Compilers and computer architecture From strings to ASTs (2): context free grammars

Martin Berger<sup>1</sup>

October 2019

<sup>1</sup>Email: M.F.Berger@sussex.ac.uk, Office hours: Wed 12-13 in Chi-2R312

## Recall the function of compilers



# Recall we are discussing parsing



Remember, we want to take a program given as a string and:

- Check if it's syntactically correct, e.g. is every opened bracket later closed?
- Produce an AST to facilitate efficient code generation.



We split that task into two phases, lexing and parsing. Lexing throws away some information (e.g. how many white-spaces) and prepares a token-list, which is used by the parser. The token-list simplifies the parser, because some detail is not important for syntactic correctness:

if x < 2+3 then P else Q

is syntactically correct exactly when

if y < 111 + 222 then P else Q

The token-list simplifies the parser, because some detail is not important for syntactic correctness:

if x < 2+3 then P else Q

is syntactically correct exactly when

if y < 111 + 222 then P else Q

So from the point of view of the next stage (parsing), all we need to know is that the input is

T\_if T\_var T\_less T\_int T\_plus T\_int T\_then ...

Of course we cannot throw away the names of variables etc completely, as the later stages (type-checking and code generation) need them. They are just irrelevant for syntax checking. We keep them and our token-lists are like this

T\_if T\_var ( "x" ) T\_less T\_int ( 2 ) T\_plus ...

#### Two tasks of syntax analysis

As with the lexical phase, we have to deal with two distinct tasks.

- Specifying that the syntactically correct programs (token lists) are.
- Checking if an input program (token list) is syntactically correct according to the specification, and output a corresponding AST.

Let's deal with specification first. What are our options? How about using regular expressions for this purpose?

Alas not every language can be expressed in these formalisms. Example:

```
Alphabet = \{'(', ')'\}.
```

Language = all **balanced** parentheses (), ()(), (()), (()), ((())), (()))), ..., note: the empty string is balanced.

## FSAs/REs can't count

Let's analyse the situation a bit more. Why can we not describe the language of all balanced parentheses using REs or FSAs.

Each FSA has only a fixed number (say *n*) of states. But what if we have more than *n* open brackets before we hit a closing bracket?

Since there are only *n* states, when we reach the *n* open bracket, we must have gone back to a state that we already visited earlier, say when we processed the *i*-th bracket with i < n. This means the automaton treats *i* as it does *n*, leading to confusion.

Summary: FSAs can't count, and likewise for REs (why?).

Lack of expressivity of regular expressions & FSAs

Why is it a problem for syntax analysis in programming languages if REs and FSAs can't count?

Because programming languages contain many bracket-like constructs that can be nested, e.g.

```
begin ... end
do ... while
if ( ... ) then { ... } else { ... }
3 + ( 3 - (x + 6) )
```

But we must formalise the syntax of our language if we want to computer to process it. So we need a formalism that can 'count'.

#### Problem

# What we are looking for is something like REs, but more powerful:

regular expression/FSA		???
lexer	_	parser

Let me introduce you to: context free grammars (CFGs).

## Context free grammars

Programs have a naturally recursive and nested structure: A program is e.g.:

▶ if P then Q else Q', where P, Q, Q' are programs.

$$\blacktriangleright$$
 x := *P*, where *P* is a program.

begin x := 1; begin ... end; y := 2; end

CFGs are a generalisation of regular expression that is ideal for describing such recursive and nested structures.

#### Context free grammar

#### A context-free grammar is a tuple (A, V, Init, R) where

- A is a finite set called **alphabet**.
- ► *V* is a finite, non-empty set of **variables**.
- $\blacktriangleright A \cap V = \emptyset.$
- Init  $\in$  V is the initial variable.
- *R* is the finite set of **reductions**, where each reduction in *R* is of the form (*I*, *r*) such that
  - $\blacktriangleright$  *I* is a variable, i.e.  $I \in V$ .
  - ▶ *r* is a string (possibly empty) over the **new** alphabet  $A \cup V$ .

We usually write  $I \rightarrow r$  for  $(I, r) \in R$ .

Note that the alphabet are often also called **terminal symbols**, reductions are also called **reduction steps** or **transitions** or **productions**, some people say **non-terminal symbol** for variable, and the initial variable is also called **start symbol**.

# Context free grammar

Example:

$$\blacktriangleright A = \{a, b\}.$$

► 
$$V = \{S\}.$$

- ► The initial variable is *S*.
- R contains only three reductions:

Recall that  $\epsilon$  is the empty string.

Now the CFG is (A, V, S, R).

The language of balanced brackets with *a* being the open bracket, and *b* being the closed bracket!

To make this intuition precise, we need to say precisely what the language of a CFG is.

## The language accepted by a CFG

# The key idea is simple: **replace the variables according to the reductions**.

Given a string *s* over  $A \cup V$ , i.e. the alphabet and variables, any occurrence of a variable *T* in *s* can be replaced by the string  $r_1...r_n$ , provided there is a reduction  $T \rightarrow r_1...r_n$ .

For example if we have a reduction

$$S 
ightarrow a \ T \ b$$

then we can rewrite the string

aaSbb

to

How do we start this rewriting of variables? With the initial variable.

When does this rewriting of variables stop? When the string we arrive at by rewriting in a finite number of steps from the initial variable contains no more variables.

## The language accepted by a CFG

Then: the **language** of a CFG is the set of all strings over the alphabet of the CFG that can be arrived at by rewriting from the initial variable.

## The language accepted by a CFG

Let's do this with the CFG for balanced brackets (A, V, S, R) where

• 
$$A = \{(,)\}.$$

$$\blacktriangleright V = \{S\}.$$

► The initial variable is *S*.

▶ Reductions *R* are *S* → (*S*), *S* → *SS*, and *S* →  $\epsilon$ 

$$S \rightarrow (S)$$
  

$$\rightarrow (SS)$$
  

$$\rightarrow ((S)S)$$
  

$$\rightarrow (((S))S)$$
  

$$\rightarrow (((S))SS)$$
  

$$\rightarrow (((S))\epsilon S) = (((S))S)$$
  

$$\rightarrow (((\epsilon))S) = (((())S)$$
  

$$\rightarrow ((())\epsilon) = ((()))$$

Question: Why / how can CFGs count?

Why / how does the CFG (A, V, S, R) with

count?

Because only  $S \to (S)$  introduces new brackets. But by construction it always introduces a closing bracket for each new open bracket.

## The language accepted by a CFG: infinite reductions

Note that many CFGs allow infinite reductions: for example with the grammar the previous slide we can do this:

:

Such infinite reductions don't affect the language of the grammar. Only sequences of rewrites that end in a string free from variables count towards the language.

#### The language accepted by a CFG If you like formal definitions ...

Given a fixed CFG  $\mathcal{G} = (A, V, S, R)$ . For arbitrary strings  $\sigma, \sigma' \in (V \cup A)^*$  we define the *one-step reduction relation*  $\Rightarrow$  which relates strings from  $(V \cup A)^*$  as follows.  $\sigma \Rightarrow \sigma'$  if and only if:

•  $\sigma = \sigma_1 I \sigma_2$  where  $I \in V$ , and  $\sigma_1, \sigma_2$  are strings from  $(V \cup A)^*$ .

• There is a reduction  $I \longrightarrow \gamma$  in R.

$$\blacktriangleright \sigma' = \sigma_1 \gamma \sigma_2.$$

The *language accepted by*  $\mathcal{G}$ , written  $lang(\mathcal{G})$  is given as follows.

$$lang(\mathcal{G}) \stackrel{\text{def}}{=} \{ \gamma_n \mid S \to \gamma_1 \to \cdots \to \gamma_n, \text{ where } \gamma_n \in A^* \}$$

The sequence  $S \rightarrow \gamma_1 \rightarrow \cdots \rightarrow \gamma_n$  is called **derivation**.

**Note:** only strings free from variables can be in  $lang(\mathcal{G})$ .

#### Example CFG

Consider the following CFG where while, if, ; etc are elements of the alphabet, and *M* is a variable.

М	$\rightarrow$	while $M$ do $M$
М	$\rightarrow$	if ${\it M}$ then ${\it M}$
М	$\rightarrow$	M; M
	÷	

If *M* is the starting variable, then we can derive

÷

$$egin{array}{rcl} M& o&MM\ o&M ext{if}\ M ext{ then}\ M\ o&M ext{if}\ M ext{ then}\ while\ M ext{ do}\ M \end{array}$$

We do this until we reach a string without variables.

## Some conventions regarding CFGs

Here is a collection of conventions for making CFGs more readable. You will find them a lot when programming languages are discussed.

Variables are CAPITALISED, the alphabet is lower case (or vice versa).

Variables are in **BOLD**, the alphabet is not (or vice versa).

Variables are written in  $\langle angle-brackets \rangle$ , the alphabet isn't.

#### Some conventions regarding CFGs

Instead of multiple reductions from the same variable, like

 $N 
ightarrow r_1$  $N 
ightarrow r_2$  $N 
ightarrow r_3$ 

we write

$$N \rightarrow r_1 \mid r_2 \mid r_3$$

Instead of

$$P 
ightarrow ext{if} P$$
 then  $P \mid$  while  $P$  do  $P$ 

We often write

$$P, Q 
ightarrow$$
 if  $P$  then  $Q \mid$  while  $P$  do  $Q$ 

Finally, many write ::= instead of  $\rightarrow$ .

## Simple arithmetic expressions

Let's do another example. Grammar:

$$E \rightarrow E+E \mid E * E \mid (E) \mid 0 \mid 1 \mid ...$$

The language contains:

- ▶ 7
- ▶ 7 ∗ 4
- ▶ 7 \* 4 + 222
- ▶ 7 \* (4 + 222) ...

A well-known context free grammar

#### $R \rightarrow \emptyset \mid \epsilon \mid 'c' \mid R+R \mid RR \mid R^* \mid (R)$

What's this?

(The syntax of) regular expressions can be described by a CFG (but not by an RE)!

Since regular expressions are a special case of CFGs, could we not do both, lexing and parsing, using only CFGs?

In principle yes, but lexing based on REs (and FSAs) is simpler and faster!

#### Example: Java grammar

Let's look at the CFG for a real language:

https://docs.oracle.com/javase/specs/jls/se13/ html/jls-19.html



#### CFGs, what's next?

Recall we were looking for this:

$$\frac{\text{regular expression/FSA}}{\text{FSA}} = \frac{???}{\text{parser}}$$

And the answer was CFGs. But what is a parser?

$$\frac{\text{regular expression/FSA}}{\text{FSA}} = \frac{CFG}{???}$$

### CFGs, what's next?

CFGs allow us to specify the grammar for programming languages. But that's not all we want. We also want:

- An algorithm that decided whether a given token list is in the language of the grammar or not.
- An algorithm that converts the list of tokens (if valid) into an AST.

The key idea to solving both problems in one go is the **parse** tree.

#### CFG vs AST

Here is a grammar that you were asked to write an AST for in the tutorials.

Ρ	::=	$x := e \mid if0 e then P else P$
		whileGt0 e do P   P; P
е	::=	e + e   e - e   e * e   (e)   e % e
		x   0   1   2



#### Here's a plausible definition of ASTs for the language: syntax.java

Do you notice something?

Looks very similar. The CFG is (almost?) a description of data type for the AST.

This is no coincidence, and we will use this similarity to construct the AST as we parse, in that we will see the parsing process as a tree. How?

#### Derivations and parse trees

Recall that a derivation in a CFG (A, V, I, R) is a sequence

 $I \rightarrow t_1 \rightarrow t_2 \rightarrow \cdots \rightarrow t_n$ 

where  $t_n$  is free from variables, and each step is goverend by a reduction from *R*.

We can drawn each derivation as a tree, called **parse tree**. The parse tree tells us how the input token list 'fits into' the grammar, e.g. which reduction we applied and when to 'consume' the input.

- ► The start symbol / is the tree's root.
- For each reduction X → ⟨y<sub>1</sub>,..., y<sub>n</sub>⟩ we add all the y<sub>i</sub> as children to node X.

I.e. nodes in the tree are elements of  $A \cup V$ . Let's look at an example.

Recall our CFG for arithmetic expressions: Let's do another example. Grammar:

$$E \rightarrow E+E \mid E * E \mid (E) \mid 0 \mid 1 \mid ...$$

Let's say we have the string "4 \* 3 + 17". Let's parse this string and build the corresponding parse tree.

#### Example parse tree



Let's do this in detail on the board.

## Derivations and parse trees

The following is important about parse trees.

- Terminal symbols are at the leaves of the tree.
- Variables symbols are at the non-leave nodes.
- An **in-order** traversal of the tree returns the input string.
- ► The parse tree reveals bracketing structure explicitly.

## Left- vs rightmost derivation

BUT ... usually there are many derivations for a chosen string, giving the same parse tree. For example:



Canonical choices:

- Left-most: always replace the left-most variable.
- **Right-most:** always replace the right-most variable.
- ► NB the examples above were neither left-nor right-most! In parsing we usually use either left- or rightmost derivations to construct a parse tree.

Question: do left- and rightmost derivations lead to the same parse tree?

Answer: For a context-free grammar: YES. It really doesn't matter in what order variables are rewritten. Why? Because the rest of the string is unaffected by rewriting a variable, so we can modify the order.

Alas there is a second degree of freedom: which rule to choose?

In constrast: it can make a big difference which **rule** we apply when rewriting a variable

This leads to an important subject: **ambiguity**.

# Ambiguity



#### Ambiguity In the grammar

#### $E \rightarrow E+E \mid E * E \mid (E) \mid 0 \mid 1 \mid \dots$

the string 4 \* 3 + 17 has two distinct parse trees!

#### Е Е Е Е Е Е Е Е Е 17 Е 17 3 3

A CFG with this property is called **ambiguous**.





# Ambiguity



More precisely: a context-free grammar is **ambiguous** if there is a string in the language of the grammar that has more than one parse tree.

Note that this has **nothing** to do with left- vs right-derivation. Each of the ambiguous parse trees has a left- and a right-derivation.





We also have ambiguity in natural language, e.g.



# Ambiguity



Ambiguity is programming languages is bad, because it leaves the meaning of a program unclear, e.g. the compiler should generate different code for 1 + 2 \* 3 when it's uderstood as (1 + 2) \* 3 than for 1 + (2 \* 3).

Can we automatically check whether a grammar is ambigouous?

Bad news: ambiguity of grammars is undecidable, i.e. no algorithm can exist that takes as input a CFG and returns "Ambiguous" or "Not ambiguous" correctly for all CFGs.

# Ambiguity



There are several ways to deal with ambiguity.

- Parser returns all possible parse trees, leaving choice to later compiler phases. Example: combinator parsers often do this, Earley parser. Downside: kicks can down the road ... need to disambiguate later (i.e. doesn't really solve the problem), and does too much work if some of the parse trees are later discarded.
- ▶ Use non-ambiguous grammar. Easier said than done ...
- Rewriting the grammar to remove ambiguity. For example by enforcing precedence that \* binds more tightly than +. We look at this now.



The problem with

$$E \rightarrow E+E \mid E * E \mid (E) \mid 0 \mid 1 \mid ...$$

is that addition and multiplication have the same status. But in our everyday understanding, we think of a \* b + c as meaning (a \* b) + c. Moreover, we evaluate a + b + c as (a + b) + c. But there's nothing in the naive grammar that ensures this.

Let's bake these preferences into the grammar.

Ambiguity: grammar rewriting



#### Let's rewrite

to  

$$E \rightarrow E + E \mid E * E \mid (E) \mid 0 \mid 1 \mid \dots$$

$$E \rightarrow F + E \mid F$$

$$F \rightarrow N * F \mid N \mid (E) * F \mid (E)$$

$$N \rightarrow 0 \mid 1 \mid \dots$$

Examples in class.



Here is a problem that often arises when specifying programming languages.

$$egin{array}{lll} M & o & ext{if} \ M ext{ then } M \ & & & \ & & \ & & \ & & \ & & \ & & \ & & \$$

# If-Then ambiguity



#### Now we can find two distinct parse trees for

if B then if B' then P else  ${\tt Q}$ 



#### If-Then ambiguity

We solved the \*/+ ambiguity by giving \* precedence. At the level of grammar that meant we had + coming 'first' in the grammar.

Let's do this for the if-then ambiguity by saying:

else always closes the nearest **unclosed** if, so if-then-else has priority over if-then.



# If-Then ambiguity, aka the dangling-else problem

 $\rightarrow$  if *M* then *ITE* else *ITE* 

 $\begin{array}{rrr} M & \rightarrow & ITE \\ & \mid & BOTH \end{array}$ 

ITE

only if-then-else both if-then and if-then-else

other reductions

 $\begin{array}{rcl} \textit{BOTH} & \to & \textit{if } \textit{M} \textit{ then } \textit{M} \\ & & | & \textit{if } \textit{M} \textit{ then } \textit{ITE} \textit{ else } \textit{BOTH} & \textit{no other reductions} \end{array}$ 



# Ambiguity: general algorithm?



Alas there is no algorithm that can rewrite all ambiguous CFGs into unambiguous CFGs with the same language, since some CFGs are **inherently ambiguous**, meaning they are only recongnised by ambiguous CFGs.

Fortunately, such languages are esoteric and not relevant for programming languages. For languages relevant in programming, it is generally straightforward to produce an unambiguous CFG.

I will not ask you in the exam to convert an ambiguous CFG into an unambiguous CFG. You should just know what ambiguity means in parsing and why it is a problem.