab"}}\,\text{concat}}{\text{"ab"} \sim ab}\,\text{seq} \quad \dfrac{c \in \Sigma}{\text{"c"} \sim c}\,\text{char} \quad \dfrac{\vdots}{\text{"ab"} +\!+ \text{"c"} = \text{"abc"}}\,\text{concat}}{\text{"abc"} \sim abc}\,\text{seq}$$

$$\frac{}{\text{“”} \sim \epsilon} \text{epsilon} \qquad \frac{c \in \Sigma}{\text{“}c\text{”} \sim c} \text{char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \mathbin{+\!\!+} s_2 = s_3}{s_3 \sim r_1 r_2} \text{seq} \qquad \frac{s \sim r}{s \sim (r)} \text{paren}$$

Judgments:   $s \sim r$   "string $s$ matches regular expression $r$"
Variables:   $c$ character   $r$ regular expression   $s$ string
Symbols:   $\epsilon$   " "   |   $\star$   ( )

"abc" $\sim$ a(bc)?

$$\frac{\dfrac{a \in \Sigma}{\text{“a”} \sim a} \text{char} \qquad \text{“bc”} \sim \text{(bc)} \qquad \dfrac{\vdots}{\text{“a”} \mathbin{+\!\!+} \text{“bc”} = \text{“abc”}} \text{concat}}{\text{“abc”} \sim \text{a(bc)}} \text{seq}$$

$$\frac{}{\text{``''} \sim \epsilon} \text{ epsilon} \qquad \frac{c \in \Sigma}{\text{``}c\text{''} \sim c} \text{ char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \mathbin{++} s_2 = s_3}{s_3 \sim r_1 r_2} \text{ seq} \qquad \frac{s \sim r}{s \sim (r)} \text{ paren}$$

Judgments: $\quad s \sim r \quad$ "string $s$ matches regular expression $r$"

Variables: $\quad c$ character $\quad r$ regular expression $\quad s$ string

Symbols: $\quad \epsilon \quad$ " " $\quad | \quad * \quad ( \quad )$

## "abc" $\sim$ a(bc)?

$$\frac{\dfrac{a \in \Sigma}{\text{``a''} \sim \text{a}} \text{ char} \quad \dfrac{\dfrac{\text{``bc''} \sim \text{bc}}{\text{``bc''} \sim \text{(bc)}} \text{ paren} \quad \dfrac{\vdots}{\text{``a''} \mathbin{++} \text{``bc''} = \text{``abc''}} \text{ concat}}{\text{``abc''} \sim \text{a(bc)}}}{} \text{ seq}$$

## Regular Expressions

$$\frac{}{\text{“”} \sim \epsilon}\,\text{epsilon} \qquad \frac{c \in \Sigma}{\text{“}c\text{”} \sim c}\,\text{char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \mathbin{++} s_2 = s_3}{s_3 \sim r_1 r_2}\,\text{seq} \qquad \frac{s \sim r}{s \sim (r)}\,\text{paren}$$

Judgments: $\quad s \sim r \quad$ "string $s$ matches regular expression $r$"

Variables: $\quad c$ character $\quad r$ regular expression $\quad s$ string

Symbols: $\quad \epsilon \quad$ “ ” $\quad | \quad * \quad ( \ )$

## "abc" $\sim$ a(bc)?

$$\cfrac{a \in \Sigma}{\text{“a”} \sim \mathsf{a}}\,\text{char} \quad \cfrac{\cfrac{\cfrac{b \in \Sigma}{\text{“b”} \sim \mathsf{b}}\,\text{char} \quad \cfrac{c \in \Sigma}{\text{“c”} \sim \mathsf{c}}\,\text{char} \quad \cfrac{\vdots}{\text{“b”} \mathbin{++} \text{“c”} = \text{“bc”}}\,\text{concat}}{\cfrac{\text{“bc”} \sim \mathsf{bc}}{\text{“bc”} \sim \mathsf{(bc)}}\,\text{paren}}\,\text{seq} \quad \cfrac{\vdots}{\text{“a”} \mathbin{++} \text{“bc”} = \text{“abc”}}\,\text{concat}}{\text{“abc”} \sim \mathsf{a(bc)}}\,\text{seq}$$

## Regular Expressions

$$\frac{}{\text{``''} \sim \epsilon}\;\text{epsilon} \qquad \frac{c \in \Sigma}{\text{``}c\text{''} \sim c}\;\text{char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \mathbin{+\!\!+} s_2 = s_3}{s_3 \sim r_1 r_2}\;\text{seq}$$

```haskell
-- An algebraic data type resolves ambiguity in RE structure (parentheses unneeded)
data RE = Epsilon    -- Epsilon/empty string
        | Ch Char    -- Single Character
        | Seq RE RE  -- Sequence, e.g., r1 r2

infix 5 ~~           -- Regular expression match operator (Haskell already uses ~)
(~~)            :: String -> RE -> Bool
""   ~~ Epsilon    = True       -- epsilon
[c1] ~~ Ch c2      = c1 == c2    -- char
s3   ~~ Seq r1 r2  =             -- What to do for seq? How do we choose s1, s2?
_    ~~ _          = False       -- default
```

$$\frac{}{\text{"" } \sim \epsilon} \text{ epsilon} \qquad \frac{c \in \Sigma}{\text{"}c\text{" } \sim c} \text{ char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \mathbin{++} s_2 = s_3}{s_3 \sim r_1 r_2} \text{ seq}$$

```
ghci> import Data.List (inits, tails)
ghci> inits "abc"      -- All prefixes, shortest first
["","a","ab","abc"]

ghci> tails "abc"      -- All suffixes, longest first
["abc","bc","c",""]

ghci> :t zipWith       -- zipWith f [a₁, a₂, ...] [b₁, b₂, ...] = [f a₁ b₁, f a₂ b₂, ...]
zipWith :: (a -> b -> c) -> [a] -> [b] -> [c]

ghci> :t or            -- Logical OR of a list of Booleans
or :: Foldable t => t Bool -> Bool
```

## Regular Expressions

$$\overline{\text{``''} \sim \epsilon}\text{ epsilon} \qquad \frac{c \in \Sigma}{\text{``}c\text{''} \sim c}\text{ char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 +\!\!+ s_2 = s_3}{s_3 \sim r_1 r_2}\text{ seq}$$

```haskell
import Data.List (inits, tails)
data RE = Epsilon | Ch Char | Seq RE RE

infix 5 ~~
(~~)              :: String -> RE -> Bool
""   ~~ Epsilon   = True       -- epsilon
[c1] ~~ Ch c2     = c1 == c2   -- char
s3   ~~ Seq r1 r2 = or $ zipWith testSplit (inits s3) (tails s3) -- seq
                    where testSplit s1 s2 = s1 ~~ r1 && s2 ~~ r2
_    ~~ _         = False      -- default
```

## Regular Expressions

$$\frac{}{\text{“”} \sim \epsilon}\,\text{epsilon} \qquad \frac{c \in \Sigma}{\text{“}c\text{”} \sim c}\,\text{char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \mathbin{++} s_2 = s_3}{s_3 \sim r_1 r_2}\,\text{seq}$$

$$\frac{s \sim r_1}{s \sim r_1 | r_2}\,\text{alt-l} \qquad \frac{s \sim r_2}{s \sim r_1 | r_2}\,\text{alt-r}$$

## Regular Expressions

$$\overline{\text{``''} \sim \epsilon} \text{ epsilon} \qquad \frac{c \in \Sigma}{\text{``}c\text{''} \sim c} \text{ char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \mathbin{++} s_2 = s_3}{s_3 \sim r_1 r_2} \text{ seq}$$

$$\frac{s \sim r_1}{s \sim r_1 | r_2} \text{ alt-l} \qquad \frac{s \sim r_2}{s \sim r_1 | r_2} \text{ alt-r}$$

```
data RE = Epsilon | Ch Char | Seq RE RE | Alt RE RE

(~~)                 :: String -> RE -> Bool
""   ~~ Epsilon    = True                          -- epsilon
[c1] ~~ Ch c2      = c1 == c2                       -- char
s3   ~~ Seq r1 r2  = or $ zipWith testSplit (inits s3) (tails s3) -- seq
                     where testSplit s1 s2 = s1 ~~ r1 && s2 ~~ r2
s    ~~ Alt r1 r2  = s ~~ r1 || s ~~ r2             -- alt-l and alt-r
_    ~~ _          = False                          -- default
```

## Regular Expressions

$$\frac{}{\text{``''} \sim \epsilon} \text{ epsilon} \qquad \frac{c \in \Sigma}{\text{``}c\text{''} \sim c} \text{ char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \mathbin{++} s_2 = s_3}{s_3 \sim r_1 r_2} \text{ seq}$$

$$\frac{s \sim r_1}{s \sim r_1 | r_2} \text{ alt-l} \qquad \frac{s \sim r_2}{s \sim r_1 | r_2} \text{ alt-r} \qquad \frac{}{\text{``''} \sim r\star} \text{ star-0} \qquad \frac{s_1 \sim r \quad s_2 \sim r\star \quad s_1 \mathbin{++} s_2 = s_3}{s_3 \sim r\star} \text{ star-1}$$

## Regular Expressions

$$\overline{\text{“”} \sim \epsilon}\ \text{epsilon} \qquad \frac{c \in \Sigma}{\text{“}c\text{”} \sim c}\ \text{char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \mathbin{++} s_2 = s_3}{s_3 \sim r_1 r_2}\ \text{seq}$$

$$\frac{s \sim r_1}{s \sim r_1 | r_2}\ \text{alt-l} \qquad \frac{s \sim r_2}{s \sim r_1 | r_2}\ \text{alt-r} \qquad \overline{\text{“”} \sim r\star}\ \text{star-0} \qquad \frac{s_1 \sim r \quad s_2 \sim r\star \quad s_1 \mathbin{++} s_2 = s_3}{s_3 \sim r\star}\ \text{star-1}$$

```haskell
data RE = Epsilon | Ch Char | Seq RE RE | Alt RE RE | Star RE
(~~)                 :: String -> RE -> Bool      -- HANGS TESTING "b" ∼ ε*
""   ~~ Epsilon   = True                          -- epsilon
[c1] ~~ Ch c2     = c1 == c2                       -- char
s3   ~~ Seq r1 r2 = or $ zipWith testSplit (inits s3) (tails s3)
                    where testSplit s1 s2 = s1 ~~ r1 && s2 ~~ r2
s    ~~ Alt r1 r2 = s ~~ r1 || s ~~ r2            -- alt-l and alt-r
""   ~~ Star _    = True                           -- star-0
s3   ~~ Star r    = or $ zipWith testSplit (inits s3) (tails s3)
                    where testSplit s1 s2 = s1 ~~ r && s2 ~~ (Star r)
_    ~~ _         = False                          -- default
```

## Regular Expressions

$$\overline{\text{""} \sim \epsilon}\ \text{epsilon} \qquad \frac{c \in \Sigma}{\text{"}c\text{"} \sim c}\ \text{char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \mathbin{+\!\!+} s_2 = s_3}{s_3 \sim r_1 r_2}\ \text{seq}$$

$$\frac{s \sim r_1}{s \sim r_1 | r_2}\ \text{alt-l} \qquad \frac{s \sim r_2}{s \sim r_1 | r_2}\ \text{alt-r} \qquad \overline{\text{""} \sim r\star}\ \text{star-0} \qquad \frac{s_1 \sim r \quad s_2 \sim r\star \quad s_1 \mathbin{+\!\!+} s_2 = s_3}{s_3 \sim r\star}\ \text{star-1}$$

```haskell
data RE = Epsilon | Ch Char | Seq RE RE | Alt RE RE | Star RE
(~~)              :: String -> RE -> Bool    -- STILL HANGS ON "b" ~ ε*
""   ~~ Epsilon   = True                     -- epsilon
[c1] ~~ Ch c2     = c1 == c2                 -- char
s3   ~~ Seq r1 r2 = or $ zipWith testSplit (inits s3) (tails s3)
                    where testSplit s1 s2 = s1 ~~ r1 && s2 ~~ r2
s    ~~ Alt r1 r2 = s ~~ r1 || s ~~ r2       -- alt-l and alt-r
""   ~~ Star _    = True                      -- star-0
s3   ~~ Star r    = or $ zipWith testSplit (inits s3) (tails s3)
                    where testSplit s1 s2 = s2 ~~ (Star r) && s1 ~~ r
_    ~~ _         = False                      -- default
```

## Regular Expressions

$$\frac{}{\text{``''} \sim \epsilon} \text{ epsilon} \qquad \frac{c \in \Sigma}{\text{``}c\text{''} \sim c} \text{ char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \,{+}{+}\, s_2 = s_3}{s_3 \sim r_1 r_2} \text{ seq}$$

$$\frac{s \sim r_1}{s \sim r_1 | r_2} \text{ alt-l} \quad \frac{s \sim r_2}{s \sim r_1 | r_2} \text{ alt-r} \quad \frac{}{\text{``''} \sim r\star} \text{ star-0} \quad \frac{s_1 \sim r \quad s_2 \sim r\star \quad s_1 \,{+}{+}\, s_2 = s_3}{s_3 \sim r\star} \text{ star-1}$$

"b" $\sim \epsilon\star$?

"b" $\sim \epsilon\star$

## Regular Expressions

$$\frac{}{\text{“”} \sim \epsilon} \text{ epsilon} \qquad \frac{c \in \Sigma}{\text{“}c\text{”} \sim c} \text{ char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \mathbin{++} s_2 = s_3}{s_3 \sim r_1 r_2} \text{ seq}$$

$$\frac{s \sim r_1}{s \sim r_1 | r_2} \text{ alt-l} \qquad \frac{s \sim r_2}{s \sim r_1 | r_2} \text{ alt-r} \qquad \frac{}{\text{“”} \sim r\star} \text{ star-0} \qquad \frac{s_1 \sim r \quad s_2 \sim r\star \quad s_1 \mathbin{++} s_2 = s_3}{s_3 \sim r\star} \text{ star-1}$$

### “b” $\sim \epsilon\star$?

$$\frac{\text{“”} \sim \epsilon \qquad \text{“b”} \sim \epsilon\star \qquad \text{“”} \mathbin{++} \text{“b”} = \text{“b”}}{\text{“b”} \sim \epsilon\star} \text{ star-1} \qquad\qquad s_1 = \text{“”} \quad s_2 = \text{“b”}$$

## Regular Expressions

$$\frac{}{\text{``''} \sim \epsilon}\,\text{epsilon} \qquad \frac{c \in \Sigma}{\text{``}c\text{''} \sim c}\,\text{char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \mathrel{++} s_2 = s_3}{s_3 \sim r_1 r_2}\,\text{seq}$$

$$\frac{s \sim r_1}{s \sim r_1 | r_2}\,\text{alt-l} \quad \frac{s \sim r_2}{s \sim r_1 | r_2}\,\text{alt-r} \quad \frac{}{\text{``''} \sim r\star}\,\text{star-0} \quad \frac{s_1 \sim r \quad s_2 \sim r\star \quad s_1 \mathrel{++} s_2 = s_3}{s_3 \sim r\star}\,\text{star-1}$$

## "b" $\sim \epsilon\star$?

$$\frac{}{\text{``''} \sim \epsilon}\,\text{epsilon} \qquad \frac{\text{``b''} \sim \epsilon\star \qquad \text{``''} \mathrel{++} \text{``b''} = \text{``b''}}{\text{``b''} \sim \epsilon\star}\,\text{star-1} \qquad\qquad s_1 = \text{``''} \quad s_2 = \text{``b''}$$

## Regular Expressions

$$\frac{}{\text{``''} \sim \epsilon}\text{ epsilon} \qquad \frac{c \in \Sigma}{\text{``}c\text{''} \sim c}\text{ char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \mathbin{++} s_2 = s_3}{s_3 \sim r_1 r_2}\text{ seq}$$

$$\frac{s \sim r_1}{s \sim r_1 | r_2}\text{ alt-l} \qquad \frac{s \sim r_2}{s \sim r_1 | r_2}\text{ alt-r} \qquad \frac{}{\text{``''} \sim r\star}\text{ star-0} \qquad \frac{s_1 \sim r \quad s_2 \sim r\star \quad s_1 \mathbin{++} s_2 = s_3}{s_3 \sim r\star}\text{ star-1}$$

## "b" $\sim \epsilon\star$?

$$\frac{\dfrac{}{\text{``''} \sim \epsilon}\text{ epsilon} \qquad \dfrac{\vdots}{\text{``b''} \sim \epsilon\star}\text{ star-1} \qquad \text{``''} \mathbin{++} \text{``b''} = \text{``b''}}{\text{``b''} \sim \epsilon\star}\text{ star-1} \qquad\qquad s_1 = \text{``''} \quad s_2 = \text{``b''}$$

### Regular Expressions

$$\frac{}{\text{""} \sim \epsilon}\text{epsilon} \qquad \frac{c \in \Sigma}{\text{``c''} \sim c}\text{char} \qquad \frac{s_1 \sim r_1 \quad s_2 \sim r_2 \quad s_1 \mathbin{++} s_2 = s_3}{s_3 \sim r_1 r_2}\text{seq}$$

$$\frac{s \sim r_1}{s \sim r_1 | r_2}\text{alt-l} \quad \frac{s \sim r_2}{s \sim r_1 | r_2}\text{alt-r} \quad \frac{}{\text{""} \sim r\star}\text{star-0} \quad \frac{s_1 \sim r \quad s_2 \sim r\star \quad s_1 \mathbin{++} s_2 = s_3 \quad s_1 \neq \text{""}}{s_3 \sim r\star}\text{star-1}$$

```
(~~)                :: String -> RE -> Bool
""   ~~ Epsilon   = True                         -- epsilon
[c1] ~~ Ch c2     = c1 == c2                      -- char
s3   ~~ Seq r1 r2 = or $ zipWith testSplit (inits s3) (tails s3)
                    where testSplit s1 s2 = s1 ~~ r1 && s2 ~~ r2
s    ~~ Alt r1 r2 = s ~~ r1 || s ~~ r2           -- alt-l and alt-r
""   ~~ Star _    = True                          -- star-0
s3   ~~ Star r    = or $ zipWith testSplit (inits s3) (tails s3)
              where testSplit [] _ = False
                    testSplit s1 s2 = s1 ~~ r && s2 ~~ (Star r)
_    ~~ _         = False                          -- default
```
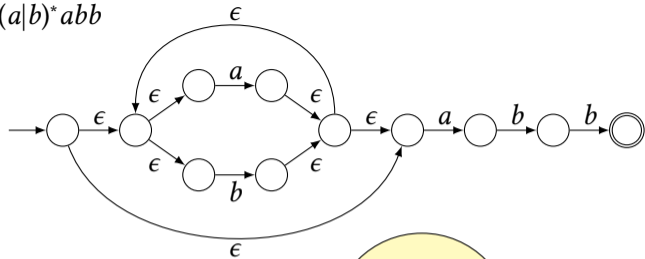
This is a backtracking algorithm that tries everything until it works
What does this do on "aaaaaaabb" $\sim$ (aaa)*b*?

```
(~~)                :: String -> RE -> Bool
""   ~~ Epsilon   = True                      -- epsilon
[c1] ~~ Ch c2     = c1 == c2                   -- char
s3   ~~ Seq r1 r2 = or $ zipWith testSplit (inits s3) (tails s3)
                    where testSplit s1 s2 = s1 ~~ r1 && s2 ~~ r2
s    ~~ Alt r1 r2 = s ~~ r1 || s ~~ r2         -- alt-l and alt-r
""   ~~ Star _    = True                       -- star-0
s3   ~~ Star r    = or $ zipWith testSplit (inits s3) (tails s3)
            where testSplit [] _ = False
                  testSplit s1 s2 = s1 ~~ r && s2 ~~ (Star r)
_    ~~ _         = False                      -- default
```

# A Better Way: Thompson's Algorithm



Ken Thompson. Programming techniques: Regular expression search algorithm.
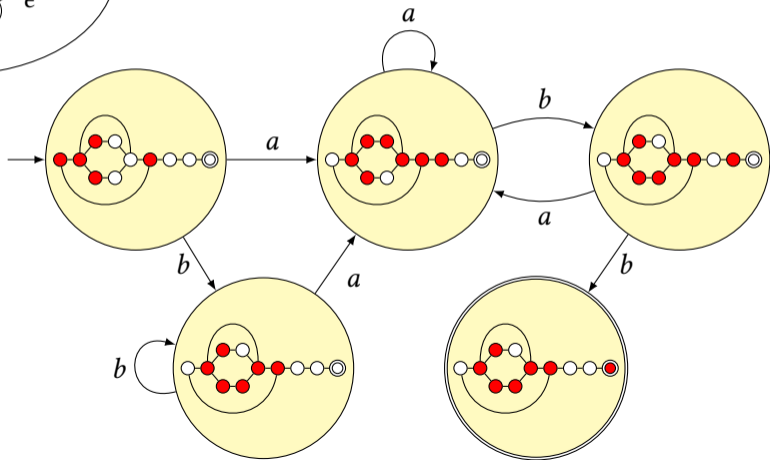*Communications of the ACM*, 11(6):419–422, June 1968.
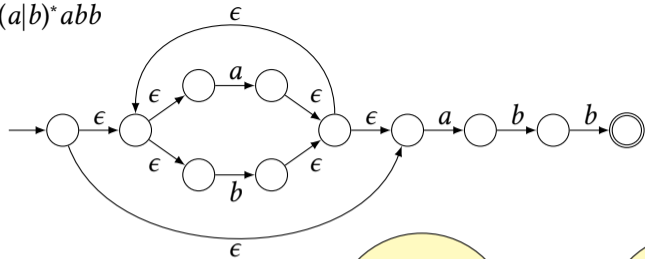
$(a|b)^* abb$

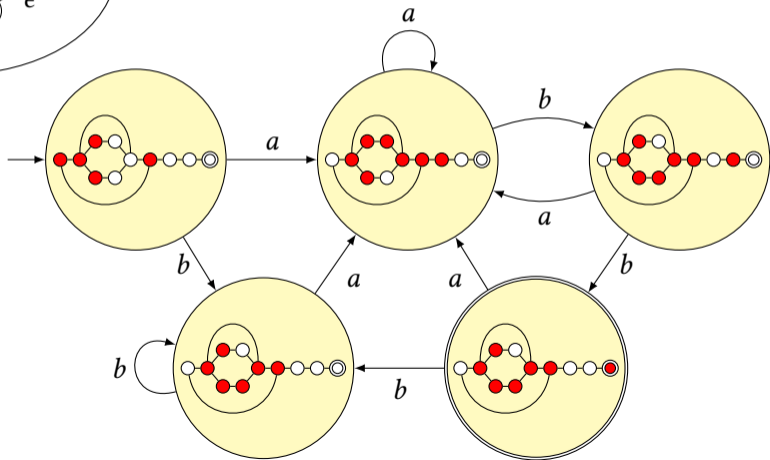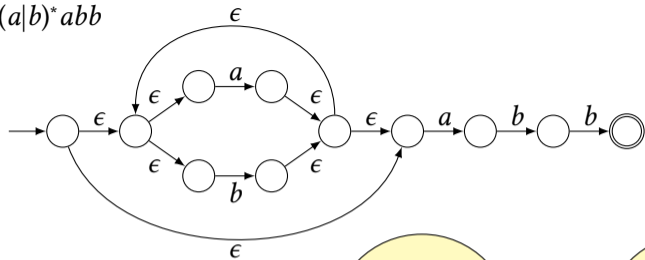$(a|b)^*abb$

$(a|b)^*abb$

$(a|b)^*abb$

$(a|b)^*abb$

$(a|b)^*abb$

## Brzozowski derivatives

$\partial_s R$ is the derivative of regular expression $R$ w.r.t. the string $s$

"Every string that can follow $s$ to match $R$"

$\partial_{s_1} R = \{ s_2 \mid s_1 s_2 \in L(R) \}$, where $L(R)$ is the language of $R$

Janusz A. Brzozowski. Derivatives of regular expressions.
*Journal of the Association for Computing Machinery*, 11(4):481–494, October 1964.

## Brzozowski derivatives

$\partial_s R$ is the derivative of regular expression $R$ w.r.t. the string $s$
"Every string that can follow $s$ to match $R$"
$\partial_{s_1} R = \{ s_2 \mid s_1 s_2 \in L(R) \}$, where $L(R)$ is the language of $R$

$\partial_a\, a = \epsilon \quad \partial_a\, aa = a \quad \partial_a\, abc = bc \quad \partial_b\, abc = \varnothing \quad \partial_a\, ab\,|\,cd = b$
$\partial_a\, abc\,|\,acd = bc\,|\,cd \quad \partial_a\, a{\star}bc = a{\star}bc \quad \partial_a\, a{\star}ac = a{\star}ac\,|\,c$

Janusz A. Brzozowski. Derivatives of regular expressions.
*Journal of the Association for Computing Machinery*, 11(4):481–494, October 1964.

## Brzozowski derivatives

$\partial_s R$ is the derivative of regular expression $R$ w.r.t. the string $s$

"Every string that can follow $s$ to match $R$"

$\partial_{s_1} R = \{ s_2 \mid s_1 s_2 \in L(R) \}$, where $L(R)$ is the language of $R$

---

$\partial_a\, a = \epsilon \quad \partial_a\, aa = a \quad \partial_a\, abc = bc \quad \partial_b\, abc = \varnothing \quad \partial_a\, ab\,|\,cd = b$

$\partial_a\, abc\,|\,acd = bc\,|\,cd \quad \partial_a\, a{*}bc = a{*}bc \quad \partial_a\, a{*}ac = a{*}ac\,|\,c$

Theorem: the derivative of a regular expression is a regular expression (including $\varnothing$)

Some subtlety when "" $\sim R$, rules otherwise look like those for polynomials

Janusz A. Brzozowski. Derivatives of regular expressions.
*Journal of the Association for Computing Machinery*, 11(4):481–494, October 1964.

### Brzozowski derivatives

$\partial_s R$ is the derivative of regular expression $R$ w.r.t. the string $s$

"Every string that can follow $s$ to match $R$"

$\partial_{s_1} R = \{s_2 \mid s_1 s_2 \in L(R)\}$, where $L(R)$ is the language of $R$

$\partial_a a = \epsilon$   $\partial_a aa = a$   $\partial_a abc = bc$   $\partial_b abc = \emptyset$   $\partial_a ab|cd = b$
$\partial_a abc|acd = bc|cd$   $\partial_a a*bc = a*bc$   $\partial_a a*ac = a*ac|c$

Theorem: the derivative of a regular expression is a regular expression (including $\emptyset$)

Some subtlety when "" $\sim R$, rules otherwise look like those for polynomials

Use "subset construction" to build a DFA: label states with regular expression derivatives

Janusz A. Brzozowski. Derivatives of regular expressions.
*Journal of the Association for Computing Machinery*, 11(4):481–494, October 1964.

Scott Owens, John Reppy, and Aaron Turon. Regular-expression derivatives re-examined.
*Journal of Functional Programming*, 19(2):173–190, March 2009.