

# CFGs for Finite Unary Sets

Lecture 11

Binghui Peng

This lecture is based on the paper:

*Simulating Finite Automata with Context-Free Grammars*

by Domaratzki, Pighizzini, and Shallit

*Information Processing Letters*, Volume 84, 2002, pp. 339-344.

# Modified Chomsky Normal Form

A context-free grammar is in **Modified Chomsky Normal Form (CNF)** if its production rules take one of the following forms:

- $A \rightarrow BCD$
- $A \rightarrow BC$
- $A \rightarrow \sigma$
- $A \rightarrow \epsilon$

In this lecture, we focus on **the number of nonterminals (NTs)** required to generate a language.

Allowing rules of the form  $A \rightarrow BCD$  and  $A \rightarrow \epsilon$  simplifies the construction for unary sets. Construction for the conventional Chomsky Normal Form will be more complicated and omitted from the lecture.

# Objective: Succinct CFGs for Unary Sets

**Goal:** Construct a small CFG for a set  $A \subseteq \{\epsilon, a, a^2, \dots, a^n\}$ .

## Plan:

- 1 Construct a CNF grammar for  $\{\epsilon, a, \dots, a^{124}\}$  with 16 NTs.
- 2 Show that for any  $A \subseteq \{\epsilon, a, \dots, a^{124}\}$ , there exists a CNF with 16 NTs.
- 3 Generalize:  $\forall A \subseteq \{\epsilon, a, \dots, a^n\}$ , there exists a CNF with  $O(n^{1/3})$  NTs.
- 4 Analyze the lower bound:  $\exists A \subseteq \{\epsilon, a, \dots, a^n\}$  requiring  $\Omega(n^{1/3})$  NTs.

# Construction Strategy for $\{\epsilon, a, \dots, a^{124}\}$

We define three sets of nonterminals:

- $E_0, \dots, E_4$  such that  $L(E_i) = \{a^i\}$ .
- $F_0, \dots, F_4$  such that  $L(F_i) = \{a^{5i}\}$ .
- $G_0, \dots, G_4$  such that  $L(G_i) = \{a^{25i}\}$ .

Every  $m \in \{0, \dots, 124\}$  can be represented in base 5 as  $m = i + 5j + 25k$ , where  $0 \leq i, j, k \leq 4$ . Thus,  $a^m = L(E_i)L(F_j)L(G_k)$ .

## Production Rules for $E_i, F_j, G_k$

- $E_0 \rightarrow \epsilon, E_1 \rightarrow a, E_2 \rightarrow E_1 E_1, E_3 \rightarrow E_2 E_1, E_4 \rightarrow E_3 E_1$
- $F_0 \rightarrow \epsilon, F_1 \rightarrow E_4 E_1, F_2 \rightarrow F_1 F_1, F_3 \rightarrow F_2 F_1, F_4 \rightarrow F_3 F_1$
- $G_0 \rightarrow \epsilon, G_1 \rightarrow F_4 F_1, G_2 \rightarrow G_1 G_1, G_3 \rightarrow G_2 G_1, G_4 \rightarrow G_3 G_1$

To generate  $\{\epsilon, a, \dots, a^{124}\}$ , include rules:

$$S \rightarrow E_i F_j G_k \quad \forall 0 \leq i, j, k \leq 4$$

## Property of the Construction

Each string  $a^m$  for  $0 \leq m \leq 124$  is generated by exactly one rule  $S \rightarrow E_i F_j G_k$ .

Removing the rule  $S \rightarrow E_i F_j G_k$  specifically omits the string  $a^{i+5j+25k}$  from the language.

# Subsets of $\{\epsilon, a, \dots, a^{124}\}$

**Theorem:** For any  $A \subseteq \{\epsilon, a, \dots, a^{124}\}$ , there exists a CNF grammar with 16 nonterminals.

**Proof:** Starting with the grammar for  $\{\epsilon, a, \dots, a^{124}\}$ , remove the rule  $S \rightarrow E_i F_j G_k$  if  $a^{i+5j+25k} \notin A$ . The number of nonterminals remains 16.

# Generalization to $A \subseteq \{\epsilon, a, \dots, a^{t^3-1}\}$

**Theorem:** For  $t \in \mathbb{N}$ , any  $A \subseteq \{\epsilon, a, \dots, a^{t^3-1}\}$  can be generated by a CNF with  $3t + 1$  nonterminals.

## Construction:

- $E_0, \dots, E_{t-1}$  with  $L(E_i) = \{a^i\}$  ( $t$  NTs)
- $F_0, \dots, F_{t-1}$  with  $L(F_i) = \{a^{ti}\}$  ( $t$  NTs)
- $G_0, \dots, G_{t-1}$  with  $L(G_i) = \{a^{t^2i}\}$  ( $t$  NTs)
- Start symbol  $S$  with rules  $S \rightarrow E_i F_j G_k$  (1 NT)

**Theorem:** For any  $A \subseteq \{\epsilon, a, \dots, a^n\}$ , there exists a CNF grammar with  $O(n^{1/3})$  nonterminals.

By choosing  $t = \lceil (n+1)^{1/3} \rceil$ , we have  $t^3 - 1 \geq n$ , and the construction yields  $3t + 1 = O(n^{1/3})$  nonterminals.

# Results for Standard Chomsky Normal Form

Using the standard definition of CNF (rules  $A \rightarrow BC$  or  $A \rightarrow \sigma$ ):

**Theorem:** For any  $A \subseteq \{\epsilon, a, \dots, a^{t^3-1}\}$ , there exists a standard CNF with  $4t + 1$  nonterminals.

**Corollary:** For any  $A \subseteq \{\epsilon, a, \dots, a^n\}$ , there exists a standard CNF with  $O(n^{1/3})$  nonterminals.

## Optimality of $O(n^{1/3})$ Nonterminals

**Question:** Can we achieve better than  $O(n^{1/3})$  nonterminals for all  $A \subseteq \{\epsilon, a, \dots, a^n\}$ ?

## Optimality of $O(n^{1/3})$ Nonterminals

**Question:** Can we achieve better than  $O(n^{1/3})$  nonterminals for all  $A \subseteq \{\epsilon, a, \dots, a^n\}$ ?

**Answer: No.** There exists a subset  $A$  such that any CNF grammar requires  $\Omega(n^{1/3})$  nonterminals.

## Lower Bound: Counting Argument

- Number of subsets of  $\{\epsilon, a, \dots, a^{n-1}\}$ :  $2^n$ .
- Number of possible CNF grammars with  $t$  nonterminals:
  - Rules of form  $A \rightarrow BC$ :  $t^3$  possibilities.
  - Rules of form  $A \rightarrow \sigma$ :  $O(t)$  possibilities.
  - Total rules:  $t^3 + O(t) \leq 2t^3$ .
  - Total grammars:  $\leq 2^{2t^3}$ .

If  $2^{2t^3} < 2^n$ , there must exist a subset  $A$  that cannot be generated by any CNF with  $t$  nonterminals.

Setting  $2t^3 < n$  implies  $t < (n/2)^{1/3}$ .

Thus, for  $t = \Omega(n^{1/3})$ , there exists a subset  $A$  requiring at least  $t$  nonterminals.