# Context-Free Languages (CFLs)

Monday, March 6, 2023

Turing-recognizable

decidable

context-free

regular

# Announcements

- HW 4 in
  - ~~due Sun 3/5 11:59pm EST~~

- HW 5 out
  - due Sun 3/19 11:59pm EST
  - (after Spring Break!)

Quiz Preview

-  What do we call the class of languages that are generated by a CFG?

*Last Time:*

**Pumping lemma** If $A$ is a regular language, then there is a number $p$ (the pumping length) where if $s$ is any string in $A$ of length at least $p$, then $s$ may be divided into three pieces, $s = xyz$, satisfying the following conditions:

1. for each $i \geq 0$, $xy^i z \in A$,
2. $|y| > 0$, and
3. $|xy| \leq p$.

Let $B$ be the language $\{0^n 1^n \mid n \geq 0\}$. We use the pumping lemma to prove that $B$ is not regular. The proof is by contradiction.

- <u>Assume:</u> language $B$ is regular

- So it <u>must satisfy</u> the **Pumping Lemma:**

*Last Time:*

**Pumping lemma** If $A$ is a regular language, then there is a number $p$ (the pumping length) where if $s$ is any string in $A$ of length at least $p$, then $s$ may be divided into three pieces, $s = xyz$, satisfying the following conditions:

1. for each $i \geq 0$, $xy^i z \in A$,
2. $|y| > 0$, and
3. $|xy| \leq p$.

Let $B$ be the language $\{0^n 1^n \mid n \geq 0\}$. We use the pumping lemma to prove that $B$ is not regular. The proof is by contradiction.

- Assume: language $B$ is regular

- So it must satisfy the **Pumping Lemma**:
  - All strings $\geq$ length $p$ …
  - … can be split into some $xyz$ … where $y$ is "pumpable"

- Get **contradiction** by finding **counterexample**: a not "pumpable string $\geq$ length $p$: $0^p 1^p$
  - Must show string cannot be pumped for all possible splittings into $xyz$
  - Use pumping lemma condition #3 to eliminate some cases

- Therefore, $B$ is not regular
  - (This is the **contrapositive** of the Pumping Lemma)
- This is a **contradiction** of the assumption!

contradiction

$$00 \ldots 011 \ldots 1$$

$y$ must be in the first $p$ 0s!

4

*Last Time:*

**Pumping lemma** If $A$ is a regular language, then there is a number $p$ (the pumping length) where if $s$ is any string in $A$ of length at least $p$, then $s$ may be divided into three pieces, $s = xyz$, satisfying the following conditions:
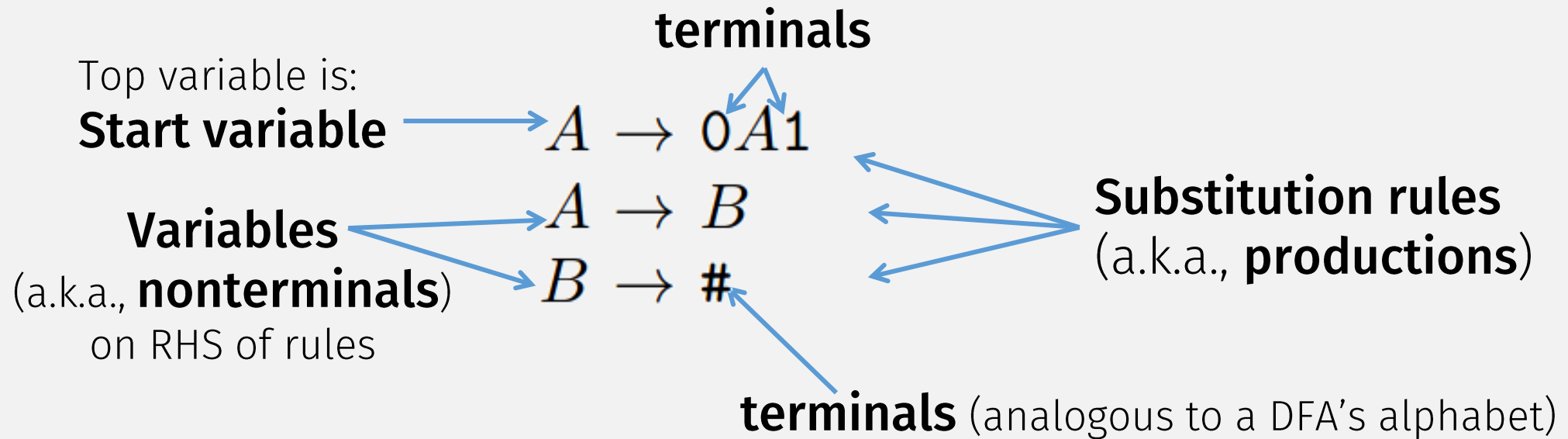
1. for each $i \geq 0$, $xy^i z \in A$,
2. $|y| > 0$, and
3. $|xy| \leq p$.

Let $B$ be the language $\{0^n 1^n \,|\, n \geq 0\}$. We use the pumping lemma to prove that $B$ is not regular. The proof is by contradiction.

If this language is not regular, then what is it???

Maybe? … a **context-free language** (CFL)?

# A Context-Free Grammar (CFG)

**terminals**

Top variable is:
**Start variable** $\longrightarrow$ $A \rightarrow 0A1$

$A \rightarrow B$

**Variables** $\longrightarrow$
(a.k.a., **nonterminals**)
on RHS of rules

$B \rightarrow \#$

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

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

# A Context-Free Grammar (CFG)

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

$R$ is this set of rules (mappings):
Top variable is:

**terminals**

**Start variable**

$$A \rightarrow 0A1$$
$$A \rightarrow B$$
$$B \rightarrow \#$$

**Variables**
(a.k.a., **nonterminals**)

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

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

**terminals** (analogous to a 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



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

A **CFG** specifies a language!

(definition of a language)

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:

```
# 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
```

(indentation checking
probably not
describable with a CFG)

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

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

**Many Other Language** (partially)

## 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



**ORACLE**

## 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

15

# Generating Strings with a CFG

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

$$G_1 =$$

1st rule $\longrightarrow$

$$A \rightarrow 0A1$$
$$A \rightarrow B$$
$$B \rightarrow \#$$

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

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

"Applying a rule"
= replace LHS variable
with RHS

At each step, can 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 applying substitution rules:

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

Start variable

After applying 1st rule

1st rule again

1st rule again

Use 2nd rule

Use last rule

# Derivations: Formally

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

**Single-step**

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

Where:

$\alpha, \beta \in (V \cup \Sigma)^*$    Strings of terminals and variables

$A \in V$    Variable

$A \to \gamma \in R$    Rule

**Extended Derivation**

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

Recursive case:     (multistep)

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

  Single step

  Recursive call

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

17

# 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.

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

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

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

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

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

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

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

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 $L$
   $$A \rightarrow 0A1$$
   $$A \rightarrow B$$
   $$B \rightarrow \varepsilon$$

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

# A String Can Have Multiple Derivations

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

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

27

# Derivations and Parse Trees

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

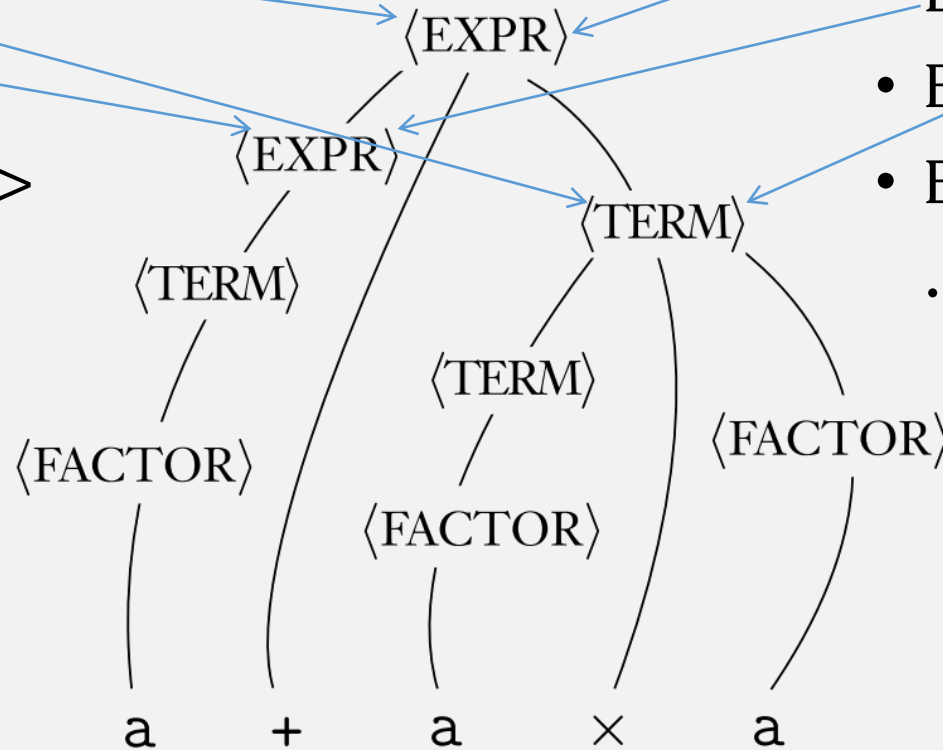A derivation may also be represented as a **parse tree**

# 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⟩

⟨FACTOR⟩

⟨TERM⟩

⟨FACTOR⟩

⟨TERM⟩

⟨FACTOR⟩

a    +    a    ×    a

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

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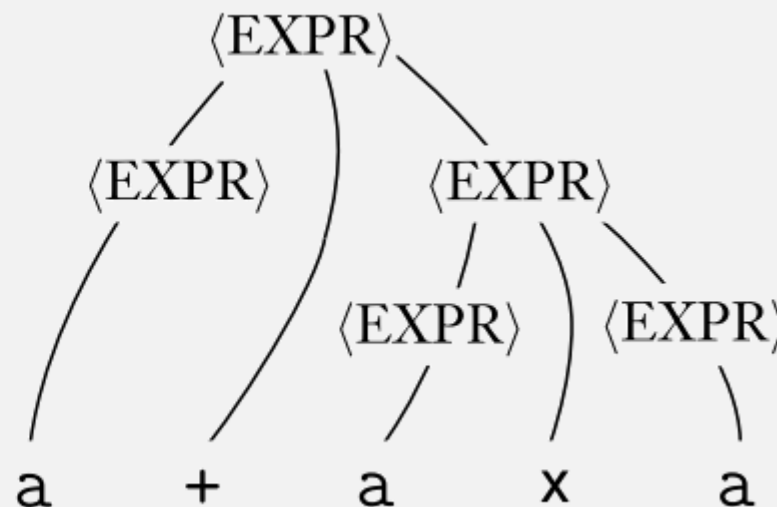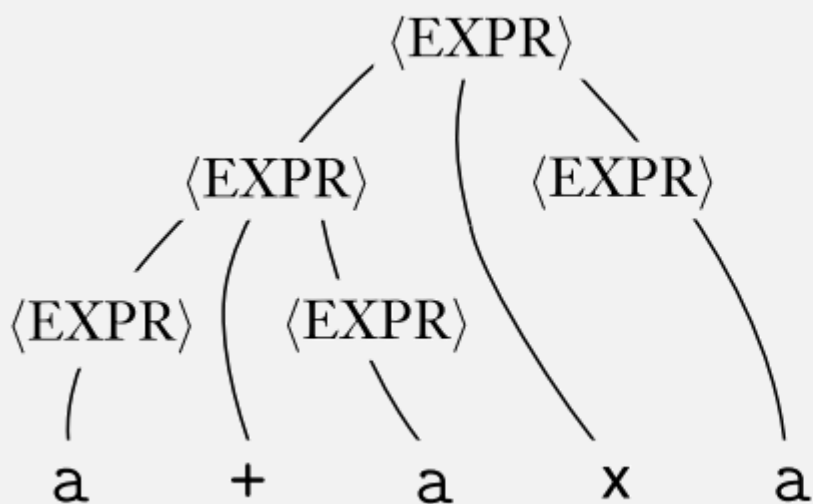a DFA computation

Same parse tree

30

# Ambiguity

grammar $G_5$:

$$\langle \text{EXPR} \rangle \rightarrow \langle \text{EXPR} \rangle + \langle \text{EXPR} \rangle \mid \langle \text{EXPR} \rangle \times \langle \text{EXPR} \rangle \mid (\langle \text{EXPR} \rangle) \mid \texttt{a}$$

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");
```

⬇

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

**VS**

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

> This string has 2 <u>parsings</u>, and thus <u>2 meanings</u>!

> **Ambiguous** grammars are confusing. A <u>computation</u> on a string should ideally have only <u>one result</u>.

> Thus <u>in practice</u>, we typically focus on the **unambiguous** <u>subset</u> 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")

# 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., XML
  - ELEMENT $\rightarrow$ <TAG>CONTENT</TAG>
  - Start and end tags are "linked"

2. Start with small grammars and then combine (just like FSMs)

# Designing Grammars: Building Up

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

  - 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 $\{0^n 1^n | n \geq 0\}$ :
    $$S_1 \rightarrow 0S_1 1 \mid \varepsilon$$

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

  - Then combine:
    $$S \rightarrow S_1 \mid S_2$$
    $$S_1 \rightarrow 0S_1 1 \mid \varepsilon$$
    $$S_2 \rightarrow 1S_2 0 \mid \varepsilon$$

New start variable and rule combines two smaller grammars

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

# (Closed) Operations on CFLs?

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

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

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

- "Repetition": $\quad S' \rightarrow S' S_1 \mid \varepsilon$

Could you write out the full proof?

# In-class Example: Designing grammars

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

$\{w|\ w$ starts and ends with the same symbol$\}$

- $S \rightarrow 0C'0\ |\ 1C'1\ |\ \varepsilon$      "string <u>starts/ends</u> with same symbol, <u>middle</u> can be anything"

- $C' \rightarrow C'C\ |\ \varepsilon$      "<u>middle</u>: all possible terminals, repeated (ie, <u>all possible </u>strings)"

- $C \rightarrow 0\ |\ 1$
   "<u>all possible </u>terminals"

*Next Time:*

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

# Next Time:

| Regular Languages | Context-Free Languages (CFLs) |
|:---:|:---:|
| Regular Expression | Context-Free Grammar (CFG) |
| A Reg Expr <u>describes</u> a Regular Lang | A CFG <u>describes</u> a CFL |
| | |
| Finite Automaton (FSM) | **Push-down Automaton** (PDA) |
| An FSM <u>recognizes</u> a Regular Lang | A **PDA** <u>recognizes</u> a CFL |
| | |
| | |
| | |

# Next Time:

| Regular Languages | Context-Free Languages (CFLs) |
|:---:|:---:|
| Regular Expression | Context-Free Grammar (CFG) |
| A Reg Expr <u>describes</u> a Regular Lang | A CFG <u>describes</u> a CFL |
| | |
| Finite Automaton (FSM) | **Push-down Automaton** (PDA) |
| An FSM <u>recognizes</u> a Regular Lang | A **PDA** <u>recognizes</u> a CFL |
| <u>DIFFERENCE:</u> | <u>DIFFERENCE:</u> |
| A Regular Lang is <u>defined</u> with a FSM | A CFL is <u>defined</u> with a CFG |
| *Proved*: Reg Expr ⇔ Reg Lang | *Must prove*: **PDA** ⇔ CFL |

thm

def

def

thm

43

# Check-in Quiz 3/6

On gradescope