

Design and Implementation of Automated Grammar Transformation and Parser Generators

文法自动变换及解析器生成器的设计与实现

Yihan Zhang

张艺瀚

Northeastern University

zhng1573@hotmail.com

2016 年 6 月 2 日

Overview

- 1 文法自动变换
 - 预备知识与记号
 - 迭代到不动点
- 2 解析器生成器
 - CYK 算法及其改进
 - Earley 算法及其改进
- 3 系统设计
- 4 系统实现
- 5 系统测试
- 6 结论

Definition

A context-free grammar(上下文无关文法) G is a quadruple(四元组) $G = (V, T, P, S)$, where V is nonterminal(variable, 非终结符号或变元) set, T is terminal(终结符号) set, P is production(产生式) set with productions shaped like $A \rightarrow \alpha$, $A \in V, \alpha \in (V \cup T)^*$, S is start variable(开始变元).

迭代到不动点

In this frame, we will introduce a technique called "iteration to a fixed point(迭代到不动点)". Take computing the generating symbols(可产生的符号) as an example. A symbol is said to be generating if $A \xRightarrow{*} \omega, \omega \in T^*$. In most textbooks, generating symbols are generally computed by the following rules:

- All terminals are generating;
- If $A \rightarrow \alpha \in P, \alpha = \varepsilon$ or every symbol in α is generating. Then A is generating.

Given any system of set equations, it's easy to get corresponding algorithm using the technique iteration to a fixed point. Algorithm for the last example is:

Algorithm 1 Compute generating variables

Require: $G = (V, T, P, S)$

Ensure: $V_{old} = \emptyset, V_{new} = \{A : A \rightarrow \omega \in P, \omega \in T^*\}$

- 1: **while** $V_{old} \neq V_{new}$ **do**
 - 2: $V_{old} = V_{new}$
 - 3: $V_{new} = V_{old} \cup \{A : A \rightarrow \alpha \in P, \alpha \in (T \cup V_{old})^*\}$
 - 4: **end while**
 - 5: $V = V_{new}$
-

原始 CYK 算法

Algorithm 2 Original CYK

Require: $G = (V, T, P, S)$ in CNF, $a_1 a_2 \cdots a_n \in T^n$

```
1: for  $i = 0$  to  $n - 1$  do
2:    $T_{i0} = \{A : A \rightarrow a_i \in P\}$ 
3: end for
4: for  $j = 1$  to  $n - 1$  do
5:   for  $i = 0$  to  $n - j - 1$  do
6:      $T_{ij} = \emptyset$ 
7:     for  $k = 0$  to  $j - 1$  do
8:        $T_{ij} = T_{ij} \cup \{A : A \rightarrow BC \in P, B \in T_{ik}, C \in T_{(i+k)(j-k)}\}$ 
9:     end for
10:  end for
11: end for
```

Remark

CYK algorithm requires the input grammar to be in CNF(乔姆斯基范式), whose productions are all in the form $A \rightarrow BC$, $B, C \in V$. Definition of T_{ij} in the algorithm is $T_{ij} = \{A : A \xRightarrow{} a_i a_{i+1} \cdots a_{i+j-2} a_{i+j-1}\}$. After running the algorithm, if $S \in T_{1n}$, then $a_1 a_2 \cdots a_n \in L(G)$.*

改进 CYK 算法

Algorithm 3 Improved CYK

Ensure: $M = \text{all false}$, $M[0, i, j] = \text{true}$, $i = 0, 1, \dots, n$, $V_j \rightarrow a_i \in P$

```
1: for  $i = 1$  to  $n - 1$  do
2:   for  $j = 0$  to  $n - i - 1$  do
3:     for  $k = 0$  to  $i - 1$  do
4:       for all  $V_A \text{ to } V_B V_C$  do
5:         if  $M[k, j, B], M[i - k - 1, j + k + 1, C]$  then
6:            $M[i, j, A] = \text{true}$ 
7:         end if
8:       end for
9:     end for
10:   end for
11: end for
```


Remark

After running the improved version, if $M[n-1, 0, S]$, then $V_S \in T_{1n}$. In fact, $V_k \in T[i, j] \iff M[i, j, k]$. Considering V is a finite set, we use a 3-dimensional boolean matrix(三维布尔矩阵) M to represent whether a variable is in an entry of the original 2-dimensional variable-set matrix(二维变元集合矩阵) T . This improvement avoids union-find operation of sets(集合的并查操作) which may have a high-complexity implementation and makes the algorithm more practical.

原始 Earley 算法

Given any augmented CFG(增广文法, with additional variable S' and production $S' \rightarrow S$) G and $a_1 a_2 \cdots a_n \in T^n$, Earley algorithm maintains a list of $n + 1$ sets and can be divided into two stages:

- ① Initialization(初始化): $S_0 = \{(S' \rightarrow \cdot S, 0)\}$;
- ② Assuming the current position is in the k -th slot of input string, perform the following operation until a fixed point:
 - Predict(预测) S_k ;
 - Scan(扫描) S_k ;
 - Complete(扩充) S_k .

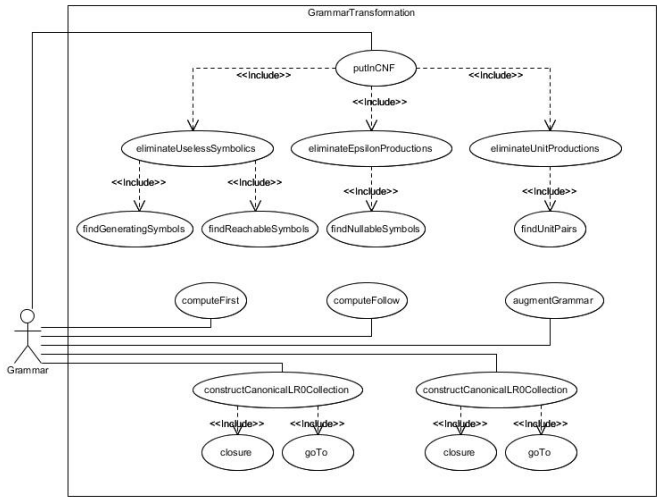
改进 Earley 算法

Algorithm 4 Earley

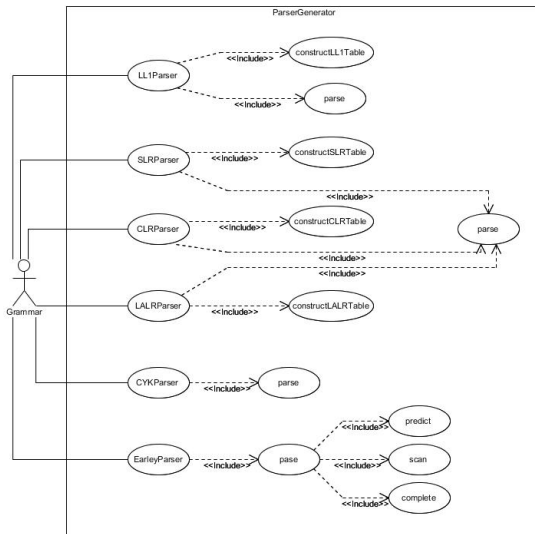
```
1: for  $k = 0$  to  $n$  do
2:   for all  $s \in S_k$  do
3:     if  $\cdot$  is at the end then
4:       if nonterminal after  $\cdot$  then
5:         Predict( $G, k, s$ )
6:       else
7:         Scan( $a_1 a_2 \cdots a_n, k, s$ )
8:       end if
9:     else
10:      Complete( $k, s$ )
11:    end if
12:  end for
```

Remark

The improved version modifies three operations to make them act on a single state rather than a state set and puts them into a wider loop, reducing repeated search(减少重复搜索).



图：文法自动变换用例图



图：解析器生成器用例图

Automated grammar transformer(simplify any given CFG and compute relevant sets automatically, 文法自动变换: 自动化简任意 CFG 并计算相关集合):

- Eliminate useless symbols(去除无用符号);
- Eliminate ε -productions(去除 ε -产生式);
- Eliminate unit production(去除单一产生式);
- Put in CNF(化为乔姆斯基范式).

Parser generators(generate parse table for any given CFG automatically, 解析器生成器: 自动生成任意 CFG 的语法分析表):

- CYK, Earley
- LL(1)
- SLR, CLR, LALR

- Hardware: ThinkPad-Edge-E530, 3.5GiB memory, Intel Core i5-3210M CPU, 2.50GHzx4, Intel Ivybridge Mobile;
- Operating system: Ubuntu 14.04 LTS, 64-bit;
- Tools: C++1y, Sublime Text 3, makefile;
- Environment requirement: g++ (Ubuntu 4.8.4-2ubuntu1 14.04.1),
Copyright (C) 2013 Free Software Foundation, Inc.

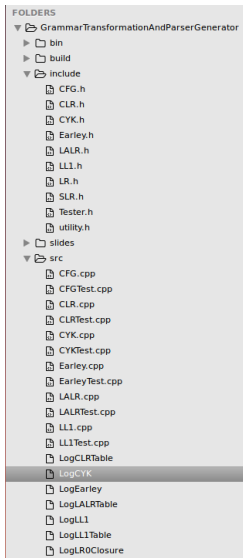


图: Directory structure

系统测试

```

1  CYKTable = {
2      (0, 0) { B, },
3      (0, 1) { A, C, },
4      (0, 2) { A, C, },
5      (0, 3) { B, },
6      (0, 4) { A, C, },
7      (1, 0) { S, A, },
8      (1, 1) { B, },
9      (1, 2) { S, C, },
10     (1, 3) { S, A, },
11     (2, 0) { },
12     (2, 1) { B, },
13     (2, 2) { B, },
14     (3, 0) { },
15     (3, 1) { S, A, C, },
16     (4, 0) { S, A, C, },
17 }
18
19 CYKTable = {
20     (0, 0) { NP, },
21     (0, 1) { VP, V, },
22     (0, 2) { Det, },
23     (0, 3) { N, },
24     (0, 4) { P, },
25     (0, 5) { Det, },
26     (0, 6) { N, },
27     (1, 0) { S, },
28     (1, 1) { },
29     (1, 2) { NP, },
30     (1, 3) { },
31     (1, 4) { },
32     (1, 5) { NP, },
33     (2, 0) { },
34     (2, 1) { VP, },
35     (2, 2) { },
36     (2, 3) { },
37     (2, 4) { PP, },
38     (3, 0) { S, },
39     (3, 1) { },
40     (3, 2) { },

```

(a) CYK log

1	$S = \{$
2	$(0, 0), (S \rightarrow S, 0)$
3	$(0, 1) (S \rightarrow S + M, 0)$
4	$(0, 2) (S \rightarrow M, 0)$
5	$(0, 3) (M \rightarrow M * T, 0)$
6	$(0, 4) (M \rightarrow T, 0)$
7	$(0, 5) (T \rightarrow 1, 0)$
8	$(0, 6) (T \rightarrow 2, 0)$
9	$(0, 7) (T \rightarrow 3, 0)$
10	$(0, 8) (T \rightarrow 4, 0)$
11	$(1, 0) (T \rightarrow 2, 0)$
12	$(1, 1) (M \rightarrow T, 0)$
13	$(1, 2) (S \rightarrow M, 0)$
14	$(1, 3) (M \rightarrow M * T, 0)$
15	$(1, 4) (S \rightarrow S, 0)$
16	$(1, 5) (S \rightarrow S + M, 0)$
17	$(2, 0) (S \rightarrow S + M, 0)$
18	$(2, 1) (M \rightarrow M * T, 2)$
19	$(2, 2) (M \rightarrow T, 2)$
20	$(2, 3) (T \rightarrow 1, 2)$
21	$(2, 4) (T \rightarrow 2, 2)$
22	$(2, 5) (T \rightarrow 3, 2)$
23	$(2, 6) (T \rightarrow 4, 2)$
24	$(3, 0) (T \rightarrow 3, 2)$
25	$(3, 1) (M \rightarrow T, 2)$
26	$(3, 2) (S \rightarrow S + M, 0)$
27	$(3, 3) (M \rightarrow M * T, 2)$
28	$(3, 4) (S \rightarrow S, 0)$
29	$(3, 5) (S \rightarrow S + M, 0)$
30	$(4, 0) (M \rightarrow M * T, 2)$
31	$(4, 1) (T \rightarrow 1, 4)$
32	$(4, 2) (T \rightarrow 2, 4)$
33	$(4, 3) (T \rightarrow 3, 4)$
34	$(4, 4) (T \rightarrow 4, 4)$
35	$(5, 0) (T \rightarrow 4, 4)$
36	$(5, 1) (M \rightarrow M * T, 2)$
37	$(5, 2) (S \rightarrow S + M, 0)$
38	$(5, 3) (M \rightarrow M * T, 2)$
39	$(5, 4) (S \rightarrow S, 0)$
40	$(5, 5) (S \rightarrow S + M, 0)$

(b) Earley log

```

1 LL1Table = {
2   (E, +) error,
3   (E, *) error,
4   (E, () E -> T E_ ,
5   (E, )) error,
6   (E, 0) E -> T E_ ,
7   (E, $) error,
8   (T, +) error,
9   (T, *) error,
10  (T, () T -> F T_ ,
11  (T, )) error,
12  (T, 0) T -> F T_ ,
13  (T, $) error,
14  (E_, +) E_ -> + T E_ ,
15  (E_, *) error,
16  (E_, () error,
17  (E_, )) E_ -> ,
18  (E_, 0) error,
19  (E_, $) E_ -> ,
20  (F, +) error,
21  (F, *) error,
22  (F, () F -> ( E ) ,
23  (F, )) error,
24  (F, 0) F -> 0 ,
25  (F, $) error,
26  (T_, +) T_ -> ,
27  (T_, *) T_ -> * F T_ ,
28  (T_, () error,
29  (T_, )) T_ -> ,
30  (T_, 0) error,
31  (T_, $) T_ -> ,
32 }
33

```

(a) LL(1) table log

```

1 |stk: { E, $, }
2 x: E
3 a: 0
4 E -> T E_
5
6 |stk: { T, E_, $, }
7 x: T
8 a: 0
9 T -> F T_
10
11 |stk: { F, T_, E_, $, }
12 x: F
13 a: 0
14 F -> 0
15
16 |stk: { 0, T_, E_, $, }
17 x: 0
18 a: 0
19 match 0
20
21 |stk: { T_, E_, $, }
22 x: T_
23 a: +
24 T_ ->
25
26 |stk: { E_, $, }
27 x: E_
28 a: +
29 E_ -> + T E_
30
31 |stk: { +, T, E_, $, }
32 x: +
33 a: +
34 match +
35
36 |stk: { T, E_, $, }
37 x: T
38 a: 0
39 T -> F T_
40

```

(b) LL(1) parsing log

```

1 |canonicalLR0Collection = {
2   I0 = {
3     E -> . E ,
4     E -> . E + T ,
5     E -> . T ,
6     T -> . T * F ,
7     T -> . F ,
8     F -> . ( E ) ,
9     F -> . 0 ,
10  }
11
12  I1 = {
13    E -> E . ,
14    E -> E . + T ,
15  }
16
17  I2 = {
18    E -> T . ,
19    T -> T . * F ,
20  }
21
22  I3 = {
23    T -> F . ,
24  }
25
26  I4 = {
27    F -> ( . E ) ,
28    E -> . E + T ,
29    E -> . T ,
30    T -> . T * F ,
31    T -> . F ,
32    F -> . ( E ) ,
33    F -> . 0 ,
34  }
35
36  I5 = {
37    F -> 0 . ,
38  }
39
40  I6 = {

```

(a) LR(0) collection log

```

1 |action = {
2   (0, +) error,
3   (0, *) error,
4   (0, () shift 4,
5   (0, )) error,
6   (0, 0) shift 5,
7   (0, $) error,
8   (1, +) shift 6,
9   (1, *) error,
10  (1, () error,
11  (1, )) error,
12  (1, 0) error,
13  (1, $) accept,
14  (2, +) reduce E -> T ,
15  (2, *) shift 7,
16  (2, () error,
17  (2, )) reduce E -> T ,
18  (2, 0) error,
19  (2, $) reduce E -> T ,
20  (3, +) reduce T -> F ,
21  (3, *) reduce T -> F ,
22  (3, () error,
23  (3, )) reduce T -> F ,
24  (3, 0) error,
25  (3, $) reduce T -> F ,
26  (4, +) error,
27  (4, *) error,
28  (4, () shift 4,
29  (4, )) error,
30  (4, 0) shift 5,
31  (4, $) error,
32  (5, +) reduce F -> 0 ,
33  (5, *) reduce F -> 0 ,
34  (5, () error,
35  (5, )) reduce F -> 0 ,
36  (5, 0) error,
37  (5, $) reduce F -> 0 ,
38  (6, +) error,
39  (6, *) error,
40  (6, () shift 4,

```

(b) SLR table log

```

1 canonicalLR1Collection = {
2     I0 = {
3         S -> . S , $ ,
4         S -> . C C , $ ,
5         C -> . C C , c ,
6         C -> . c C , d ,
7         C -> . d , c ,
8         C -> . d , d ,
9     }
10
11     I1 = {
12         _S -> S . , $ ,
13     }
14
15     I2 = {
16         S -> C . C , $ ,
17         C -> . c C , $ ,
18         C -> . d , $ ,
19     }
20
21     I3 = {
22         C -> c . C , c ,
23         C -> c . C , d ,
24         C -> . c C , c ,
25         C -> . d , c ,
26         C -> . c C , d ,
27         C -> . d , d ,
28     }
29
30     I4 = {
31         C -> d . , c ,
32         C -> d . , d ,
33     }
34
35     I5 = {
36         S -> C C . , $ ,
37     }
38
39     I6 = {
40         C -> c . C , $ ,

```

(a) LR(1) collection log

```

1 action = {
2     (0, c) shift 3,
3     (0, d) shift 4,
4     (0, $) error,
5     (1, c) error,
6     (1, d) error,
7     (1, $) accept,
8     (2, c) shift 6,
9     (2, d) shift 7,
10    (2, $) error,
11    (3, c) shift 3,
12    (3, d) shift 4,
13    (3, $) error,
14    (4, c) reduce C -> d ,
15    (4, d) reduce C -> d ,
16    (4, $) error,
17    (5, c) error,
18    (5, d) error,
19    (5, $) reduce S -> C C ,
20    (6, c) shift 6,
21    (6, d) shift 7,
22    (6, $) error,
23    (7, c) error,
24    (7, d) error,
25    (7, $) reduce C -> d ,
26    (8, c) reduce C -> c C ,
27    (8, d) reduce C -> c C ,
28    (8, $) error,
29    (9, c) error,
30    (9, d) error,
31    (9, $) reduce C -> c C ,
32 }
33
34 goTo = {
35     (0, S) 1,
36     (0, C) 2,
37     (1, S) error,
38     (1, C) error,
39     (2, S) error,
40     (2, C) 5,

```

(b) CLR table log

```

1 action = {
2     (0, c) shift 3,
3     (0, d) shift 4,
4     (0, $) error,
5     (1, c) error,
6     (1, d) error,
7     (1, $) accept,
8     (2, c) shift 3,
9     (2, d) shift 4,
10    (2, $) error,
11    (3, c) shift 3,
12    (3, d) shift 4,
13    (3, $) error,
14    (4, c) reduce C -> d ,
15    (4, d) reduce C -> d ,
16    (4, $) reduce C -> d ,
17    (5, c) error,
18    (5, d) error,
19    (5, $) reduce S -> C C ,
20    (6, c) reduce C -> c C ,
21    (6, d) reduce C -> c C ,
22    (6, $) reduce C -> c C ,
23 }
24
25 goTo = {
26     (0, S) 1,
27     (0, C) 2,
28     (1, S) error,
29     (1, C) error,
30     (2, S) error,
31     (2, C) 5,
32     (3, S) error,
33     (3, C) 6,
34     (4, S) error,
35     (4, C) error,
36     (5, S) error,
37     (5, C) error,
38     (6, S) error,
39     (6, C) error,
40 }

```

(c) LALR table log

结论

We design a C-like grammar(类 C 文法) which supports most basic syntax of C and use (part of) it to test the system. Conclusions are shown below:

- Any grammar which $L(G) \setminus \{\varepsilon\} \neq \emptyset$ can be put in CNF;
- Complexity of CYK and Earley are both $O(n^3)$, while $O(n^2)$ for Earley when handling unambiguous grammars(无二义性文法);
- Construction of parser generators(解析器生成器) is equivalent to construction of parse tables(语法分析表);
- LL(1) parser is equivalent to recursive descent parser without backtrack(无回溯递归下降语法分析器);
- Range of applicable grammars: $SLR \subset LALR \subset CLR$, scale of generated parse tables: $SLR = LALR < CLR$.

Q & A