

# Frontend: Lexing & Parsing

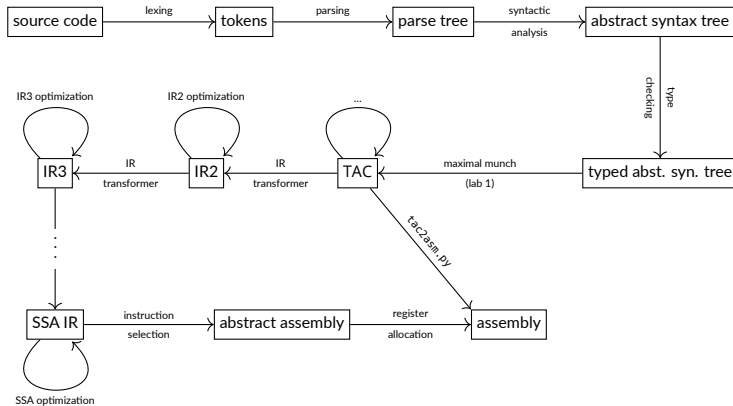
## (CSC 3F002 EP) Compilers

**Pierre-Yves Strub**

(Slide deck author: **Kaustuv Chaudhuri**)

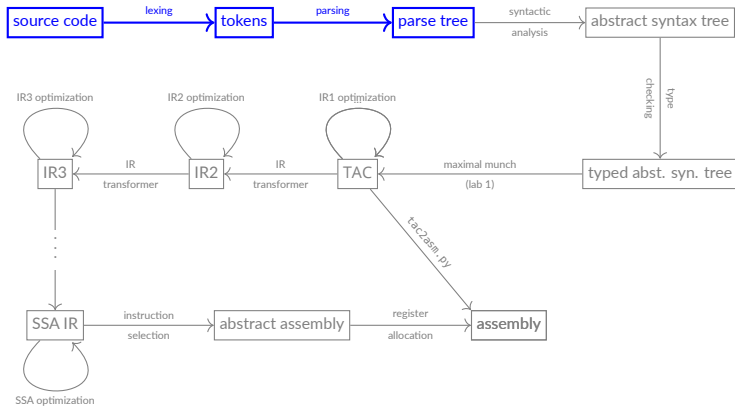
2025-09-04 // Week 1

# Compiler Stages: Block Diagram



# Today's Focus

(will be part of lab 1)



# BX Grammar (Straightline Fragment)

(we'll use this for the first labs)

```
⟨program⟩ ::= "def main() { " (⟨vardecl⟩ | ⟨stmt⟩)* "}"
```

```
⟨vardecl⟩ ::= "var" ⟨ident⟩ "=" ⟨expr⟩ ": int;"
```

```
⟨stmt⟩ ::= ⟨assign⟩ | ⟨print⟩
```

```
⟨assign⟩ ::= IDENT "=" ⟨expr⟩ ";"
```

```
⟨print⟩ ::= "print(" ⟨expr⟩ ");"
```

```
⟨expr⟩ ::= IDENT | NUMBER
```

```
      | ⟨expr⟩ "+" ⟨expr⟩ | ⟨expr⟩ "-" ⟨expr⟩
```

```
      | ⟨expr⟩ "*" ⟨expr⟩ | ⟨expr⟩ "/" ⟨expr⟩ | ⟨expr⟩ "%" ⟨expr⟩
```

```
      | ⟨expr⟩ "&" ⟨expr⟩ | ⟨expr⟩ "|" ⟨expr⟩ | ⟨expr⟩ "^" ⟨expr⟩
```

```
      | ⟨expr⟩ "<<" ⟨expr⟩ | ⟨expr⟩ ">>" ⟨expr⟩
```

```
      | "-" ⟨expr⟩ | "~" ⟨expr⟩
```

```
      | "(" ⟨expr⟩ ")"
```

```
IDENT ::= / [A-Za-z][A-Za-a0-9_]* /
```

```
NUMBER ::= / 0 | [1-9][0-9]* /
```

(except reserved words)

(value must fit in 63 bits)

# Compiler Frontend: Reading the Source

- Steps:
  - 1 Convert BX0 source code (text) to **token** sequence (lexing)
  - 2 Parse token sequence (parsing)
  - 3 generate an **abstract syntax tree** (AST)
  - 4 Perform **syntactic correctness** checks on the AST
  - 5 Perform **type checking** to compute and store type information
  - 6 (Optional) Simplify the AST when necessary  
(e.g., convert **for** loops to **while** loops)
- Lexing and Parsing
  - Can be written **by hand** in e.g. Python
    - Flexible – can parse very complicated grammars
    - Labor intensive, esp. for maintenance
    - Difficult to make them fast
  - Lexer/parser **generators** transform the grammar into efficient lexers/parsers
    - Python: the PLY library
    - C/C++: (F)Lex and Yacc (or Bison)
    - Java: Antlr

## Step 1: Lexical Scanning

- Source code contains many **kinds** of things:  
variables, numbers, operators, keywords, whitespace,  
comments, punctuation, file inclusion, ...
- When defining a grammar, it makes sense to categorize these kinds of things into **tokens**
  - By convention, tokens are written in UPPERCASE
  - All variables are the same token: IDENT ("identifier")
  - Each keyword is a separate token: PRINT, IF, WHILE, ...
  - Operators and punctuations are also separate tokens: PLUS, STAR, SEMICOLON, LPAREN, RPAREN, ...
  - Usually **not** tokens: whitespace, comments
- A **lexical scanner (lexer)** converts source code into tokens

# A Simple Scanner

```
1  import re
2
3  ident_re = re.compile(r'[A-Za-z_][A-Za-z0-9_]*')
4  wsp_re   = re.compile(r'[\t\f\v\r\n]+')
5
6  def lex(source, pos=0):
7      while pos < len(source):
8          # using Python 3.8+ syntax:
9          if (match := wsp_re.match(source, pos)):
10             pos += len(match.group(0)) # skip whitespace
11         elif (match := ident_re.match(source, pos)):
12             ident = match.group(0)
13             yield ('IDENT', ident, pos)
14             pos += len(ident)
15         else:
16             print(f'Unknown character at offset {pos}: {source[pos]}')
17             pos += 1 # skip it
```

```
>>> print(*lex('mary had a little lambda'), sep='\n')
('IDENT', 'mary', 0)
('IDENT', 'had', 5)
('IDENT', 'a', 9)
('IDENT', 'little', 11)
('IDENT', 'lambda', 18)
>>>
```

# Adding Numbers

```
1  # ...
2  number_re = re.compile(r'[0-9]+')
3
4  def lex(source, pos=0):
5      while pos < len(source):
6          # ...
7          elif (match := number_re.match(source, pos)):
8              numlit = match.group(0)
9              yield ('NUMBER', int(numlit), pos)
10             pos += len(numlit)
11         # ...
```

```
>>> print(*lex('mary had 3 lambdas'), sep='\n')
('IDENT', 'mary', 0)
('IDENT', 'had', 5)
('NUMBER', 3, 9)
('IDENT', 'lambdas', 11)
>>>
```



# Handling Line-Comments

```
1  # ...
2  comment_re = re.compile(r'//.*\n?')
3
4  def lex(source, pos=0):
5      while pos < len(source):
6          # ...
7          elif (match := comment_re.match(source, pos)):
8              pos += len(match.group(0))
9          # ...
```

```
>>> print(*lex('mary had 3 // snow white\nlambdas'), sep='\n')
('IDENT', 'mary', 0)
('IDENT', 'had', 5)
('NUMBER', 3, 9)
('IDENT', 'lambdas', 25)
>>>
```

# Tracking Line Numbers

```
1  # ...
2  wsp_re = re.compile(r'[\t\f\v\r]+') # r'\n' is gone!
3  nl_re = re.compile(r'\n')
4
5  def lex(source, line=1, pos=0):
6      while pos < len(source):
7          # ...
8          elif nl_re.match(source, pos):
9              pos += 1
10             line += 1
11         # ...
12         elif (match := ident_re.match(source, pos)):
13             ident = match.group(0)
14             yield ('IDENT', ident, line, pos)
15             pos += len(ident)
16         # ...
```

```
>>> print(*lex('// a pythonic rhyme\nmary had 3 // snow white\nlambdas'), sep='\n')
('IDENT', 'mary', 2, 20)
('IDENT', 'had', 2, 25)
('NUMBER', 3, 2, 29)
('IDENT', 'lambdas', 3, 45)
>>>
```

# Operators and Punctuation

```
1  # ...
2  lparen_re    = re.compile(r'\(')
3  rparen_re    = re.compile(r'\)')
4  plus_re      = re.compile(r'\+')
5  minus_re     = re.compile(r'\-')
6  semicolon_re = re.compile(r';')
7
8  def lex(source, line=1, pos=0):
9      while pos < len(source):
10         # ...
11         elif lparen_re.match(source, pos):
12             pos += 1
13             yield ('LPAREN', None, line, pos)
14         elif rparen_re.match(source, pos):
15             pos += 1
16             yield ('RPAREN', None, line, pos)
17         elif plus_re.match(source, pos):
18             pos += 1
19             yield ('PLUS', None, line, pos)
20         elif minus_re.match(source, pos):
21             pos += 1
22             yield ('MINUS', None, line, pos)
23         elif semicolon_re.match(source, pos):
24             pos += 1
25             yield ('SEMICOLON', None, line, pos)
26         # ...
```

# Handling Keywords

```
1  # ...
2  reserved = {
3      'print': 'PRINT',
4      'while': 'WHILE',
5  }
6
7  def lex(source, line=1, pos=0):
8      while pos < len(source):
9          # ...
10         elif (match := ident_re.match(source, pos)):
11             ident = match.group(0)
12             yield (reserved.get(ident, 'IDENT'), ident, line, pos)
13             pos += len(ident)
14         # ...
```

```
>>> print(*lex('print(x + y);'), sep='\n')
('PRINT', 'print', 1, 0)
('LPAREN', None, 1, 6)
('IDENT', 'x', 1, 6)
('PLUS', None, 1, 9)
('IDENT', 'y', 1, 10)
('RPAREN', None, 1, 12)
('SEMICOLON', None, 1, 13)
>>>
```

# Using PLY/Lex To Streamline Lexical Scanners

In a file called, say, `scanner.py`:

```
1  import ply.lex as lex
2
3  reserved = {'print': 'PRINT', 'while': 'WHILE'}
4
5  # The 'tokens' tuple must be present and list all the valid tokens
6  tokens = (
7      'PLUS', 'MINUS', 'SEMICOLON', 'LPAREN', 'RPAREN', 'IDENT', 'NUMBER'
8  ) + tuple(reserved.values())
9
10 # Regexp strings definitions beginning with 't_' define simple tokens
11 t_LPAREN = r'\('
12 t_RPAREN = r'\)'
13 t_PLUS = r'\+'
14 t_MINUS = r'\-'
15 t_SEMICOLON = ';'
16
17 # Functions beginning with 't_' define complex token processing code.
18 # The docstrings of the functions contain the regexp that is matched for the token
19 def t_IDENT(t):
20     r'[A-Za-z_][A-Za-z0-9_]*' # docstring contains the regexp
21     t.type = reserved.get(t.value, 'IDENT')
22     return t
23
24 def t_NUMBER(t):
25     r'[1-9][0-9]*'
26     t.value = int(t.value)
27     return t
```

# Using a PLY/Lex Scanner

Still in scanner.py:

```
1  # error handling with t_error()
2  def t_error(t):
3      print(f"Illegal character '{t.value[0]}' on line {t.lexer.lineno}")
4      t.lexer.skip(1) # skip one character
5
6  # characters to ignore
7  t_ignore = ' \t\f\v'
8
9  def t_newline(t):
10     r'\n'
11     t.lexer.lineno += 1
12     # no return, signifying ignore
13
14  lexer = lex.lex()
15  # This will use Python introspection (reflection) to find out all the
16  # 'tokens' and 't_stuff' in this module and create a suitable lexer from it
```

```
>>> from scanner import lexer
>>> lexer.input('print(x + y);')      # give it some input
>>> lexer.token()                     # get next token
LexToken(PRINT, 'print', 1, 0)
>>> lexer.token()
LexToken(LPAREN, '(', 1, 5)
>>> t = lexer.token()
>>> (t.type, t.value, t.lineno, t.lexpos)
('IDENT', 'x', 1, 6)
```

## Step 2: Parsing

- Using the token stream generated by the lexer, the parser will convert it into an **abstract syntax tree** (AST)
- The structure of the AST is given by a **grammar**, usually specified in Backus-Naur Form (BNF).
- You have already seen BNF grammars in CSE 206.
- Example:

```
 $\langle \text{expr} \rangle ::= \text{IDENT} \mid \text{NUMBER}$   
                   $\mid \langle \text{expr} \rangle \text{ PLUS } \langle \text{expr} \rangle \mid \langle \text{expr} \rangle \text{ MINUS } \langle \text{expr} \rangle$   
                   $\mid \text{LPAREN } \langle \text{expr} \rangle \text{ RPAREN}$ 
```

# Using PLY/Yacc to Build a LR Parser

In a file called, say, `parser.py`:

```
1  import ply.yacc as yacc
2  import bxast          # We will write this module later
3
4  from scanner import tokens
5
6  ## Every parser function is written with a 'p_' prefix.
7  ## The docstring of the function is the portion of the grammar it handles
8
9  def p_expr_ident(p):
10     """expr : IDENT"""
11     # p[0]    p[1]
12     #
13     # The parser function must define p[0] in terms of p[1], p[2], ...
14     p[0] = bxast.Variable(p[1])
15
16  def p_expr_number(p):
17     """expr : NUMBER"""
18     p[0] = bxast.Number(int(p[1]))
19
20  def p_expr_plus(p):
21     """expr : expr PLUS expr"""
22     # p[0]    p[1] p[2] p[3]
23     p[0] = BinopApp('PLUS', p[1], p[3])
24
25  def p_expr_minus(p):
26     """expr : expr MINUS expr"""
27     p[0] = BinopApp('MINUS', p[1], p[3])
28
29  def p_expr_parens(p):
30     """expr : LPAREN expr RPAREN"""
31     p[0] = p[2]
```



# Using the PLY/Yacc Parser

Still in `parser.py`:

```
1 parser = yacc.yacc(start='expr')
2 # Build the parser by introspection from all the 'p_' functions.
3 # The 'start' parameter determines the start symbol, i.e., the overall
4 # output of the parser
```

```
>>> from scanner import lexer
>>> lexer.input('(x + y)')      # give it some input
>>> from parser import parser
>>> p = parser.parse(lexer=lexer)
>>> p
<bxast.BinopApp object at 0x7f57550e5310>
>>> print(p.to_json())          # assuming you wrote such a function
[[ "<binop>", [[ "<var>", "x" ], [] ],
  [[ "PLUS" ], [] ], [[ "<var>", "y" ], [] ] ], [] ]
```

# Operator Precedence

- Some operators can be used in multiple modes.  
E.g.:  $-$  can be used *infix* (subtraction) or *prefix* (negation)
- Infix operators have an associativity.  
E.g.:  $(x - y + z)$  is  $((x - y) + z)$ , not  $(x - (y + z))$
- Most operators also have a **binding strength** or **precedence**.  
E.g.:  $(x + y * z)$  is  $(x + (y * z))$ , not  $((x + y) * z)$
- Resolving precedence by hand is a lot of fun, but a bit out of scope for CSC 3F002 EP, at least at present
  - Look up: *Pratt parsing* aka *precedence climbing*
  - Also see: *shunting yard algorithm*
  - May revisit later if the need arises
- For now: we will use the built-in precedence declaration of PLY/Yacc.

# Operator Precedence in PLY/Yacc Parser

Still in `parser.py`:

```
1  # To declare precedence, there needs to be a 'precedence' table (tuple)
2  # that is listed in order of increasing precedence
3  precedence = (
4      ('left', 'PLUS', 'MINUS'),          # left-assoc., low precedence
5      ('left', 'TIMES', 'DIV', 'MODULUS'), # left-assoc., medium precedence
6      ('right', 'UMINUS'),                 # right-assoc., high precedence
7  )
8
9  # ...
10
11 parser = yacc.yacc(start='expr')
```

```
>>> from scanner import lexer
>>> from parser import parser
>>> parser.parse('x + y * z', lexer=lexer)
<...> # BinopApp('PLUS', Variable('x'), BinopApp('TIMES', Variable('y'), Variable('z')))
>>> parser.parse('(x + y) * z', lexer=lexer)
<...> # BinopApp('TIMES', BinopApp('PLUS', Variable('x'), Variable('y')), Variable('z'))
>>>
```

## Sequences (and Separators)

- Core BNF has no built in support for sequences, but:
- To get a sequence of  $\langle \text{item} \rangle$ s, one can do:
  - Zero or more  $\langle \text{item} \rangle$ s as an auxiliary nonterminal  $\langle \text{items} \rangle$

$$\langle \text{items} \rangle ::= \epsilon \mid \langle \text{item} \rangle \langle \text{items} \rangle$$

- One or more  $\langle \text{item} \rangle$ s as an auxiliary nonterminal  $\langle \text{items1} \rangle$

$$\langle \text{items1} \rangle ::= \langle \text{item} \rangle \langle \text{items} \rangle$$

- Easy to modify to separate individual  $\langle \text{item} \rangle$ s by COMMA, say:

$$\begin{aligned} \langle \text{items-comma} \rangle &::= \epsilon \mid \langle \text{items-comma1} \rangle \\ \langle \text{items-comma1} \rangle &::= \langle \text{item} \rangle \mid \langle \text{item} \rangle \text{ COMMA } \langle \text{items-comma1} \rangle \end{aligned}$$

## Extended BNF

- Instead of tokens names (e.g. COMMA), use their textual form (e.g. `","`)
- Group items with parentheses – note: the parentheses are not part of the language being described
- Indicate zero-or-more repetitions with `*`  
One-or-more repetitions with `+`  
Zero-or-one occurrences with `?`

```
⟨program⟩ ::= "def main() {" ⟨vardecl⟩* ⟨stmt⟩* "}"
```

```
⟨procdef⟩ ::= "def" IDENT "(" (⟨param⟩ ("," ⟨param⟩)*)? ")" ... (~ lab 4)
```

## Step 3: Syntactic Analysis

- Grammars specified in BNF are **context free**, so:
  - Cannot detect if every variable has been declared
  - Cannot detect if no variable has multiple declarations
  - Similarly, cannot detect that every procedure has a unique definition
- The output of the parser is sometimes called a **parse tree** rather than an AST
- Syntactic checks are implemented by traversing the parse tree.
- For the body of a procedure:
  - Whenever a  $\langle \text{vardecl} \rangle$  is encountered, check that:
    - ① the variable has not already been declared
    - ② the initializer is a well formed expression
  - Whenever an  $\langle \text{expr} \rangle$  expression is encountered, check that:
    - ① Every variable has been declared earlier
  - More complex checks as the language grows...

## Lab 1: Preview

- You will write:
  - A lexer and parser for a fragment of BX using PLY
- Fragment of BX:
  - Will be a single function (main) program, straightline code, and **integer** vars

$\langle \text{program} \rangle ::= \text{"def main() \{ " } \langle \text{stmt} \rangle^* \text{" \} "}$

$\langle \text{stmt} \rangle ::= \langle \text{vardecl} \rangle \mid \langle \text{assign} \rangle \mid \langle \text{print} \rangle$

$\langle \text{vardecl} \rangle ::= \text{"var" IDENT "=" } \langle \text{expr} \rangle \text{" : int;"}$