

Announcements

- P2 posted
- HW2 posted

Context-free grammars (CFGs)

Roadmap

Last time

- Regex == DFA
- JLex: a tool for generating (Java code for) Lexers/Scanners

This time

- CFGs, the underlying abstraction for Parsers

Next week

- Java CUP: A tool for generating (Java code for) parser

regular expression ==> scanner, JLex

CFG ==> parser, Java CUP

RegExs Are Great!

Perfect for tokenizing a language

RegExs are great for scanner and tokenization, but not sufficient for parsing.

They do have some limitations

- Limited class of language that cannot specify all programming constructs we need
- No notion of structure

Let's explore both of these issues

Limitations of RegExps

Regex and DFAs are equivalent.
For every regexp, there's an equivalent DFA that represents the same language, and vice versa.

Cannot handle “matching”

E.g., language of balanced parentheses

$$L_{()} = \{ ({}^n)^n \text{ where } n > 0 \}$$

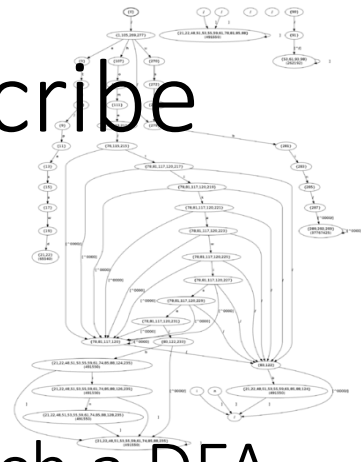
No DFA exists for this language

Intuition: A given FSM only has a fixed, finite amount of memory

- For an FSM, memory = the states
- With a fixed, finite amount of memory, how could an FSM remember how many “(” characters it has seen?

Theorem: No RegEx/DFA can describe the language $L_{()}()$

Can be asked in the exam.

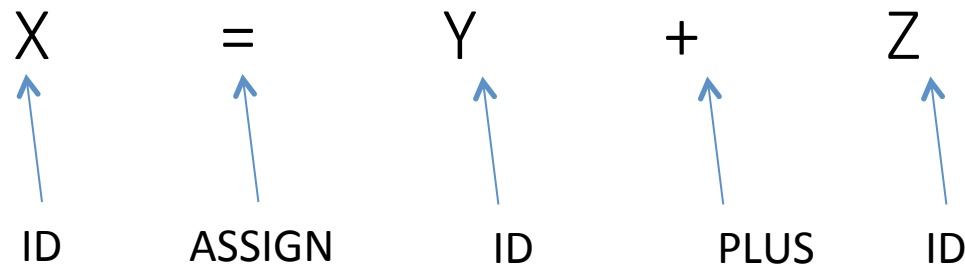


By contradiction:

- Let's say there exists a DFA A for $L_{()}()$ and such a DFA has N states
- A has to accept the string $(^N)^N$ with some path $q_0q_1...q_N...q_{2N}$
 q_0 is the start state, so $q_1...q_{2N}$ takes $2N$ characters in total.
- By *pigeonhole principle* some state has repeated:
 $q_i = q_j$ for some $i < j \leq N$
- Therefore the run $q_0q_1...q_iq_{j+1}...q_N...q_{2N}$ is also accepting
- A accepts the string $(^{N-(j-i)})^N$ not in $L_{()}()$ → contradiction!

Limitations of RegEx: Structure

Our Enhanced-RegEx scanner can emit a stream of tokens:



... but this doesn't really enforce any order of operations

Regular expression cannot indicate the precedence or associativity of the operators.

The Chomsky Hierarchy



Turing machine

LANGUAGE CLASS:

Recursively enumerable

Context-Sensitive

Context-Free

Regular

Happy medium?

FSM



Noam Chomsky

power

efficiency



Context Free Grammars (CFGs)

A set of (recursive) rewriting rules to generate patterns of strings

To build a tree, you need recursion.

Can envision a “parse tree” that keeps structure

CFG: Intuition

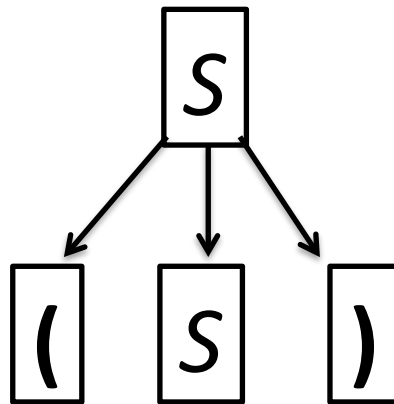
$$S \rightarrow \underbrace{(' S ')}_{\text{single set of parenthesis}}$$

A rule that says that you
can rewrite S to be an S surrounded by
a single set of parenthesis

Before applying rule



After applying rule



Context Free Grammars (CFGs)

A CFG is a 4-tuple (N, Σ, P, S)

- N is a set of non-terminals, e.g., $A, B, S...$
- Σ is the set of terminals
- P is a set of production rules
- S (in N) is the initial non-terminal symbol

Context Free Grammars (CFGs)

A CFG is a 4-tuple (N, Σ, P, S)

- N is a set of non-terminals, e.g., $A, B, S \dots$
- Σ is the set of terminals
- P is a set of production rules
- S (in N) is the initial non-terminal symbol

Placeholder / interior nodes
in the parse tree

Leaf of the tree.

Tokens from
scanner

Rules for deriving strings

Compare math representation of FSM and CFG:

$FSM = (Q, \Sigma, \delta, q, F)$

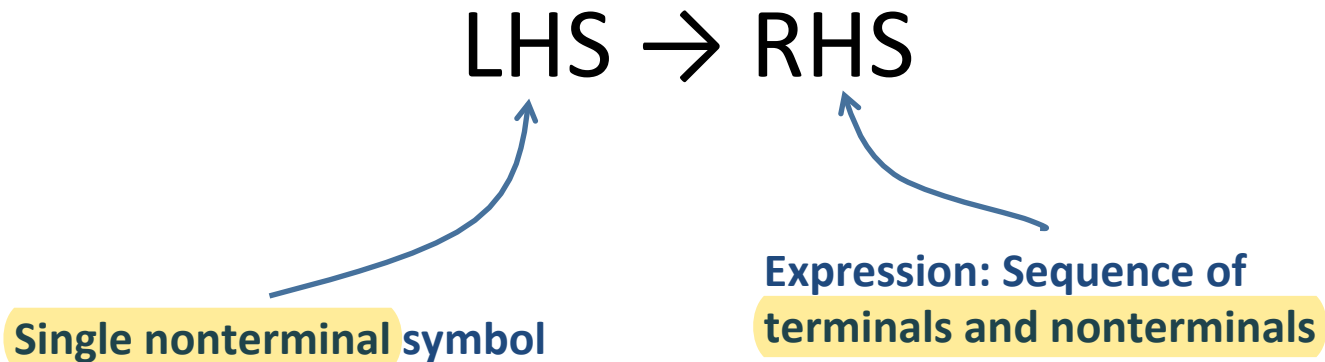
Q : finite set of states, Σ : alphabets, δ : transition functions, q : start state, F : final states.

$CFG = (N, \Sigma, P, S)$

N : non-terminals, Σ : terminals, P : production rules, S : initial non-terminal symbol.

If not otherwise specified, use the
non-terminal that appears on the LHS
of the first production as the start

Production Syntax



Examples:

$$S \rightarrow '(' S ')'$$
$$S \rightarrow \epsilon$$

Production Shorthand

Nonterm \rightarrow expression

$S \rightarrow '(' S ')'$

Nonterm $\rightarrow \varepsilon$

$S \rightarrow \varepsilon$

equivalently:

Nonterm \rightarrow expression
 $\quad \quad \quad | \varepsilon$

$S \rightarrow '(' S ')'$
 $\quad \quad \quad | \varepsilon$

equivalently:

Nonterm \rightarrow expression $| \varepsilon$

$S \rightarrow '(' S ')'$ $| \varepsilon$

Derivations

Now we are generating the syntax tree by randomly picking a production rule.

To derive a string:

- Start by setting “*Current Sequence*” to the start symbol
- Repeat:
 - Find a Nonterminal X in the Current Sequence
 - Find a production of the form $X \rightarrow \alpha$
 - “Apply” the production: create a new “current sequence” in which α replaces X
- Stop when there are no more non-terminals
- This process derives a string of terminal symbols

Derivation Syntax

- We'll use the symbol $\Rightarrow \perp$ for *derives*
- We'll use the symbol $\Rightarrow + \perp$ for *derives in one or more steps*
- We'll use the symbol $\Rightarrow^* \perp$ for *derives in zero or more steps*

An Example Grammar

An Example Grammar

Terminals

begin

end

semicolon

assign

id

plus

An Example Grammar

For readability, bold and lowercase

Terminals

begin

end

semicolon

assign

id

plus

An Example Grammar

For readability, bold and lowercase

Terminals

begin } **Program**
end } **boundary**
semicolon
assign
id
plus

An Example Grammar

For readability, bold and lowercase

Terminals

begin } **Program**
end } **boundary**

semicolon  **Represents ";"**
Separates statements

assign

id

plus

An Example Grammar

For readability, bold and lowercase

Terminals

begin } **Program**
end } **boundary**
semicolon — **Represents “;”**
assign — **Separates statements**
id
plus — **Represents “=” statement**

An Example Grammar

For readability, bold and lowercase

Terminals

begin } **Program**
end } **boundary**
semicolon — **Represents “;”**
assign — **Separates statements**
id — **Represents “=” statement**
plus — **Identifier / variable name**

An Example Grammar

For readability, bold and lowercase

Terminals

begin } **Program**
end } **boundary**
semicolon — **Represents “;”**
assign — **Separates statements**
id — **Represents “=” statement**
plus — **Identifier / variable name**
— **Represents “+” expression**

An Example Grammar

For readability, bold and lowercase

Terminals

begin
end
semicolon
assign
id
plus

Nonterminals

Prog
Stmts
Stmt
Expr

An Example Grammar

For readability, bold and lowercase

Terminals

begin
end
semicolon
assign
id
plus

For readability, Italics and UpperCamelCase

Nonterminals

Prog
Stmts
Stmt
Expr

An Example Grammar

For readability, bold and lowercase

Terminals

begin
end
semicolon
assign
id
plus

For readability, Italics and UpperCamelCase

Nonterminals

Prog ————— **Root of the parse tree**
Stmts
Stmt
Expr

An Example Grammar

For readability, bold and lowercase

Terminals

begin
end
semicolon
assign
id
plus

For readability, Italics and UpperCamelCase

Nonterminals

Prog ————— **Root of the parse tree**
Stmts ————— **List of statements**
Stmt
Expr

An Example Grammar

For readability, bold and lowercase

Terminals

begin
end
semicolon
assign
id
plus

For readability, Italics and UpperCamelCase

Nonterminals

<i>Prog</i>	—————	Root of the parse tree
<i>Stmts</i>	—————	List of statements
<i>Stmt</i>	—————	A single statement
<i>Expr</i>		

An Example Grammar

For readability, bold and lowercase

Terminals

begin
end
semicolon
assign
id
plus

For readability, Italics and UpperCamelCase

Nonterminals

<i>Prog</i>	—————	Root of the parse tree
<i>Stmts</i>	—————	List of statements
<i>Stmt</i>	—————	A single statement
<i>Expr</i>	—————	A mathematical expression

An Example Grammar

For readability, bold and lowercase

Terminals

begin
end
semicolon
assign
id
plus

For readability, Italics and UpperCamelCase

Nonterminals

Prog
Stmts
Stmt
Expr

Defines the syntax of legal programs

Productions

$Prog \rightarrow \mathbf{begin} \text{ Stmts } \mathbf{end}$

$Stmts \rightarrow \text{ Stmts } \mathbf{semicolon} \text{ Stmt}$
 $\quad \quad \quad | \text{ Stmt}$

$Stmt \rightarrow \mathbf{id} \mathbf{assign} \text{ Expr}$

$Expr \rightarrow \mathbf{id}$

$\quad \quad \quad | \text{ Expr } \mathbf{plus} \text{ id}$

An Example Grammar

For readability, bold and lowercase

Terminals

begin } **Program**
end } **boundary**
semicolon — **Represents “;”**
assign — **Separates statements**
id — **Represents “=” statement**
plus — **Identifier / variable name**
— **Represents “+” expression**

For readability, Italics and UpperCamelCase

Nonterminals

Prog — **Root of the parse tree**
Stmts — **List of statements**
Stmt — **A single statement**
Expr — **A mathematical expression**

Defines the syntax of legal programs

Productions

Prog → **begin** *Stmts* **end**

Stmts → *Stmts* **semicolon** *Stmt*
| *Stmt*

Stmt → **id** **assign** *Expr*

Expr → **id**

| *Expr* **plus** *id*

Productions

1. *Prog* → **begin** *Stmts* **end**
2. *Stmts* → *Stmts* **semicolon** *Stmt*
3. | *Stmt*
4. *Stmt* → **id assign** *Expr*
5. *Expr* → **id**
6. | *Expr* **plus id**

Productions

1. $Prog \rightarrow \mathbf{begin} \textit{Stmts} \mathbf{end}$
2. $Stmts \rightarrow \textit{Stmts} \mathbf{semicolon} \textit{Stmt}$
3. $\quad \quad \quad | \textit{Stmt}$
4. $Stmt \rightarrow \mathbf{id} \mathbf{assign} \textit{Expr}$
5. $Expr \rightarrow \mathbf{id}$
6. $\quad \quad \quad | \textit{Expr} \mathbf{plus} \mathbf{id}$

Derivation Sequence

Productions

Parse Tree

1. *Prog* → **begin** *Stmts* **end**
2. *Stmts* → *Stmts* **semicolon** *Stmt*
3. | *Stmt*
4. *Stmt* → **id assign** *Expr*
5. *Expr* → **id**
6. | *Expr* **plus id**

Derivation Sequence

Productions

- 1. *Prog* → **begin** *Stmts* **end**
- 2. *Stmts* → *Stmts* **semicolon** *Stmt*
- 3. | *Stmt*
- 4. *Stmt* → **id assign** *Expr*
- 5. *Expr* → **id**
- 6. | *Expr* **plus id**

Derivation Sequence

Parse Tree

Key

terminal

Nonterminal

Rule
used

Productions

- 1. *Prog* → **begin** *Stmts* **end**
- 2. *Stmts* → *Stmts* **semicolon** *Stmt*
- 3. | *Stmt*
- 4. *Stmt* → **id** **assign** *Expr*
- 5. *Expr* → **id**
- 6. | *Expr* **plus** **id**

Derivation Sequence

Prog

Parse Tree



Key

terminal

Nonterminal

Rule
used

Productions

- 1. *Prog* → **begin** *Stmts* **end**
- 2. *Stmts* → *Stmts* **semicolon** *Stmt*
- 3. | *Stmt*
- 4. *Stmt* → **id** **assign** *Expr*
- 5. *Expr* → **id**
- 6. | *Expr* **plus** **id**

Derivation Sequence

Prog ⇒ **begin** *Stmts* **end** 

Parse Tree



Key

terminal

Nonterminal

Rule
used

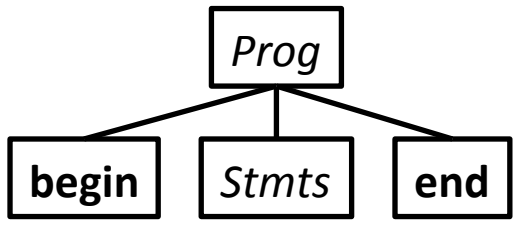
Productions

- 1. *Prog* → **begin** *Stmts* **end**
- 2. *Stmts* → *Stmts* **semicolon** *Stmt*
- 3. | *Stmt*
- 4. *Stmt* → **id** **assign** *Expr*
- 5. *Expr* → **id**
- 6. | *Expr* **plus** **id**

Derivation Sequence

Prog ⇒ **begin** *Stmts* **end** 1

Parse Tree



Key

terminal

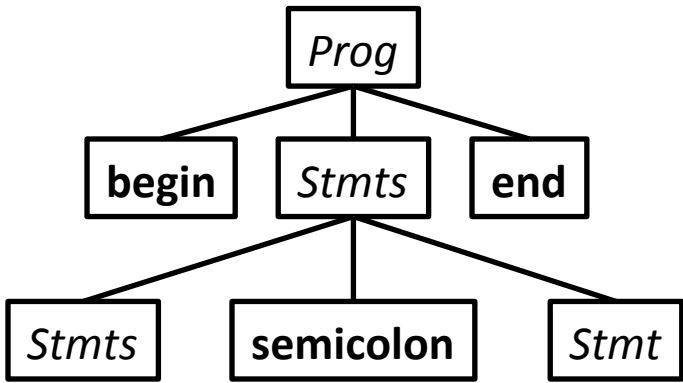
Nonterminal

Rule used

Productions

- 1. *Prog* → **begin** *Stmts* **end**
- 2. *Stmts* → *Stmts* **semicolon** *Stmt*
- 3. | *Stmt*
- 4. *Stmt* → **id** **assign** *Expr*
- 5. *Expr* → **id**
- 6. | *Expr* **plus** **id**

Parse Tree



Derivation Sequence

Prog ⇒ **begin** *Stmts* **end** ①
 ⇒ **begin** *Stmts* **semicolon** *Stmt* **end** ②

Key

terminal

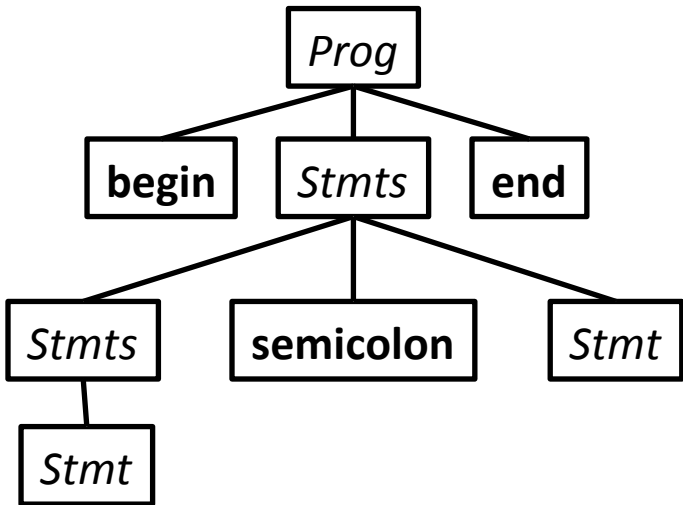
Nonterminal

Rule used

Productions

- 1. *Prog* → **begin** *Stmts* **end**
- 2. *Stmts* → *Stmts* **semicolon** *Stmt*
- 3. | *Stmt*
- 4. *Stmt* → **id** **assign** *Expr*
- 5. *Expr* → **id**
- 6. | *Expr* **plus** **id**

Parse Tree



Derivation Sequence

- Prog* ⇒ **begin** *Stmts* **end** 1
- ⇒ **begin** *Stmts* **semicolon** *Stmt* **end** 2
- ⇒ **begin** *Stmt* **semicolon** *Stmt* **end** 3

Key

terminal

Nonterminal

Rule used

Productions

- 1. *Prog* → **begin** *Stmts* **end**
- 2. *Stmts* → *Stmts* **semicolon** *Stmt*
- 3. | *Stmt*
- 4. *Stmt* → **id** **assign** *Expr*
- 5. *Expr* → **id**
- 6. | *Expr* **plus** **id**

Derivation Sequence

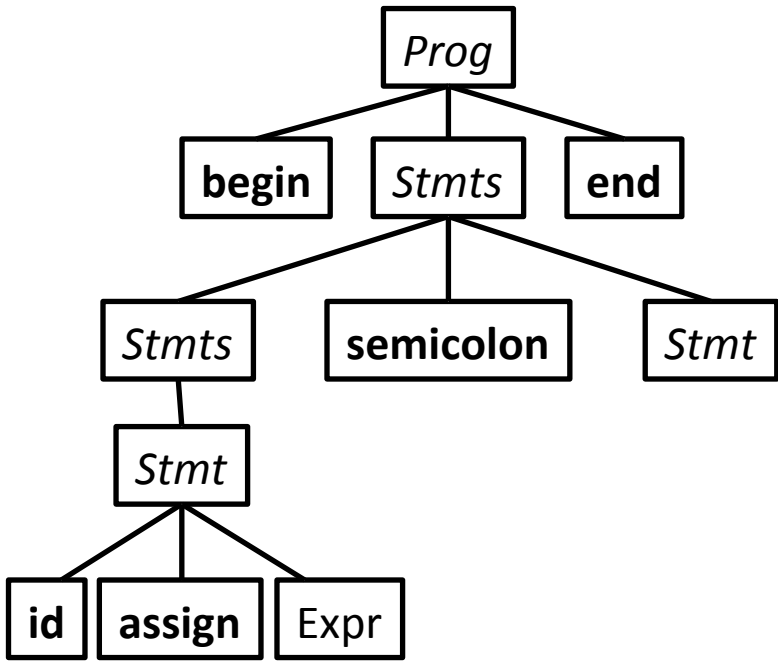
Prog ⇒⊥ **begin** *Stmts* **end** 1

⇒⊥ **begin** *Stmts* **semicolon** *Stmt* **end** 2

⇒⊥ **begin** *Stmt* **semicolon** *Stmt* **end** 3

⇒⊥ **begin** **id** **assign** *Expr* **semicolon** *Stmt* **end** 4

Parse Tree



Key

terminal

Nonterminal

Rule used

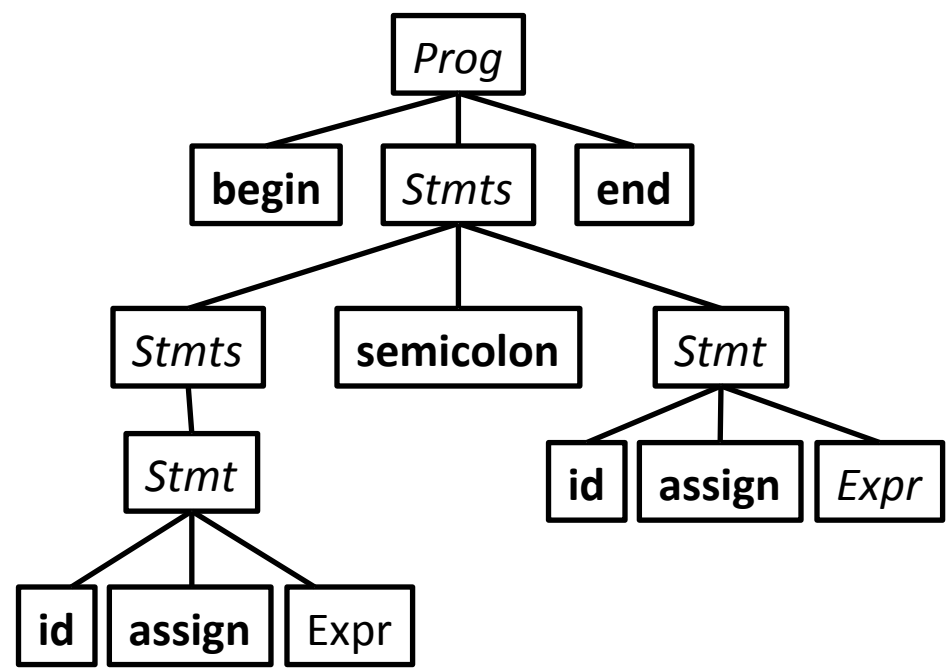
Productions

- 1. *Prog* → **begin** *Stmts* **end**
- 2. *Stmts* → *Stmts* **semicolon** *Stmt*
- 3. | *Stmt*
- 4. *Stmt* → **id** **assign** *Expr*
- 5. *Expr* → **id**
- 6. | *Expr* **plus** **id**

Derivation Sequence

Prog ⇒ **begin** *Stmts* **end** ①
⇒ **begin** *Stmts* **semicolon** *Stmt* **end** ②
⇒ **begin** *Stmt* **semicolon** *Stmt* **end** ③
⇒ **begin** **id** **assign** *Expr* **semicolon** *Stmt* **end** ④
⇒ **begin** **id** **assign** *Expr* **semicolon** **id** **assign** *Expr* **end** ④

Parse Tree



Key

terminal

Nonterminal

Rule used

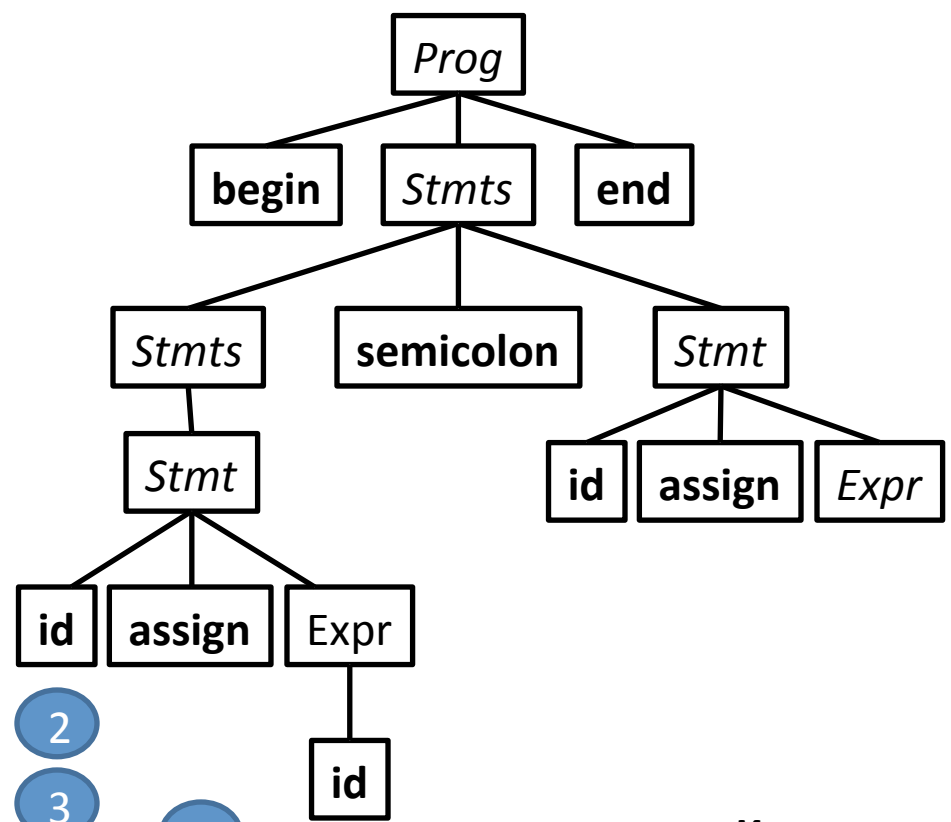
Productions

- 1. *Prog* → **begin** *Stmts* **end**
- 2. *Stmts* → *Stmts* **semicolon** *Stmt*
- 3. | *Stmt*
- 4. *Stmt* → **id** **assign** *Expr*
- 5. *Expr* → **id**
- 6. | *Expr* **plus** **id**

Derivation Sequence

Prog ⇒ **begin** *Stmts* **end** ①
⇒ **begin** *Stmts* **semicolon** *Stmt* **end** ②
⇒ **begin** *Stmt* **semicolon** *Stmt* **end** ③
⇒ **begin** **id** **assign** *Expr* **semicolon** *Stmt* **end** ④
⇒ **begin** **id** **assign** *Expr* **semicolon** **id** **assign** *Expr* **end** ④
⇒ **begin** **id** **assign** **id** **semicolon** **id** **assign** *Expr* **end** ⑤

Parse Tree



Key

terminal

Nonterminal

Rule used

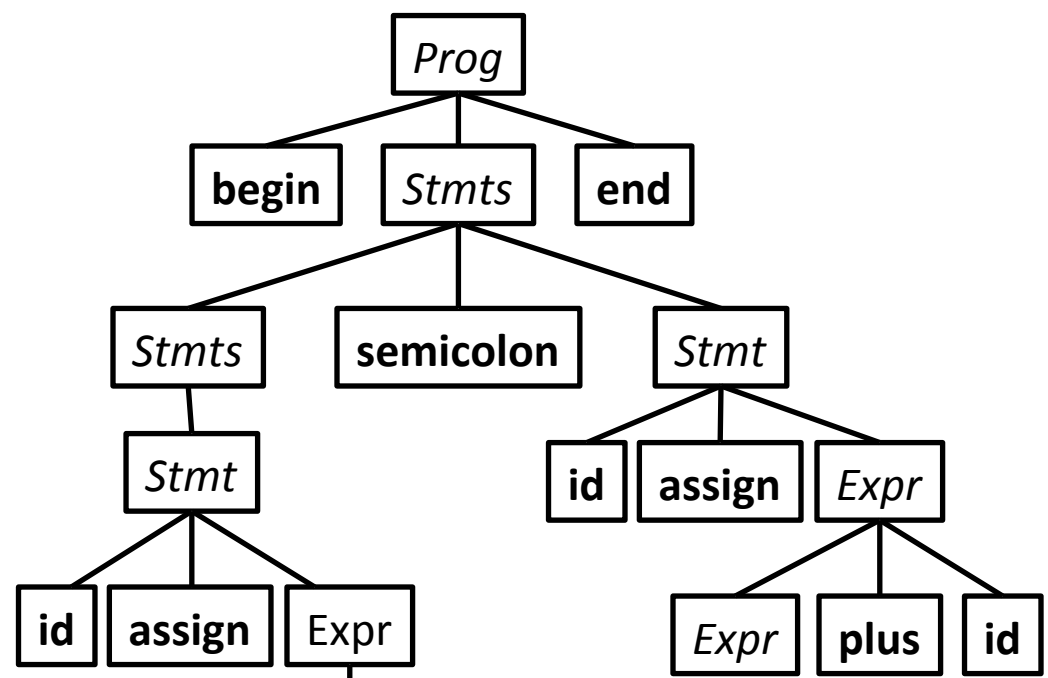
Productions

- 1. *Prog* → **begin** *Stmts* **end**
- 2. *Stmts* → *Stmts* **semicolon** *Stmt*
- 3. | *Stmt*
- 4. *Stmt* → **id** **assign** *Expr*
- 5. *Expr* → **id**
- 6. | *Expr* **plus** **id**

Derivation Sequence

Prog ⇒ **begin** *Stmts* **end** ①
⇒ **begin** *Stmts* **semicolon** *Stmt* **end** ②
⇒ **begin** *Stmt* **semicolon** *Stmt* **end** ③
⇒ **begin** **id** **assign** *Expr* **semicolon** *Stmt* **end** ④
⇒ **begin** **id** **assign** *Expr* **semicolon** **id** **assign** *Expr* **end** ④
⇒ **begin** **id** **assign** **id** **semicolon** **id** **assign** *Expr* **end** ⑤
⇒ **begin** **id** **assign** **id** **semicolon** **id** **assign** *Expr* **plus** **id** **end** ⑥

Parse Tree



Key

terminal

Nonterminal

Rule used

Productions

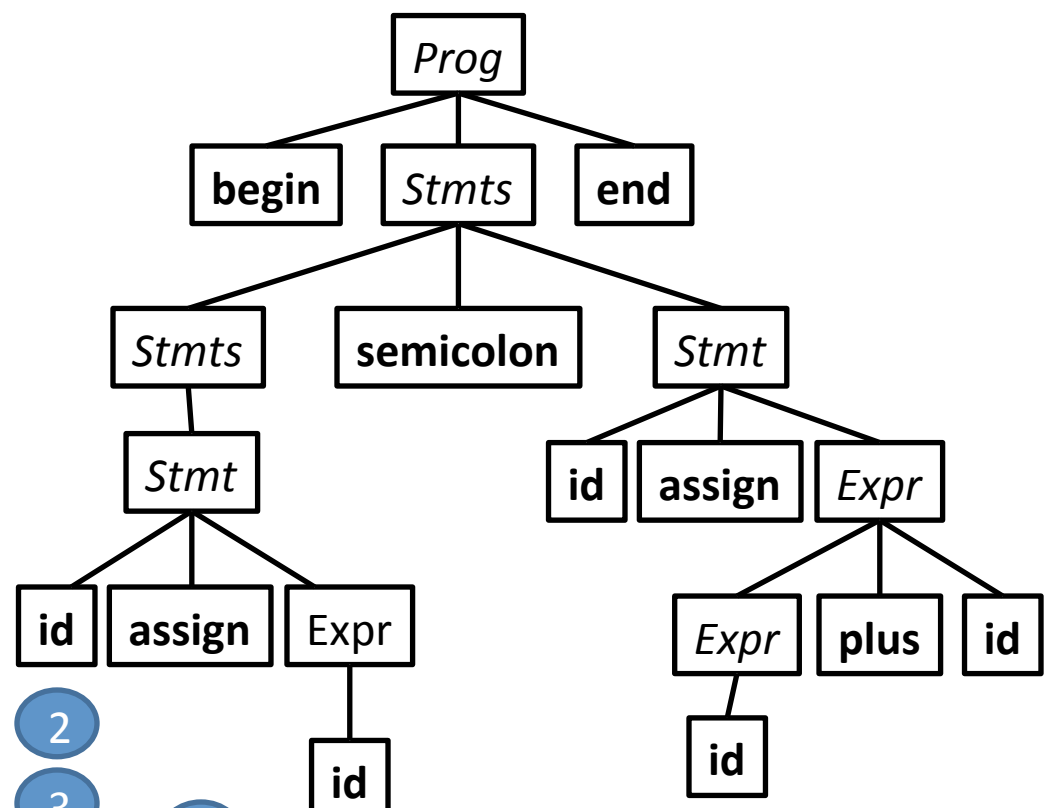
- 1. *Prog* → **begin** *Stmts* **end**
- 2. *Stmts* → *Stmts* **semicolon** *Stmt*
- 3. | *Stmt*
- 4. *Stmt* → **id** **assign** *Expr*
- 5. *Expr* → **id**
- 6. | *Expr* **plus** **id**

Derivation Sequence

Prog ⇒ **begin** *Stmts* **end** ①
⇒ **begin** *Stmts* **semicolon** *Stmt* **end** ②
⇒ **begin** *Stmt* **semicolon** *Stmt* **end** ③
⇒ **begin** **id** **assign** *Expr* **semicolon** *Stmt* **end** ④
⇒ **begin** **id** **assign** *Expr* **semicolon** **id** **assign** *Expr* **end** ④
⇒ **begin** **id** **assign** **id** **semicolon** **id** **assign** *Expr* **end** ⑤
⇒ **begin** **id** **assign** **id** **semicolon** **id** **assign** *Expr* **plus** **id** **end** ⑥
⇒ **begin** **id** **assign** **id** **semicolon** **id** **assign** **id** **plus** **id** **end** ⑤

One of the sequences accepted by the CFG.

Parse Tree



Key

terminal

Nonterminal

Rule used

A five minute introduction

MAKEFILE

Makefiles: Motivation

- Typing the series of commands to generate our code can be tedious
 - Multiple steps that depend on each other
 - Somewhat complicated commands
 - May not need to rebuild everything
- Makefiles solve these issues
 - Record a series of commands in a script-like DSL
 - Specify dependency rules and Make generates the results

Makefiles: Basic Structure

<target>: <dependency list>

(tab) <command to satisfy target>

Makefiles: Basic Structure

<target>: <dependency list>

(tab) <command to satisfy target>

Example

```
Example.class: Example.java IO.class  
    javac Example.java
```

```
IO.class: IO.java  
    javac IO.java
```

Makefiles: Basic Structure

<target>: <dependency list>

(tab) <command to satisfy target>

Example

Example.class depends on example.java and IO.class

```
Example.class: Example.java IO.class
    javac Example.java
```

```
IO.class: IO.java
    javac IO.java
```

Makefiles: Basic Structure

<target>: <dependency list>

(tab) <command to satisfy target>

Example

Example.class depends on example.java and IO.class

```
Example.class: Example.java IO.class
```

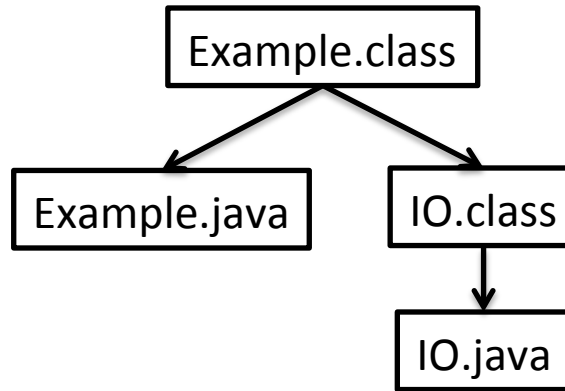
```
    javac Example.java
```

Example.class is generated by
javac Example.java

```
IO.class: IO.java
```

```
    javac IO.java
```

Makefiles: Dependencies



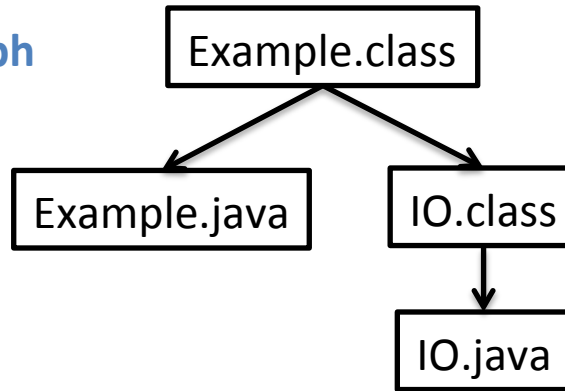
Example

```
Example.class: Example.java IO.class  
    javac Example.java
```

```
IO.class: IO.java  
    javac IO.java
```

Makefiles: Dependencies

Internal Dependency graph



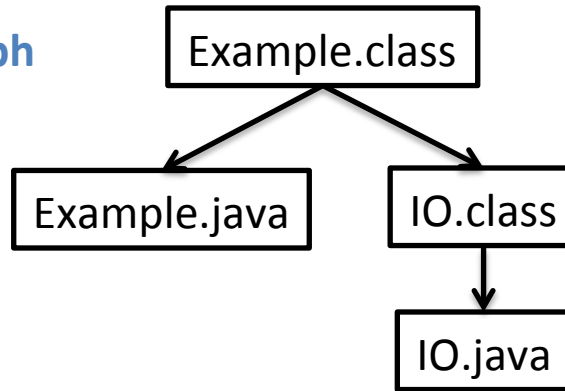
Example

```
Example.class: Example.java IO.class  
    javac Example.java
```

```
IO.class: IO.java  
    javac IO.java
```

Makefiles: Dependencies

Internal Dependency graph



A file is rebuilt if one of its dependencies changes

Example

```
Example.class: Example.java IO.class
               javac Example.java
```

```
IO.class: IO.java
          javac IO.java
```

Makefiles: Variables

You can thread common configuration values through your makefile

Makefiles: Variables

You can thread common configuration values through your makefile

Example

JC = /s/std/bin/javac

JFLAGS = -g

Makefiles: Variables

You can thread common configuration values through your makefile

Example

JC = /s/std/bin/javac

JFLAGS = -g **Build for debug**

Makefiles: Variables

You can thread common configuration values through your makefile

Example

```
JC = /s/std/bin/javac
```

```
JFLAGS = -g Build for debug
```

```
Example.class: Example.java IO.class  
    $(JC) $(JFLAGS) Example.java
```

```
IO.class: IO.java  
    $(JC) $(JFLAGS) IO.java
```

Makefiles: Phony Targets

- You can run commands through make.
 - Write a target with no dependencies (called phony)
 - Will cause it to execute the command every time



Makefiles: Phony Targets

- You can run commands through make.
 - Write a target with no dependencies (called phony)
 - Will cause it to execute the command every time

Example

```
clean:  
    rm -f *.class
```



Makefiles: Phony Targets

- You can run commands through make.
 - Write a target with no dependencies (called phony)
 - Will cause it to execute the command every time

Example

`clean:`

```
rm -f *.class
```

`test:`

```
java -cp . Test.class
```



Recap

- We've defined context-free grammars
 - More powerful than regular expressions
- Learned a bit about makefile
- Next time we'll look at grammars in more detail