

Обобщённый синтаксический анализ

Generalized parsing

Автор: Григорьев Семён Вячеславович
rsdpisuy@gmail.com

Математико-Механический факультет
Санкт-Петербургского государственного университета

21 ноября 2013г.

- Управляющие таблицы содержат не более одной инструкции в ячейке.
 - ▶ Можно считать, что ровно одно, если все пустые ячейки заполнить инструкцией *ERROR*
- А если в процессе построения таблицы невозможно однозначно определить очередное действие?
 - ▶ Конфликты
 - ★ Shift-Reduce
 - ★ Reduce-Reduce

- В большинстве практических ситуаций нам нужен единственный вывод данной строки в данной грамматике.

Конфликты на практике

- В большинстве практических ситуаций нам нужен единственный вывод данной строки в данной грамматике.
- Или факт его отсутствия.

Конфликты на практике

- Конфликты приводят к недетерминированному поведению распознавателя.

Конфликты на практике

- Конфликты приводят к недетерминированному поведению распознавателя.
- Для одной входной цепочки можно получить несколько выводов.

Конфликты на практике

- Конфликты приводят к недетерминированному поведению распознавателя.
- Для одной входной цепочки можно получить несколько выводов.
- А на сколько это плохо?

Конфликты на практике

- Конфликты приводят к недетерминированному поведению распознавателя.
- Для одной входной цепочки можно получить несколько выводов.
- А на сколько это плохо?
- А всегда ли нужен ровно один вывод?

Конфликты на практике

- Конфликты приводят к недетерминированному поведению распознавателя.
- Для одной входной цепочки можно получить несколько выводов.
- А на сколько это плохо?
- А всегда ли нужен ровно один вывод?
- Какой из них “правильный”?

Definition (Неоднозначная грамматика)

Будем называть формальную грамматику неоднозначной, если существует входная последовательность для которой может быть найдено более одного левостороннего вывода в этой грамматике.

- “Аккуратный” способ.
 - ▶ Преобразования грамматики.
 - ▶ В некоторых случаях можно преобразовать грамматику языка так, чтобы она не содержала неоднозначностей.

Разрешение неоднозначностей

- “Аккуратный” способ.
 - ▶ Преобразования грамматики.
 - ▶ В некоторых случаях можно преобразовать грамматику языка так, чтобы она не содержала неоднозначностей.
- Практический способ.
 - ▶ Автоматический детерминированный выбор одного из вариантов при построении таблиц.

Разрешение неоднозначностей

- “Аккуратный” способ.
 - ▶ Преобразования грамматики.
 - ▶ В некоторых случаях можно преобразовать грамматику языка так, чтобы она не содержала неоднозначностей.
- Практический способ.
 - ▶ Автоматический детерминированный выбор одного из вариантов при построении таблиц.
 - ▶ Shift-Reduce → Shift
 - ▶ Reduce-Reduce → “magic first” Reduce

Пример преобразования грамматики

Оригинальная
грамматика G_1

① $S_s \rightarrow E\$$

② $E \rightarrow E + E$

③ $E \rightarrow E * E$

④ $E \rightarrow d$

Преобразованная
грамматика G_2

① $S_s \rightarrow S\$$

② $S \rightarrow S + T$

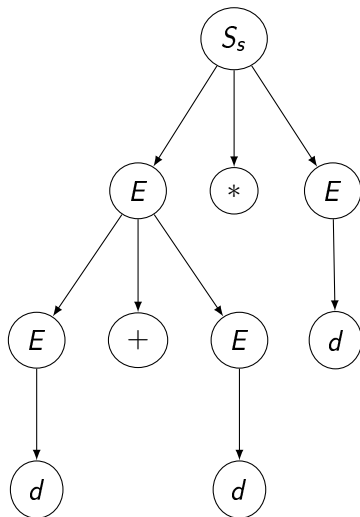
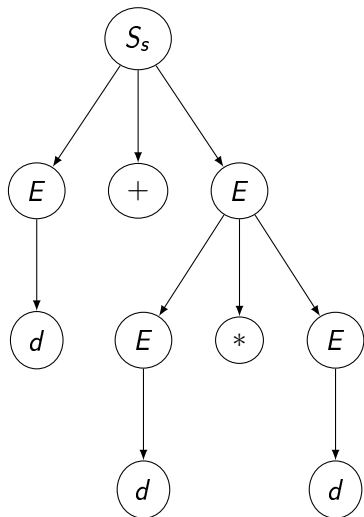
③ $S \rightarrow T$

④ $T \rightarrow T * F$

⑤ $T \rightarrow F$

⑥ $F \rightarrow d$

$d + d + d$ в грамматике G_1



- Не всегда удобны для применения на практике.
 - ▶ Не всегда просты.
 - ▶ Результирующая грамматика сложна.
 - ▶ Не все языки можно описать однозначными грамматиками.

Существенно неоднозначные языки

Definition (Существенно неоднозначный язык)

Язык L называется *существенно неоднозначным* если он может быть задан только неоднозначной грамматикой.

Существенно неоднозначные языки

Definition (Существенно неоднозначный язык)

Язык L называется *существенно неоднозначным* если он может быть задан только неоднозначной грамматикой.

Пример:

$$L_1 = \{a^n b^m c^m d^n \mid n, m > 0\}$$

$$L_2 = \{a^n b^n c^m d^m \mid n, m > 0\}$$

$$L_3 = L_1 \cup L_2$$

- А почему бы не разрешить обрабатывать все варианты?

- А почему бы не разрешить обрабатывать все варианты?
- Записывать в ячейки таблицы больше одной инструкции.
- Во время анализа в момент возникновения неоднозначности дклать нужное количество копий стека и работать с каждым из них независимо.

- А почему бы не разрешить обрабатывать все варианты?
- Записывать в ячейки таблицы больше одной инструкции.
- Во время анализа в момент возникновения неоднозначности дклать нужное количество копий стека и работать с каждым из них независимо.
- Много лишних копий данных.

Объединение общих частей стека

- Копирование появляется при возникновении конфликтов (Reduce/Reduce, Shift/Reduce).
- Его можно заменить ветвлением.

Объединение общих частей стека

- Копирование появляется при возникновении конфликтов (Reduce/Reduce, Shift/Reduce).
- Его можно заменить ветвлением.
- Если вершины двух веток соответствуют одному состоянию, то их можно объединить.

“Историческая справка”

Идеи и алгоритмы известны достаточно давно.

- 1974. Bernard Lang. Не признали.
- 1984. Masaru Tomita. Книга в 200 страниц.

Graph Structuref Stack (GSS)

- Стек можно организовать в виде графа.
- В общем случае граф произвольный.

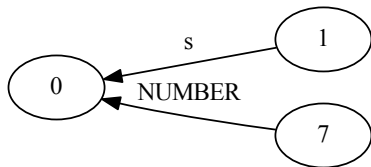
Пример

- Грамматика

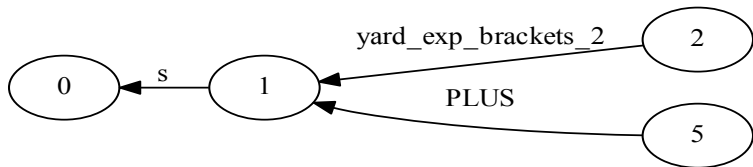
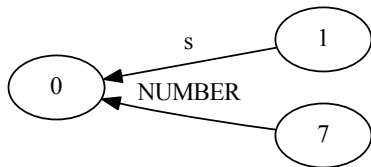
$s: (s(+|*)s) | \text{NUM}$

- Входная строка: NUM+NUM+NUM

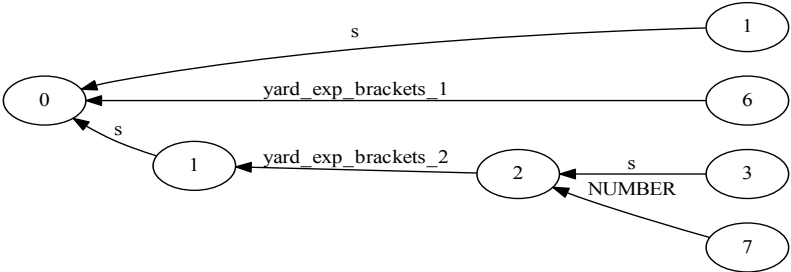
Пример: GSS



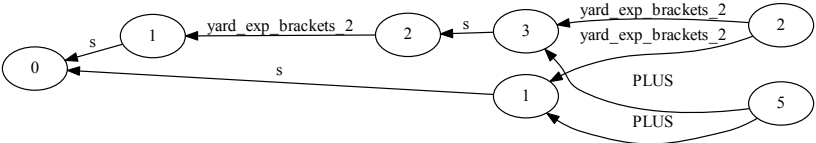
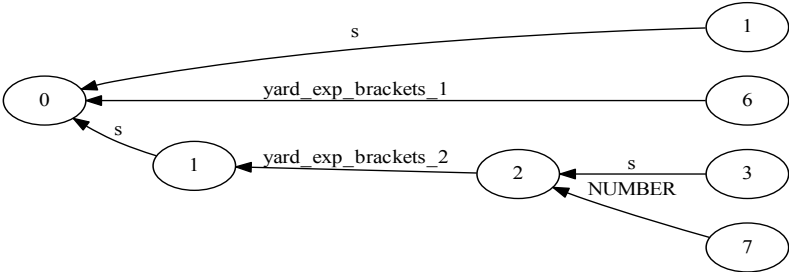
Пример: GSS



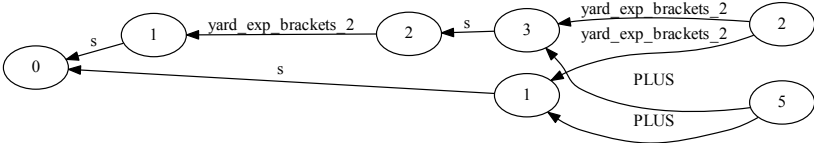
Пример: GSS



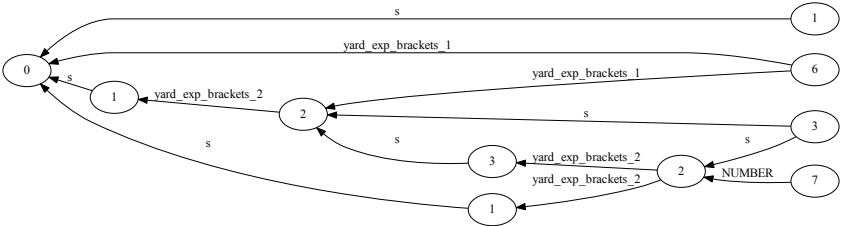
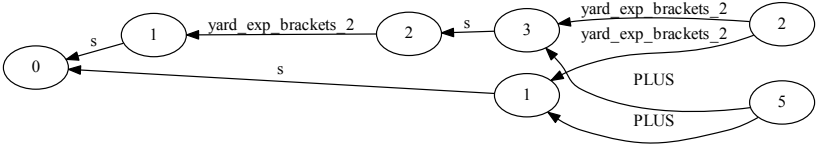
Пример: GSS



Пример: GSS



Пример: GSS



Обобщение GLR-анализа

Немного про то, когда нам могут понадобиться все варианты разбора.

- Динамический SQL

```
IF @X = @Y
    SET @TABLE = '#table1'
ELSE
    SET @TABLE = 'table2'
EXECUTE
    ('SELECT x FROM' + @TABLE + ' WHERE ISNULL(n,0) > 1')
```

- JavaScript в Java

```
String script =
    "function hello(name) print('Hello, ' + name); ";
engine.eval(script);
Invocable inv = (Invocable) engine;
inv.invokeFunction("hello", "Scripting!!!" );
```

- Ошибки в динамически формируемых выражениях обнаруживаются лишь во время выполнения.
- Динамически формируемые выражения – код на некотором языке. Его нужно соответствующим образом поддерживать и обрабатывать.
 - ▶ Поддержка в IDE.
 - ▶ Реинжиниринг ПО, разработанного с использованием встроенных языков.

- Статический анализ для встроенных языков. Абстрактный лексический и синтаксический анализы.
 - ▶ Поддержка в IDE:
 - ★ Многие ошибки можно искать без запуска программы.
 - ★ Автодополнение, рефакторинг.
 - ★ Возможны трансформации.
 - ▶ Реинжиниринг:
 - ★ Статический анализ.
 - ★ Трансформация (трансляция).

- Для каждого выражения строится конструкция, аппроксимирующая множество его возможных значений.
 - ▶ Data-flow уравнение
 - ▶ **Граф**
 - ▶ Регулярное выражение
- Выполнение лексического, синтаксического анализа над графом – абстрактный анализ.

Пример

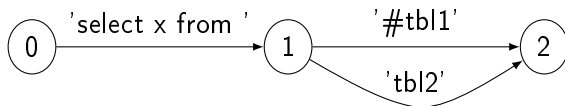
- Код:

```
var tbl1 = "#tbl1"
```

```
var tbl2 = "tbl2"
```

```
execute ("select x from " + if cond then tbl1 else tbl2)
```

- Входной “поток” для анализатора:



Shift-Shift псевдоконфликты

- Ветвление во входном потоке соответствует возможности в текущем состоянии сделать несколько разных переносов (Shift).
- Конечно, это не конфликты в строгом смысле слова, но давайте по аналогии называть их Shift-Shift конфликтами.

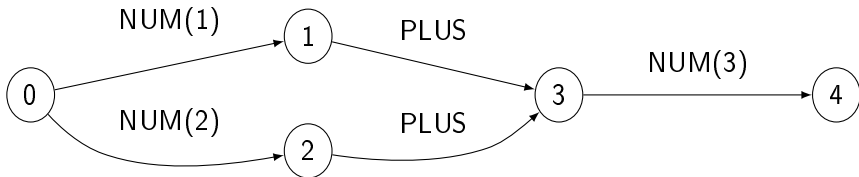
Пример

- Грамматика:

s : NUM PLUS e

e : NUM

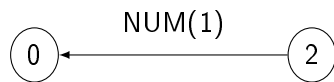
- Входной граф:



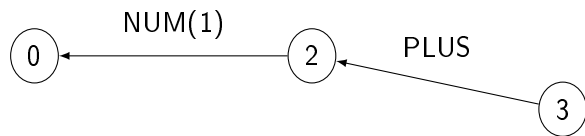
Пример: стек

0

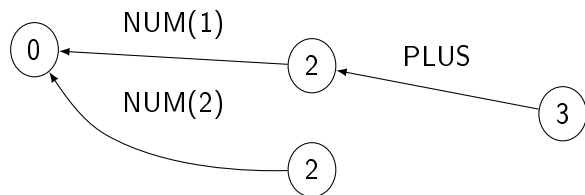
Пример: стек



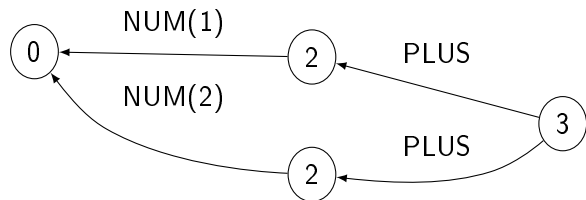
Пример: стек



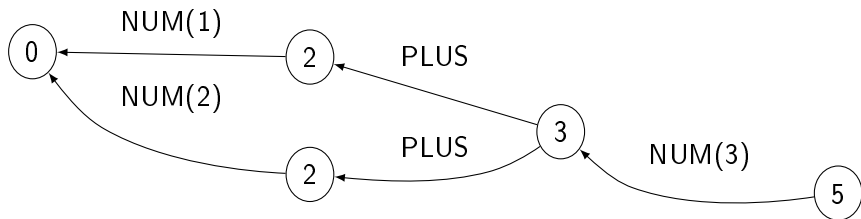
Пример: стек



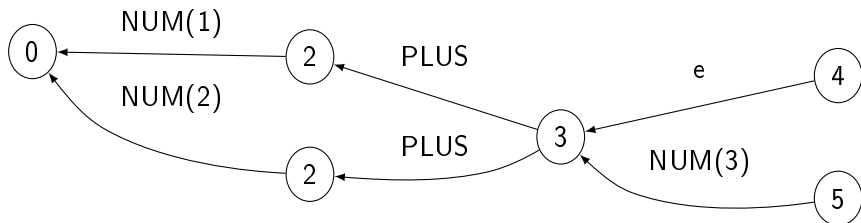
Пример: стек



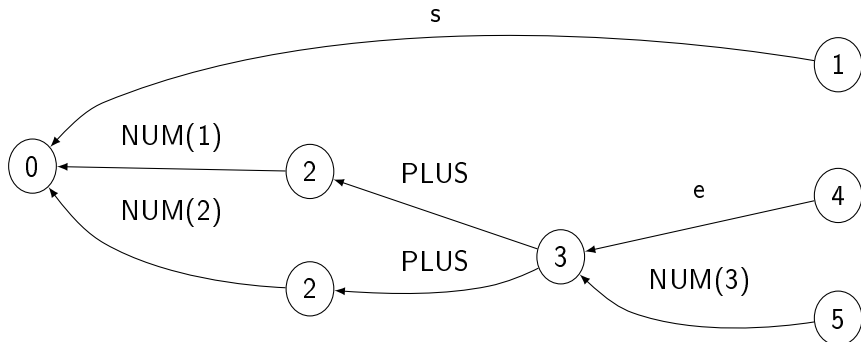
Пример: стек



Пример: стек



Пример: стек



Припер: результат разбора

