# Automata and Formal Languages (2)

| | |
|---|---|
| Email: | christian.urban at kcl.ac.uk |
| Office: | S1.27 (1st floor Strand Building) |
| Slides: | KEATS |

# Languages

A language is a set of strings.

A regular expression specifies a set of strings, or language.

# Strings

Different ways of writing strings:

$$"hello" \qquad [h, e, l, l, o] \qquad h :: e :: l :: l :: o :: Nil$$

$$"" \qquad\qquad [] \qquad\qquad\qquad Nil$$

# Strings

Different ways of writing strings:

$$"hello" \qquad [h, e, l, l, o] \qquad h :: e :: l :: l :: o :: Nil$$

$$"" \qquad\qquad [] \qquad\qquad\qquad Nil$$

The concatenation operation on strings and sets of strings:

$$"foo" \ @ \ "bar" = "foobar"$$

$$A \ @ \ B \stackrel{\text{def}}{=} \{s_1 @ s_2 \mid s_1 \in A \land s_2 \in B\}$$

# Regular Expressions

Their inductive definition:

$$
\begin{array}{llll}
r & ::= & \varnothing & \text{null} \\
& | & \epsilon & \text{empty string / '''' / []} \\
& | & c & \text{character} \\
& | & r_1 \cdot r_2 & \text{sequence} \\
& | & r_1 + r_2 & \text{alternative / choice} \\
& | & r^* & \text{star (zero or more)}
\end{array}
$$

# Re...

Their indu

```scala
abstract class Rexp

case object NULL extends Rexp
case object EMPTY extends Rexp
case class CHAR(c: Char) extends Rexp
case class ALT(r1: Rexp, r2: Rexp) extends Rexp
case class SEQ(r1: Rexp, r2: Rexp) extends Rexp
case class STAR(r: Rexp) extends Rexp
```

$$
\begin{array}{llll}
r & ::= & \varnothing & \text{null} \\
  & | & \epsilon & \text{empty string / ''''' / []} \\
  & | & c & \text{character} \\
  & | & r_1 \cdot r_2 & \text{sequence} \\
  & | & r_1 + r_2 & \text{alternative / choice} \\
  & | & r^* & \text{star (zero or more)}
\end{array}
$$

# The Meaning of a Regular Expression

$$L(\varnothing) \stackrel{\text{def}}{=} \varnothing$$

$$L(\epsilon) \stackrel{\text{def}}{=} \{""\}$$

$$L(c) \stackrel{\text{def}}{=} \{"c"\}$$

$$L(r_1 + r_2) \stackrel{\text{def}}{=} L(r_1) \cup L(r_2)$$

$$L(r_1 \cdot r_2) \stackrel{\text{def}}{=} L(r_1) \, @ \, L(r_2)$$

$$L(r^*) \stackrel{\text{def}}{=} \bigcup_{n \geq 0} L(r)^n$$

$L$ is a function from regular expressions to sets of strings

$L : \text{Rexp} \Rightarrow \text{Set}[\text{String}]$

# The Meaning of a Regular Expression

$$L(\varnothing) \overset{\text{def}}{=} \varnothing$$
$$L(\epsilon) \overset{\text{def}}{=} \{""\}$$
$$L(c) \overset{\text{def}}{=} \{"c"\}$$
$$L(r_1 + r_2) \overset{\text{def}}{=} L(r_1) \cup L(r_2)$$
$$L(r_1 \cdot r_2) \overset{\text{def}}{=} L(r_1) @ L(r_2)$$
$$L(r^*) \overset{\text{def}}{=} \bigcup_{n \geq 0} L(r)^n$$

$$L(r)^0 \overset{\text{def}}{=} \{""\}$$
$$L(r)^{n+1} \overset{\text{def}}{=} L(r) @ L(r)^n$$

$L$ is a function from regular expressions to sets of strings

$L : \text{Rexp} \Rightarrow \text{Set}[\text{String}]$

What is $L(a^*)$?

# Reg Exp Equivalences

$$(a + b) + c \equiv^? a + (b + c)$$

$$a + a \equiv^? a$$

$$(a \cdot b) \cdot c \equiv^? a \cdot (b \cdot c)$$

$$a \cdot a \equiv^? a$$

$$\epsilon^* \equiv^? \epsilon$$

$$\varnothing^* \equiv^? \varnothing$$

$$\forall \text{ r.} \qquad r \cdot \epsilon \equiv^? r$$

$$\forall \text{ r.} \qquad r + \epsilon \equiv^? r$$

$$\forall \text{ r.} \qquad r + \varnothing \equiv^? r$$

$$\forall \text{ r.} \qquad r \cdot \varnothing \equiv^? r$$

$$c \cdot (a + b) \equiv^? (c \cdot a) + (c \cdot b)$$

$$a^* \equiv^? \epsilon + (a \cdot a^*)$$

# Reg Exp Equivalences

$$(a + b) + c \;\equiv^?\; a + (b + c) \qquad \text{yes}$$

$$a + a \;\equiv^?\; a$$

$$(a \cdot b) \cdot c \;\equiv^?\; a \cdot (b \cdot c)$$

$$a \cdot a \;\equiv^?\; a$$

$$\epsilon^* \;\equiv^?\; \epsilon$$

$$\varnothing^* \;\equiv^?\; \varnothing$$

$\forall$ r. $\qquad r \cdot \epsilon \;\equiv^?\; r$

$\forall$ r. $\qquad r + \epsilon \;\equiv^?\; r$

$\forall$ r. $\qquad r + \varnothing \;\equiv^?\; r$

$\forall$ r. $\qquad r \cdot \varnothing \;\equiv^?\; r$

$$c \cdot (a + b) \;\equiv^?\; (c \cdot a) + (c \cdot b)$$

$$a^* \;\equiv^?\; \epsilon + (a \cdot a^*)$$

# Reg Exp Equivalences

$$(a + b) + c \equiv^? a + (b + c) \qquad \text{yes}$$

$$a + a \equiv^? a \qquad \text{yes}$$

$$(a \cdot b) \cdot c \equiv^? a \cdot (b \cdot c)$$

$$a \cdot a \equiv^? a$$

$$\epsilon^* \equiv^? \epsilon$$

$$\varnothing^* \equiv^? \varnothing$$

$\forall$ r. $\qquad r \cdot \epsilon \equiv^? r$

$\forall$ r. $\qquad r + \epsilon \equiv^? r$

$\forall$ r. $\qquad r + \varnothing \equiv^? r$

$\forall$ r. $\qquad r \cdot \varnothing \equiv^? r$

$$c \cdot (a + b) \equiv^? (c \cdot a) + (c \cdot b)$$

$$a^* \equiv^? \epsilon + (a \cdot a^*)$$

# Reg Exp Equivalences

$$(a + b) + c \equiv^? a + (b + c) \quad \text{yes}$$

$$a + a \equiv^? a \quad \text{yes}$$

$$(a \cdot b) \cdot c \equiv^? a \cdot (b \cdot c) \quad \text{yes}$$

$$a \cdot a \equiv^? a$$

$$\epsilon^* \equiv^? \epsilon$$

$$\varnothing^* \equiv^? \varnothing$$

$$\forall \, r. \quad r \cdot \epsilon \equiv^? r$$

$$\forall \, r. \quad r + \epsilon \equiv^? r$$

$$\forall \, r. \quad r + \varnothing \equiv^? r$$

$$\forall \, r. \quad r \cdot \varnothing \equiv^? r$$

$$c \cdot (a + b) \equiv^? (c \cdot a) + (c \cdot b)$$

$$a^* \equiv^? \epsilon + (a \cdot a^*)$$

# Reg Exp Equivalences

$$(a + b) + c \quad \equiv^? \quad a + (b + c) \qquad \text{yes}$$

$$a + a \quad \equiv^? \quad a \qquad \text{yes}$$

$$(a \cdot b) \cdot c \quad \equiv^? \quad a \cdot (b \cdot c) \qquad \text{yes}$$

$$a \cdot a \quad \equiv^? \quad a \qquad \text{no}$$

$$\epsilon^* \quad \equiv^? \quad \epsilon$$

$$\varnothing^* \quad \equiv^? \quad \varnothing$$

$$\forall\, r. \qquad r \cdot \epsilon \quad \equiv^? \quad r$$

$$\forall\, r. \qquad r + \epsilon \quad \equiv^? \quad r$$

$$\forall\, r. \qquad r + \varnothing \quad \equiv^? \quad r$$

$$\forall\, r. \qquad r \cdot \varnothing \quad \equiv^? \quad r$$

$$c \cdot (a + b) \quad \equiv^? \quad (c \cdot a) + (c \cdot b)$$

$$a^* \quad \equiv^? \quad \epsilon + (a \cdot a^*)$$

# Reg Exp Equivalences

$$(a + b) + c \equiv^? a + (b + c) \quad \text{yes}$$

$$a + a \equiv^? a \quad \text{yes}$$

$$(a \cdot b) \cdot c \equiv^? a \cdot (b \cdot c) \quad \text{yes}$$

$$a \cdot a \equiv^? a \quad \text{no}$$

$$\epsilon^* \equiv^? \epsilon \quad \text{yes}$$

$$\varnothing^* \equiv^? \varnothing$$

$$\forall\, r. \quad r \cdot \epsilon \equiv^? r$$

$$\forall\, r. \quad r + \epsilon \equiv^? r$$

$$\forall\, r. \quad r + \varnothing \equiv^? r$$

$$\forall\, r. \quad r \cdot \varnothing \equiv^? r$$

$$c \cdot (a + b) \equiv^? (c \cdot a) + (c \cdot b)$$

$$a^* \equiv^? \epsilon + (a \cdot a^*)$$

# Reg Exp Equivalences

$$(a + b) + c \equiv^? a + (b + c) \qquad \text{yes}$$

$$a + a \equiv^? a \qquad \text{yes}$$

$$(a \cdot b) \cdot c \equiv^? a \cdot (b \cdot c) \qquad \text{yes}$$

$$a \cdot a \equiv^? a \qquad \text{no}$$

$$\epsilon^* \equiv^? \epsilon \qquad \text{yes}$$

$$\varnothing^* \equiv^? \varnothing \qquad \text{no}$$

$$\forall\, r. \qquad r \cdot \epsilon \equiv^? r$$

$$\forall\, r. \qquad r + \epsilon \equiv^? r$$

$$\forall\, r. \qquad r + \varnothing \equiv^? r$$

$$\forall\, r. \qquad r \cdot \varnothing \equiv^? r$$

$$c \cdot (a + b) \equiv^? (c \cdot a) + (c \cdot b)$$

$$a^* \equiv^? \epsilon + (a \cdot a^*)$$

# Reg Exp Equivalences

$$(a + b) + c \equiv^? a + (b + c) \qquad \text{yes}$$

$$a + a \equiv^? a \qquad \text{yes}$$

$$(a \cdot b) \cdot c \equiv^? a \cdot (b \cdot c) \qquad \text{yes}$$

$$a \cdot a \equiv^? a \qquad \text{no}$$

$$\epsilon^* \equiv^? \epsilon \qquad \text{yes}$$

$$\varnothing^* \equiv^? \varnothing \qquad \text{no}$$

$$\forall\, r. \qquad r \cdot \epsilon \equiv^? r \qquad \text{yes}$$

$$\forall\, r. \qquad r + \epsilon \equiv^? r$$

$$\forall\, r. \qquad r + \varnothing \equiv^? r$$

$$\forall\, r. \qquad r \cdot \varnothing \equiv^? r$$

$$c \cdot (a + b) \equiv^? (c \cdot a) + (c \cdot b)$$

$$a^* \equiv^? \epsilon + (a \cdot a^*)$$

# Reg Exp Equivalences

$$(a + b) + c \equiv^? a + (b + c) \quad \text{yes}$$

$$a + a \equiv^? a \quad \text{yes}$$

$$(a \cdot b) \cdot c \equiv^? a \cdot (b \cdot c) \quad \text{yes}$$

$$a \cdot a \equiv^? a \quad \text{no}$$

$$\epsilon^* \equiv^? \epsilon \quad \text{yes}$$

$$\varnothing^* \equiv^? \varnothing \quad \text{no}$$

$$\forall\, r. \quad r \cdot \epsilon \equiv^? r \quad \text{yes}$$

$$\forall\, r. \quad r + \epsilon \equiv^? r \quad \text{no}$$

$$\forall\, r. \quad r + \varnothing \equiv^? r$$

$$\forall\, r. \quad r \cdot \varnothing \equiv^? r$$

$$c \cdot (a + b) \equiv^? (c \cdot a) + (c \cdot b)$$

$$a^* \equiv^? \epsilon + (a \cdot a^*)$$

# Reg Exp Equivalences

$$(a + b) + c \equiv^? a + (b + c) \quad \text{yes}$$
$$a + a \equiv^? a \quad \text{yes}$$
$$(a \cdot b) \cdot c \equiv^? a \cdot (b \cdot c) \quad \text{yes}$$
$$a \cdot a \equiv^? a \quad \text{no}$$
$$\epsilon^* \equiv^? \epsilon \quad \text{yes}$$
$$\varnothing^* \equiv^? \varnothing \quad \text{no}$$
$$\forall\, r. \quad r \cdot \epsilon \equiv^? r \quad \text{yes}$$
$$\forall\, r. \quad r + \epsilon \equiv^? r \quad \text{no}$$
$$\forall\, r. \quad r + \varnothing \equiv^? r \quad \text{yes}$$
$$\forall\, r. \quad r \cdot \varnothing \equiv^? r$$
$$c \cdot (a + b) \equiv^? (c \cdot a) + (c \cdot b)$$
$$a^* \equiv^? \epsilon + (a \cdot a^*)$$

# Reg Exp Equivalences

$$(a + b) + c \equiv^? a + (b + c) \quad \text{yes}$$

$$a + a \equiv^? a \quad \text{yes}$$

$$(a \cdot b) \cdot c \equiv^? a \cdot (b \cdot c) \quad \text{yes}$$

$$a \cdot a \equiv^? a \quad \text{no}$$

$$\epsilon^* \equiv^? \epsilon \quad \text{yes}$$

$$\varnothing^* \equiv^? \varnothing \quad \text{no}$$

$$\forall \, r. \quad r \cdot \epsilon \equiv^? r \quad \text{yes}$$

$$\forall \, r. \quad r + \epsilon \equiv^? r \quad \text{no}$$

$$\forall \, r. \quad r + \varnothing \equiv^? r \quad \text{yes}$$

$$\forall \, r. \quad r \cdot \varnothing \equiv^? r \quad \text{no}$$

$$c \cdot (a + b) \equiv^? (c \cdot a) + (c \cdot b)$$

$$a^* \equiv^? \epsilon + (a \cdot a^*)$$

# Reg Exp Equivalences

$$(a + b) + c \equiv^? a + (b + c) \qquad \text{yes}$$

$$a + a \equiv^? a \qquad \text{yes}$$

$$(a \cdot b) \cdot c \equiv^? a \cdot (b \cdot c) \qquad \text{yes}$$

$$a \cdot a \equiv^? a \qquad \text{no}$$

$$\epsilon^* \equiv^? \epsilon \qquad \text{yes}$$

$$\varnothing^* \equiv^? \varnothing \qquad \text{no}$$

$$\forall \, r. \qquad r \cdot \epsilon \equiv^? r \qquad \text{yes}$$

$$\forall \, r. \qquad r + \epsilon \equiv^? r \qquad \text{no}$$

$$\forall \, r. \qquad r + \varnothing \equiv^? r \qquad \text{yes}$$

$$\forall \, r. \qquad r \cdot \varnothing \equiv^? r \qquad \text{no}$$

$$c \cdot (a + b) \equiv^? (c \cdot a) + (c \cdot b) \qquad \text{yes}$$

$$a^* \equiv^? \epsilon + (a \cdot a^*)$$

# Reg Exp Equivalences

$$(a + b) + c \equiv^? a + (b + c) \quad \text{yes}$$

$$a + a \equiv^? a \quad \text{yes}$$

$$(a \cdot b) \cdot c \equiv^? a \cdot (b \cdot c) \quad \text{yes}$$

$$a \cdot a \equiv^? a \quad \text{no}$$

$$\epsilon^* \equiv^? \epsilon \quad \text{yes}$$

$$\varnothing^* \equiv^? \varnothing \quad \text{no}$$

$$\forall\, r. \quad r \cdot \epsilon \equiv^? r \quad \text{yes}$$

$$\forall\, r. \quad r + \epsilon \equiv^? r \quad \text{no}$$

$$\forall\, r. \quad r + \varnothing \equiv^? r \quad \text{yes}$$

$$\forall\, r. \quad r \cdot \varnothing \equiv^? r \quad \text{no}$$

$$c \cdot (a + b) \equiv^? (c \cdot a) + (c \cdot b) \quad \text{yes}$$

$$a^* \equiv^? \epsilon + (a \cdot a^*) \quad \text{yes}$$

# The Specification for Matching

a regular expression $r$ matches a string $s$ if and only if

$$s \in L(r)$$

# The Specification for Matching

a regular expression $r$ matches a string $s$ if and only if

$$s \in L(r)$$

# $(a?\{n\}) \cdot a\{n\}$

secs — as

- ○ Python
- ⬠ Ruby

# Evil Regular Expressions

- **R**egular **e**xpression **D**enial **o**f **S**ervice (ReDoS)

- Evil regular expressions
  - $(a?\{n\}) \cdot a\{n\}$
  - $(a^+)^+$
  - $([a\text{-}z]^+)^*$
  - $(a + a \cdot a)^+$
  - $(a + a?)^+$

# A Matching Algorithm

...whether a regular expression can match the empty string:

$$nullable(\varnothing) \stackrel{\text{def}}{=} false$$

$$nullable(\epsilon) \stackrel{\text{def}}{=} true$$

$$nullable(c) \stackrel{\text{def}}{=} false$$

$$nullable(r_1 + r_2) \stackrel{\text{def}}{=} nullable(r_1) \vee nullable(r_2)$$

$$nullable(r_1 \cdot r_2) \stackrel{\text{def}}{=} nullable(r_1) \wedge nullable(r_2)$$

$$nullable(r^*) \stackrel{\text{def}}{=} true$$

# A Matching Algorithm

...whether a regular expression can match the empty string:

$$nullable(\varnothing) \quad \stackrel{\text{def}}{=} \quad false$$

$$nullable(\epsilon) \quad \stackrel{\text{def}}{=} \quad true$$

$$nullable(c) \quad \stackrel{\text{def}}{=} \quad false$$

$$nullable(r_1 + r_2) \quad \stackrel{\text{def}}{=} \quad nullable(r_1) \vee nullable(r_2)$$

$$nullable(r_1 \cdot r_2) \quad \stackrel{\text{def}}{=} \quad nullable(r_1) \wedge nullable(r_2)$$

$$nulla \quad \stackrel{\text{def}}{=}$$

```scala
1  def nullable (r: Rexp) : Boolean = r match {
2    case NULL => false
3    case EMPTY => true
4    case CHAR(_) => false
5    case ALT(r1, r2) => nullable(r1) || nullable(r2)
6    case SEQ(r1, r2) => nullable(r1) && nullable(r2)
7    case STAR(_) => true
8  }
```

# The Derivative of a Rexp

If $r$ matches the string $c :: s$, what is a regular expression that matches $s$?

*der c r* gives the answer

# The Derivative of a Rexp (2)

$$der\ c\ (\varnothing) \overset{\text{def}}{=} \varnothing$$

$$der\ c\ (\epsilon) \overset{\text{def}}{=} \varnothing$$

$$der\ c\ (d) \overset{\text{def}}{=} \text{if } c = d \text{ then } \epsilon \text{ else } \varnothing$$

$$der\ c\ (r_1 + r_2) \overset{\text{def}}{=} der\ c\ r_1 + der\ c\ r_2$$

$$der\ c\ (r_1 \cdot r_2) \overset{\text{def}}{=} \text{if } nullable(r_1)$$
$$\text{then } (der\ c\ r_1) \cdot r_2 + der\ c\ r_2$$
$$\text{else } (der\ c\ r_1) \cdot r_2$$

$$der\ c\ (r^*) \overset{\text{def}}{=} (der\ c\ r) \cdot (r^*)$$

# The Derivative of a Rexp (2)

$$der\ c\ (\varnothing) \quad \overset{\text{def}}{=}\ \varnothing$$

$$der\ c\ (\epsilon) \quad \overset{\text{def}}{=}\ \varnothing$$

$$der\ c\ (d) \quad \overset{\text{def}}{=}\ \text{if}\ c = d\ \text{then}\ \epsilon\ \text{else}\ \varnothing$$

$$der\ c\ (r_1 + r_2) \quad \overset{\text{def}}{=}\ der\ c\ r_1 + der\ c\ r_2$$

$$der\ c\ (r_1 \cdot r_2) \quad \overset{\text{def}}{=}\ \text{if}\ nullable(r_1)$$
$$\text{then}\ (der\ c\ r_1) \cdot r_2 + der\ c\ r_2$$
$$\text{else}\ (der\ c\ r_1) \cdot r_2$$

$$der\ c\ (r^*) \quad \overset{\text{def}}{=}\ (der\ c\ r) \cdot (r^*)$$

$$ders\ [\,]\ r \quad \overset{\text{def}}{=}\ r$$

$$ders\ (c :: s)\ r \quad \overset{\text{def}}{=}\ ders\ s\ (der\ c\ r)$$

# The Derivative of a Rexp (2)

$$der\, c\, (\varnothing) \qquad \stackrel{\text{def}}{=} \varnothing$$

$$der\, c\, (\epsilon) \qquad \stackrel{\text{def}}{=} \varnothing$$

```scala
def der (r: Rexp, c: Char) : Rexp = r match {
  case NULL => NULL
  case EMPTY => NULL
  case CHAR(d) => if (c == d) EMPTY else NULL
  case ALT(r1, r2) => ALT(der(r1, c), der(r2, c))
  case SEQ(r1, r2) =>
    if (nullable(r1)) ALT(SEQ(der(r1, c), r2), der(r2, c))
    else SEQ(der(r1, c), r2)
  case STAR(r) => SEQ(der(r, c), STAR(r))
}

def ders (s: List[Char], r: Rexp) : Rexp = s match {
  case Nil => r
  case c::s => ders(s, der(c, r))
}
```

# Examples

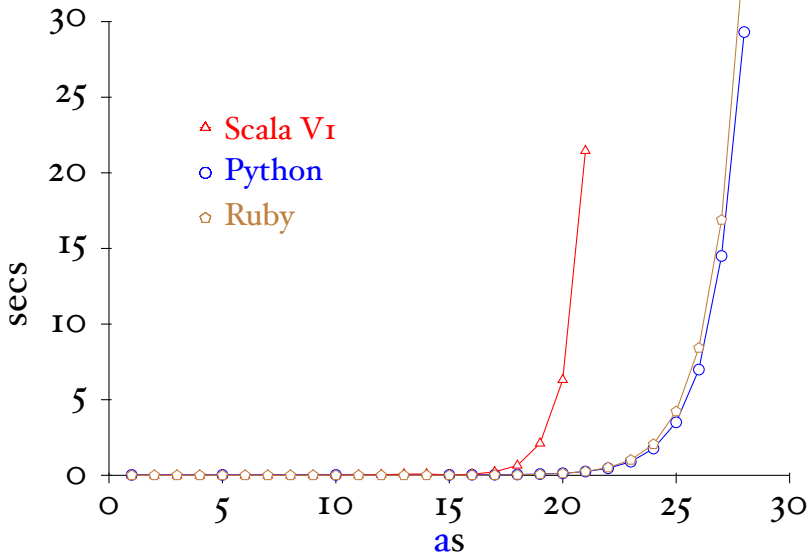Given $r \stackrel{\text{def}}{=} ((a \cdot b) + b)^*$ what is

$$der\ a\ r$$
$$der\ b\ r$$

$$(a?\{n\}) \cdot a\{n\}$$

# Proofs about Rexps

Remember their inductive definition:

$$
\begin{aligned}
r \quad ::= \quad & \varnothing \\
| \quad & \epsilon \\
| \quad & c \\
| \quad & r_1 \cdot r_2 \\
| \quad & r_1 + r_2 \\
| \quad & r^*
\end{aligned}
$$

If we want to prove something, say a property $P(r)$, for all regular expressions $r$ then ...

# Proofs about Rexp (2)

- $P$ holds for $\varnothing$, $\epsilon$ and $c$

- $P$ holds for $r_1 + r_2$ under the assumption that $P$ already holds for $r_1$ and $r_2$.

- $P$ holds for $r_1 \cdot r_2$ under the assumption that $P$ already holds for $r_1$ and $r_2$.

- $P$ holds for $r^*$ under the assumption that $P$ already holds for $r$.

# Proofs about Rexp (3)

Assume $P(r)$ is the property:

$$nullable(r) \text{ if and only if } "" \in L(r)$$

# Proofs about Rexp (4)

Let $Der\ c\ A$ be the set defined as

$$Der\ c\ A \stackrel{\text{def}}{=} \{s \mid c{::}s \in A\}$$

We can prove

$$L(der\ c\ r) = Der\ c\ (L(r))$$

by induction on $r$.

# Proofs about Strings

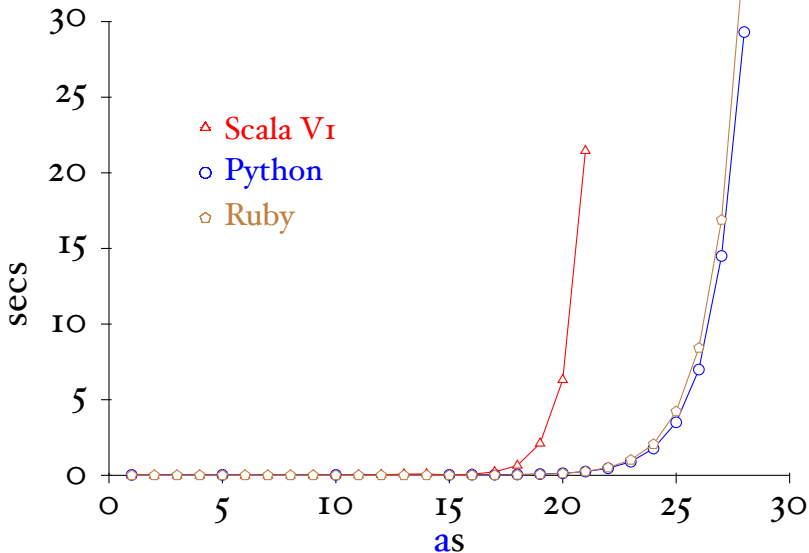If we want to prove something, say a property $P(s)$, for all strings $s$ then ...

- $P$ holds for the empty string, and

- $P$ holds for the string $c :: s$ under the assumption that $P$ already holds for $s$

# Proofs about Strings (2)

We can finally prove

$$matcher(r, s) \text{ if and only if } s \in L(r)$$

# $(a?\{n\}) \cdot a\{n\}$

# A Problem

We represented the "n-times" $a\{n\}$ as a sequence regular expression:

  1:    $a$

  2:    $a \cdot a$

  3:    $a \cdot a \cdot a$

        ...

 13:    $a \cdot a \cdot a \cdot a \cdot a \cdot a \cdot a \cdot a \cdot a \cdot a \cdot a \cdot a \cdot a$

        ...

20:

This problem is aggravated with $a?$ being represented as $\epsilon + a$.
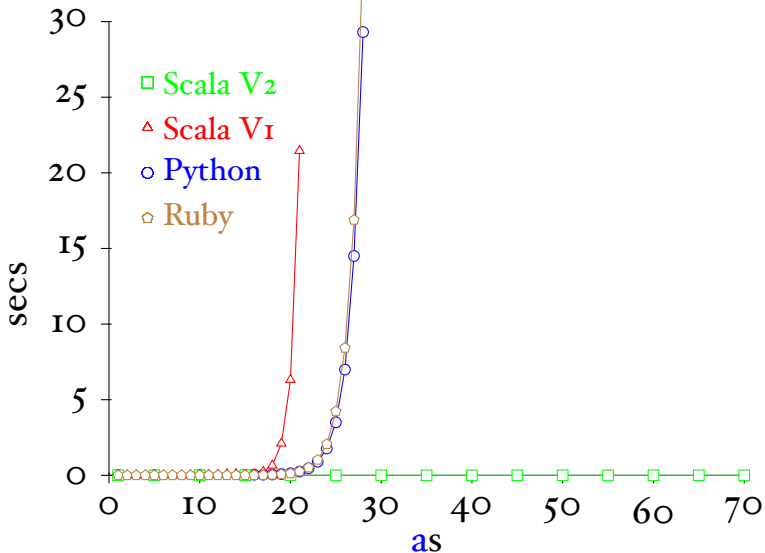
# Solving the Problem
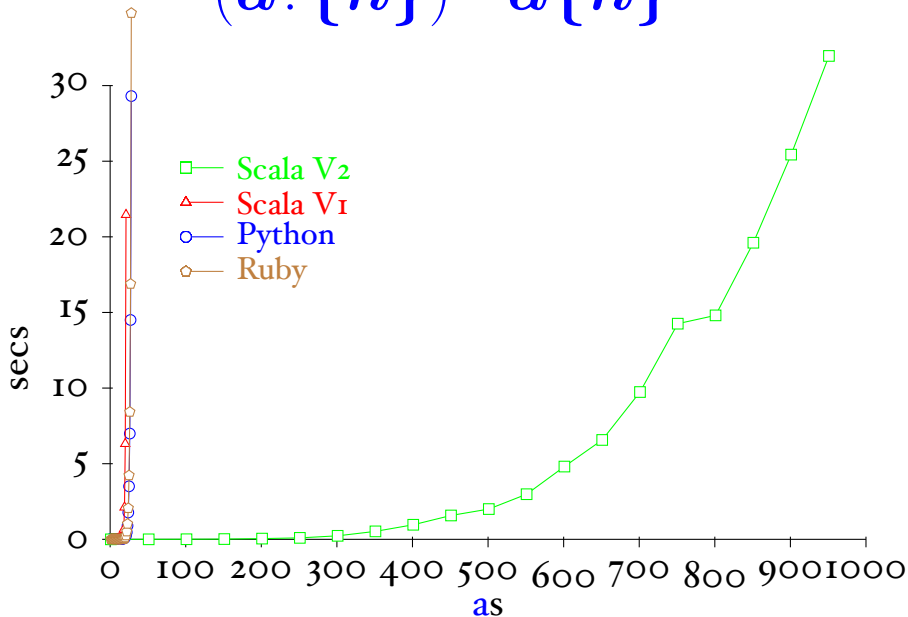
What happens if we extend our regular expressions

$$r \quad ::= \quad ...$$
$$| \quad r\{n\}$$
$$| \quad r?$$

What is their meaning? What are the cases for *nullable* and *der*?

# $(a?\{n\}) \cdot a\{n\}$

# $(a?\{n\}) \cdot a\{n\}$

# Examples

Recall the example of $r \stackrel{\text{def}}{=} ((a \cdot b) + b)^*$ with

$$der\ a\ r = ((\epsilon \cdot b) + \varnothing) \cdot r$$
$$der\ b\ r = ((\varnothing \cdot b) + \epsilon) \cdot r$$

What are these regular expressions equal to?

$(a?\{n\}) \cdot a\{n\}$