# CS 536 Announcements for Thursday, February 13, 2025

**Programming Assignment 2** – due Tuesday, February 18

**Homework 2** – is available

**Last Time**
- why regular expressions aren't enough
- CFGs
    - formal definition
    - examples
    - language defined by a CFG
- parse trees
- Makefiles
- ambiguous grammars
- grammars for expressions
    - precedence
    - associativity
- grammars for lists

**Today**
- syntax-directed translation
- abstract syntax trees
- implementing ASTs    *end of Midterm 1 material*

# CFG review

prog  → BEGIN stmts END

stmts → stmts SEMICOLON stmt
    | stmt

stmt  → ID ASSIGN expr

expr  → expr PLUS term
    | term

term  → term TIMES factor
    | factor

factor → expon POW factor
    | expon

expon → INTLIT
    | LPAREN expr RPAREN

*precedence low*

*high*

# Overview of CFGs

Scanner $\rightarrow$ stream of tokens $\rightarrow$ parser $\rightarrow$ AST

**CFGs for language definition**

- the CFGs we've discussed can generate/define languages of valid strings

Start by building parse tree
& end with some valid string $w \in L(G)$

**CFGs for language recognition**

Start with string $w$

& end with yes/no answer depending on whether $w \in L(G)$

**CFGs for parsing**

Start with string $w$

& end with parse tree for $w$ if $w \in L(G)$
  ↳ generally use AST instead of parse tree

— need to translate sequence of tokens ($w$)

# Syntax-directed translation (SDT)

=  translating from a sequence of tokens into a sequence of actions/other form, based on underlying syntax

Could be: AST, value, type, etc.

**To define a syntax-directed translation**

Augment CFG with ***translation rules*** (at most 1 rule per production)

↳ LHS → RHS

- define translation of LHS non-terminal as a function of

  - Constants

  - translations of RHS non-terminals

  - values of token (terminals) on RHS

**To translate a sequence of tokens using SDT** ✸

- build parse tree

- use translation rules to compute translation of each non-terminal in parse tree
  bottom-up ⟵ handle children of node before node

- translation of sequence of tokens is the translation of the parse tree's root non-terminal (ie, start symbol)

The **type** of the translation can be anything: numeric, string, set, tree, ...

✸ Note: above is how to understand the translation, not how a compiler actually does it

# Example: grammar for language of binary numbers

CFG

b → 0
b → 1
b → b 0
b → b 1

translation rules

b.trans = 0
b.trans = 1
$b_1$.trans = $b_2$.trans * 2
$b_1$.trans = $b_2$.trans * 2 + 1

SDT to compute
the decimal equivalent
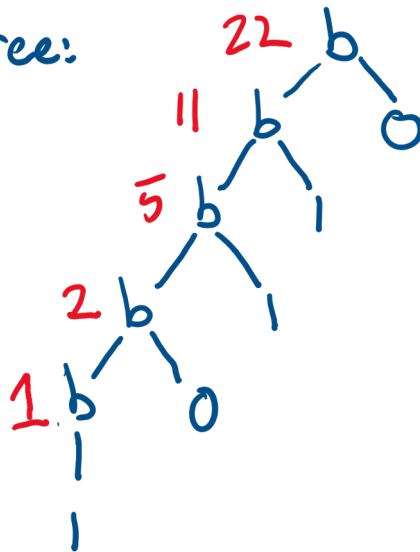of a binary number

Example: in put string  10110

Parse tree:

steps

④
③
②
①

22 b
11 b    0
5 b    1
2 b    1
1 b    0
b
1

Translation is 22

# Example: grammar for language of variable declarations

CFG

declList → ε

       | decl declList$_2$

decl → type ID ;

type → INT

       | BOOL

Translation rules

string
concatenation

declList.trans = " "

declList$_1$.trans = decl.trans + " " + declList$_2$.trans

decl.trans = ID.value

Write a syntax-directed translation for the CFG given above so that the translation of a sequence of tokens is a string containing the ID's that have been declared.
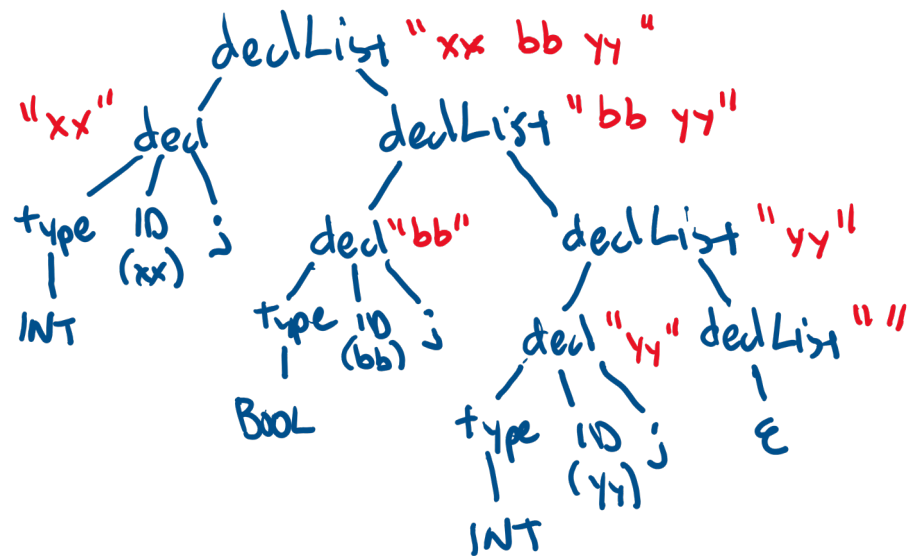
Example input

int xx;
bool bb;
int yy;

Output

"bb xx yy"
(in any order)

Parse tree

Translation is

"xx bb yy"

declList "xx bb yy"

"xx" decl
    type ID ;
     (xx)
    INT

declList "bb yy"

decl "bb"
  type ID ;
    (bb)
  BOOL

declList "yy"

decl "yy"
  type ID ;
   (yy)
  INT

declList " "
  ε

# Example: grammar for language of variable declarations

CFG

declList → ε

     |    decl declList

decl → type ID ;

type → INT

     |    BOOL

Translation rules

$declList.trans = ""$

$declList_1.trans = decl.trans + " " + declList_2.trans$

$decl.trans = type.trans ? ID.value : ""$   ✶

$type.trans = true$

$type.trans = false$

Modify the previous syntax-directed translation so that only declarations of type `int` are added to the output string.

Example input
int xx;
bool bb;
int yy;

Output
"xx yy"
(in any order)

Note:
1) different non-terms can have different types as their translation
2) translation rules can be conditional

✶ $x = a ? b : c;$ equiv to if (a)
            x = b;
        else
            x = c;

# SDT for parsing

Previous examples showed SDT process assigning different types to the translation

- translate tokenized stream to an integer value
- translate tokenized stream to a string

For parsing, we'll need to translate a tokenized stream to an **abstract-syntax tree (AST)**
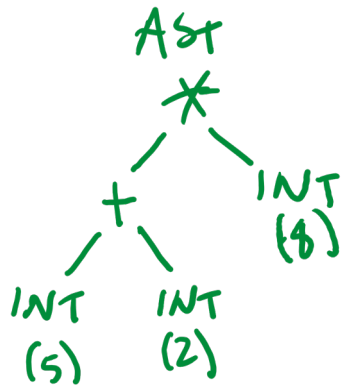
# Abstract syntax trees

**AST** = condensed form of parse tree

- operators are internal nodes (not leaves)
- chains of productions are collapsed
- lists are flattened
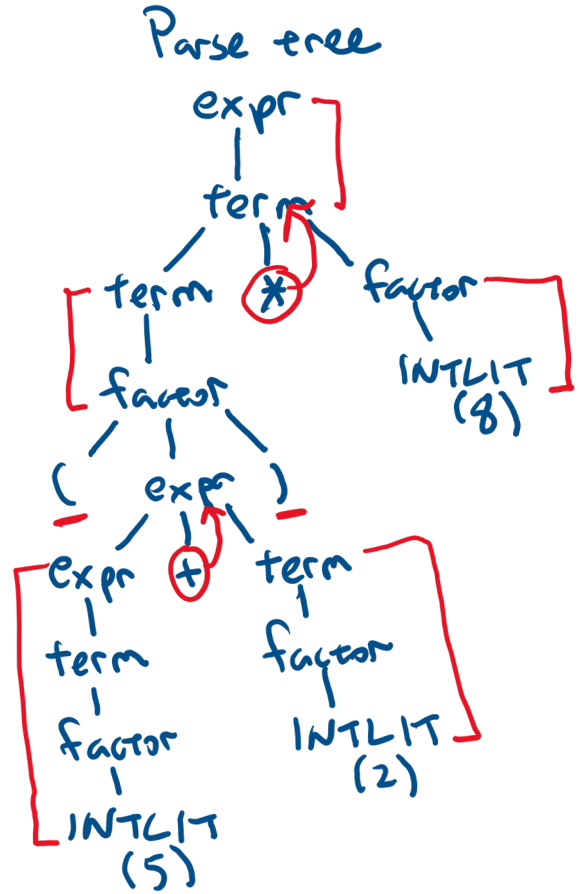- syntactic details are omitted
  - ↳ eg   ; parens

# AST Example

CFG

expr → expr PLUS term
| term

term → term TIMES factor
| factor

factor → INTLIT
| LPAREN expr RPAREN

**AST**



$(5+2) * 8$

**Parse tree**



Résume
at 7:30

# SDT review

**SDT** = translating from a <u>sequence of tokens</u> into a <u>sequence of actions/other form</u>, based on <u>underlying syntax</u>

**To define a syntax-directed translation**

- augment CFG with **_translation rules_**  lhs → rhs    *(contains terms, non-terms, ε)*
  - define translation of LHS non-terminal as a function of:
    - constants  2, " "
    - translations of RHS non-terminals  rhs.trans
    - values of terminals (tokens) on RHS  TOKEN.value

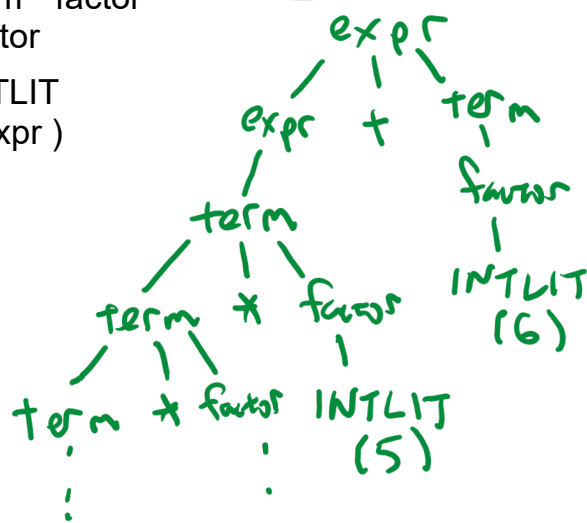**To translate a sequence of tokens using SDT (conceptually)**

- build parse tree
- use translation rules to compute translation of each non-terminal (bottom-up)
- translation of sequence of tokens = translation of parse tree's root non-terminal

For parsing, we'll need to translate tokenized stream to **abstract-syntax tree (AST)**
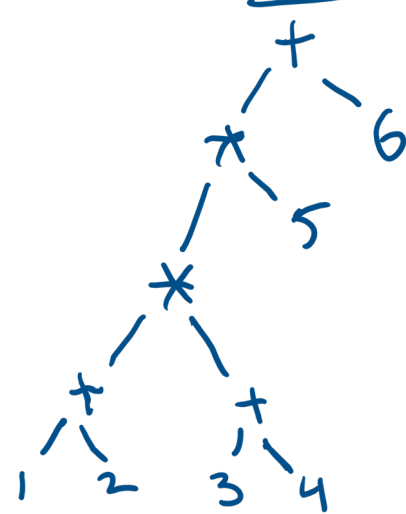
**Example** $(1+2) * (3+4) * 5 + 6$

expr → expr + term
| term
term → term * factor
| factor
factor → INTLIT
| ( expr )

Parse tree



AST



translation rule:

$expr_1.trans = MkPlusNode(expr_2.trans, term.trans)$

## AST for parsing

We've been showing the translation in two steps:

token stream → parse tree → AST then throw away parse tree

In practice we'll do

token stream → AST

**Why have an AST?**

– captures essential structure
– easier to work with

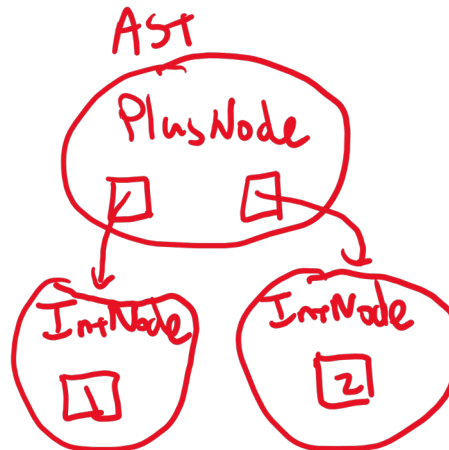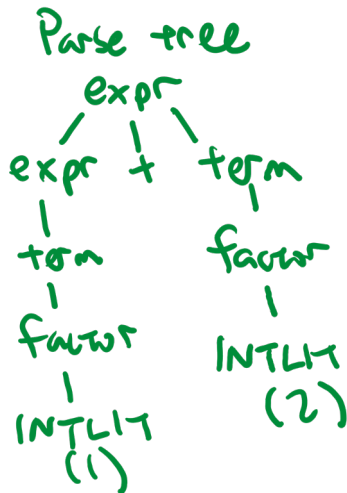# AST implementation

$expr \rightarrow expr + term$        $expr_1.trans = Mk PlusNode(expr_2.trans, term.trans)$

Define a class for each kind of AST node

Create a new node object in some rules

- new node object is the value of LHS.trans
- fields of node object come from translations of RHS non-terminals

Given 1+2

Parse tree

```
        expr
       / |  \
   expr  +   term
    |         |
   term     factor
    |         |
  factor    INTLIT
    |        (2)
  INTLIT
   (1)
```

AST

PlusNode
 □   □
 ↓    ↓
IntNode  IntNode
  1        2

```
class PlusNode {
    IntNode left;
    IntNode right;
}
class IntNode {
    int value;
}
```

Need class hierarchy
& make subclasses of ExpNode

```
class PlusNode extends ExpNode {
    ExpNode left;
    ExpNode right;      → put into ExpNode
}
class IntNode extends Expnode {
    int value;
}
```

# Translation rules to build ASTs for expressions

CFG                                         Translation rules

$expr \rightarrow expr + term$         $expr_1.trans$ = new PlusNode($expr_2.trans$, $term.trans$)

$\quad\quad |\quad term$                    $expr.trans$ = term.trans

$term \rightarrow term * factor$        $term_1.trans$ = new TimesNode($term_2.trans$, $factor.trans$)

$\quad\quad |\quad factor$                 $term.trans$ = factor.trans

$factor \rightarrow INTLIT$             $factor.trans$ = new IntNode(INTLIT.value)

$\quad\quad |\quad ( expr )$               $factor.trans$ = expr.trans
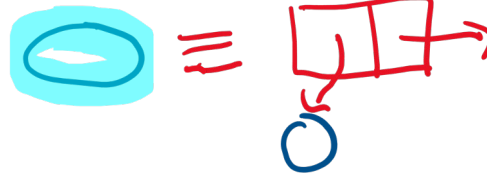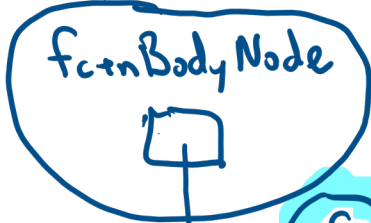
Example 1+2

**Example**
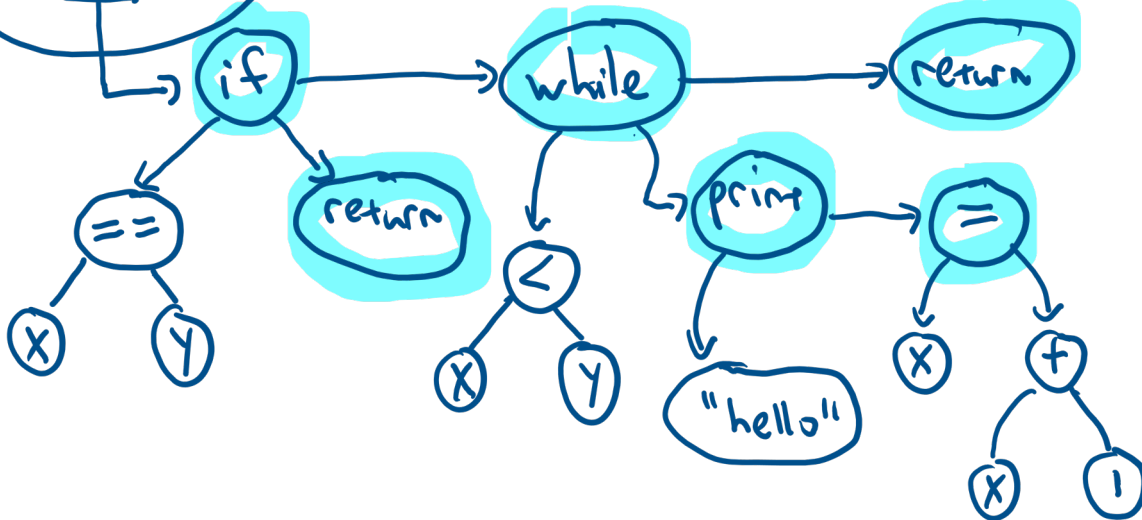
```
void foo(int x, int y) {
    if (x == y) {
        return;
    }
    while (x < y) {
        cout << "hello";
        x = x + 1;
    }
    return;
}
```
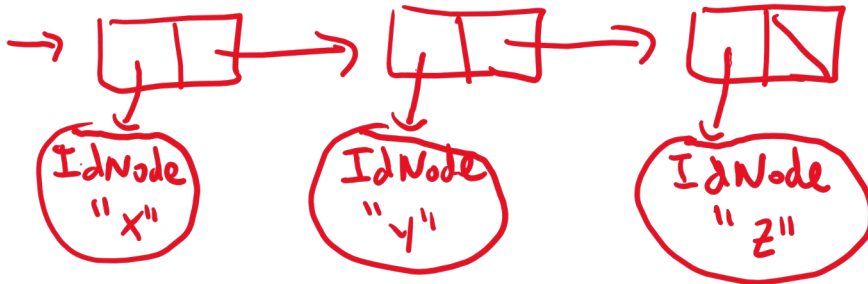
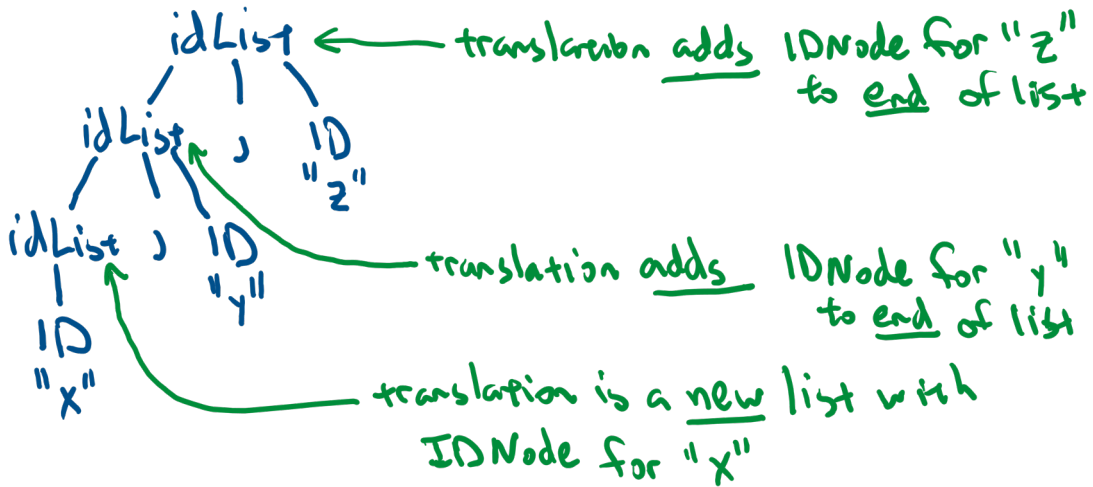AST for function body



fctnBody Node

ie linked list

<u>CFG</u>

idList → idList COMMA ID
| ID

Want AST to be



Parse tree

idList ← translation adds IDNode for "z" to end of list

idList , ID "z"

idList , ID "y" ← translation adds IDNode for "y" to end of list

idList , ID "x"

ID "x" ← translation is a new list with IDNode for "x"
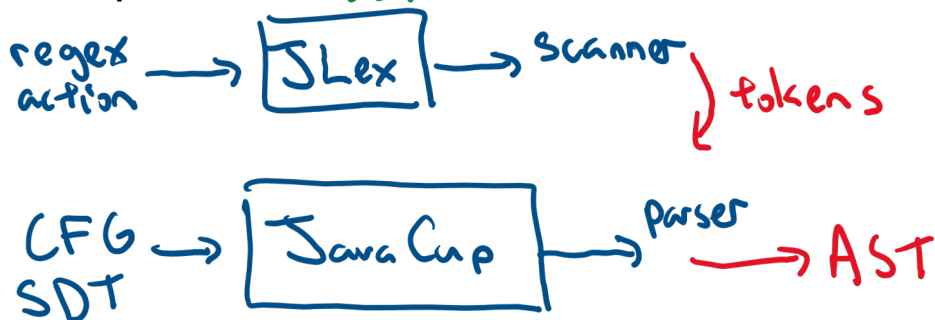
Input x, y, z

# The bigger picture

**Scanner**

- **Language abstraction**: regular expressions

- **Output**: token stream

- **Tool**: JLex

- **Implementation**: interpret DFA using table (for δ),
  recording most_recent_accepted_position & most_recent_token

**Parser**

- **Language abstraction**: CFG

- **Output**: AST (by way of a syntax-directed translation)

- **Tool**: Java CUP ← next time

- **Implementation**: ??? ← next time

regex action ⟶ JLex ⟶ Scanner ⟩ tokens

CFG SDT ⟶ Java Cup ⟶ Parser ⟶ AST

**Next Time**
- Java CUP
- approaches to parsing
- bottom-up parsing
- CFG transformations
- CYK algorithm