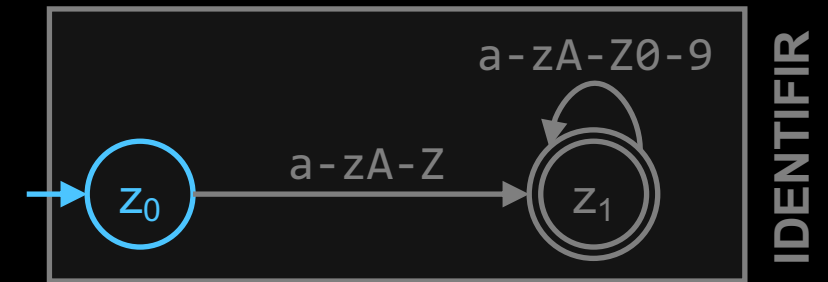


## Brainstorming: Beispielhafte Grammatik für Java

- 1) `<program>` ::= `<modifier>` `<type>` `<identifier>` `"{"` ... `"}"`
- 2) `<modifier>` ::= `"public"`
- 3) `<modifier>` ::= `"private"`
- 4) `<modifier>` ::= `"protected"`
- 5) `<type>` ::= `"class"`
- 6) `<type>` ::= `"interface"`
- 7) `<type>` ::= `"enum"`
- 8) `<identifier>` ::= `<letter>`
- 9) `<identifier>` ::= `<letter>` `<suffix>`
- 10) `<suffix>` ::= `<letter>`
- 11) `<suffix>` ::= `<digit>`
- 12) `<letter>` ::= `"a"` | `"b"` | `"c"` | ... | `"x"`
- 13) `<digit>` ::= `"0"` | `"1"` | ... | `"8"` | `"9"`

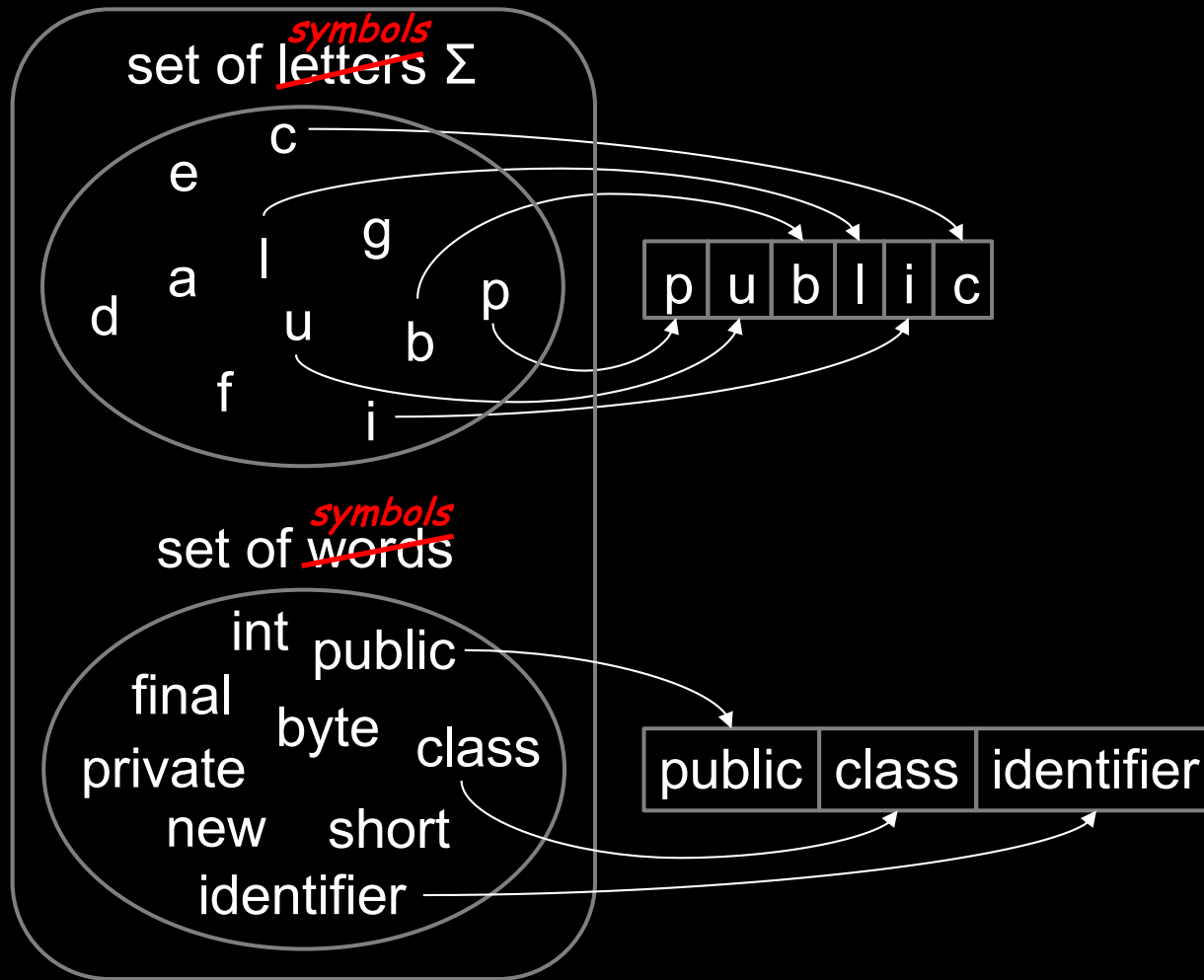
## Brainstorming: Beispielhafte Grammatik für Java

- 1) `<program>` ::= `<modifier>` `<type>` `<identifier>` "{" ... "}"
- 2) `<modifier>` ::= "public"
- 3) `<modifier>` ::= "private"
- 4) `<modifier>` ::= "protected"
- 5) `<type>` ::= "class"
- 6) `<type>` ::= "interface"
- 7) `<type>` ::= "enum"
- 8) `<identifier>` ::= `<letter>`
- 9) `<identifier>` ::= `<letter>` `<suffix>`
- 10) `<suffix>` ::= `<letter>`
- 11) `<suffix>` ::= `<digit>`
- 12) `<letter>` ::= "a" | "b" | "c" | ... | "x"
- 13) `<digit>` ::= "0" | "1" | ... | "8" | "9"



Hatten wir für  
Bezeichner nicht schon  
einen endlichen  
Automaten definiert?

Feststellung: ohne Flexion sind Wörter und Sätze methodisch gleichartig



word  $\in \Sigma^*$

~~language~~ <sup>language</sup> ~~vocabulary~~  $\subseteq \Sigma^*$

~~word~~ sentence  $\in \Sigma^*$

language  $\subseteq \Sigma^*$

## Offene Fragen

- ~~1. Nicht jedes Token darf an jeder Stelle stehen. Wie lässt sich das regulieren?~~
2. Sind die von endlichen Automaten akzeptierten Sprachen identisch zu den von Grammatiken erzeugten Sprachen?
- ~~3. Kann man die Erstellung eines endlichen Automaten automatisieren?~~
- ~~4. Lassen sich alle NEAs in DEAs umwandeln?~~

Zeit das gelernte festzuhalten...

## concatenation

Let  $A$  be a non-empty, finite set.

The set of all strings on  $A$  is the set  $A^* = \{\langle a_1, \dots, a_n \rangle : n \geq 0 \wedge a_1, \dots, a_n \in A\}$

The operation  $*$  is called the Kleene closure.

Concatenation  $\circ$  is a two-place operation on  $A^*$  defined as  $\langle a_1, \dots, a_n \rangle \circ \langle b_1, \dots, b_m \rangle = \langle a_1, \dots, a_n, b_1, \dots, b_m \rangle$

$A^*$  also contains the empty string  $\langle \rangle$ . It is denoted as  $\varepsilon$ . In particular,  $w \circ \varepsilon = w = \varepsilon \circ w$  is valid.

Concatenation is associate, but not commutative.

$$\begin{aligned}(\alpha \circ \beta) \circ \gamma &= \alpha \circ (\beta \circ \gamma) \\ \alpha \circ \beta &\neq \beta \circ \alpha\end{aligned}$$

## substrings &amp; atoms

If  $\alpha$  is  $a_1 \dots \mathbf{b_1} \dots \mathbf{b_n} \dots a_m$ , we call  $b_1 \dots b_n$  a **substring** of  $\alpha$ .

A **prefix** of  $\alpha$  is an initial substring of  $\alpha$ , a **suffix** of  $\alpha$  is a final substring of  $\alpha$ .

$\varepsilon$  is a substring of every string.

$\alpha$  is an **atom** in  $A^*$  if  $\alpha \neq \varepsilon$  and the **only substrings of  $\alpha$  in  $A^*$  are  $\varepsilon$  and  $\alpha$  itself**.

## alphabet

An alphabet is a **non-empty, finite set**  $A$  such that every element  $a \in A$  is an **atom in  $A^*$** .

We call the elements of alphabet  $A$  **symbols** or **lexical items**.

An alphabet is **often denoted by  $\Sigma$** , yet can be given any name.

$B = \{0,1\}$                 Says  $B$  is an alphabet of two symbols, 0 and 1.

$C = \{a,b,c\}$             Says  $C$  is an alphabet of three symbols, a, b and c.

Sometimes **space and comma** are in an alphabet while other times they are **meta symbols** used for descriptions.



## symbol / lexical items

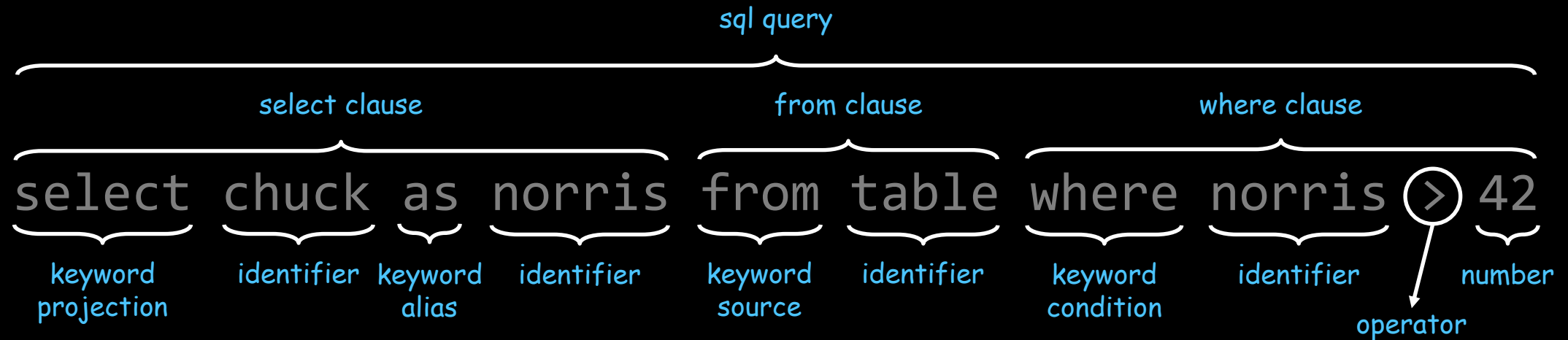
A symbol is a character, glyph, mark. It is an **abstract entity** that has **no meaning by itself**, often called uninterpreted. Letters from various alphabets, digits and special characters are the most commonly used symbols.

## formal language

A language  $L$  in alphabet  $\Sigma$  is a set of strings in  $\Sigma^*$ . We write  $L \subseteq \Sigma^*$ .

The expression 'aaa'  $\longrightarrow$  word  $w \in L$   
 belongs to the  
 language of all expressions that  
 consist only of the letter 'a'.  
 } language  $L = \{a, aa, aaa, aaaa, \dots\} = \Sigma^*$   
 alphabet  
 $\Sigma = \{a\}$   
 dazu gehört auch Epsilon!

# Aufbau formaler Sprachen



**Beobachtung:** Bei natürlichen Sprachen sind Satzzeichen üblicherweise keine Elemente des Alphabets, obwohl sie in fast jedem korrekten Satz auftreten. Das Alphabet beinhaltet eben nur jene Buchstaben, die benötigt werden, um alle Wörter der Sprache zu bilden.

Formale Sprachen verwenden dagegen häufig einige Sonderzeichen als Abkürzung. Anstelle des größer-Zeichens im oben abgebildeten Beispiel hätte man sich auch auf GT oder „greater than“ festlegen können. Das größer-Zeichen kann demnach ebenfalls als Wort angesehen werden und ist damit Teil des Alphabets.

## formal grammar

A formal grammar is specified by a tuple  $(V_N, V_T, S, P)$  where

$V_N$  is a finite set of non-terminal symbols,

$V_T$  is a finite set of terminal symbols,

$S \in V_N$  is the start symbol,

$P \subset V^*NV^* \times V^*$  is a finite set of production rules,

We write rules  $(v, w) \in P$  as  $v \rightarrow w$ . A derivation step has the form  $\alpha v \beta \Rightarrow \alpha w \beta$  where  $\alpha, \beta \in V^*$ .

A derivation has the form  $w_1 \Rightarrow w_2 \Rightarrow \dots \Rightarrow w_n$  which we abbreviate as

$$\mathbf{w}_1 \stackrel{*}{\Rightarrow} \mathbf{w}_n \Leftrightarrow \begin{cases} \mathbf{w}_1 = \mathbf{w}_n \\ \exists \mathbf{w} \in V^*: \mathbf{w}_1 \stackrel{*}{\Rightarrow} \mathbf{w} \wedge \mathbf{w} \Rightarrow \mathbf{w}_n \end{cases}$$

The language generated by  $G$  is defined as

$$L(G) = \{w \in \Sigma^* \mid S \stackrel{*}{\Rightarrow} w\}.$$

## deterministic finite acceptor

A deterministic finite acceptor (DFA) is specified by a tuple  $(Q, \Sigma, \delta, q_0, F)$  where

$Q$	is a finite set of states,
$\Sigma$	is a finite alphabet,
$\delta: Q \times \Sigma \rightarrow Q$	is the transition function,
$q_0 \in Q$	is the initial state,
$F \subseteq Q$	is the set of final states.

The function  $f^*: (Q \times \Sigma^*) \rightarrow Q$  specifies the state that is reached last when the word  $w$  is entered.

$$\begin{aligned} f^*(q, \varepsilon) &= q \\ \forall w \in \Sigma^* \forall x \in \Sigma: f^*(q, wx) &= \delta(f^*(q, w), x) \end{aligned}$$

A deterministic finite acceptor  $D = (Q, \Sigma, \delta, q_0, F)$  recognizes a language  $L \subseteq \Sigma^*$  if

$$L(D) = \{w \in \Sigma^* \mid f^*(q_0, w) \in F\}.$$