

## 15-312 Lecture on Concrete and Abstract Syntax

### Solution to Exercise 1

In a deductive system containing the rules for  $\_ \text{nat}$  and  $\text{sum}(\_, \_, \_)$ :

$$\frac{}{z \text{ nat}} \text{z\_nat} \qquad \frac{n \text{ nat}}{s \text{ nat}} \text{s\_nat}$$

$$\frac{n \text{ nat}}{\text{sum}(z, n, n)} \text{sum\_z} \qquad \frac{\text{sum}(m, n, p)}{\text{sum}(s \text{ } m, n, s \text{ } p)} \text{sum\_s}$$

prove that:

*If  $\mathcal{D} :: \text{sum}(m, n, p)$ , then there exists derivations  $\mathcal{D}_m :: m \text{ nat}$ ,  $\mathcal{D}_n :: n \text{ nat}$  and  $\mathcal{D}_p :: p \text{ nat}$ ,*

**Proof:** The proof proceeds by induction on the structure of the given derivation  $\mathcal{D}$ . There are two cases to examine:

1. Case:

$$\mathcal{D} = \frac{\mathcal{D}'_n \quad n' \text{ nat}}{\text{sum}(z, n', n')} \text{sum\_z}$$

Then, it must be the case that  $m = z$ ,  $n = n'$  and  $p = n'$ . We need to build derivations of  $m \text{ nat}$ ,  $n \text{ nat}$  and  $p \text{ nat}$ . These are easily obtained as follows:

$$\mathcal{D}_m = \frac{}{z \text{ nat}} \text{z\_nat} \qquad \mathcal{D}_n = \mathcal{D}'_n \qquad \mathcal{D}_p = \mathcal{D}'_n$$

This concludes this case of the proof.

2. Case:

$$\mathcal{D} = \frac{\mathcal{D}' \quad \text{sum}(m', n', p')}{\text{sum}(s \text{ } m', n', s \text{ } p')} \text{sum\_s}$$

Then, it must be the case that  $m = s m'$ ,  $n = n'$  and  $p = s p'$ . Since we have a derivation  $\mathcal{D}'$  of  $\text{sum}(m', n', p')$ , by induction hypothesis on  $\mathcal{D}'$ , we can assume that there exist derivations  $\mathcal{D}'_m :: m' \text{ nat}$ ,  $\mathcal{D}'_n :: n' \text{ nat}$  and  $\mathcal{D}'_p :: p' \text{ nat}$ .

We need to build derivations of  $m \text{ nat}$ ,  $n \text{ nat}$  and  $p \text{ nat}$ . We obtain them as follows:

$$\mathcal{D}_m = \frac{\mathcal{D}'_m}{\frac{m' \text{ nat}}{s m' \text{ nat}} \text{ s\_nat}} \quad \mathcal{D}_n = \mathcal{D}'_n \quad \mathcal{D}_p = \frac{\mathcal{D}'_p}{\frac{p' \text{ nat}}{s p' \text{ nat}} \text{ s\_nat}}$$

This concludes this case of the proof.

This concludes the proofs since we have proved each case.  $\square$

## Solution to Exercise 2

In the following deductive system for transition sequences:

$$\frac{}{s \mapsto^* s} \text{ id} \quad \frac{s \mapsto s' \quad s' \mapsto^* s''}{s \mapsto^* s''} \text{ it\_r}$$

show that the following rule is admissible:

$$\frac{s \mapsto^* s' \quad s' \mapsto s''}{s \mapsto^* s''} \text{ it\_l}$$

Is it derivable?

Let us rephrase this question as a property. Prove that

*For every derivations  $\mathcal{D} :: s_1 \mapsto^* s_2$  and  $\mathcal{E} :: s_2 \mapsto s_3$ , there exists a derivation  $\mathcal{F} :: s_1 \mapsto^* s_3$*

(We have renamed the states for clarity.)

**Proof:** The proof proceeds by induction on the structure of derivation  $\mathcal{D}$ . (Note that we have not given any rule for the judgment  $s \mapsto s'$  and therefore we have no idea what a derivation of  $\mathcal{E}$  looks like. For this reason, we cannot proceed by induction on  $\mathcal{E}$ ). There are two cases to examine, one for each rule defining the judgment  $s \mapsto^* s'$ :

1. Case:

$$\mathcal{D} = \frac{}{s \mapsto^* s} \text{ id}$$

where  $s_1 = s_2 = s$ .

Since the derivation  $\mathcal{E} :: s \mapsto s_3$  is given to us (recall that  $s_2 = s$ ), we build the desired derivation  $\mathcal{F}$  as follows:

$$\mathcal{F} = \frac{\mathcal{E} \quad \frac{}{s_3 \mapsto^* s_3} \text{ id}}{s \mapsto^* s_3} \text{ it\_r}$$

Note how we are using  $\mathcal{E}$  on the left of the construction for  $\mathcal{F}$  while it would appear on the right in the admissible rule. Note also how we needed to reconstruct from scratch the right subderivation of rule **it\_r**.

2. Case:

$$\mathcal{D} = \frac{\mathcal{D}' \quad \mathcal{D}''}{s \mapsto s' \quad s' \mapsto^* s''} \text{it\_r}$$

where  $s_1 = s$ ,  $s_2 = s''$  and  $s'$  is some intermediate state that must exist for this rule to have been applied.

Since  $\mathcal{D}''$  is smaller than  $\mathcal{D}$ , we can apply the induction hypothesis to  $\mathcal{D}''$  and  $\mathcal{E}$  and this allows us to postulate the existence of a derivation  $\mathcal{F}'' :: s' \mapsto^* s_3$ . We can then combine derivations  $\mathcal{D}'$  and  $\mathcal{F}''$  using rule **it\_r** into the desired derivation  $\mathcal{F}$  as follows:

$$\mathcal{F} = \frac{\mathcal{D}' \quad \mathcal{F}''}{s \mapsto s' \quad s' \mapsto^* s_3} \text{it\_r}$$

which is what we wanted since  $s_1 = s$ .

This concludes the proof of this statement, and this also allows us to conclude that rule

$$\frac{s \mapsto^* s' \quad s' \mapsto s''}{s \mapsto^* s''} \text{it\_l}$$

is admissible. □

To understand why rule **it\_l** is not derivable, note the shape of a proof of  $s \mapsto s'$  using rules **id** and **it\_r**: it is a tree that is completely unbalanced to the right with all the left branches consisting of subderivations of the judgment  $\_ \mapsto \_$ . So, every derivable rule may have premises for  $\_ \mapsto \_$  on the left but the rightmost premise will have to be for the  $\_ \mapsto^* \_$  judgment.

Instead, a derivation that uses **it\_l** has a subderivation for the  $\_ \mapsto \_$  judgment as its rightmost premise. For this reason, it cannot be derived from rules **id** and **it\_r**. Note that it builds a derivation that grows left. Indeed, if only **id** and **it\_l** are ever used, the resulting derivation tree will be completely unbalanced toward the left.

To reinforce these concepts, prove that the following deductive system is equivalent to :

$$\frac{}{s \mapsto^* s} \text{id}' \quad \frac{s \mapsto s'}{s \mapsto^* s'} \text{it}'_{\text{step}} \quad \frac{s \mapsto^* s' \quad s' \mapsto^* s''}{s \mapsto^* s''} \text{it}'_{\text{lr}}$$

(Note that both premises of **it'\_lr** use the transition sequence — not step — judgment.)

## Strings from First Principles

Given an alphabet  $\Sigma$ , we define the judgment  $c \text{ char}_\Sigma$  to indicate that  $c$  is a character in  $\Sigma$ . It is extensionally defined by giving one rule for each  $c \in \Sigma$ :

$$\frac{}{c \text{ char}_\Sigma} c$$

Then, strings over  $\Sigma$  are just lists of characters with  $\epsilon$  as the empty string (corresponding to the empty list of characters) and the constructor  $_ \cdot _$ , where  $c \cdot s$  indicates the extension of string  $s$  with character  $c$  (on the left):

$$\frac{}{\epsilon \text{ string}_\Sigma} \text{str}_\epsilon \qquad \frac{c \text{ char}_\Sigma \quad s \text{ string}_\Sigma}{c \cdot s \text{ string}_\Sigma} \text{str}_\cdot$$

From now on, we will assume that the alphabet  $\Sigma$  is fixed, and omit it as a subscript of these judgments.

Traversing strings left to right one character at a time is rather tedious. We'd like to manipulate strings as a concatenation of substrings. Let's define the judgment  $s \text{ as } s_1 \hat{\ } s_2$  that splits a string  $s$  as the concatenation of two string  $s_1$  and  $s_2$  (or dually builds  $s$  as the concatenation of  $s_1$  and  $s_2$ ):

$$\frac{s \text{ string}}{s \text{ as } \epsilon \hat{\ } s} \hat{\ }_r \qquad \frac{s \text{ as } s_1 \hat{\ } s_2}{c \cdot s \text{ as } (c \cdot s_1) \hat{\ } s_2} \hat{\ }_l$$

Note that the mode of these rules is  $(\exists!, \forall, \forall)$ , thus implementing a concatenation function. However, if we think about them as splitting a given string, they are non-deterministic because a string can be split in many ways (it also has mode  $(\forall, \exists, \exists)$ ).

## Lexing

The following grammar provides the raw concrete syntax for simple arithmetic expressions:

$$\begin{aligned} \text{Spaces} \quad \_ & ::= \epsilon \mid \_ \langle \text{space} \rangle \mid \_ \langle \text{tab} \rangle \mid \_ \langle \text{new line} \rangle \\ \text{Digits} \quad D & ::= 0 \mid 1 \mid \dots \mid 9 \\ \text{Numbers} \quad N & ::= D \mid D N \\ \text{Expressions} \quad E & ::= \_ N \_ \mid E \_ ++ \_ E \mid E \_ * \_ E \end{aligned}$$

We have written addition as  $++$  rather than  $+$  for illustration purposes.

Raw concrete syntax is what we type on a computer and write on paper, but it is contains a number of distracting details when we want to reason about a language as a mathematical object. We will eventually rely on *terms*, or *abstract syntax trees*, for this purpose.

A first step in this direction is to isolate the tokens that constitute a phrase, without bothering yet about the tree structure. We want in particular to ignore spaces, interpret numbers as numbers rather than sequences of digits, and give composite keywords (for

example ++ above) an atomic representation. *Lexing* will reduce that grammar to this one:

$$\text{Expressions } E ::= \text{num}[n] \mid E + E \mid E * E$$

Here,  $\text{num}[n]$ ,  $+$  and  $*$  are tokens, where  $n$  is a natural number (in decimal notation for simplicity). We used a different font to distinguish them from characters. We will construct a token stream as a string over these three types of tokens. For convenience, we will build them left-to-right rather than right-to-left as we did earlier — this is done in the same way and we still use  $\epsilon$  and  $\cdot$  as our constructors. We will also use  $\hat{\phantom{x}}$  for token stream concatenation (or splitting).

We will now implement lexing as a series of judgments, that model the way an actual lexer works, although in a very simplified and specialized way. We can do so in many ways. We will use two judgments:

$$\begin{array}{ll} t \triangleright s \text{ lex\_to } t' & \text{Given token prefix } t, \text{ string } s \text{ lexes to tokens } t' \\ t \triangleright_n s \text{ lex\_to } t' & \text{Given token prefix } t \text{ and number } n, \text{ string } s \text{ lexes to tokens } t' \end{array}$$

The first argument of these judgments is used as an accumulator where tokens are stored up to the point where the input string has been totally consumed. It is then output as the third argument. We build a token stream for a concrete string  $s$  by attempting to find a derivation for the judgment  $\epsilon \triangleright s \text{ lex\_to } t$ : if such a derivation exists, then the object  $t$  built along the way is a token stream for  $s$ . Otherwise, no token stream for  $s$  exists.

Although they are mutually recursive, we will now give the defining rules one judgment at a time:

$$\begin{array}{c} \frac{}{t \triangleright \epsilon \text{ lex\_to } t} \text{lex\_}\epsilon \qquad \frac{t \triangleright_d s \text{ lex\_to } t'}{t \triangleright d \cdot s \text{ lex\_to } t'} \text{lex\_}d \\ \frac{t \triangleright s \text{ lex\_to } t'}{t \triangleright \_ \cdot s \text{ lex\_to } t'} \text{lex\_}\_ \qquad \frac{t \cdot + \triangleright s \text{ lex\_to } t'}{t \triangleright + \cdot \_ \cdot s \text{ lex\_to } t'} \text{lex\_}+ \qquad \frac{t \cdot * \triangleright s \text{ lex\_to } t'}{t \triangleright * \cdot s \text{ lex\_to } t'} \text{lex\_}* \end{array}$$

Rule  $\text{lex\_}\epsilon$  returns the accumulated token stream when the input string is empty. When seeing a digit  $d$ , rule  $\text{lex\_}d$  passes it to the second judgment (below), which will convert it into a number once all of its digits have been collected. Rule  $\text{lex\_}\_$  ignores a space. Rules  $\text{lex\_}+$  and  $\text{lex\_}*$  recognize the strings  $++$  and  $*$  respectively and extend the token stream with the tokens  $+$  and  $*$  respectively.

$$\begin{array}{c} \frac{t \triangleright_{n'} s \text{ lex\_to } t' \quad n' = 10 \times n + d}{t \triangleright_n d \cdot s \text{ lex\_to } t'} \text{lex\_}dd \\ \frac{t \cdot \text{num}[n] \triangleright \epsilon \text{ lex\_to } t'}{t \triangleright_n \epsilon \text{ lex\_to } t'} \text{lex\_}d\epsilon \qquad \frac{t \cdot \text{num}[n] \triangleright s \text{ lex\_to } t'}{t \triangleright_n \_ \cdot s \text{ lex\_to } t'} \text{lex\_}d\_ \\ \frac{t \cdot \text{num}[n] \triangleright + \cdot s \text{ lex\_to } t'}{t \triangleright_n + \cdot s \text{ lex\_to } t'} \text{lex\_}d+ \qquad \frac{t \cdot \text{num}[n] \triangleright * \cdot s \text{ lex\_to } t'}{t \triangleright_n * \cdot s \text{ lex\_to } t'} \text{lex\_}d+ \end{array}$$

The one interesting rule here is `lex_dd`, which recognizes an additional digit of a numeric string. Defining rules for the judgment in the right premise is rather simple but quite tedious. All the other rules simply call the main judgment after adding the token `num[n]` to the token stream.

## Exercise

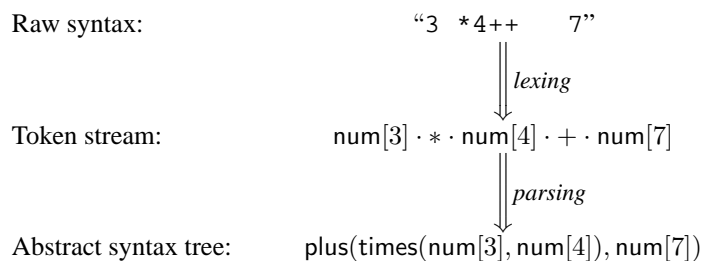
Define a syntax for regular expressions  $r$  and give the rules for a judgment  $s$  matches  $r$  which is derivable only when concrete string  $s$  matches regular expression  $r$ . It is easy to define this judgment non-deterministically by using string concatenation to split strings at appropriate places. You may want to model a more efficient version after some of the material you saw in 15-212 (a more general judgment is needed).

Assume you are given token forms  $t \in T$  (for example all the tokens that can appear in a Java program) and rules implementing a judgment  $s \triangleright_r t$  that builds the token  $t$  out of the string  $s$  which matches regular expression  $r$  (assume you already know that  $s$  matches  $r$ ; in the example above, one such judgment would build the token `num[n]` given a string of digits  $s$  and the regular expression  $(0|\dots|9)^+$  and another would build `+` when it sees `++`), generalize the above judgment  $t \triangleright s \text{ lex\_to } t'$  to produce a stream of these generic tokens.

## Parsing

Moving from raw strings to token streams hides a substantial amount of irrelevant details for the purpose of examining a program as a mathematical object. Parsing exposes the structure of this flat sequence of tokens as defined by the grammar of the language. In particular, it turns tokens that are meant to connect other tokens (such as `*` in `num[3]·*·num[4]`) into operators applied to arguments (here `times(num[3], num[4])`). What we obtain as a result of parsing is an *abstract syntax tree*, often simply called a *term*. Programs written in this way are said to be in *abstract syntax*, as opposed to the *concrete syntax* of what is entered in a computer.

Consider the following example, starting from raw strings:



(the raw string has been compacted for clarity, and the trailing  $\epsilon$  of the token stream has been dropped.)

The abstract syntax is defined inductively in a way that is coherent with the interpretation of a program or expression (recall that raw strings are also defined inductively,

but this definition has nothing to do with the higher-level meaning of a string, rather with the sequencing of characters). This inductive definition comes directly from the grammar for its language. For example, the grammar for (tokenized) expressions given earlier can be automatically transformed into the following inference rules for the judgment  $t \text{ exp}$ , which interprets a token stream  $t$  as an expression:

$$\frac{}{\text{num}[n] \text{ exp}} \text{ num}[n] \qquad \frac{t_1 \text{ exp} \quad t_2 \text{ exp}}{t_1 \hat{+} t_2 \text{ exp}} \text{ plus} \qquad \frac{t_1 \text{ exp} \quad t_2 \text{ exp}}{t_1 \hat{*} t_2 \text{ exp}} \text{ times}$$

(here, we have taken some liberties: a fully formal definition would rely on the string splitting judgment  $s$  as  $s_1 \hat{+} s_2$  defined earlier, adding one premise to the two rightmost rules.)

Consider the following derivation for the judgment

$$\text{num}[3] \cdot * \cdot \text{num}[4] \cdot + \cdot \text{num}[7] \text{ exp}$$

involving the above example:

$$\frac{\frac{\frac{}{\text{num}[3] \text{ exp}} \text{ num}[3]}{\text{num}[3] \cdot * \cdot \text{num}[4] \text{ exp}} \text{ times} \quad \frac{\frac{}{\text{num}[4] \text{ exp}} \text{ num}[4]}{\text{num}[7] \text{ exp}} \text{ times}}{\text{num}[3] \cdot * \cdot \text{num}[4] \cdot + \cdot \text{num}[7] \text{ exp}} \text{ plus}$$

This object is a parse tree for this expression, given as a derivation rather than in the usual way.

Now, the abstract syntax for this expression can be read off directly from this derivation by simply taking the names of the rules as operators. Indeed, taking the terms corresponding to the premises of a rule as the argument of this rule's name viewed as an operator, we obtain:

$$\text{plus}(\text{times}(\text{num}[3], \text{num}[4]), \text{num}[7])$$

Therefore, the abstract syntax is just a simplified way of writing the parse tree for a (token) string. Here, we have justified it in terms of derivation. This technique is completely general, and scales to more complicated forms of abstract syntax.

This yields the following grammar for abstract expressions:

$$\text{Abstract Expressions:} \quad E ::= \text{num}[n] \mid \text{plus}(E, E) \mid \text{times}(E, E)$$

We will generally omit the tag  $\text{num}$ .

This technique is even more general than this: we can read any derivation as a term (although we often need a little bit more than just the rule names). This is the way proofs are entered into a theorem prover.

Harper's book (chapter 5) mentions an alternative way of building abstract syntax trees: rather than reading them off a parse tree in derivational form, they can be given an independent definition, and an additional judgment (parsing) is defined to relate strings (raw or tokenized) and abstract syntax trees. The net effect is the same, as this construction essentially defines the abstract syntax tree to be isomorphic to the derivation. This is another general approach — but sometimes error prone — to introducing terms to describe derivations.

## Ambiguity

Of course, the above grammar for expressions is ambiguous, and the shown derivation tree is not the only one for our example (the other one takes `*` to be the main operator, and therefore would build the term `times(num[3], plus(num[4], num[7]))`). The ambiguity here derives from the fact that the string concatenation judgment  $s \text{ as } s_1 \hat{=} s_2$ , implicitly used in rules `plus` and `times`, is non-deterministic when used for splitting  $s$  into  $s_1$  and  $s_2$ .

There are standard techniques for writing a grammar that ensure that it is non-ambiguous. Then, there is at most one parse tree (i.e., derivation in the sense above) in correspondence with any string in the language. See Harper's book (chapter 5) for a discussion on this.

Making a grammar non-ambiguous often involves introducing new non-terminals (which for us means adding judgments) and maybe parentheses. The construction given above clearly works also in this case, but doing so often carries over disambiguating machinery to the abstract syntax, which is not directly useful since the abstract syntax is not ambiguous (by definition). For this reason, the abstract syntax is generally based on a very simple (but usually ambiguous) grammar for a language.