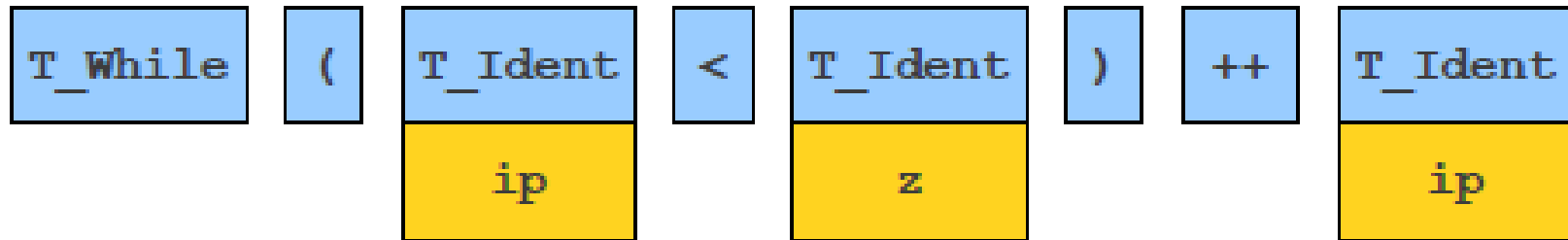
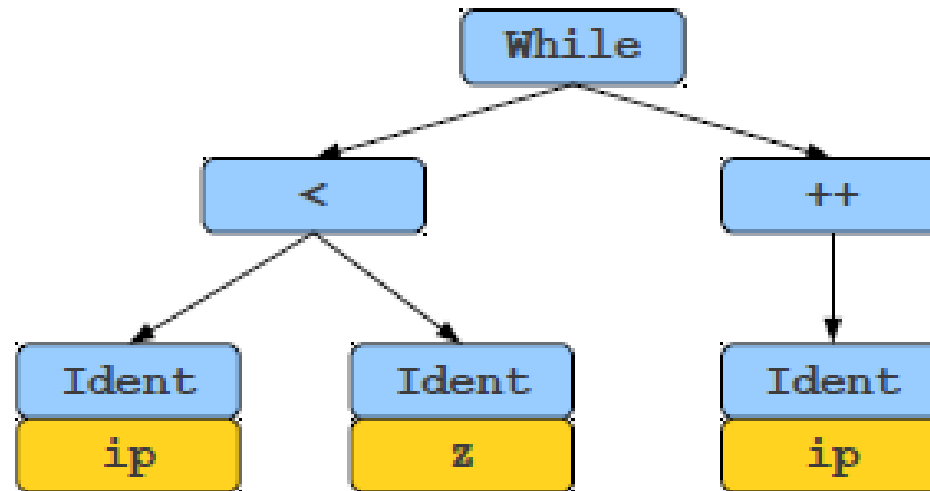


# طراحی کامپایلرها

تحلیل لغوی (بخش 1)

قمرناز تدین تبریزی



w	h	i	l	e		(	i	p		<		z	)	\n	\t	+	+	i	p	;
---	---	---	---	---	--	---	---	---	--	---	--	---	---	----	----	---	---	---	---	---

```
while (ip < z)
    ++ip;
```



T_Do	[	T_For	]	=	T_New	T_IntConst
						0

d	o	[	f	o	r	]		=		n	e	w		0	;
---	---	---	---	---	---	---	--	---	--	---	---	---	--	---	---

do[for] = new 0;

# پویش فایل مبدأ

4

w	h	i	l	e		(	1	3	7		<		i	)	\n	\t	+	+	i	;
---	---	---	---	---	--	---	---	---	---	--	---	--	---	---	----	----	---	---	---	---

w	h	i	l	e		(	1	3	7		<		i	)	\n	\t	+	+	i	;
---	---	---	---	---	--	---	---	---	---	--	---	--	---	---	----	----	---	---	---	---

w	h	i	l	e		(	1	3	7		<		i	)	\n	\t	+	+	i	;
---	---	---	---	---	--	---	---	---	---	--	---	--	---	---	----	----	---	---	---	---

w	h	i	l	e		(	1	3	7		<		i	)	\n	\t	+	+	i	;
---	---	---	---	---	--	---	---	---	---	--	---	--	---	---	----	----	---	---	---	---

w	h	i	l	e		(	1	3	7		<		i	)	\n	\t	+	+	i	;
---	---	---	---	---	--	---	---	---	---	--	---	--	---	---	----	----	---	---	---	---

w	h	i	l	e		(	1	3	7		<		i	)	\n	\t	+	+	i	;
---	---	---	---	---	--	---	---	---	---	--	---	--	---	---	----	----	---	---	---	---

w	h	i	l	e		(	1	3	7		<		i	)	\n	\t	+	+	i	;
---	---	---	---	---	--	---	---	---	---	--	---	--	---	---	----	----	---	---	---	---

w	h	i	l	e		(	1	3	7		<		i	)	\n	\t	+	+	i	;
---	---	---	---	---	--	---	---	---	---	--	---	--	---	---	----	----	---	---	---	---

□ این بخش از برنامه اصلی، که توکن بر اساس آن ساخته می شود، لغت (Lexeme) گفته می شود

T\_While

□ توکن: یک نوع عددی شده است که نشان میدهد چه موجودیت منطقی از کد مبدأ خوانده شده است

w	h	i	l	e		(	1	3	7		<		i	)	\n	\t	+	+	i	;
---	---	---	---	---	--	---	---	---	---	--	---	--	---	---	----	----	---	---	---	---

گاهی اوقات یک لغت ذخیره نمیشود و  
 ندیده گرفته میشود. فضای خالی هیچ  
 معنایی برای برنامه ندارد بنابراین  
 نادیده گرفته می شود. □

T\_While

```
w h i l e   ( 1 3 7   <   i ) \n\t + + i ;
```

```
w h i l e   ( 1 3 7   <   i ) \n\t + + i ;
```

```
w h i l e   ( 1 3 7   <   i ) \n\t + + i ;
```

```
w h i l e   ( 1 3 7   <   i ) \n\t + + i ;
```

```
w h i l e   ( 1 3 7   <   i ) \n\t + + i ;
```

```
w h i l e   ( 1 3 7   <   i ) \n\t + + i ;
```

```
T_While
```

```
(
```

```
T_IntConst
137
```

برخی از توکن ها خصیصه ای دارند که اطلاعات اضافه درباره توکن را نگهداری میکند.

# هدف از تحلیل لغوی

8

- تبدیل از توصیف فیزیکی برنامه به دنباله ای از توکنها
  - هر توکن یک بخش منطقی از فایل مبدأ را نشان میدهد  
مثل یک کلمه کلیدی، نام متغیر و ...
- هر توکن با یک لغت مرتبط است.
  - متن واقعی توکن مثل 137, "int", ...
- هر توکن میتواند خصیصه ای داشته باشد که از متن به دست می آید مثل مقدار عدد
- دنباله توکن ها در مرحله بعد در پارسر به کار می رود.



# انتخاب توکن

9

چه توکنهایی مفید هستند؟ □

```
for (int k = 0; k < myArray[5]; ++k) {  
  cout << k << endl;  
}
```

```
For      {  
int      }  
<<      ;  
=        <  
(        [  
)        ]  
++
```

Identifier

IntegerConstant

# انتخاب توکن های خوب

10

□ بسیار وابسته به زبان است.

□ معمولا:

- کلمات کلیدی، توکن خودشان را دارند.
- هر نماد نقطه گذاری، توکن خودش را دارد.
- لغتهای نمایشگر شناسه ها، ثابتهای عددی، رشته ها و غیره با هم گروه بندی میشوند.
- اطلاعات نامربوط کنارگذاشته می شوند (فضای خالی، توضیحات)

- در fortran همه فضاهاى خالى حذف مى شوند
- **DO 5 I = 1,25**
- **DO5I = 1.25**
- تعیین این که ورودی کجا تفکیک شود مشکل است.

□ در زبان PL/1 کلمات کلیدی می توانند به عنوان شناسه استفاده شوند.

□ IF THEN THEN THEN = ELSE; ELSE ELSE = IF

□ **IF THEN THEN THEN = ELSE; ELSE ELSE = IF**

□ برچسب گذاری لغات مشکل است

- چگونه تعیین کنیم کدام لغات با هر توکن مرتبط هستند.
- اگر چند روش برای پویش ورودی وجود دارد چگونه بدانیم کدام را انتخاب کنیم.
- چگونه میتوان موارد فوق را با کارایی در نظر گرفت

# ارتباط لغات و توکنها

14

- توکنها روشی برای گروه بندی لغات فراهم میکنند.
- برخی توکنها ممکن است فقط با یک لغت مرتبط شوند.
  - توکنها برای کلمات کلیدی مثل if و While میتوانند جزء این دسته باشند.
- بعضی توکنها ممکن است با لغات زیادی مرتبط شوند:
  - نام متغیرها، اعداد، رشته ها و...

# مجموعه لغات

15

- ارتباط لغات با هر توکن
- توکن "number" با مجموعه  $\{0, 1, 2, \dots, 10, 11, 12, \dots\}$
- توکن "string" با مجموعه  $\{\dots, "a", "b", " ", \dots\}$
- توکن While با مجموعه  $\{\text{while}\}$

# Formal Languages

16

- زبانهای رسمی، مجموعه ای از رشته ها هستند.
- بسیاری از زبانهای نامتناهی دارای تعریف متناهی هستند.
  - تعریف زبان با استفاده از اتوماتا
  - تعریف زبان با استفاده از گرامر
  - تعریف زبان با استفاده از عبارت باقاعده



# عبارت باقاعده

17

- مجموعه ای از توصیفها برای تعریف زبانهای معینی است.
- توصیف فشرده و قابل فهمی از زبان
- اساسی برای بسیاری سیستمهای نرم افزاری مثل Flex که بعدا بررسی خواهد شد

# عبارات باقاعده اتمیک

18

- نماد  $\frac{\circ}{\circ}$  یا رشته تهی را نشان میدهد.
- برای هر نماد  $a$ ، یک عبارت باقاعده است که با کاراکتر  $a$  منطبق می شود.

# عبارات باقاعده ترکیبی

19

□ اگر  $R1$  و  $R2$  عبارت باقاعده باشند،  
 $R1R2$  (الحاق)،  $R1|R2$  (اجتماع) و  $R1^*$  بستار  
ستاره ای آنها هم باقاعده است.

□ اولویت عملگرها:

$(R)$

$R^*$

$R_1R_2$

$R_1 | R_2$

# مثال

20

□ الفبای 0 و 1:

□ رشته های شامل 00:  $(1|0)^*00(1|0)^*$

11011100101  
0000  
1111101111001111

# مثال

21

- الفبای 0 و 1:
- رشته های با طول 4:

$(0|1)(0|1)(0|1)(0|1)$        $(0|1)\{4\}$

0000  
1010  
1111  
1000

# مثال

22

□ الفبای 0 و 1:

□ رشته هایی که حداکثر یک صفر دارند:

$1^*(0 \mid \epsilon)1^*$        $1^*0?1^*$

11110111

111111

0111

0

# مثال کاربردی

23

- فرض کنید الفبا شامل  $a$  ,  $@$  و  $.$  باشد که  $a$  به معنی هر حرف الفبا باشد.
- عبارت باقاعده برای آدرس ایمیل:

$aa^* (.aa^*)^* @ aa^*.aa^* (.aa^*)^*$

$a^+ (.a^+)^* @ a^+ (.a^+)^*$

cs143@cs.stanford.edu  
first.middle.last@mail.site.org  
barack.obama@whitehouse.gov

# مثال کاربردی

24

- الفبا: همه کاراکترهای اسکی
- عبارت باقاعده برای اعداد زوج

$(+|-)?(0|1|2|3|4|5|6|7|8|9)^*(0|2|4|6|8)$

$(+|-)?[0-9]^*[02468]$

42

+1370

-3248

-9999912



# پیاده سازی عبارتهای باقاعده

25

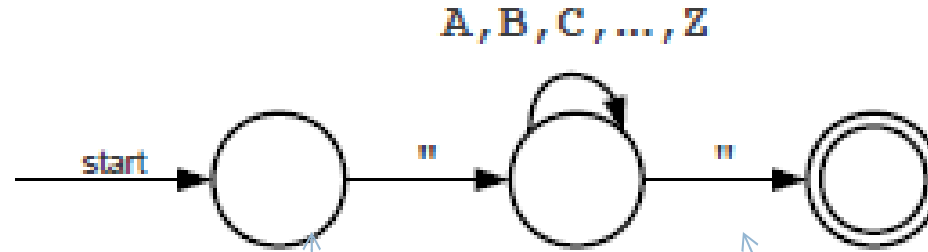
□ عبارات باقاعده با اتوماتای متنهای قابل پیاده سازی هستند.

□ دو نوع اصلی:

- NFA(Nondeterministic Finite Automata)
- DFA(Deterministic Finite Automata)

# یک اتوماتای ساده

26

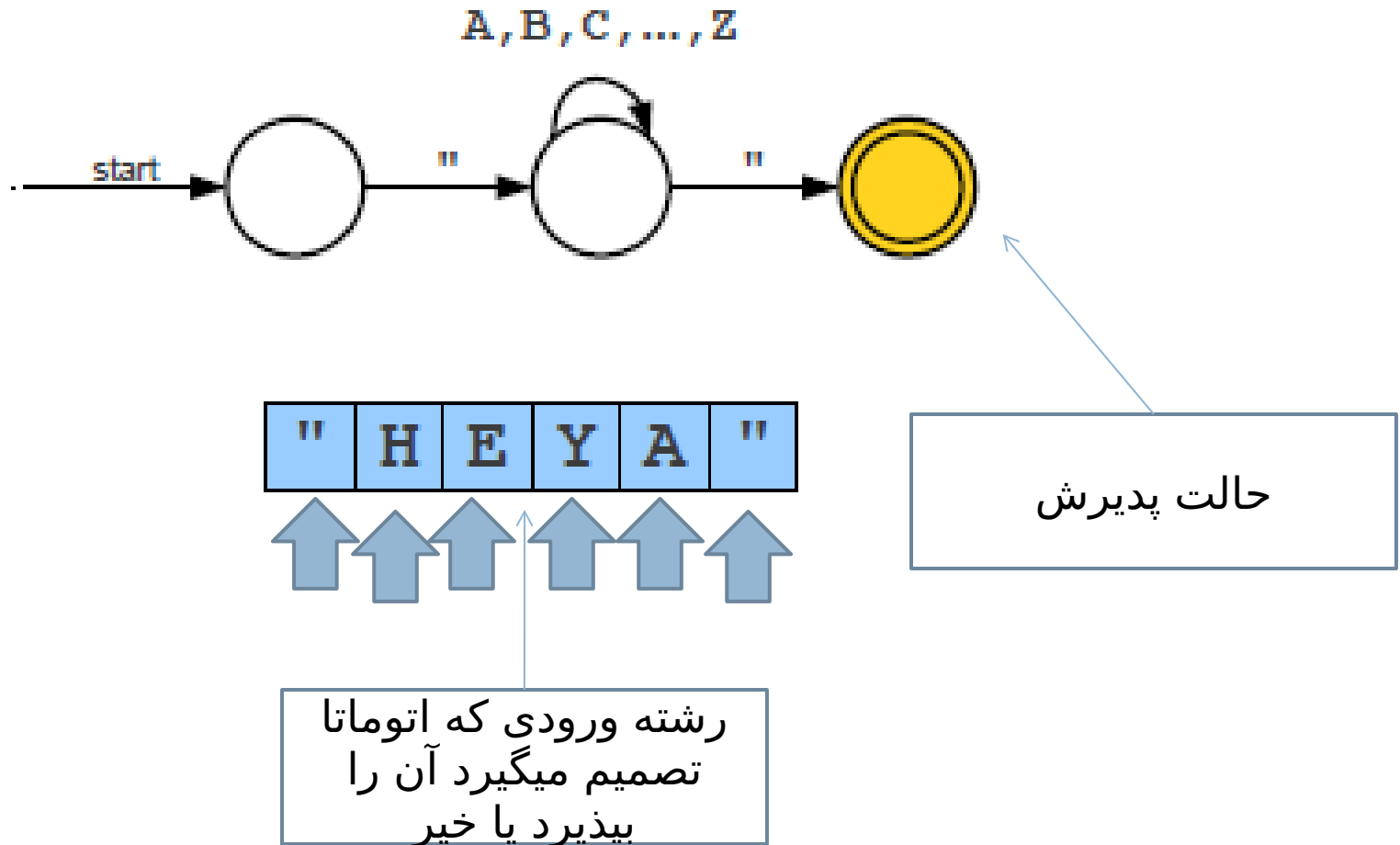


هر دایره یک وضعیت  
اتوماتا است.

کمانها تغییر وضعیت  
هستند

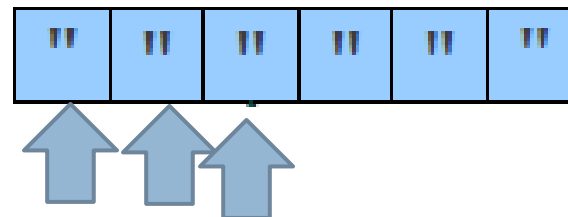
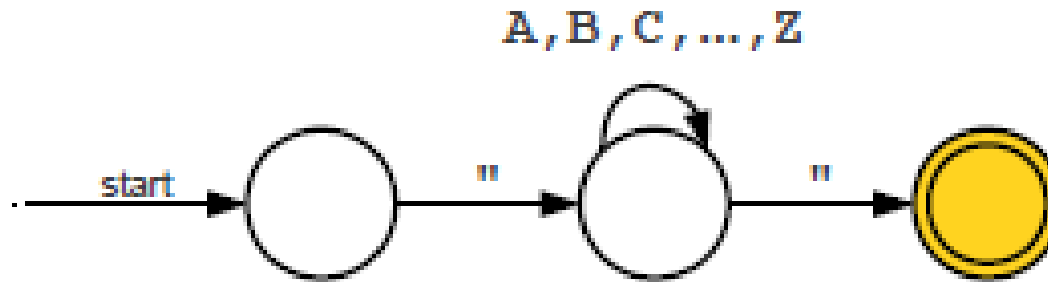
# یک اتوماتای ساده، ادامه....

27



# یک اتوماتای ساده، ادامه....

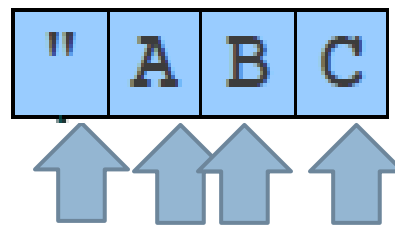
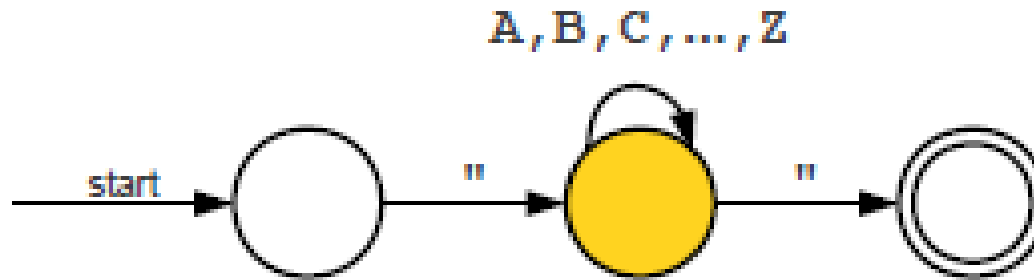
28



هیچ تغییر حالتی وجود ندارد. ورودی رد می شود.

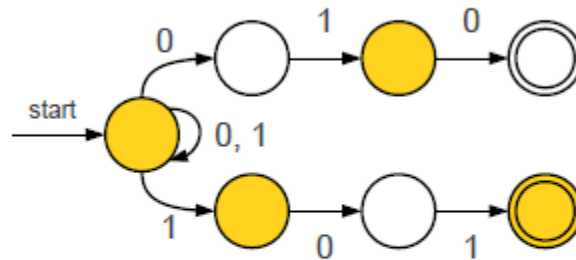
# یک اتوماتای ساده، ادامه....

29



این حالت، پذیرش نیست  
پس اتوماتا رشته را رد  
میکند.

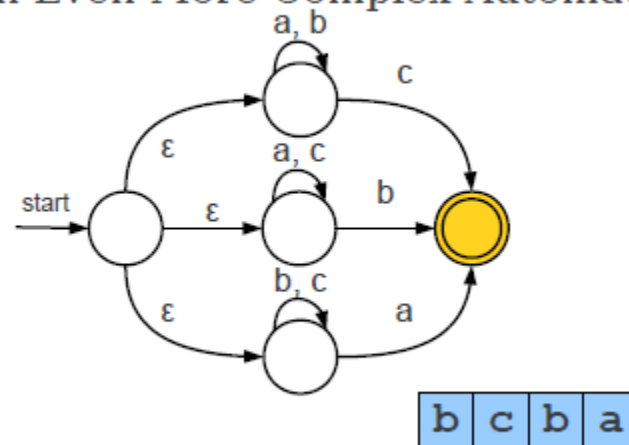
### A More Complex Automaton



پذیرش

0	1	1	1	0	1
---	---	---	---	---	---

### An Even More Complex Automaton



گذر  $\epsilon$  به صورت خودکار  
و بدون استفاده از  
ورودی دنبال می شود

# شبیه سازی NFA

32

- وضعیتهای دنبال می شوند. پردازش از وضعیت شروع و هر وضعیتی که با تغییر حالت  $\delta$  قابل دسترسی باشد آغاز می شود.
- برای هر کاراکتر ورودی:
  - مجموعه ای از وضعیتهای بعدی نگهداری می شوند. (مقدار اولیه تهی است)
  - برای هر وضعیت جاری:
    - همه انتقالهایی که برچسب آنها حرف جاری است دنبال می شوند.
    - این وضعیتهای به مجموعه وضعیتهای جدید اضافه می شوند.
  - هر وضعیتی که با حرکت  $\delta$  قابل دسترسی باشد به مجموعه وضعیت های بعدی اضافه میشود.
- پیچیدگی:  $O(mn^2)$  برای رشته هایی به طول  $m$  و اتوماتایی با  $n$  وضعیت



- هر عبارت باقاعده با طول  $n$  را می توان به یک NFA با  $O(n)$  وضعیت تبدیل کرد.
- در زمان  $O(mn^2)$  میتوان تعیین کرد که آیا رشته ای با طول  $m$  با عبارت باقاعده با طول  $n$  منطبق می شود یا خیر.
- بعدا نشان میدهیم چگونه این زمان را به  $O(m)$  کاهش میدهیم.