



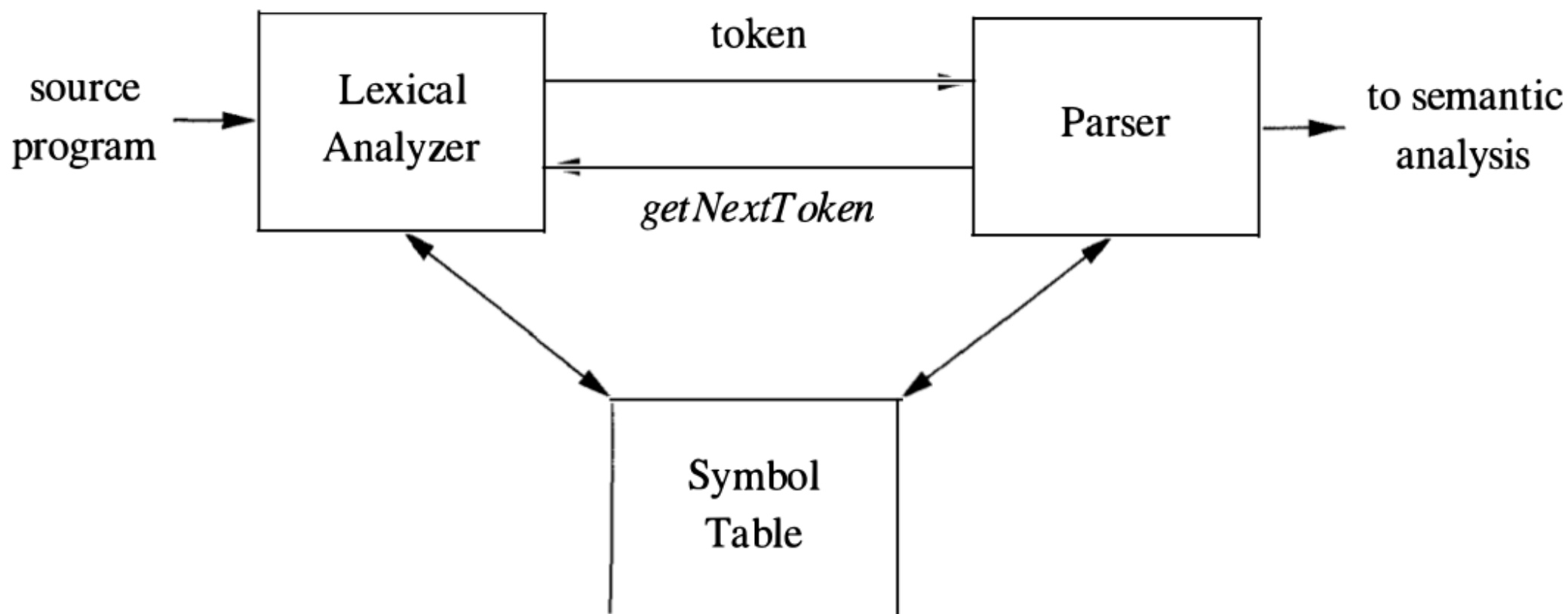
دانشگاه صنعتی شاهرود
دانشکده مهندسی کامپیوتر و فناوری اطلاعات

درس اصول طراحی کامپایلر

تحلیل لغوی Lexical Analysis

مدرس:
علیرضا تجری

موقعیت تحلیل لغوی در فرآیند کامپایل



عبارت‌های معنادار در زبان‌های برنامه‌نویسی

- int float if else while (1) کلمه‌های کلیدی
- ; , () { } [] (2) علائم
- > >= ! + - * ** / = (3) عملگرها
- variableName max ControllerClass (4) شناسه‌ها
- 15.25 "Hello World" 15 (5) ثوابت
- /*This is a comment*/ // # (6) توضیحات
- space tab \t new line \n carriage return \r (7) فضاهاى خالى

تحلیلگر معنایی به دنبال کدام یک از موارد بالا است؟

ورودی و خروجی تحلیلگر لغوی

```
int main(){  
    return 0;  
}
```

کد نوشته شده در فایل

```
int main(){\n\treturn 0;\n}
```

آنچه تحلیلگر لغوی می بیند

دنباله توکن هایی که برای پارسر ارسال می شود

```
<INT> <ID,1> <OPP> <CLP> <OPB> <RETURN>  
<NUMBER,0> <SC> <CLB>
```

توکن، الگو و lexeme

- توکن (نشانه، Token): رکوردی شامل یک نام و یک خصوصیت اختیاری

<name, attribute>

- الگو (Pattern): توصیفی از دنباله کاراکترهایی که معادل یک توکن است.

- lexeme: رشته کاراکتری در برنامه ورودی که دارای یک الگو است و توسط تحلیلگر لغوی به عنوان یک توکن شناخته می‌شود.

```
int main(
```

```
<INT> <ID,1> <OPP>
```

مثال توکن، الگو و lexeme

TOKEN	INFORMAL DESCRIPTION	SAMPLE LEXEMES
if	characters i, f	if
else	characters e, l, s, e	else
comparison	< or > or <= or >= or == or !=	<=, !=
id	letter followed by letters and digits	pi, score, D2
number	any numeric constant	3.14159, 0, 6.02e23
literal	anything but ", surrounded by "'s	"core dumped"

انواع توکن‌ها

در زبان‌های برنامه نویسی معمولاً ۵ دسته توکن وجود دارد:

- (1) برای هر کلمه‌کلیدی یک توکن وجود دارد. الگوی این توکن، دنباله کاراکترهای کلمه‌کلیدی است.
- (2) برای عملگرها توکن وجود دارد. می‌توان برای هر عملگر یک توکن داشت و یا برای عملگرهای مشابه یک توکن در نظر گرفت. الگوی این توکن نیز دنباله کاراکتری عملگر است.
- (3) برای شناسه‌ها (identifier) – مانند نام متغیر، نام تابع و نام کلاس – یک توکن وجود دارد. الگوی مربوط به این توکن وابسته به زبان برنامه‌نویسی است.
- (4) یک یا چند توکن نشان دهنده ثابت‌ها (ثابت‌های عددی و رشته‌ای) وجود دارد.
- (5) توکن‌های مربوط به علائم (, ; : { } []) نیز وجود دارد.

مثال توکن‌ها

IF, WHILE, RETURN

(1) کلمه‌کلیدی

GT > , LTEQ <=, EQ ==, RELOP

(2) عملگرها

ID var1 max

(3) شناسه

NUMBER 10.25, STRING_LITERAL "Hi"

(4) ثابت‌ها

OPP (, CLP), OPB {, CLB }

(5) علائم

خصوصیات توکن‌ها

- تحلیلگر لغوی در برخی موارد، علاوه بر نام توکن، باید داده‌های بیشتری را همراه آن قرار دهد.
- مثلاً هنگامی که یک lexeme از نوع NUMBER است، باید مقدار آن نیز مشخص باشد.
- در مرحله تحلیل گرامر (Syntax analyser)، از نام توکن‌ها استفاده می‌شود.
- در مرحله تحلیل معنا (Semantic analyser)، از خصوصیت توکن‌ها هم استفاده می‌شود.
- برخی از توکن‌ها (مانند شناسه) دارای بیش از یک خصوصیت هستند.
 - مثلاً علاوه بر lexeme مربوط به شناسه، باید نوع آن را هم تعیین کرد.
- اطلاعات این توکن‌ها در جدول نمادها (Symbol Table) ذخیره می‌شود.

مثال از دنباله توکن‌های تولید شده

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

<**id**, pointer to symbol-table entry for E>

<**assign_op**>

<**id**, pointer to symbol-table entry for M>

<**mult_op**>

<**id**, pointer to symbol-table entry for C>

<**exp_op**>

<**number**, integer value 2>

الگوی توکن‌ها را چگونه بیان کنیم؟

یک کاراکتر حرفی و بعد از آن
یک یا چند کاراکتر حرفی یا عددی

توصیفی

$[a-zA-Z][a-zA-Z0-9]^+$

با عبارت‌های منظم

عبارت‌های منظم

EXPRESSION	MATCHES	EXAMPLE
c	the one non-operator character c	a
$\backslash c$	character c literally	$\backslash *$
$"s"$	string s literally	$"**"$
$.$	any character but newline	$a.*b$
$^$	beginning of a line	abc
$\$$	end of a line	$abc\$$
$[s]$	any one of the characters in string s	$[abc]$
$[^s]$	any one character not in string s	$[^abc]$
r^*	zero or more strings matching r	a^*
r^+	one or more strings matching r	a^+
$r^?$	zero or one r	$a^?$
$r\{m,n\}$	between m and n occurrences of r	$a[1,5]$
r_1r_2	an r_1 followed by an r_2	ab
$r_1 \mid r_2$	an r_1 or an r_2	$a \mid b$
(r)	same as r	$(a \mid b)$

مقدمه‌ای بر عبارتهای منظم

■ اولویت عملگرها

■ $*$ $+$ $?$

■ $()$

■ concat

■ $|$

$ad+b$

$ab?c$

$ab | cd$

$a(b | c)d$

$(ab) | (cd)$

$(ab^*) | (cd+)$

$a | bc | d+e$

نمونه الگو

- یک رقم
- اعداد صحیح
- عدد اعشاری
- نماد علمی
- شناسه
- کلمه کلیدی while
- عملگرهای رابطه‌ای
- عملگر تقسیم
- علامت ;

تعریف منظم

نامگذاری الگوهای ساده برای استفاده در الگوهای پیچیده

$$d1 ==> r1$$

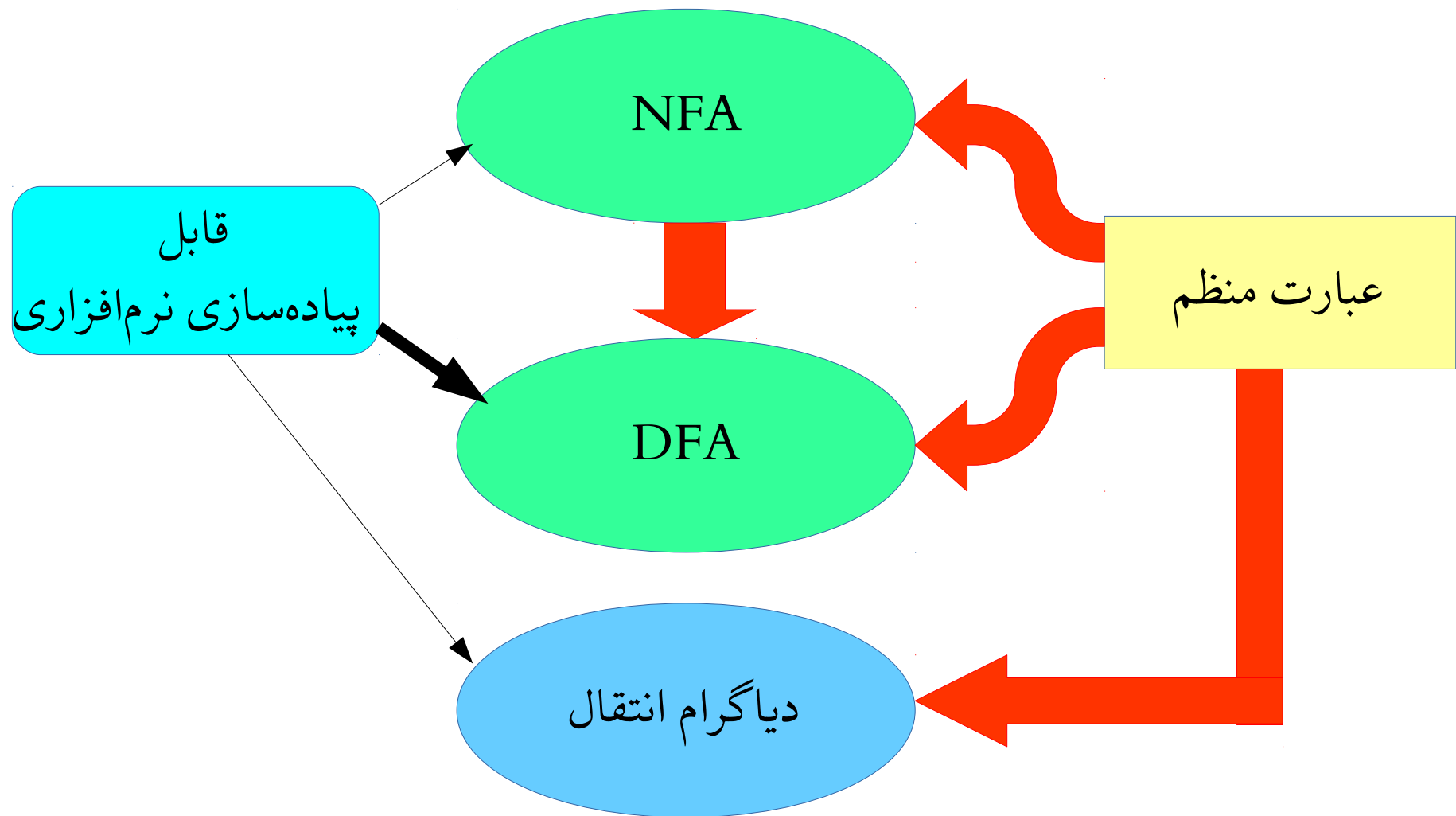
$$d2 ==> r2$$

$$dn ==> d1 \ d2$$

تعریف منظم

<i>digit</i>	→	[0-9]
<i>digits</i>	→	<i>digit</i> ⁺
<i>number</i>	→	<i>digits</i> (. <i>digits</i>)? (E [+-]? <i>digits</i>)?
<i>letter</i>	→	[A-Za-z]
<i>id</i>	→	<i>letter</i> (<i>letter</i> <i>digit</i>)*
<i>if</i>	→	<i>if</i>
<i>then</i>	→	<i>then</i>
<i>else</i>	→	<i>else</i>
<i>relop</i>	→	< > <= >= = <>

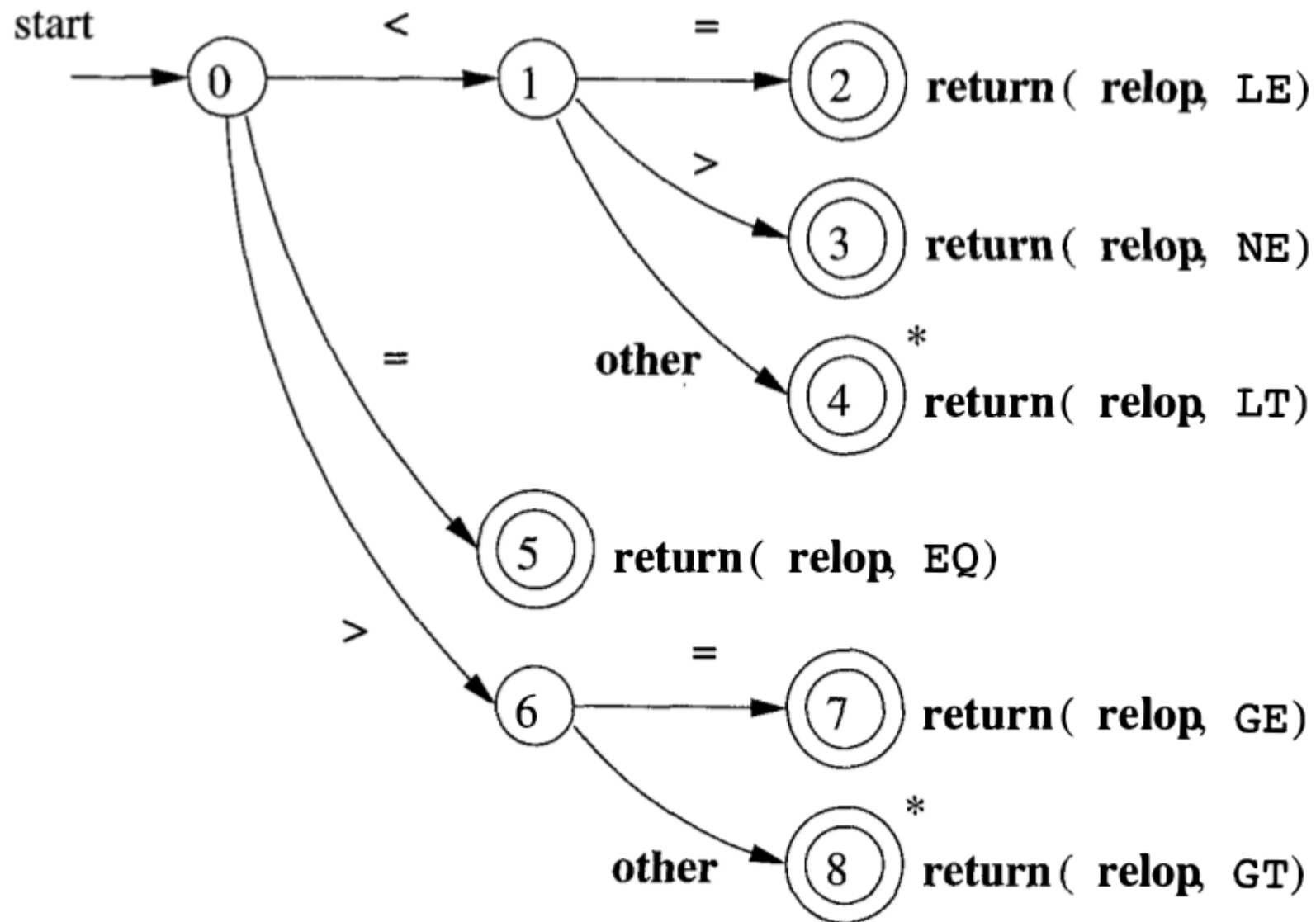
چطور عبارت‌های منظم را پیاده‌سازی کنیم؟



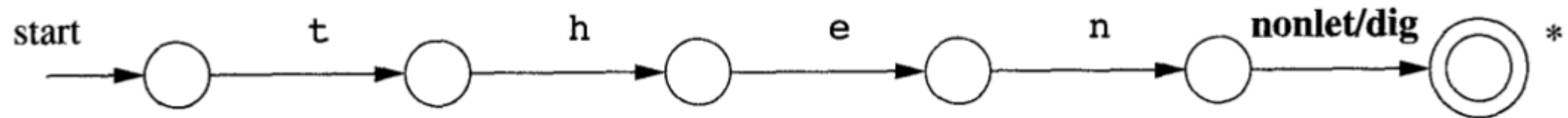
دیاگرام انتقال Transition Diagram

- عبارتهای منظم را می توان به دیاگرام انتقال تبدیل کرد.
- دیاگرام انتقال یک گراف جهت دار است.
- به راس های دیاگرام، حالت (state) گفته می شود. یک حالت شروع نیز وجود دارد.
- یال های بین راس ها دارای برچسب هستند.
 - یک یا چند کاراکتر
- حالت (state) ها، قطعی (deterministic) هستند.
 - با یک ورودی خاص، از یک حالت فقط به یک حالت دیگر می توان رفت.
- برخی از حالت ها، پذیرنده (accepting) هستند.
 - حالت پذیرنده نشان دهنده کشف یک lexeme است.
- در برخی از حالت های پذیرنده، باید ورودی برگردد. در این موارد از * استفاده می شود.

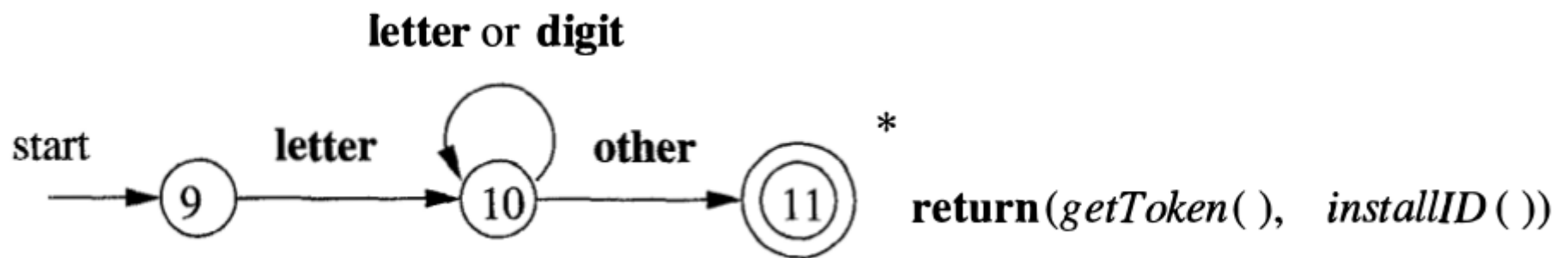
نمونه دیاگرام انتقال



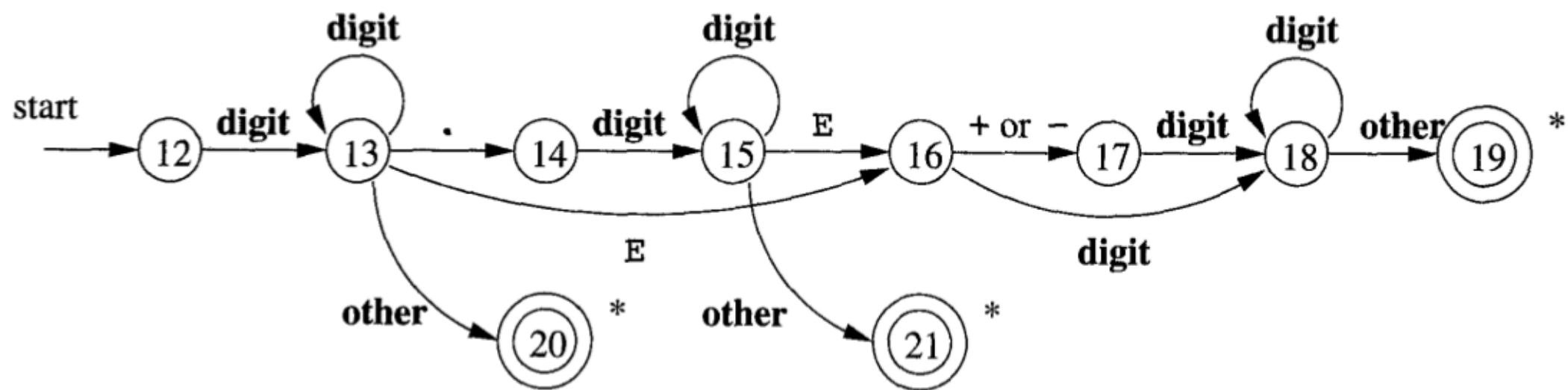
نمونه دیاگرام انتقال



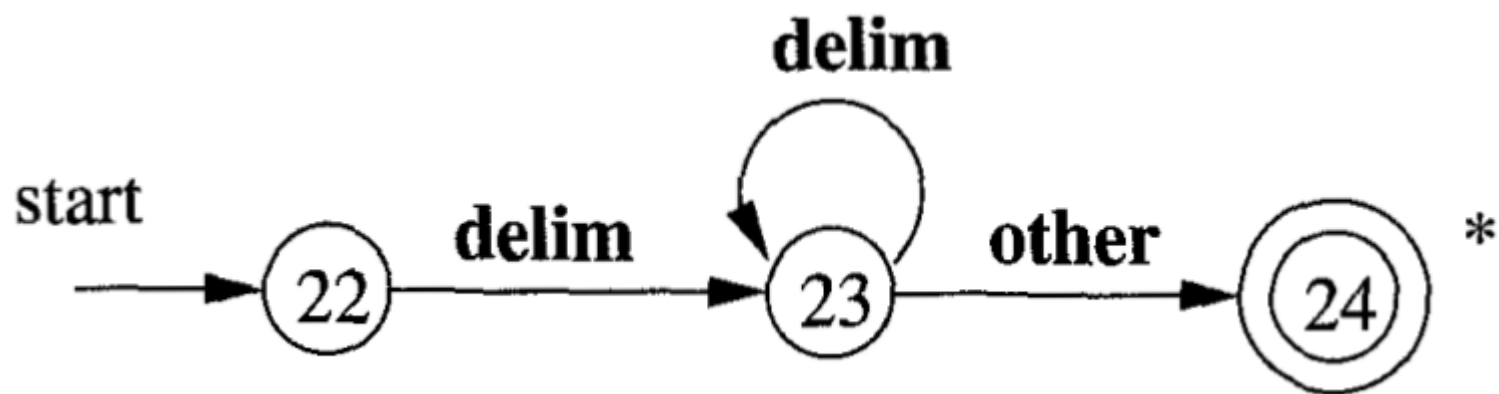
نمونه دیاگرام انتقال



نمونه دیاگرام انتقال



نمونه دیاگرام انتقال



نمونه دیاگرام انتقال

```
TOKEN getRelop()
{
    TOKEN retToken = new(RELOP);
    while(1) { /* repeat character processing until a return
                or failure occurs */
        switch(state) {
            case 0: c = nextChar();
                    if ( c == '<' ) state = 1;
                    else if ( c == '=' ) state = 5;
                    else if ( c == '>' ) state = 6;
                    else fail(); /* lexeme is not a relop */
                    break;
            case 1: ...
                    ...
            case 8: retract();
                    retToken.attribute = GT;
                    return(retToken);
        }
    }
}
```