# UMB CS 622
# Chomsky Normal Form

Wednesday, April 17, 2024

Turing-recognizable

decidable

context-free

regular

Halting TMs,
a.k.a., "algorithms"

... that analyze CFLs

# Announcements

- HW 8 in
  - ~~Due Wed April 17 12pm noon~~

- HW 9 out
  - Due Wed April 24 12pm noon

# *Last Time:* Decider Turing Machines

- 2 classes of Turing Machines
  - **Recognizers** (all TMs): may loop forever
    - TM that **loops** on an input does **not accept** that input
  - **Deciders** (subset of TMs) (**algorithms**) always halt
    - Must **accept** or **reject**


- **Decider** definitions must include a **termination argument:**
  - Explains (informally) **why every step in the TM halts**
  - (Pay special attention to loops)

# *Last Time:* Algorithms About Regular Langs

- $A_{\mathsf{DFA}} = \{\langle B, w\rangle|\ B \text{ is a DFA that accepts input string } w\}$
  - **Decider**: Simulates DFA by implementing extended $\delta$ function

- $A_{\mathsf{NFA}} = \{\langle B, w\rangle|\ B \text{ is an NFA that accepts input string } w\}$
  - **Decider**: Uses **NFA→DFA** decider + $A_{\mathsf{DFA}}$ decider

- $A_{\mathsf{REX}} = \{\langle R, w\rangle|\ R \text{ is a regular expression that generates string } w\}$
  - **Decider**: Uses **RegExpr→NFA** decider + $A_{\mathsf{NFA}}$ decider

- $E_{\mathsf{DFA}} = \{\langle A\rangle|\ A \text{ is a DFA and } L(A) = \emptyset\}$
  - **Decider**: Reachability algorithm · Lang of the DFA

Remember:
**TMs ~ programs**
**Creating TM ~ programming**
**Previous theorems ~ library**

- $EQ_{\mathsf{DFA}} = \{\langle A, B\rangle|\ A \text{ and } B \text{ are DFAs and } L(A) = L(B)\}$
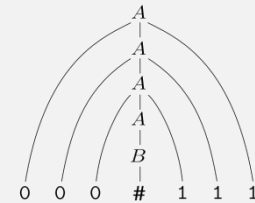  - **Decider**: Uses complement and intersection closure construction + $E_{\mathsf{DFA}}$ decider

# *Next:* Algorithms (Decider TMs) for CFLs?

- What can we predict about CFGs or PDAs?

# Thm: $A_{\mathsf{CFG}}$ is a decidable language

$$A_{\mathsf{CFG}} = \{\langle G, w \rangle \mid G \text{ is a CFG that generates string } w\}$$

- This is a very practically important problem ...
- ... equivalent to:
  - **Algorithm** to **parse** "program" $w$ for a programming language with grammar $G$?

- A Decider for this problem could ... ?
  - Try every possible derivation of $G$, and check if it's equal to $w$?
  - But this might never halt
    - E.g., what if there are rules like: $S \rightarrow 0S$ or $S \rightarrow S$
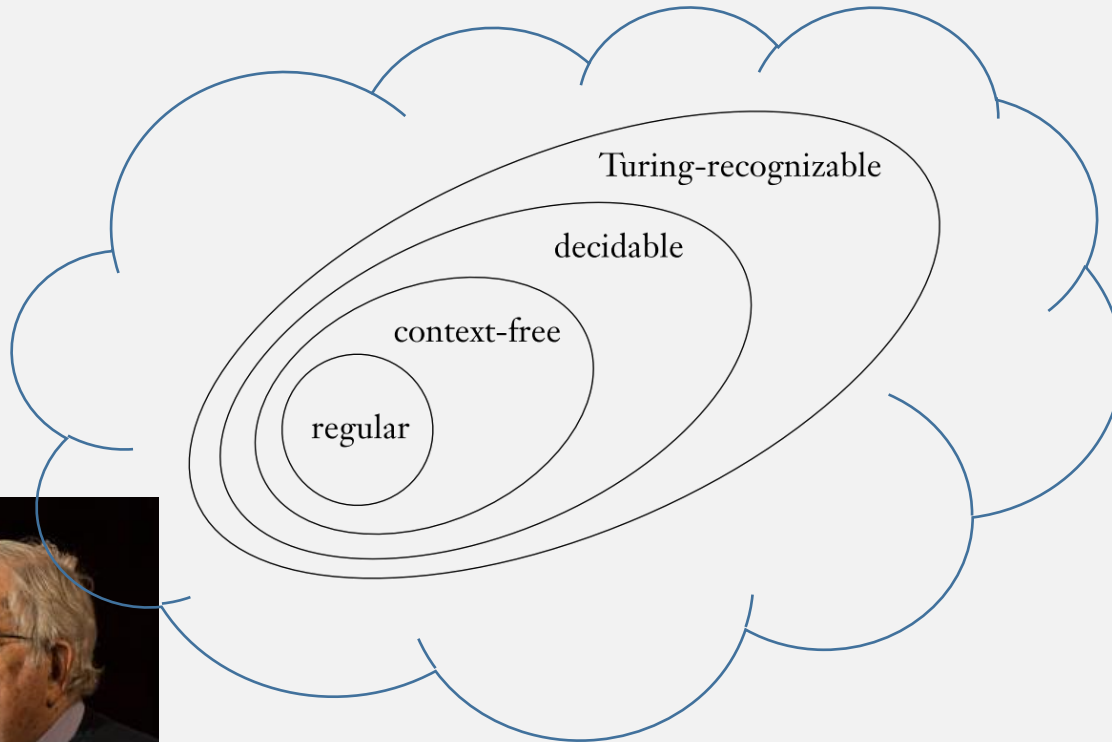  - This TM would be a recognizer but not a decider

Idea: can the TM stop checking after some length?
  - I.e., Is there upper bound on the number of derivation steps?

# Chomsky Normal Form

# Noam Chomsky



Turing-recognizable

decidable

context-free

regular

He came up with this <u>hierarchy</u> of languages

# Chomsky Normal Form

A context-free grammar is in **Chomsky normal form** if every rule is of the form

(non-start) Variables only

$$A \rightarrow BC$$

2 rule shapes

$$A \rightarrow a$$

Terminals only

where $a$ is any terminal and $A$, $B$, and $C$ are any variables—except that $B$ and $C$ may not be the start variable. In addition, we permit the rule $S \rightarrow \varepsilon$, where $S$ is the start variable.

# Chomsky Normal Form Example

- $S \rightarrow AB$
- $B \rightarrow AB$
- $A \rightarrow$ a
- $B \rightarrow$ b

- To generate string of length: 2
  - Use $S$ rule: **1 time**; Use $A$ or $B$ rules: **2 times**
  - $S \Rightarrow AB \Rightarrow \text{a}B \Rightarrow \text{ab}$
  - Derivation total steps: 1 + 2 = $\boxed{3}$
- To generate string of length: 3
  - Use $S$ rule: **1 time**; $A$ rule: **1 time**; $A$ or $B$ rules: **3 times**
  - $S \Rightarrow AB \Rightarrow AAB \Rightarrow \text{a}AB \Rightarrow \text{aa}B \Rightarrow \text{aab}$
  - Derivation total steps: 1 + 1 + 3 = $\boxed{5}$
- To generate string of length: 4
  - Use $S$ rule: **1 time** ; $A$ rule: **2 times**; $A$ or $B$ rules: **4 times**
  - $S \Rightarrow AB \Rightarrow AAB \Rightarrow AAAB \Rightarrow \text{a}AAB \Rightarrow \text{aa}AB \Rightarrow \text{aaa}B \Rightarrow \text{aaab}$
  - Derivation total steps: 3 + 4 = $\boxed{7}$
- ...

A context-free grammar is in ***Chomsky normal form*** if every rule is of the form

☑    $A \rightarrow BC$
     $A \rightarrow a$    2 rule shapes

where $a$ is any terminal and $A$, $B$, and $C$ are any variables—except that $B$ and $C$ may not be the start variable. In addition, we permit the rule $S \rightarrow \varepsilon$, where $S$ is the start variable.

# Chomsky Normal Form: Number of Steps

To generate a string of length $n$:

$n - 1$ steps: to generate $n$ variables

$+\, n$ steps: to turn each variable into a terminal

Total: $2n - 1$ steps

(A *finite* number of steps!)

Makes the string long enough

Convert string to terminals

**Chomsky normal form**

$A \rightarrow BC$   Use $n$-1 times

$A \rightarrow a$   Use $n$ times

<u>Thm</u>: $A_{\mathsf{CFG}}$ is a decidable language

$$A_{\mathsf{CFG}} = \{\langle G, w\rangle \mid G \text{ is a CFG that generates string } w\}$$

<u>Proof</u>: create the decider:

$S$ = "On input $\langle G, w\rangle$, where $G$ is a CFG and $w$ is a string:
1. Convert $G$ to an equivalent grammar in Chomsky normal form.
2. List all derivations with $2n-1$ steps, where $n$ is the length of $w$; except if $n = 0$, then instead list all derivations with one step.
3. If any of these derivations generate $w$, *accept*; if not, *reject*."

We first need to prove this is <u>true for all</u> CFGs!

Step 1: Conversion to Chomsky Normal Form is an algorithm …
Step 2:
Step 3:

Termination argument?

# Thm: Every CFG has a Chomsky Normal Form

Proof: Create algorithm to convert any CFG into Chomsky Normal Form

*Chomsky normal form*

$$A \rightarrow BC$$
$$A \rightarrow a$$

1. Add new start variable $S_0$ that does not appear on any RHS
   - I.e., add rule $S_0 \rightarrow S$, where $S$ is old start var

$$S \rightarrow ASA \mid aB$$
$$A \rightarrow B \mid S$$
$$B \rightarrow b \mid \varepsilon$$

$\Rightarrow$

$$S_0 \rightarrow S$$
$$S \rightarrow ASA \mid aB$$
$$A \rightarrow B \mid S$$
$$B \rightarrow b \mid \varepsilon$$

# <u>Thm</u>: Every CFG has a Chomsky Normal Form

*Chomsky normal form*

$$A \rightarrow BC$$
$$A \rightarrow a$$

1. Add new start variable $S_0$ that does not appear on any RHS
   - I.e., add rule $S_0 \rightarrow S$, where $S$ is old start var

2. Remove all "empty" rules of the form $A \rightarrow \varepsilon$
   - $A$ must not be the start variable
   - Then for every rule with $A$ on RHS, add new rule with $A$ deleted
     - E.g., If $R \rightarrow uAv$ is a rule, add $R \rightarrow uv$
   - Must cover all combinations if $A$ appears more than once in a RHS
     - E.g., if $R \rightarrow uAvAw$ is a rule, add 3 rules: $R \rightarrow uvAw$, $R \rightarrow uAvw$, $R \rightarrow uvw$

$$S_0 \rightarrow S$$
$$S \rightarrow ASA \mid aB \mid \mathbf{a}$$
$$A \rightarrow B \mid S \mid \boldsymbol{\varepsilon}$$
$$B \rightarrow b \mid \varepsilon$$

Then, add

First, remove

$$S_0 \rightarrow S$$
$$S \rightarrow ASA \mid aB \mid a \mid \boldsymbol{SA} \mid \boldsymbol{AS} \mid \boldsymbol{S}$$
$$A \rightarrow B \mid S \mid \varepsilon$$
$$B \rightarrow b$$

Then add, to account for possibly empty $A$

Then, remove

# Thm: Every CFG has a Chomsky Normal Form

*Chomsky normal form*

$$A \rightarrow BC$$
$$A \rightarrow a$$

1. Add new start variable $S_0$ that does not appear on any RHS
   - I.e., add rule $S_0 \rightarrow S$, where $S$ is old start var

2. Remove all "empty" rules of the form $A \rightarrow \varepsilon$
   - $A$ must not be the start variable
   - Then for every rule with $A$ on RHS, add new rule with $A$ deleted
     - E.g., If $R \rightarrow uAv$ is a rule, add $R \rightarrow uv$
   - Must cover all combinations if $A$ appears more than once in a RHS
     - E.g., if $R \rightarrow uAvAw$ is a rule, add 3 rules: $R \rightarrow uvAw, R \rightarrow uAvw, R \rightarrow uvw$

3. Remove all "unit" rules of the form $A \rightarrow B$
   - Then, for every rule $B \rightarrow u$, add rule $A \rightarrow u$

$$S_0 \rightarrow S$$
$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS \mid S$$
$$A \rightarrow B \mid S$$
$$B \rightarrow b$$

Remove, no add (same variable)

$$S_0 \rightarrow S \mid ASA \mid aB \mid a \mid SA \mid AS$$
$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$
$$A \rightarrow B \mid S$$
$$B \rightarrow b$$

Remove, then add $S$ RHSs to $S_0$

$$S_0 \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$
$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$
$$A \rightarrow S \mid b \mid ASA \mid aB \mid a \mid SA \mid AS$$
$$B \rightarrow b$$

Remove, then add $S$ RHSs to $A$

# Thm: Every CFG has a Chomsky Normal Form

***Chomsky normal form***

$$A \to BC$$
$$A \to a$$

1. Add new start variable $S_0$ that does not appear on any RHS
   - I.e., add rule $S_0 \to S$, where $S$ is old start var

2. Remove all "empty" rules of the form $A \to \varepsilon$
   - $A$ must not be the start variable
   - Then for every rule with $A$ on RHS, add new rule with $A$ deleted
     - E.g., If $R \to uAv$ is a rule, add $R \to uv$
   - Must cover all combinations if $A$ appears more than once in a RHS
     - E.g., if $R \to uAvAw$ is a rule, add 3 rules: $R \to uvAw, R \to uAvw, R \to uvw$

$$S_0 \to \boxed{ASA} \mid \boxed{aB} \mid a \mid SA \mid AS$$
$$S \to ASA \mid aB \mid a \mid SA \mid AS$$
$$A \to b \mid ASA \mid aB \mid a \mid SA \mid AS$$
$$B \to b$$

3. Remove all "unit" rules of the form $A \to B$
   - Then, for every rule $B \to u$, add rule $A \to u$

4. Split up rules with RHS longer than length 2
   - E.g., $A \to wxyz$ becomes $A \to wB, B \to xC, C \to yz$

5. Replace all terminals on RHS with new rule
   - E.g., for above, add $W \to w, X \to x, Y \to y, Z \to z$

$$S_0 \to \boxed{AA_1} \mid \boxed{UB} \mid a \mid SA \mid AS$$
$$S \to AA_1 \mid UB \mid a \mid SA \mid AS$$
$$A \to b \mid AA_1 \mid UB \mid a \mid SA \mid AS$$
$$A_1 \to \boxed{SA}$$
$$\boxed{U \to a}$$
$$B \to b$$

# Thm: $A_{\mathsf{CFG}}$ is a decidable language

$$A_{\mathsf{CFG}} = \{\langle G, w\rangle \mid G \text{ is a CFG that generates string } w\}$$

Proof: create the decider:

$S = $ "On input $\langle G, w\rangle$, where $G$ is a CFG and $w$ is a string:
1. Convert $G$ to an equivalent grammar in Chomsky normal form.
2. List all derivations with $2n-1$ steps, where $n$ is the length of $w$; except if $n = 0$, then instead list all derivations with one step.
3. If any of these derivations generate $w$, *accept*; if not, *reject*."

We first need to prove this is <u>true for all</u> CFGs!

Termination argument:
**Step 1**: any CFG has only a finite # rules
**Step 2**: $2n$-1 = finite # of derivations to check
**Step 3**: checking finite number of derivations

# Thm: $E_{\text{CFG}}$ is a decidable language.

$$E_{\text{CFG}} = \{\langle G \rangle \mid G \text{ is a } \boxed{\text{CFG}} \text{ and } L(G) = \emptyset\}$$

Recall:

$$E_{\text{DFA}} = \{\langle A \rangle \mid A \text{ is a } \boxed{\text{DFA}} \text{ and } L(A) = \emptyset\}$$

$T$ = "On input $\langle A \rangle$, where $A$ is a DFA:

1. Mark the start state of $A$.
2. Repeat until no new states get marked:
3. Mark any state that has a transition coming into it from any state that is already marked.
4. If no accept state is marked, *accept*; otherwise, *reject*."

"Reachability" (of accept state from start state) algorithm

Can we compute "reachability" for a CFG?

# Thm: $E_{\text{CFG}}$ is a decidable language.

$$E_{\text{CFG}} = \{\langle G \rangle \mid G \text{ is a CFG and } L(G) = \emptyset\}$$

Proof: create **decider** that calculates reachability for grammar $G$
- Go <u>backwards</u>, start from **terminals**, to <u>avoid getting stuck in looping rules</u>

> Loop marks 1 new variable on each iteration or stops: it eventually terminates because there are a finite # of variables

$R =$ "On input $\langle G \rangle$, where $G$ is a CFG:
1. Mark all terminal symbols in $G$.
2. Repeat until no new variables get marked:
3.     Mark any variable $A$ where $G$ has a rule $A \rightarrow U_1 U_2 \cdots U_k$ and each symbol $U_1, \ldots, U_k$ has already been marked.
4. If the start variable is not marked, *accept*; otherwise, *reject*."

Termination argument?
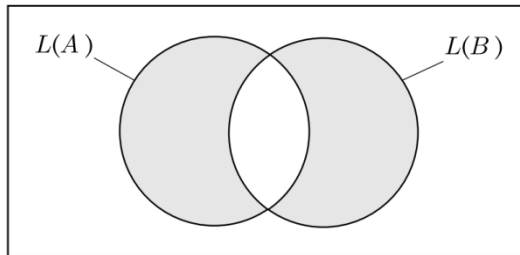
# Thm: $EQ_{\mathsf{CFG}}$ is a decidable language?

$$EQ_{\mathsf{CFG}} = \{\langle G, H \rangle \mid G \text{ and } H \text{ are CFGs and } L(G) = L(H)\}$$

Recall: $$EQ_{\mathsf{DFA}} = \{\langle A, B \rangle \mid A \text{ and } B \text{ are DFAs and } L(A) = L(B)\}$$

- Used Symmetric Difference

$L(A)$      $L(B)$

$$L(C) = \emptyset \text{ iff } L(A) = L(B)$$

  - where $C$ = complement, union, intersection of machines $A$ and $B$

- Can't do this for CFLs!
  - Intersection and complement are not closed for CFLs!!!

# Intersection of CFLs is <u>Not</u> Closed!

<u>Proof</u> (by contradiction), <u>Assume</u> intersection is closed for CFLs
- Then **intersection of these CFLs should be a CFL:**

$$A = \{\mathbf{a}^m \mathbf{b}^n \mathbf{c}^n \mid m, n \geq 0\}$$

$$B = \{\mathbf{a}^n \mathbf{b}^n \mathbf{c}^m \mid m, n \geq 0\}$$

- But $A \cap B = \{\mathbf{a}^n \mathbf{b}^n \mathbf{c}^n \mid n \geq 0\}$

- **… which is not a CFL!** (So we have a contradiction)

# Complement of a CFL is not Closed!

- <u>Assume</u> **CFLs closed under complement,** then:

$$\text{if } G_1 \text{ and } G_2 \text{ context-free}$$

$$\overline{L(G_1)} \text{ and } \overline{L(G_2)} \text{ context-free}$$ From the assumption

$$\overline{L(G_1)} \cup \overline{L(G_2)} \text{ context-free}$$ Union of CFLs is closed

$$\overline{\overline{L(G_1)} \cup \overline{L(G_2)}} \text{ context-free}$$ From the assumption

$$L(G_1) \cap L(G_2) \text{ context-free}$$ DeMorgan's Law!

But intersection is not closed for CFLS (prev slide)

# Thm: $EQ_{\mathsf{CFG}}$ is a decidable language?

$$EQ_{\mathsf{CFG}} = \{\langle G, H \rangle \mid G \text{ and } H \text{ are CFGs and } L(G) = L(H)\}$$
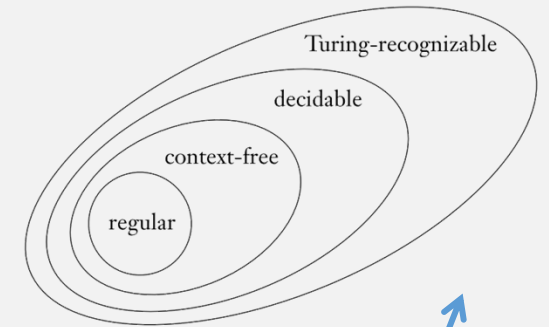
- No!
  - There's no algorithm to decide whether two grammars are equivalent!

- It's not recognizable either! (Can't create <u>any</u> TM to do this!!!)
  - (details later)

- I.e., **this is an <u>impossible computation</u>!**

# *Summary* Algorithms About CFLs

- $A_{\mathsf{CFG}} = \{\langle G, w\rangle|\ G \text{ is a CFG that generates string } w\}$
  - **Decider**: Convert grammar to Chomsky Normal Form
  - Then **check all possible derivations up to length 2|w| - 1 steps**

- $E_{\mathsf{CFG}} = \{\langle G\rangle|\ G \text{ is a CFG and } L(G) = \emptyset\}$
  - **Decider**: Compute "reachability" of start variable from terminals

- $EQ_{\mathsf{CFG}} = \{\langle G, H\rangle|\ G \text{ and } H \text{ are CFGs and } L(G) = L(H)\}$
  - We couldn't prove that this is decidable!
  - (So you cant use this theorem when creating another decider)

# The Limits of Turing Machines?

- **TMs represent all possible "computations"**
  - I.e., **any** (`Python`, `Java`, ...) program you write is a TM

- But **some things are** **not** computable? I.e., some langs are out here **?**

- To explore the limits of computation, we have been studying ...
  ... computation about other computation ...
  - <u>Thought</u>: Is there a decider (algorithm) to determine whether a TM is an decider?

Hmmm, this doesn't feel right ...

*Next time:* Is $A_{\mathsf{TM}}$ decidable?

$$A_{\mathsf{TM}} = \{\langle M, w\rangle \mid M \text{ is a TM and } M \text{ accepts } w\}$$