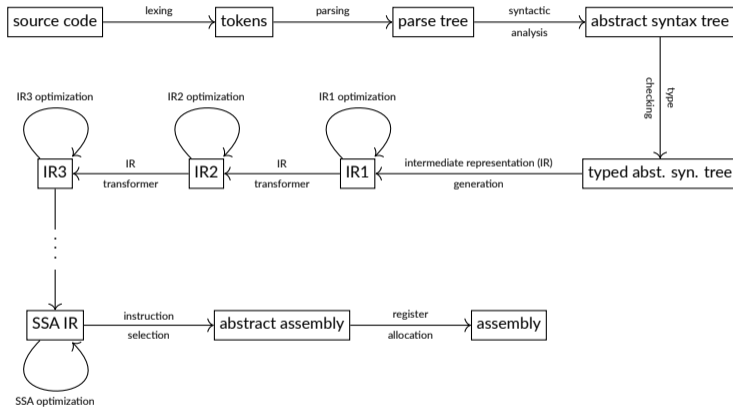# IR Generation
## (CSC 3F002 EP) Compilers

**Pierre-Yves Strub**
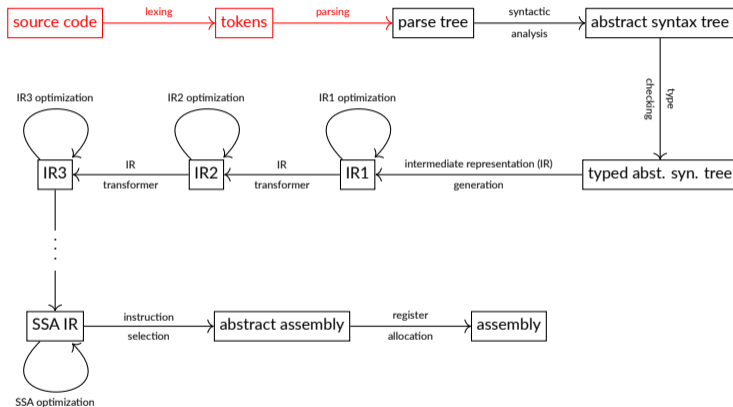
(Slide deck author: **Kaustuv Chaudhuri**)

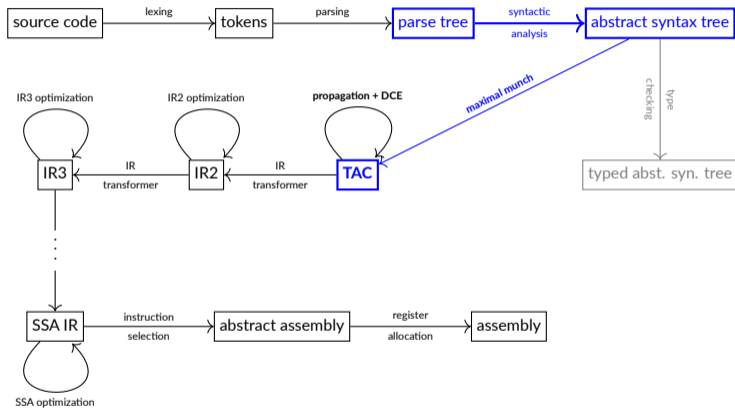2025-09-08 // Week 2

# Compiler Stages: Block Diagram

source code — lexing → tokens — parsing → parse tree — syntactic analysis → abstract syntax tree

type checking

IR3 optimization
IR2 optimization
IR1 optimization

IR3 ← IR transformer ← IR2 ← IR transformer ← IR1 ← intermediate representation (IR) generation ← typed abst. syn. tree

SSA IR — instruction selection → abstract assembly — register allocation → assembly

SSA optimization

# Last lecture

# Today

# Reminder: straightline BX

- A simplification of BX with:
  - Only one type of variables, signed 64-bit **int**s
  - No control structures
  - Only one function, `main()`
  - Only I/O: printing **int**s
- BX expressions:
  - Arithmetic (+, -, *, /, %)
  - Bitwise ops (~, &, |, ^)
  - Arithmetic shifting (<<, >>)
- BX statements:
  - Assignment: `x = expression;`
  - Print: `print(expression);`

# Reminder: BX Expressions and Assignments

```
x = 10;              // assign immediate

y = 2 * x;           // compute and assign
z = y / 2;           // compute and assign

w = z - x - y;       // compound expression

v = - w - - w;       // - is both binary and unary
                     // parsed as (- w) - (- w), i.e., 0

m = x + y * 2;       // precedence: * binds tighter than +
                     //    so, parsed as x + (y * 2)
n = (x + y) * 2;     // parentheses to force evaluation orders
```
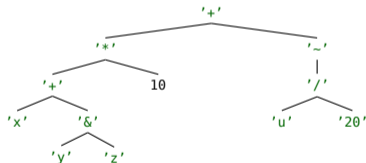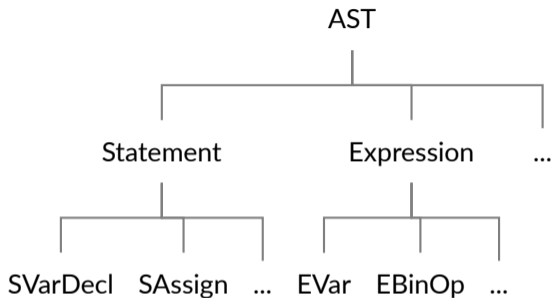
# Abstract Syntax Tree

- Programs are represented using an Parse Tree (PST) / Abstract Syntax Tree (AST)
- PST is what is given by the parser. AST is more abstract, can contain extra information (e.g. typing information). The line between PST/AST is thin and sometimes the parser outputs an AST directly.
- For example, ⟨expr⟩s are tree-like with arbitrary depth

```
                      '+'
            ┌──────────┴──────────┐
          '*'                    '~'
       ┌───┴───┐                  │
     '+'      10                 '/'
   ┌──┴──┐                    ┌───┴───┐
  'x'   '&'                  'u'    '20'
      ┌─┴─┐
    'y'  'z'
```

- Order of evaluation of an ⟨expr⟩ depends on traversal (preorder, postorder, …)
- In OOP, PST/AST are defined using a *class hierarchy*

# Example of an AST hierarchy



- Each class has several data attributes that store the relevement information. For example, the SVarDecl class (Variable Declaration) contains the variable name (a str), its type and its initializer (an Expression)

- This hierarchy is recursive. For example, the EBinOp (Binary Operator Applicaton) contains two Expression (the operands).

# AST definition in Python

```python
1   import abc
2   import dataclasses as dc
3
4   @dc.dataclass
5   class AST(abc.ABC):
6       pass
7
8   @dc.dataclass
9   class Statement(AST):
10      pass
11
12  @dc.dataclass
13  class Expression(AST):
14      pass
```

```python
1   @dc.dataclass
2   class SVarDecl(Statement):
3       name: str
4       type: str
5       init: Expression
6
7   @dc.dataclass
8   class SAssignment(Statement):
9       lvalue: str
10      rvalue: Expression
11
12  @dc.dataclass
13  class EVar(Expression):
14      name: str
15
16  @dc.dataclass
17  class EBinOp(Expression):
18      binop: str
19      left: Expression
20      right: Expression
```

# AST transformers

There are different ways to work with an AST. The ugly one:

```python
def do_something(ast : AST):
    if isinstance(ast, EVar):
        # Do something
    elif isinstance(ast, EBinOp):
        # Do something else
    # ...
```

- Error prone, hard to read, huge functions, ...
- Pass your way

# AST transformers

The OOP way:

```
1    @dc.dataclass
2    class AST(abc.ABC):
3      @abc.abstractmethod
4      def transformer1(self):
5        pass
6
7      @abc.abstractmethod
8      def transformer2(self, thearg):
9        pass
```

```
1    @dc.dataclass
2    class SAssignment:
3      lvalue: str
4      rvalue: Expression
5
6      def transformer1(self):
7        # Do something (1)
8
9      def transformer2(self, thearg):
10       # Do something (2)
```

- Classic representation in OOP, but...
- ...transformers' bodies are spread across multiple classes.
- it is also possible to use the visitor design pattern.

# AST transformers

Python has now a notion of pattern matching that plays well with the `dataclasses` module:

```
1    def transformer1_expr(e : Expression):
2      match e:
3        case EVar(name):
4          # Do something. Can use the 'name'
             variable
5        case EBinOp(_, left, right):
6          # Do something. '_' means that the
7          # first data attribute has been
8          # ignored (not bound)
```

```
1    def transformer1_stmt(s : Statement):
2      match e:
3        case SVarDecl(name = x, init = e):
4          # Data attributes can be matched
5          # using their names (more robust).
6        case SAssign(lvalue, rvalue):
7          # Do something.
```

- Close to algebraic pattern matching
- Widely used in functional programming
- Very effective when dealing with an AST
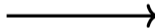
# AST Linearization

- Order of evaluation of an ⟨expr⟩ depends on traversal (preorder, postorder, …)

---

- Instructions are shallow –
  operands can be registers, memory addresses, or numbers (depth 0)
- Instructions are sequence-like
- Order of evaluation is explicit – top-to-bottom

---

- Instruction Generation: going from tree-like to list-like structures

# BX (AST) to TAC: Example

BX

```
def main() {
  // variable declarations
  var x = 0 : int;
  var y = 0 : int;

  // straightline code
  x = 10;
  y = 2 * x;
  x = y * y / 2;
  print(9 * x * x + 3 * x - 8);
}
```

TAC

```
proc @main:
  %0 = const 0;
  %1 = const 0;
  %0 = const 10;
  %2 = const 2;
  %1 = mul %2, %0;
  %4 = mul %1, %1;
  %7 = const 2;
  %0 = div %4, %7;
  %12 = const 9;
  %11 = mul %12, %0;
  %10 = mul %11, %0;
  %16 = const 3;
  %15 = mul %16, %0;
  %9 = add %10, %15;
  %18 = const 8;
  %8 = sub %9, %18;
  print %8;
```

# Three Address Code (TAC)
(our first intermediate representation, IR)

- A (low-level) language that is closer to assembly than BX

- Temporaries:
  - Think: *abstract* register
  - Infinite supply – can use arbitrarily many
  - Anonymous temporaries: %0, %1, %42, ...                    (*we will use these mostly*)
  - Named temporaries: %x, %foo, ...

- Instructions:
  - Unique opcode
  - Zero, one, or two argument temporaries                    (*with one exception:* const)
  - Zero or one result temporary that is written
  - Prototype syntax:     %42 = add %10, %3
                          result  opcode   args

# TAC Instructions Reference

(%1, %2 = examples of reads;   %0 = example of write;   42 = example of immediate)

| Instruction | Description |
|---|---|
| %0 = const 42; | Set temporary %0 to the value 42. |
| %0 = copy %1; | Copy the value of %1 to %0. |
| %0 = binop %1, %2; | Compute the value of binary operator binop $\in$ {add, sub, mul, div, mod, and, or, xor, shl, shr} applied to %1 and %2 and store in %0. |
| %0 = unop %1; | Compute the value of a unary operator unop $\in$ {neg, not} applied to %1 and store in %0. |
| print %1; | Print the value of %1 to the standard output. |

# Algorithm 1: Top-Down Maximal Munch (TMM)

High-level overview

- Start from the root of the AST and do a postorder (i.e., children before root) traversal
- Allocate fresh temporaries for each child node (if any) and generate instructions (recursively!) to output the value fo that child to that temporary
- Finally, generate the instruction for the root using the temporaries of the children

# Algorithm 1: Top-Down Maximal Munch (TMM)

High-level overview

- Start from the root of the AST and do a postorder (i.e., children before root) traversal
- Allocate fresh temporaries for each child node (if any) and generate instructions (recursively!) to output the value fo that child to that temporary
- Finally, generate the instruction for the root using the temporaries of the children

We will now flesh this out.

# Top-Down Maximal Munch: Expressions

## code $(e, x)$

- Generate code for the expression $e$.
- Storing the result in the temporary $x$.
- Returns: list of instructions.

| $e$ | code $(e, x)$ | proviso |
|---|---|---|
| 42 | $[x =$ `const 42;`$]$ | — |
| $y$ | $[x =$ `copy y;`$]$ | — |
| $e_1 + e_2$ | code $(e_1, y)$ +<br>code $(e_2, z)$ +<br>$[x =$ `add y, z;`$]$ | $y, z$   fresh |
| $- e_1$ | code $(e_1, y)$ +<br>$[x =$ `neg y;`$]$ | $y$   fresh |

# Top-Down Maximal Munch: Statements

### code $(s)$

- Generate code for the statement $s$.
- Returns: list of instructions.

| $s$ | code $(s)$ | proviso |
|-----|-----------|---------|
| $x = e;$ | code $(e, x)$ | — |
| `print(e);` | code $(e, x)$ + [`print x;`] | $x$ fresh |

# Top-Down Maximal Munch — Example

TAC

```
%0 = const 10;
%2 = const 2;
%3 = copy %0;
%1 = mul %2, %3;
%5 = copy %1;
%6 = copy %1;
%4 = mul %5, %6;
%7 = const 2;
%0 = div %4, %7;
%12 = const 9;
%13 = copy %0;
%11 = mul %12, %13;
%14 = copy %0;
%10 = mul %11, %14;
%16 = const 3;
%17 = copy %0;
%15 = mul %16, %17;
%9 = add %10, %15;
%18 = const 8;
%8 = sub %9, %18;
print %8;
```

BX

```
x = 10;
y = 2 * x;
x = y * y / 2;
print(9 * x * x + 3 * x - 8);
```

TMM →

# Top-Down Maximal Munch — Example

TAC

```
%0 = const 10;
%2 = const 2;
%3 = copy %0;
%1 = mul %2, %3;
%5 = copy %1;
%6 = copy %1;
%4 = mul %5, %6;
%7 = const 2;
%0 = div %4, %7;
%12 = const 9;
%13 = copy %0;
%11 = mul %12, %13;
%14 = copy %0;
%10 = mul %11, %14;
%16 = const 3;
%17 = copy %0;
%15 = mul %16, %17;
%9 = add %10, %15;
%18 = const 8;
%8 = sub %9, %18;
print %8;
```

BX

```
x = 10;
y = 2 * x;
x = y * y / 2;
print(9 * x * x + 3 * x - 8);
```

TMM $\longrightarrow$

TMM generates many (often redundant) copies!

# Algorithm 2: Bottom-Up Maximal Munch (BMM)

Key idea: make the destination temporary an output instead of an input.

---

$(x, L) = \text{code}\,(e)$

- Generate code ($L$) for the expression $e$.
- Storing the result in the temporary $x$.
- Returns: both $x$ and $L$

---

| $e$ | code $(e)$ | proviso |
|---|---|---|
| 42 | $x, [x = \texttt{const } 42;]$ | $x$   fresh |
| $y$ | $y, [\,]$ | — |
| $e_1 + e_2$ | $x, L_1 + L_2 + [x = \texttt{add } y, z;]$ <br> where $(y, L_1) = \text{code}\,(e_1)$ <br> and $(z, L_2) = \text{code}\,(e_2)$ | $x$   fresh |
| - $e_1$ | $x, L_1 + [x = \texttt{neg } y;]$ <br> where $(y, L_1) = \text{code}\,(e_1)$ | $x$   fresh |

# Bottom-Up Maximal Munch: Statements

## code $(s) = L$

- Generate code for the statement $s$
- Returns: list $L$ of instructions

| $s$ | code $(s)$ | proviso |
|---|---|---|
| $x = e;$ | $L + [x = \text{copy } y;]$ <br> where $(y, L) = \text{code}(e)$ | — |
| print$(e);$ | $L + [\text{print } x;]$ <br> where $(x, L) = \text{code}(e)$ | — |

# Both Maximal Munches — Example

TAC

```
%0 = const 10;
%2 = const 2;
%3 = copy %0;
%1 = mul %2, %3;
%5 = copy %1;
%6 = copy %1;
%4 = mul %5, %6;
%7 = const 2;
%0 = div %4, %7;
%12 = const 9;
%13 = copy %0;
%11 = mul %12, %13;
%14 = copy %0;
%10 = mul %11, %14;
%16 = const 3;
%17 = copy %0;
%15 = mul %16, %17;
%9 = add %10, %15;
%18 = const 8;
%8 = sub %9, %18;
print %8;
```

← TMM

BX

```
x = 10;
y = 2 * x;
x = y * y / 2;
print(9 * x * x + 3 * x -
    8);
```

BMM →

TAC

```
%0 = const 10;
%1 = copy %0;
%2 = const 2;
%3 = mul %2, %1;
%4 = copy %3;
%5 = mul %4, %4;
%6 = const 2;
%7 = div %5, %6;
%1 = copy %7;
%8 = const 9;
%9 = mul %8, %1;
%10 = mul %9, %1;
%11 = const 3;
%12 = mul %11, %1;
%13 = add %10, %12;
%14 = const 8;
%15 = sub %13, %14;
print %15;
```

# Summary of Lecture 1

- The frontend produces an parse/abstract syntax tree (AST)
- Parse/abstract syntax tree are represented using tree-like structures. The concrete representation depends on the programming language
  - Class hierarchy,
  - Algebraic datatypes,
  - ...
- AST is linearized into instructions with maximal munch
  - Top-Down: `copy` on reads
  - Bottom-Up: `copy` on writes