**Definition:** [Grammar] A *grammar* is a triple $(\Gamma, \rightarrow, \gamma)$ such that,

- $\Gamma = T \cup N$ with $T \cap N = \emptyset$, where $T$ is a set of symbols called the *terminals* and $N$ is a set of symbols called the *non-terminals*,[1]

- $\rightarrow$ is a set of rules of the form $u \rightarrow v$ with $u, v \in \Gamma^*$,

- $\gamma$ is called the *start symbol* and $\gamma \in N$.

---

[1]The fact that $T$ and $N$ are non-overlapping means that there will never be confusion between terminals and non-terminals.

## Grammars

**Example:** Grammar for arithmetic expressions. We define the grammar $(\Gamma, \rightarrow, \gamma)$ as follows:

- $\Gamma = T \cup N$ with $T = \{a, b, c, +, *, (,)\}$ and $N = \{E\}$,
- $\rightarrow$ is is defined as,

$$
\begin{array}{rcl}
E & \rightarrow & E + E \\
E & \rightarrow & E * E \\
E & \rightarrow & (E) \\
E & \rightarrow & a \\
E & \rightarrow & b \\
E & \rightarrow & c
\end{array}
$$

- $\gamma = E$ (clearly this satisfies $\gamma \in N$).

With grammars, derivations always start with the start symbol. Consider,

$$E \Rightarrow E*E \Rightarrow (E)*E \Rightarrow (E+E)*E \Rightarrow (a+E)*E \Rightarrow (a+b)*E \Rightarrow (a+b)*c.$$

Here, $(a + b) * c$ is a normal form often also called a *terminal* or *derived string*.

**Exercise:** Identify the rule that was applied at each rewrite step in the above derivation.

**Exercise:** Derive the string $((a))$.

**Exercise:** Derive the string $a + b * c$. Is the derivation unique? Why? Why not?

We are now in the position to define exactly what we mean by a *language*.

**Definition:**[Language] Let $G = (\Gamma, \rightarrow, \gamma)$ be a grammar, then we define the *language of grammar G* as the set of all terminal strings that can be derived from the start symbol $s$ by rewriting using the rules in $\rightarrow$. Formally,
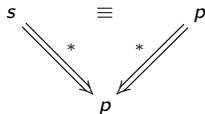
$$L(G) = \{q \mid \gamma \Rightarrow^* q \wedge q \in T^*\}.$$

**Example:** Let $J = (\Gamma, \rightarrow, \gamma)$ be the grammar of Java, then $L(J)$ is the set of all possible Java programs. The Java language is defined as the set of all possible Java programs.

# Grammars

**Observations:**

- With the concept of a language we can now ask interesting questions. For example, given a grammar $G$ and some sentence $p \in T^*$, does $p$ belong to $L(G)$?

- If we let $J$ be the grammar of Java, then asking whether some string $p \in T^*$ is in $L(J)$ is equivalent to asking whether $p$ is a *syntactically correct program*.

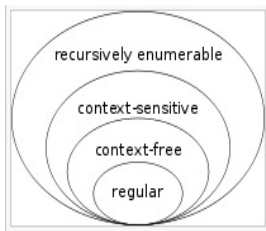- We can prove language membership by showing that the start symbol is equivalent to the sentence in question,

$$s \quad\searrow\quad \equiv \quad\swarrow\quad p$$
$$\searrow^* \qquad \swarrow^*$$
$$p$$

**Observations:**

- By restricting the shape of the rewrite rules in a grammar we obtain different language *classes*.
- The most famous set of language classes is the *Chomsky Hierarchy*.

# Grammars

### Table: The Chomsky Hierarchy

| Rules | Grammar | Language | Machine |
|---|---|---|---|
| $\alpha \to \beta$ | Type-0 | Recursively Enumerable | Turing machine |
| $\alpha A \beta \to \alpha \gamma \beta$ | Type-1 | Context-sensitive | Linear-bounded Turing machine |
| $A \to \gamma$ | Type-2 | Context-free | Pushdown automaton |
| $A \to a$ and $A \to aB$ | Type-3 | Regular | Finite state automaton |

where $\alpha, \beta, \gamma \in \Gamma^*, A, B \in N, a \in T$. In Type-1 $\gamma$ is not allowed to be the empty string.

**Observation:** The most convenient language class for programming language specification are the context-free languages – they are decidable – pushdown automata can be efficiently implemented in order to prove language membership.

# Grammars

**Example:** A simple imperative language. We define grammar $G = (\Gamma, \rightarrow, \gamma)$ as follows:

- $\Gamma = T \cup N$ where

  $T = \{0, \ldots, 9, a, \ldots, z, \textbf{true}, \textbf{false}, \textbf{skip}, \textbf{if}, \textbf{then}, \textbf{else}, \textbf{while}, \textbf{do}, \textbf{end} +, -, *, =, \leq, !, \&\&, ||, :=, ;, (,)\}$

  and

  $$N = \{A, B, C, D, L, V\}.$$

- The rule set $\rightarrow$ is defined by,

  $$
  \begin{array}{rcl}
  A & \rightarrow & D \mid V \mid A + A \mid A - A \mid A * A \mid (A) \\
  B & \rightarrow & \textbf{true} \mid \textbf{false} \mid A = A \mid A \leq A \mid !B \mid B\&\&B \mid B||B \mid (B) \\
  C & \rightarrow & \textbf{skip} \mid V := A \mid C \,;\, C \mid \textbf{if } B \textbf{ then } C \textbf{ else } C \textbf{ end} \mid \textbf{while } B \textbf{ do } C \textbf{ end} \\
  D & \rightarrow & L \mid -L \\
  L & \rightarrow & 0\,L \mid \ldots \mid 9\,L \mid 0 \mid \ldots \mid 9 \\
  V & \rightarrow & a\,V \mid \ldots \mid z\,V \mid a \mid \ldots z
  \end{array}
  $$

- $\gamma = C$.

Observe that this is a context-free grammar!

Here are some elements in $L(G)$,

$x := 1; y := x$
$v := 1;$ **if** $v \leq 0$ **then** $v := (-1) * v$ **else skip end**
$n := 5; f := 1;$ **while** $2 \leq n$ **do** $f := n * f; n := n - 1$ **end**

**Exercise:** Prove that they belong to $L(G)$.

HW#1 – see website