

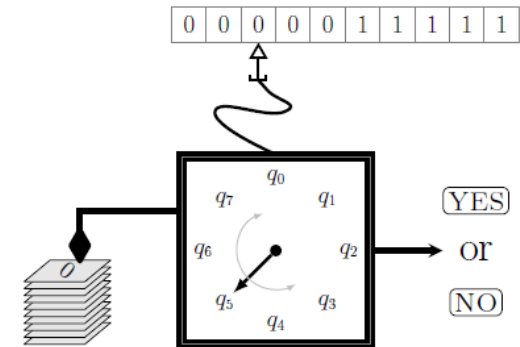
Context Free Grammars



- Malik Magdon-Ismael. Discrete Mathematics and Computing.
 - Chapter 25

Adding Memory

- DFAs have no scratch paper. It's hard to compute entirely in your head.
- **Stack Memory.** Think of a file-clerk with a stack of papers.
- The clerk's capabilities:
 - see the top sheet;
 - remove the top sheet (*pop*)
 - *push* something new onto the top of the stack.
 - no access to inner sheets without removing top.
- DFA with a stack is a *pushdown automaton (PDA)*
- How does the stack help to solve $\mathcal{L}_{0^n 1^n} = \{0^n 1^n \mid n \geq 0\}$?
 1. When you read in each 0, write it to the stack.
 2. For each 1, pop the stack. At the end if the stack is empty, accept.
- The memory allows the automaton to “remember” n .



- Solving a problem by listing out the language.
- Rules for Context Free Grammars (CFG).
- Examples of Context Free Grammars.
 - English.
 - Programming.
- Proving a CFG solves a problem.
- Parse Trees.
- Pushdown Automata and non context free languages.

Recursively Defined Language: Listing a Language

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

- How do we define this language recursively?

$\varepsilon \in \mathcal{L}_{0^n 1^n}$ [basis]

$x \in \mathcal{L}_{0^n 1^n} \rightarrow 0x1 \in \mathcal{L}_{0^n 1^n}$ [constructor rule]

Nothing else is in $\mathcal{L}_{0^n 1^n}$ [minimality]

- To test if $0010 \in \mathcal{L}_{0^n 1^n}$: generate strings in order of length and test each for a match:

$\varepsilon \rightarrow 01 \rightarrow 0011 \rightarrow 000111$

NO

- A Context Free Grammar is like a recursive definition

1. $S \rightarrow \varepsilon$

2. $S \rightarrow 0S1$

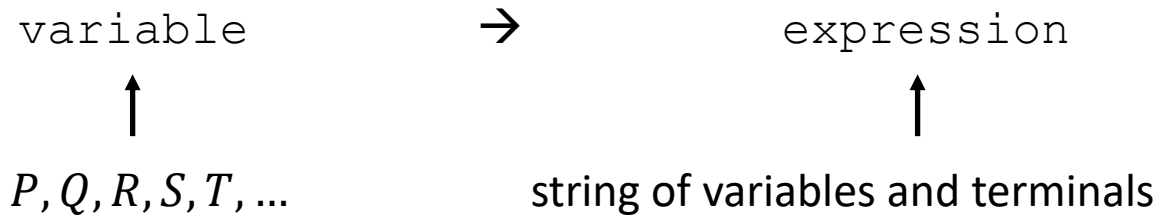
Rules for Context Free Grammars (CFGs)

- Production rules of CFGs

1. $S \rightarrow \varepsilon$

2. $S \rightarrow 0S1$

- Each production rule has the form



$$\begin{array}{ccccccccc} S & \xRightarrow{2} & 0S1 & \xRightarrow{2} & 00S11 & \xRightarrow{2} & 000S111 & \xRightarrow{2} & 0000S1111 \\ \begin{array}{c} 1 \Downarrow \\ \varepsilon \end{array} & & \begin{array}{c} 1 \Downarrow \\ 01 \end{array} & & \begin{array}{c} 1 \Downarrow \\ 0011 \end{array} & & \begin{array}{c} 1 \Downarrow \\ 000111 \end{array} & & \begin{array}{c} 1 \Downarrow \\ 00001111 \end{array} \end{array}$$

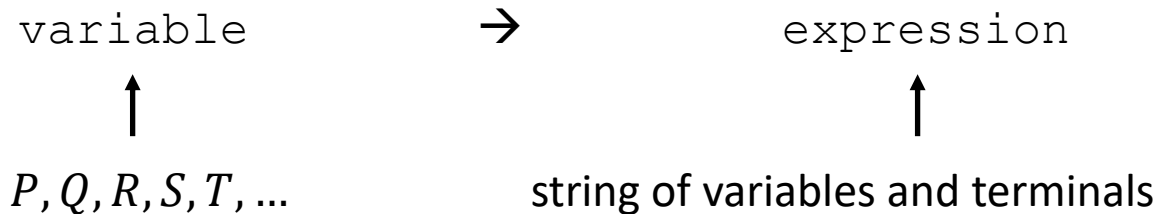


- Production rules of CFGs

1. $S \rightarrow \varepsilon$

2. $S \rightarrow 0S1$

- Each production rule has the form



1. Write down the start variable (from the first production rule, typically S).
 2. Replace *one* variable in the current string with the expression from a production rule that *starts* with that variable on the left.
 3. Repeat step 2 until no variables remain in the string.
- “Replace variable with expression, no matter where (independent of context)”
 - Shorthand: 1: $S \rightarrow \varepsilon|0S1$

- How do we write the CFG for equality?

$$\text{CFG}_{\text{bal}} \quad 1: S \rightarrow \varepsilon \mid 0S1S \mid 1S0S$$

- A derivation of 0110 in CFG_{bal} (each step is called an inference).

$$S \xrightarrow{1} 0S1S \xrightarrow{1} 0S11S0S \xrightarrow{1} 0\varepsilon11S0S \xrightarrow{*} 0110$$

- Notation

$$S \xRightarrow{*} 0110$$

$\xRightarrow{*}$ means “A derivation starting from S yields 0110”

- Distinguish S from a *mathematical* variable (e.g. x)

$$0S1S \quad \text{versus} \quad 0x1x$$

- Two S 's are replaced independently. Two x 's must be the same, e.g. $x = 11$

1. $S \rightarrow \langle \text{phrase} \rangle \langle \text{verb} \rangle$
2. $\langle \text{phrase} \rangle \rightarrow \langle \text{article} \rangle \langle \text{noun} \rangle$
3. $\langle \text{article} \rangle \rightarrow A_ | \text{The}_$
4. $\langle \text{noun} \rangle \rightarrow \text{cat}_ | \text{dog}_$
5. $\langle \text{verb} \rangle \rightarrow \text{walks.} | \text{runs.} | \text{walks.}_ S | \text{runs.}_ S$

$S \stackrel{1}{\Rightarrow} \langle \text{phrase} \rangle \langle \text{verb} \rangle$

$\stackrel{5}{\Rightarrow} \langle \text{phrase} \rangle \text{walks}$

$\stackrel{2}{\Rightarrow} \langle \text{article} \rangle \langle \text{noun} \rangle \text{walks}$

$\stackrel{3}{\Rightarrow} A_ \langle \text{noun} \rangle \text{walks}$

$\stackrel{4}{\Rightarrow} A_ \text{cat}_ \text{walks}$

A CFG for Programming



1. $S \rightarrow \langle \text{stmt} \rangle ; S \mid \langle \text{stmt} \rangle ;$
2. $\langle \text{stmt} \rangle \rightarrow \langle \text{assign} \rangle \mid \langle \text{declare} \rangle$
3. $\langle \text{declare} \rangle \rightarrow \text{int } _ \langle \text{variable} \rangle$
4. $\langle \text{assign} \rangle \rightarrow \langle \text{variable} \rangle = \langle \text{integer} \rangle$
5. $\langle \text{integer} \rangle \rightarrow \langle \text{integer} \rangle \langle \text{digit} \rangle \mid \langle \text{digit} \rangle$
6. $\langle \text{digit} \rangle \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$
7. $\langle \text{variable} \rangle \rightarrow x \mid x \langle \text{variable} \rangle$

Constructing a CFG to Solve a Problem

- Going back to our language of balanced 0s and 1s

$$\mathcal{L}_{bal} = \{\text{strings with an equal number of 1s and 0s}\}$$

$$001011010110 = 0 \cdot \mathbf{0101} \cdot 1 \cdot \mathbf{010110}$$



in \mathcal{L}_{bal}



in \mathcal{L}_{bal}

$$= 0S1S$$

- Every large string in \mathcal{L}_{bal} can be obtained (recursively) from smaller ones.

$$S \rightarrow \varepsilon \mid 0S1S \mid 1S0S$$

- We must prove that:
 - Every string generated by this CFG is in \mathcal{L}_{bal}
 - Every string in \mathcal{L}_{bal} can be derived by this grammar

Proving a CFG Solves a Problem



- Going back to our language of balanced 0s and 1s

$$\mathcal{L}_{bal} = \{\text{strings with an equal number of 1s and 0s}\}$$

- Every large string in \mathcal{L}_{bal} can be obtained (recursively) from smaller ones.

$$S \rightarrow \varepsilon \mid 0S1S \mid 1S0S$$

- Proof. [Every derivation in CFG_{bal} generates a string in \mathcal{L}_{bal}]

- Use strong induction on the length of the derivation

- i.e., number of production rules invoked

- Base Case. length-1 derivation gives ε .

- Induction step. The derivation starts in one of two ways:

$$1. S \rightarrow 0S1S$$

$$\begin{array}{c} * \Downarrow \quad * \Downarrow \\ x \quad y \end{array}$$

$$2. S \rightarrow 1S0S$$

$$\begin{array}{c} * \Downarrow \quad * \Downarrow \\ x \quad y \end{array}$$

- The derivations of x and y are shorter.

- By the induction hypothesis, $x, y \in \mathcal{L}_{bal}$, so the final strings are in \mathcal{L}_{bal}

Proving a CFG Solves a Problem

- Going back to our language of balanced 0s and 1s

$$\mathcal{L}_{bal} = \{\text{strings with an equal number of 1s and 0s}\}$$

- Every large string in \mathcal{L}_{bal} can be obtained (recursively) from smaller ones.

$$S \rightarrow \varepsilon \mid 0S1S \mid 1S0S$$

- Proof. [Every derivation in CFG_{bal} generates a string in \mathcal{L}_{bal}]

- Proof. [Every string in \mathcal{L}_{bal} can be derived within CFG_{bal}]

- Use strong induction on the length of the derivation

- i.e., number of production rules invoked

- Base Case. length-1 derivation gives ε .

- Induction step. Any string $w \in \mathcal{L}_{bal}$ has one of two forms:

1. $w = 0w_11w_2$

2. $w = 1w_10w_2$

- where w_1, w_2 have same number of 1s and 0s and have smaller length

- Why must w_1 and w_2 exist?

- String w must have a balanced substring (0 and 1 may be endpoints of w)

- By the induction hypothesis, $S \xRightarrow{*} w_1$ and $S \xRightarrow{*} w_2$, so $S \xRightarrow{*} w$



- Exercise 25.5

- Consider these two languages and their CFGs:

$$\mathcal{L}_1 = \{0^n 1^n \mid n \geq 0\}$$

$$A \rightarrow \varepsilon \mid 0A1$$

$$\mathcal{L}_2 = \{1^n 0^n \mid n \geq 0\}$$

$$B \rightarrow \varepsilon \mid 1B0$$

- What is $\mathcal{L}_1 \cup \mathcal{L}_2$?
 - All strings of equal 0s and 1s, where either all 0s come first or all 1s come first
- What is the CFG?

$$S \rightarrow A \mid B$$

$$A \rightarrow \varepsilon \mid 0A1$$

$$B \rightarrow \varepsilon \mid 1B0$$

- Consider these two languages and their CFGs:

$$\mathcal{L}_1 = \{0^n 1^n \mid n \geq 0\}$$

$$A \rightarrow \varepsilon \mid 0A1$$

$$\mathcal{L}_2 = \{1^n 0^n \mid n \geq 0\}$$

$$B \rightarrow \varepsilon \mid 1B0$$

- What is $\mathcal{L}_1 \bullet \mathcal{L}_2$?
 - All strings of equal 0s and 1s, where the first part comes from \mathcal{L}_1 and the second part comes from \mathcal{L}_2 😊
- What is the CFG?

$$S \rightarrow AB$$

$$A \rightarrow \varepsilon \mid 0A1$$

$$B \rightarrow \varepsilon \mid 1B0$$

- Consider these two languages and their CFGs:

$$\mathcal{L}_1 = \{0^n 1^n \mid n \geq 0\}$$

$$A \rightarrow \varepsilon \mid 0A1$$

$$\mathcal{L}_2 = \{1^n 0^n \mid n \geq 0\}$$

$$B \rightarrow \varepsilon \mid 1B0$$

- What is \mathcal{L}_1^* ?
 - All concatenations of strings in \mathcal{L}_1
- What is the CFG?

$$S \rightarrow \varepsilon \mid SA$$

$$A \rightarrow \varepsilon \mid 0A1$$

- **Example 25.2.** CFGs can implement DFAs, and so are strictly more powerful.

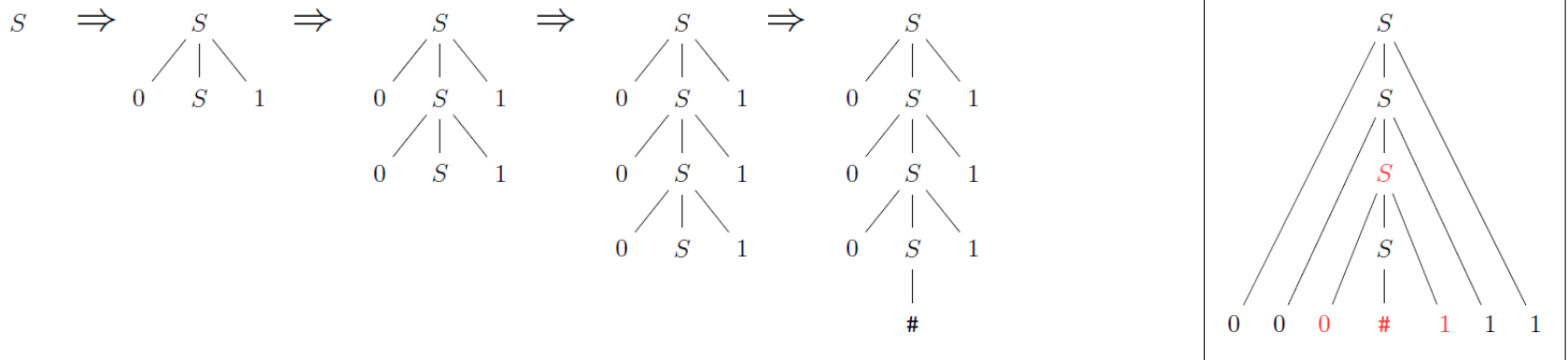
- Consider the CFG

$$S \rightarrow \#|0S1$$

- What is the derivation of 000#111?

$$S \Rightarrow 0S1 \Rightarrow 00S11 \Rightarrow 000S111 \Rightarrow 000\#111$$

- The *parse tree* gives us more information than a derivation



- Clearly shows how a substring was derived from its parent variable.

- Here is the CFG for arithmetic

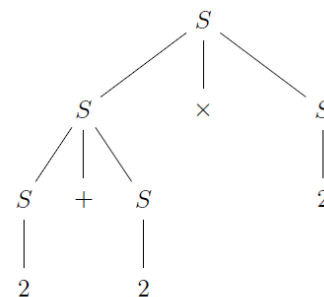
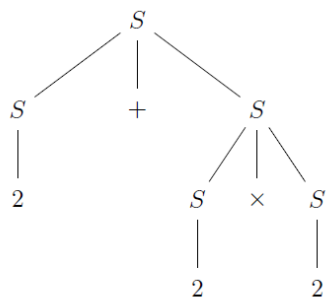
$$S \rightarrow S + S \mid S \times S \mid (S) \mid 2$$

– (the terminals are **+**, **×**, **(**, **)** and **2**)

- What are two derivations of $2 + 2 \times 2$

$$S \Rightarrow S + S \Rightarrow S + S \times S \stackrel{*}{\Rightarrow} 2 + 2 \times 2$$

$$S \Rightarrow S \times S \Rightarrow S + S \times S \stackrel{*}{\Rightarrow} 2 + 2 \times 2$$



- Parse tree \leftrightarrow How you interpret the string.
- Different parse trees \leftrightarrow different meanings.
- BAD!** We want unambiguous meaning
 - programs, html-code, math, English, . . .

- Here is the CFG for arithmetic

$$S \rightarrow S + S | S \times S | (S) | 2$$

– (the terminals are **+**, **×**, **(**, **)** and **2**)

- What are two derivations of $2 + 2 \times 2$

$$S \Rightarrow S + S \Rightarrow S + S \times S \overset{*}{\Rightarrow} 2 + 2 \times 2$$

$$S \Rightarrow S \times S \Rightarrow S + S \times S \overset{*}{\Rightarrow} 2 + 2 \times 2$$

- How do we write an unambiguous version of the arithmetic CFG?

– Use parenthesis and order of operations!

$$S \rightarrow P | S + P$$

$$P \rightarrow T | P \times T$$

$$T \rightarrow 2 | (S)$$

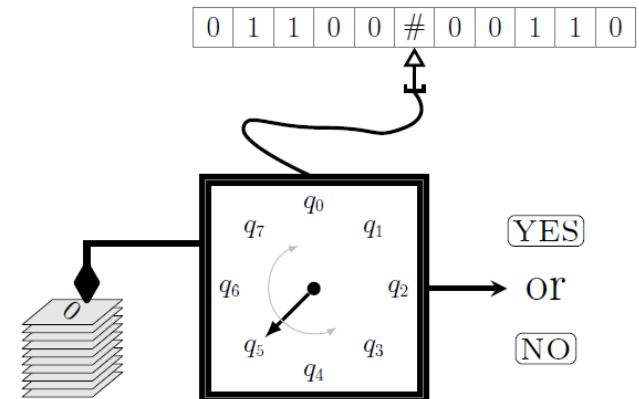
- Consider the palindrome-like language:

$$\mathcal{L} = \{w\#w^R \mid w \in \{0,1\}^*\}$$

- What is the CFG?

$$S \rightarrow \#|0S0|1S1$$

- Why can't a DFA decide whether words are in this language?
 - It needs to remember the first string w
- We can come with a DFA for deciding this language if we give it memory!
- DFA with stack memory (push, pop, read)
- Push the first half of the string (before #).
- For each bit in the second half, pop the stack and compare.
- DFAs with stack memory closely related to CFGs.



Non Context Free Expressions

- What languages do you think are hard to solve with a simple stack?

- Repetition (a stack only counts in reverse)

$$\{w#w\}$$

- Multiple-equality (similar to repetition)

$$\{0^n 1^n 0^n\}$$

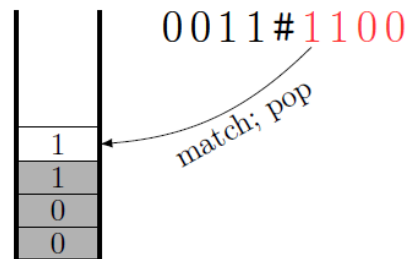
- Squaring

$$\{0^{n^2}\}, \{0^n 1^n 0^n\}$$

- Exponentiation

$$\{0^{2^n}\}, \{0^n 1^{2^n}\}$$

- Let's see why $w#w^R$ is context-free:



0011 is pushed. DFA matches 1100 by popping

Non Context Free Expressions, cont'd

- What languages do you think are hard to solve with a simple stack?

- Repetition (a stack only counts in reverse)

$$\{w#w\}$$

- Multiple-equality (similar to repetition)

$$\{0^n 1^n 0^n\}$$

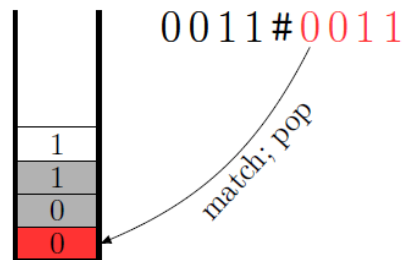
- Squaring

$$\{0^{n^2}\}, \{0^n 1^n 0^{n^2}\}$$

- Exponentiation

$$\{0^{2^n}\}, \{0^n 1^{2^n}\}$$

- Let's see why $w#w$ is not context-free:



0011 is pushed. DFA needs bottom-access to match

Non Context Free Expressions, cont'd

- What languages do you think are hard to solve with a simple stack?

- Repetition (a stack only counts in reverse)

$$\{w#w\}$$

- Multiple-equality (similar to repetition)

$$\{0^n 1^n 0^n\}$$

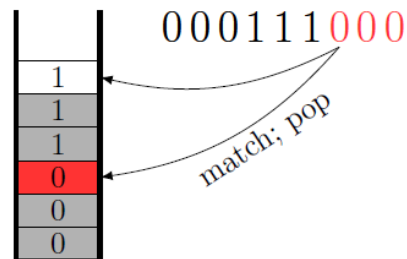
- Squaring

$$\{0^{n^2}\}, \{0^n 1^n 0^n\}$$

- Exponentiation

$$\{0^{2^n}\}, \{0^n 1^{2^n}\}$$

- Let's see why $0^n 1^n 0^n$ is not context-free:



000111 is pushed. DFA needs random access to match

Non Context Free Expressions, cont'd

- The file clerk who only has access to the top of the *stack* of papers has fundamentally less power than the file clerk who has a *filing cabinet* with access to all papers.
- **We need a new model**, one with Random Access Memory (RAM).
 - Random doesn't mean probabilistic. Means the machine can access any part of the memory