

COMP130014编译

第三讲：上下文无关文法

徐辉

xuh@fudan.edu.cn



主要内容

一、上下文无关文法

二、扩展BNF范式

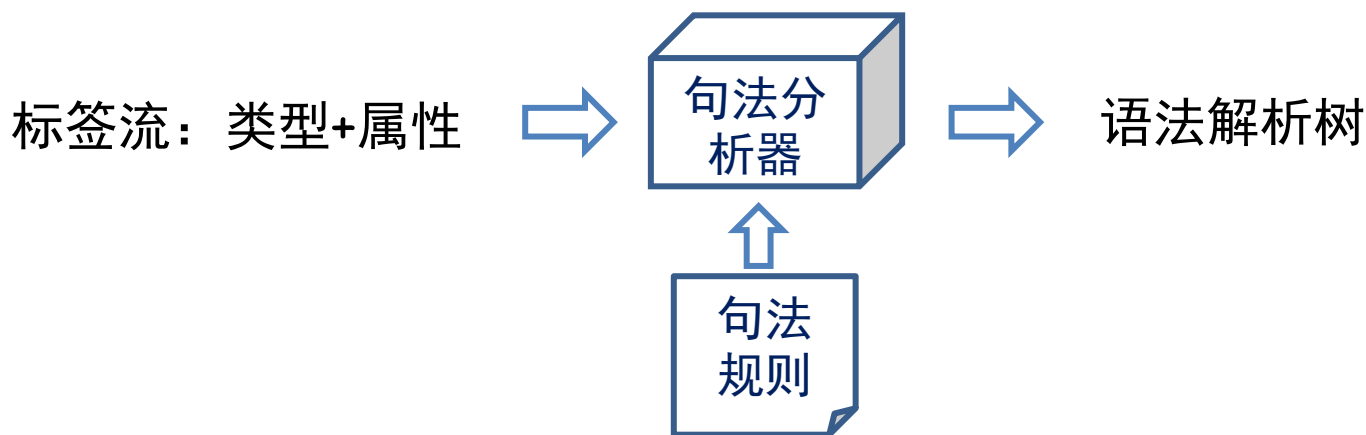
三、TeaPL文法定义

四、语言分析问题

一、上下文无关文法

句法解析问题

- 给定一个句子和句法规则，找到可生成该句子的一个推导
- 句法规则定义了句法分析器可接受的标签序列及其推导方式
- 文法/语法 = 词法 + 句法



语法推导举例

语法规则

- [1] $E \rightarrow E + E$
- [2] $E \rightarrow E - E$
- [3] $E \rightarrow E \times E$
- [4] $E \rightarrow E / E$
- [5] $E \rightarrow \langle \text{UNUM} \rangle$

目标句子: $1 + 2 \times 3$

$\langle \text{UNUM}(1) \rangle + \langle \text{UNUM}(2) \rangle \times \langle \text{UNUM}(3) \rangle$

推导步骤

- [1] $E \rightarrow E + E$
- [5] $E \rightarrow \langle \text{UNUM}(1) \rangle + E$
- [3] $E \rightarrow \langle \text{UNUM}(1) \rangle + E \times E$
- [5] $E \rightarrow \langle \text{UNUM}(1) \rangle + \langle \text{UNUM}(2) \rangle \times E$
- [5] $E \rightarrow \langle \text{UNUM}(1) \rangle + \langle \text{UNUM}(2) \rangle \times \langle \text{UNUM}(3) \rangle$

基本概念和符号

- 一门语言是多个句子的集合
- 句子是由终结符组成的序列
- 字符串是包含终结符和非终结符的序列
 - 非终结符：X、Y、Z
 - 终结符（标签）：<BINOP>、<NUM>
 - 字符串符号： α 、 β 、 γ
- 语法包含一个开始符号和多条推导规则
 - $S \rightarrow \beta$
- 语法G的语言L(G)是该语法可推导的所有句子的集合

上线文无关文法 (Context-Free Grammar)

- 上下文无关语法是一个四元组(T, NT, S, P)
- T : 终结符
- NT : 非终结符
- S : 起始符号
- P : 推导规则集合: $\{X \rightarrow \gamma\}$
 - X 是非终结符
 - γ 是字符串
 - 规则左侧只能有一个非终结符

括号匹配问题

- 用CFG语法设计一套括号匹配规则
- 验证： $()()()$ 是该语法的一个推导吗？

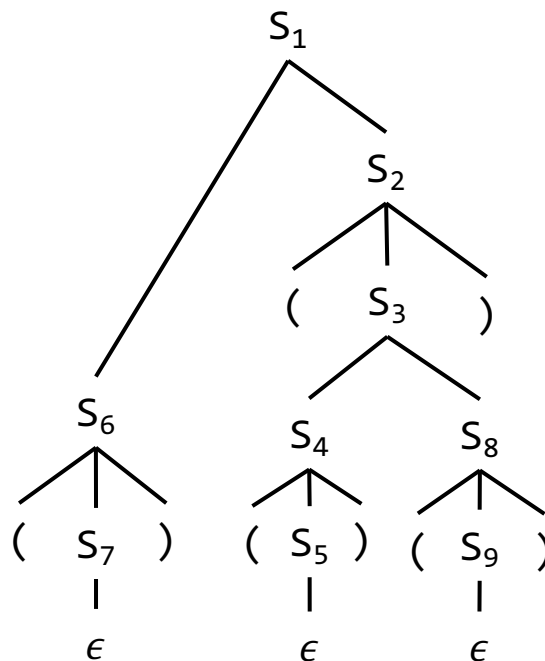
语法规则

- | | |
|-----|--------------------------|
| [1] | $S \rightarrow \epsilon$ |
| [2] | $S \rightarrow (S)$ |
| [3] | $S \rightarrow SS$ |

推导

- [3] $S \rightarrow SS$
[2] $S \rightarrow S(S)$
[3] $S \rightarrow S(SS)$
[2] $S \rightarrow S(S(S))$
[1] $S \rightarrow S(S())$
[2] $S \rightarrow S((S)())$
[1] $S \rightarrow S()()$
[2] $S \rightarrow (S)()()$
[1] $S \rightarrow ()()()$

语法解析树



写出计算器的CFG文法

```
[1] E → E <ADD> E
[2] E → E <SUB> E
[3] E → E <MUL> E
[4] E → E <DIV> E
[5] E → E <POW> E
[6] E → <LPAR> E <RPAR>
[7] E → NUM
[8] NUM → <UNUM>
[9] NUM → <SUB> <UNUM>
```

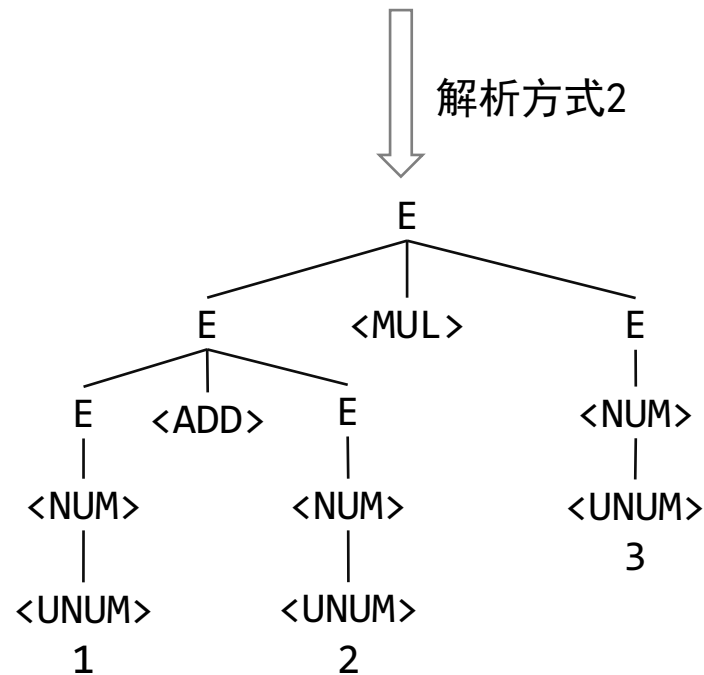
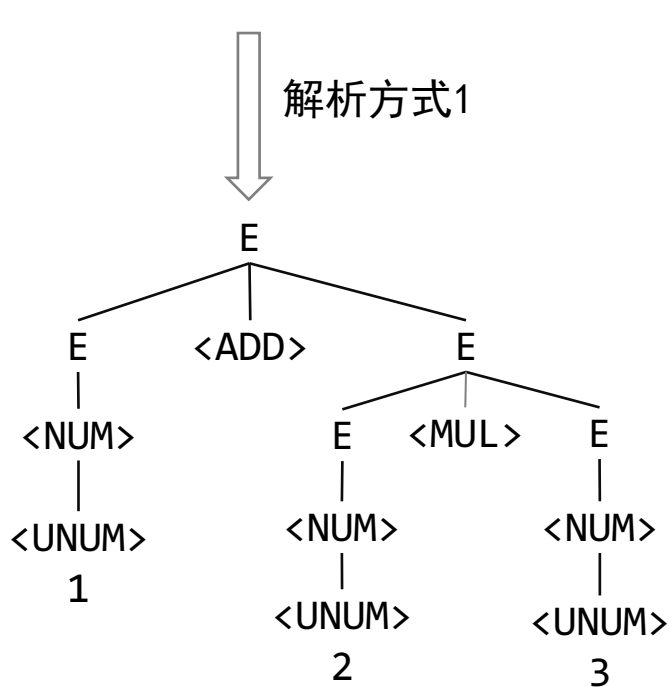
=

```
[1] E → E '+' E
[2] E → E '/' E
[3] E → E '*' E
[4] E → E '^' E
[5] E → '(' E ')'
[6] E → NUM
[7] E → NUM
[8] NUM → <UNUM>
[9] NUM → '-' <UNUM>
```

二义性问题

- $L(G)$ 中的某个句子存在一个以上的最左（或最右）推导
- 语法解析树不同

标签流: $\langle \text{UNUM}(1) \rangle \langle \text{ADD} \rangle \langle \text{UNUM}(2) \rangle \langle \text{MUL} \rangle \langle \text{UNUM}(3) \rangle$







A programmer's wife asks him to go to the grocery. She says "Get a gallon of milk. If they have eggs, get 12."

The programmer returns with 12 gallons of milk.

消除二义性：将运算符特性融入语法规则

- 优先级： $^ > \times / \div > + / -$
- 结合性： $\times / \div / + / -$ 左结合， $^$ 右结合

```
[1] E → E '+' E
[2] E → E '/' E
[3] E → E '*' E
[4] E → E '/' E
[5] E → E '^' E
[6] E → '(' E ')'
[7] E → NUM
[8] NUM → <UNUM>
[9] NUM → '-' <UNUM>
```



```
[1] E → E OP1 E1
[2] E → E1
[3] E1 → E1 OP2 E2
[4] E1 → E2
[5] E2 → E3 OP3 E2
[6] E2 → E3
[7] E3 → NUM
[8] E3 → '(' E ')'
[9] NUM → <UNUM>
[10] NUM → '-' <UNUM>
[11] OP1 → '+'
[12] OP1 → '-'
[13] OP2 → '*'
[14] OP2 → '/'
[15] OP3 → '^'
```

练习：为下列语言设计语法规则：

- 1) 所有0和1组成的字符串，每个0后面紧跟若干个1
- 2) 所有0和1组成的字符串，0和1的个数相同
- 3) 所有0和1组成的字符串，0和1的个数不相同

练习：语法设计

- 为正则语言设计CFG（用于解析正则表达式）
 - 支持字符 [A-Za-z0-9]
 - 支持连接、选择|、闭包*
 - 支持()
- 检查语法是否有二义性？

二、扩展BNF范式

CFG的问题

- 规则条目多且复杂，易写易读性差
- 语法解析树复杂

```
[1] E → E OP1 E1
[2] E → E1
[3] E1 → E1 OP2 E2
[4] E1 → E2
[5] E2 → E3 OP3 E2
[6] E2 → E3
[7] E3 → NUM
[8] E3 → '(' E ')'
```

[9] NUM → <UNUM>

[10] NUM → '-' <UNUM>

[11] OP1 → '+'

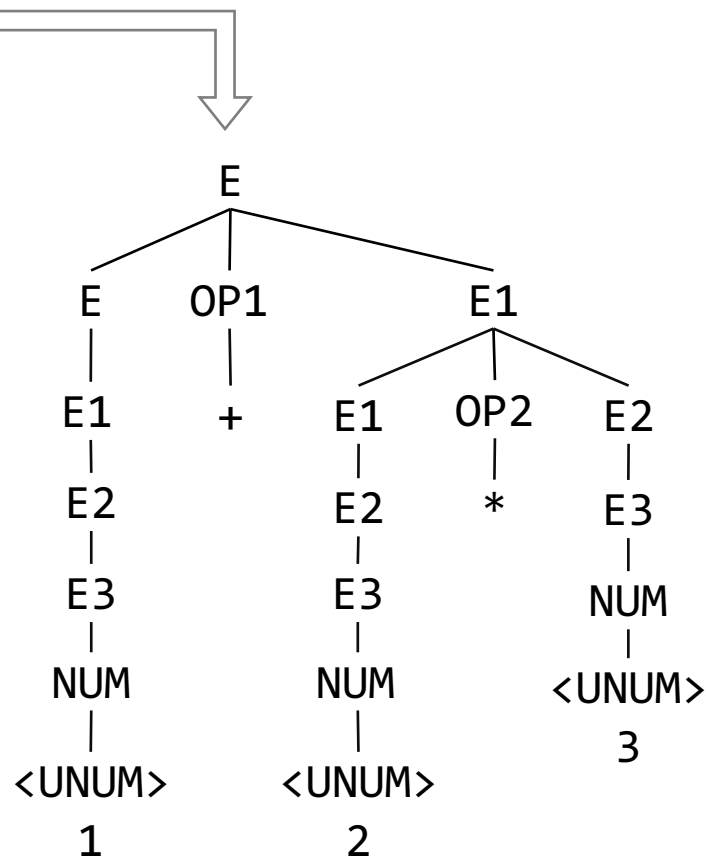
[12] OP1 → '-'

[13] OP2 → '*'

[14] OP2 → '/'

[15] OP3 → '^'

应用：1+2*3的语法解析树



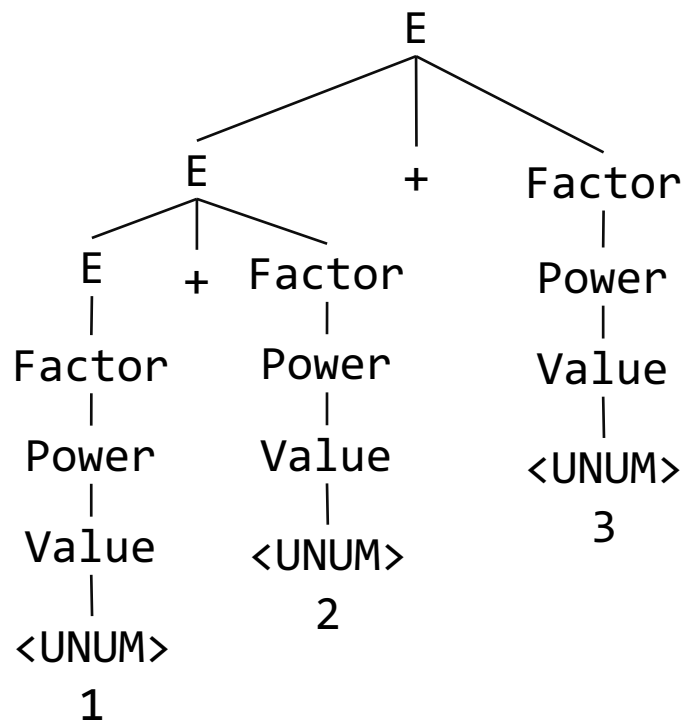
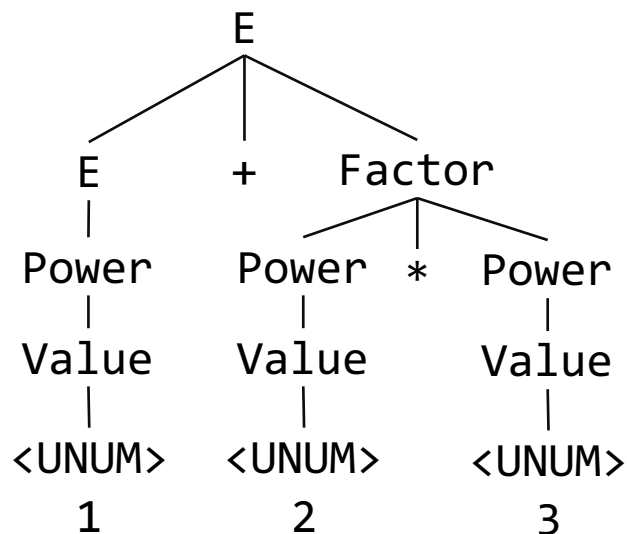
扩展BNF范式（Extended Backus-Naur Form）

- 引入更多运算符提升规则描述效率
- EBNF运算符有多种表示方法，我们沿用正则文法符号

构造方式	符号	优先级	示例
字符串	' '	5	'a'
通配符	.	5	
字符集	[]	5	[a-z]
可选匹配	?	4	$\beta?$
闭包	*	4	β^*
正闭包	+	4	β^+
非	!	3	$!\beta$
连接		2	$\alpha\beta$
选择		1	$\alpha \beta$

EBNF及其对应的语法解析树

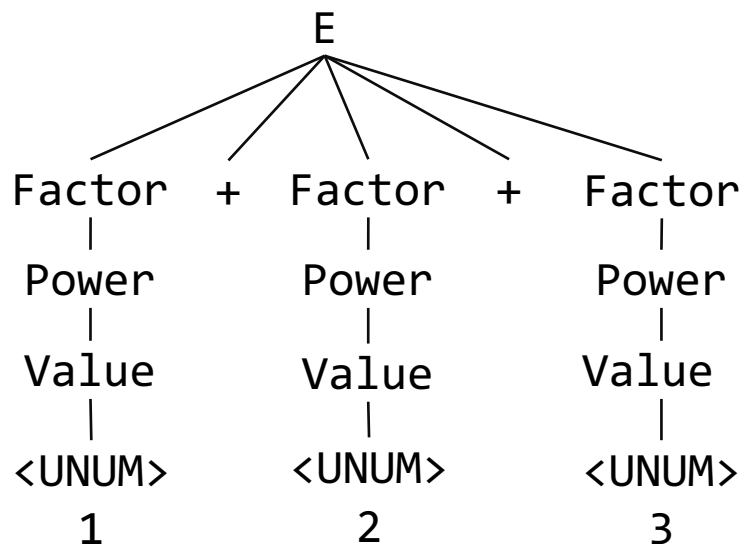
- [1] $E \rightarrow (E ('+' | '-'))? \text{ Factor}$
- [2] $\text{Factor} \rightarrow (\text{Factor} ('*' | '/'))? \text{ Power}$
- [3] $\text{Power} \rightarrow \text{Value} ('^' \text{ Power})?$
- [4] $\text{Value} \rightarrow \langle \text{UNUM} \rangle \mid ('-' \langle \text{UNUM} \rangle) \mid ('(' E ')')$



EBNF对应的语法解析树

- [1] $E \rightarrow \text{Factor } (('+' | '-') \text{Factor})^*$
 - [2] $\text{Factor} \rightarrow \text{Power } (('*' | '/') \text{Power})^*$
 - [3] $\text{Power} \rightarrow \text{Value } ('^' \text{Power})?$
 - [4] $\text{Value} \rightarrow \langle \text{UNUM} \rangle \mid ('-' \langle \text{UNUM} \rangle) \mid ('(' E ')')$

?



练习：使用EBNF设计/改写正则语言文法

- 为正则语言设计CFG（用于解析正则表达式）
 - 支持字符 [A-Za-z0-9]
 - 支持连接、或|、闭包*运算
 - 支持()
- 检查语法是否有二义性？

思考

- 扩展EBNF和CFG表达能力是否等价？

三、TeaPL文法定义

使用EBNF定于TeaPL：程序组成

$\text{program} \mapsto (\text{varDeclStmt} \mid \text{structDef} \mid \text{fnDeclStmt} \mid \text{fnDef} \mid \text{macro} \mid \text{comment})^*$

<code>varDeclStmt</code>	全局变量声明
<code>structDef</code>	数据结构定义
<code>fnDeclStmt</code>	函数声明
<code>fnDef</code>	函数定义
<code>macro</code>	宏
<code>comment</code>	注释

变量声明形式

```
let a:int;
```

→ 变量声明

```
let a:int = 0;
```

→ 声明时初始化

```
let a;
```

→ 类型可省略

```
let a = 0;
```

```
let a[5]:int;
```

→ 支持数组类型

```
let a[n]:int;
```

```
let a[n];
```

```
let a[2]:int = {0};
```

→ 数组声明时初始化

```
let a[2]:int = {1, 2};
```

→ 长度应保持一致（类型检查）

不支持：

- 二维数组：let a[m][n];
- 一条语句同时声明多个变量：let i,j;

变量声明

$\text{varDeclStmt} \mapsto \text{'let' (varDecl | varDef) ';'}$

$\text{varDecl} \mapsto \text{id (':' type)?}$
 $\quad \quad \quad | \text{id '[' (id | num) ']' (':' type)?}$

$\text{varDef} \mapsto \text{id (':' type)? '=' rightVal}$
 $\quad \quad \quad | \text{id '[' (id | num) ']' (':' type)? '=' '{' num '}'}$

类型

$\text{type} \mapsto \text{primitiveType} \mid \text{structType} \mid \text{ptrType}$

$\text{primitiveType} \mapsto \text{int} \mid \text{bool} \mid \text{char} \mid \text{long} \mid \text{float} \mid \text{double}$

$\text{structType} \mapsto \text{id}$

$\text{ptrType} \mapsto '*' \text{ type}$

$\text{structDef} \mapsto \text{'struct' id '{' fieldDecl (, fieldDecl)* '}'}$

$\text{fieldDecl} \mapsto \text{id ':' type}$

$\mid \text{id '[' (id \mid \text{num}) ']' ':' type}$

```
struct MyStruct {  
    x: int,  
    y[10]: int  
}
```

→ 结构体内可以包含数组

右值表达式

`rightVal` \mapsto `arithExpr` | `boolExpr`

`arithExpr` \mapsto (`arithExpr` ('+' | '-'))? `factor`

`factor` \mapsto (`factor` ('*' | '/'))? `exprUnit`

`exprUnit` \mapsto `num` | `id` | `fnCall` | '(' `rightVal` ')'
| `id` '.' `id` | `id` '[' (`id` | `num`) ']'
| `deref` | `addr` | `string`

`num` \mapsto `unum` | ('-' `unum`)

`deref` \mapsto '*' `id`

`addr` \mapsto '&' `id`

`string` \mapsto '"' (!'')* '"'

函数声明和定义

```
fn foo(a:int, b:int) -> int; → 函数声明  
fn foo(a:int, b:int) -> int { → 函数定义  
    return a + b;  
}
```

$\text{fnDeclStmt} \mapsto \text{'fn' fnSign ';'}$

$\text{fnSign} \mapsto \text{id '(' params? ')' '->' type?}$

$\text{params} \mapsto \text{id ':' type (',' id ':' type)*}$

$\text{fnDef} \mapsto \text{'fn' fnSign codeBlock}$

$\text{codeBlock} \mapsto \text{'{' (stmt | codeBlock)* '}'}$

基本语句

$\text{stmt} \mapsto \text{varDeclStmt} \mid \text{assignStmt} \mid \text{callStmt}$
 $\mid \text{retStmt} \mid \text{ifStmt} \mid \text{whileStmt}$
 $\mid \text{breakStmt} \mid \text{continueStmt} \mid ';' \mid \text{forStmt} \mid \text{matchStmt}$

$\text{assignStmt} \mapsto \text{leftVal} '=' \text{rightVal} ';' \mid$

$\text{leftVal} \mapsto \text{id} \mid \text{id} '[' (\text{num} \mid \text{id}) ']' \mid \text{id} '.' \text{id}$
 $\mid \text{deref}$

基本语句

`callStmt` \mapsto `fnCall` `';`

`fnCall` \mapsto `id` `'('` `(rightVal` `(,` `rightVal)*``)? '`

`retStmt` \mapsto `'ret'` `rightVal?` `';`

`ifStmt` \mapsto `'if'` `'('` `boolExpr` `')'` `codeBlock` `(else` `codeBlock)?`

`whileStmt` \mapsto `'while'` `'('` `boolExpr` `')'` `codeBlock`

`breakStmt` \mapsto `'break'` `';`

`continueStmt` \mapsto `'continue'` `';`

```
if(a>=b) {  
    ...  
} else {  
    ...  
}
```

→ 强制使用括号，避免歧义

条件表达式

$\text{boolExpr} \mapsto (\text{boolExpr} \text{ '||' })? \text{ andExpr}$

$\text{andExpr} \mapsto (\text{andExpr} \text{ '&&' }) \text{ boolUnit}$

$\text{boolUnit} \mapsto \text{cmpExpr} \mid \text{'!(' cmpExpr ')'} \mid \text{'!'? '(' boolExpr ')'}$

$\text{cmpExpr} \mapsto \text{exprUnit} \text{ ('>' \mid '>=' \mid '<' \mid '<=' \mid '==' \mid '!=')}$

exprUnit

```
if(!(a>=b)) {...}
```

!后强制使用括号，避免歧义

```
if((a>b) && (b>c)){...}
```

```
if(a>b>c) {...}
```

不允许，易产生歧义

```
if(a+b>c+d){...}
```

不允许，易产生歧义

```
if((a+b)>(c+d)){...}
```

```
if(a>b==c){...}
```

不允许，易产生歧义

运算符优先级总结：与C语言标准兼容

优先级	运算符	含义	结合性	TeaPL限制
8	-, !	单目运算符	右	'-'后只跟数字, '!'只能跟'('
7	*, /	乘除号	左	
6	+, -	加减号	左	
5	>, >=, <, <=	比大小	左	不支持连续比较 比较对象不支持逻辑运算
4	==, !=	等价性	左	
3	&&	逻辑与	左	
2		逻辑非	左	
1	=	赋值	右	不支持连续赋值

数字、标识符和注释

$\text{unum} \mapsto [1-9][0-9]^* | 0$
 $\quad \quad \quad | ([1-9][0-9]^* | 0) \text{'.'} [0-9]^+$

$\text{id} \mapsto [a-z_A-Z] ([a-z_A-Z0-9])^*$

$\text{comment} \mapsto \text{'//'} (!\text{newline})^* \text{newline}$

$\text{comment} \mapsto \text{'/*'} (!\text{'*/'})^* \text{'*/'}$

$\text{newline} \mapsto \text{'\n'}$

问题： 下列哪些代码符合TeaPL语法？

```
fn foo() {  
    ret 1;  
}
```

```
fn foo() {  
    ret 2>1;  
}
```

```
fn foo() -> int {  
    ret 1;  
}
```

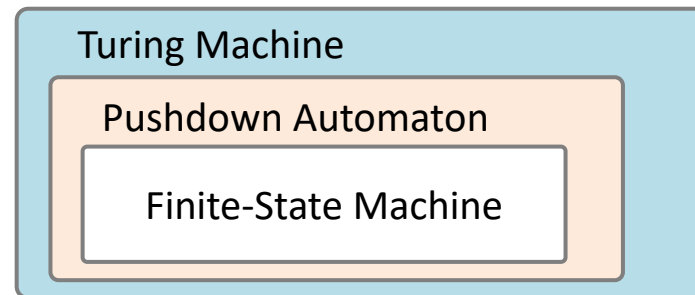
```
fn foo() -> bool {  
    ret 2>1;  
}
```

四、文法能力分类

按表达能力分类

Chomsky Hierarchy

类型	文法名称	计算模型	规则形式	语言示例
0 型	递归枚举	图灵机	无限制	
1 型	上下文敏感	线性有界图灵机	左侧可以多个符号 $\alpha S \rightarrow \beta$	$a^n b^n c^n$
2 型	上下文无关	下推自动机	左侧仅一个符号 $S \rightarrow \beta$	$a^n b^n$
3 型	正则	有穷自动机	右侧全部为终结符 $S \rightarrow ab$	a^n



正则语言 VS 上下文无关语言

- 正则语言也可以用CFG规则形式表示：
 - $X \rightarrow \gamma$
 - $\gamma \rightarrow \gamma_1$
 - ...
- 特点：右侧的非终结符均可替换为终结符

[1]	$S \rightarrow A B$
[2]	$A \rightarrow (0?1)^*$
[3]	$B \rightarrow (1?0)^*$

 $\Longrightarrow S \rightarrow (0?1)^*|(1?0)^*$

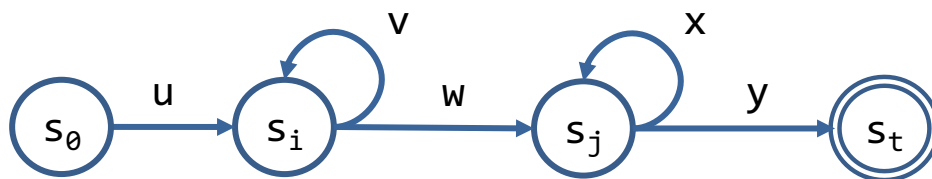
非CFG语言：上下文敏感语法

- $L = \{a^n b^n c^n, n > 0\}$ 不是CFG语言
- 上下文敏感文法规则形式： $aS \rightarrow \beta$

[1]	$S \rightarrow aBC$
[2]	$S \rightarrow aSBC$
[3]	$CB \rightarrow BC$
[4]	$aB \rightarrow ab$
[5]	$bB \rightarrow bb$
[6]	$bC \rightarrow bc$
[7]	$cC \rightarrow cc$

非CFG语言的泵引理

- CFG语言的泵引理（必要条件）：
 - 任意长度超过 p 的句子可以被拆分为 $uvwxy$ 的形式
 - v 和 x 被重复任意次后得到的新句子（如 $uvvwxxy$ ）仍属于该语言
- 正则属于CFG: $uv^nwx^n\epsilon$



练习：下列语言是否为正则语言？

- 集合表示

1) $L = \{a^n b^n | n \leq 100\}$

2) $L = \{a^n | n \geq 1\}$

3) $L = \{a^{2^n} | n \geq 1\}$

4) $L = \{a^p | p \text{ is prime}\}$

- Regex/CFG语法表示

1) $S \rightarrow (0?1)^*$

2) $S \rightarrow aT | \epsilon, T \rightarrow Sb$

3) $S \rightarrow 0S1S | 1S0S | \epsilon$

思考

- 1) 用正则表达式可以定义任意正则语言吗？
- 2) 有穷自动机可以解析任意正则表达式吗？
- 3) 用CFG可以定义任意正则语言吗？
- 4) 用CFG可以定义任意上下文无关语言吗？
- 5) 用下推自动机可以解析任意正则表达式吗？
- 6) 用下推自动机可以解析任意CFG规则吗？
- 7) 用通用图灵机可以解析任意CFG规则吗？
- 8) 用通用图灵机可以解析任意程序吗？