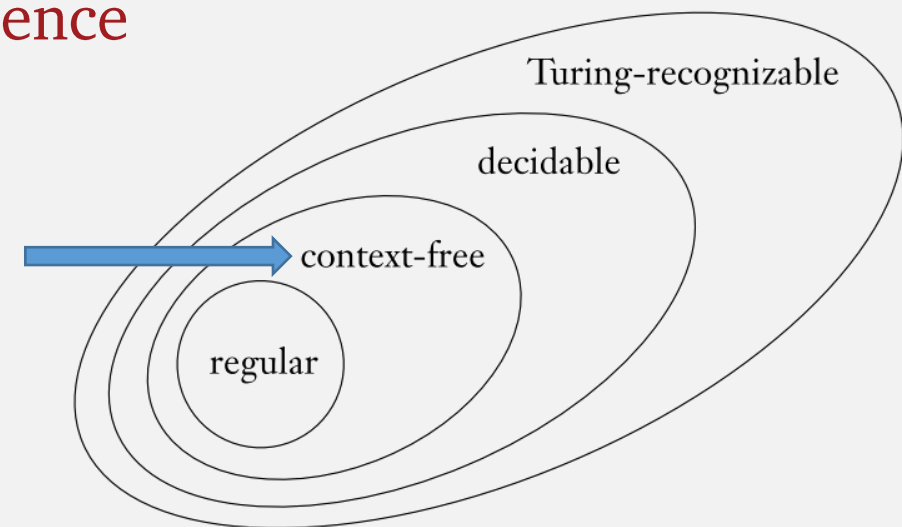# CS 420 / CS 620

# Context-Free Grammars (CFGs) and Context-Free Languages (CFLs)
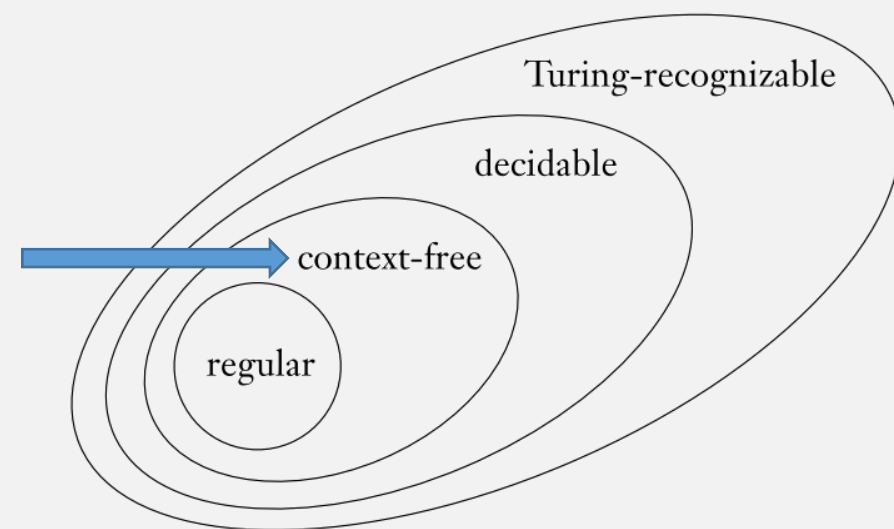
Monday October 20, 2025

UMass Boston Computer Science

# Announcements

- HW 6
  - ~~Due: Mon 10/20 12pm (noon)~~

- HW 7
  - Out: Mon 10/20 12pm (noon)
  - Due: Mon 10/27 12pm (noon)

# Non-Regular Languages

Example:

An arbitrary count

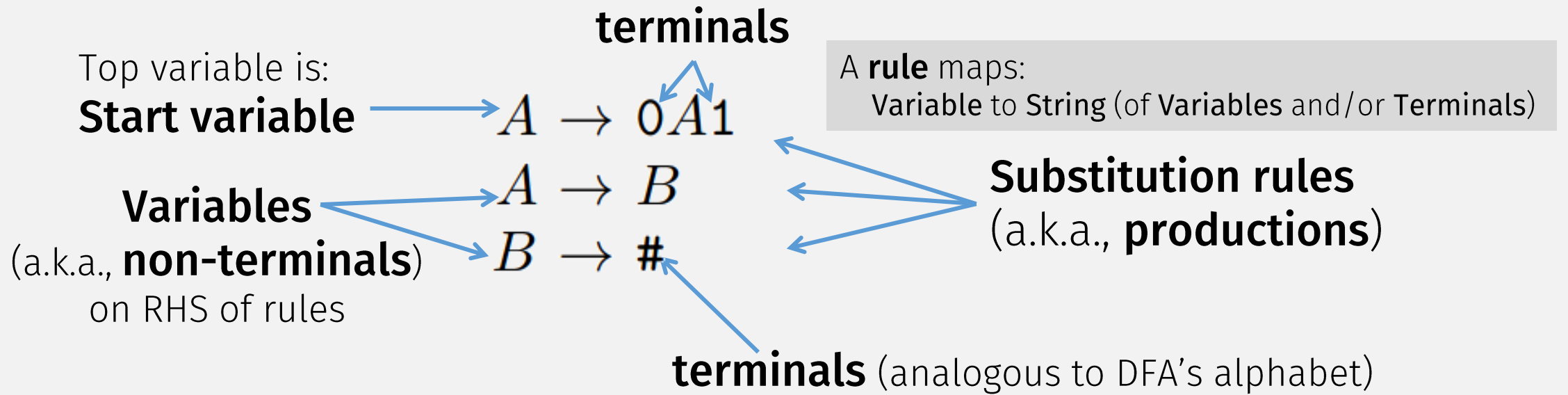$L = \{\ \mathtt{0}^n\mathtt{1}^n \mid n \geq 0\ \}$

- **A DFA recognizing** $L$ would **require infinite states!** (impossible)
  - States representing: **zero 0s seen, one 0 seen, two 0s, …**

- **This language is the same as many PLs,** e.g., HTML!
  - To better see this replace:
    - "**0**" with "`<tag>`" or "**(**"
    - "**1**" with "`</tag>`" or "**)**"

- The Problem: **remembering** <u>nestedness</u>
  - **Need to count arbitrary nesting depths**
    - E.g., **{ { { … } } }**
  - Thus: **most programming language syntax** is <u>not regular</u>!

We can
_prove_ **non-regularness** …
with the **Pumping Lemma**
(and **proof by contradiction**)

But … **what kind of
language is it then?**

# A Context-Free Grammar (CFG)

**terminals**

Top variable is:
**Start variable**

A **rule** maps:
Variable to String (of Variables and/or Terminals)

$A \rightarrow 0A1$

$A \rightarrow B$

**Variables**
(a.k.a., **non-terminals**)
on RHS of rules

$B \rightarrow \#$

**Substitution rules**
(a.k.a., **productions**)

**terminals** (analogous to DFA's alphabet)

# Context-Free Grammar (CFG): Formally

Grammar $G_1 = (V, \Sigma, R, S)$

$R$ is this <u>set</u> of rules: Var->String (of vars+terminals) mappings:

**CFG** <u>Practical Application</u>:
Used to describe
<u>programming language
syntax</u>!

Top variable is:
**Start variable**

terminals

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

**Variables**
(a.k.a., **non-terminals**)

$$B \rightarrow \#$$

**Substitution rules**
(a.k.a., **productions**)

**terminals** (analogous to DFA's alphabet)

A **context-free grammar** is a 4-tuple $(V, \Sigma, R, S)$, where

1. $V$ is a finite set called the **variables**,
2. $\Sigma$ is a finite set, disjoint from $V$, called the **terminals**,
3. $R$ is a finite set of **rules**, with each rule being a variable and a string of variables and terminals, and
4. $S \in V$ is the start variable.

$V = \{A, B\},$

$\Sigma = \{0, 1, \#\},$

$S = A,$

# Java Syntax: Described with CFGs

**ORACLE**

Chapter 2

## Chapter 2. Grammars

This chapter describes the context-free grammars used in this specification

### 2.1. Context-Free Grammars

"productions" = rules

"nonterminal" = variable

A *context-free grammar* consists of a number of *productions*. Each production has an abstract symbol called a *nonterminal* as its *left hand side*, and a sequence of one or more nonterminal and *terminal* symbols as its *right-hand side*. For each grammar, the terminal symbols are drawn from a specified *alphabet*.

"goal symbol" = Start variable

A **CFG** specifies a language!

Starting from a sentence consisting of a single distinguished nonterminal, called the *goal symbol*, a given context-free grammar specifies a language, namely, the set of possible sequences of terminal symbols that can result from repeatedly replacing any nonterminal in the sequence with a right-hand side of a production for which the nonterminal is the left-hand side.

(definition of a **language:** set of sequences of symbols)

### 2.2. The Lexical Grammar

A *lexical grammar* for the Java programming language is given in §3. This grammar has as its terminal symbols the characters of the Unicode character set. It defines a set of productions, starting from the goal symbol *Input* (§3.5), that describe how sequences of Unicode characters (§3.1) are translated into a sequence of input elements (§3.5).

Definition:
A **CFG** describes a **context-free language!**

(Could you write it in our typical IF-THEN form?)

`https://docs.oracle.com/javase/specs/jls/se7/html/jls-2.html`

# Analogies

| Regular Language | Context-Free Language (CFL) |
|:---:|:---:|
| Regular Expression | Context-Free Grammar (CFG) |
| A **Reg Expr** <u>describes</u> a **Regular lang** | A **CFG** <u>describes</u> a **CFL** |
| | |
| | |
| | |
| | |
| | |
| | |

thm

def

# (partially)
# Python Syntax: Described with a CFG

## 10. Full Grammar specification

This is the full Python grammar, as it is read by the parser generator and used to parse Python source files:

(indentation checking not expressible with CFG?)

```
# Grammar for Python

# NOTE WELL: You should also follow all the steps listed at
# https://devguide.python.org/grammar/

# Start symbols for the grammar:
#       single_input is a single interactive statement;
#       file_input is a module or sequence of commands read from an input file;
#       eval_input is the input for the eval() functions.
#       func_type_input is a PEP 484 Python 2 function type comment
# NB: compound_stmt in single_input is followed by extra NEWLINE!
# NB: due to the way TYPE_COMMENT is tokenized it will always be followed by a NEWLINE
single_input: NEWLINE | simple_stmt | compound_stmt NEWLINE
file_input: (NEWLINE | stmt)* ENDMARKER
eval_input: testlist NEWLINE* ENDMARKER
```

https://docs.python.org/3/reference/grammar.html

**Many Other Language** (partially)

# ~~Python~~ Syntax: Described with a CFG

## 10. Full Grammar specification

This is the full Python grammar, as it is read by the parser generator and used to parse Python source files:

```
# Grammar for Python

# NOTE WELL: You should also follow all the steps listed at
# https://devguide.python.org/grammar/

# Start symbols for the grammar:
#       single_input is a single interactive statement;
#       file_input is a module or sequence of commands read from an input file;
#       eval_input is the input for the eval() functions.
#       func_type_input is a PEP 484 Python 2 function type comment
# NB: compound_stmt in single_input is followed by extra NEWLINE!
# NB: due to the way TYPE_COMMENT is tokenized it will always be followed by a NEWLINE
single_input: NEWLINE | simple_stmt | compound_stmt NEWLINE
file_input: (NEWLINE | stmt)* ENDMARKER
eval_input: testlist NEWLINE* ENDMARKER
```

https://docs.python.org/3/reference/grammar.html

# Java Syntax: Described with CFGs

**Definition:**
A **CFG** describes a **context-free language**!
but **what** <u>strings</u> are **in the language?**

ORACLE

Java SE > Java SE Specifications > Java Language Specification

Prev

## Chapter 2. Grammars

This chapter describes the context-free grammars used in this specification to define the lexical and syntactic structure of a program

### 2.1. Context-Free Grammars

A *context-free grammar* consists of a number of *productions*. Each production has an abstract symbol called a *nonterminal* as its *left hand side*, and a sequence of one or more nonterminal and *terminal* symbols as its *right-hand side*. For each grammar, the terminal symbols are drawn from a specified *alphabet*.

Starting from a sentence consisting of a single distinguished nonterminal, called the *goal symbol*, a given context-free grammar specifies a language, namely, the set of possible sequences of terminal symbols that can result from repeatedly replacing any nonterminal in the sequence with a right-hand side of a production for which the nonterminal is the left-hand side.

### 2.2. The Lexical Grammar

A *lexical grammar* for the Java programming language is given in §3. This grammar has as its terminal symbols the characters of the Unicode character set. It defines a set of productions, starting from the goal symbol *Input* (§3.5), that describe how sequences of Unicode characters (§3.1) are translated into a sequence of input elements (§3.5)

https://docs.oracle.com/javase/specs/jls/se7/html/jls-2.html

# Generating Strings with a CFG

*In-class exercise:*

**Write 3 more strings that can be generated by this grammar**

Definition:
A **CFG** describes a **context-free language!**
but what <u>strings</u> are in the language?

1st rule ⟶ $A \rightarrow 0A1$

2nd rule ⟶ $A \rightarrow B$

Last rule ⟶ $B \rightarrow \#$

"Applying a rule" = replace LHS variable with RHS sequence

At each step, *arbitrarily* choose <u>any</u> variable to replace, and <u>any</u> rule to apply

Stop when: string is all terminals

A CFG **generates** a string, by repeatedly <u>applying</u> substitution rules:

Example:

$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000B111 \Rightarrow 000\#111$

Start with: Start variable

Apply 1st rule

1st rule again

1st rule again

Apply 2nd rule

Apply last rule

# Generating Strings with a CFG

$$G_1 =$$

$$A \rightarrow 0A1$$
$$A \rightarrow B$$
$$B \rightarrow \texttt{\#}$$

Strings in CFG's language
= all possible **generated** / **derived** strings

$$L(G_1) \text{ is } \{0^n \texttt{\#} 1^n \mid n \geq 0\}$$

A CFG **generates** a string, by repeatedly <u>applying</u> substitution rules:

Example:
$$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000B111 \Rightarrow 000\texttt{\#}111$$

This sequence of steps is called a **derivation**

# Derivations: Formally

Let $G = (V, \Sigma, R, S)$

**Single-step**

$$\alpha A \beta \underset{G}{\Rightarrow} \alpha \gamma \beta$$

Where:

sequence of:
variables or terminals

$\alpha, \beta \in (V \cup \Sigma)^*$

$A \in V$ ← Variable

$A \rightarrow \gamma \in R$ ← Rule

# Derivations: Formally

**Let** $G = (V, \Sigma, R, S)$

**Single-step**

$$\alpha A \beta \underset{G}{\Rightarrow} \alpha \gamma \beta$$

Where:

$$\alpha, \beta \in (V \cup \Sigma)^*$$

$$A \in V \longleftarrow \boxed{\text{Variable}}$$

$$A \rightarrow \gamma \in R \longleftarrow \boxed{\text{Rule}}$$

**Multi-step** (recursively defined, on # steps)

<u>Base case</u>:   $\alpha \overset{*}{\underset{G}{\Rightarrow}} \alpha$     (0 steps)

<u>Recursive case</u>:        (1 or more steps)

$\boxed{\text{Single step}}$                                                   (smaller)
                                                                                 Recursive "call"

- If: $\alpha \underset{G}{\Rightarrow} \beta$   and   $\beta \overset{*}{\underset{G}{\Rightarrow}} \gamma$

- Then: $\alpha \overset{*}{\underset{G}{\Rightarrow}} \gamma$

# Formal Definition of a CFL

A **context-free grammar** is a 4-tuple $(V, \Sigma, R, S)$, where

1. $V$ is a finite set called the **variables**,
2. $\Sigma$ is a finite set, disjoint from $V$, called the **terminals**,
3. $R$ is a finite set of **rules**, with each rule being a variable and a string of variables and terminals, and
4. $S \in V$ is the start variable.

$$G = (V, \Sigma, R, S)$$

"all possible sequences of terminal symbols ..."

... "that can be **generated** with rules of grammar $G$"

"the language of a grammar $G$ is ..."

$$L(G) = \left\{ w \in \Sigma^* \mid S \underset{G}{\overset{*}{\Rightarrow}} w \right\}$$

Any language that can be generated by some context-free grammar is called a **context-free language**

Alternatively (an easier form to use in a proof is),
IF a **language** can be **generated by some CFG**,
THEN that **language** is a **CFL**

Or IF a **CFG** describes a language, THEN that **language** is a **CFL**

*Flashback:* $\{0^n 1^n \mid n \geq 0\}$

- Pumping Lemma says: not a regular language
- It's a **context-free language!**
  - Proof?
  - Key step: Come up with CFG describing it ...
  - <u>Hint</u>: It's similar to:

$$A \rightarrow 0A1$$
$$A \rightarrow B$$
$$B \rightarrow \cancel{\#}\ \varepsilon$$

$$L(G_1) \text{ is } \{0^n \cancel{\#} 1^n \mid n \geq 0\}$$

Statements and Justifications?

*Proof:* $L = \{0^n 1^n \mid n \geq 0\}$ is a CFL

**Statements**

1. If a CFG describes a language, then it is a CFL

2. CFG $G_1$ describes $L$
$$A \to 0A1$$
$$A \to B$$
$$B \to \varepsilon$$

3. $L = \{0^n 1^n \mid n \geq 0\}$ is a CFL

**Justifications**

1. Definition of CFL

2. (Did you come up with examples???)

3. By Statements #1 and #2

Must be the same "$P$" to use modus ponens

# A String Can Have Multiple Derivations

"|" symbol = Shorthand for multiple rules with same LHS variable

$$\langle EXPR \rangle \rightarrow \langle EXPR \rangle + \langle TERM \rangle \mid \langle TERM \rangle$$
$$\langle TERM \rangle \rightarrow \langle TERM \rangle \times \langle FACTOR \rangle \mid \langle FACTOR \rangle$$
$$\langle FACTOR \rangle \rightarrow ( \langle EXPR \rangle ) \mid a$$

Want to generate this string: **a + a × a**

- EXPR ⇒
- EXPR + TERM ⇒
- EXPR + TERM × FACTOR ⇒
- EXPR + TERM × a ⇒
- …

**RIGHTMOST** DERIVATION

- EXPR ⇒
- EXPR + TERM ⇒
- TERM + TERM ⇒
- FACTOR + TERM ⇒
- **a** + TERM

  …

**LEFTMOST** DERIVATION

# Derivations and Parse Trees

$$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000B111 \Rightarrow 000\#111$$

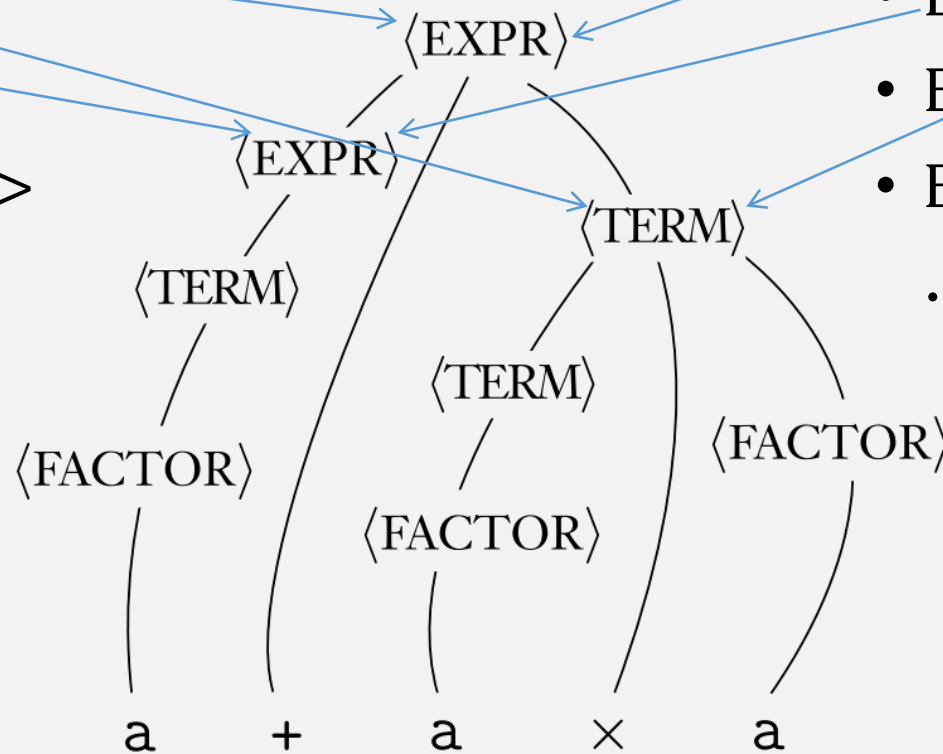A derivation may also be represented as a **parse tree**

# Multiple Derivations, Single Parse Tree

- EXPR =>
- EXPR + TERM =>
- TERM + TERM =>
- FACTOR + TERM =>
- a + TERM

  ...

- EXPR =>
- EXPR + TERM =>
- EXPR + TERM x FACTOR =>
- EXPR + TERM x a=>

  ...

⟨EXPR⟩
⟨EXPR⟩
⟨TERM⟩
⟨TERM⟩
⟨FACTOR⟩
⟨TERM⟩
⟨FACTOR⟩
⟨FACTOR⟩

a + a × a

Same parse tree

A **parse tree** represents a CFG computation ... like a **sequence of states** represents a DFA computation

A **Parse Tree** gives "meaning" to a **string**

# Ambiguity

grammar $G_5$:

$\langle EXPR \rangle \rightarrow \langle EXPR \rangle + \langle EXPR \rangle \mid \langle EXPR \rangle \times \langle EXPR \rangle \mid (\langle EXPR \rangle) \mid a$

Same **string**,
different **derivation**,
and different **parse tree!**

So this string has
two meanings!

# Ambiguity

A string $w$ is derived **ambiguously** in context-free grammar $G$ if it has two or more different leftmost derivations. Grammar $G$ is **ambiguous** if it generates some string ambiguously.

An **ambiguous grammar** can give a
string <u>multiple meanings</u>, ie **represent
<u>two different computations</u>!**
(why is this <u>bad</u>?)

# Real-life Ambiguity ("Dangling" `else`)

- What is the result of this C program?

```
if (1) if (0) printf("a"); else printf("2");
```

This **string has 2 parsings,** and thus 2 meanings!

```
if (1)
  if (0)
    printf("a");
  else
    printf("2");
```

**VS**

```
if (1)
  if (0)
    printf("a");
else
  printf("2");
```

**Ambiguous** grammars are confusing. A computation (represented by a string) should ideally have only one possible result.

Thus in practice, we typically focus on the **unambiguous** subset of CFGs (CFLs) (more on this later)

Problem is, there's no easy way to create an **unambiguous** grammar (it's up to language designers to "be careful")

# Subclasses of CFLs



DCFLs

Programming language parsers / compilers are ideally in here

Unambiguous Grammars

Ambiguous Grammars

LL(k)  LR(k)

LL(1)  LR(1)

LALR(1)

SLR

LL(0)  LR(0)

All CFLS

2) choose "look ahead" amount

2 parser design decisions:
1) Parse from left (L),
or from right (R)

# Designing Grammars : Basics

1. Think about what you want to "link" together

- E.g., $0^n 1^n$
  - $A \rightarrow 0A1$
  - # 0s and # 1s are "linked"

- E.g., HTML
  - ELEMENT $\rightarrow$ <TAG>CONTENT</TAG>
  - Start and end tags are "linked"

2. Start with small grammars (computation) and then combine
   - just like with DFAs, NFAs, and programming!

# In-class exercise: Creating CFG

alphabet $\Sigma$ is $\{0,1\}$

$$L = \{w \mid w \text{ starts and ends with the same symbol}\}$$

1) come up with examples:    In the language:  **010, 101, 11011**        **1, 0 ?**  ☑

   Not in the language:  **10, 01, 110**                $\varepsilon$ ?     ☒

2) Create CFG:

Needed Rules:

$S \rightarrow \mathbf{0}M\mathbf{0} \mid \mathbf{1}M\mathbf{1} \mid \mathbf{0} \mid \mathbf{1}$     "start/end symbol are "linked" (ie, same); middle can be anything"

$M \rightarrow MT \mid \varepsilon$         "middle: all possible terminals, repeated (ie, all possible strings)"

$T \rightarrow \mathbf{0} \mid \mathbf{1}$     "all possible terminals"

3) Check CFG: generates examples in the language; doesn't generate examples not in language

Examples Table!
(justifies that the "CFG describes $L$")

# Designing Grammars: Building Up

- **Start** with **small grammars** and then **combine** (just like **programming**)

  - To create a **grammar for the language** $\{0^n 1^n | n \geq 0\} \cup \{1^n 0^n | n \geq 0\}$

  - First create **grammar for lang** $\{0^n 1^n | n \geq 0\}$:
  $$S_1 \rightarrow 0 S_1 1 \mid \varepsilon$$

  - Then create **grammar for lang** $\{1^n 0^n | n \geq 0\}$:
  $$S_2 \rightarrow 1 S_2 0 \mid \varepsilon$$

  - Then **combine:** $S \rightarrow S_1 \mid S_2$

  $$S_1 \rightarrow 0 S_1 1 \mid \varepsilon$$
  $$S_2 \rightarrow 1 S_2 0 \mid \varepsilon$$

> New start variable and rule combines two smaller grammars

> "|" = "or" = union (combines 2 rules with same left side)

# (Closed) Operations for CFLs?

- **Start** with **small grammars** and then **combine** (just like **programming**)

- "Or": $\qquad\qquad S \rightarrow S_1 \mid S_2$

- "Concatenate": $S \rightarrow S_1 S_2$

- "Star" (repetition): $S' \rightarrow S' S_1 \mid \varepsilon$

Could you write out the full proof?

# Regular Language vs CFL Comparison

| Regular Languages | Context-Free Languages (CFLs) |
|---|---|
| Regular Expression | Context-Free Grammar (CFG) |
| <u>describes</u> a **Regular Lang** | <u>describes</u> a **CFL** |
| | |
| | |
| | |
| | |
| | |
| | |

# Regular Language vs CFL Comparison

| Regular Languages | Context-Free Languages (CFLs) |
| --- | --- |
| Regular Expression | Context-Free Grammar (CFG) |
| <u>describes</u> a **Regular Lang** | <u>describes</u> a **CFL** |
| | |
| Finite State Automaton (FSM) | **???** |
| <u>recognizes</u> a **Regular Lang** | <u>recognizes</u> a **CFL** |
| | |
| | |
| | |

# Regular Language vs CFL Comparison

| Regular Languages | Context-Free Languages (CFLs) |
|---|---|
| Regular Expression describes a **Regular Lang** | Context-Free Grammar (CFG) describes a **CFL** |
| | |
| Finite State Automaton (FSM) recognizes a **Regular Lang** | **Push-down Automata** (PDA) recognizes a **CFL** |
| | |
| | |

thm

def

def

thm

# Regular Language vs CFL Comparison

| Regular Languages | Context-Free Languages (CFLs) |
|---|---|
| Regular Expression | Context-Free Grammar (CFG) |
| describes a **Regular Lang** | describes a **CFL** |
| | |
| Finite State Automaton (FSM) | **Push-down Automata** (PDA) |
| recognizes a **Regular Lang** | recognizes a **CFL** |
| | |
| Proved: | To prove: |
| **Regular Lang ⇔ Regular Expr** | CFL ⇔ PDA |

thm

def

def

thm