

Notes on formal grammars

Handout written by Maggie Johnson and revised by me.

What is a grammar?

A *grammar* is a powerful tool for describing and analyzing languages. It is a set of rules by which valid sentences in a language are constructed. Here's a trivial example of English grammar:

sentence	->	<subject> <verb-phrase> <object>
subject	->	This Computers I
verb-phrase	->	<adverb> <verb> <verb>
adverb	->	never
verb	->	is run am tell
object	->	the <noun> a <noun> <noun>
noun	->	university world cheese lies

Using the above rules or *productions*, we can derive simple sentences such as these:

This is a university.
Computers run the world.
I am the cheese.
I never tell lies.

Here is a *leftmost derivation* of the first sentence using these productions.

sentence	->	<subject> <verb-phrase> <object>
	->	This <verb-phrase> <object>
	->	This <verb> <object>
	->	This is <object>
	->	This is a <noun>
	->	This is a university

In addition to several reasonable sentences, we can also derive nonsense like “Computers run cheese” and “This am a lies”. These sentences doesn't make semantic sense, but they are syntactically correct because they are of the sequence of subject, verb-phrase, and object. Grammars are a tool for syntax, not semantics. We worry about semantics at a later point in the compiling process. In the syntax analysis phase, we verify the correctness of the structure only.

Vocabulary

We need to review some definitions before we can proceed:

<i>grammar</i>	a set of rules by which valid sentences in a language are constructed.
<i>nonterminal</i>	a grammar symbol that can be replaced/expanded to a sequence of symbols.
<i>terminal</i>	an actual word in a language; these are the symbols in a grammar that cannot be replaced by anything else. “terminal” is supposed to conjure up the idea that it is a dead-end— no further expansion is possible from here.
<i>production</i>	a grammar rule that describes how to replace/exchange symbols. The general form of a production for a nonterminal is:

$$X \rightarrow Y_1 Y_2 Y_3 \dots Y_n$$

The nonterminal X is declared equivalent to the concatenation of the symbols $Y_1 Y_2 Y_3 \dots Y_n$. The production means that anywhere where we encounter X , we may

replace it by the string $Y_1Y_2Y_3\dots Y_n$. Eventually we will have a string containing nothing that can be expanded further, i.e., it will consist of only terminals. Such a string is called a *sentence*. In the context of programming languages, a sentence is a syntactically correct and complete program.

derivation a sequence of applications of the rules of a grammar that produces a finished string of terminals. A *leftmost derivation* is where we always substitute for the leftmost nonterminal as we apply the rules (we can similarly define a rightmost derivation). A derivation is also called a *parse*.

start symbol a grammar has a single nonterminal (the start symbol) from which all sentences derive:

$$S \rightarrow X_1X_2X_3\dots X_n$$

All sentences are derived from S by successive replacement using the productions of the grammar.

null symbol ϵ it is sometimes useful to specify that a symbol can be replaced by nothing at all. To indicate this, we use the null symbol ϵ , e.g., $A \rightarrow B \mid \epsilon$.

BNF a way of specifying programming languages using formal grammars and production rules with a particular form of notation (Backus-Naur form).

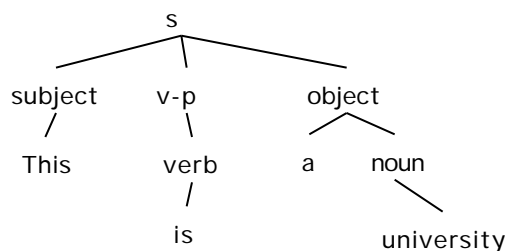
A few grammar exercises to try on your own, the alphabet in each case is $\{a,b\}$

- Define a grammar for the language of strings with one or more a's followed by zero or more b's.
- Define a grammar for even-length palindromes.
- Define a grammar for strings where the number of a's is equal to the number b's.
- Define a grammar where the number of a's is not equal to the number b's. (Hint: think about it as two separate cases...)
- Can you write regular expressions for these languages? Why or why not?

Parse representation

In working with grammars, we can represent the application of the rules to derive a sentence in two ways. The first is a derivation as shown earlier for “This is a university” where the rules are applied step-by-step and we substitute for one nonterminal at a time. Think of a derivation as a history of how the sentence was parsed because it not only includes which productions were applied, but also the order they were applied (i.e. which nonterminal was chosen for expansion at each step). There can many different derivations for the same sentence (the leftmost, the rightmost, and so on).

A *parse tree* is the second method for representation. It diagrams how each symbol derives from other symbols in a hierarchical manner. Here is a parse tree for “This is a university”:



Although the parse tree includes all of the productions that were applied, it does not encode the order they were applied. For an unambiguous grammar (we'll define ambiguity in a minute), there is exactly one parse tree for a particular sentence.

More formal definitions

Here are some other definitions we will need, described in reference to this example grammar:

$$\begin{array}{lcl} S & \rightarrow & AB \\ A & \rightarrow & Ax \mid y \\ B & \rightarrow & z \end{array}$$

alphabet

The alphabet is $\{S, A, B, x, y, z\}$. It is divided into two disjoint sets. The *terminal alphabet* consists of terminals, which appear in the sentences of the language: $\{x, y, z\}$. The remaining symbols are the *nonterminal alphabet*; these are the symbols that appear on the left side of productions and can be replaced during the course of a derivation: $\{S, A, B\}$. Formally, we use \mathbf{V} for the alphabet, \mathbf{T} for the terminal alphabet and \mathbf{N} for the nonterminal alphabet giving us: $\mathbf{V} = \mathbf{T} \cup \mathbf{N}$, and $\mathbf{T} \cap \mathbf{N} = \emptyset$.

The convention used in our lecture notes are a sans-serif font for grammar elements, lowercase for terminals, uppercase for nonterminals, and underlined lowercase (e.g., \underline{u} , \underline{v}) to denote arbitrary strings of terminal and nonterminal symbols (possibly null). In some textbooks, Greek symbols are used for arbitrary strings of terminal and nonterminal symbols (e.g., α , β).

context-free grammar

To define a language, we need a set of productions, of the general form: $\underline{u} \rightarrow \underline{v}$. In a *context-free grammar*, \underline{u} is a single nonterminal and \underline{v} is an arbitrary string of terminal and nonterminal symbols. When parsing, we can replace \underline{u} by \underline{v} wherever it occurs. We shall refer to this set of productions symbolically as \mathbf{P} .

formal grammar

We formally define a grammar as a 4-tuple $\{S, \mathbf{P}, \mathbf{N}, \mathbf{T}\}$. S is the start symbol and $S \in \mathbf{N}$, \mathbf{P} is the set of productions, and \mathbf{N} and \mathbf{T} are the nonterminal and terminal alphabets. A sentence is a string of symbols in \mathbf{T} derived from S using one or more applications of productions in \mathbf{P} . A string of symbols derived from S but possibly including nonterminals is called a *sentential form* or a *working string*.

A production $\underline{u} \rightarrow \underline{v}$ is used to replace an occurrence of \underline{u} by \underline{v} . Formally, if we apply a production $p \in \mathbf{P}$ to a string of symbols \underline{w} in \mathbf{V} to yield a new string of symbols \underline{z} in \mathbf{V} , we say that \underline{z} derived from \underline{w} using p , written as follows: $\underline{w} \Rightarrow^p \underline{z}$. We also use:

$$\begin{array}{ll} \underline{w} \Rightarrow \underline{z} & \underline{z} \text{ derives from } \underline{w} \text{ (production unspecified)} \\ \underline{w} \Rightarrow^* \underline{z} & \underline{z} \text{ derives from } \underline{w} \text{ using zero or more productions} \\ \underline{w} \Rightarrow^+ \underline{z} & \underline{z} \text{ derives from } \underline{w} \text{ using one or more productions} \end{array}$$

equivalence

The language $L(G)$ defined by grammar G is the set of sentences derivable using G . Two grammars G and G' are said to be *equivalent* if the languages they generate $L(G)$ and $L(G')$ are the same.

A hierarchy of grammars

We owe a lot of our understanding of grammars to the work of the American linguist Noam Chomsky (yes, the Noam Chomsky known for his politics). There are four categories of formal grammars in the *Chomsky Hierarchy*, they span from Type 0, the most general, to Type 3, the most restrictive. More restrictions on the grammar make it easier to describe and efficiently parse, but reduce the expressive power.

Type 0: free or unrestricted grammars

These are the most general. Productions are of the form $\underline{u} \rightarrow \underline{v}$ where both \underline{u} and \underline{v} are arbitrary strings of symbols in \mathbf{V} , with \underline{u} non-null. There are no restrictions on what appears on the left or right-hand side other than the left-hand side must be non-empty.

Type 1: context-sensitive grammars

Productions are of the form $\underline{u}X\underline{w} \rightarrow \underline{u}v\underline{w}$ where \underline{u} , \underline{v} and \underline{w} are arbitrary strings of symbols in \mathbf{V} , with \underline{v} non-null, and X a single nonterminal. In other words, X may be replaced by \underline{v} but only when it is surrounded by \underline{u} and \underline{w} . (i.e. in a particular context).

Type 2: context-free grammars

Productions are of the form $X \rightarrow \underline{v}$ where \underline{v} is an arbitrary string of symbols in \mathbf{V} , and X is a single nonterminal. Wherever you find X , you can replace with \underline{v} (regardless of context).

Type 3: regular grammars

Productions are of the form $X \rightarrow a$ or $X \rightarrow aY$ where X and Y are nonterminals and a is a terminal. That is the left-hand side must be a single nonterminal and the right-hand side can be either a single terminal by itself or with a single nonterminal. These grammars are the most limited in terms of expressive power.

Every type 3 grammar is a type 2 grammar, and every type 2 is a type 1 and so on. Type 3 grammars are particularly easy to parse because of the lack of recursive constructs. Efficient parsers exist for many classes of Type 2 grammars. Although Type 1 and Type 0 grammars are more powerful than Type 2 and 3, they are far less useful since we cannot create efficient parsers for them. In designing programming languages using formal grammars, we will use Type 2 or context-free grammars, often just abbreviated as CFG.

Issues in parsing context-free grammars

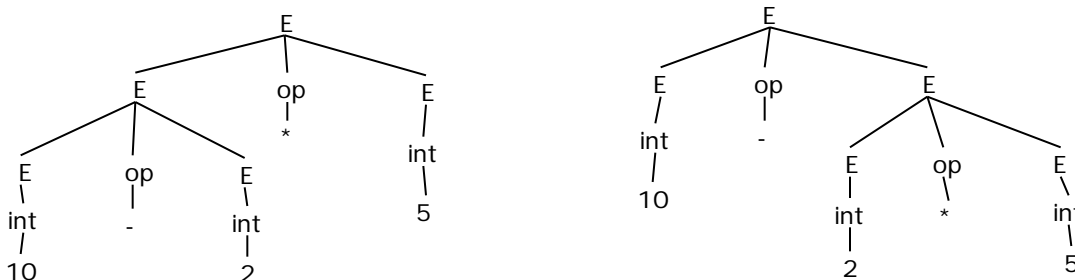
There are several efficient approaches to parsing most Type 2 grammars and we will talk through them over the next few lectures. However, there are some issues that can interfere with parsing that we must take into consideration when designing the grammar. Let's take a look at three of them: ambiguity, recursive rules, and left-factoring.

Ambiguity

If a grammar permits more than one parse tree for some sentences, it is said to be *ambiguous*. For example, consider the following classic arithmetic expression grammar:

$$\begin{array}{ll} E & \rightarrow E \text{ op } E \mid (E) \mid \text{int} \\ \text{op} & \rightarrow + \mid - \mid * \mid / \end{array}$$

This grammar denotes expressions that consist of integers joined by binary operators and possibly including parentheses. As defined above, this grammar is ambiguous because for certain sentences we can construct more than one parse tree. For example, consider the expression $10 - 2 * 5$. We parse by first applying the production $E \rightarrow E \text{ op } E$. The parse tree on the left chooses to expand that first op to $*$, the one on the right to $-$. We have two completely different parse trees. Which one is correct?



Both trees are legal in the grammar as stated and thus either interpretation is valid. Although natural languages can tolerate some kind of ambiguity (puns, plays on words, etc.), it is not acceptable in computer languages. We don't want the compiler just haphazardly deciding which way to interpret our expressions! Given our expectations from algebra concerning precedence, only one of the trees seems right. The right-hand tree fits our expectation that * "binds tighter" and for that result to be computed first then integrated in the outer expression which has a lower precedence operator.

It's fairly easy for a grammar to become ambiguous if you are not careful in its construction. Unfortunately, there is no magical technique that can be used to resolve all varieties of ambiguity. It is an undecidable problem to determine whether any grammar is ambiguous, much less attempt to mechanically remove all ambiguity. However, that doesn't mean in practice that we cannot detect ambiguity or can't do something about it. For programming language grammars, we usually take pains to construct an unambiguous grammar or introduce additional disambiguating rules to throw away the undesirable parse trees, leaving only one for each sentence.

Using the above ambiguous expression grammar, one technique would leave the grammar as is, but add disambiguating rules into the parser implementation. We could code into the parser knowledge of precedence and associativity to break the tie and force the parser to build the tree on the right rather than the left. The advantage of this is that the grammar remains simple and less complicated. But as a downside, the syntactic structure of the language is no longer given by the grammar alone.

Another approach is to change the grammar to only allow the one tree that correctly reflects our intention and eliminate the others. For the expression grammar, we can separate expressions into multiplicative and additive subgroups and force them to be expanded in the desired order.

```

E      ->  E t_op E | T
t_op   ->  + | -
T      ->  T f_op T | F
f_op   ->  * | /
F      ->  (E) | int

```

Terms are addition/subtraction expressions and factors used for multiplication and division. Since the base case for expression is a term, addition and subtraction will appear higher in the parse tree, and thus receive lower precedence.

After verifying that the above re-written grammar has only one parse tree for the earlier ambiguous expression, you might think we were home free, but now consider the expression $10 - 2 - 5$. The recursion on both sides of the binary operator allows either side to match repetitions. The arithmetic operators usually associate to the left, so by replacing the right-hand side with the base case will force the repetitive matches onto the left side. The final result is:

```

E      ->  E t_op T | T
t_op   ->  + | -
T      ->  T f_op F | F
f_op   ->  * | /
F      ->  (E) | int

```

Whew! The obvious disadvantage of changing the grammar to remove ambiguity is that it may complicate and obscure the original grammar definitions. There is no mechanical means to change any ambiguous grammar into an unambiguous one—it is known to be an undecidable problem (in fact, even determining that a CFG is ambiguous is undecidable). However, most programming languages have only limited issues with ambiguity that can be resolved using ad hoc techniques.

Recursive productions

Productions are often defined in terms of themselves. For example a list of variables in a programming language grammar could be specified by this production:

$$\text{variable_list} \rightarrow \text{variable} \mid \text{variable_list} , \text{variable}$$

Such productions are said to be *recursive*. If the recursive nonterminal is at the left of the right-side of the production, e.g. $A \rightarrow \underline{u} \mid A\underline{v}$, we call the production *left-recursive*. Similarly, we can define a *right-recursive* production: $A \rightarrow \underline{u} \mid \underline{v}A$. Some parsing techniques have trouble with one or the other variants of recursive productions and so sometimes we have to massage the grammar into a different but equivalent form. Left-recursive productions can be especially troublesome in the top-down parsers (we'll see why a bit later). Handily, there is a simple technique for rewriting the grammar to move the recursion to the other side. For example, consider this left-recursive rule:

$$X \rightarrow Xa \mid Xb \mid AB \mid C \mid DEF$$

To convert the rule, we introduce a new nonterminal X' that we append to the end of all non-left-recursive productions for X . The expansion for the new nonterminal is basically the reverse of the original left-recursive rule. The re-written productions are:

$$\begin{aligned} X &\rightarrow ABX' \mid CX' \mid DEFX' \\ X' &\rightarrow aX' \mid bX' \mid \end{aligned}$$

It appears we just exchanged the left-recursive rules for an equivalent right-recursive version. This might seem pointless, but some parsing algorithms prefer or even require only left or right recursion.

Left-factoring

The parser usually reads tokens from left to right and it is convenient if upon reading a token it can make an immediate decision about which production from the grammar to expand. However, this can be trouble if there are productions that have common first symbol(s) on the right side of the productions. Here is an example we often see in programming language grammars:

$$\text{Stmt} \rightarrow \text{if Cond then Stmt else Stmt} \mid \text{if Cond then Stmt} \mid \text{Other} \mid \dots$$

The common prefix is `if Cond then Stmt`. This causes problems because when a parser encounter an “if”, it does not know which production to use. A useful technique called *left-factoring* allows us to restructure the grammar to avoid this situation. We rewrite the productions to defer the decision about which of the options to choose until we have seen enough of the input to make the appropriate choice. We factor out the common part of the two options into a shared rule that both will use and then add a new rule that picks up where the tokens diverge.

$$\begin{aligned} \text{Stmt} &\rightarrow \text{if Cond then Stmt OptElse} \mid \text{Other} \mid \dots \\ \text{OptElse} &\rightarrow \text{else S} \mid \end{aligned}$$

In the re-written grammar, upon reading an “if” we expand first production and wait until `if Cond then Stmt` has been seen to decide whether to expand `OptElse` to `else` or `.`

Hidden left-factors and hidden left recursion

A grammar may not appear to have left recursion or left factors, yet still have issues that will interfere with parsing. This may be because the issues are hidden and need to be first exposed via substitution.

For example, consider this grammar:

$$\begin{aligned} A &\rightarrow da \mid acB \\ B &\rightarrow abB \mid daA \mid Af \end{aligned}$$

A cursory examination of the grammar may not detect that the first and second productions of B overlap with the third. We substitute the expansions for A into the third production to expose this:

$$\begin{aligned} A &\rightarrow da \mid acB \\ B &\rightarrow abB \mid daA \mid daf \mid acBf \end{aligned}$$

This exchanges the original third production of B for several new productions, one for each of the productions for A. These directly show the overlap, and we can then left-factor:

$$\begin{aligned} A &\rightarrow da \mid acB \\ B &\rightarrow aM \mid daN \\ M &\rightarrow bB \mid cBf \\ N &\rightarrow A \mid f \end{aligned}$$

Similarly, the following grammar does not appear to have any left-recursion:

$$\begin{aligned} S &\rightarrow Tu \mid wx \\ T &\rightarrow Sq \mid vvS \end{aligned}$$

Yet after substitution of S into T, the left-recursion comes to light:

$$\begin{aligned} S &\rightarrow Tu \mid wx \\ T &\rightarrow Tuq \mid wxq \mid vvS \end{aligned}$$

If we then eliminate left-recursion, we get:

$$\begin{aligned} S &\rightarrow Tu \mid wx \\ T &\rightarrow wxqT' \mid vvST' \\ T' &\rightarrow uqT' \mid \end{aligned}$$

Programming language case study: ALGOL

Algol is of interest to us because it was the first programming language to be defined using a grammar. It grew out of an international effort in the late 1950's to create a "universal programming language" that would run on all machines. At that time, FORTRAN and COBOL were the prominent languages, with new languages sprouting up all over the place. Programmers became increasingly concerned about portability of programs, and being able to communicate with one another on programming topics.

Consequently the ACM and GAMM (Gesellschaft für angewandte Mathematik und Mechanik) decided to come up with a single programming language that all could use on their computers, and in whose terms, programs could be communicated between the users of all machines. Their first decision was not to use FORTRAN as their universal language. This may seem surprising to us today, since it was the most commonly used language back then. However, as Alan J. Perlis, one of the original committee members, puts it:

"Today, FORTRAN is the property of the computing world, but in 1957, it was an IBM creation and closely tied to IBM hardware. For these reasons, FORTRAN was unacceptable as a universal language."

ALGOL-58 was the first version of the language, followed up very soon after by ALGOL-60, which is the version that had the most impact. As a language, it introduced the following features:

- block structure and nested structures
- strong typing
- scoping
- procedures and functions
- call by value, call by reference
- side effects (is this good or bad?)
- recursion

It may seem surprising that recursion was not present in the original FORTRAN or COBOL. You probably know that to implement recursion we need a runtime stack to store the activation records as functions are called. In FORTRAN and COBOL, activation records were created at compile time, not runtime. Thus, only one activation record per subroutine was created. No stack was used. The parameters for the subroutine were copied into the activation record and that data area was used for subroutine processing.

The ALGOL report was the first time we see BNF to describe a programming language. Both John Backus and Peter Naur were on the ALGOL committees. They derived this description technique from an earlier paper written by Backus. The technique was adopted because they needed a machine-independent method of description. If one looks at the early definitions of FORTRAN, one can see the links to the IBM hardware. With ALGOL, the machine was not relevant. BNF had a huge impact on programming language design and compiler construction. First, it stimulated a large number of studies on the formal structure of programming languages laying the groundwork for a theoretical approach to language design. Second, a formal syntactic description could be used to drive a compiler directly (as we shall see).

ALGOL had a tremendous impact on programming language design, compiler construction, and language theory, but the language itself was a commercial failure. Partly this was due to design decisions (overly complex features, no IO) along with the politics of the time (popularity of Fortran, lack of support from the all-powerful IBM, resistance to BNF).

Bibliography

- A. Aho, R. Sethi, J. Ullman, *Compilers: Principles, Techniques, and Tools*. Reading, MA: Addison-Wesley, 1986.
- J. Backus, "The Syntax and Semantics of the Proposed International Algebraic Language of the Zurich ACM-GAMM Conference," *Proceedings of the International Conference on Information Processing*, 1959, pp. 125-132.
- N. Chomsky, "On Certain Formal Properties of Grammars," *Information and Control*, Vol. 2, 1959, pp. 137-167.
- J.P. Bennett, *Introduction to Compiling Techniques*. Berkshire, England: McGraw-Hill, 1990.
- D. Cohen, *Introduction to Computer Theory*. New York: Wiley, 1986.
- J.C. Martin, *Introduction to Languages and the Theory of Computation*. New York, NY: McGraw-Hill, 1991.
- P. Naur, "Programming Languages, Natural Languages, and Mathematics," *Communications of the ACM*, Vol 18, No. 12, 1975, pp. 676-683.
- J. Sammet, *Programming Languages: History and Fundamentals*. Englewood-Cliffs, NJ: Prentice-Hall, 1969.
- R.L. Wexelblat, *History of Programming Languages*. London: Academic Press, 1981.