

Лекции по теории формальных языков

Лекция 13.

LR(0)-, SLR(1)-, LR(1)- и LALR(1)-анализ

Александр Сергеевич Герасимов

<http://gas-teach.narod.ru>

Кафедра математических и информационных технологий
Санкт-Петербургского академического университета
Российской академии наук.

Весенний семестр 2010/11 учебного года

13 мая 2011 г.

План

- 1 LR(0)-анализ
- 2 SLR(1)-анализ
- 3 LR(1)-анализ
- 4 LALR(1)-анализ

План

- 1 LR(0)-анализ
- 2 SLR(1)-анализ
- 3 LR(1)-анализ
- 4 LALR(1)-анализ

Анализ с помощью LR(0)-автомата: неформальное описание

- Подадим цепочку w на вход LR(0)-автомату \mathcal{A}_G грамматики G .
- После каждого такта работы автомата \mathcal{A}_G кладём новое состояние автомата в стек и, анализируя это состояние, до перехода к следующему такту производим описанные ниже действия.
- Если в состоянии есть пункт $[A \rightarrow \alpha_1 \cdot \alpha_2]$ и $\alpha_2 \neq \varepsilon$, то выполняем перенос следующего символа, т. е. переходим к следующему такту работы автомата \mathcal{A}_G .
- Если в состоянии есть пункт $[A \rightarrow \alpha \cdot]$, то выполняем свёртку по правилу $A \rightarrow \alpha$ (при этом активный префикс $\gamma\alpha$ в стеке заменяется на активный префикс γA):
 - ▶ выталкиваем из стека $|\alpha|$ состояний;
 - ▶ на вершине стека оказывается состояние I после прочтения γ ;
 - ▶ вернём автомат в состояние I и дадим прочитать ему на следующем такте A .
- Если таким образом автомат прочтёт w и придёт в состоянии, содержащее пункт $[S' \rightarrow S \cdot]$, то $w \in L(G)$, иначе $w \notin L(G)$.

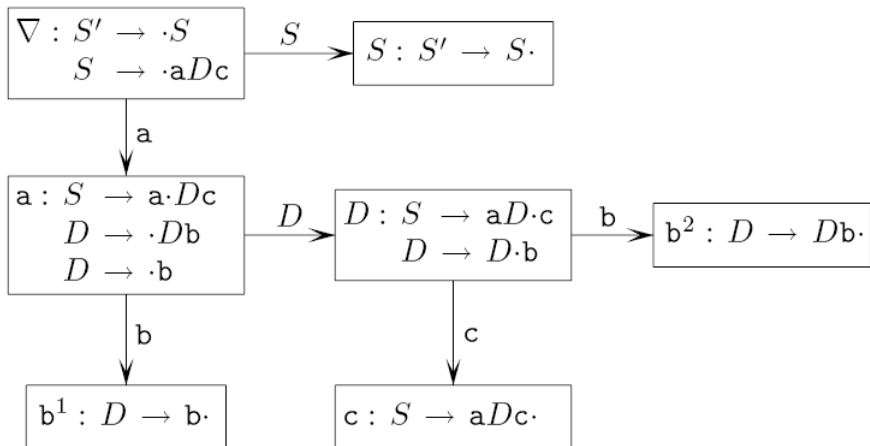
Конфликты

- *Конфликт перенос-свёртка*: в одном состоянии LR(0)-автомата имеются пункты
 - ▶ $[A \rightarrow \alpha_1 \cdot \alpha_2]$, где $\alpha_2 \neq \varepsilon$, и
 - ▶ $[B \rightarrow \beta \cdot]$.
- *Конфликт свёртка-свёртка*: в одном состоянии LR(0)-автомата имеются различные пункты
 - ▶ $[A \rightarrow \alpha \cdot]$ и
 - ▶ $[B \rightarrow \beta \cdot]$.

LR(0)-грамматика: определение и пример

Грамматика называется *LR(0)-грамматикой*, если каждое состояние её LR(0)-автомата, содержащее пункт вида $[A \rightarrow \alpha \cdot]$, состоит из единственного пункта.

Пример. $G_1 : (0) S' \rightarrow S, (1) S \rightarrow aDc, (2) D \rightarrow Db, (3) D \rightarrow b.$



Алгоритм построения LR(0)-анализатора

Вход. Расширенная LR(0)-грамматика $G = (\Sigma, \Gamma, P, S')$.

Выход. Таблица LR-анализа для грамматики G .

1. построить LR(0)-автомат $\mathcal{A}_G = (Q, \Sigma \cup \Gamma, \delta, I_0, Q)$;
2. для каждого $I \in Q$
3. для каждого $i \in I$
4. если $(i = [S' \rightarrow S \cdot])$ ACTION(I, \cdot) := \checkmark ;
5. иначе если $(i = [A \rightarrow \alpha \cdot])$
6. для каждого $a \in \Sigma \cup \{\cdot\}$
7. ACTION(I, a) := (\otimes номер($A \rightarrow \alpha$));
8. иначе если $(i = [A \rightarrow \beta_1 \cdot a \beta_2])$
9. ACTION(I, a) := ($\leftarrow \delta(I, a)$);
10. иначе (т. е. $i = [A \rightarrow \beta_1 \cdot B \beta_2]$)
11. GOTO(I, B) := $\delta(I, B)$

Почему в каждой клетке построенной таблицы не более одной записи?

Пример построения LR(0)-анализатора

$G_1 : (0) S' \rightarrow S, (1) S \rightarrow aDc, (2) D \rightarrow Db, (3) D \rightarrow b.$

| | ACTION | | | | GOTO | |
|----------|----------------|------------------|----------------|-------------|------|---|
| | a | b | c | \dagger | S | D |
| S | | | | ✓ | | |
| D | | $\leftarrow b^2$ | $\leftarrow c$ | | | |
| a | | $\leftarrow b^1$ | | | | D |
| b^1 | $\otimes 3$ | $\otimes 3$ | $\otimes 3$ | $\otimes 3$ | | |
| b^2 | $\otimes 2$ | $\otimes 2$ | $\otimes 2$ | $\otimes 2$ | | |
| c | $\otimes 1$ | $\otimes 1$ | $\otimes 1$ | $\otimes 1$ | | |
| ∇ | $\leftarrow a$ | | | | S | |

Предложение о корректности алгоритма построения LR(0)-анализатора

Предложение

LR(0)-анализатор, построенный алгоритмом на слайде 7 по грамматике G , допускает цепочку w тогда и только тогда, когда $w \in L(G)$.

Доказательство.

- Этот анализатор работает как LR(0)-автомат с запуском из неочередного состояния в случае свёртки и складыванием в состояний стек.
- Анализатор выполняет свёртку в том и только в том случае, когда он находится в состоянии I вида $\{[A \rightarrow \alpha \cdot]\}$.
- Если LR(0)-автомат находится в состоянии I , а в стеке — $\gamma = \gamma' \alpha$, то по основной теореме LR-анализа γ является активным префиксом, для которого допустим пункт $[A \rightarrow \alpha \cdot]$, а α является основой некоторой r -формы, начинающейся с γ .

Предложение о корректности алгоритма построения LR(0)-анализатора: окончание доказательства

- В состоянии I имеется единственный пункт $[A \rightarrow \alpha \cdot]$, поэтому α является основой любой r -формы, начинающейся с γ .
- Таким образом, анализатор выполняет свёртку тогда и только тогда, когда наверху стека находится основа текущей r -формы.
- Иначе анализатор пытается осуществить перенос. Если конкатенация содержимого стека и необработанного суффикса входной цепочки является r -формой, то такой перенос возможен, поскольку LR(0)-автомат распознаёт все активные префиксы.
- Наконец, допуск производится тогда и только тогда, когда входная цепочка свёрнута в аксиому грамматики.



Теорема об оценке числа шагов LR(0)-анализатора

Теорема

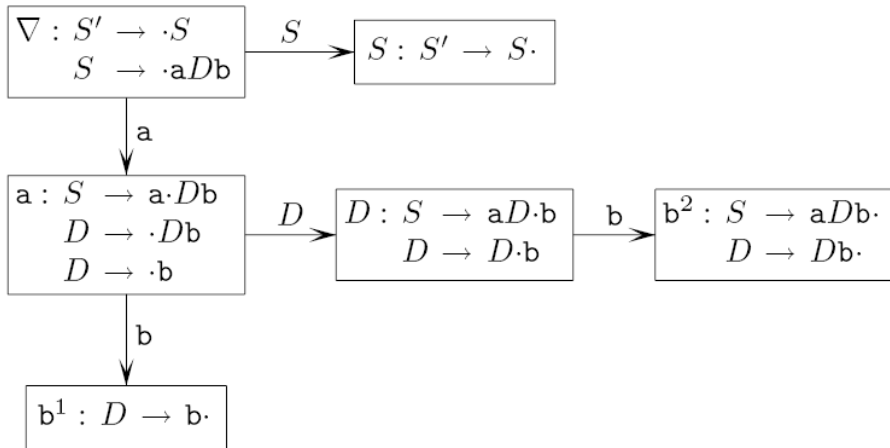
Число шагов (переносов и свёрток) LR(0)-анализатора, построенного алгоритмом на слайде 7, линейно зависит от длины входной цепочки.

План

- 1 LR(0)-анализ
- 2 SLR(1)-анализ**
- 3 LR(1)-анализ
- 4 LALR(1)-анализ

Пример конфликта «свёртка-свёртка»

$G_2 : (0) S' \rightarrow S, (1) S \rightarrow aDb, (2) D \rightarrow Db, (3) D \rightarrow b.$



Разрешение конфликта:

- если очередной символ \dagger , то $\otimes 1$;
- если очередной символ b , то $\otimes 2$.

Пример конфликта «перенос-свёртка»

$G_3 : (0) S' \rightarrow S, (1) S \rightarrow ac, (2) S \rightarrow aDbc, (3) D \rightarrow Db, (4) D \rightarrow \varepsilon.$

$$\nabla = \{[S' \rightarrow \cdot S], [S \rightarrow \cdot ac], [S \rightarrow \cdot aDbc]\} \xrightarrow{S} S = \{[S' \rightarrow S \cdot]\}$$

$\downarrow a$

$$a = \{[S \rightarrow a \cdot c], [S \rightarrow a \cdot Dbc], [D \rightarrow \cdot Db], [D \rightarrow \cdot]\} \xrightarrow{c} c^1 = \{[S \rightarrow ac \cdot]\}$$

$\downarrow D$

$$D = \{[S \rightarrow aD \cdot bc], [D \rightarrow D \cdot b]\} \xrightarrow{b} b = \{[S \rightarrow aDb \cdot c], [D \rightarrow Db \cdot]\} \xrightarrow{c} c^2 = \{[S \rightarrow aDbc \cdot]\}$$

Разрешение конфликтов:

- если очередной символ c , то перенос;
- если очередной символ b , то свёртка.

Разрешение некоторых конфликтов при помощи множеств FOLLOW(A)

- Свёртка по правилу $A \rightarrow \alpha$ является правильным действием LR-анализатора, только если очередной входной символ принадлежит FOLLOW(A).
- Пусть пункт $[A \rightarrow \alpha \cdot]$ ($A \neq S'$) принадлежит состоянию I LR(0)-автомата.
- Тогда достаточно положить $ACTION(I, a) := (\otimes \text{номер}(A \rightarrow \alpha))$ лишь для всех $a \in FOLLOW(A)$.

Алгоритм построения SLR(1)-анализатора

Вход. Расширенная грамматика $G = (\Sigma, \Gamma, P, S')$.

Выход. Таблица LR-анализа для грамматики G .

1. построить LR(0)-автомат $\mathcal{A}_G = (Q, \Sigma \cup \Gamma, \delta, I_0, Q)$;
2. построить FOLLOW(A) для каждого $A \in \Gamma$;
3. для каждого $I \in Q$
4. **для** каждого $i \in I$
5. **если** ($i = [S' \rightarrow S \cdot]$) ACTION(I, \cdot) := \checkmark ;
6. **иначе если** ($i = [A \rightarrow \alpha \cdot]$)
7. **для** каждого $a \in \text{FOLLOW}(A)$
8. ACTION(I, a) := (\otimes номер($A \rightarrow \alpha$));
9. **иначе если** ($i = [A \rightarrow \beta_1 \cdot a \beta_2]$)
10. ACTION(I, a) := ($\leftarrow \delta(I, a)$);
11. **иначе** (т. е. $i = [A \rightarrow \beta_1 \cdot B \beta_2]$)
12. GOTO(I, B) := $\delta(I, B)$

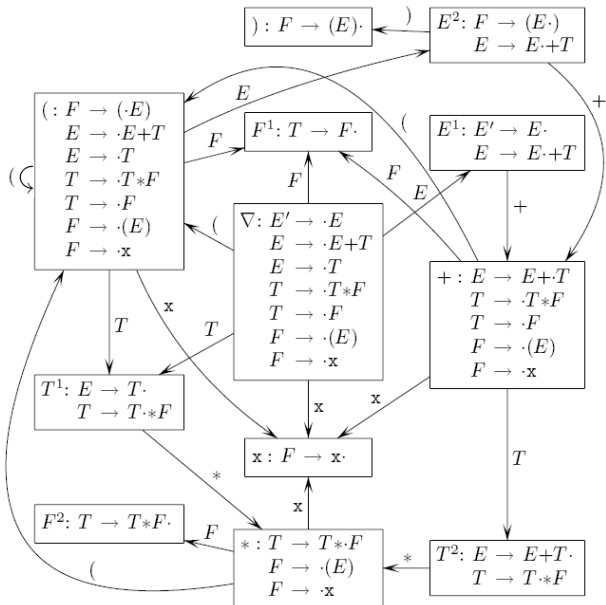
SLR(1)-грамматика: определение и пример

Грамматика называется *SLR(1)-грамматикой*, если в таблице ACTION, построенной алгоритмом на предыдущем слайде, нет конфликтов (т. е. в каждой клетке не более одной записи).

Пример. $G_2 : (0) S' \rightarrow S, (1) S \rightarrow aDb, (2) D \rightarrow Db, (3) D \rightarrow b.$
 $FOLLOW(S) = \{ \vdash \}, FOLLOW(D) = \{ b \}.$

| | ACTION | | | GOTO | |
|----------|----------------|------------------|-------------|------|---|
| | a | b | \vdash | S | D |
| S | | | ✓ | | |
| D | | $\leftarrow b^2$ | | | |
| a | | $\leftarrow b^1$ | | | D |
| b^1 | | $\otimes 3$ | | | |
| b^2 | | $\otimes 2$ | $\otimes 1$ | | |
| ∇ | $\leftarrow a$ | | | S | |

- (0) $E' \rightarrow E$, (1) $E \rightarrow E + T$, (2) $E \rightarrow T$, (3) $T \rightarrow T * F$,
 (4) $T \rightarrow F$, (5) $F \rightarrow (E)$, (6) $F \rightarrow x$.



Пример построения SLR(1)-анализатора: окончание

$FOLLOW(E) = \{+,), \vdash\}$, $FOLLOW(T) = FOLLOW(F) = \{*, +,), \vdash\}$.

| | ACTION | | | | | | GOTO | | |
|----------|----------------|----------------|----------------|----------------|----------------|--------------|-------|-------|-------|
| | + | * | x | (|) | \vdash | E | T | F |
| E^1 | $\leftarrow +$ | | | | | \checkmark | | | |
| T^1 | $\otimes 2$ | $\leftarrow *$ | | | $\otimes 2$ | $\otimes 2$ | | | |
| F^1 | $\otimes 4$ | $\otimes 4$ | | | $\otimes 4$ | $\otimes 4$ | | | |
| (| | | $\leftarrow x$ | $\leftarrow ($ | | | E^2 | T^1 | F^1 |
| x | $\otimes 6$ | $\otimes 6$ | | | $\otimes 6$ | $\otimes 6$ | | | |
| + | | | $\leftarrow x$ | $\leftarrow ($ | | | | T^2 | F^1 |
| * | | | $\leftarrow x$ | $\leftarrow ($ | | | | | F^2 |
| E^2 | $\leftarrow +$ | | | | $\leftarrow)$ | | | | |
| T^2 | $\otimes 1$ | $\leftarrow *$ | | | $\otimes 1$ | $\otimes 1$ | | | |
| F^2 | $\otimes 3$ | $\otimes 3$ | | | $\otimes 3$ | $\otimes 3$ | | | |
|) | $\otimes 5$ | $\otimes 5$ | | | $\otimes 5$ | $\otimes 5$ | | | |
| ∇ | | | $\leftarrow x$ | $\leftarrow ($ | | | E^1 | T^1 | F^1 |

План

- 1 LR(0)-анализ
- 2 SLR(1)-анализ
- 3 LR(1)-анализ**
- 4 LALR(1)-анализ

LR(1)-пункты

Пусть $G = (\Sigma, \Gamma, P, S')$ — расширенная грамматика.

- LR(1)-пункт: $[A \rightarrow \beta_1 \cdot \beta_2, a]$, где
 - ▶ $[A \rightarrow \beta_1 \cdot \beta_2]$ — LR(0)-пункт, называемый *ядром* этого LR(1)-пункта,
 - ▶ a — терминал или символ \dagger .
- LR(1)-пункт $[A \rightarrow \beta_1 \cdot \beta_2, a]$ допустим для активного префикса $\alpha\beta_1$, если существует (правый) вывод

$$S' \Rightarrow^* \alpha A w \Rightarrow \alpha \beta_1 \beta_2 w \Rightarrow^* u w$$

и цепочка $w \dagger$ начинается символом a .

- Автоматом LR(1)-пунктов грамматики G назовём ε -НКА

$$\mathcal{I}_G^{(1)} = (Q, \Sigma \cup \Gamma, \delta, \{i_0\}, Q), \text{ где}$$

- ▶ Q — множество всех LR(1)-пунктов грамматики G ,
 - ▶ $i_0 = [S' \rightarrow \cdot S, \dagger]$,
 - ▶ отношение переходов состоит из всех базисных переходов вида $([A \rightarrow \beta_1 \cdot X \beta_2, a], X, [A \rightarrow \beta_1 X \cdot \beta_2, a])$ и всех ε -переходов вида $([A \rightarrow \beta_1 \cdot B \beta_2, a], \varepsilon, [B \rightarrow \cdot \beta, b])$, где $b \in \text{FIRST}(\beta_2 a)$ и при построении $\text{FIRST}(\beta_2 a)$ мы считаем \dagger терминалом.
- Для $\mathcal{I}_G^{(1)}$ верны основная теорема LR-анализа и её следствия.

LR(1)-автомат

- ДКА (неполный) $\mathcal{A}_G^{(1)}$, построенный по ε -НКА $\mathcal{I}_G^{(1)}$ алгоритмом из лекции 2, будем называть LR(1)-автоматом грамматики G .

Пусть M — произвольное множество LR(1)-пунктов расширенной грамматики G .

- Положим, что $\text{CLOSURE}_1(M)$ есть минимальное по включению множество LR(1)-пунктов такое, что

- ▶ $M \subseteq \text{CLOSURE}_1(M)$ и
- ▶ $[A \rightarrow \beta_1 \cdot B\beta_2, a] \in \text{CLOSURE}_1(M)$ влечёт $[B \rightarrow \cdot \gamma, b] \in \text{CLOSURE}_1(M)$ для каждого правила вывода вида $B \rightarrow \gamma$ и каждого $b \in \text{FIRST}(\beta_2 a)$.

- Пусть X — терминал или нетерминал грамматики G . Положим

$$\text{GOTO}_1(M, X) = \text{CLOSURE}_1(\{[A \rightarrow \beta_1 X \cdot \beta_2, a] \mid [A \rightarrow \beta_1 \cdot X \beta_2, a] \in M\})$$

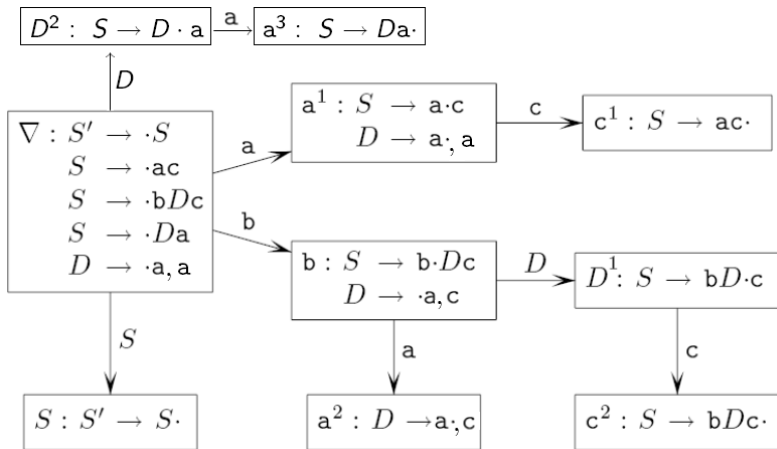
- Алгоритм построения LR(1)-автомата.

Пример построения LR(1)-автомата

G_4 : (0) $S' \rightarrow S$, (1) $S \rightarrow ac$, (2) $S \rightarrow bDc$, (3) $S \rightarrow Da$, (4) $D \rightarrow a$.

Переходы в $\mathcal{I}_{G_4}^{(1)}$ имеют вид $([A \rightarrow \beta_1 \cdot X\beta_2, a], X, [A \rightarrow \beta_1 X \cdot \beta_2, a])$ и $([A \rightarrow \beta_1 \cdot B\beta_2, a], \varepsilon, [B \rightarrow \cdot \beta, b])$, где $b \in \text{FIRST}(\beta_2 a)$.

На этом рисунке опущенный второй компонент в LR(1)-пунктах есть \perp .



Алгоритм построения LR(1)-анализатора.

Определение LR(1)-грамматики

Вход. Расширенная грамматика $G = (\Sigma, \Gamma, P, S')$.

Выход. Таблица LR-анализа для грамматики G .

1. построить LR(1)-автомат $\mathcal{A}_G^{(1)} = (Q, \Sigma \cup \Gamma, \delta, I_0, Q)$;
2. для каждого $I \in Q$
3. для каждого $i \in I$
4. если $(i = [S' \rightarrow S \cdot, \neg])$ ACTION(I, \neg) := \checkmark ;
5. иначе если $(i = [A \rightarrow \alpha \cdot, a])$
6. ACTION(I, a) := (\otimes номер($A \rightarrow \alpha$));
7. иначе если $(i = [A \rightarrow \beta_1 \cdot a \beta_2, c])$
8. ACTION(I, a) := ($\leftarrow \delta(I, a)$);
9. иначе (т. е. $i = [A \rightarrow \beta_1 \cdot B \beta_2, c]$)
10. GOTO(I, B) := $\delta(I, B)$

Грамматика называется *LR(1)-грамматикой*, если в построенной этим алгоритмом таблице ACTION нет конфликтов (т. е. в каждой клетке не более одной записи).

Пример построения LR(1)-анализатора

$G_4 : (0) S' \rightarrow S, (1) S \rightarrow ac, (2) S \rightarrow bDc, (3) S \rightarrow Da, (4) D \rightarrow a.$

| | ACTION | | | | GOTO | |
|----------|------------------|----------------|------------------|-------------|------|-------|
| | a | b | c | ⊥ | S | D |
| S | | | | ✓ | | |
| D^1 | | | $\leftarrow c^2$ | | | |
| D^2 | $\leftarrow a^3$ | | | | | |
| a^1 | $\otimes 4$ | | $\leftarrow c^1$ | | | |
| a^2 | | | $\otimes 4$ | | | |
| a^3 | | | | $\otimes 3$ | | |
| b | $\leftarrow a^2$ | | | | | D^1 |
| c^1 | | | | $\otimes 1$ | | |
| c^2 | | | | $\otimes 2$ | | |
| ∇ | $\leftarrow a^1$ | $\leftarrow b$ | | | S | D^2 |

План

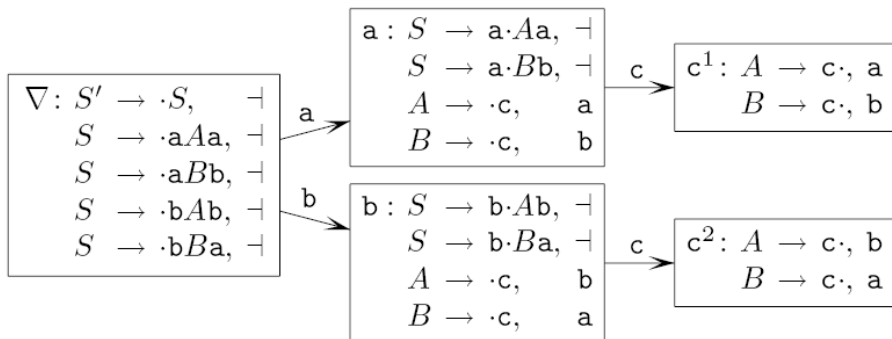
- 1 LR(0)-анализ
- 2 SLR(1)-анализ
- 3 LR(1)-анализ
- 4 LALR(1)-анализ

LARL(1)-анализ

- Для языка программирования Паскаль LR(1)-автомат имеет несколько тысяч состояний, тогда как LR(0)-автомат — несколько сотен состояний.
- LALR — Look Ahead LR, LR-анализ с предпросмотром.
- *LALR(1)-автомат* грамматики получается из LR(1)-автомата этой грамматики объединением всех состояний, имеющих одно и то же множество ядер LR(1)-пунктов, в одно состояние.
- Для LR(1)-грамматики LALR(1)-автомат не имеет конфликтов перенос-свёртка.
 - ▶ Пусть $[A \rightarrow \alpha \cdot, a]$ и $[B \rightarrow \beta_1 \cdot a \beta_2, b]$ принадлежат одному состоянию LALR(1)-автомата.
 - ▶ Тогда $[A \rightarrow \alpha \cdot, a]$ и $[B \rightarrow \beta_1 \cdot a \beta_2, c]$ (для некоторого c) принадлежат одному состоянию LR(1)-автомата.(Но конфликты свёртка-свёртка могут появиться.)
- Грамматика называется *LALR(1)-грамматикой*, если нет конфликтов в таблице ACTION, построенной алгоритмом на слайде 25, но по LALR(1)-автомату этой грамматики.

Пример LR(1)-, но не LALR(1)-грамматики

$G_5 : S' \rightarrow S, S \rightarrow aAa|aBb|bAb|bBa, A \rightarrow c, B \rightarrow c.$



$$c^{1,2} = c^1 \cup c^2 = \{[A \rightarrow c \cdot, a/b], [B \rightarrow c \cdot, a/b]\}.$$

($[A \rightarrow \alpha \cdot \beta, a/b]$ — сокращение $[A \rightarrow \alpha \cdot \beta, a], [A \rightarrow \alpha \cdot \beta, b].$)

Литература

Основная литература

- Замятин А. П., Шур А. М. Языки, грамматики, распознаватели: Учебное пособие. Екатеринбург : Изд-во Урал. ун-та, 2007 (электронный вариант книги — на <http://elar.usu.ru>, поиск).

Дополнительная литература

- Ахо А., Лам М., Сети Р., Ульман Дж. Компиляторы: принципы, технологии и инструментарий. М.: ООО "И.Д. Вильямс", 2008.
- Ахо А., Ульман Дж. Теория синтаксического анализа, перевода и компиляции. М.: Мир, 1978.
- Мартыненко Б. К. Языки и трансляции: Учеб. пособие. СПб.: Издательство С.-Петербургского университета, 2004 (электронный вариант книги — на <http://www.math.spbu.ru/user/mbk>).