



1. Overview

1.1 Motivation

1.2 Structure of a Compiler

1.3 Grammars

1.4 Syntax Tree and Ambiguity

1.5 Chomsky's Classification of Grammars

1.6 The Z# Language

Short History of Compiler Construction



Formerly "a mystery", today one of the best-known areas of computing

- | | | |
|-------------|----------------|---|
| 1957 | Fortran | first compilers
(arithmetic expressions, statements, procedures) |
| 1960 | Algol | first formal language definition
(grammars in Backus-Naur form, block structure, recursion, ...) |
| 1970 | Pascal | user-defined types, virtual machines (P-code) |
| 1985 | C++ | object-orientation, exceptions, templates |
| 1995 | Java | just-in-time compilation |

We only look at imperative languages

Functional languages (e.g. Lisp) and logical languages (e.g. Prolog) require different techniques.

Why should I learn about compilers?



It's part of the general background of a software engineer

- How do compilers work?
- How do computers work?
(instruction set, registers, addressing modes, run-time data structures, ...)
- What machine code is generated for certain language constructs?
(efficiency considerations)
- What is good language design?
- Opportunity for a non-trivial programming project

Also useful for general software development

- Reading syntactically structured command-line arguments
- Reading structured data (e.g. XML files, part lists, image files, ...)
- Searching in hierarchical namespaces
- Interpretation of command codes
- ...



1. Overview

1.1 Motivation

1.2 Structure of a Compiler

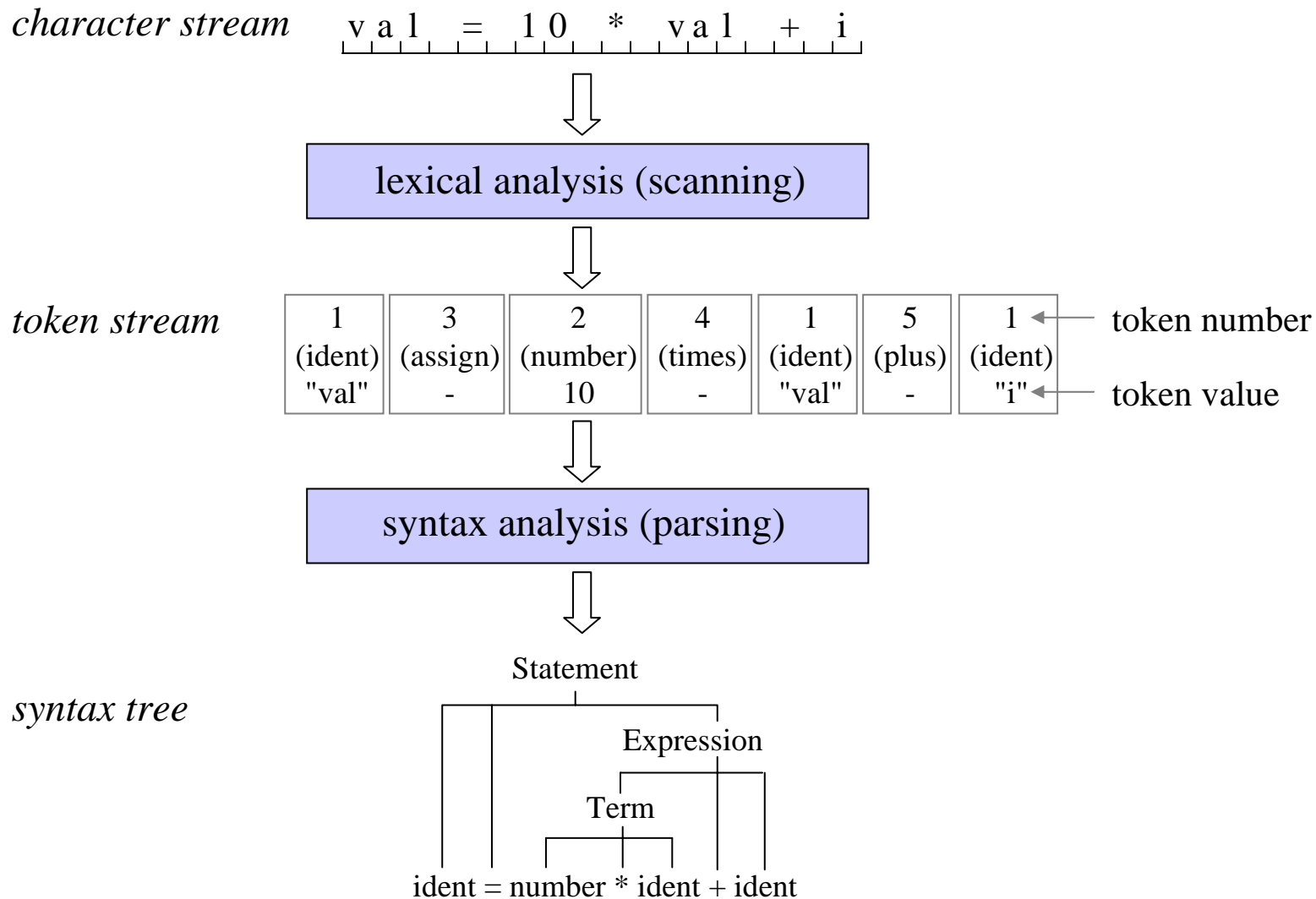
1.3 Grammars

1.4 Syntax Tree and Ambiguity

1.5 Chomsky's Classification of Grammars

1.6 The Z# Language

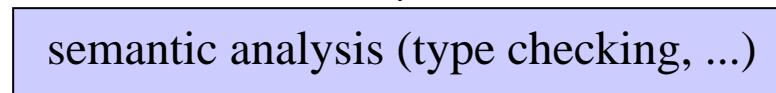
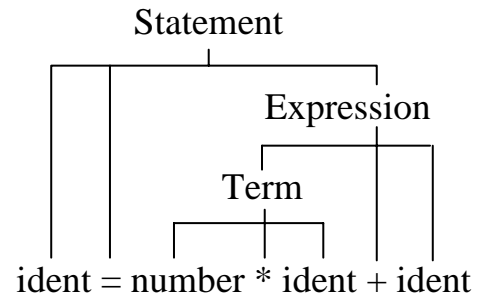
Dynamic Structure of a Compiler



Dynamic Structure of a Compiler

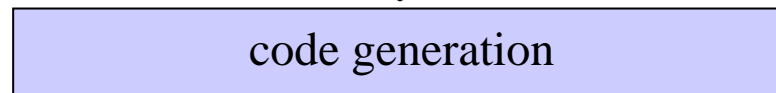
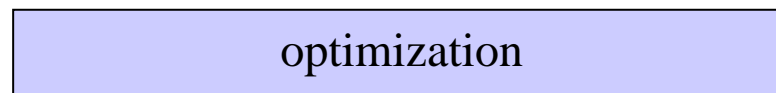


syntax tree



intermediate representation

syntax tree, symbol table, ...



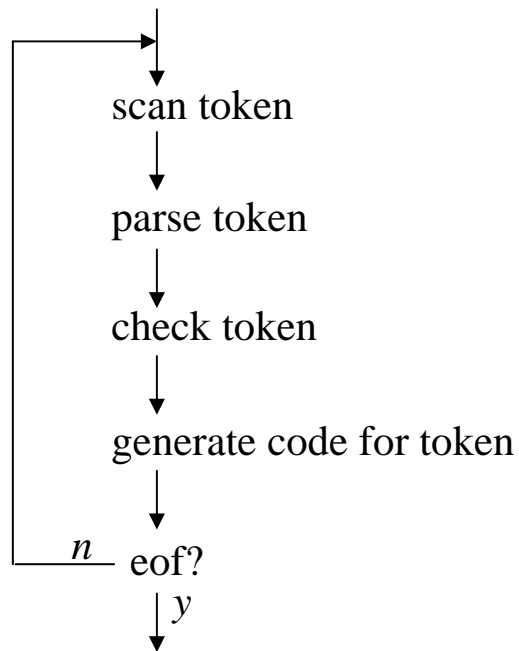
machine code

ld.i4.s 10
ldloc.1
mul
...

Single-Pass Compilers



Phases work in an interleaved way

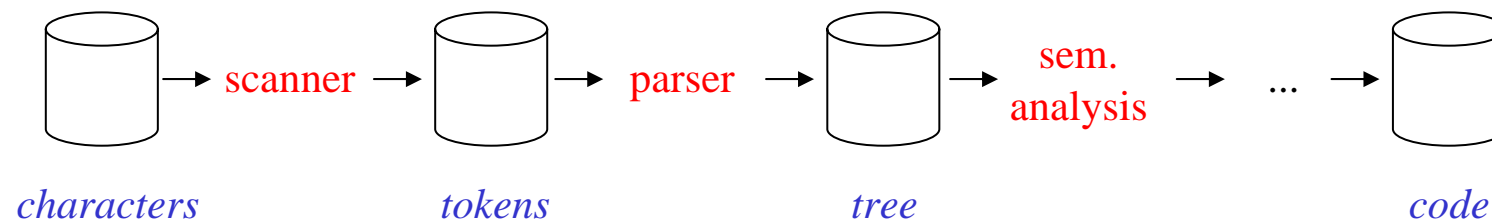


The target program is already generated while the source program is read.

Multi-Pass Compilers



Phases are separate "programs", which run sequentially

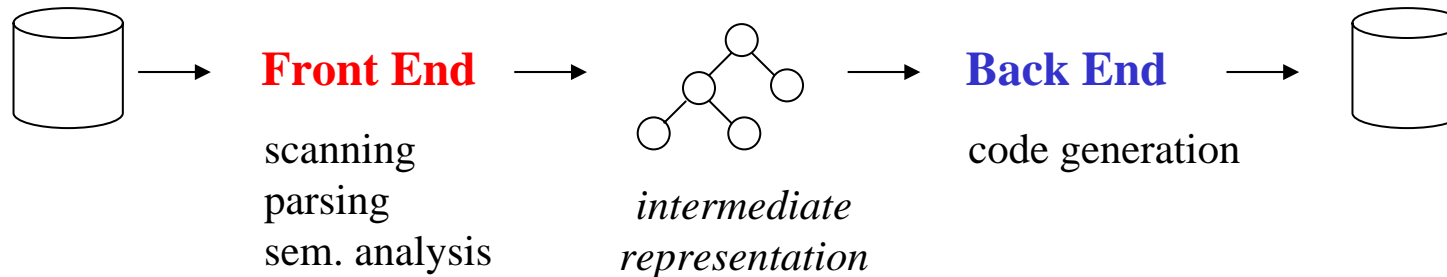


Each phase reads from a file and writes to a new file

Why multi-pass?

- if memory is scarce (irrelevant today)
- if the language is complex
- if portability is important

Today: Often Two-Pass Compilers



language-dependent

Java

C

Pascal

machine-dependent

Pentium

PowerPC

SPARC

any combination possible

Advantages

- better portability
- many combinations between front ends and back ends possible
- optimizations are easier on the intermediate representation than on source code

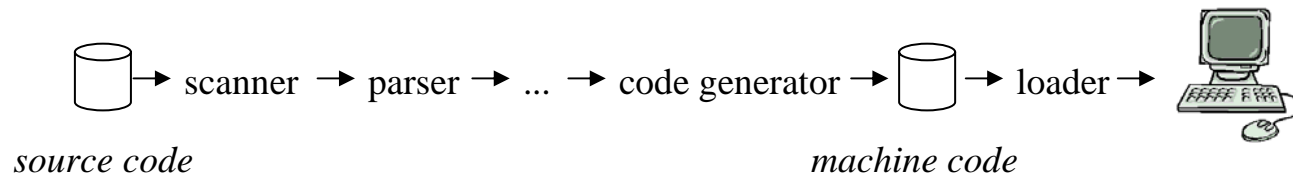
Disadvantages

- slower
- needs more memory

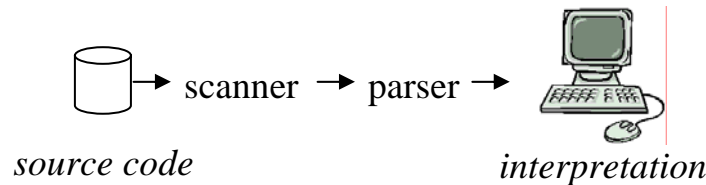
Compiler versus Interpreter



Compiler translates to machine code

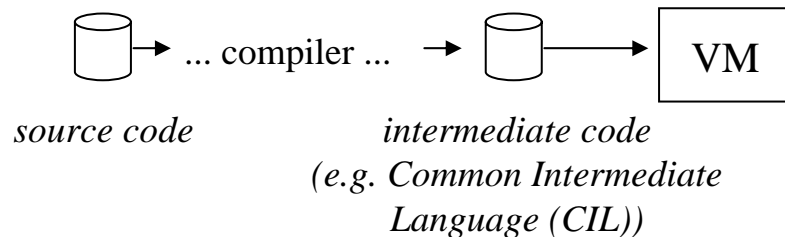


Interpreter executes source code "directly"



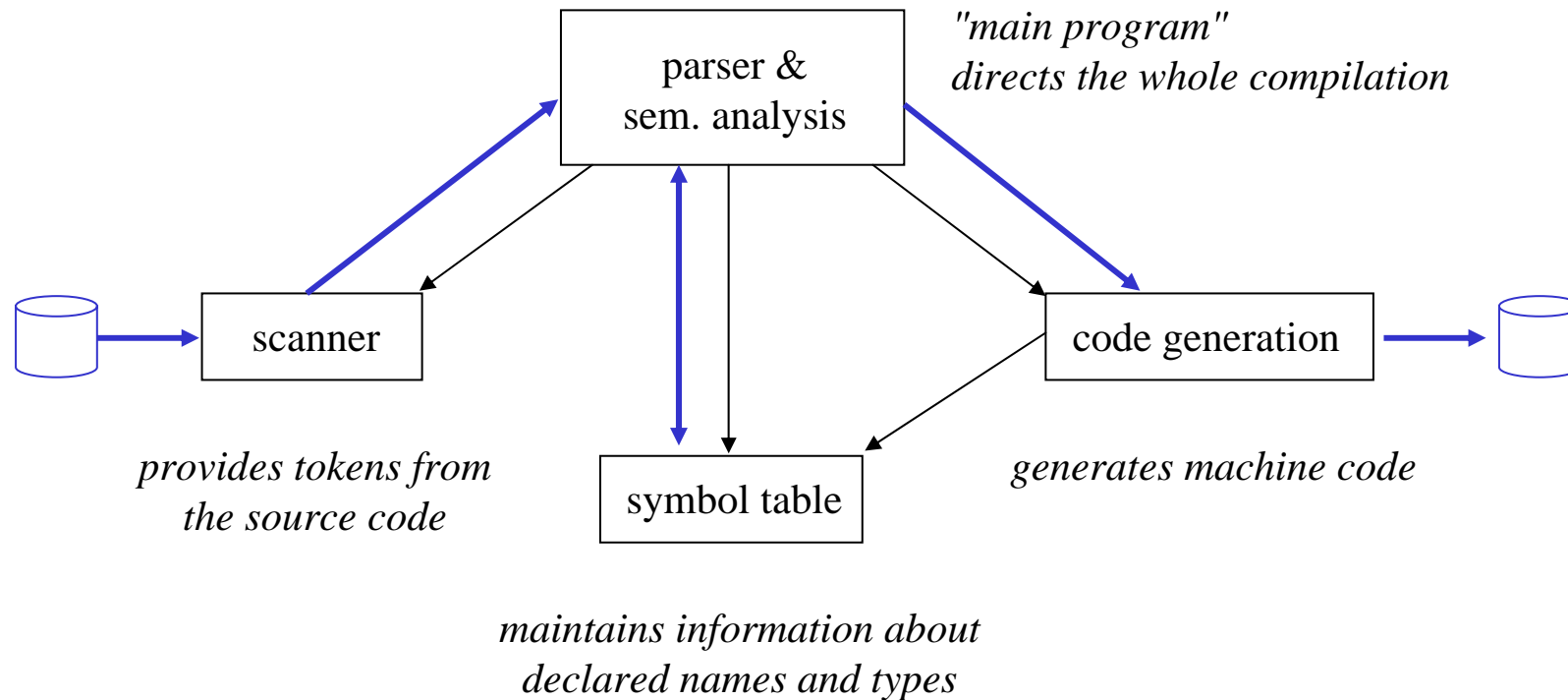
- statements in a loop are scanned and parsed again and again

Variant: interpretation of intermediate code



- source code is translated into the code of a *virtual machine* (VM)
- VM interprets the code simulating the physical machine

Static Structure of a Compiler



→ uses

→ data flow



1. Overview

1.1 Motivation

1.2 Structure of a Compiler

1.3 Grammars

1.4 Syntax Tree and Ambiguity

1.5 Chomsky's Classification of Grammars

1.6 The Z# Language

What is a grammar?



Example Statement = "if" "(" Condition ")" Statement ["else" Statement].

Four components

terminal symbols	are atomic	"if", ">=", ident, number, ...
nonterminal symbols	are derived into smaller units	Statement, Expr, Type, ...
productions	rules how to decompose nonterminals	Statement = Designator "=" Expr ";" Designator = ident ["." ident] ...
start symbol	topmost nonterminal	CSharp

EBNF Notation



Extended Backus-Naur form

John Backus: developed the first Fortran compiler

Peter Naur: edited the Algol60 report

<i>symbol</i>	<i>meaning</i>	<i>examples</i>
string	denotes itself	"=", "while"
name	denotes a T or NT symbol	ident, Statement
=	separates the sides of a production	A = b c d .
.	terminates a production	
	separates alternatives	a b c ≡ a or b or c
(...)	groups alternatives	a (b c) ≡ ab ac
[...]	optional part	[a] b ≡ ab b
{...}	repetitive part	{ a } b ≡ b ab aab aaab ...

Conventions

- terminal symbols start with lower-case letters (e.g. ident)
- nonterminal symbols start with upper-case letters (e.g. Statement)

Example: Grammar for Arithmetic Expressions



Productions

Expr = ["+" | "-"] Term { ("+" | "-") Term }.
Term = Factor { ("*" | "/") Factor }.
Factor = ident | number | "(" Expr ")".

Terminal symbols

simple TS: "+" , "-" , "*" , "/" , "(" , ")"
(just 1 instance)

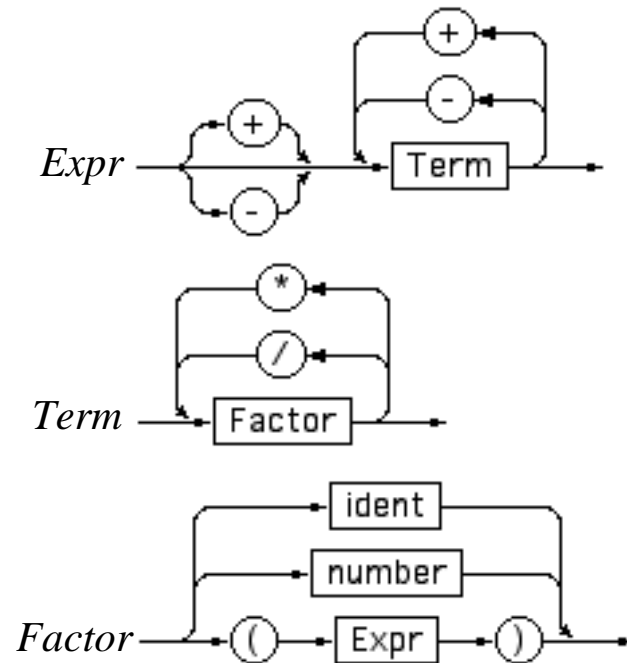
terminal classes: ident, number
(multiple instances)

Nonterminal symbols

Expr, Term, Factor

Start symbol

Expr



Operator Priority



Grammars can be used to define the priority of operators

```
Expr = [ "+" | "-" ] Term { ( "+" | "-" ) Term }.  
Term = Factor { ( "*" | "/" ) Factor }.  
Factor = ident | number | "(" Expr ")".
```

input: $-a * 3 + b / 4 - c$

\Rightarrow - ident * number + ident / number - ident
 \Rightarrow - $\underbrace{\text{Factor} * \text{Factor}} + \underbrace{\text{Factor} / \text{Factor}} - \underbrace{\text{Factor}}$
 \Rightarrow - $\underbrace{\text{Term} + \text{Term} - \text{Term}}$
 \Rightarrow $\underbrace{\hspace{10em}}_{\text{Expr}}$

"*" and "/" have higher priority than "+" and "-"
"- " does not refer to a , but to $a*3$

How must the grammar be transformed, so that "-" refers to a ?

Terminal Start Symbols of Nonterminals



Which terminal symbols can a nonterminal start with?

```
Expr = ["+" | "-"] Term {"+" | "-"} Term}.  
Term = Factor {"*" | "/" } Factor}.  
Factor = ident | number | "(" Expr ")".
```

First(Factor) = **ident, number, "("**

First(Term) = First(Factor)
= **ident, number, "("**

First(Expr) = "+", "-", First(Term)
= **+", "-", ident, number, "("**

Terminal Successors of Nonterminals



Which terminal symbols can follow after a nonterminal in the grammar?

```
Expr = [ "+" | "-" ] Term { ( "+" | "-" ) Term }.  
Term = Factor { ( "*" | "/" ) Factor }.  
Factor = ident | number | "(" Expr ")".
```

Follow(Expr) = **)", eof**

Follow(Term) = "+", "-", Follow(Expr)
= **+", "-",)", eof**

Follow(Factor) = "*", "/", Follow(Term)
= ***, /, +, -,)", eof**

Where does *Expr* occur on the right-hand side of a production?
What terminal symbols can follow there?

Some Terminology



Alphabet

The set of terminal and nonterminal symbols of a grammar

String

A finite sequence of symbols from an alphabet.

Strings are denoted by greek letters (α , β , γ , ...)

e.g: $\alpha = \text{ident} + \text{number}$

$\beta = - \text{Term} + \text{Factor} * \text{number}$

Empty String

The string that contains no symbol (denoted by ϵ).

Derivations and Reductions



Derivation

$$\alpha \Rightarrow \beta \quad (\text{direct derivation}) \quad \underbrace{\text{Term} + \underbrace{\text{Factor}}_{\text{NTS}} * \text{Factor}}_{\alpha} \Rightarrow \underbrace{\text{Term} + \underbrace{\text{ident}}_{\text{right-hand side of a production of NTS}} * \text{Factor}}_{\beta}$$

$$\alpha \Rightarrow^* \beta \quad (\text{indirect derivation}) \quad \alpha \Rightarrow \gamma_1 \Rightarrow \gamma_2 \Rightarrow \dots \Rightarrow \gamma_n \Rightarrow \beta$$

$$\alpha \Rightarrow^L \beta \quad (\text{left-canonical derivation}) \quad \text{the } \underline{\text{leftmost}} \text{ NTS in } \alpha \text{ is derived first}$$

$$\alpha \Rightarrow^R \beta \quad (\text{right-canonical derivation}) \quad \text{the } \underline{\text{rightmost}} \text{ NTS in } \alpha \text{ is derived first}$$

Reduction

The converse of a derivation:

If the right-hand side of a production occurs in β it is replaced with the corresponding NTS

Deletability



A string α is called deletable, if it can be derived to the empty string.

$$\alpha \Rightarrow^* \varepsilon$$

Example

A = B C.
B = [b].
C = c | d | .

B is deletable: $B \Rightarrow \varepsilon$

C is deletable: $C \Rightarrow \varepsilon$

A is deletable: $A \Rightarrow B C \Rightarrow C \Rightarrow \varepsilon$



More Terminology

Phrase

Any string that can be derived from a nonterminal symbol.

Term phrases: Factor
 Factor * Factor
 ident * Factor
 ...

Sentential form

Any string that can be derived from the start symbol of the grammar.

e.g.: Expr
 Term + Term + Term
 Term + Factor * ident + Term
 ...

Sentence

A sentential form that consists of terminal symbols only.

e.g.: ident * number + ident

Language (formal language)

The set of all sentences of a grammar (usually infinitely large).

e.g.: the C# language is the set of all valid C# programs



Recursion

A production is recursive if $A \Rightarrow^* \omega_1 A \omega_2$

Can be used to represent repetitions and nested structures

Direct recursion $A \Rightarrow \omega_1 A \omega_2$

Left recursion $A = b \mid A a.$ $A \Rightarrow A a \Rightarrow A a a \Rightarrow A a a a \Rightarrow b a a a a \dots$

Right recursion $A = b \mid a A.$ $A \Rightarrow a A \Rightarrow a a A \Rightarrow a a a A \Rightarrow \dots a a a a b$

Central recursion $A = b \mid "(" A ")".$ $A \Rightarrow (A) \Rightarrow ((A)) \Rightarrow (((A))) \Rightarrow (((... (b)...)))$

Indirect recursion $A \Rightarrow^* \omega_1 A \omega_2$

Example

Expr = Term { "+" Term }.
Term = Factor { "*" Factor }.
Factor = id | "(" Expr ")".

Expr \Rightarrow Term \Rightarrow Factor \Rightarrow "(" Expr ")"



How to Remove Left Recursion

Left recursion cannot be handled in topdown syntax analysis

$A = b \mid A a.$ Both alternatives start with b .
The parser cannot decide which one to choose

Left recursion can be transformed to iteration

$E = T \mid E "+" T.$

What sentences can be derived?

T
T + T
T + T + T
...

From this one can deduce the iterative EBNF rule:

$E = T \{ "+" T \}.$



1. Overview

1.1 Motivation

1.2 Structure of a Compiler

1.3 Grammars

1.4 Syntax Tree and Ambiguity

1.5 Chomsky's Classification of Grammars

1.6 The Z# Language



Plain BNF Notation

- terminal symbols* are written without quotes (ident, +, -)
nonterminal symbols are written in angle brackets (<Expr>, <Term>)
sides of a production are separated by ::=

BNF grammar for arithmetic expressions

```
<Expr> ::= <Sign> <Term>
<Expr> ::= <Expr> <Addop> <Term>
<Sign> ::= +
<Sign> ::= -
<Sign> ::=
<Addop> ::= +
<Addop> ::= -
<Term> ::= <Factor>
<Term> ::= <Term> <Mulop> <Factor>
<Mulop> ::= *
<Mulop> ::= /
<Factor> ::= ident
<Factor> ::= number
<Factor> ::= ( <Expr> )
```

- Alternatives are transformed into separate productions
- Repetition must be expressed by recursion

Advantages

- fewer meta symbols (no |, (), [], {})
- it is easier to build a syntax tree

Disadvantages

- more clumsy

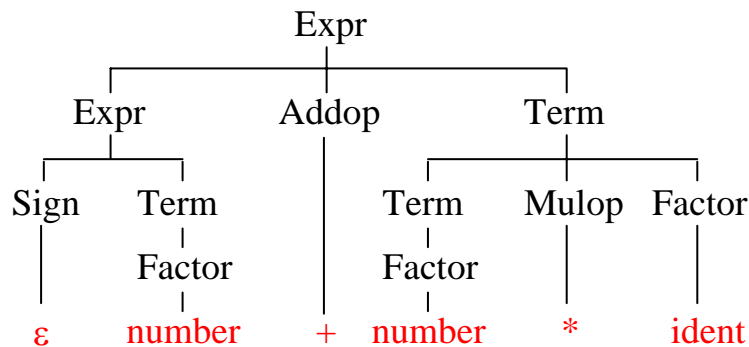
Syntax Tree



Shows the structure of a particular sentence

e.g. for $10 + 3 * i$

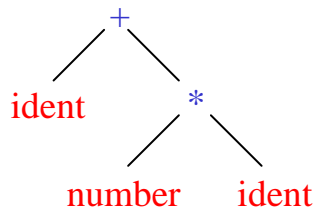
Concrete Syntax Tree (parse tree)



Would not be possible with EBNF because of [...] and {...}, e.g.:
 $\text{Expr} = [\text{Sign}] \text{Term} \{ \text{Addop} \text{Term} \}.$

Also reflects operator priorities: operators further down in the tree have a higher priority than operators further up in the tree.

Abstract Syntax Tree (leaves = operands, inner nodes = operators)



often used as an internal program representation;
used for optimizations

Ambiguity



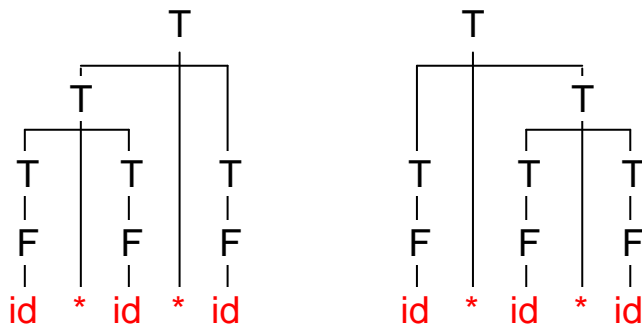
A grammar is ambiguous, if more than one syntax tree can be built for a given sentence.

Example

$T = F \mid T "*" T.$
 $F = id.$

sentence: $id * id * id$

Two syntax trees can be built for this sentence:



Ambiguous grammars cause problems in syntax analysis!



Removing Ambiguity

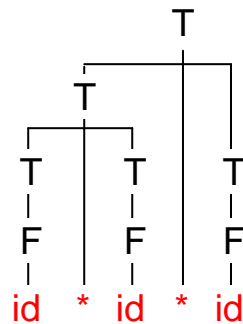
Example

$T = F \mid T "*" T.$
 $F = \text{id}.$

Only the grammar is ambiguous, not the language.

The grammar can be transformed to

$T = F \mid T "*" F.$
 $F = \text{id}.$



i.e. T has priority over F

only this syntax tree is possible

Even better: transformation to EBNF

$T = F \{ "*" F \}.$
 $F = \text{id}.$

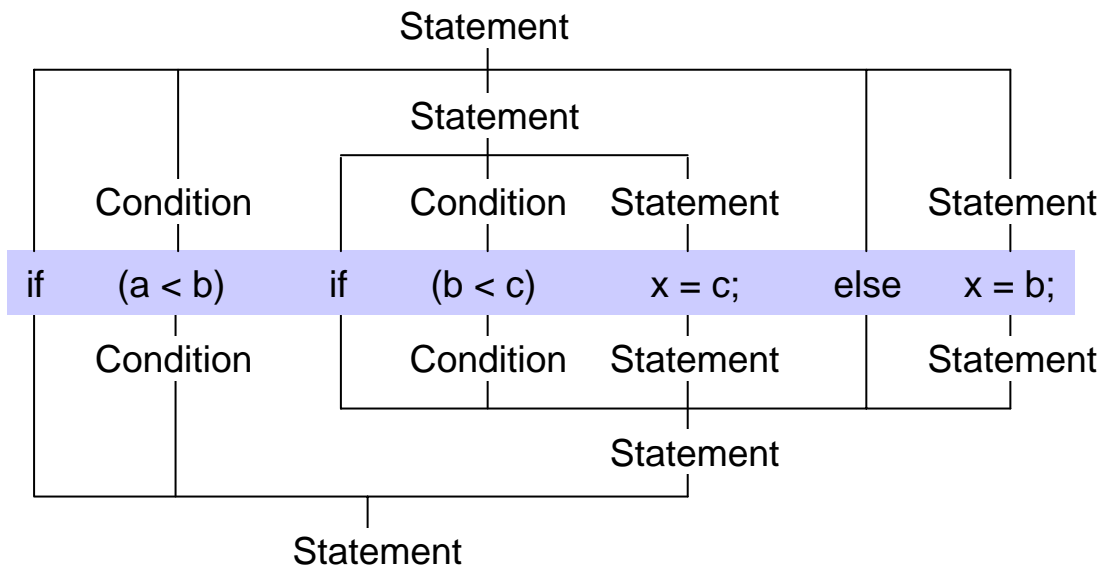


Inherent Ambiguity

There are languages which are inherently ambiguous

Example: *Dangling Else*

```
Statement = Assignment  
          | "if" Condition Statement  
          | "if" Condition Statement "else" Statement  
          | ... .
```



There is no unambiguous grammar for this language!

C# solution

Always recognize the longest possible right-hand side of a production

⇒ leads to the lower of the two syntax trees



1. Overview

1.1 Motivation

1.2 Structure of a Compiler

1.3 Grammars

1.4 Syntax Tree and Ambiguity

1.5 Chomsky's Classification of Grammars

1.6 The Z# Language

Classification of Grammars



Due to Noam Chomsky (1956)

Grammars are sets of productions of the form $\alpha = \beta$.

class 0 **Unrestricted grammars** (α and β arbitrary)

e.g: $A = a A b \mid B c B$.

$aBc = d$.

$dB = bb$.

$A \Rightarrow aAb \Rightarrow aBcBb \Rightarrow dBb \Rightarrow bbb$

Recognized by Turing machines

class 1 **Contex-sensitive grammars** ($|\alpha| \leq |\beta|$)

e.g: $a A = a b c$.

Recognized by linear bounded automata

class 2 **Context-free grammars** ($\alpha = NT, \beta \neq \varepsilon$)

e.g: $A = a b c$.

Recognized by push-down automata

class 3 **Regular grammars** ($\alpha = NT, \beta = T \mid T NT$)

e.g: $A = b \mid b B$.

Recognized by finite automata

} Only these two classes
are relevant in compiler
construction



1. Overview

1.1 Motivation

1.2 Structure of a Compiler

1.3 Grammars

1.4 Syntax Tree and Ambiguity

1.5 Chomsky's Classification of Grammars

1.6 The Z# Language

Sample Z# Program



```
class P
  const int size = 10;
  class Table {
    int[] pos;
    int[] neg;
  }
  Table val;
{
  void Main ()
    int x, i;
  { /*----- initialize val -----*/
    val = new Table;
    val.pos = new int[size];
    val.neg = new int[size];
    i = 0;
    while (i < size) {
      val.pos[i] = 0; val.neg[i] = 0; i++;
    }
    /*----- read values -----*/
    read(x);
    while (-size < x && x < size) {
      if (0 <= x) { val.pos[x]++; }
      else { val.neg[-x]++; }
      read(x);
    }
  }
}
```

main program class; no separate compilation

inner classes (without methods)

global variables

local variables

Lexical Structure of Z#



Names ident = letter { letter | digit | '_' }.

Numbers number = digit { digit }. all numbers are of type *int*

Char constants charConst = \" char \". all character constants are of type *char*
no strings (may contain \r and \n)

Keywords class
if else while read write return break
void const new

Operators + - * / % ++ --
== != > >= < <=
&& ||
() [] { }
= ; , .

Comments /* ... */ may be nested

Types *int* *char* arrays classes

Syntactical Structure of Z# (1)



Programs

```
Program = "class" ident
         { ConstDecl | VarDecl | ClassDecl }
         "{" { MethodDecl } "}"
```

```
class P
  ... declarations ...
{ ... methods ...
}
```

Declarations

```
ConstDecl = "const" Type ident "=" ( number | charConst ) ";"
VarDecl   = Type ident { "," ident } ";"
MethodDecl = (Type | "void") ident "(" [ FormPars ] ")" Block
Type       = ident [ "[" "]" ]
FormPars   = Type ident { "," Type ident }
```

only one-dimensional arrays

Syntactical Structure of Z# (2)



Statements

```
Block      = "{" {Statement} }".
Statement  = Designator ( "=" Expr ";"
                    | "(" [ActPars] ")" ";"
                    | "++" ";"
                    | "--" ";"
                    )
            | "if" "(" Condition ")" Block [ "else" Block ]
            | "while" "(" Condition ")" Block
            | "break" ";"
            | "return" [ Expr ] ";"
            | "read" "(" Designator ")" ";"
            | "write" "(" Expr [ "," number ] ")" ";"
            | ";".
ActPars    = Expr { "," Expr }.
```

- input from *System.Console*
- output to *System.Console*

Syntactical Structure of Z# (3)



Expressions

Condition = CondTerm { "|" CondTerm }.
CondTerm = CondFact { "&&" CondFact }.
CondFact = Expr Relop Expr.
Relop = "==" | "!=" | ">" | ">=" | "<" | "<=".

Expr = ["-"] Term { Addop Term }.
Term = Factor { Mulop Factor }.
Factor = Designator ["(" [ActPars] ")"]
| number
| charConst
| "new" ident ["[" Expr "]"]
| "(" Expr ")".
Designator = ident ["[" Expr "]"] { "." ident ["[" Expr "]"] }.
Addop = "+" | "-".
Mulop = "*" | "/" | "%".

no constructors