

Compiladores

(IF688)

Leopoldo Teixeira
lmt@cin.ufpe.br | @leopoldomt

Estendendo para identificadores e números maiores que 9

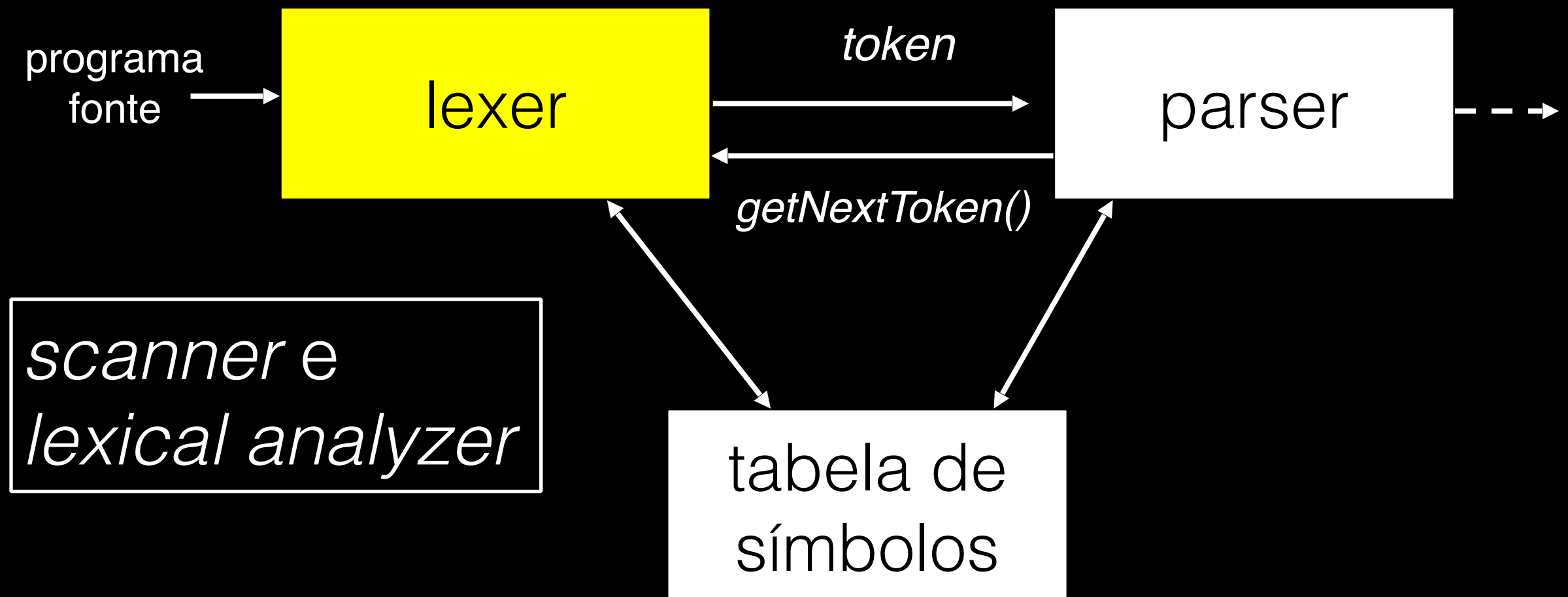
$expr \rightarrow$	$expr + term$	{print ('+') }
$expr \rightarrow$	$expr - term$	{print ('-') }
$expr \rightarrow$	$term$	
$term \rightarrow$	$term * factor$	{print ('*') }
$term \rightarrow$	$term / factor$	{print ('/') }
$term \rightarrow$	$factor$	
$factor \rightarrow$	$(expr)$	
$factor \rightarrow$	num	{print (num.value) }
$factor \rightarrow$	id	{print (id.lexeme) }

Reescrivendo...

expr	→	term rest
rest	→	+ term {print ('+') } rest
rest	→	- term {print ('-') } rest
rest	→	ϵ
term	→	factor resto
resto	→	* factor {print ('*') } resto
resto	→	/ factor {print ('/') } resto
resto	→	ϵ
factor	→	0 {print ('0') }
...		
factor	→	9 {print ('9') }

Antes de traduzir um programa, compiladores precisam entender a sua estrutura e significado.

Analizador léxico



Construindo um analisador léxico

- Especificar tokens e lexemas da linguagem
- Implementar (ou gerar através de ferramenta) o analisador a partir da especificação

Quais as razões para
separarmos análise
léxica de sintática?

Lexer vs. Parser

- Design mais simples – misturar análise léxica com sintática torna o parser bem mais complicado
- Eficiência – separação possibilita construir processadores léxicos e sintáticos mais eficientes.
- Portabilidade – variações de dispositivos podem ficar restritas ao analisador léxico.

Tokens, Padrões e Lexemas

- Token: par com o nome do token e atributos opcionais
- Padrão: descrição dos possíveis lexemas associados a um tipo de token
- Lexema: sequência de caracteres que casam com o padrão de um tipo de token.

Exemplo

- Para a entrada: “var1 = 23”;
- Primeiro **token** da entrada: <ID, “var1”>
- Para este token
 - **Padrão**: “seq. de caracteres seguida por dígito”
 - ID (para identificador) é o tipo do token
 - **Lexema** “var1”

Observações

- Existem conjuntos de strings na entrada que geram o mesmo tipo de token
 - Exemplo: var1, xyz2 são tokens do mesmo tipo ID
- Símbolos terminais de uma gramática correspondem a tokens
 - Exemplo: Palavras reservadas, operadores, identificadores, constantes, parenteses, etc.

Exemplos de Tokens

TOKEN	Descrição Informal	Ex.: Lexemas
<code>if</code>	caracteres <code>i, f</code>	<code>if</code>
<code>else</code>	caracteres <code>e, l, s, e</code>	<code>else</code>
<code>comparison</code>	<code><, >, <=, >=, ==, !=</code>	<code><=, !=</code>
<code>id</code>	letra seguida de letras e dígitos	<code>pi, score, D2</code>
<code>number</code>	qualquer constante numérica	<code>3.14159, 0, 6.02e23</code>
<code>literal</code>	qualquer coisa exceto <code>"</code> , envolvida por aspas	<code>"exemplo"</code>

Atributos de um token

- Ao reconhecer um token **num** é relevante saber se seu valor é zero ou um, por exemplo.
- Geralmente associado a um token existe um atributo.
- Uso da tabela de símbolos para guardar informações auxiliares sobre tokens.

Exemplo

$$E = M * C ** 2$$

Exemplo

E	=	M	*	C	**	2
---	---	---	---	---	----	---

Exemplo

E	=	M	*	C	**	2
----------	----------	----------	----------	----------	-----------	----------

```
<id, apontador p/posição de E na tabela de símbolos>  
<assign_op,>  
<id, apontador p/posição de M na tabela de símbolos>  
<mult_op,>  
<id, apontador p/posição de C na tabela de símbolos>  
<exp_op,>  
<num, valor 2>
```


Capturando erros com análise léxica

- É difícil para um analisador léxico determinar, sem a ajuda de outros componentes, que há um erro no código-fonte
- Por exemplo:
 - `fi (a == f(x)) ...`
- Qual o token que deve ser retornado?

Especificação de Tokens

- Baseada em Teoria de Linguagens:
 - Alfabetos, strings e linguagens
- Operações sobre linguagens
- Seja L o conjunto $\{A, B, \dots, Z, a, b, \dots, z\}$ e D o conjunto $\{0, 1, \dots, 9\}$.

Expressões Regulares

Expressões

- Na matemática, podemos usar os operadores para construir expressões a serem avaliadas, como:

$$(5+3)^*4$$

Estas expressões retornam um resultado numérico.

Operações Regulares

- Em linguagens formais podemos utilizar as operações regulares para criar expressões envolvendo linguagens:

$$(0+1)0^*$$

- Operações: união, concatenação, fecho reflexivo.
- São chamadas de expressões regulares.

Expressão Regular

- Formalismo denotacional
- Definida a partir de conjuntos (linguagens) básicos, concatenação e união.
- Adequadas para a comunicação
 - humano \times humano
 - humano \times máquina
- Importantes em diversas aplicações onde é necessário encontrar padrões em textos.

Exemplos

- $0^*10^* = ?$
- $(0+1)^*1(0+1)^* = ?$
- $(0+1)^*001(0+1)^* = ?$

Exemplos

- $0^*10^* =$
 $\{w \mid w \text{ contém um único } 1\}$
- $(0+1)^*1(0+1)^* =$
 $\{w \mid w \text{ tem pelo menos um } 1\}$
- $(0+1)^*001(0+1)^*$
 $\{w \mid w \text{ contém a sub-sentença } 001\}$

Mais Exemplos

ER	Linguagem Gerada
aa	
ba^*	
$(a+b)^*$	
$(a+b)^*aa(a+b)^*$	
$a^*ba^*ba^*$	
$(a+b)^*(aa+bb)$	
$(a+\varepsilon)(b+ba)^*$	

Mais Exemplos

ER	Linguagem Gerada
aa	somente a palavra aa
ba*	todas as palavras que iniciam por b, seguido por zero ou mais a
(a+b)*	todas as palavras sobre { a, b }
(a+b)*aa(a+b)*	todas as palavras contendo aa como subpalavra
a*ba*ba*	todas as palavras contendo exatamente dois b
(a+b)*(aa+bb)	todas as palavras que terminam com aa ou bb
(a+ε)(b+ba)*	todas as palavras que não possuem dois a consecutivos

Definições Regulares

- Por conveniência de notação, nomeamos algumas expressões regulares
- Estes nomes são usados como símbolos em expressões subsequentes

Identificadores em C

letter $\rightarrow A | B | \dots | Z | a | b | \dots | z$

digit $\rightarrow 0 | 1 | \dots | 9$

id $\rightarrow letter (letter | digit)^*$

Generalizando...

Em um alfabeto Σ

$$d_1 \rightarrow r_1$$

$$d_2 \rightarrow r_2$$

...

$$d_n \rightarrow r_n$$

onde:

1. cada d_i é um novo símbolo, que não está em Σ e é diferente dos demais d
2. cada r_i é uma expressão regular sobre $\Sigma \cup \{d_1, d_2, \dots, d_{i-1},\}$

Numerais

digit $\rightarrow 0 \mid 1 \mid \dots \mid 9$

digits $\rightarrow \textit{digit digit}^*$

optionalFraction $\rightarrow . \textit{digits} \mid \varepsilon$

optionalExponent $\rightarrow (E (+ \mid - \mid \varepsilon) \textit{digits}) \mid \varepsilon$

number $\rightarrow \textit{digits optionalFraction optionalExponent}$

Operador +

- Um outro operador regular comumente utilizado é o + que significa:

$$r^+ = rr^*$$

- ou seja, enquanto r^* representa toda sentença formada por 0 ou mais concatenações, r^+ representa sentenças formada por pelo menos 1 concatenação de r .

Operador ?

- Outro operador regular comumente utilizado é o ? que significa:

$$r? = (r + \epsilon)$$

- ou seja, a linguagem formada por 0 ou 1 concatenação sucessiva de r.

Classes de Caracteres

- Uma expressão regular $a_1|a_2|\dots|a_n$ pode ser substituída pela abreviação $[a_1a_2a_n]$
- Se a_1, a_2, \dots, a_n formam uma sequência, podemos substituir por a_1-a_n
 - **$[abc] = a|b|c$**
 - **$[a-z] = a|b|c|\dots|z$**

Reescrevendo...

letter \rightarrow [A-Za-z_]

digit \rightarrow [0-9]

id \rightarrow *letter* (*letter* | *digit*)*

Reescrevendo...

digit \rightarrow [0-9]

digits \rightarrow *digit*⁺

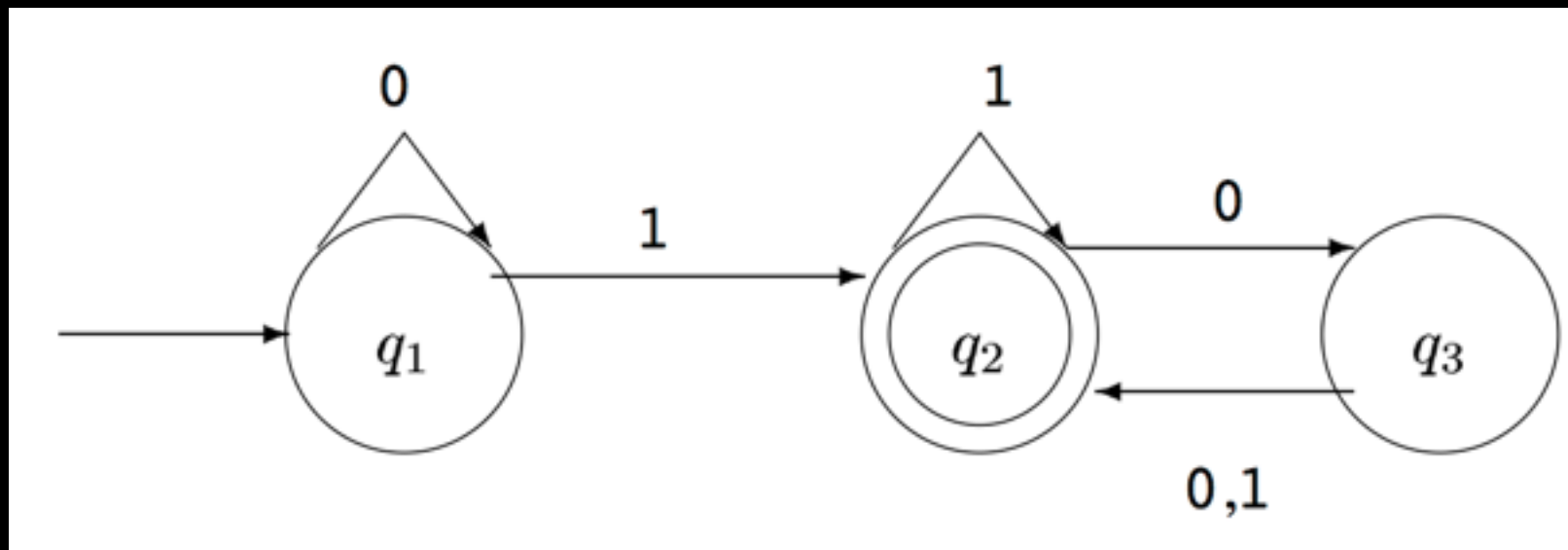
number \rightarrow *digits* (*.* *digits*)? (*E* [*+-*]? *digits*)?

Como reconhecer os
tokens?

Reconhecimento de Tokens

- Gerar diagramas de transição e depois implementar uma máquina de estados.
- Autômatos finitos determinísticos e não-determinísticos.

AFD



$M_1 = (\Sigma, Q, \delta, q_0, F)$, onde

$\Sigma = \{0, 1\}$

$Q = \{q_1, q_2, q_3\}$

$q_0 = q_1$

$F = \{q_3\}$

$\delta =$

	0	1
q ₁	q ₁	q ₂
q ₂	q ₃	q ₂
q ₃	q ₂	q ₂

AFD

$$M_1 = (\Sigma, Q, \delta_1, q_0, F)$$

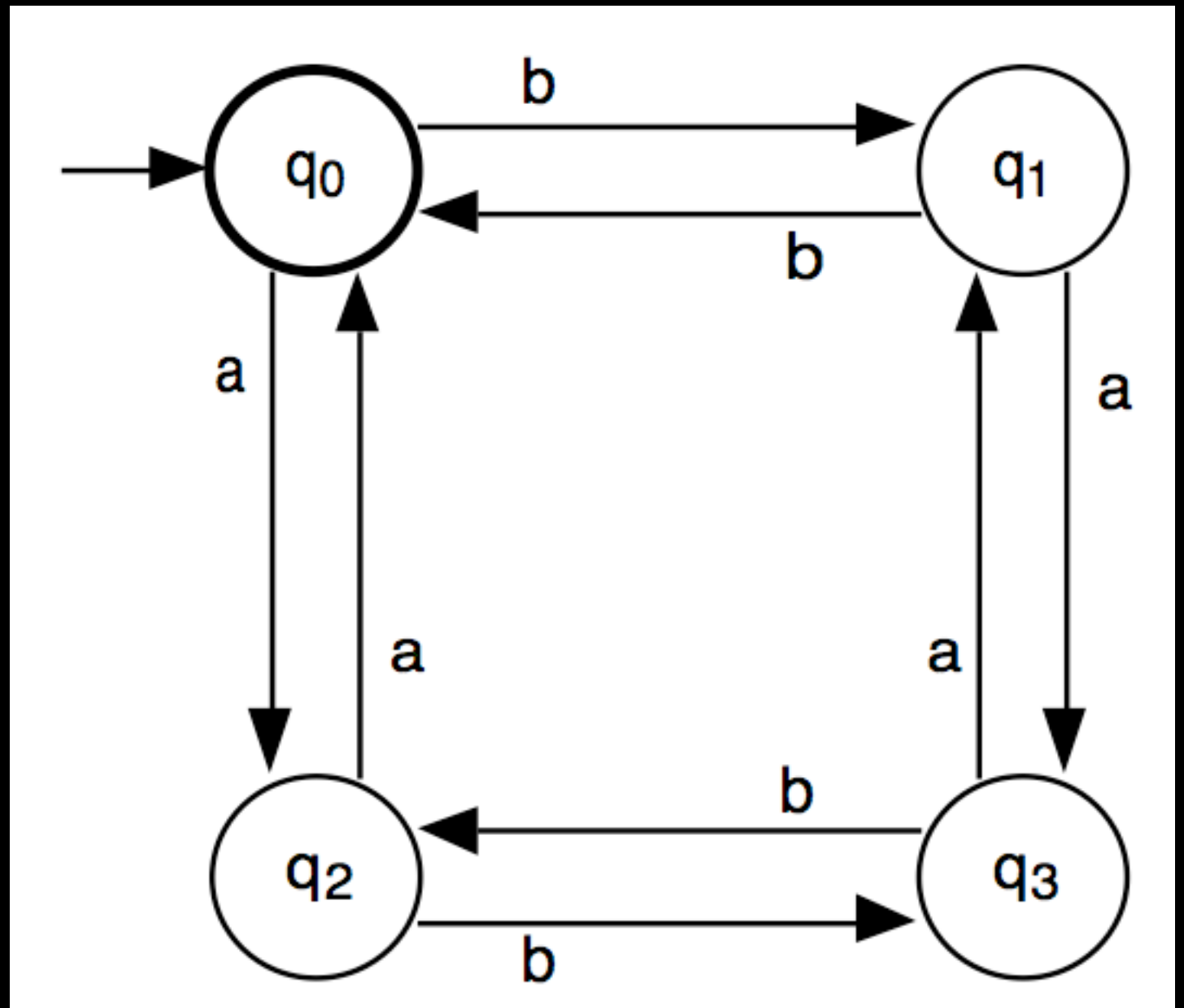
$$\Sigma = \{a, b\}$$

$$Q = \{q_0, q_1, q_2, q_3\}$$

$$F = \{q_0\}$$

$\delta_1 =$

	a	b
q ₀	q ₂	q ₁
q ₁	q ₃	q ₀
q ₂	q ₀	q ₃
q ₃	q ₁	q ₂

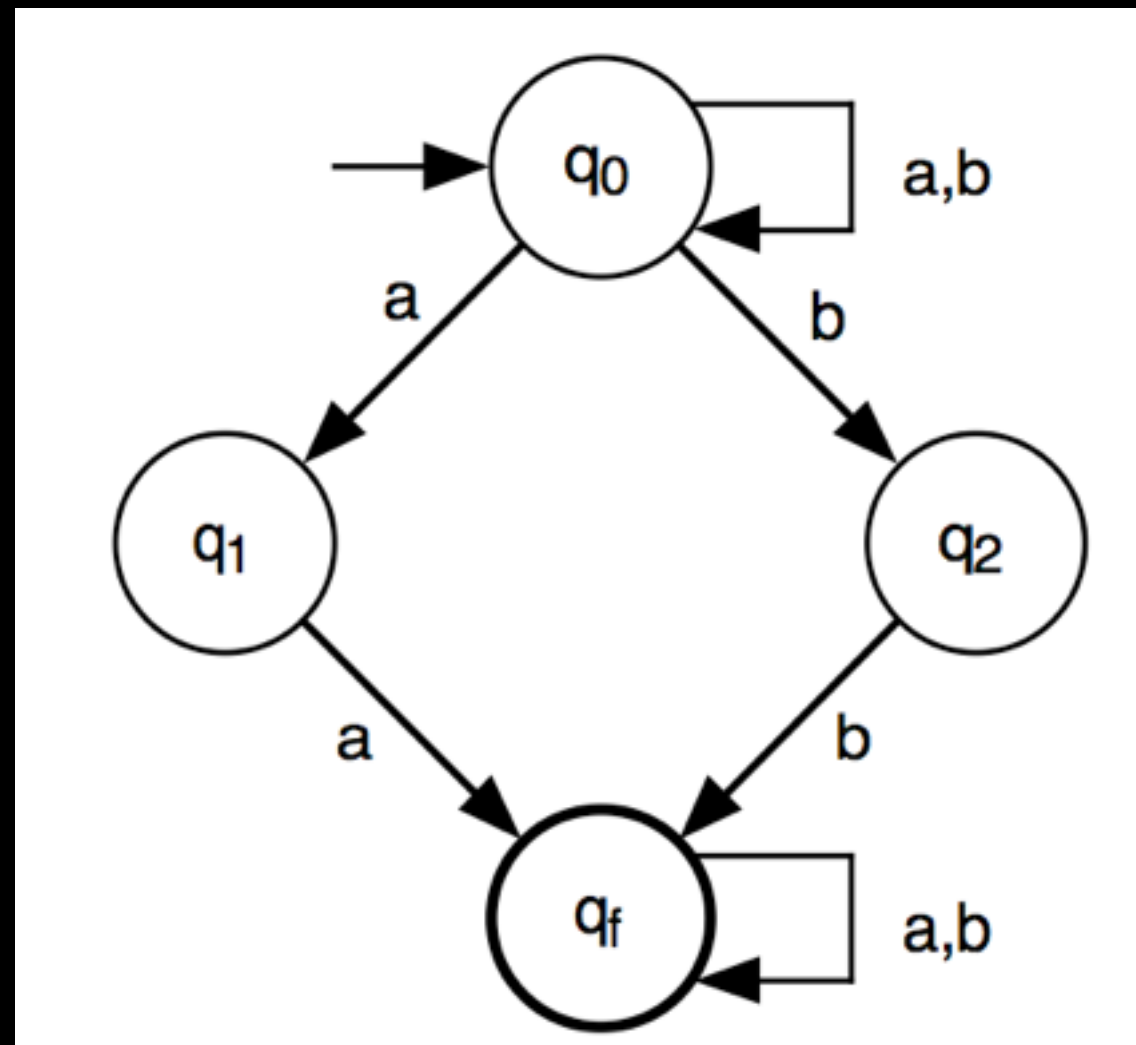


AFN

$$M_2 = (\{a,b\}, \{q_0, q_1, q_2, q_f\}, \bar{\delta}_2, q_0, \{q_f\})$$

$\bar{\delta}_2 =$

	a	b
q₀	$\{q_0, q_1\}$	$\{q_0, q_2\}$
q₁	$\{q_f\}$	-
q₂	-	$\{q_f\}$
q_f	$\{q_f\}$	$\{q_f\}$

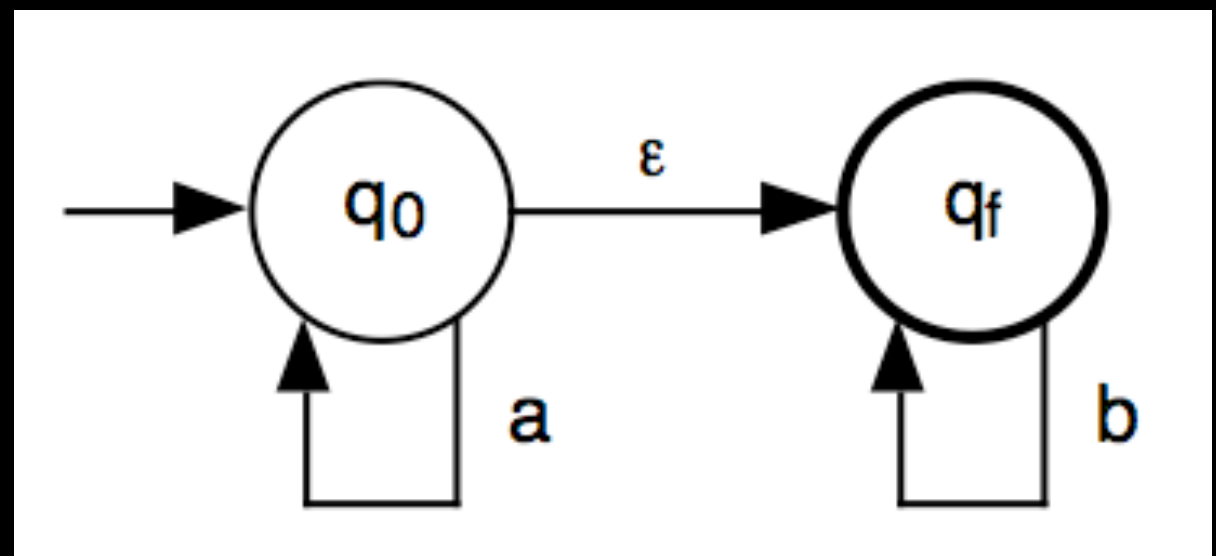


AFN ϵ

$$M_3 = (\{a,b\}, \{q_0, q_f\}, \delta_3, q_0, \{q_f\})$$

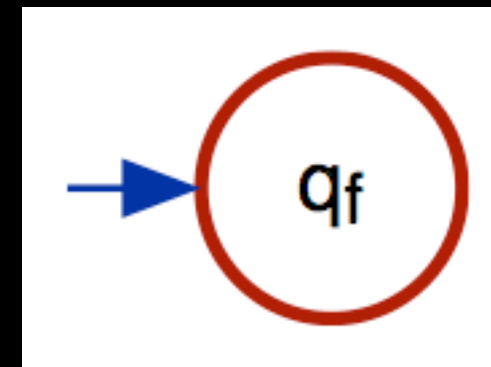
$$\delta_3 =$$

	a	b	ϵ
q₀	$\{q_0\}$	-	$\{q_f\}$
q_f	-	$\{q_f\}$	$\{q_f\}$



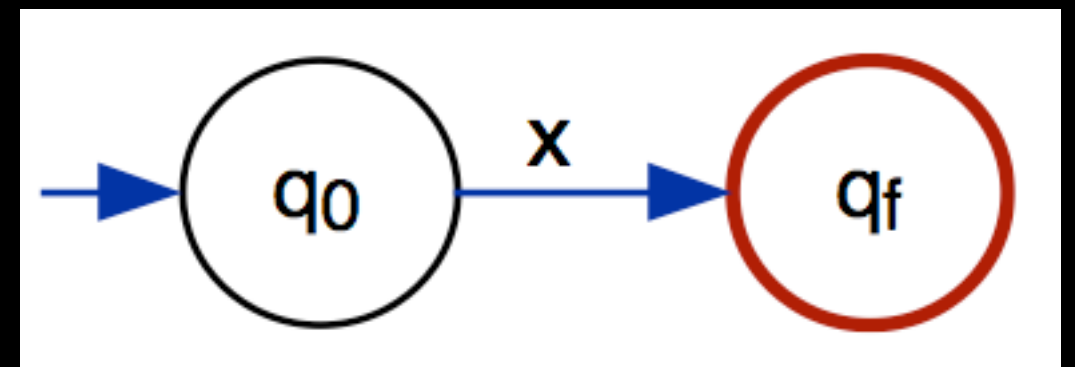
ERs como autômatos

$r = \varepsilon,$



$$M_1 = (\emptyset, \{q_f\}, \delta_1, q_f, \{q_f\})$$

$r = x,$



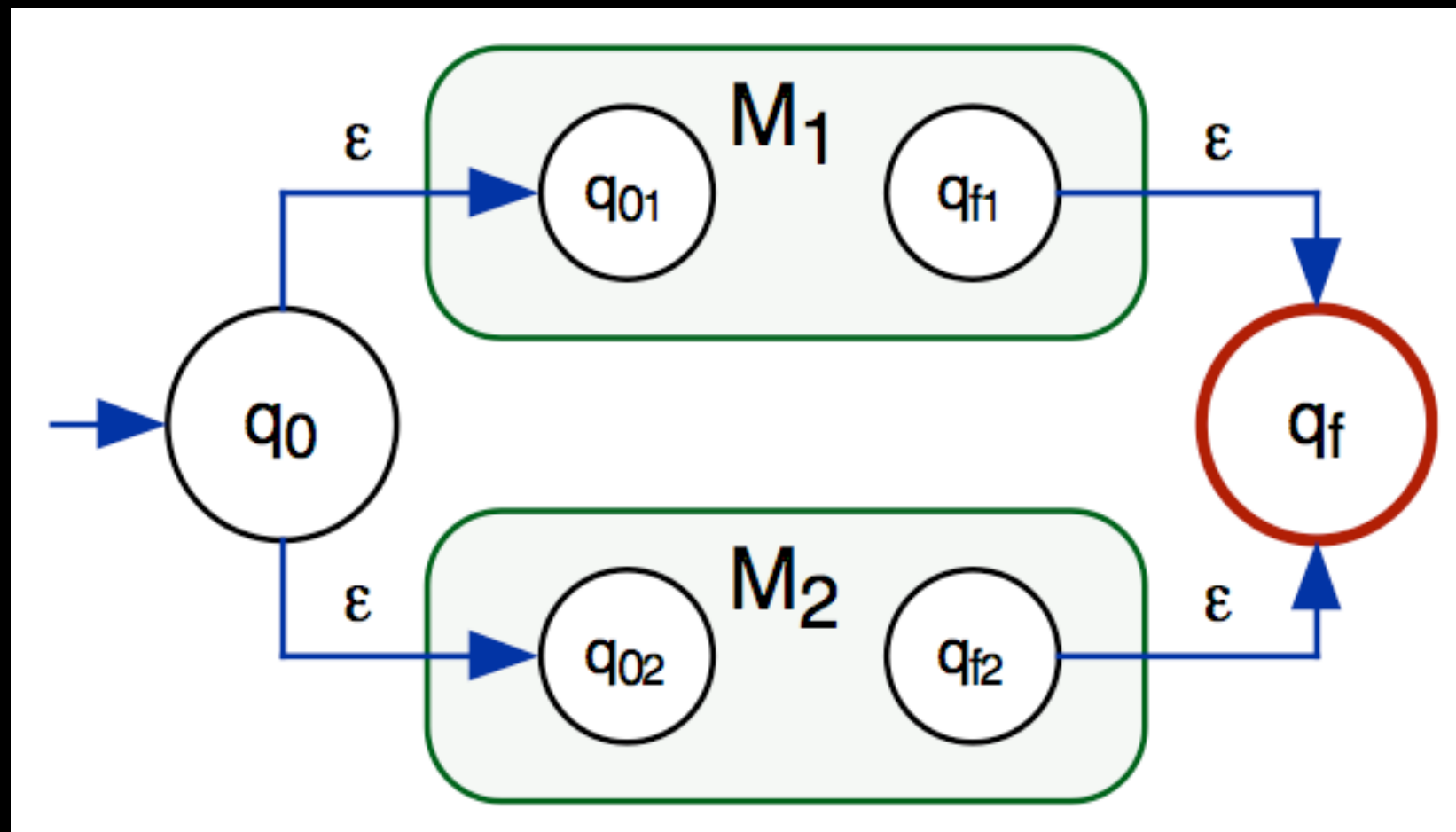
$$M_2 = (\{x\}, \{q_0, q_f\}, \delta_2, q_0, \{q_f\})$$

Passo de indução

- Se r é uma ER com $n+1$ operadores, pode ser representada por:
 - $r = r_1 + r_2$
 - $r = r_1 r_2$
 - $r = r_1^*$
- É possível representar cada uma destas com autômatos?

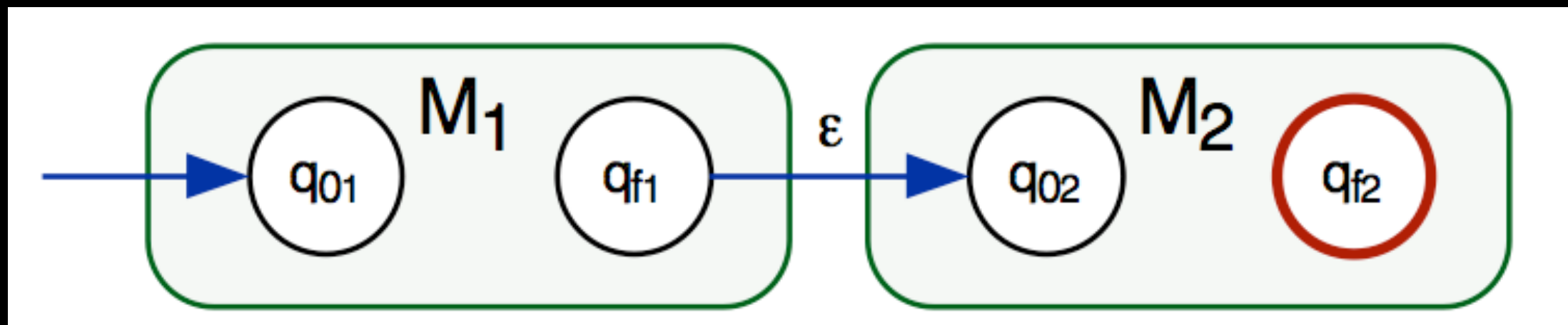
$$r = r_1 + r_2$$

$$M_3 = (\Sigma_1 \cup \Sigma_2, Q_1 \cup Q_2 \cup \{q_0, q_f\}, \delta, q_0, \{q_f\})$$



$$r = r_1 r_2$$

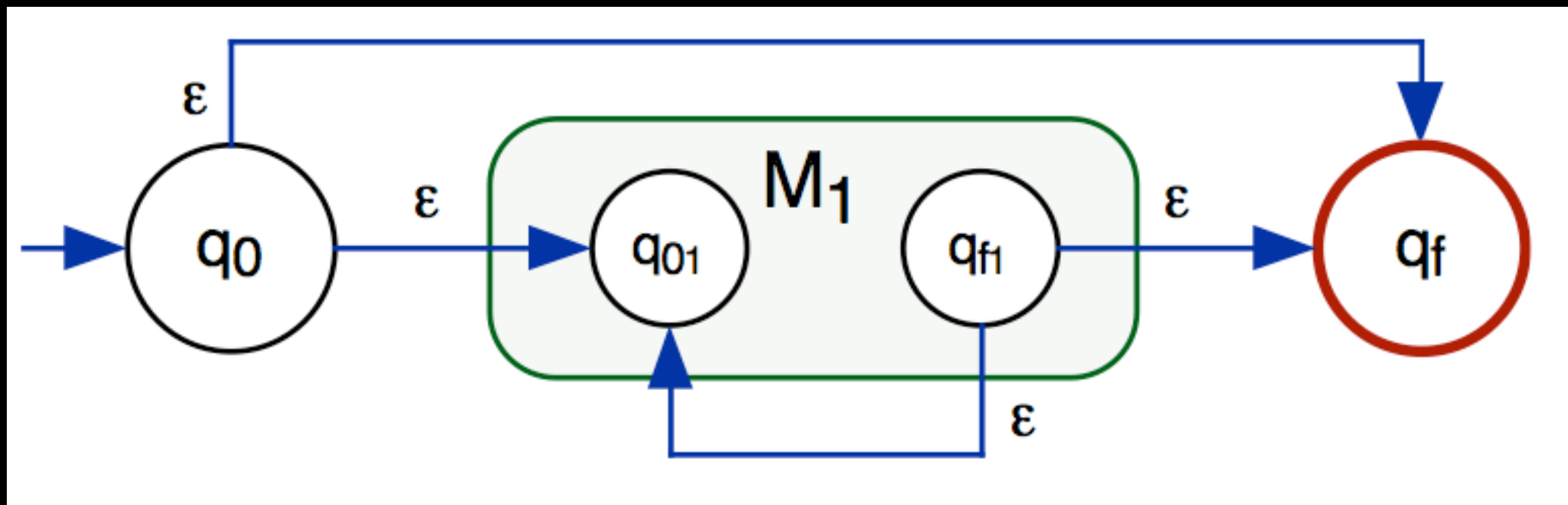
$$M_4 = (\Sigma_1 \cup \Sigma_2, Q_1 \cup Q_2 \cup \{q_0, q_f\}, \delta, q_0, \{q_f\})$$



$$r = r_1^*$$

(suponha $q_0 \notin Q_1$, $q_f \notin Q_1$)

$$M_5 = (\Sigma_1, Q_1 \cup \{q_0, q_f\}, \delta, q_0, \{q_f\})$$



Exemplo

if → if

id → [a-z][a-z0-9]*

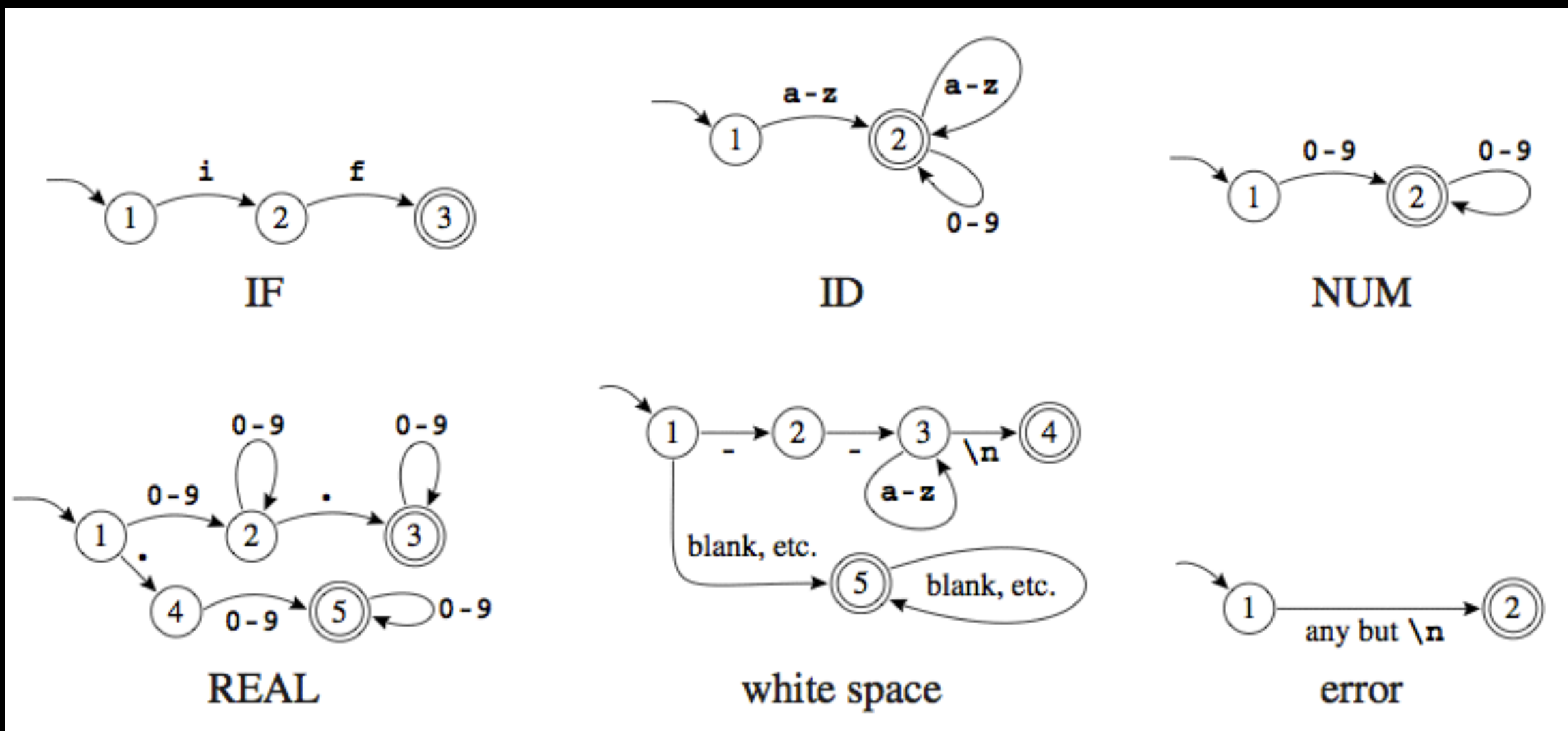
num → [0-9]+

real → ([0-9]+ "." [0-9]*) | ([0-9]* "." [0-9]+)

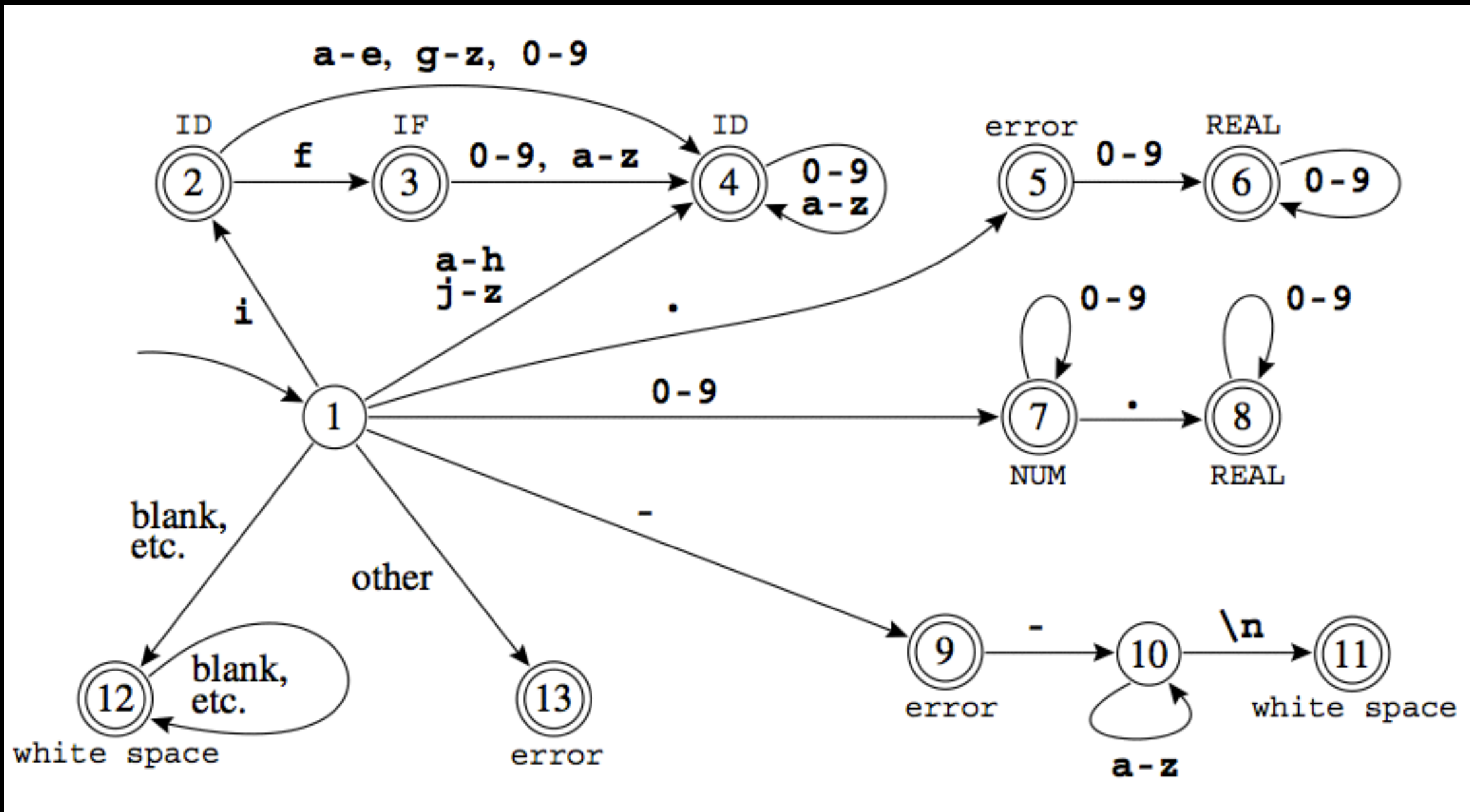
ws → ("--"[a-z]* "\n") | (" "[a-z]* "\n" | "\t")+

error → .

Autômatos Exemplo



Combinando os Autômatos



Próxima aula:
implementando
analisadores léxicos