# Guaranteed Parsing

*The CYK algorithm for parsing works with*
*any context-free language.*

## CYK

The CYK algorithm is named after Cocke, Younger and Kasami. It assumes CFG $G$ in Chomsky Normal Form.

It uses dynamic programming...

## Generalized Problem

Say input $w = w_1 w_2 \ldots w_n$. Then, let $w_{i,j}$ be substring $w_i w_{i+1} \ldots w_j$. The problem one solves is:

*which variables produce which substrings.*

## A Recursive Formula

In general, suppose we want to know if variable $A \stackrel{*}{\Longrightarrow} w_{i,j}$ where $|w_{i,j}| \geq 2$:

The first step in derivation must be production of form $A \rightarrow EF$. So, $w_{i,j}$ can be split into two pieces: the first generated from $E$ and the second from $F$. (But we don't know where split occurs.) Hence the recursive formula:

*Consider all productions $A \rightarrow EF$. For all possible $k$ from $i$ up to $j - 1$, ask whether $E \stackrel{*}{\Longrightarrow} w_{i,k}$ and $F \stackrel{*}{\Longrightarrow} w_{k+1,j}$.*

## The Overall Algorithm

To make efficient, answer question for smaller strings first, keeping results in a table.

**CYK algorithm.** *1. Start by answering for each $i$ and each variable $A$ whether $A \stackrel{*}{\Longrightarrow} w_{i,i}$. (Look at unit productions.)*

*2. Then answer for each $i$ and each variable $A$ whether $A \stackrel{*}{\Longrightarrow} w_{i,i+1}$. (Use recursive formula.)*

*3. Repeat for all $w_{i,i+2}$, then all $w_{i,i+3}$, and so on.*

Eventually, we determine variables for $w = w_{1,n}$.

Consider CFG with start variable $S$:

$$S \rightarrow ST \mid TU \mid \texttt{b}$$
$$T \rightarrow SU \mid \texttt{a}$$
$$U \rightarrow SS \mid \texttt{b}$$

Consider input string $w = \texttt{aababb}$.

## Example Table for aababb

finish

| | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| 1 | $T$ | . | . | . | $S$ | $S, T, U$ |
| 2 | | $T$ | $S$ | $S$ | $S, T, U$ | $S, T, U$ |
| 3 | | | $S, U$ | $S$ | $T, U$ | $S, T, U$ |
| 4 | | | | $T$ | $S$ | $S, T, U$ |
| 5 | | | | | $S, U$ | $T, U$ |
| 6 | | | | | | $S, U$ |

start

For example, entry in row 3 column 5 says that variables $T$ and $U$ generate $w_{3,5}$: $T$ is here since $T \to SU$ and $S \stackrel{*}{\Longrightarrow} w_{3,4}$, $U \stackrel{*}{\Longrightarrow} w_{5,5}$.

For the earlier grammar

> 1: $S \to \mathtt{r}L$
>
> 2: $L \to L\,\mathtt{,}\,I$
>
> 3: $L \to I$
>
> 4: $I \to \mathtt{v}$

Convert to Chomsky Normal Form, and then apply the CYK algorithm to the string $\mathtt{rv,v,v}$.

# Solution to Practice

$S \rightarrow RL$

$L \rightarrow LF \mid$ v

$F \rightarrow CI$

$I \rightarrow$ v

$C \rightarrow$ ,

$R \rightarrow$ r

|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| 1 | $R$ | $S$ | . | $S$ | . | $S$ |
| 2 |   | $L, I$ | . | $L$ | . | $L$ |
| 3 |   |   | $C$ | $F$ | . | . |
| 4 |   |   |   | $L, I$ | . | $L$ |
| 5 |   |   |   |   | $C$ | $F$ |
| 6 |   |   |   |   |   | $L, I$ |

**Summary**

The CYK algorithm can be used to parse any context-free language.