# Context-Free Grammars

# Describing Languages

- We've seen two models for the regular languages:

  - *Finite automata* accept precisely the strings in the language.

  - *Regular expressions* describe precisely the strings in the language.

- Finite automata *recognize* strings in the language.

  - Perform a computation to determine whether a specific string is in the language.

- Regular expressions *match* strings in the language.

  - Describe the general shape of all strings in the language.

# Context-Free Grammars

- A ***context-free grammar*** (or ***CFG***) is an entirely different formalism for defining a class of languages.

- Goal: Give a procedure for listing off all strings in the language.

- CFGs are best explained by example...

# Arithmetic Expressions

- Suppose we want to describe all legal arithmetic expressions using addition, subtraction, multiplication, and division.

- Here is one possible CFG:

$E \rightarrow$ int

$E \rightarrow E$ Op $E$

$E \rightarrow (E)$

Op $\rightarrow +$

Op $\rightarrow -$

Op $\rightarrow *$

Op $\rightarrow /$

$E$
$\Rightarrow E$ Op $E$
$\Rightarrow E$ Op $(E)$
$\Rightarrow E$ Op $(E$ Op $E)$
$\Rightarrow E * (E$ Op $E)$
$\Rightarrow$ int $* (E$ Op $E)$
$\Rightarrow$ int $* ($int Op $E)$
$\Rightarrow$ int $* ($int Op int$)$
$\Rightarrow$ int $* ($int $+$ int$)$

# Arithmetic Expressions

- Suppose we want to describe all legal arithmetic expressions using addition, subtraction, multiplication, and division.

- Here is one possible CFG:

E → int

E → E Op E

E → (E)

Op → +

Op → -

Op → *

Op → /

E
⇒ E Op E
⇒ E Op int
⇒ int Op int
⇒ int / int

# Context-Free Grammars

- Formally, a context-free grammar is a collection of four objects:

  - A set of ***nonterminal symbols*** (also called ***variables***),

  - A set of ***terminal symbols*** (the ***alphabet*** of the CFG)

  - A set of ***production rules*** saying how each nonterminal can be replaced by a string of terminals and nonterminals, and

  - A ***start symbol*** (which must be a nonterminal) that begins the derivation.

$$E \rightarrow \texttt{int}$$

$$E \rightarrow E \; Op \; E$$

$$E \rightarrow \texttt{(}E\texttt{)}$$

$$Op \rightarrow \texttt{+}$$

$$Op \rightarrow \texttt{-}$$

$$Op \rightarrow \texttt{*}$$

$$Op \rightarrow \texttt{/}$$

# Some CFG Notation

- Capital letters in **Bold Red Uppercase** will represent nonterminals.

  - i.e. **A**, **B**, **C**, **D**

- Lowercase letters in `blue monospace` will represent terminals.

  - i.e. `t`, `u`, `v`, `w`

- Lowercase Greek letters in *gray italics* will represent arbitrary strings of terminals and nonterminals.

  - i.e. *α*, *γ*, *ω*

# A Notational Shorthand

E → int

E → E Op E

E → (E)

Op → +

Op → −

Op → *

Op → /

# A Notational Shorthand

$$E \rightarrow \texttt{int} \mid E \; Op \; E \mid (E)$$
$$Op \rightarrow \texttt{+} \mid \texttt{-} \mid \texttt{*} \mid \texttt{/}$$

# Derivations

$$E \rightarrow E\ Op\ E \mid \texttt{int} \mid (E)$$
$$Op \rightarrow \texttt{+} \mid \texttt{*} \mid \texttt{-} \mid \texttt{/}$$

E

$\Rightarrow$ E Op E

$\Rightarrow$ E Op (E)

$\Rightarrow$ E Op (E Op E)

$\Rightarrow$ E * (E Op E)

$\Rightarrow$ int * (E Op E)

$\Rightarrow$ int * (int Op E)

$\Rightarrow$ int * (int Op int)

$\Rightarrow$ int * (int + int)

- A sequence of steps where nonterminals are replaced by the right-hand side of a production is called a *derivation*.

- If string $\alpha$ derives string $\omega$, we write $\alpha \Rightarrow^* \omega$.

- In the example on the left, we see $E \Rightarrow^* $ int * (int + int).

# The Language of a Grammar

- If *G* is a CFG with alphabet $\Sigma$ and start symbol **S**, then the ***language of G*** is the set

$$\mathscr{L}(G) = \{\ \omega \in \Sigma^* \mid S \Rightarrow^* \omega\ \}$$

- That is, $\mathscr{L}(G)$ is the set of strings derivable from the start symbol.

- Note: $\omega$ must be in $\Sigma^*$, the set of strings made from terminals. Strings involving nonterminals aren't in the language.

# Context-Free Languages

- A language $L$ is called a ***context-free language*** (or CFL) if there is a CFG $G$ such that $L = \mathscr{L}(G)$.

- Questions:

  - What languages are context-free?

  - How are context-free and regular languages related?

# From Regexes to CFGs

- CFGs consist purely of production rules of the form **A** → *ω*. They do not have the regular expression operators * or ∪.

- However, we can convert regular expressions to CFGs as follows:

$$\textbf{S} \rightarrow \texttt{a*b}$$

# From Regexes to CFGs

- CFGs consist purely of production rules of the form **A** → *ω*. They do not have the regular expression operators * or ∪.

- However, we can convert regular expressions to CFGs as follows:

$$S \rightarrow Ab$$
$$A \rightarrow Aa \mid \varepsilon$$

# From Regexes to CFGs

- CFGs consist purely of production rules of the form **A** → *ω*. They do not have the regular expression operators * or ∪.

- However, we can convert regular expressions to CFGs as follows:

$$\textbf{S} \rightarrow \texttt{a(b} \cup \texttt{c*)}$$

# From Regexes to CFGs

- CFGs consist purely of production rules of the form $A \to \omega$. They do not have the regular expression operators * or ∪.

- However, we can convert regular expressions to CFGs as follows:

$$S \to aX$$
$$X \to b \mid C$$
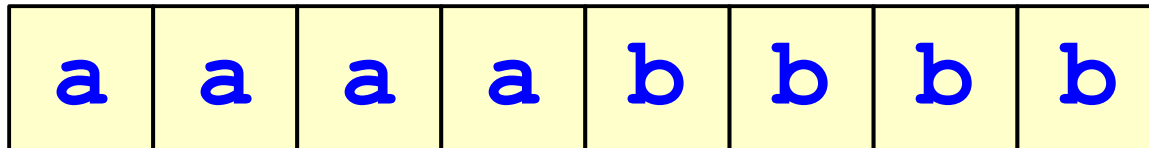$$C \to Cc \mid \varepsilon$$

# Regular Languages and CFLs

- ***Theorem:*** Every regular language is context-free.

- ***Proof Idea:*** Use the construction from the previous slides to convert a regular expression for $L$ into a CFG for $L$. ∎

- ***Problem Set Exercise:*** Instead, show how to convert a DFA/NFA into a CFG.
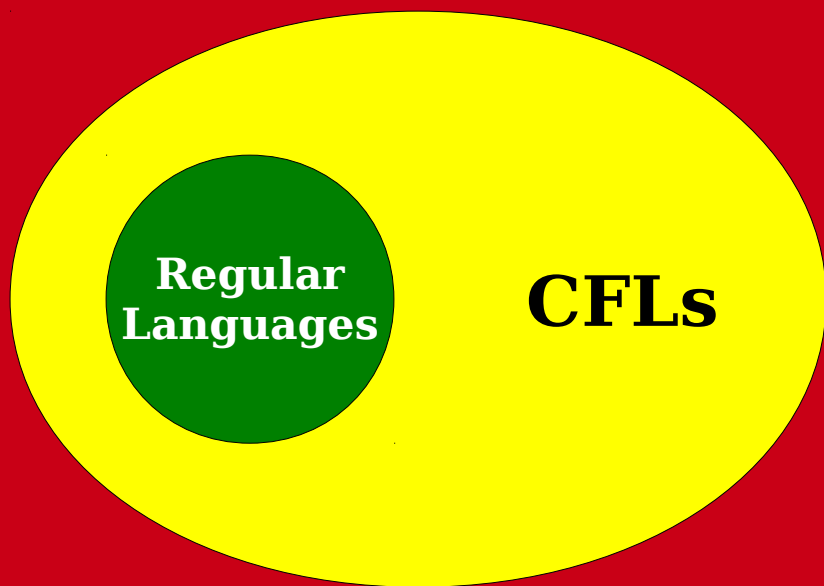
# The Language of a Grammar

- Consider the following CFG *G*:

$$S \rightarrow aSb \mid \varepsilon$$

- What strings can this generate?

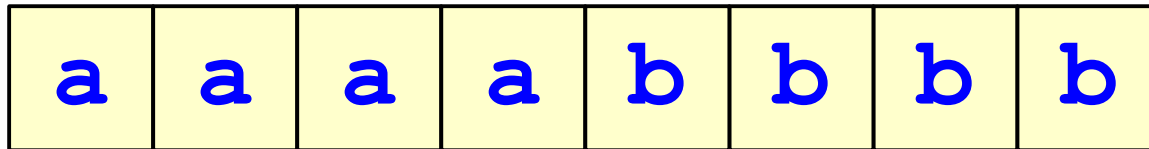| a | a | a | a | b | b | b | b |
|---|---|---|---|---|---|---|---|

$$\mathcal{L}(G) = \{\ a^n b^n \mid n \in \mathbb{N}\ \}$$

# Why the Extra Power?

- Why do CFGs have more power than regular expressions?

- *Intuition:* Derivations of strings have unbounded "memory."

$$S \rightarrow aSb \mid \varepsilon$$

| a | a | a | a | b | b | b | b |
|---|---|---|---|---|---|---|---|

# Time-Out for Announcements!

# Problem Set Seven

- Problem Set Six was due at the start of today's lecture.

  - Want to use late days? Submit up to Monday at 3:00PM.

- Problem Set Seven goes out now. It's due next Friday.

  - Play around with the Myhill-Nerode theorem and the limits of regular languages!

  - Play around with your very own CFGs!

# Midterms Graded

- Midterms have been graded. They're available for pickup in the Gates building.

  - SCPD students: we've sent the exams back to the SCPD office. You should hear back from them soon.

- Solutions and stats are available in the Gates building in the normal handout filing cabinet.

# Midterm Regrades

- If you believe that we made a grading error on the exam, you can submit it for a regrade. To do so, fill out the form online, staple it to your exam, and hand it to Keith by next Friday.

- Please only submit regrades if you
  - believe that we actually graded your exam incorrectly, and
  - you've talked about the exam with the course staff and they agree with you.

- ***Your score can go down if you ask for a regrade***. Please be sure you really want to ask for it before you submit a regrade request.

# Back to CS103!

# Designing CFGs

- Like designing DFAs, NFAs, and regular expressions, designing CFGs is a craft.

- When thinking about CFGs:

  - ***Think recursively:*** Build up bigger structures from smaller ones.

  - ***Have a construction plan:*** Know in what order you will build up the string.

  - ***Store information in nonterminals:*** Have each nonterminal correspond to some useful piece of information.

# Designing CFGs

- Let $\Sigma = \{$**a**, **b**$\}$ and let $L = \{w \in \Sigma^* \mid w$ is a palindrome $\}$

- We can design a CFG for $L$ by thinking inductively:

  - Base case: $\varepsilon$, **a**, and **b** are palindromes.

  - If $\omega$ is a palindrome, then **a**$\omega$**a** and **b**$\omega$**b** are palindromes.

$$S \rightarrow \varepsilon \mid a \mid b \mid aSa \mid bSb$$

# Designing CFGs

- Let Σ = { **(**, **)** } and let $L = \{w \in \Sigma^* \mid w$ is a string of balanced parentheses $\}$

- Some sample strings in $L$:

$$\text{((()))}$$

$$\text{(())()}$$

$$\text{(()())(()())}$$

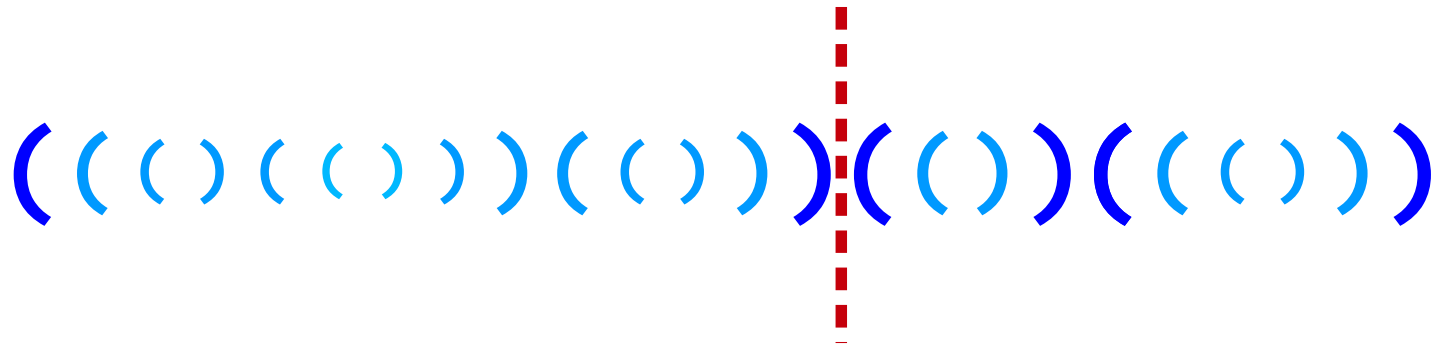$$\text{(((()))(()))}$$

$$\varepsilon$$

$$\text{()()}$$

# Designing CFGs

- Let $\Sigma = \{$ (, ) $\}$ and let $L = \{w \in \Sigma^* \mid w$ is a string of balanced parentheses $\}$

- Let's think about this recursively.

  - Base case: the empty string is a string of balanced parentheses.

  - Recursive step: Look at the closing parenthesis that matches the first open parenthesis.

$$((\,)(\,))(\,(\,))(\,)((\,))$$

# Designing CFGs

- Let Σ = { **(**, **)** } and let $L$ = {$w$ ∈ Σ* | $w$ is a string of balanced parentheses }

- Let's think about this recursively.

  - Base case: the empty string is a string of balanced parentheses.

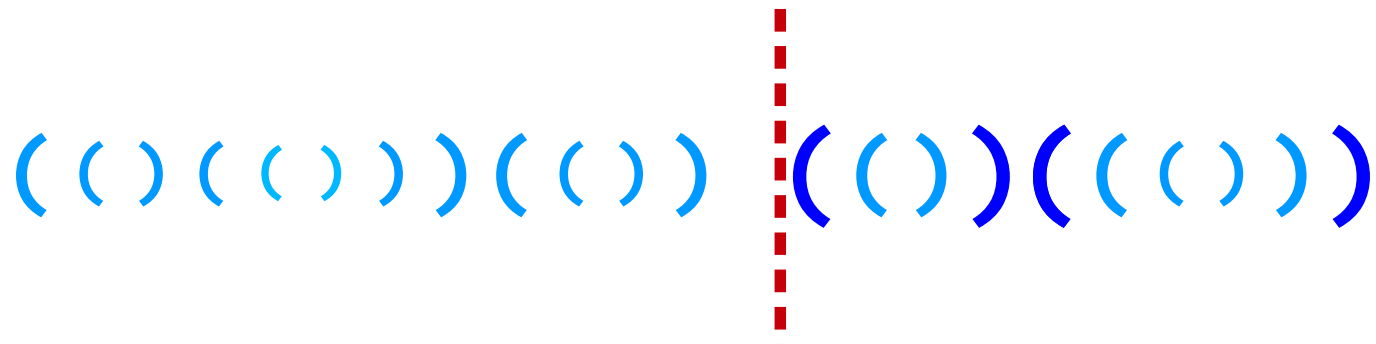  - Recursive step: Look at the closing parenthesis that matches the first open parenthesis.

# Designing CFGs

- Let $\Sigma = \{$ (, ) $\}$ and let $L = \{ w \in \Sigma^* \mid w$ is a string of balanced parentheses $\}$

- Let's think about this recursively.

  - Base case: the empty string is a string of balanced parentheses.

  - Recursive step: Look at the closing parenthesis that matches the first open parenthesis.

( ) ( ( ) ) ) ( ( ) ) ( ( ) ) ( ( ( ) ) )

# Designing CFGs

- Let Σ = { **(**, **)** } and let $L$ = { $w \in \Sigma^*$ | $w$ is a string of balanced parentheses }

- Let's think about this recursively.

  - Base case: the empty string is a string of balanced parentheses.

  - Recursive step: Look at the closing parenthesis that matches the first open parenthesis. Removing the first parenthesis and the matching parenthesis forms two new strings of balanced parentheses.

$$S \rightarrow (S)S \mid \varepsilon$$

# Designing CFGs: A Caveat

- Let $\Sigma = \{$a, b$\}$ and let $L = \{w \in \Sigma^* \mid w$ has the same number of a's and b's $\}$

- Is this a CFG for $L$?

$$S \to aSb \mid bSa \mid \varepsilon$$

- Can you derive the string abba?

# Designing CFGs: A Caveat

- When designing a CFG for a language, make sure that it
    - generates all the strings in the language and
    - never generates a string outside the language.
- The first of these can be tricky – make sure to test your grammars!
- You'll design your own CFG for this language on the next problem set.

# CFG Caveats II

- Is the following grammar a CFG for the language { $a^n b^n \mid n \in \mathbb{N}$ }?

$$S \to aSb$$

- What strings can you derive?

  - Answer: **None!**

- What is the language of the grammar?

  - Answer: **Ø**

- When designing CFGs, make sure your recursion actually terminates!

# CFG Caveats III

- When designing CFGs, remember that each nonterminal can be expanded out independently of the others.

- Let $\Sigma = \{\mathbf{a}, \overset{?}{=}\}$ and let $L = \{\mathbf{a}^n \overset{?}{=} \mathbf{a}^n \mid n \in \mathbb{N}\}$.

- Is the following a CFG for $L$?

$$\mathbf{S} \to \mathbf{X}\overset{?}{=}\mathbf{X}$$

$$\mathbf{X} \to \mathbf{a}\mathbf{X} \mid \varepsilon$$

$\mathbf{S}$
$\Rightarrow \mathbf{X}\overset{?}{=}\mathbf{X}$
$\Rightarrow \mathbf{a}\mathbf{X}\overset{?}{=}\mathbf{X}$
$\Rightarrow \mathbf{a}\mathbf{a}\mathbf{X}\overset{?}{=}\mathbf{X}$
$\Rightarrow \mathbf{a}\mathbf{a}\overset{?}{=}\mathbf{X}$
$\Rightarrow \mathbf{a}\mathbf{a}\overset{?}{=}\mathbf{a}\mathbf{X}$
$\Rightarrow \mathbf{a}\mathbf{a}\overset{?}{=}\mathbf{a}$

# Finding a Build Order

- Let $\Sigma = \{\mathtt{a}, \stackrel{?}{=}\}$ and let $L = \{\mathtt{a}^n \stackrel{?}{=} \mathtt{a}^n \mid n \in \mathbb{N}\}$.

- To build a CFG for $L$, we need to be more clever with how we construct the string.

  - If we build the strings of $\mathtt{a}$'s independently of one another, then we can't enforce that they have the same length.

  - *Idea:* Build both strings of $\mathtt{a}$'s at the same time.

- Here's one possible grammar based on that idea:

$$S \to \stackrel{?}{=} \mid \mathtt{a}S\mathtt{a}$$

$$S$$
$$\Rightarrow \mathtt{a}S\mathtt{a}$$
$$\Rightarrow \mathtt{aa}S\mathtt{aa}$$
$$\Rightarrow \mathtt{aaa}S\mathtt{aaa}$$
$$\Rightarrow \mathtt{aaa}\stackrel{?}{=}\mathtt{aaa}$$

# Function Prototypes

- Let Σ = {**void**, **int**, **double**, **name**, **(**, **)**, **,**, **;**}.

- Let's write a CFG for C-style function prototypes!

- Examples:

  - **void name(int name, double name);**

  - **int name();**

  - **int name(double name);**

  - **int name(int, int name, int);**

  - **void name(void);**

# Function Prototypes

- Here's one possible grammar:

  - **S** → **Ret** `name` **(Args);**

  - **Ret** → **Type** | `void`

  - **Type** → `int` | `double`

  - **Args** → ε | `void` | **ArgList**

  - **ArgList** → **OneArg** | **ArgList**, **OneArg**

  - **OneArg** → **Type** | **Type** `name`

- Fun question to think about: what changes would you need to make to support pointer types?

# Summary of CFG Design Tips

- Look for recursive structures where they exist: they can help guide you toward a solution.

- Keep the build order in mind – often, you'll build two totally different parts of the string concurrently.

    - Usually, those parts are built in opposite directions: one's built left-to-right, the other right-to-left.

- Use different nonterminals to represent different structures.

# Applications of Context-Free Grammars

# CFGs for Programming Languages

$$
\begin{array}{lll}
\textbf{BLOCK} & \rightarrow & \textbf{STMT} \\
& | & \{ \textbf{STMTS} \} \\
\\
\textbf{STMTS} & \rightarrow & \boldsymbol{\varepsilon} \\
& | & \textbf{STMT STMTS} \\
\\
\textbf{STMT} & \rightarrow & \textbf{EXPR}; \\
& | & \texttt{if (}\textbf{EXPR}\texttt{) }\textbf{BLOCK} \\
& | & \texttt{while (}\textbf{EXPR}\texttt{) }\textbf{BLOCK} \\
& | & \texttt{do }\textbf{BLOCK}\texttt{ while (}\textbf{EXPR}\texttt{)}; \\
& | & \textbf{BLOCK} \\
& | & \ldots \\
\\
\textbf{EXPR} & \rightarrow & \texttt{identifier} \\
& | & \texttt{constant} \\
& | & \textbf{EXPR + EXPR} \\
& | & \textbf{EXPR -- EXPR} \\
& | & \textbf{EXPR * EXPR} \\
& | & \ldots
\end{array}
$$

# Grammars in Compilers

- One of the key steps in a compiler is figuring out what a program "means."

- This is usually done by defining a grammar showing the high-level structure of a programming language.

- There are certain classes of grammars (LL(1) grammars, LR(1) grammars, LALR(1) grammars, etc.) for which it's easy to figure out how a particular string was derived.

- Tools like `yacc` or `bison` automatically generate parsers from these grammars.

- Curious to learn more? Take CS143!

# Natural Language Processing

- By building context-free grammars for actual languages and applying statistical inference, it's possible for a computer to recover the likely meaning of a sentence.
  - In fact, CFGs were first called ***phrase-structure grammars*** and were introduced by Noam Chomsky in his seminal work *Syntactic Structures*.
  - They were then adapted for use in the context of programming languages, where they were called ***Backus-Naur forms***.
- Stanford's <span style="color:orangered">CoreNLP project</span> is one place to look for an example of this.
- Want to learn more? Take CS124 or CS224N!

# Next Time

- **Turing Machines**

  - What does a computer with unbounded memory look like?

  - How do you program them?