# The Chomsky-Schützenberger Theorem

Deepali Nemade and Nikhil Panwar

Department of Computer Science and Automation,
Indian Institute of Science, Bangalore.

24th November, 2011

## Homomorphism

A *Homomorphism* is a map $h : \Gamma^* \to \Sigma^*$ such that for all $x, y \in \Gamma^*$

$$h(xy) = h(x)h(y)$$
$$h(\epsilon) = \epsilon$$

## Context Free Grammar

A *Context Free Grammar* (CFG) is a quadruple
$$G = (N, \Sigma, P, S)$$

## Context Free Grammar

A *Context Free Grammar* (CFG) is a quadruple
$$G = (N, \Sigma, P, S)$$
where,

## Context Free Grammar

A *Context Free Grammar* (CFG) is a quadruple
$$G = (N, \Sigma, P, S)$$
where,

N is a finite set of Non terminal Symbols,

## Context Free Grammar

A *Context Free Grammar* (CFG) is a quadruple
$$G = (N, \Sigma, P, S)$$
where,

N is a finite set of Non terminal Symbols,
$\Sigma$ is a finite set of Terminal Symbols,

## Context Free Grammar

A *Context Free Grammar* (CFG) is a quadruple
$$G = (N, \Sigma, P, S)$$
where,

        N is a finite set of Non terminal Symbols,

        $\Sigma$ is a finite set of Terminal Symbols,

        P is a finite subset of $N \times (N \cup \Sigma)^*$ (set of productions)

## Context Free Grammar

A *Context Free Grammar* (CFG) is a quadruple
$$G = (N, \Sigma, P, S)$$
where,

      N is a finite set of Non terminal Symbols,

      $\Sigma$ is a finite set of Terminal Symbols,

      P is a finite subset of $N \times (N \cup \Sigma)^*$ (set of productions)

      $S \in N$ is the start symbol

# Chomsky Normal Form

A CFG is in "Chomsky normal form(CNF)" if all productions are of
the form

# Chomsky Normal Form

A CFG is in "Chomsky normal form(CNF)" if all productions are of the form

$$A \rightarrow BC$$

## Chomsky Normal Form

A CFG is in "Chomsky normal form(CNF)" if all productions are of the form

$$A \rightarrow BC$$
or

## Chomsky Normal Form

A CFG is in "Chomsky normal form(CNF)" if all productions are of the form

$$A \rightarrow BC$$
$$\text{or}$$
$$A \rightarrow a$$

## Chomsky Normal Form

A CFG is in "Chomsky normal form(CNF)" if all productions are of the form

$$A \rightarrow BC$$
$$\text{or}$$
$$A \rightarrow a$$

where,

$A, B, C \in N$ and $a \in \Sigma$

## $PAREN_n$

$PAREN_n$ denote the language consisting of all balanced strings of parentheses of $n$ distinct types.

## $PAREN_n$

$PAREN_n$ denote the language consisting of all balanced strings of parentheses of $n$ distinct types.

This language is generated by the grammar-

$$S \to [^1 S]^1 | [^2 S]^2 | \ldots\ldots\ldots | [^n S]^n | \varepsilon$$

# $PAREN_n$

$PAREN_n$ denote the language consisting of all balanced strings of parentheses of $n$ distinct types.

This language is generated by the grammar-

$$S \to [^1 S]^1 | [^2 S]^2 | .............. | [^n S]^n | \varepsilon$$

The languages $PAREN_n$ are called *Dyck Languages* in the literature.

## $PAREN_n$

$PAREN_n$ denote the language consisting of all balanced strings of parentheses of $n$ distinct types.

This language is generated by the grammar-

$$S \rightarrow [^1S]^1|[^2S]^2|.............|[^nS]^n|\varepsilon$$

The languages $PAREN_n$ are called *Dyck Languages* in the literature.

**Example:** $[^1[^2]^2]^1[^3]^3$

## Theorem

### The Chomsky-Schützenberger Theorem

*Every context-free language is a homomorphic image of the intersection of a parenthesis language and a regular set. In other words, for every CFL A, there is an $n >= 0$, a regular set R, and a homomorphism h such that*

$$A = h(PAREN_n \cap R)$$

## Proof

Let $G = (N, \Sigma, P, S)$ be an arbitrary CFG in Chomsky Normal Form (CNF).

Let the productions in P are denoted by $\pi, \rho, \sigma \ldots$

## Proof

Let $G = (N, \Sigma, P, S)$ be an arbitrary CFG in Chomsky Normal Form (CNF).

Let the productions in P are denoted by $\pi, \rho, \sigma \ldots$

For $\pi \in P$, define $\pi^{'}$

$$A \to [^1_\pi B]^1_\pi [^2_\pi C]^2_\pi \qquad \text{if } \pi = A \to BC,$$

$$A \to [^1_\pi]^1_\pi [^2_\pi]^2_\pi \qquad \text{if } \pi = A \to a$$

## Proof

Let $G = (N, \Sigma, P, S)$ be an arbitrary CFG in Chomsky Normal
Form (CNF).
Let the productions in P are denoted by $\pi, \rho, \sigma \ldots$

For $\pi \in P$, define $\pi^{'}$

$$A \to [^1_\pi B]^1_\pi [^2_\pi C]^2_\pi \qquad \text{if } \pi = A \to BC,$$
$$A \to [^1_\pi]^1_\pi [^2_\pi]^2_\pi \qquad \text{if } \pi = A \to a$$

and define the grammar $G^{'} = (N, \Sigma, P^{'}, S)$ with

$$\Gamma = \{[^1_\pi, ]^1_\pi, [^2_\pi, ]^2_\pi | \pi \in P\}$$
$$P^{'} = \{\pi^{'} | \pi \in P\}$$

## Proof continued...

$L(G^{'}) \subseteq PAREN_{\Gamma}$

## Proof continued...

$L(G^{'}) \subseteq PAREN_{\Gamma}$
Properties satisfied by strings in $L(G^{'})$ that are not satisfied by
strings in $PAREN_{\Gamma}$ in general:

## Proof continued...

$L(G^{'}) \subseteq PAREN_\Gamma$

Properties satisfied by strings in $L(G^{'})$ that are not satisfied by strings in $PAREN_\Gamma$ in general:

- **Property 1 :** Every $]_\pi^1$ is immediately followed by a $[_\pi^2$.

## Proof continued...

$L(G') \subseteq PAREN_\Gamma$
Properties satisfied by strings in $L(G')$ that are not satisfied by strings in $PAREN_\Gamma$ in general:

- **Property 1 :** Every $]_\pi^1$ is immediately followed by a $[_\pi^2$.
- **Property 2 :** No $]_\pi^2$ is immediately followed by a left parenthesis.

## Proof continued...

$L(G') \subseteq PAREN_\Gamma$

Properties satisfied by strings in $L(G')$ that are not satisfied by strings in $PAREN_\Gamma$ in general:

- **Property 1 :** Every $]_\pi^1$ is immediately followed by a $[_\pi^2$.

- **Property 2 :** No $]_\pi^2$ is immediately followed by a left parenthesis.

- **Property 3 :** If $\pi = A \to BC$, then every $[_\pi^1$ is immediately followed by $[_\rho^1$ for some $\rho \in P$ with left hand side B, and every $[_\pi^2$ is immediately followed by $[_\sigma^1$ for some $\sigma \in P$ with left-hand side $C$.

## Proof continued...

$L(G^{'}) \subseteq PAREN_\Gamma$
Properties satisfied by strings in $L(G^{'})$ that are not satisfied by
strings in $PAREN_\Gamma$ in general:

- **Property 1 :** Every $]_\pi^1$ is immediately followed by a $[_\pi^2$.

- **Property 2 :** No $]_\pi^2$ is immediately followed by a left
  parenthesis.

- **Property 3 :** If $\pi = A \to BC$, then every $[_\pi^1$ is immediately
  followed by $[_\rho^1$ for some $\rho \in P$ with left hand side B, and every
  $[_\pi^2$ is immediately followed by $[_\sigma^1$ for some $\sigma \in P$ with left-hand
  side $C$.

- **Property 4 :** If $\pi = A \to a$, then every $[_\pi^1$ is immediately
  followed by $]_\pi^1$ and every $[_\pi^2$ is immediately followed by $]_\pi^2$.

## Proof continued...

In addition, all strings $x$ such that $A \xrightarrow[G']{*} x$ satisfy the property

- **Property** $(v_A)$: The string $x$ begins with $[^1_\pi$ for some $\pi \in P$ with left-hand side A.

## Proof continued...

In addition, all strings $x$ such that $A \xrightarrow[G']{*} x$ satisfy the property

- **Property** ($v_A$): The string $x$ begins with $[^1_\pi$ for some $\pi \in P$ with left-hand side A.

Now we can define a regular expression that satisfies all the above properties as:

$$R_A = \{ \ x \in \Gamma^* \mid x \text{ satisfies } property\,1 \text{ through } (v_A) \ \}$$

## Proof continued...

### Lemma

$$A \xrightarrow[G']{*} x \iff x \in (PAREN_\Gamma \cap R_A)$$

## Proof Continued...

**Proof of Lemma:**

Prove $\Rightarrow$: $A \xrightarrow[G']{*} x \Rightarrow x \in (PAREN_\Gamma \cap R_A)$

## Proof Continued...

**Proof of Lemma:**

Prove $\Rightarrow$: $A \xrightarrow[G']{*} x \Rightarrow x \in (PAREN_\Gamma \cap R_A)$

Applying induction on the length of derivation:

## Proof Continued...

**Proof of Lemma:**

Prove $\Rightarrow$: $A \xrightarrow[G']{*} x \Rightarrow x \in (PAREN_\Gamma \cap R_A)$

Applying induction on the length of derivation:

**Basis:** n=1

$$A \to [_\pi^1 B]_\pi^1 [_\pi^2 C]_\pi^2$$
$$A \to [_\pi^1]_\pi^1 [_\pi^2]_\pi^2$$

Since RHS satisfies all properties so it is true for $n = 1$ .

## Proof Continued...

**Proof of Lemma:**

Prove $\Rightarrow$: $A \xrightarrow[G']{*} x \Rightarrow x \in (PAREN_\Gamma \cap R_A)$

Applying induction on the length of derivation:

**Basis:** n=1

$$A \rightarrow [^1_\pi B]^1_\pi [^2_\pi C]^2_\pi$$
$$A \rightarrow [^1_\pi]^1_\pi [^2_\pi]^2_\pi$$

Since RHS satisfies all properties so it is true for $n = 1$ .

**Induction Hypothesis:** Let $A \xrightarrow[G']{*} \alpha$

where, $\alpha$ is a sentential form of length n that satisfies all properties.

## Proof Continued...

**Proof of Lemma:**

Prove $\Rightarrow$: $A \xrightarrow[G']{*} x \Rightarrow x \in (PAREN_\Gamma \cap R_A)$

Applying induction on the length of derivation:

**Basis:** n=1

$$A \to [_\pi^1 B]_\pi^1 [_\pi^2 C]_\pi^2$$
$$A \to [_\pi^1 ]_\pi^1 [_\pi^2 ]_\pi^2$$

Since RHS satisfies all properties so it is true for $n = 1$ .

**Induction Hypothesis:** Let $A \xrightarrow[G']{*} \alpha$

where, $\alpha$ is a sentential form of length n that satisfies all properties.

**Induction Step:** Proving it for $n + 1$ length of derivation

$$A \xrightarrow[G']{1} [_\pi^1 B]_\pi^1 [_\pi^2 C]_\pi^2 \xrightarrow[G']{*} \alpha$$

## Proof Continued...

Prove $\Leftarrow$: $x \in (PAREN_\Gamma \cap R_A) \Rightarrow A \xrightarrow[G']{*} x$

## Proof Continued...

Prove $\Leftarrow$: $x \in (PAREN_\Gamma \cap R_A) \Rightarrow A \xrightarrow[G']{*} x$

Applying induction on the length of x:

## Proof Continued...

Prove $\Leftarrow$: $x \in (PAREN_\Gamma \cap R_A) \Rightarrow A \xrightarrow[G']{*} x$

Applying induction on the length of x:

It follows from properties that x is a string of balanced parentheses of the form

$$x = [^1_\pi y]^1_\pi [^2_\pi z]^2_\pi$$

for some $y, z \in \Gamma^*$ and $\pi$ with left hand side A.

## Proof continued...

If $\pi = A \to BC$ , then

## Proof continued...

If $\pi = A \rightarrow BC$ , then

From property 3, y satisfies $(v_B)$ and z satisfies $(v_C)$.
Also y and z are balanced.

Thus $y \in PAREN_\Gamma \cap R_B$ and $z \in PAREN_\Gamma \cap R_c$ .

## Proof continued...

If $\pi = A \rightarrow BC$ , then

From property 3, y satisfies $(v_B)$ and z satisfies $(v_C)$.
Also y and z are balanced.

Thus $y \in PAREN_\Gamma \cap R_B$ and $z \in PAREN_\Gamma \cap R_c$ .

By induction hypothesis, $B \xrightarrow[G']{*} y$ and $C \xrightarrow[G']{*} z$ therefore,

$$A \xrightarrow[G']{1} [^1_\pi B]^1_\pi [^2_\pi C]^2_\pi \xrightarrow[G']{*} [^1_\pi x]^1_\pi [^2_\pi y]^2_\pi = x$$

## Proof continued...

If $\pi = A \rightarrow a$ , then

## Proof continued...

If $\pi = A \rightarrow a$ , then

From property 4, $y = z = \epsilon$ , and

$$A \rightarrow [^1_\pi]^1_\pi[^2_\pi]^2_\pi = x$$

## Proof continued...

If $\pi = A \to a$ , then

From property 4, $y = z = \epsilon$ , and

$$A \to [^1_\pi]^1_\pi [^2_\pi]^2_\pi = x$$

It follows from Lemma that $L(G') = PAREN_\Gamma \cap R_S$ .

## Proof continued...

**Applying Homomorphism**

Define homomorphism $h : \Gamma^* \rightarrow \Sigma^*$ as follows:

## Proof continued...

**Applying Homomorphism**

Define homomorphism $h : \Gamma^* \to \Sigma^*$ as follows:

For $\pi$ of the form $A \to BC$, take
$$h([^1_\pi) = h(]^1_\pi) = h([^2_\pi) = h(]^2_\pi) = \epsilon,$$

## Proof continued...

**Applying Homomorphism**
Define homomorphism $h : \Gamma^* \to \Sigma^*$ as follows:

For $\pi$ of the form $A \to BC$, take
$$h([_\pi^1) = h(]_\pi^1) = h([_\pi^2) = h(]_\pi^2) = \epsilon,$$

For $\pi$ of the form $A \to a$ , take
$$h(]_\pi^1) = h([_\pi^2) = h(]_\pi^2) = \epsilon,$$
$$h([_\pi^1) = a$$

## Proof continued...

**Applying Homomorphism**

Define homomorphism $h : \Gamma^* \to \Sigma^*$ as follows:

For $\pi$ of the form $A \to BC$, take
$$h([^1_\pi) = h(]^1_\pi) = h([^2_\pi) = h(]^2_\pi) = \epsilon,$$

For $\pi$ of the form $A \to a$ , take
$$h(]^1_\pi) = h([^2_\pi) = h(]^2_\pi) = \epsilon,$$
$$h([^1_\pi) = a$$

Applying $h$ to the production $\pi$ of $P'$ gives the production $\pi$ of P
thus $L(G) = h(L(G')) = h(PAREN_\Gamma \cap R_S)$.

## Proof continued...

**Applying Homomorphism**
Define homomorphism $h : \Gamma^* \to \Sigma^*$ as follows:

For $\pi$ of the form $A \to BC$, take
$$h([^1_\pi) = h(]^1_\pi) = h([^2_\pi) = h(]^2_\pi) = \epsilon,$$

For $\pi$ of the form $A \to a$ , take
$$h(]^1_\pi) = h([^2_\pi) = h(]^2_\pi) = \epsilon,$$
$$h([^1_\pi) = a$$

Applying $h$ to the production $\pi$ of $P'$ gives the production $\pi$ of P
thus $L(G) = h(L(G')) = h(PAREN_\Gamma \cap R_S)$.

This completes the proof of the Chomsky-Schützenberger theorem.

### Example

Apply the theorem on $\{ a^n b^n \mid n >= 0 \} = h(L(G^{'}) \cap R)$

## Example continued...

Let $G = (N, \Sigma, P, S)$ be the CFG corresponding to our CFL $\{a^n b^n | n = 0\}$ where,

$$S \rightarrow aSb$$
$$S \rightarrow \epsilon$$

Converting it to CNF we get

$$
\begin{array}{ll}
\pi & S \rightarrow AC \\
\sigma & C \rightarrow SB \\
\rho & A \rightarrow a \\
\lambda & B \rightarrow b \\
\gamma & S \rightarrow \epsilon
\end{array}
$$

## Example continued...

Now define grammar $G' = (N, \Gamma, P', S)$ with

$$\Gamma = \{[_\pi^1, ]_\pi^1, [_\pi^2, ]_\pi^2 | \pi \in P\} \ ,$$
$$P' = \{\pi' | \pi \in P\} \ .$$

where production $P'$ are,

$$\pi \qquad S \rightarrow [_\pi^1 A]_\pi^1 [_\pi^2 C]_\pi^2$$

$$\sigma \qquad C \rightarrow [_\sigma^1 S]_\sigma^1 [_\sigma^2 B]_\sigma^2$$

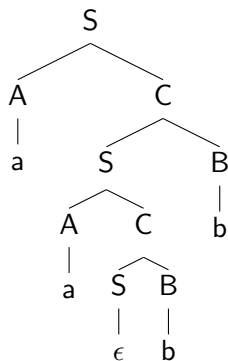$$\rho \qquad A \rightarrow [_\rho^1]_\rho^1$$

$$\lambda \qquad B \rightarrow [_\lambda^1]_\lambda^1$$

$$\gamma \qquad S \rightarrow [_\gamma^1]_\gamma^1$$

## Example continued...

Consider a string x generated by grammar G: $x = aabb$

The parse tree generated by G is

THANK YOU