

Week 5
Parsing

Hassam



Outline

Compilers

Parsing



Section 1

Compilers



C to executable?

How do you go from "code" to something the computer understands?



C to executable?

How do you go from "code" to something the computer understands?

- What does the computer even understand?



C to executable?

How do you go from "code" to something the computer understands?

- What does the computer even understand?
- What is code?

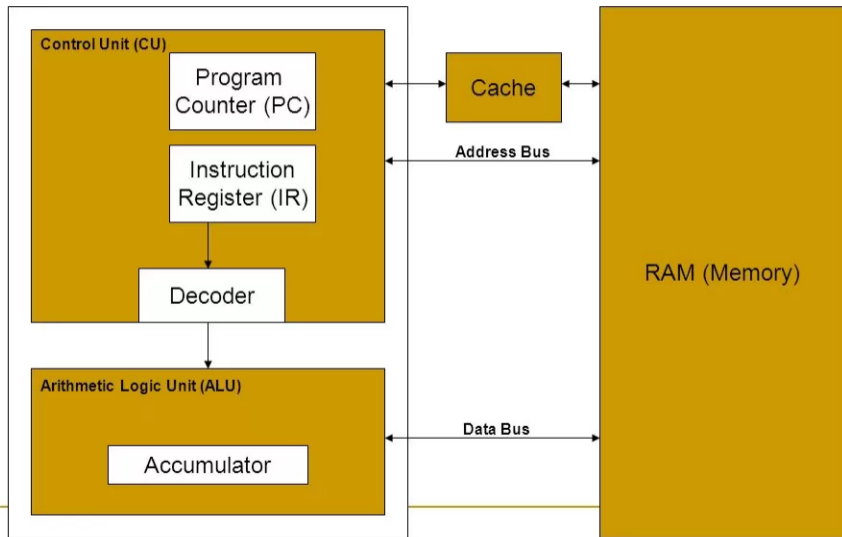


Assembly

@nebu addressed this two weeks ago, but let's review.



CPU Diagram

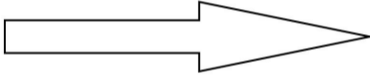


Assembly Language

```
mov ecx, ebx  
mov esp, edx  
mov edx, r9d  
mov rax, rdx
```

Programmer

Assembler + Linker



Machine Language

```
100101011001  
010011111011  
111010101101  
01010101010
```

Processor



Code?

```
1  int main(int argc, char** argv) {  
2      puts("Hello, World!");  
3      return 0;  
4  }
```



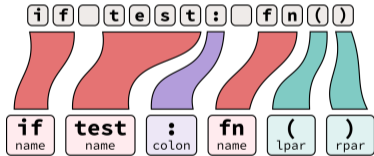
Code?

```
1 def main():  
2     print("Hello, world!")
```

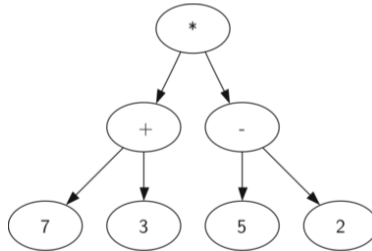
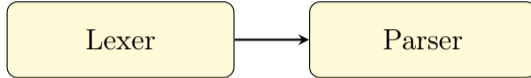


What are the steps in between?

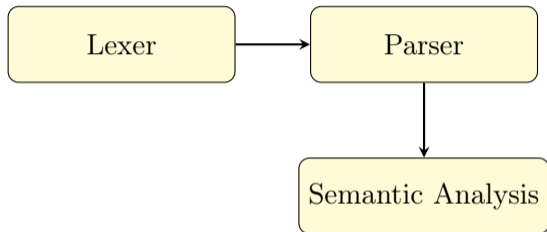
Lexer



What are the steps in between?



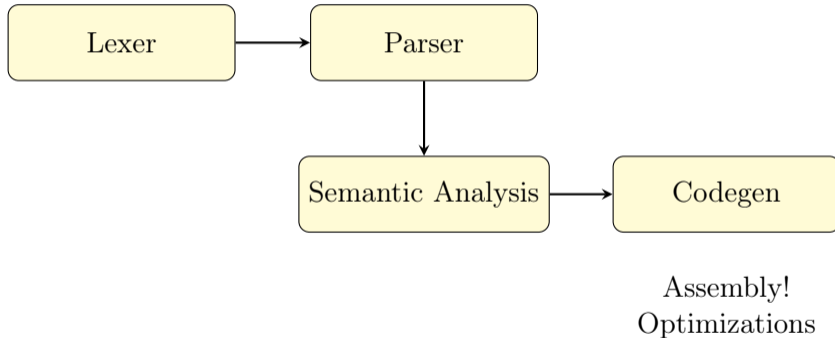
What are the steps in between?



Syntax Errors
Type Checking
Lifetime Analysis



What are the steps in between?



Questions?



Section 2

Parsing



Remember CFGs?

- So far, we have been using CFGs without any reason.
- Soon, we'll use them to prove properties of computation!



Remember CFGs?

- So far, we have been using CFGs without any reason.
- Soon, we'll use them to prove properties of computation!
- But, sometimes you just want to specify a grammar and parse it, CFGs are useful here too!
- CFGs are ugly though, and we're leaving the world of theory, so we have something better.

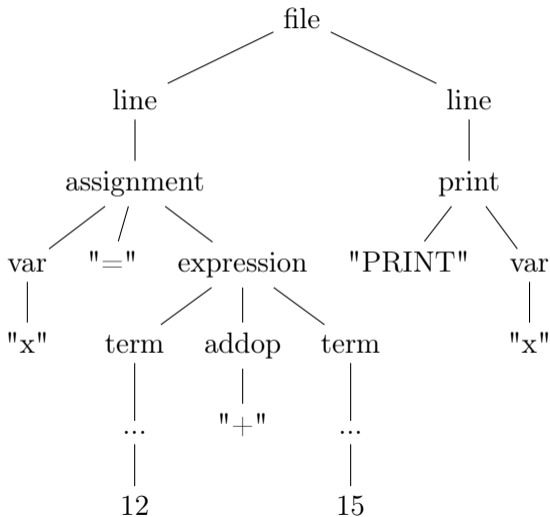


Extended Backus–Naur form

```
1 file = line { line } <EOF>.
2 line = [ assignment | print | reset ] <NL>.
3 assignment = var " :=" expression.
4 print = "PRINT" var.
5 reset = "RESET".
6 expression = term { addop term }.
7 term = factor { mulop factor }.
8 factor = "(" expression ")" | var | number.
9 addop = "+" | "-".
10 mulop = "*".
11 var = letter { letter | digit }.
12 number = [ "-" ] digit { digit }.
13 letter = ('A'-'Z') | ('a'-'z').
14 digit = ('0'-'9').
```



Parse Tree



Parse Tree

```
1  x = 12 + 15
2  PRINT x
```



So, PDAs?

- This "language" is representable by a CFG, so should we use a PDA to parse it?



So, PDAs?

- This "language" is representable by a CFG, so should we use a PDA to parse it?
- Let's not.



LL Parsing

- Reading **L**eft to right, performing a **L**eftmost processing.
- Formally, this type of parser is equivalent to a deterministic pushdown automata. But, much easier to write.
- Requires our grammar to follow some rules, but all grammars (parseable by DPDAs) can be (annoyingly) converted to do so.



Easy to write?

So easy to write:

- 1 `start = F | '(' start '+' F ')'`
- 2 `F = 'a'`

This parses expressions of the form: $(a + a)$, $((a + a) + a)$, etc.



How do we code this?

```
1 parseF 'a':xs = F, xs
2 parseF [] = error
3
4 parseS '(':xs = {
5     S1, '+':restS = parseS(xs)
6     F, ')':rest = parseF(restS)
7     return (S (S1, F)), rest
8 }
9 parseS x = {
10     F, rest = parseF(x)
11     return S(F), rest
12 }
13 parseS [] = error
```

```
1 start = F
2   | '(' start '+' F ')'
3   F = 'a'
```



What goes wrong?

- Our code from before is valid and generalizable, but it makes some assumptions.



What goes wrong?

- Our code from before is valid and generalizable, but it makes some assumptions.
- It assumes that there is no "left recursion"



What goes wrong?

- Our code from before is valid and generalizable, but it makes some assumptions.
- It assumes that there is no "left recursion"
- It assumes that there are no common prefixes



What goes wrong?

- Our code from before is valid and generalizable, but it makes some assumptions.
- It assumes that there is no "left recursion"
- It assumes that there are no common prefixes
- It assumes that you can look one character ahead and determine what to do next. This is formally known as LL(1).



What goes wrong?

- Our code from before is valid and generalizable, but it makes some assumptions.
- It assumes that there is no "left recursion"
- It assumes that there are no common prefixes
- It assumes that you can look one character ahead and determine what to do next. This is formally known as LL(1).
- It assumes there's no ambiguity.



What goes wrong?

- Our code from before is valid and generalizable, but it makes some assumptions.
- It assumes that there is no "left recursion"
- It assumes that there are no common prefixes
- It assumes that you can look one character ahead and determine what to do next. This is formally known as LL(1).
- It assumes there's no ambiguity.

Which of these can we fix?



Ambiguity?

```
1 void func() {  
2     a < b , c > d;  
3 }
```



Ambiguity?

```
1 void func() {  
2     a < b , c > d;  
3 }
```

Is this a template instantiation? Or are we using `operator<` and `operator>` on two variables?



Ambiguity?

```
1 void func() {  
2     a < b , c > d;  
3 }
```

Is this a template instantiation? Or are we using `operator<` and `operator>` on two variables?

It's both.



Fixing bad grammars

It turns out, we can fix them all, except for the last one. We need a more powerful tool for ambiguity.

- It is possible to rewrite anything left-recursive to be non-left-recursive.
- Removing common prefixes is called "left factoring".
- We can change from an LL(1) parser to an LL(k) parser, and "look ahead" k steps. This becomes messier and messier the further we look ahead though.

Nearly every mainstream programming language is LL(1). They are convenient to parse, and even more convenient to write.



Fixing left recursion

```
1 start = start '+' | 'i'
```

This parses `i+++++`, and onwards. How do we get rid of the left recursion?



Fixing left recursion

1 `start = start '+' | 'i'`

This parses `i+++++`, and onwards. How do we get rid of the left recursion?

The general strategy is:

- `S = 'x'` becomes `S = 'x'S_new`
- `S = S'x'` becomes `S_new = 'x'S_new`
- Add `S_new = ϵ`



Fixing left recursion

```
1 start = start '+' | 'i'
```

This parses `i+++++`, and onwards. How do we get rid of the left recursion?

The general strategy is:

- $S = 'x'$ becomes $S = 'x'S_{\text{new}}$
- $S = S'x'$ becomes $S_{\text{new}} = 'x'S_{\text{new}}$
- Add $S_{\text{new}} = \epsilon$

```
1 start = 'i' start_new
```

```
2 start_new = '+'start_new |  $\epsilon$ 
```



Fixing common prefixes

```
1 start = 'a' B | 'a' 'd'  
2 B = 'c'
```



Fixing common prefixes

1 start = 'a' B | 'a' 'd'

2 B = 'c'

Becomes:

1 start = 'a' D

2 D = B | 'd'

3 B = 'c'



Solving Ambiguity

Remember our CFG from last week?

$$S \rightarrow \varepsilon \mid Sab \mid aSb \mid aSbS \\ \mid Sba \mid bSa \mid bSaS$$

We can try to fix the left recursion and the common prefixes, but we'll get stuck in an infinite loop. This grammar is "ambiguous", we would need to look ahead infinitely many steps to parse it. We need non-determinism to parse this.



Backtracking to the rescue?

```
1  def parseS(inp):
2      if inp == '': return S(), ''
3      try:
4          s, rest = parseS(inp)
5          if rest[0] == 'a':
6              assert rest[1] == 'b'
7              return S(s, 'a', 'b'), rest[2:]
8          elif rest[0] == 'b':
9              assert rest[1] == 'a'
10             return S(s, 'b', 'a'), rest[2:]
11         else:
12             assert False
13     except Exception:
14         assert inp[0] == 'a'
15         try:
16             s, rest = parseS(inp)
17             ...
```



Yikes

This is messy, and it's slow! Try finishing it by yourself.

Most modern compilers for non-LL(1) languages [C and C++ :(] use handwritten backtracking parsers like this one.



Writing parsers for fun and profit

Writing parsers is actually very fun! I have had four interviews this semester that required writing a custom parser.

- Try writing parsers for some of the simple examples we discussed today in your favorite language.
- Go meta. Write a "parser compiler" that takes some basic version of EBNF and turns it into code that runs a parser.
- So much more to explore. LR(k), LALR(k), etc.
- Parsing is "solved", more interesting problems in compilers now.



Goodbye

Any sufficiently complicated C or Fortran program contains an ad hoc, informally-specified, bug-ridden, slow implementation of half of Common Lisp.

— Greenspun's tenth rule of programming (1993)

