CSE110A: Compilers

April 28, 2028

Topics:

- Top down parsing
 - Lookahead sets
 - Recursive descent parsers
- Symbol Tables



Announcements

- HW 2 is out
 - You should have everything you need for it after today
 - Due in 6 days (May 4)
 - Plenty of time for help
- We are working on grading HW 1
- Midterm will be given on May 8
 - Taken during class
 - Study material is homeworks, slides, and book readings
 - 3 pages of notes (front and back, handwritten or typed)

Quiz

We'll revisit a few of the questions from the last quiz

Quiz

To prepare a grammar for a top-down parser, you must ensure that there is no recursion, except in the right-most element of any production rule.

⊖ True

 \bigcirc False

What is the issue with left recursion?

```
root = start symbol;
focus = root;
push(None);
                                  What could a demonic
to match = s.token();
                                  choice do?
while (true):
  if (focus is a nonterminal)
    pick next rule (A ::= B1,B2,B3...BN);
    push(BN... B3, B2);
    focus = B1
  else if (focus == to match)
    to match = s.token()
    focus = pop()
```

```
else if (to_match == None and focus == None)
Accept
```

Variable	Value
focus	
to_match	
s.istring	
stack	

```
1: Expr ::= Expr '+' ID
2: | ID
```

Can we derive the string a

Expanded Rule	Sentential Form
start	Expr

```
root = start symbol;
focus = root;
push(None);
                                  What could a demonic
to match = s.token();
                                  choice do?
while (true):
  if (focus is a nonterminal)
    pick next rule (A ::= B1,B2,B3...BN);
    push(BN... B3, B2);
    focus = B1
  else if (focus == to match)
    to match = s.token()
    focus = pop()
```

```
else if (to_match == None and focus == None)
Accept
```

Variable	Value
focus	
to_match	
s.istring	
stack	

```
1: Expr ::= Expr '+' ID
2: | ID
```

Can we derive the string a

Expanded Rule	Sentential Form
start	Expr
1	Expr + ID
1	Expr + ID + ID
1	Expr + ID + ID + ID

infinite recursion

Eliminating direct left recursion

A and B can be any sequence of non-terminals and terminals

Eliminating direct left recursion

```
1: Expr ::= Expr Op Unit
2:  | Unit
3: Unit ::= '(' Expr ')'
4:  | ID
5: Op ::= '+'
6:  | '*'
```

Lets do this one as an example:

Eliminating direct left recursion

A = Op Unit B = Unit

1:	Expr	::=	Uni	it Exp	or2
2:	Expr2	::=	Op	Unit	Expr2
3:			11 11		

Lets do this one as an example:

How about indirect left recursion?

Identify indirect left left recursion

 $Expr_base \rightarrow_{lhs} Expr_op \rightarrow_{lhs} Expr_base$

How about indirect left recursion?

Identify indirect left left recursion

 $Expr_base \rightarrow_{lhs} Expr_op \rightarrow_{lhs} Expr_base$

inline indirect non-terminal

It is always possible to eliminate left recursion

Quiz

It is only possible to write a top-down parser if you can determine exactly which production rule to apply at each step.

⊖ True

 \bigcirc False

```
else if (to_match == None and focus == None)
    Accept
```

```
else if (focus == to_match)
  to_match = s.token()
  focus = pop()
```

else if (we have a cached state)
 backtrack();

else

```
parser_error()
```

```
1: Expr ::= ID Expr2
2: Expr2 ::= '+' ID Expr2
| ""
```

Can we match: "a"?

Expanded Rule	Sentential Form	
start	Expr	

Backtracking gets complicated...

- Do we need to backtrack?
 - In the general case, **yes**
 - In many useful cases, no

```
root = start symbol;
focus = root;
push(None);
to_match = s.token();
while (true):
    if (focus is a nonterminal)
    pick next rule (A ::= B1,B2,B3...BN);
    if B1 == "": focus=pop(); continue;
    push(BN... B3, B2);
    focus = B1
```

```
else if (focus == to_match)
  to_match = s.token()
  focus = pop()
```

```
else if (to_match == None and focus == None)
Accept
```

Variable	Value
focus	Expr2
to_match	None
s.istring	un
stack	None

1: Expr ::= ID Expr2 2: Expr2 ::= '+' Expr2 3: | ""

Can we match: "a"?

Expanded Rule	Sentential Form	
start	Expr	
1	ID <mark>Expr2</mark>	

The First Set

For each production choice, find the set of tokens that each production can start with

			First s	sets:
1:	Expr	::= Unit Expr2	1: {}	
2:	Expr2	::= Op Unit Expr2	2: {}	
3:		<i>и п</i>	3: {}	
4:	Unit	::= '(' Expr ')'	4: {}	
5:		ID	5: {}	
6:	Op	::= '+'	6 : {}	
7:		/ * /	7: {}	

The First Set

For each production choice, find the set of tokens that each production can start with

```
First sets:
1: { '(', ID}
2: { '+', '*' }
3: { "" }
4: { '(' }
5: { ID }
6: { '+' }
7: { '*' }
```

We can use first sets to decide which rule to pick!

```
root = start symbol;
focus = root;
push(None);
to match = s.token();
```

while (true):

```
if (focus is a nonterminal)
  pick next rule (A ::= B1,B2,B3...BN);
  push(BN... B3, B2);
  focus = B1
```

```
else if (focus == to_match)
  to_match = s.token()
  focus = pop()
```

```
else if (to_match == None and focus == None)
Accept
```

Variable	Value
focus	
<mark>to_match</mark>	
s.istring	
stack	

```
First sets:
1: {'(', ID}
2: {'+', '*'}
3: {'''}
4: {'('}
5: {ID}
6: {'+'}
7: {'*'}
```

We simply use to_match and compare it to the first sets for each choice

For example, Op and Unit

Quiz

In many cases, a top-down parser requires the grammar to be re-written. Write a few sentences about why this might be an issue when developing a compiler and how the issues might be addressed.

Class discussion

New material

- The Follow set
- The First+ set
- Recursive descent parser

The Follow Set

Rules with "" in their First set need special attention

			First sets:	Follow sets:
1:	Expr	::= Unit Expr2	1: {'(', ID}	1: NA
2:	Expr2	::= Op Unit Expr2	2: { '+', '*' }	2: NA
3:	_		3: { " " }	3: {}
4:	Unit	::= '(' Expr ')'	4: { ' (' }	4: NA
5:		ID	5: {ID}	5: NA
6:	Op	::= '+'	6: { '+' }	6: NA
7:		/ * /	7: { ' * ' }	7: NA

We need to find the tokens that any string that follows the production can start with.

The Follow Set

Rules with "" in their First set need special attention

			First sets:	Follow sets:
1:	Expr	::= Unit Expr2	1: {'(', ID}	1: NA
2:	Expr2	::= Op Unit Expr2	2: { '+', '*' }	2: NA
3:	_		3: {""}	3: {None, ')'}
4:	Unit	::= '(' Expr ')'	4: { ' (' }	4: NA
5:		ID	5: {ID}	5: NA
6:	Op	::= '+'	6: { '+' }	6: NA
7:		/ * /	7: { ' * ' }	7: NA

We need to find the tokens that any string that follows the production can start with.

The First+ Set

The First+ set is the combination of First and Follow sets

			First sets:	Follow sets:	First+ sets:
1:	Expr	::= Unit Expr2	1: {'(', ID}	1: NA	1: {'(', ID}
2:	Expr2	::= Op Unit Expr2	2: { '+', '*' }	2: NA	2: { '+', '*' }
3:			3: { " " }	3: {None, ')'}	3: {None, ')'}
4:	Unit	::= '(' Expr ')'	4: { ' (' }	4: NA	4: { (' }
5:		ID	5: {ID}	5: NA	5: {ID}
6:	Op	::= '+'	6: { '+' }	6: NA	6: { '+' }
7:		/*/	7: { ' * ' }	7: NA	7: { ' * ' }

Do we need backtracking?

The First+ set is the combination of First and Follow sets

```
First+ sets:
1: Expr ::= Unit Expr2
                               1: {'(', ID}
                          2: { '+', '*' }
2: Expr2 ::= Op Unit Expr2
                               3: {None, ')'}
3:
          11 11
4: Unit ::= '(' Expr ')' 4: {'('}
5:
                             5: {ID}
              ID
6: Op ::= '+'
                               6: { '+' }
                               7: { '*' }
             1 * 1
7:
```

For each non-terminal: if every production has a disjoint First+ set then we do not need any backtracking!

Do we need backtracking?

The First+ set is the combination of First and Follow sets





For each non-terminal: if every production has a disjoint First+ set then we do not need any backtracking!

Do we need backtracking?

The First+ set is the combination of First and Follow sets



These grammars are called LL(1)

- L scanning the input left to right
- L left derivation
- 1 how many look ahead symbols

They are also called predictive grammars

Many programming languages are LL(1)

For each non-terminal: if every production has a disjoint First+ set then we do not need any backtracking!

1: Factor ::= ID 2: | ID '[' Args ']' 3: | ID '(' Args ')'

							Fir	rst
1:	Factor	::=	ID				1:	{}
2:			ID	'['	Args	']'	2:	{}
3:		Ì	ID	'('	Args	')′	3:	{}
	•							

1:	Factor	::=	ID				
2:			ID	'['	Args	']′	
3:			ID	'('	Args	')′	

First						
1:	{ID}					
2:	{ID}					
<mark>3:</mark>	{ID}					
• • •						

We cannot select the next rule based on a single look ahead token!

T 1 - - - - 1

	•						• • •	•
3:			ID	'('	Args	')′	3:	{ID}
2:			ID	'['	Args	']′	2:	{ID}
1:	Factor	::=	ID				1:	{ID}
							FII	ST

We can refactor

			First
1:	Factor	::= ID Option_args	1: {}
2:	Option_args	::= '[' Args ']'	2: {}
3:		'(' Args ')'	3: {}
4:		11 11	4: {}

							Η	fir	st
1:	Factor	::=	ID]	L :	{ID}
2:			ID	'['	Args	']′		2:	{ID}
3:		Ì	ID	'('	Args	')′		3:	{ID}
	•								

We can refactor

1:	Factor	::=	ID	Optior	n_args
2:	Option_args	::=	'['	Args	']′
3:			'('	Args	')′
4:			<i> </i>		

First



// We will need to compute the follow set

Firet

								- D C
1:	Factor	::=	ID				1:	{ID]
2:			ID	'['	Args	']′	2:	{ID]
3:			ID	'('	Args	')′	3:	{ID]
	-							_

It is not always possible to rewrite grammars into a predictive form, but many programming languages can be.

We can refactor

1:	Factor	::=	ID	Optior	n_args
2:	Option_args	::=	'['	Args	']′
3:			'('	Args	')′
4:					

Fi	rst
1:	{ID}
2:	{ ' [' }
3:	{ ' (' }
4:	{ " " }

// We will need to compute the follow set

We now have a full top-down parsing algorithm!

```
root = start symbol;
focus = root;
push(None);
to match = s.token();
while (true):
  if (focus is a nonterminal)
    pick next rule (A ::= B1,B2,B3...BN);
    push(BN... B3, B2);
    focus = B1
```

else if (focus == to match)

to match = s.token()

focus = pop()

```
1: Expr := Unit Expr2
2: { '+', '*' }
                   2: Expr2 ::= Op Unit Expr2
                             11 11
3: {None, ')'}
                   3:
                   4: Unit ::= '(' Expr ')'
                   5:
                                 ID
                            ::= '+'
                   6: Op
                   7:
                                1 * 1
```

First+ sets for each production rule

First+ sets:

4: { ' (' }

6: $\{'+'\}$

7: { '*' }

5: {ID}

1: {'(', ID}

input grammar, refactored to remove *left recursion*

```
else if (to match == None and focus == None)
 Accept
```

To pick the next rule, compare to match with the possible first+ sets. Pick the rule whose first+ set contains to match.

If there is no such rule then it is a parsing error.
Moving on to a simpler implementation:

Recursive Descent Parser

How do we parse an Expr?

How do we parse an Expr? We parse a Unit followed by an Expr2

How do we parse an Expr? We parse a Unit followed by an Expr2

We can just write exactly that!

```
def parse_Expr(self):
    self.parse_Unit();
    self.parse_Expr2();
    return
```

How do we parse an Expr2?

1:	Expr	::= Unit Expr2
2:	Expr2	::= Op Unit Expr2
3:		<i>11 11</i>
4:	Unit	::= '(' Expr ')'
5:		ID
6:	Op	::= '+'
7:		/*/

How do we parse an Expr2?

First+ sets:
1: {'(', ID}
2: {'+', '*'}
3: {None, ')'}
4: {'('}
5: {ID}
6: {'+'}
7: {'*'}

```
1: Expr ::= Unit Expr2
                                                          How do we parse an Expr2?
2: Expr2 ::= Op Unit Expr2
              11 11
3:
4: Unit ::= '(' Expr ')'
5:
                   ID
6: Op
            ::= '+'
                                   def parse Expr2(self):
7:
                 1 * 1
                                      token id = get token id(self.to match)
                                      # Expr2 ::= Op Unit Expr2
                                      if token id in ["PLUS", "MULT"]:
                                          self.parse Op()
First+ sets:
                                          self.parse Unit()
                                          self.parse_Expr2()
1: { '(', ID}
                                          return
2: { '+', '*' }
                                      # Expr2 ::= ""
3: {None, ')'}
                                      if token id in [None, "RPAR"]:
4: { ' ( ' }
                                          return
5: {ID}
                                      raise ParserException(-1,
                                                                                  # line number (for you to do)
6: \{'+'\}
                                                                                  # observed token
                                                          self.to match,
                                                          ["PLUS", "MULT", "RPAR"]) # expected token
7: { '*' }
```

How do we parse a Unit?

First+ sets:
1: { '(', ID}
2: { '+', '*'}
3: {None, ')'}
4: { '('}
5: {ID}
6: { '+'}
7: { '*'}

```
1: Expr ::= Unit Expr2
                                                         How do we parse a Unit?
2: Expr2 ::= Op Unit Expr2
3:
              11 11
4: Unit ::= '(' Expr ')'
                                         def parse Unit(self):
5:
                  ID
            ::= '+'
6: Op
                                             token_id = get_token_id(self.to_match)
7:
                 1 * 1
                                             # Unit ::= '(' Expr ')'
                                             if token id == "LPAR":
                                                self.eat("LPAR")
                                                self.parse Expr()
                                                self.eat("RPAR")
First+ sets:
                                                return
1: { '(', ID}
                                             # Unit :: = ID
2: { '+', '*' }
                                             if token id == "ID":
3: {None, ')'}
                                                self.eat("ID")
                                                return
4: { ' ( ' }
                                             raise ParserException(-1, # line number (for you to do)
5: {ID}
                                                                self.to match, # observed token
6: { '+' }
                                                                 ["LPAR", "ID"]) # expected token
7: { '*' }
```

```
1: Expr ::= Unit Expr2
                                                           How do we parse a Unit?
2: Expr2 ::= Op Unit Expr2
3:
               11 11
4: Unit ::= '(' Expr ')'
                                           def parse Unit(self):
5:
                   ID
            ::= '+'
6: Op
                                               token_id = get_token_id(self.to_match)
7:
                  1 * 1
                                               # Unit ::= '(' Expr ')'
                                               if token id == "LPAR":
                                                                         ensure that to_match has token ID of "LPAR"
                                                  self.eat("LPAR")
                                                                         and get the next token
                                                  self.parse_Expr()
                                                  self.eat("RPAR")
First+ sets:
                                                  return
1: { '(', ID}
                                               # Unit :: = ID
2: { '+', '*' }
                                               if token id == "ID":
                                                  self.eat("ID")
3: {None, ')'}
                                                  return
4: { ' ( ' }
                                               raise ParserException(-1, # line number (for you to do)
5: {ID}
                                                                   self.to match, # observed token
6: \{'+'\}
                                                                   ["LPAR", "ID"]) # expected token
7: { '*' }
```

How do we parse an Op?

First+ sets:
1: {'(', ID}
2: {'+', '*'}
3: {None, ')'}
4: {'('}
5: {ID}
6: {'+'}
7: {'*'}

```
1: Expr ::= Unit Expr2
2: Expr2 ::= Op Unit Expr2
3:
            11 11
4: Unit ::= '(' Expr ')'
5:
                ID
          ::= '+'
6: Op
7:
               1 * 1
First+ sets:
1: { '(', ID}
2: { '+', '*' }
3: {None, ')'}
4: { ' ( ' }
5: {ID}
6: { '+' }
7: { '*' }
```

```
How do we parse an Op?
```

```
def parse_0p(self):
```

```
token_id = get_token_id(self.to_match)
```

```
# Op ::= '+'
if token_id == "PLUS":
    self.eat("PLUS")
    return
```

```
# Op ::= '*'
if token_id == "MULT":
    self.eat("MULT")
    return
```

Recursive Descent Parsers

- Requires some work on the grammar
 - Remove left recursion
 - Always possible
 - Identify first+ sets
 - Not always possible
 - Might require rewriting the grammar
- Allows a straightforward and efficient parser implementation
 - Many industry contacts have told me they use recursive descent parsers internally
- Next homework will have you implement a recursive descent parser for the grammar you are working on now!
- Very interesting use case of recursion!

Moving on: Scope

Scope

- What is scope?
- Can it be determined at compile time? Can it be determined at runtime?
- C vs. Python
- Anyone have any interesting scoping rules they know of?

One consideration: Scope

• Lexical scope example

int x = 0; int y = 0; { int y = 0; x+=1; y+=1; } x+=1; y+=1;

What are the final values in x and y?

- Symbol table object
- two methods:
 - lookup(id) : lookup an id in the symbol table. Returns None if the id is not in the symbol table.
 - insert(id, info) : insert a new id (or overwrite an existing id) into the symbol table along with a set of information about the id.

a very simple programming language

ID = [a-z]+ INCREMENT = "\+\+" TYPE = "int" LBRAC = "{" RBRAC = "}" SEMI = ";" int x; x++; int y; y++;

statements are either a declaration or an increment

a very simple programming language

ID = [a - z] +	int x;
INCREMENT = " + + "	{ int y; x++;
TYPE = "int"	
LBRAC = "{"	у++;
RBRAC = "}"	} v++:
SEMI = ";"	

statements are either a declaration or an increment

a very simple programming language





error!

statements are either a declaration or an increment

• SymbolTable ST;

Say we are matched the statement: int x;

declare_statement ::= TYPE ID SEMI { }

lookup(id) : lookup an id in the symbol table. Returns None if the id is not in the symbol table.

insert(id,info) : insert a new id (or overwrite an existing id) into the symbol table along with a set of information about the id.

• SymbolTable ST;

Say we are matched the statement: int x;

```
declare_statement ::= TYPE ID SEMI
{
   self.eat(TYPE)
   variable_name = self.to_match.value
   self.eat(ID)
   ST.insert(variable_name,None)
   self.eat(SEMI)
}
```

• SymbolTable ST;

Say we are matched string: x++;

inc_statement ::= ID INCREMENT SEMI { }

lookup(id) : lookup an id in the symbol table. Returns None if the id is not in the symbol table.

insert(id,info) : insert a new id (or overwrite an existing id) into the symbol table along with a set of information about the id.

```
• SymbolTable ST;
```

Say we are matched string: x++;

```
inc_statement ::= ID INCREMENT SEMI
{
    variable_name = self.to_match.value
    if ST.lookup(variable_name) is None:
        raise SymbolTableException(variable_name)
    self.eat(ID)
    self.eat(INCREMENT)
    self.eat(SEMI)
```

• SymbolTable ST;

statement : LBRAC statement_list RBRAC



• SymbolTable ST;

statement : LBRAC statement_list RBRAC



start a new scope S

remove the scope S

- Symbol table
- four methods:
 - lookup(id) : lookup an id in the symbol table. Returns None if the id is not in the symbol table.
 - insert(id, info) : insert a new id into the symbol table along with a set of information about the id.
 - push_scope() : push a new scope to the symbol table
 - pop_scope() : pop a scope from the symbol table

• SymbolTable ST;

statement : LBRAC statement_list RBRAC

You will be adding the functions to push and pop scopes in your homework

- Thoughts? What data structures are good at mapping strings?
- Symbol table
- four methods:
 - lookup(id) : lookup an id in the symbol table. Returns None if the id is not in the symbol table.
 - insert(id, info) : insert a new id into the symbol table along with a set of information about the id.
 - push_scope() : push a new scope to the symbol table
 - **pop_scope()** : pop a scope from the symbol table

- Many ways to implement:
- A good way is a stack of hash tables:

base scope

HT 0

- Many ways to implement:
- A good way is a stack of hash tables:

push_scope()

HT 0

- Many ways to implement:
- A good way is a stack of hash tables:



- Many ways to implement:
- A good way is a stack of hash tables:





- Many ways to implement:
- A good way is a stack of hash tables:

insert (id -> data) at top hash table



Stack of hash tables

insert(id,data)

- Many ways to implement:
- A good way is a stack of hash tables:

HT 1

HT 0

lookup(id)
- Many ways to implement:
- A good way is a stack of hash tables:



- Many ways to implement:
- A good way is a stack of hash tables:

HT 1

lookup(id)

HT 0

then check

here

- Many ways to implement:
- A good way is a stack of hash tables:



pop_scope()

HT 0

- Many ways to implement:
- A good way is a stack of hash tables:

HT 0

• Example int x = 0; int y = 0; { int y = 0; x++; y++; } x++;

y++;

HT 0

See you on Friday!

• We will discuss parser generators