

## Supplementary Lecture G

### The Chomsky–Schützenberger Theorem

Let  $\text{PAREN}_n$  denote the language consisting of all balanced strings of parentheses of  $n$  distinct types. This language is generated by the grammar

$$S \rightarrow [S]_1 \mid [S]_2 \mid \cdots \mid [S]_n \mid SS \mid \epsilon.$$

The languages  $\text{PAREN}_n$  are sometimes called *Dyck languages* in the literature.

The following theorem shows that the parenthesis languages  $\text{PAREN}_n$  play a special role in the theory of context-free languages: *every* CFL is essentially a parenthesis language modified in some relatively simple way. In a sense, balanced parentheses capture the essential structure of CFLs that differentiates them from the regular sets.

**Theorem G.1 (Chomsky–Schützenberger)** *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  $n \geq 0$ , a regular set  $R$ , and a homomorphism  $h$  such that*

$$A = h(\text{PAREN}_n \cap R).$$

Recall from Lecture 10 that a *homomorphism* is a map  $h : \Gamma^* \rightarrow \Sigma^*$  such that  $h(xy) = h(x)h(y)$  for all  $x, y \in \Gamma^*$ . It follows from this property that  $h(\epsilon) = \epsilon$  and that  $h$  is completely determined by its values on  $\Gamma$ . The

homomorphic image of a set  $B \subseteq \Gamma^*$  under  $h$  is the set  $\{h(x) \mid x \in B\} \subseteq \Sigma^*$ , denoted  $h(B)$ .

*Proof.* Let  $G = (N, \Sigma, P, S)$  be an arbitrary CFG in Chomsky normal form. Denote productions in  $P$  by  $\pi, \rho, \sigma, \dots$

For  $\pi \in P$ , define

$$\pi' = \begin{cases} A \rightarrow \begin{matrix} 1 & 1^2 & 2 \\ [B] & [C] & \\ \pi & \pi & \pi \end{matrix} & \text{if } \pi = A \rightarrow BC, \\ A \rightarrow \begin{matrix} 1 & 1^2 & 2 \\ [ ] & [ ] & \\ \pi & \pi & \pi \end{matrix} & \text{if } \pi = A \rightarrow a; \end{cases}$$

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

$$\Gamma = \left\{ \begin{matrix} 1 & 1 & 2 \\ [ , ] , [ , ] & \mid & \pi \in P \end{matrix} \right\},$$

$$P' = \{\pi' \mid \pi \in P\}.$$

The idea here is that a balanced string of parentheses generated by  $G'$  encodes a corresponding string generated by  $G$  along with its parse tree.

Let  $\text{PAREN}_\Gamma$  be the parenthesis language over parentheses  $\Gamma$ . Surely  $L(G') \subseteq \text{PAREN}_\Gamma$ , since the productions of  $G'$  generate parentheses in well-nested matched pairs. However, not all strings in  $\text{PAREN}_\Gamma$  are generated by  $G'$ . Here are some properties satisfied by strings in  $L(G')$  that are not satisfied by strings in  $\text{PAREN}_\Gamma$  in general:

- (i) Every  $\begin{matrix} 1 \\ \pi \end{matrix}$  is immediately followed by a  $\begin{matrix} 2 \\ \pi \end{matrix}$ .
- (ii) No  $\begin{matrix} 2 \\ \pi \end{matrix}$  is immediately followed by a left parenthesis.
- (iii) If  $\pi = A \rightarrow BC$ , then every  $\begin{matrix} 1 \\ \pi \end{matrix}$  is immediately followed by  $\begin{matrix} 1 \\ \rho \end{matrix}$  for some  $\rho \in P$  with left-hand side  $B$ , and every  $\begin{matrix} 2 \\ \pi \end{matrix}$  is immediately followed by  $\begin{matrix} 1 \\ \sigma \end{matrix}$  for some  $\sigma \in P$  with left-hand side  $C$ .
- (iv) If  $\pi = A \rightarrow a$ , then every  $\begin{matrix} 1 \\ \pi \end{matrix}$  is immediately followed by  $\begin{matrix} 1 \\ \pi \end{matrix}$  and every  $\begin{matrix} 2 \\ \pi \end{matrix}$  is immediately followed by  $\begin{matrix} 1 \\ \pi \end{matrix}$ .

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

- (v<sub>A</sub>) The string  $x$  begins with  $\begin{matrix} 1 \\ \pi \end{matrix}$  for some  $\pi \in P$  with left-hand side  $A$ .

Each of the properties (i) through (v<sub>A</sub>) can be described by a regular expression; thus the sets

$$R_A = \{x \in \Gamma^* \mid x \text{ satisfies (i) through (v}_A)\}$$

are regular. We claim

Lemma G.2  $A \xrightarrow{G'}^* x \iff x \in \text{PAREN}_\Gamma \cap R_A.$

*Proof.* The direction ( $\Rightarrow$ ) is a straightforward proof by induction on the length of the derivation. For the direction ( $\Leftarrow$ ), suppose  $x \in \text{PAREN}_\Gamma \cap R_A.$  We proceed by induction on the length of  $x.$  It follows from properties (i) through ( $v_A$ ) and the fact that  $x$  is a string of balanced parentheses that  $x$  is of the form

$$x = \begin{matrix} 1 & 12 & 2 \\ \pi & \pi\pi & \pi \end{matrix} [y] [z]$$

for some  $y, z \in \Gamma^*$  and  $\pi$  with left-hand side  $A.$  If  $\pi = A \rightarrow BC,$  then from property (iii),  $y$  satisfies ( $v_B$ ) and  $z$  satisfies ( $v_C$ ). Also,  $y$  and  $z$  satisfy (i) through (iv) and are balanced. Thus  $y \in \text{PAREN}_\Gamma \cap R_B$  and  $z \in \text{PAREN}_\Gamma \cap R_C.$  By the induction hypothesis,  $B \xrightarrow{G'}^* y$  and  $C \xrightarrow{G'}^* z;$  therefore,

$$A \xrightarrow{G'}^* \begin{matrix} 1 & 12 & 2 \\ \pi & \pi\pi & \pi \end{matrix} [B] [C] \xrightarrow{G'}^* \begin{matrix} 1 & 12 & 2 \\ \pi & \pi\pi & \pi \end{matrix} [y] [z] = x.$$

If  $\pi = A \rightarrow a,$  then from property (iv),  $y = z = \epsilon,$  and

$$A \xrightarrow{G'}^* \begin{matrix} 1 & 12 & 2 \\ \pi & \pi\pi & \pi \end{matrix} [ ] [ ] = x. \quad \square$$

It follows from Lemma G.2 that  $L(G') = \text{PAREN}_\Gamma \cap R_S.$  Now define the homomorphism  $h : \Gamma^* \rightarrow \Sigma^*$  as follows. For  $\pi$  of the form  $A \rightarrow BC,$  take

$$h\left(\begin{matrix} 1 \\ \pi \end{matrix}\right) = h\left(\begin{matrix} 1 \\ \pi \end{matrix}\right) = h\left(\begin{matrix} 2 \\ \pi \end{matrix}\right) = h\left(\begin{matrix} 2 \\ \pi \end{matrix}\right) = \epsilon.$$

For  $\pi$  of the form  $A \rightarrow a,$  take

$$h\left(\begin{matrix} 1 \\ \pi \end{matrix}\right) = h\left(\begin{matrix} 2 \\ \pi \end{matrix}\right) = h\left(\begin{matrix} 2 \\ \pi \end{matrix}\right) = \epsilon,$$

$$h\left(\begin{matrix} 1 \\ \pi \end{matrix}\right) = a.$$

Applying  $h$  to the production  $\pi'$  of  $P'$  gives the production  $\pi$  of  $P;$  thus  $L(G) = h(L(G')) = h(\text{PAREN}_\Gamma \cap R_S).$  This completes the proof of the Chomsky–Schützenberger theorem.  $\square$

### Historical Notes

The pivotal importance of balanced parentheses in the theory of context-free languages was recognized quite early on. The Chomsky–Schützenberger theorem is due to Chomsky and Schützenberger [19, 22].