# Programming Languages
## 2nd edition
### Tucker and Noonan

Chapter 2
Syntax

*A language that is simple to parse for the compiler is also simple to parse for the human programmer.*

*N. Wirth*

# Contents

## 2.3 Syntax of a Small Language: *Clite*

Motivation for using a subset of C:

| Language | Grammar (pages) | Reference |
|---|---|---|
| Pascal | 5 | Jensen & Wirth |
| C | 6 | Kernighan & Richie |
| C++ | 22 | Stroustrup |
| Java | 14 | Gosling, et. al. |

The *Clite* grammar fits on one page (next 3 slides),
so it's a far better tool for studying language design.

# Fig. 2.7 *Clite* Grammar: Statements

$Program \rightarrow$ `int main ( ) {` *Declarations Statements* `}`

$Declarations \rightarrow$ **{** *Declaration* **}**

$Declaration \rightarrow Type\ Identifier$ **[** **[** *Integer* **]** **]** **{** **,** *Identifier* **[** **[** *Integer* **]** **]** **}**

$Type \rightarrow$ `int | bool | float | char`

$Statements \rightarrow$ **{** *Statement* **}**

$Statement \rightarrow$ **;** | *Block* | *Assignment* | *IfStatement* | *WhileStatement*

$Block \rightarrow$ **{** *Statements* **}**

$Assignment \rightarrow Identifier$ **[** **[** *Expression* **]** **]** = *Expression* **;**

$IfStatement \rightarrow$ `if` **(** *Expression* **)** *Statement* **[** `else` *Statement* **]**

$WhileStatement \rightarrow$ `while` **(** *Expression* **)** *Statement*

# Fig. 2.7 *Clite* Grammar: Expressions

*Expression* → *Conjunction* **{** || *Conjunction* **}**

*Conjunction* → *Equality* **{** && *Equality* **}**

*Equality* → *Relation* **[** *EquOp Relation* **]**

*EquOp* → == | !=

*Relation* → *Addition* **[** *RelOp Addition* **]**

*RelOp* → < | <= | > | >=

*Addition* → *Term* **{** *AddOp Term* **}**

*AddOp* → + | −

*Term* → *Factor* **{** *MulOp Factor* **}**

*MulOp* → * | / | %

*Factor* → **[** *UnaryOp* **]** *Primary*

*UnaryOp* → − | !

*Primary* → *Identifier* **[** [ *Expression* ] **]** | *Literal* | ( *Expression* ) |
Type ( *Expression* )

# Fig. 2.7 *Clite* grammar: lexical level

*Identifier* → *Letter* **{** *Letter* | *Digit* **}**

*Letter* → `a | b | ... | z | A | B | ... | Z`

*Digit* → `0 | 1 | ... | 9`

*Literal* → *Integer* | *Boolean* | *Float* | *Char*

*Integer* → *Digit* **{** *Digit* **}**

*Boolean* → `true | False`

*Float* → *Integer* `.` *Integer*

*Char* → ' *ASCII Char* '

# Issues Not Addressed by this Grammar

- Comments

- Whitespace

- Distinguishing one token <= from two tokens <  =

- Distinguishing identifiers from keywords like if

These issues are addressed by identifying two levels:

- *lexical level*

- *syntactic level*

# 2.3.1 Lexical Syntax

*Input*: a stream of characters from the ASCII set, keyed by a programmer.

*Output*: a stream of *tokens* or basic symbols, classified as follows:

- *Identifiers*      e.g., Stack, x, i, push
- *Literals*         e.g., 123, 'x', 3.25, true
- *Keywords*         bool char else false float if int
                     main true while
- *Operators*        = || && == != < <= > >= + - * / !
- *Punctuation*      ; , { } ( )

# Whitespace

Whitespace is any space, tab, end-of-line character (or characters), or character sequence inside a comment

No token may contain embedded whitespace (unless it is a character or string literal)

Example:

>=   *one token*

>  = *two tokens*

# Whitespace Examples in Pascal

while  a  <  b  do          *legal* - spacing between tokens

while  a<b  do                    spacing not needed for <


whilea<bdo                  *illegal* - can't tell boundaries

whilea  <  bdo                    between tokens

# Comments

Not defined in grammar

*Clite* uses // comment style of C++

# Identifier

Sequence of letters and digits, starting with a letter

if is both an identifier and a keyword

Most languages require identifiers to be distinct from keywords

In some languages, identifiers are merely predefined (and thus can be redefined by the programmer)

# Redefining Identifiers can be dangerous

```
program confusing;
const true = false;
begin
    if  (a<b) = true then
f(a)
    else …
```

# Should Identifiers be case-sensitive?

Older languages: no.   Why?

- *Pascal: no.*

- *Modula: yes*

- *C, C++: yes*

- *Java: yes*

- *PHP: partly yes, partly no.  What about orthogonality?*

# 2.3.2 Concrete Syntax

Based on a parse of its *Tokens*

*; is a statement terminator*

*(Algol-60, Pascal use ; as a separator)*

Rule for *IfStatement* is ambiguous:

"The else ambiguity is resolved by connecting
    an **else** with the last encountered else-less if."

[Stroustrup, 1991]

# Expressions in *Clite*

13 grammar rules

Use of meta braces – operators are left associative

C++ expressions require 4 pages of grammar rules [Stroustrup]

C uses an ambiguous expression grammar [Kernighan and Ritchie]

# Associativity and Precedence

| Clite Operator | Associativity |
|---|---|
| Unary - ! | none |
| * / | left |
| + - | left |
| < <= > >= | none |
| == != | none |
| && | left |
| \|\| | left |

# *Clite* Equality, Relational Operators

… are non-associative.

(an idea borrowed from Ada)

Why is this important?

In C++, the expression:

```
if (a < x < b)
```
is *not* equivalent to
```
if (a < x && x < b)
```
But it is error-free!
So, what does it mean?