# Discrete Mathematics, Lecture 27.

2024-02-07

## Abstract Syntax Trees

# **Concrete syntax:** strings of symbols

▶ possibly including symbols to disambiguate the semantics (brackets, white space, *etc*),

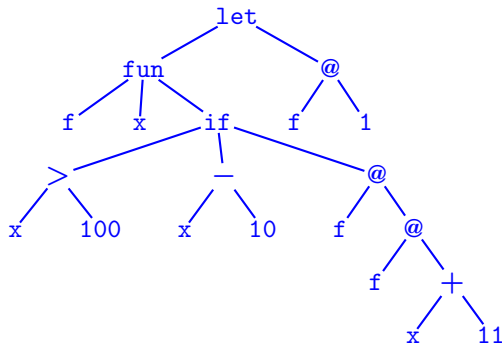▶ or that have no semantic content (*e.g.* syntax for comments).

For example, an ML expression:

```
let fun f x =
    if x > 100 then x − 10
    else f ( f ( x + 11 ) )
in f 1 end
(* value   is   99 *)
```

# Abstract syntax: finite rooted trees

- ▶ vertexes with **n** children are labelled by operators expecting **n** arguments (**n**-**ary** operators) – in particular leaves are labelled with **0**-ary (nullary) operators (constants, variables, *etc*)

- ▶ label of the root gives the 'outermost form' of the whole phrase

E.g. for the ML expression on Slide 25:

# Regular expressions (concrete syntax)

over a given alphabet $\Sigma$.

Let $\Sigma'$ be the 6-element set $\{\epsilon, \varnothing, |, *, (, )\}$ (assumed disjoint from $\Sigma$)

$$U = (\Sigma \cup \Sigma')^*$$

axioms: $\dfrac{\overset{(a \in \Sigma)}{\rule{2cm}{0.4pt}}}{a} \qquad \dfrac{\rule{1.5cm}{0.4pt}}{\epsilon} \qquad \dfrac{\rule{1.5cm}{0.4pt}}{\varnothing}$

rules: $\dfrac{r}{(r)} \qquad \dfrac{r \quad s}{r|s} \qquad \dfrac{r \quad s}{rs} \qquad \dfrac{r}{r^*}$

(where $a \in \Sigma$ and $r, s \in U$)

# Some derivations of regular expressions
## (assuming $a, b \in \Sigma$)

$$\cfrac{\widetilde{\epsilon} \quad \cfrac{\widehat{a} \quad \cfrac{\overline{b}}{b^*}}{ab^*}}{\epsilon | ab^*}$$

$$\cfrac{\cfrac{\widehat{\epsilon} \quad \widehat{a}}{\epsilon | a} \quad \cfrac{\overline{b}}{b^*}}{\epsilon | ab^*}$$

$$\cfrac{\widehat{\epsilon} \quad \cfrac{\cfrac{\overline{a} \quad \overline{b}}{ab}}{ab^*}}{\epsilon | ab^*}$$

$$\cfrac{\widetilde{\epsilon} \quad \cfrac{\widetilde{a} \quad \cfrac{\cfrac{\overline{b}}{b^*}}{(b^*)}}{a(b^*)}}{\epsilon | (a(b^*))}$$

$$\cfrac{\cfrac{\cfrac{\widehat{\epsilon} \quad \widehat{a}}{\epsilon | a}}{(\epsilon | a)} \quad \cfrac{\cfrac{\overline{b}}{b^*}}{(b^*)}}{(\epsilon | a)(b^*)}$$

$$\cfrac{\widetilde{\epsilon} \quad \cfrac{\cfrac{\cfrac{\cfrac{\overline{a} \quad \overline{b}}{ab}}{(ab)}}{(ab)^*}}{((ab)^*)}}{\epsilon | ((ab)^*)}$$

28

# Regular expressions (abstract syntax)

The 'signature' for regular expression abstract syntax trees (over an alphabet $\Sigma$) consists of

- ▶ binary operators *Union* and *Concat*
- ▶ unary operator *Star*
- ▶ nullary operators (constants) *Null*, *Empty* and *Sym$_a$* (one for each $a \in \Sigma$).

$(ab) \mid c^*$    Union (Concat (Sym$_a$, Sym$_b$), Star(Sym$_c$))

# Regular expressions (abstract syntax)

The 'signature' for regular expression abstract syntax trees (over an alphabet $\Sigma$) as an ML datatype declaration:

$$
\begin{aligned}
\texttt{datatype } {'a}\,\texttt{RE} = {} & \texttt{Union of } ({'a}\,\texttt{RE}) * ({'a}\,\texttt{RE}) \\
\mid {} & \texttt{Concat of } ({'a}\,\texttt{RE}) * ({'a}\,\texttt{RE}) \\
\mid {} & \texttt{Star of } {'a}\,\texttt{RE} \\
\mid {} & \texttt{Null} \\
\mid {} & \texttt{Empty} \\
\mid {} & \texttt{Sym of } {'a}
\end{aligned}
$$

(the type $'a\,\texttt{RE}$ is parameterised by a type variable $'a$ standing for the alphabet $\Sigma$)

Some abstract syntax trees of regular expressions
(assuming $a, b \in \Sigma$)

$\varepsilon \mid a(b^*)$



1.    2.    3.

(*cf.* examples on Slide 28)

We will use a textual representation of trees, for example:

1. $Union(Null, Concat(Sym_a, Star(Sym_b)))$

2. $Concat(Union(Null, Sym_a), Star(Sym_b))$

3. $Union(Null, Star(Concat(Sym_a, Sym_b)))$

# Relating concrete and abstract syntax

for regular expressions over an alphabet $\Sigma$, via an inductively defined relation $\sim$ between strings and trees:
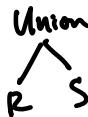
$$\overline{a \sim Sym_a} \qquad \overline{\epsilon \sim Null} \qquad \overline{\emptyset \sim Empty}$$

$$\frac{r \sim R}{(r) \sim R} \qquad \frac{r \sim R \qquad s \sim S}{r|s \sim Union(R, S)}$$

Union

R    S

$$\frac{r \sim R \qquad s \sim S}{rs \sim Concat(R, S)} \qquad \frac{r \sim R}{r^* \sim Star(R)}$$

For example:

$$\epsilon | (a(b^*)) \sim Union(Null, Concat(Sym_a, Star(Sym_b)))$$
$$\epsilon | ab^* \sim Union(Null, Concat(Sym_a, Star(Sym_b)))$$
$$\epsilon | ab^* \sim Concat(Union(Null, Sym_a), Star(Sym_b))$$

Thus $\sim$ is a 'many-many' relation between strings and trees.

- **Parsing:** algorithms for producing abstract syntax trees $parse(r)$ from concrete syntax $r$, satisfying $r \sim parse(r)$.
- **Pretty printing:** algorithms for producing concrete syntax $pp(R)$ from abstract syntax trees $R$, satisfying $pp(R) \sim R$.

(See CST IB Compiler construction course.)

We can introduce **operator precedence** and **associativity** conventions for concrete syntax and cut down the pairs in the $\sim$ relation to make it single-valued (that is, $r \sim R$ & $r \sim R' \Rightarrow R = R'$). For example, for regular expressions we decree:

$\_^*$ binds more tightly than $\_\_$, binds more tightly than $\_|\_$

So, for example, the only parse of $\epsilon|ab^*$ is the tree

$$Union(Null, Concat(Sym_a, Star(Sym_b)))$$

Concat ( Unon ( ε, Syn a ), Star (Syn b))

We also decree that the binary operations of concatenation $\_\_$ and union $\_|\_$ are left associative, so that for example $abc$ parses as

$$Concat(Concat(Sym_a, Sym_b), Sym_c)$$

(However, the union and concatenation operators for regular expressions will always be given a semantics that is associative, so the left-associativity convention is less important than the operator-precedence convention.)

**Note: for the rest of the course we adopt these operator-precedence and associativity conventions for regular expressions and refer to abstract syntax trees using concrete syntax.**

# Matching

Each regular expression $r$ over an alphabet $\Sigma$ determines a language $L(r) \subseteq \Sigma^*$. The strings $u$ in $L(r)$ are by definition the ones that **match** $r$, where

▶ $u$ matches the regular expression $a$ (where $a \in \Sigma$) iff $u = a$ *(Sigma)*

▶ $u$ matches the regular expression $\epsilon$ iff $u$ is the null string $\varepsilon$

▶ no string matches the regular expression $\emptyset$

▶ $u$ matches $r|s$ iff it either matches $r$ or it matches $s$ *(Union(r,s))*

▶ $u$ matches $rs$ iff it can be expressed as the concatenation of two strings, $u = vw$, with $v$ matching $r$ and $w$ matching $s$ *(Concat(r,s))*

▶ $u$ matches $r^*$ iff either $u = \varepsilon$, or $u$ matches $r$, or $u$ can be expressed as the concatenation of two or more strings, each of which matches $r$.

Concat(Sym$_a$, Union(Sym$_b$, Sym$_c$))

ab
ac

(ab)*
/|\
ε,   ab,   abab,   ababab, ....

# Inductive definition of matching

$$U = \Sigma^* \times \{\text{regular expressions over } \Sigma\}$$

abstract syntax trees

axioms:
$$\frac{}{(a, a)} \quad \frac{}{(\varepsilon, \epsilon)} \quad \frac{}{(\varepsilon, r^*)}$$

rules:  (a, Sym a)   (ε, Null)   (ε, Star (r))

$$\frac{(u, r)}{(u, r|s)} \qquad \frac{(u, s)}{(u, r|s)}$$

$$\frac{(v, r) \quad (w, s)}{(vw, rs)} \qquad \frac{(u, r) \quad (v, r^*)}{(uv, r^*)}$$

Concat(r, s)

(No axiom/rule involves the empty regular expression $\varnothing$ – why?)

# Examples of matching

Assuming $\Sigma = \{a, b\}$, then:

- $a|b$ is matched by each symbol in $\Sigma$

- $b(a|b)^*$ is matched by any string in $\Sigma^*$ that starts with a '$b$'

- $((a|b)(a|b))^*$ is matched by any string of even length in $\Sigma^*$

- $(a|b)^*(a|b)^*$ is matched by any string in $\Sigma^*$

- $(\epsilon|a)(\epsilon|b)|bb$ is matched by just the strings $\varepsilon$, $a$, $b$, $ab$, and $bb$

- $\varnothing b|a$ is just matched by $a$

# Some questions

(a) Is there an algorithm which, given a string $u$ and a regular expression $r$, computes whether or not $u$ matches $r$?

(b) In formulating the definition of regular expressions, have we missed out some practically useful notions of pattern? $\subset Sym_a = \varepsilon \mid b(a|b)^* \mid a(a|b)(a|b)^*$

(c) Is there an algorithm which, given two regular expressions $r$ and $s$, computes whether or not they are **equivalent**, in the sense that $L(r)$ and $L(s)$ are equal sets?

(d) Is every language (subset of $\Sigma^*$) of the form $L(r)$ for some $r$?

(a) The answer to question (a) on Slide 38 is 'yes': algorithms for deciding such pattern-matching questions make use of finite automata. We will see this next.

(b) If you have used the UNIX utility `grep`, or a text editor with good facilities for regular expression based search, like `emacs`, you will know that the answer to question (b) on Slide 38 is also 'yes'—the regular expressions defined on Slide 27 leave out some forms of pattern that one sees in such applications. However, the answer to the question is also 'no', in the sense that (for a fixed alphabet) these extra forms of regular expression are definable, up to equivalence, from the basic forms given on Slide 27. For example, if the symbols of the alphabet are ordered in some standard way, it is common to provide a form of pattern for naming ranges of symbols—for example [a-z] might denote a pattern matching any lower-case letter. It is not hard to see how to define a regular expression (albeit a rather long one) which achieves the same effect. However, some other commonly occurring kinds of pattern are much harder to describe using the rather minimalist syntax of Slide 27. The principal example is **complementation**, $\sim(r)$:

$u$ matches $\sim(r)$     iff     $u$ *does not* match $r$.

It will be a corollary of the work we do on finite automata (and a good measure of its power) that every pattern making use of the complementation operation $\sim(-)$ can be replaced by an equivalent regular expression just making use of the operations on Slide 27. But why do we stick to the minimalist syntax of regular expressions on that slide? The answer is that it reduces the amount of work we will have to do to show that, in principle, matching strings against patterns can be decided via the use of finite automata.