

Context-Free Grammars

An **alphabet** is a set of symbols —finite or infinite.

Given an alphabet A the set of all *finite* sequences with items in A is written A^* .

A **language (over A)** is a subset —finite or infinite— of A^* .

Some languages:

- $\{[0, 0], [0, 1], [1, 0], [1, 1]\}$ a finite language over $\{0, 1\}$
- $\{0, 1\}^* = \{[], [0], [1], [0, 0], [0, 1], [1, 0], [1, 1], [0, 0, 0], \dots\}$,
an infinite language over $\{0, 1\}$
- $\{[], [0], [1], [0, 0], [1, 1], [0, 0, 0], [0, 1, 0], [1, 0, 1], [1, 1, 1], \dots\}$
the infinite set of palindromes over $\{0, 1\}$
- $\{[1], [2], [1, '+', 1], [1, '+', 2], [2, '+', 2], [1, '+', 1, '+', 1], \dots\}$
arithmetic expressions using 1, 2, and +

As languages are often infinite sets, we need *finite* descriptions of *infinite* languages.

A **context-free language (CFL)** is a language defined by a *context-free grammar* (below).

Context free languages have numerous applications in data formats, programming languages, communication protocols, etc.

Languages that are too complex to be context free are nevertheless often described by first describing a context free language and then imposing additional restrictions

- E.g. the set of all syntactically correct Java classes is context free

```
class C { int i = 13.5 ; }
```

- The set of all type correct Java classes is *not* context free.

A **context free grammar (CFG)** $(A, N, P, n_{\text{start}})$ consists of

- An alphabet A (i.e. a set of symbols)
- A finite set of **nonterminal symbols** N disjoint from A .
- A finite set of **production rules** of the form $n \rightarrow \alpha$, where $n \in N$ and $\alpha \in (N \cup A)^*$
- A **starting nonterminal** $n_{\text{start}} \in N$.

Example

$G_0 = (\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, (,), -, \#\}, \{pn, d\}, P_0, pn)$

where P_0 is¹

$$\begin{aligned} & \{pn \rightarrow (ddd)\#ddd - dddd, \\ & d \rightarrow 0, d \rightarrow 1, d \rightarrow 2, d \rightarrow 3, d \rightarrow 4, \\ & d \rightarrow 5, d \rightarrow 6, d \rightarrow 7, d \rightarrow 8, d \rightarrow 9\} \end{aligned}$$

¹ In formal language theory, it is usual to write sequences just by listing the items. E.g. ddd instead of $[d, d, d]$. This creates a usually harmless ambiguity between symbols and sequences of length 1. Catenation is written as st rather than $s^{\wedge}t$. The empty sequence is written as either ϵ or as nothing at all rather than $[]$.

For a given grammar $G = (A, N, P, n_{\text{start}})$, define a relation \Longrightarrow on sequences as follows²

$$\alpha \Longrightarrow \beta$$

if and only if

$$\alpha, \beta \in (N \cup A)^* \text{ and}$$

there are sequences $\gamma, \delta, \theta \in (N \cup A)^*$ and an $n \in N$ such that $\alpha = \gamma n \theta$ and $\beta = \gamma \delta \theta$ and $(n \rightarrow \delta) \in P$

If $\alpha \Longrightarrow \beta$, we say that we can **derive** β from α in one **step**.

Each **derivation step** $\gamma n \theta \Longrightarrow \gamma \delta \theta$ replaces one occurrence of some nonterminal symbol n with δ where $(n \rightarrow \delta) \in P$. For example the following are derivation steps for G_0

$$d d d \Longrightarrow 8 d d$$

$$d d d \Longrightarrow d 0 d$$

$$p n \Longrightarrow (d d d) \# d d d - d d d d$$

Thus each grammar defines a directed graph in which the nodes are elements of $(A \cup N)^*$ and, for each α and β , there is an edge from α to β iff $\alpha \Longrightarrow \beta$.

² Using the notation of the rest of the course, we would write the last line of this definition as “such that $\alpha = \gamma \hat{[n]} \hat{\theta}$ and $\beta = \gamma \hat{\delta} \hat{\theta}$ and $(n \rightarrow \delta) \in P$ ”.

A **derivation** is a finite path in this graph. We write $\alpha \xRightarrow{*} \beta$ to mean there is a derivation that starts at α and ends at β . I.e. $\alpha \xRightarrow{*} \beta$ means that we can transform α into β via 0 or more derivation steps.

For example

$$pn \xRightarrow{*} (dd9)\#ddd - d d 0 d$$

The **language defined by** a CFG $(A, N, P, n_{\text{start}})$ is the set of sequences in $\alpha \in A^*$ such that $n_{\text{start}} \xRightarrow{*} \alpha$.

For example, the language defined by G_0 includes the sequence

$$(709)\#867 - 5309$$

To prove this, all we need to do is show 1 derivation (of the many) from pn .

$$\begin{aligned} pn &\implies (ddd)\#ddd - dddd \implies (dd9)\#ddd - dddd \\ &\implies (dd9)\#ddd - dd0d \implies (d09)\#ddd - dd0d \\ &\implies (709)\#ddd - dd0d \implies (709)\#ddd - dd0d \\ &\implies (709)\#ddd - dd09 \implies (709)\#ddd - d309 \\ &\implies (709)\#8dd - d309 \implies (709)\#86d - d309 \\ &\implies (709)\#867 - d309 \implies (709)\#867 - 5309 \end{aligned}$$

Another example. Let $G_1 = (A_1, N_1, P_1, \text{block})$

- $A_1 = \{+, *, /, -, (,), <, :=, \text{if, then, while, do, else, end}\} \cup \mathcal{I} \cup \mathcal{N}$, where \mathcal{I} is a finite set of identifiers disjoint from $\{\text{if, then, while, do, else, end}\}$ and \mathcal{N} is some finite subset of \mathbb{N} .
- $N_1 = \{\text{block, command, exp, comparand, term, factor}\}$
- and P_1 contains all of the following production rules³

$$\text{block} \rightarrow \epsilon$$

$$\text{block} \rightarrow \text{command block}$$

$$\text{command} \rightarrow i := \text{exp} \quad \text{for all } i \in \mathcal{I}$$

$$\text{command} \rightarrow \text{if exp then block else block end if}$$

$$\text{command} \rightarrow \text{while exp do block end while}$$

$$\text{exp} \rightarrow \text{comparand}$$

$$\text{exp} \rightarrow \text{comparand} < \text{comparand}$$

$$\text{comparand} \rightarrow \text{term}$$

$$\text{comparand} \rightarrow \text{term} + \text{comparand}$$

$$\text{comparand} \rightarrow \text{term} - \text{comparand}$$

$$\text{term} \rightarrow \text{factor}$$

$$\text{term} \rightarrow \text{factor} * \text{term}$$

$$\text{term} \rightarrow \text{factor} / \text{term}$$

$$\text{factor} \rightarrow n \quad \text{for all } n \in \mathcal{N}$$

$$\text{factor} \rightarrow i \quad \text{for all } i \in \mathcal{I}$$

$$\text{factor} \rightarrow (\text{exp})$$

An example of a string in this language is

while $i < n$ do $j := j + i$ $i := i + 1$ end while

³ Recall that ϵ means an empty sequence.

We can show that

while $i < n$ **do** $j := j + i$ $i := i + 1$ **end while**

is in the language by showing a derivation.

block

\Rightarrow command block

\Rightarrow **while** exp **do** block **end while** block

⋮

\Rightarrow **while** $i < n$ **do** $j := j + i$ command block **end while** block

\Rightarrow **while** $i < n$ **do** $j := j + i$ $i :=$ exp block **end while** block

\Rightarrow **while** $i < n$ **do** $j := j + i$ $i :=$ comparand block \dots

\Rightarrow **while** $i < n$ **do** $j := j + i$ $i :=$ term + comparand block \dots

\Rightarrow **while** $i < n$ **do** $j := j + i$ $i :=$ factor + comparand block \dots

\Rightarrow **while** $i < n$ **do** $j := j + i$ $i := i +$ comparand block \dots

\Rightarrow **while** $i < n$ **do** $j := j + i$ $i := i +$ term block \dots

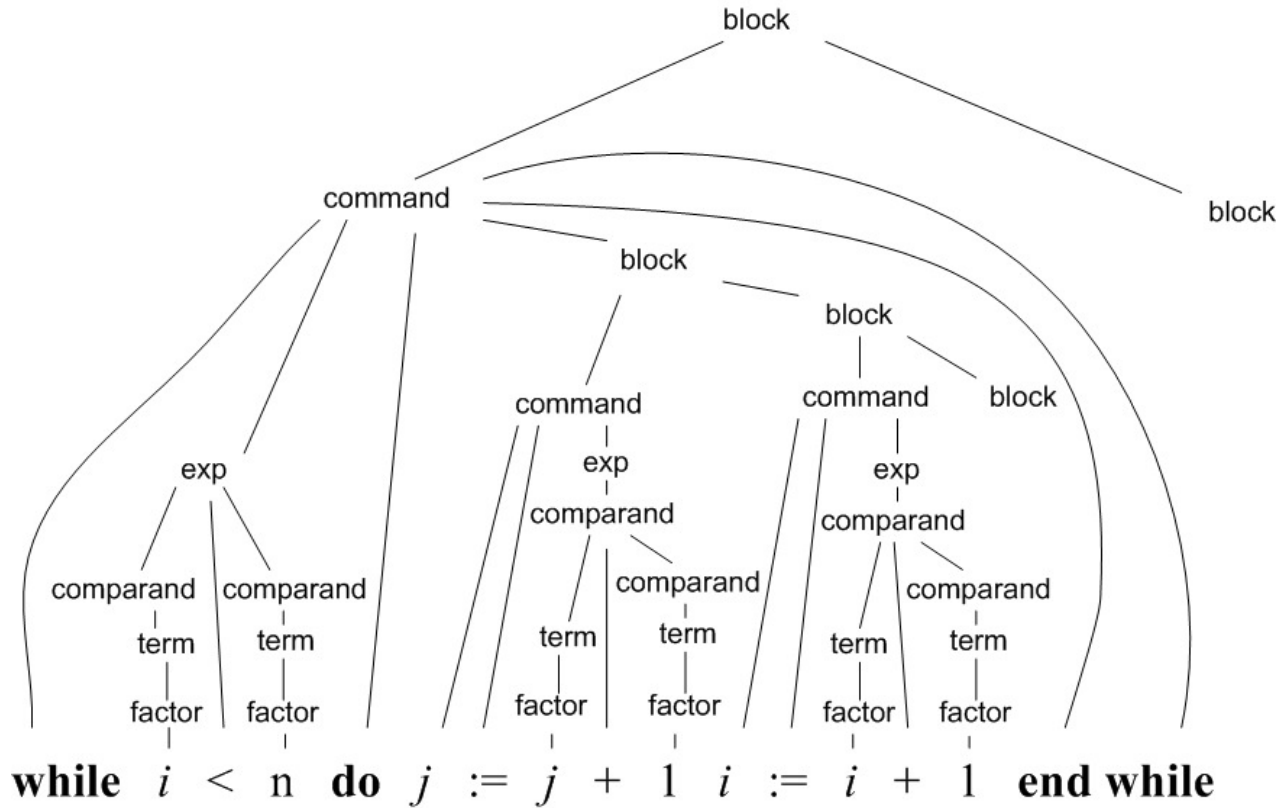
\Rightarrow **while** $i < n$ **do** $j := j + i$ $i := i +$ factor block \dots

\Rightarrow **while** $i < n$ **do** $j := j + i$ $i := i + 1$ block **end while** block

\Rightarrow **while** $i < n$ **do** $j := j + i$ $i := i + 1$ **end while** block

\Rightarrow **while** $i < n$ **do** $j := j + i$ $i := i + 1$ **end while**

Another way to show that a sequence is in a context free language is with a **parse tree**.



Let's define parse trees.

An ordered tree is a directed tree whose nodes are either leaves (with no children) or branches whose (0 or more) children are arranged in a sequence $children(t)$

A parse tree for a CFG (A, N, P, n_{start}) is a finite ordered tree whose nodes are labelled with symbols from $A \cup N$, such that

- the root is labelled with n_{start} and
- each branch node is labelled with a nonterminal symbol $n \in N$ and has a sequence of children whose labels form a finite sequence α such that $(n \rightarrow \alpha) \in P$.

Given a parse tree t the sequence of leaves is given by

```
proc fringe(  $t$  )
```

```
  if  $t$  is a leaf then return [ $label(t)$ ]
```

```
  else //  $t$  is a branch
```

```
    var  $\alpha := []$ 
```

```
    for  $u \leftarrow children(t)$  do  $\alpha := \alpha \hat{\ } fringe(u)$  end for
```

```
    return  $\alpha$ 
```

```
  end if
```

```
end fringe
```

Exercise: Show that, for any parse tree t , if $fringe(t) = \alpha$, then $n_{start} \xRightarrow{*} \alpha$.

Exercise: Show that if $n_{start} \xRightarrow{*} \alpha$, there is a parse tree t such that $\alpha = fringe(t)$.

Thus, for any $s \in A^*$, s is in the language of G if and only if there is a parse tree t where $s = fringe(t)$.

Exercise: Design a CFG for the language of palindromes in $\{0, 1\}^*$.

Exercise: Design a CFG for the language of valid boolean expressions in $\{p, q, r, \wedge, \vee, \Rightarrow, \neg, (,)\}^*$

Exercise: Design a CFG for the language of valid C++ variable declarations in $\{\text{int}, *, [,], (,), ;, a, b, c, 0, 1\}^*$ where a , b , and c are identifiers. E.g.

```
int a(int, int*);
```

declares a to be a function, while

```
int (*b[10])();
```

declares b to be an array of 10 pointers to functions returning int results.

Exercise: Look up the definition of regular language.

- Show that every regular language is a context-free language.
- Show that some context free languages are not regular languages.