# Context-Free Grammars

# A Motivating Question

```
Terminal — python — 66×21

>>> (26 + 42) * 2 + 1
```

How does my computer know what this sequence of characters means? How can it determine whether or not this expression is even syntactically valid?

# An Analogy: Mad Libs

**THE MAGIC COMPUTERS**

Today, every student has a computer small enough to fit into his

_____. He can solve any math problem by simply
        NOUN

pushing the computer's little _____. Computers
                                PLURAL NOUN

can add, multiply, divide, and _____. They
                                VERB (PRESENT TENSE)

can also _____ better than a human. Some com-
          VERB (PRESENT TENSE)

puters are _____. Others have a/an
            PART OF BODY (PLURAL)

                                        'inds of _____
                                                  PLURAL NOUN

When you're filling out Mad Libs, you have these **placeholders** for different parts of speech.

# Mad Libs for Arithmetic Expressions

Imagine I have a template like this:

**(** ___ ___ ___ **)** ___ ___ ___ ___
$\quad\ $ **Int** $\ $ **Op** $\ $ **Int** $\qquad\ $ **Op** $\ $ **Int** $\ $ **Op** $\ $ **Int**

# Mad Libs for Arithmetic Expressions

Here's one way I could fill it out:

( 26 + 42 ) * 2 + 1

$$\underset{\text{Int}}{26} \quad \underset{\text{Op}}{+} \quad \underset{\text{Int}}{42} \quad \underset{\text{Op}}{*} \quad \underset{\text{Int}}{2} \quad \underset{\text{Op}}{+} \quad \underset{\text{Int}}{1}$$

# Mad Libs for Arithmetic Expressions

Here's another:

$$( \quad \underline{7} \quad \underline{*} \quad \underline{5} \quad ) \quad \underline{/} \quad \underline{5} \quad \underline{-} \quad \underline{49}$$

$$\text{Int} \quad \text{Op} \quad \text{Int} \qquad \text{Op} \quad \text{Int} \quad \text{Op} \quad \text{Int}$$

Imagine you have a computer that's pre-programmed with this template.

You could then enter a string and be able to check whether it is valid. You can also understand what individual pieces of the string mean based on which part of the template they're filling in.

# Mad Libs for Arithmetic Expressions

This is nice but I can only make expressions of the form **(Int Op Int) Op Int Op Int**

**(** ___ ___ ___ **)** ___ ___ ___ ___
    **Int   Op   Int**      **Op   Int   Op   Int**

But there are many valid arithmetic expressions that don't follow this pattern!

# Mad Libs for Arithmetic Expressions

Idea: could we come up with a set of rules for generating valid arithmetic Mad Libs templates?

Eg. **Int Op Int**, **(Int Op (Int Op Int))**, **(Int Op Int) Op (Int Op Int)** ...

# 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 description of a language by recursively describing the structure of the 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

$E \rightarrow E\ Op\ E$

$E \rightarrow (E)$

$Op \rightarrow +$

$Op \rightarrow -$

$Op \rightarrow \times$

$Op \rightarrow /$

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

# Arithmetic Expressions

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

- Here is one possible CFG:

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

$$E$$
$$\Rightarrow E \; Op \; E$$
$$\Rightarrow E \; Op \; \texttt{int}$$
$$\Rightarrow \texttt{int} \; Op \; \texttt{int}$$
$$\Rightarrow \texttt{int / int}$$

# Context-Free Grammars

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

    - 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. By convention, the start symbol is the one on the left-hand side of the first production.

E → int
E → E Op E
E → (E)
Op → +
Op → −
Op → ×
Op → /

# Some CFG Notation

- In today's slides, capital letters in **Bold Red Uppercase** will represent nonterminals.
  - e.g. **A**, **B**, **C**, **D**
- Lowercase letters in `blue monospace` will represent terminals.
  - e.g. `t`, `u`, `v`, `w`
- Lowercase Greek letters in *gray italics* will represent arbitrary strings of terminals and nonterminals.
  - e.g. *α*, *γ*, *ω*
- You don't need to use these conventions on your own; just make sure whatever you do is readable. ☺

# A Notational Shorthand

$E \rightarrow$ int

$E \rightarrow E\ Op\ E$

$E \rightarrow (E)$

$Op \rightarrow +$

$Op \rightarrow -$

$Op \rightarrow \times$

$Op \rightarrow /$

# A Notational Shorthand

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

# Derivations

```
E → E Op E | int | (E)
Op → + | × | - | /
```

    E

⇒ E Op E

⇒ E Op (E)

⇒ E Op (E Op E)

⇒ E × (E Op E)

⇒ int × (E Op E)

⇒ int × (int Op E)

⇒ int × (int Op int)

⇒ 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

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

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

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

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

Consider the following CFG $G$ over $\Sigma = \{a, b, c, d\}$:

$$S \to Sa \mid dT$$
$$T \to bTb \mid c$$

How many of the following strings are in $\mathcal{L}(G)$?

dca

cad

bcb

dTaa

# Context-Free Languages

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

- 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:

$$S \rightarrow a*b$$

# 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:

$$S \rightarrow a\text{*}b$$

$$A \rightarrow Aa \mid \varepsilon$$

# 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:

$$\mathbf{S} \rightarrow \texttt{a*b}$$
$$\mathbf{A} \rightarrow \mathbf{A}\texttt{a} \mid \varepsilon$$

# 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:

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

# From Regexes to CFGs

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

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

$$S \rightarrow a(b \cup c^*)$$

# 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:

$$S \to a(b \cup c*)$$
$$X \to b \mid c*$$

# From Regexes to CFGs

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

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

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

# 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:

$$S \to aX$$
$$X \to b \mid c*$$

# 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:

$$S → aX$$
$$X → b \mid c*$$

# 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:

$$S \rightarrow aX$$
$$X \rightarrow b \mid c*$$
$$C \rightarrow Cc \mid \varepsilon$$

# From Regexes to CFGs

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

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

$$S \rightarrow aX$$
$$X \rightarrow b \mid c*$$
$$C \rightarrow Cc \mid \varepsilon$$

# 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 $\cup$.

- 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 $L$. ∎

- ***Great 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?

# The Language of a Grammar

- Consider the following CFG $G$:

$$S \rightarrow aSb \mid \varepsilon$$
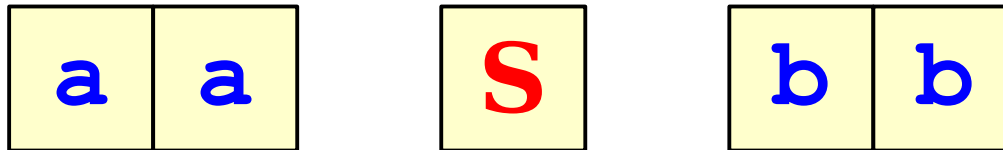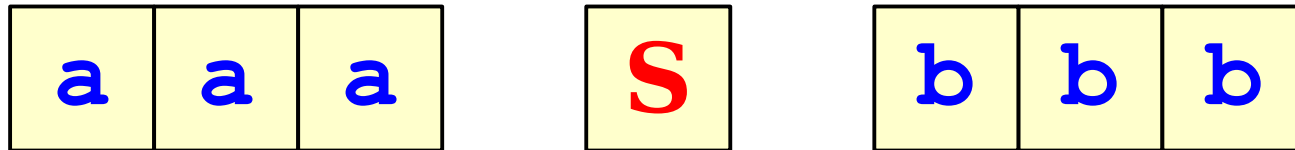
- What strings can this generate?

S

# The Language of a Grammar

- Consider the following CFG *G*:

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

- What strings can this generate?

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

# The Language of a Grammar

- Consider the following CFG *G*:

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

- What strings can this generate?

a  S  b

# The Language of a Grammar

- Consider the following CFG $G$:

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

- What strings can this generate?

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

# The Language of a Grammar

- Consider the following CFG *G*:

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

- What strings can this generate?

| a | a |   | S |   | b | b |

# The Language of a Grammar

- Consider the following CFG $G$:

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

- What strings can this generate?

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

# The Language of a Grammar

- Consider the following CFG *G*:

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

- What strings can this generate?

| a | a | a | | S | | b | b | b |

# The Language of a Grammar

- Consider the following CFG $G$:

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

- What strings can this generate?

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

# The Language of a Grammar

- Consider the following CFG *G*:

$$\mathbf{S} \to \mathbf{aSb} \mid \varepsilon$$

- What strings can this generate?

| a | a | a | a |
|---|---|---|---|

| b | b | b | b |
|---|---|---|---|

# 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 |
|---|---|---|---|---|---|---|---|

# 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."

# Why the Extra Power?

- Why do CFGs have more power than regular expressions?

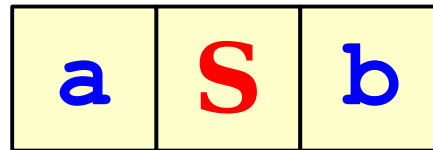- *Intuition:* Derivations of strings have unbounded "memory."

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

# Why the Extra Power?

- Why do CFGs have more power than regular expressions?

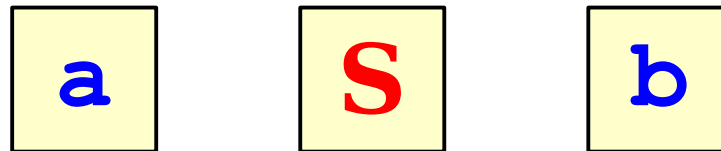- ***Intuition:*** Derivations of strings have unbounded "memory."

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

S

# Why the Extra Power?

- Why do CFGs have more power than regular expressions?

- *Intuition:* Derivations of strings have unbounded "memory."
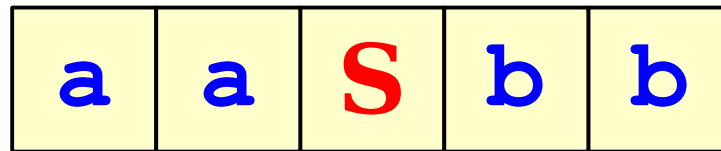
$$S \rightarrow aSb \mid \varepsilon$$

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

# Why the Extra Power?

- Why do CFGs have more power than regular expressions?

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

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

a  S  b

# Why the Extra Power?

- Why do CFGs have more power than regular expressions?

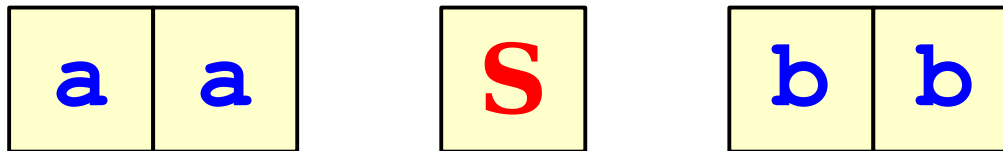- ***Intuition:*** Derivations of strings have unbounded "memory."

$$S \rightarrow aSb \mid \varepsilon$$
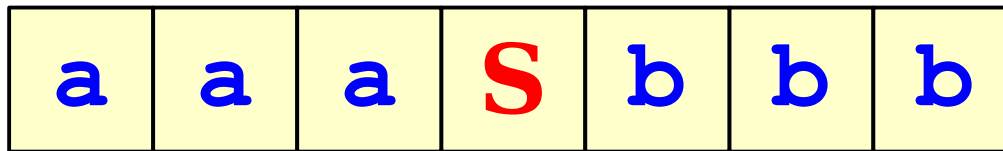
| a | a | S | b | b |
|---|---|---|---|---|

# 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 |
|---|---|

| S |
|---|

| b | b |
|---|---|

# Why the Extra Power?

- Why do CFGs have more power than regular expressions?

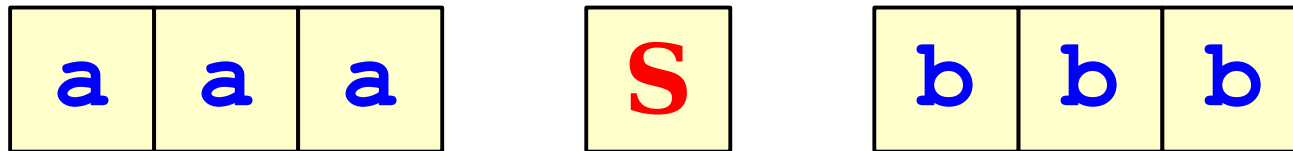- *Intuition:* Derivations of strings have unbounded "memory."

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

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

# Why the Extra Power?

- Why do CFGs have more power than regular expressions?

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

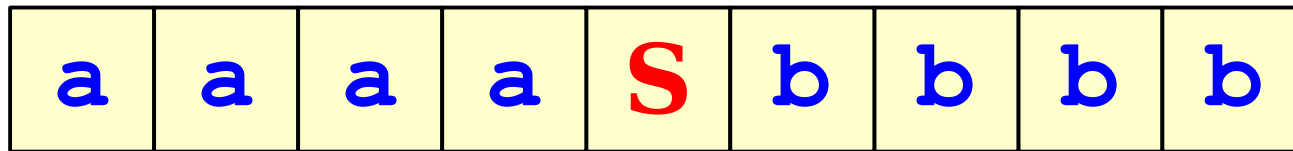$$S \rightarrow aSb \mid \varepsilon$$

| a | a | a |
|---|---|---|

| S |
|---|

| b | b | b |
|---|---|---|

# 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 | S | b | b | b | b |
|---|---|---|---|---|---|---|---|---|

# Why the Extra Power?

- Why do CFGs have more power than regular expressions?

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

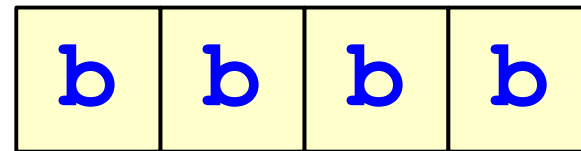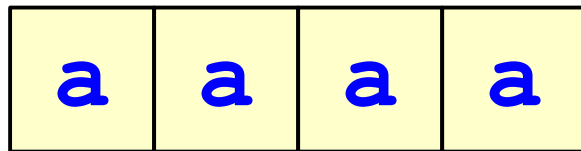$$S \rightarrow aSb \mid \varepsilon$$

| a | a | a | a |
|---|---|---|---|

| b | b | b | b |
|---|---|---|---|

# Why the Extra Power?

- Why do CFGs have more power than regular expressions?

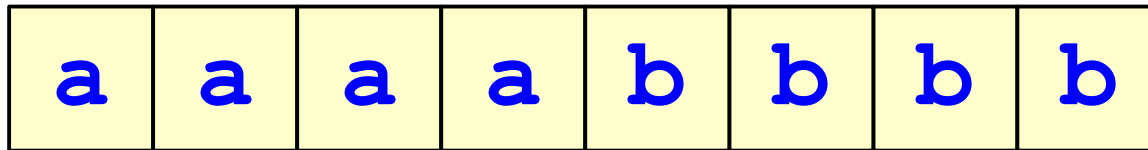- *Intuition:* Derivations of strings have unbounded "memory."

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

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

Let's take a five minute break!

# 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: ε, a, and b are palindromes.

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

  - No other strings are palindromes.

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

# 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:

**S → ε | a | b | aSa | bSb**
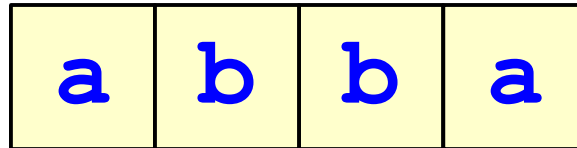
# Designing CFGs

- 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:

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

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

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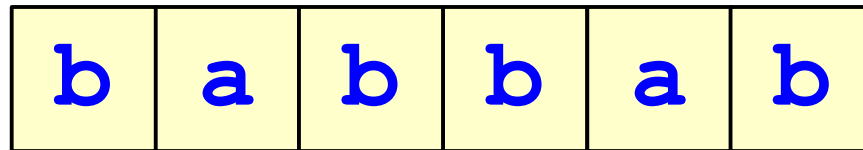
- Let $\Sigma = \{$ a, b $\}$ and let $L = \{w \in \Sigma^* \mid w$ is a palindrome $\}$

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

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

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

# Designing CFGs

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

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

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

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

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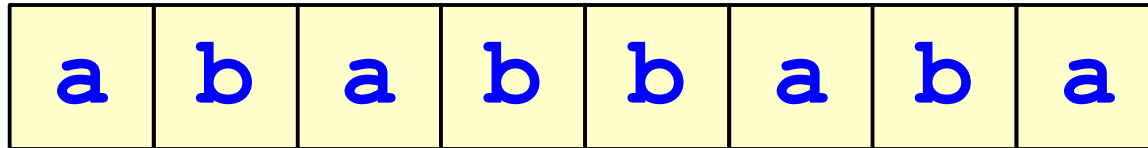
- Let $\Sigma = \{$a, b$\}$ and let $L = \{w \in \Sigma^* \mid w$ is a palindrome $\}$

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

Inductive (building up) perspective: you can take any palindrome and build a larger one by adding the same character to both ends.

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

# 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 recursively:

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

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

# Designing CFGs

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

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



$$\textcolor{red}{S} \rightarrow \textcolor{blue}{\varepsilon} \mid \textcolor{blue}{a} \mid \textcolor{blue}{b} \mid \textcolor{blue}{a}\textcolor{red}{S}\textcolor{blue}{a} \mid \textcolor{blue}{b}\textcolor{red}{S}\textcolor{blue}{b}$$
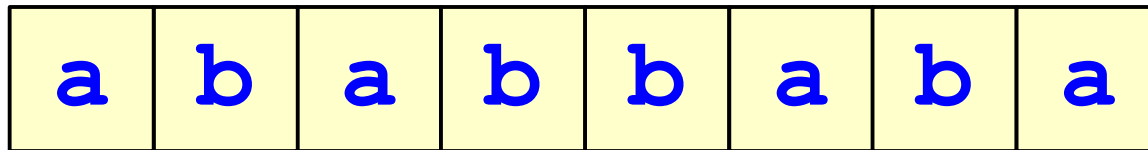
# Designing CFGs

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

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

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

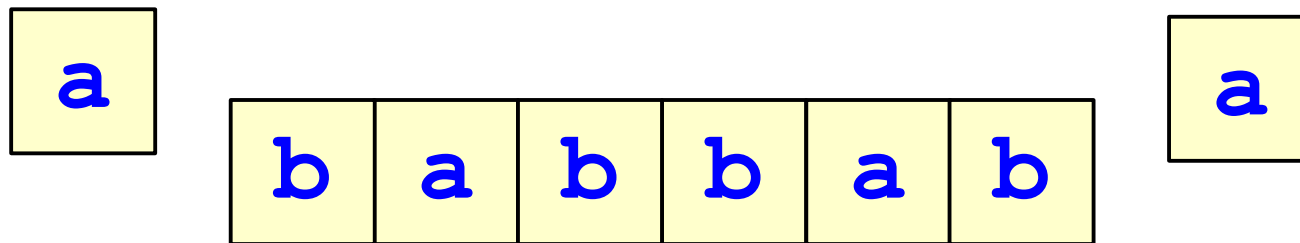$$S \rightarrow \varepsilon \mid a \mid b \mid aSa \mid bSb$$

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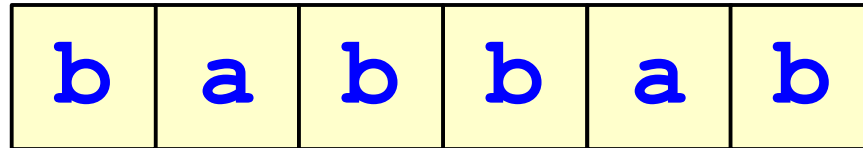
- Let $\Sigma = \{$a, b$\}$ and let $L = \{w \in \Sigma^* \mid w$ is a palindrome $\}$

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

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

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

# Designing CFGs

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

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

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

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

# Designing CFGs

- Let $\Sigma = \{$**a**, **b**$\}$ and let $L = \{w \in \Sigma^* \mid w$ is a palindrome $\}$
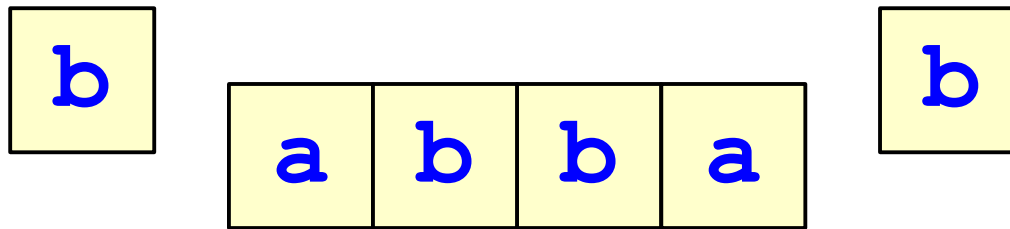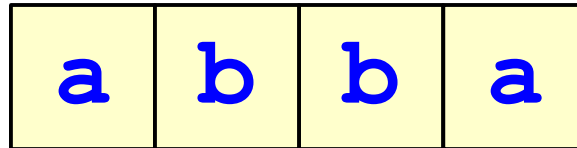
- We can design a CFG for $L$ by thinking recursively:

> Recursive (building down) perspective: you can take any palindrome and repeatedly remove the same character from both ends, leaving behind a palindrome.

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

# Designing CFGs

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

- Some sample strings in $L$:

<p align="center">{{{}}}</p>

<p align="center">{{}}{}</p>

<p align="center">{{}{}}{{}{}}</p>

<p align="center">{{{{}}}{{}}}}</p>

<p align="center">ε</p>

<p align="center">{}{}</p>

# Designing CFGs

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

- Let's think about this recursively.

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

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

$$\textbf{\{\{\{\}\{\{\}\}\}\{\{\}\}\}\{\{\}\}\{\{\{\}\}\}}$$

# Designing CFGs

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

- Let's think about this recursively.

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

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

**{{{}{{}}}{{}}}{{}}{{{}}}**

# Designing CFGs

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

- Let's think about this recursively.

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

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

**{ { { } { { } } } { { } } } { { } } { { { } } }**

# Designing CFGs

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

- Let's think about this recursively.

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

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

$$\{\{\}\{\{\}\}\}\{\{\}\} \mid \{\{\}\}\{\{\}\}\}$$

# Designing CFGs

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

- Let's think about this recursively.

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

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

$$S \rightarrow \textbf{\{}S\textbf{\}}S \mid \varepsilon$$

# Designing CFGs

- Here's the derivation from class today:

  **S**

  ⇒ **{S}S**

  ⇒ **{{S}S}S**

  ⇒ **{{{S}S}S}S**

  ⇒ **{{{S}{S}S}S}S**

  ⇒ **{{{ε}{S}S}S}S**

  ⇒ **{{{ε}{ε}S}S}S**

  ⇒ **{{{ε}{ε}ε}S}S**

  ⇒ **{{{ε}{ε}ε}ε}S**

  ⇒ **{{{ε}{ε}ε}ε}ε**

# Designing CFGs

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

How many of the following CFGs have language $L$?

**S** → **aSb** | **bSa** | **ε**

**S** → **abS** | **baS** | **ε**

**S** → **abSba** | **baSab** | **ε**

**S** → **SbaS** | **SabS** | **ε**

# Designing CFGs

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

How many of the following CFGs have language $L$?

**S → aSb | bSa | ε**

**S → abS | baS | ε**

**S → abSba | baSab | ε**

**S → SbaS | SabS | ε**

# Designing CFGs

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

How many of the following CFGs have language $L$?

**S → aSb | bSa | ε**

**S → abS | baS | ε**

**S → abSba | baSab | ε**

**S → SbaS | SabS | ε**

# Designing CFGs

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

How many of the following CFGs have language $L$?

$S \to \textcolor{blue}{a}\textcolor{red}{S}\textcolor{blue}{b} \mid \textcolor{blue}{b}\textcolor{red}{S}\textcolor{blue}{a} \mid \varepsilon$

$S \to \textcolor{blue}{ab}\textcolor{red}{S} \mid \textcolor{blue}{ba}\textcolor{red}{S} \mid \varepsilon$

$S \to \textcolor{blue}{ab}\textcolor{red}{S}\textcolor{blue}{ba} \mid \textcolor{blue}{ba}\textcolor{red}{S}\textcolor{blue}{ab} \mid \varepsilon$

$S \to \textcolor{red}{S}\textcolor{blue}{ba}\textcolor{red}{S} \mid \textcolor{red}{S}\textcolor{blue}{ab}\textcolor{red}{S} \mid \varepsilon$

# Designing CFGs

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

How many of the following CFGs have language $L$?

S → aSb | bSa | ε

S → abS | baS | ε

S → abSba | baSab | ε

S → SbaS | SabS | ε

# 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 Problem Set 8.

# CFG Caveats II

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

$$\textbf{S} \rightarrow \textbf{aSb}$$

- What strings in {**a**, **b**}* can you derive?

  - Answer: *None!*

- What is the language of the grammar?

  - Answer: Ø

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

# Designing CFGs

- 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 $L$?

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

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

$$
\begin{aligned}
&S \\
\Rightarrow\ & X \overset{?}{=} X \\
\Rightarrow\ & \mathbf{a}X \overset{?}{=} X \\
\Rightarrow\ & \mathbf{aa}X \overset{?}{=} X \\
\Rightarrow\ & \mathbf{aa} \overset{?}{=} X \\
\Rightarrow\ & \mathbf{aa} \overset{?}{=} \mathbf{a}X \\
\Rightarrow\ & \mathbf{aa} \overset{?}{=} \mathbf{a}
\end{aligned}
$$

# Finding a Build Order

- Let $\Sigma = \{$ **a**, $\stackrel{?}{=}$ $\}$ and let $L = \{$ **a**$^n$$\stackrel{?}{=}$**a**$^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 **a**'s independently of one another, then we can't enforce that they have the same length.

  - ***Idea:*** Build both strings of **a**'s at the same time.

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

$$\mathbf{S} \rightarrow \stackrel{?}{=} \mid \mathbf{aSa}$$

| |
|---|
| **S** |
| $\Rightarrow$ **aSa** |
| $\Rightarrow$ **aaSaa** |
| $\Rightarrow$ **aaaSaaa** |
| $\Rightarrow$ **aaa**$\stackrel{?}{=}$**aaa** |

# Storing Information in Nonterminals

- ***Key idea:*** Different non-terminals should represent different states or different types of strings.

  - For example, different phases of the build, or different possible structures for the string.

  - Think like the same ideas from DFA/NFA design where states in your automata represent pieces of information.

# Storing Information in Nonterminals

- Let $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ and let $L = \{w \in \Sigma^* \mid |w| \equiv_3 0$ and all the characters in the first third of $w$ are the same $\}$.

- Examples:

| | |
|---|---|
| $\boldsymbol{\varepsilon} \in L$ | $\mathbf{a} \notin L$ |
| $\mathbf{abb} \in L$ | $\mathbf{b} \notin L$ |
| $\mathbf{bab} \in L$ | $\mathbf{ababab} \notin L$ |
| $\mathbf{aababa} \in L$ | $\mathbf{aabaaaaaa} \notin L$ |
| $\mathbf{bbbbbb} \in L$ | $\mathbf{bbbb} \notin L$ |

# Storing Information in Nonterminals

- Let $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ and let $L = \{ w \in \Sigma^* \mid |w| \equiv_3 0$ and all the characters in the first third of $w$ are the same $\}$.

- Examples:

$\boldsymbol{\varepsilon} \in L$          $\mathbf{a} \notin L$

$\mathbf{a} \vdots \mathbf{bb} \in L$      $\mathbf{b} \notin L$

$\mathbf{b} \vdots \mathbf{ab} \in L$      $\mathbf{ab} \vdots \mathbf{abab} \notin L$

$\mathbf{aa} \vdots \mathbf{baba} \in L$      $\mathbf{aab} \vdots \mathbf{aaaaaa} \notin L$

$\mathbf{bb} \vdots \mathbf{bbbb} \in L$      $\mathbf{bbbb} \notin L$

# Storing Information in Nonterminals

- Let $\Sigma = \{$**a**, **b**$\}$ and let $L = \{w \in \Sigma^* \mid |w| \equiv_3 0$ and all the characters in the first third of $w$ are the same $\}$.

- One approach:

  **aaa**      **bab**

  **abb**      **bbb**

  **aaabab**    **bbabbb**

  **aababa**    **bbbaaaaaa**

  **aaaaaaaaa**  **bbbbbabaa**

***Observation 1:***

Strings in this language are either: the first third is **a**s or the first third is **b**s.

# Storing Information in Nonterminals

- Let $\Sigma = \{$**a**, **b**$\}$ and let $L = \{w \in \Sigma^* \mid |w| \equiv_3 0$ and all the characters in the first third of $w$ are the same $\}$.

- One approach:

| | |
|---|---|
| **aaa** | bab |
| **abb** | bbb |
| **aaabab** | bbabbb |
| **aababa** | bbbaaaaa |
| **aaaaaaaaa** | bbbbbabaa |

# Storing Information in Nonterminals

- Let $\Sigma = \{$**a**, **b**$\}$ and let $L = \{w \in \Sigma^* \mid |w| \equiv_3 0$ and all the characters in the first third of $w$ are the same $\}$.

- One approach:

| | | |
|---|---|---|
| **aaa** | bab | |
| **abb** | bbb | |
| **aaabab** | bbabbb | |
| **aababa** | bbbaaaaaa | |
| **aaaaaaaa** | bbbbbabaa | |

***Observation 2:***

Amongst these strings, for every **a** I have in the first third, I need two other characters in the last two thirds.

# Storing Information in Nonterminals

- Let $\Sigma = \{$ **a**, **b** $\}$ and let $L = \{w \in \Sigma^* \mid |w| \equiv_3 0$ and all the characters in the first third of $w$ are the same $\}$.

- One approach:

| | |
|---|---|
| **aaa** | bab |
| **abb** | bbb |
| **aaabab** | bbabbb |
| aaaaaa | |
| babaa | |

**Observation 2:**

Amongst these strings, for every **a** I have in the first third, I need two other characters in the last two thirds.

This pattern of "for every x I see here, I need a y somewhere else in the string" is very common in CFGs!

# Storing Information in Nonterminals

- Let $\Sigma = \{$**a**, **b**$\}$ and let $L = \{w \in \Sigma^* \mid |w| \equiv_3 0$ and all the characters in the first third of $w$ are the same $\}$.

- One approach:

| | | **Observation 2:** |
|---|---|---|
| **aaa** | bab | Amongst these strings, for every **a** I have in the first third, I need two other characters in the last two thirds. |
| **abb** | bbb | |
| **aaabab** | bbabbb | |
| **aababa** | bbbaaaaa | |
| **aaaaaaaa** | bbbbbabaa | |

$$A \rightarrow aAXX \mid \varepsilon \qquad X \rightarrow a \mid b$$

# Storing Information in Nonterminals

- Let $\Sigma = \{$**a**, **b**$\}$ and let $L = \{w \in \Sigma^* \mid |w| \equiv_3 0$ and all the characters in the first third of $w$ are the same $\}$.

- One approach:

**aaa**

bab

**abb**

**aaabab**

Here the nonterminal **A** represents "a string where the first third is **a**'s" and the nonterminal **X** represents "any character"

**aababa**

**aaaaaaaaa**

bbbbbabaa

**A → aAXX | ε**        **X → a | b**

# Storing Information in Nonterminals

- Let $\Sigma$ = {**a**, **b**} and let $L$ = {$w \in \Sigma^*$ | $|w| \equiv_3 0$ and all the characters in the first third of $w$ are the same }.

- One approach:

| | |
|---|---|
| **aaa** | bab |
| **abb** | bbb |
| **aaabab** | bbabbb |
| **aababa** | bbbaaaaa |
| **aaaaaaaa** | bbbbbabaa |

$\color{red}{A} \rightarrow \textbf{a}\color{red}{AXX}$ | $\textbf{ε}$     $\color{red}{X} \rightarrow \textbf{a}$ | $\textbf{b}$

# Storing Information in Nonterminals

- Let $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ and let $L = \{w \in \Sigma^* \mid |w| \equiv_3 0$ and all the characters in the first third of $w$ are the same $\}$.

- One approach:

| | |
|---|---|
| aaa | **bab** |
| abb | **bbb** |
| aaabab | **bbabbb** |
| aababa | **bbbaaaaaa** |
| aaaaaaaaa | **bbbbbabaa** |

$$\mathbf{B} \rightarrow \mathbf{bBXX} \mid \varepsilon \qquad \mathbf{X} \rightarrow \mathbf{a} \mid \mathbf{b}$$

# Storing Information in Nonterminals

- Let $\Sigma = \{$**a**, **b**$\}$ and let $L = \{ w \in \Sigma^* \mid |w| \equiv_3 0$ and all the characters in the first third of $w$ are the same $\}$.

- Tying everything together:

    **S** → **A** | **B**

    **A** → **a**AXX | ε

    **B** → **b**BXX | ε

    **X** → **a** | **b**

# Storing Information in Nonterminals

- Let $\Sigma = \{\mathbf{a}, \mathbf{b}\}$ and let $L = \{w \in \Sigma^* \mid |w| \equiv_3 0$ and all the characters in the first third of $w$ are the same $\}$.

- Tying everything together:

**S** → **A** | **B**

**A** → **aAXX** | ε

**B** → **bBXX** | ε

**X** → **a** | **b**

Overall strings in this language either follow the pattern of **A** or **B**.

# Storing Information in Nonterminals

- Let $\Sigma = \{$**a**, **b**$\}$ and let $L = \{w \in \Sigma^* \mid |w| \equiv_3 0$ and all the characters in the first third of $w$ are the same $\}$.

- Tying everything together:

$$S \to A \mid B$$

$$A \to \mathbf{aAXX} \mid \varepsilon$$

$$B \to bBXX \mid \varepsilon$$

$$X \to a \mid b$$

A represents "strings where the first third is **a**'s"

# Storing Information in Nonterminals

- Let $\Sigma = \{$**a**, **b**$\}$ and let $L = \{w \in \Sigma^* \mid |w| \equiv_3 0$ and all the characters in the first third of $w$ are the same $\}$.

- Tying everything together:

  **S** → **A** | **B**

  **A** → **aAXX** | **ε**

  **B** → **bBXX** | **ε**

  **X** → **a** | **b**

> **B** represents "strings where the first third is **b**'s"

# Storing Information in Nonterminals

- Let $\Sigma = \{$**a**, **b**$\}$ and let $L = \{w \in \Sigma^* \mid |w| \equiv_3 0$ and all the characters in the first third of $w$ are the same $\}$.

- Tying everything together:

  S → A | B

  A → aAXX | ε

  B → bBXX | ε

  **X** → **a** | **b**

  > **X** represents "either an **a** or a **b**"

# 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`

# 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

```
>>> (26 + 42) * 2 + 1
```

How does my computer know what this sequence of characters means? How can it determine whether or not this expression is even syntactically valid?

# Applications of CFGs

E → E Op E | int | (E)
Op → + | × | – | /

E

⇒ E Op E

⇒ E Op (E)

⇒ E Op (E Op E)

⇒ E × (E Op E)

⇒ int × (E Op E)

⇒ int × (int Op E)

⇒ int × (int Op int)

⇒ int × (int + int)

Given a set of production rules and an expression,

If I can somehow reverse engineer the derivation, I can ascribe meaning to the pieces of my string.

Exact details of how to do this are beyond the scope of this class – *Take CS143!*

# CFGs for Programming Languages

BLOCK → STMT
| { STMTS }

STMTS → ε
| STMT STMTS

STMT → EXPR;
| if (EXPR) BLOCK
| while (EXPR) BLOCK
| do BLOCK while (EXPR);
| BLOCK
| ...

EXPR → identifier
| constant
| EXPR + EXPR
| EXPR – EXPR
| EXPR * EXPR
| ...

# 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 **CoreNLP project** 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 would you program it?