

Lexical Analysis

Scanning of the characters and Checking of the language words and symbols used in the syntactic structures
(identifiers, numerals, keywords, separators, delimiters, ...)

The **Lexical Languages** are, in general, quite simple.

- They are:

Regular Languages: \mathfrak{R}

- They can be defined using:

Regular Expressions

Regular (or Linear) Grammars

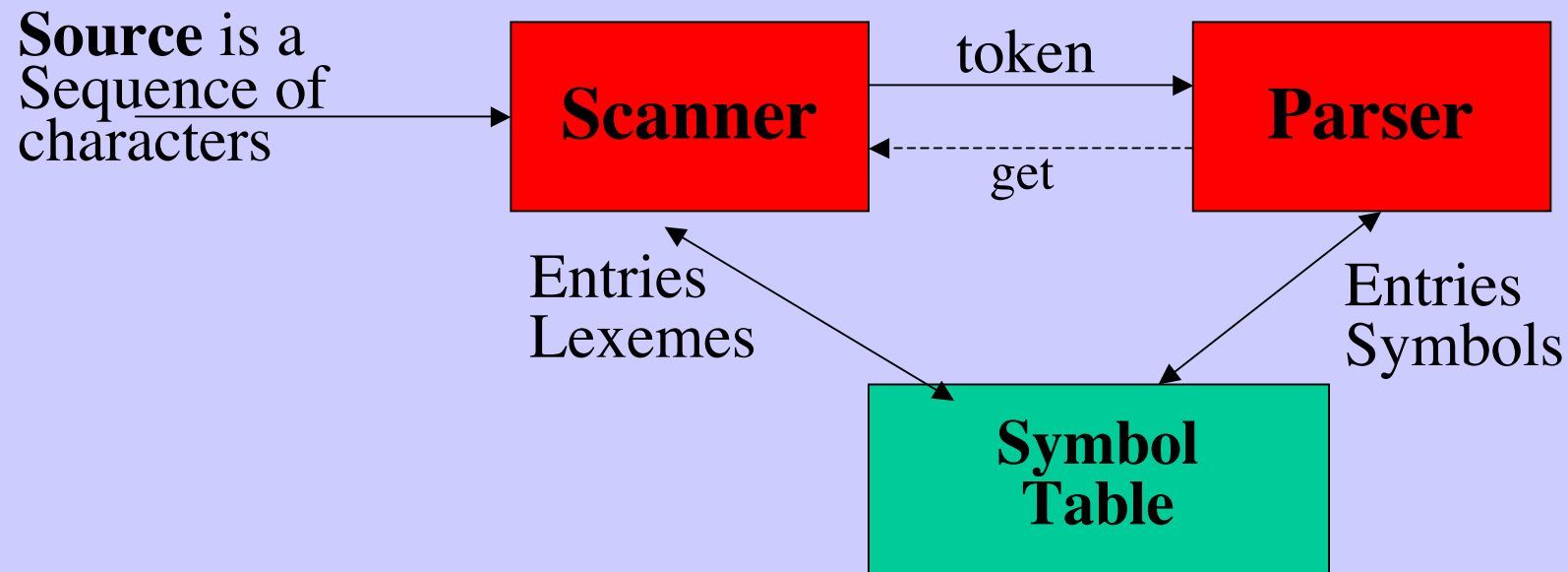
Transition Diagrams

- They can be analysed using:

N/D Automata

Where and How in a Compiler:

One Pass Structure - First Phase



Lexical Analysis (**Scanner**) is driven from Syntactic Analysis (**Parser**) which is asking for next token

Definitions: Token, Lexeme, Pattern

Token = a Lexical *Category*

lexeme = a *Value* of a Category

pattern = Lexical *rules* defining Tokens and Lexemes

Exemple

source: **const** pigreco = 3.1416

scanning: **const** **id** **rel** **num**

Lexeme:

Where source symbols are stored ?

- Source Symbols are maintained in Information Columns
- Scanner does not generate token but: $\langle \text{token}, \text{attribute} \rangle$
 $\langle \text{const}, \rangle \langle \text{id}, k \rangle \langle \text{rel}, = \rangle \langle \text{num}, k+1 \rangle$

	lexeme	
k	pigreco	
k+1	3.1416	

Information Columns

SymbolTable

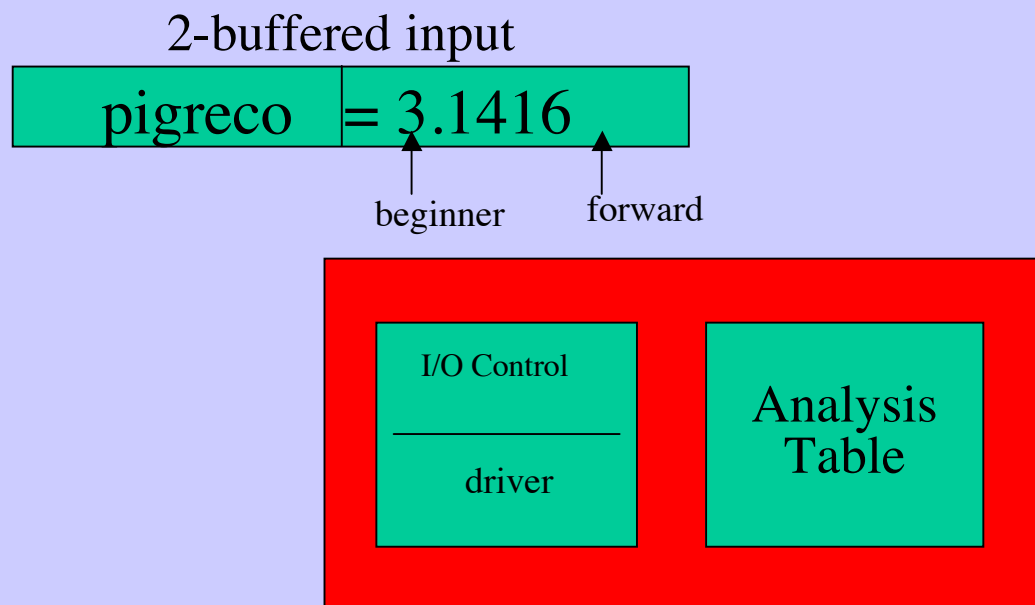
Patterns:

How Tokens are *unambiguously* defined ?

- **Patterns** may be defined using many different Formalisms:
 - Regular Expressions
 - Regular Grammars
 - Transition Diagrams
- Though equivalent, such Formalisms:
 - Have different expressiveness.
 - Hence the use of some is easier, clearer, and neater than that of the other.
 - Have different working frameworks.
 - Hence, the mechanization of the underlying analysis process is different for the different formalisms

Scanner Generators

- Today scanners are automatically generated using Regular Grammars, Regular Expressions, DFAutomata, in this order;
- Meaning preserving transformations:
 - allow to pass from one formulation to the next one, and
 - result into the structure below:



Regular Grammars

(but before) **Regular Expressions: E_Σ**

The Syntax

Regular Expressions E_Σ on an alphabet Σ are:

The minimum set recursively defined by:

1. $\epsilon \in E_\Sigma$
2. $a \in E_\Sigma, \quad \forall a \in \Sigma$
3. $e_1 \cdot e_2 \in E_\Sigma, \quad \forall e_1, e_2 \in E_\Sigma$
4. $e_1 \mid e_2 \in E_\Sigma, \quad \forall e_1, e_2 \in E_\Sigma$
5. $e^* \in E_\Sigma, \quad \forall e \in E_\Sigma$
6. $e^i \in E_\Sigma, \quad \forall e \in E_\Sigma$

Regular Expressions: The Meaning

The meaning of E_Σ is $\mathfrak{R} \subseteq 2^{\Sigma^*}$

Exponentiation

$$\Sigma = \{a,b\}, \quad \Sigma^i = \Sigma \times \Sigma^{i-1}, \quad \Sigma^0 = \{\lambda\}, \quad \Sigma^2 = \{aa, ab, bb, ba\}$$

$$\Sigma^* = \bigcup_{i \in \mathbb{N}} \Sigma^i = \{\lambda, a, b, aa, ab, bb, ba, aaa, \dots\}$$

Power Set

$$2^{\Sigma^*} = \{u \mid u \subseteq \Sigma^*\} = \{\{\lambda\}, \{a\}, \{b\}, \dots, \{\lambda, a\}, \dots, \{\lambda, a, b, ab\}, \dots\}$$

$[_]: E_\Sigma \rightarrow \mathfrak{R}$

1. $[\varepsilon] = \{\lambda\}$

2. $[a] = \{a\}$

3. $[e_1.e_2] = \{uv \mid u \in [e_1], v \in [e_2]\} = [e_1] \times [e_2]$

4. $[e_1|e_2] = \{u \mid u \in [e_1] \cup [e_2]\}$

5. $[e^*] = \{u \mid u \in \bigcup_{i \in \mathbb{N}} [e]^i\}$

6. $[e^i] = \{u \mid u \in [e]^i\}$ *shortening*

- singleton with nullary string

- singleton with 1 char string

- juxtaposition of the product set

- set union

- exponentiation set

- i^{th} -power

Regular Expressions: Examples

$$0 \mid 1 \mid \dots \mid 9$$
$$[0 \mid 1 \mid \dots \mid 9] = \{0, 1, \dots, 9\}$$
$$(0 \mid 1 \mid \dots \mid 9)^*$$
$$[(0 \mid 1 \mid \dots \mid 9)^*] = \{0, 1, \dots, 9, 00, 01, \dots, 3808, \dots\}$$
$$(1 \mid \dots \mid 9).(0 \mid 1 \mid \dots \mid 9)^*$$
$$[(1 \mid \dots \mid 9).(0 \mid 1 \mid \dots \mid 9)^*] = \{1, \dots, 9, 10, 11, \dots, 3808, \dots\}$$

Usually, dot notation in concatenation operator is omitted. Then:

$(1 \mid \dots \mid 9)(0 \mid 1 \mid \dots \mid 9)^*$ is written instead of

Grammars on Σ

$G = \langle V, \Sigma, s \in V, P \rangle$

- V *non-terminal set* (Tokens and auxiliarys categories)
- Σ *terminal set*
- s *initial symbol* (of V)
- $P \subseteq V \times E_{\Sigma \cup V}$ *production set*

Grammars

When it is a Regular Grammar

G is Regular, if in addition:

- V has a complete ordering $<$ such that:
- $\forall v ::= e \in P, \quad e \in E_{\Sigma \cup \{v' < v\}}$

Grammar

Example: Relations between grammars, expressions

The use of a regular grammar G in a lexic for numbers

```
num ::= simple | fract | exp
simple ::= digit digit*
fract ::= simple . simple
exp ::= fract E simple |
       fract E (+|-) simple
digit ::= 0 | 1 | ... | 9
```

How proving that G is regular?:

- Look for a complete ordering $<$ on V
- Prove that $<$ satisfies the second property

Only grammar productions are shown:

- What is Σ ?
- And, the other components of G ?
- Easily inferable: V (look left), s (look top)

How proving that G defines the language we are looking for?:

- Another Story Begins.

Grammars - Meaning $L(G)$: Equations on Languages

Meaning associates to:

- **Each** non-terminal, V , **One** language, $L(V) \in 2^{\Sigma^*}$.
- **Each** production, $v ::= e$, **One** equation, $L(v) = [e]$

see below on how do it

$$e \in E_{\Sigma} \implies [e] \in 2^{\Sigma^*}$$

$$v ::= e \implies L(v) = [e] \in (V \times 2^{\Sigma^*})$$

$$\{v_1 ::= e_1, \dots, v_k ::= e_k\} \implies$$

$$\{\underline{L}(v_1), \dots, \underline{L}(v_k) \mid \forall i, \underline{L}(v_i) \in 2^{\Sigma^*} \text{ and } \underline{L}(v_i) \equiv [[\underline{L}(v_j)/v_j \mid v_j < v_i]e_i]\}$$

$$G = \langle V, \Sigma, s \in V, P \rangle \implies \underline{L}(G) = \underline{L}(s)$$

Meaning of (Regular) Grammars:

An Example of Computation

$n ::= s \mid f \mid e$
 $s ::= d d^*$
 $f ::= s.s$
 $e ::= f E s \mid$
 $\quad f E (+|-) s$
 $d ::= 0 \mid 1 \mid \dots \mid 9$

$$L(0|1|\dots|9) = \{0, \dots, 9\} \in (V \times 2^{\Sigma^*})$$

$$\underline{L}(d) = \{0, \dots, 9\},$$

$$\underline{L}(s) = L([\{0, \dots, 9\}/d] d d^*)$$

$$= \{0, \dots, 9\} \{0, \dots, 9\}^* = \{0, \dots, 9\}^+$$

$$\underline{L}(e) = \dots$$

$$= \{u.vEw \mid u, v, w \in \{0, \dots, 9\}^+\}$$

$$L(n) = \{0, \dots, 9\}^+ \cup$$

$$\{u.v \mid u, v \in \{0, \dots, 9\}^+\} \cup$$

$$\{u.vEw \mid u, v, w \in \{0, \dots, 9\}^+\}$$

...

Regular Grammars: Meaning

Compare the formulation in the previous slide,

$$\{v_1 ::= e_1, \dots, v_k ::= e_k\} \implies$$

$$\{\underline{L}(v_1), \dots, \underline{L}(v_k) \mid \forall i, \underline{L}(v_i) \in 2^{\Sigma^*} \text{ and } \underline{L}(v_i) \equiv [[\underline{L}(v_j)/v_j \mid v_j < v_i]e_i]\}$$

with the the one below.

$$\{v_1 ::= e_1, \dots, v_k ::= e_k\} \implies$$

$$\{\underline{L}(v_1), \dots, \underline{L}(v_k) \mid \forall i, \underline{L}(v_i) \in 2^{\Sigma^*} \text{ and } \underline{L}(v_i) \equiv [[e_i^1, \dots, e_i^{n_i}/v_i^1, \dots, v_i^{n_i}]e_i]\}$$

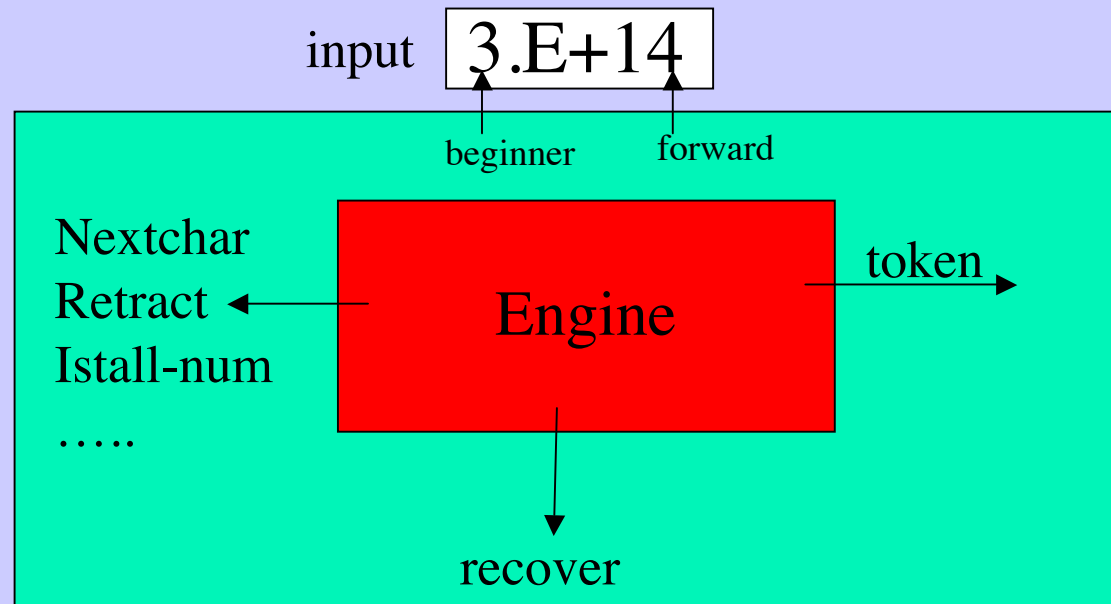
where: 1) $\{v_i^1, \dots, v_i^{n_i}\} = \{v_j \in V \mid v_j < v_i\}$;
 2) $v_i^1 < \dots < v_i^{n_i}$;
 3) $[e_i^1, \dots, e_i^{n_i}/v_i^1, \dots, v_i^{n_i}]e_i = [e_i^1, \dots, e_i^{n_i-1}/v_i^1, \dots, v_i^{n_i-1}][e_i^{n_i}/v_i^{n_i}]e_i$

Comment the differences.

Justify why the last one is more correct than the other.

How to Recognize the Lexic defined by a Grammar

```
n ::= s | f | e
s ::= d d*
f ::= s.s
e ::= f E s |
      f E (+|-) s
d ::= 0 | 1 | ... | 9
```



Scanners do it