

# 1 Первая домашка

**Язык программирования** (ЯП) — множество программ на этом языке. ЯП — частный случай всех языков. Любой язык состоит из синтаксиста (правила построения предложений), семантики (смысл предложений) и прагматики (взаимодействие с носителем).

В языках программирования разделяют **абстрактный и конкретный синтаксис**.

**Абстрактный синтаксис языка выражений:**

$$\begin{aligned}\mathcal{X} &= \{x, y, z, \dots\} \\ \otimes &= \{+, -, \times, /, \%, <, \leq, >, \geq, =, \neq, \vee, \wedge\} \\ \mathcal{E} &= \mathcal{X} \\ &\quad \mathbb{N} \\ &\quad \mathcal{E} \otimes \mathcal{E}\end{aligned}$$

Задано множество переменных, множество бинарных операций и множество **абстрактных синтаксических деревьев** (АСД).

АСД является:

- изолированным узлом, который помечен переменной или числом;
- узлом бинарной операции с двумя дочерними деревьями.

Следовательно, дерево упорядочено и определяется по индукции (базой является первый случай).



Множество АСД определяет язык, а каждая программа определяется АСД. Пример АСД:



Пример конкретного синтаксиста:  $(a + b) * c$ . Чтобы абстрактный синтаксис выразить в конкретный потребуются дополнительные понятия (в данном случае старшинство операций).

**Семантика** — это всюду определенное отображение языка в семантический домен:

$$\llbracket \bullet \rrbracket_{\mathcal{L}} : \mathcal{L} \mapsto \mathcal{D}$$

Функциональный семантический домен — множество функций.

Введем **функцию состояния**  $\sigma : \mathcal{X} \rightarrow \mathbb{Z}$ , которая связывает переменные с их значениями.

Рассмотрим **операционную семантику большого шага**. Она определяется тернарным отношением:

$$\Rightarrow_{\mathcal{E}} \subseteq \Sigma \times \mathcal{E} \times \mathbb{Z}.$$

Запись  $\sigma \xRightarrow{\mathcal{E}} n$  означает, что вычисление выражения  $e$  в состоянии  $\sigma$  дает значение  $n$ .

Три семантических правила:

$$\begin{array}{ll} \frac{n \in \mathbb{N}}{\sigma \xRightarrow{\mathcal{E}} n} & [\text{CONST}] \\ \frac{x \in \mathcal{X}}{\sigma \xrightarrow{x}_{\mathcal{E}} \sigma x} & [\text{VAR}] \\ \frac{\sigma \xRightarrow{l}_{\mathcal{E}} x, \sigma \xRightarrow{r}_{\mathcal{E}} y}{\sigma \xRightarrow{l \otimes r}_{\mathcal{E}} x \oplus y} & [\text{BINOP}] \end{array}$$

В случае с переменной мы должны извлечь из переменной число, чтобы получить результат.

Линия отделяет **посылку** (правило сверху) от **заключения** (правило снизу). К заключению можно приступить только если выполнена посылка. Если у правила отсутствует посылка, то оно является **аксиомой**.

**Базовые операции:**

$\otimes$	$\oplus$ in $\mathcal{AM}^a$
+	+
−	−
×	*
/	/
%	%
<	<
>	>
≤	<=
≥	>=
=	=
≠	!=
∧	&&
∨	!!

Пример программы:  $1 \wedge (1/0)$ . Пример описания:

$$\frac{\frac{\frac{\wedge}{\Rightarrow} \wedge}{\wedge \vee (\wedge / 0)} \quad \frac{\wedge / 0}{\Rightarrow}}{\wedge \vee (\wedge / 0)}$$

Результатом будет неопределенность, так как результат деления на ноль не определен.

Мы рассмотрели **строгую (strict) семантику**. Функция называется неопределенной, если функция от неопределенности есть неопределенность.

Добавим уточнение, что бинарная операция  $\neq \wedge \vee$  и добавим два новых правила, которые сделают семантику нестрогой по второму аргументу:

$$\frac{c \stackrel{e}{\Rightarrow} 1}{c \stackrel{e \vee r}{\Rightarrow} 1} \quad \frac{c \stackrel{e}{\Rightarrow} z+1, c \stackrel{e}{\Rightarrow} z'}{c \stackrel{e \vee r}{\Rightarrow} z \parallel z'}$$

Если левый аргумент — единица, то результат будет единицей. Если левый аргумент не равен единице, а правый аргумент — целое число, то результат вычисляется из этих двух чисел.

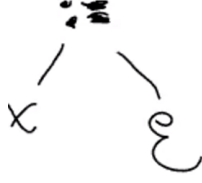
**Обобщение правил  $\Rightarrow_{\mathcal{E}}$  в семантику  $\llbracket \bullet \rrbracket_{\mathcal{E}}$ :**

$$\frac{\sigma \stackrel{e}{\Rightarrow}_{\mathcal{E}} n}{\llbracket e \rrbracket_{\mathcal{E}} \sigma = n}$$

Язык состоит из множества синтаксических категорий. Существует несколько синтаксических категорий. Категорию выражений (expression) мы рассмотрели. Следующая категория — операторов (statement). Рассмотрим язык операторов:

$\mathcal{S} =$  **skip**  
 $\mathcal{X} := \mathcal{E}$   
**read**( $\mathcal{X}$ )  
**write**( $\mathcal{E}$ )  
 $\mathcal{S}; \mathcal{S}$

Второе правило определяет операцию присвоения, где слева находится переменная, справа — выражение (рассмотрено ранее).



Последнее правило определяет композицию операторов.

Определим семантику этого языка. Пусть наш язык описывает программы, которые принимают строки целых чисел и возвращают строки целых чисел:

$$\llbracket \bullet \rrbracket : \mathcal{S} \mapsto \mathbb{Z}^* \rightarrow \mathbb{Z}^*$$

Семантика большого шага связывает две конфигурации и программу:

$$\Rightarrow_{\mathcal{S}} \subseteq \mathcal{C} \times \mathcal{S} \times \mathcal{C}.$$

Значит, что при запуске программы в одной конфигурации, через несколько шагов получаем другую конфигурацию. Семантика большого шага так называется, так как при переходе выполняется цепочка шагов.

**Конфигурация** — это пара состояния и мира:

$$\mathcal{C} : \Sigma \times \mathcal{W}$$

**Мир** — это инкапсуляция входного и выходного потока:

$$\mathcal{W} : \mathbb{Z}^* \times \mathbb{Z}^*$$

**Определим операции:**

$$\begin{aligned} \mathbf{read} \langle xi, o \rangle &= \langle x, \langle i, o \rangle \rangle \\ \mathbf{write} \ x \langle i, o \rangle &= \langle i, ox \rangle \\ \mathbf{out} \langle i, o \rangle &= o \end{aligned}$$

Здесь  $xi \in \mathbb{Z}^*$  — конкатенация слов  $x$  и  $i$ . Создает новый мир, где входной поток определяется  $x$ , а выходной прошлым миром. Запись добавляет  $x$  в конец выходного потока. Вывод выводит выходной поток.

$$\begin{array}{c}
c \xRightarrow{\text{skip}}_{\mathcal{S}} c \quad [\text{SKIP}] \\
\langle \sigma, \omega \rangle \xRightarrow{x := e}_{\mathcal{S}} \langle \sigma[x \leftarrow \llbracket e \rrbracket_{\mathcal{E}} \sigma], \omega \rangle \quad [\text{ASSIGN}] \\
\frac{\langle z, \omega' \rangle = \text{read } \omega}{\langle \sigma, \omega \rangle \xRightarrow{\text{read}(x)}_{\mathcal{S}} \langle \sigma[x \leftarrow z], \omega' \rangle} \quad [\text{READ}] \\
\langle \sigma, \omega \rangle \xRightarrow{\text{write}(e)}_{\mathcal{S}} \langle \sigma, \text{write}(\llbracket e \rrbracket_{\mathcal{E}} \sigma) \omega \rangle \quad [\text{WRITE}] \\
\frac{c_1 \xRightarrow{S_1}_{\mathcal{S}} c' \quad c' \xRightarrow{S_2}_{\mathcal{S}} c_2}{c_1 \xRightarrow{S_1; S_2}_{\mathcal{S}} c_2} \quad [\text{SEQ}]
\end{array}$$

Во время присвоения конфигурация явно представляется в виде пары состояния и мира  $\langle \sigma, \omega \rangle$ . В результате переменной  $x$  в состоянии  $\sigma$  присваивается результат выполнения выражения, а мир не изменяется.

Операционная семантика языка определяет эталонный интерпретатор. Интерпретатор — это программа, которая принимает программу на данном языке и для любого входа возвращает выход этой программы. Семантика интерпретатора:

$$\begin{aligned}
\llbracket \text{int} \rrbracket : \mathcal{L} &\rightarrow I \rightarrow O \\
\llbracket \text{int} \rrbracket_{\mathcal{L}}(p, x) &= \llbracket p \rrbracket_{\mathcal{M}}(x)
\end{aligned}$$

Обычно скобки не ставят, чтобы запись не была перегружена разными символами и читалась легче.

Обобщим семантику оператора и определим  $\llbracket \bullet \rrbracket_{\mathcal{S}}$ :

$$\frac{\langle \Lambda, \langle i, \epsilon \rangle \rangle \xRightarrow{S}_{\mathcal{S}} \langle \sigma, \omega \rangle}{\llbracket S \rrbracket_{\mathcal{S}} i = \text{out } \omega}$$

Семантика принимает числа  $i$ , после чего из когфигурации, где сосояние пусто (ни одна переменна не определена), выполняется переход по соответствующему правилу в конфигурацию с каким-то состоянием и миром. Результатом является выходной поток полученного мира.

Наш компилятор является многопроходным. Сначала получаем программу в виде синтаксического дерева, затем интерпретируем ее в соответствии с семантикой большого шага. Далее мы компилируем программу в стековую машину, а из стековой машины можем выполнить компиляцию в x86.

Фактически стековая машина является языком программирования со своей семантикой.

Синтаксис стековой машины состоит из двух синтаксических категорий: инструкции (instructions) и программы (programs)

$$\begin{aligned}
\mathcal{I} &= \text{BINOP} \otimes \\
&\quad \text{CONST } \mathbb{N} \\
&\quad \text{READ} \\
&\quad \text{WRITE} \\
&\quad \text{LD } \mathcal{X} \\
&\quad \text{ST } \mathcal{X} \\
\mathcal{P} &= \varepsilon \\
&\quad \mathcal{I} \mathcal{P}
\end{aligned}$$

Программа стековой машины является списком (возможно пустым) инструкций.

Конфигурация стековой машины включает стек. Множество конфигураций можно определить следующим образом:

$$\mathcal{C}_{SM} = \mathbb{Z}^* \times \mathcal{C}$$

Правила семантики большого шага для СМ:

$$\begin{aligned}
c &\xRightarrow{\varepsilon}_{\mathcal{SM}} c && [\text{STOP}_{SM}] \\
\frac{\langle (x \oplus y)s, c \rangle \xRightarrow{p}_{\mathcal{SM}} c'}{\langle yxs, c \rangle \xRightarrow{[\text{BINOP } \otimes] p}_{\mathcal{SM}} c'} &&& [\text{BINOP}_{SM}] \\
\frac{\langle zs, c \rangle \xRightarrow{p}_{\mathcal{SM}} c'}{\langle s, c \rangle \xRightarrow{[\text{CONST } z] p}_{\mathcal{SM}} c'} &&& [\text{CONST}_{SM}] \\
\frac{\langle z, \omega' \rangle = \mathbf{read } \omega, \langle zs, \langle \sigma, \omega' \rangle \rangle \xRightarrow{p}_{\mathcal{SM}} c'}{\langle s, \langle \sigma, \omega \rangle \rangle \xRightarrow{\text{READ } p}_{\mathcal{SM}} c'} &&& [\text{READ}_{SM}] \\
\frac{\langle s, \langle \sigma, \mathbf{write } z \omega \rangle \rangle \xRightarrow{p}_{\mathcal{SM}} c'}{\langle zs, \langle \sigma, \omega \rangle \rangle \xRightarrow{\text{WRITE } p}_{\mathcal{SM}} c'} &&& [\text{WRITE}_{SM}] \\
\frac{\langle (\sigma x)s, \langle \sigma, \omega \rangle \rangle \xRightarrow{p}_{\mathcal{SM}} c'}{\langle s, \langle \sigma, \omega \rangle \rangle \xRightarrow{[\text{LD } x] p}_{\mathcal{SM}} c'} &&& [\text{LD}_{SM}] \\
\frac{\langle s, \langle \sigma[x \leftarrow z], \omega \rangle \rangle \xRightarrow{p}_{\mathcal{SM}} c'}{\langle zs, \langle \sigma, \omega \rangle \rangle \xRightarrow{[\text{ST } x] p}_{\mathcal{SM}} c'} &&& [\text{ST}_{SM}]
\end{aligned}$$

$p$  на стрелке означает выполнение оставшейся части программы.

LD кладет значение переменной  $(\sigma x)$  в стек. ST кладет значение из стека  $z$  в переменную  $x$  ( $\sigma[x \leftarrow z]$ ).

Семантика СМ:

$$\frac{\langle \varepsilon, \langle \Lambda, \langle i, \varepsilon \rangle \rangle \rangle \xRightarrow{p}_{\mathcal{SM}} \langle s, \langle \sigma, \omega \rangle \rangle}{\llbracket p \rrbracket_{\mathcal{SM}} i = \mathbf{out } \omega}$$

Компилятор из языка операторов в стековую машину:

$$\llbracket \bullet \rrbracket^{comp} : \mathcal{S} \mapsto \mathcal{P}$$

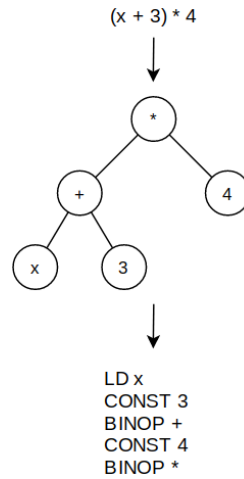
Компилятор — это функция, которая отображает программу на одном языке в программу на другом, сохраняя при этом семантику.

Компилятор можно рассматривать как семантику языка программирования. У компилятора в качестве семантического домена выступает другой язык. Таким образом, у языка может быть больше одной семантики.

Опишем семантику в виде. Просто потому что так проще. Пусть исходный язык состоит из двух синтаксических категорий: выражения и операторы. Компилятор будет состоять из двух функций, которые генерируют код для каждой категории. Компилятор выражений (получает выражение, возвращает последовательность инструкций):

$$\begin{aligned} \llbracket x \rrbracket_{\mathcal{E}}^{comp} &= [\text{LD } x] \\ \llbracket n \rrbracket_{\mathcal{E}}^{comp} &= [\text{CONST } n] \\ \llbracket A \otimes B \rrbracket_{\mathcal{E}}^{comp} &= \llbracket A \rrbracket_{\mathcal{E}}^{comp} \llbracket B \rrbracket_{\mathcal{E}}^{comp} [\text{BINOP } \otimes] \end{aligned}$$

Пример компиляции:



Компилятор операторов:

$$\begin{aligned} \llbracket x := e \rrbracket^{comp} &= \llbracket e \rrbracket_{\mathcal{E}}^{comp} [\text{ST } x] \\ \llbracket \text{read}(x) \rrbracket^{comp} &= [\text{READ}] [\text{ST } x] \\ \llbracket \text{write}(e) \rrbracket^{comp} &= \llbracket e \rrbracket_{\mathcal{E}}^{comp} [\text{WRITE}] \\ \llbracket S_1; S_2 \rrbracket^{comp} &= \llbracket S_1 \rrbracket^{comp} \llbracket S_2 \rrbracket^{comp} \end{aligned}$$

$$\llbracket \text{skip} \rrbracket^{comp} = \epsilon$$

В случае с присвоением мы при помощи компилятора для выражений получаем значение, затем кладем его в переменную.

Ветка с домашкой начинается на A01. Если в корне проекта выполнить `make`, то начнут выполняться регрессионные тесты. Чтобы их пройти нужно

<code>src/World.lama</code>	Реализация семантики мира
<code>src/State.lama</code>	Реализация состояния
<code>src/Expr.lama,</code>	Реализация правил и компиляции. Нужно дописать
<code>src/Stmt.lama,</code>	функции, которые возвращают failure
<code>src/SM.lama</code>	

В компиляторе лучше использовать буфер для списка инструкций

Делать пул реквест в тот же бранч. Сначала убедиться, что тесты работают локально. Если его отклонили, то все ок. В табличке результаты. Одна ветка — одно ДЗ. В комментариях к пул реквесту указать кто я.

Установка компилятора языка Lama:

- Ubuntu 18 с 4 ГиБ памяти,
- пакетный менеджер OPAM,
- язык программирования OCaml 4.10.1.

```
$ vagrant init ubuntu/bionic64
```

```
$ vagrant up
```

```
$ vagrant ssh
```

```
# apt-get update
```

```
# apt-get -y upgrade
```

```
# apt-get -y install gcc-multilib make m4 unzip bubblewrap gdb
```

```
# sh <(curl -sL https://raw.githubusercontent.com/ocaml/opam/master/shell/install.sh)
```

```
$ opam init -y
```

```
$ eval $(opam env)
```

```
$ opam switch create lama ocaml-variants.4.10.1+fp+flambda -y
```

```
$ eval $(opam env)
```

```
$ opam pin add Lama https://github.com/JetBrains-Research/Lama-devel.git\#1.10+ocaml4.10 --no-action -y
```

```
$ opam depext Lama -y
```

```
$ opam install Lama -y
```

regression

## 2 Вторая домашка

Задача: реализовать компилятор CM в GAS x86.

Ветка для решения: A02-straight-line-x86. Нужно использовать код из предыдущей домашки.

```
git remote rename origin fork
```

```
git remote add original https://github.com/danyaberezun/compilers-2021-spring
```

```
git fetch original
```

```
git checkout -b A02-straight-line-x86 original/A02-straight-line-x86
```



Mov (opnd, opnd)	Копирует значение из первого во второй операнд
Binop	делает бинарную операцию
IDiv	Деление целого
Cltld	s

Разницы между X86 и стековой машиной практически нет. У обоих есть стек, мир и состояние. Здесь входным потоком будет stdin, выходным — stdout. Память будет состоянием.

Регистров в стековой машине нет. Но их нужно использовать при генерации кода. Будем использовать символическую интерпретацию. Будем генерировать код с помощью символического интерпретатора для стековой машины. На стеке будут лежать не данные, а размещения (регистр, позиция на стеке, либо переменная в памяти).

Раньше компилятор получал конфигурацию и код. Сейчас компилятор будет получать окружение и код. Окружение можно считать объектом.

Размещения:

R (id)	регистр
S (int)	позиция на стеке
M (string)	именованная локация в памяти
L (int)	k

Регистры:

ebx	регистр (ebx, ecx, esi,
ecx	делает бинарную операцию
esi	Деление целого
edi	k
eax	Деление
edx	Деление
ebp	
esp	

READ – отводим новую позицию на символическом стеке (сначала выделяются регистры, затем стек). Получаем новую позицию и новое окружение. Возвращается пара — новое окружение и сгенерированный код

env.allocate – выделить новую позицию на стеке.

ST (x) снять с символического стека, засунуть в переменную

Важно создавать переменную во время чтения, если она не определена.

### 3 Третья домашка

Частично корректные компиляторы расширяют область определения программ. Они не строго соответствуют семантике. Ее поведение может отличаться после компиляции. Эталонный интерпретатор полностью повторяет операционную семантику.

Компиляторы как правило являются частично корректными. Просто потому что это сложно. Частичная коррекция достигается не специально.

Результат какой-то программы определяется ее семантикой или эталонным интерпретатором.

**Тьюринг полнота** — это возможность реализовать любую вычислимую функцию.

**Вычислимая функция:**

$$f : \mathbb{N}^k \rightarrow \mathbb{N}$$

На данный момент наш язык позволяет реализовать только полиномы, он не является тьюринг-полным.

Тут написать про оптимальность компилятора, перерисовать. Минимум шагов.

Задача оптимальной компиляции даже нашего языка неразрешимая.

10-я проблема Гильберта. Пусть есть многочлен с  $n$  переменными и целыми коэффициентами  $p(x_1, x_2, \dots, x_k)$ . Корни ищатся перебором. Если к нулю не сходится, то будет выполняться бесконечно. Процедуры, которая позволит найти корни без перебора не существует.

Синтаксический анализ занимает меньшую часть компилятора (четверть или десятую).

Синтаксически центричный принцип создания компилятора. Конкретный синтаксис языка программирования задается КС грамматикой. Конкретных синтаксисов у одного языка может быть несколько (например поддержка ключевых слов на русском языке).

Синтаксический анализатор преобразует конкретный синтаксис в абстрактный.

Сейчас у нас есть абстрактный синтаксис, но нет конкретного.

Дерево вывода КС грамматики является деревом абстрактного синтаксиса.

Проблемы:

На самом деле большинство языков (даже простых) синтаксис не является КС. Пример для языка C (конструкция `a * b`):

```
typedef int a;

a * b;

int f (int a, int b) {
    return a * b;
}
```

Во-вторых, дерево вывода не является АСД, так как один язык может иметь несколько грамматик, а АСД одно.

В-третьих, грамматика может быть контекстно зависима. Например,

```
<A <B>> f () {}

a >> b
```

Терминал » означает разное в зависимости от своей синтаксической позиции.

Отступы не значащие, но не строго. Например, в си нельзя удалить пробелы здесь (это еще одно подтверждение, что грамматика не КС):

```
void f () { ... }
```

В реальности используются монадические парсер-комбинаторы.

Атрибутные грамматики — это КС грамматика, где для каждого нетерминала указано 2 множества атрибутов. Атрибуты бывают наследуемые и синтезируемые. Если воспринимать нетерминал как процедуру, которая распознает что-то во входном потоке. Наследуемые атрибуты — это факормальные параметры распознающей процедуры, а синтезируемые атрибуты — это результат.

То есть раньше КС-грамматика распознавала строку и возвращала 0 или 1 (результат распознавания), а сейчас грамматика параметризована и может возвращать больше результатов.

L-атрибутные грамматики.

$L \rightarrow R1 R2 \dots R3$

Схема L-атрибутного анализа:

$L[a] \rightarrow x=R1[a] \ y=R2[x] \ z=R3[y] \ \{z\}$

Наследуемый атрибут можно использовать везде. Терминалы могут использовать синтезируемые атрибуты предыдущих терминалов. Результатом является z

Атрибуты позволяют выразить контекстную зависимость.

Будем использовать левоатрибутный анализ.

Как реализовывать такие грамматики?

Есть монадические парсер-комбинаторы. Комбинатор — это функция высшего порядка, она комбинирует функции. Парсер принимает входной поток и что-то простое распознает в нем.

Функция высшего порядка — это функция, принимающая в качестве аргументов другие функции или возвращающая другую функцию в качестве результата.

```
fun (stream) -> Ok (result, stream) | fail (desc)
```

Получает список, возвращает результат и остаток списка.

Парсер комбинатор принимает что-то и возвращает парсер token принимает t и возвращает функцию, которая распознает это t.

atl — парсер альтернативы. Является левопредпочтительным, так как возвращает результат левого парсера, если распознавание с его помощью удалось. Иначе применяется правый парсер.

Монадические комбинаторы: seq.

Данный комбинатор комбинирует парсер и функцию, которая получает результат предыдущего парсера и возвращает парсер.

Ostap улучшенная реализация.

Лама содержит в себе DSL (предмето ориентированный язык) для описания синтаксического анализатора.

Начинается символом `syntax`. Является расширенной формулой бэкуса наура.

Lexer.lama — лексический анализатор. Набор примитивных парсеров. Содержит регулярки для терминалов (число, отступ)

Нужно дописать Parser.lama

memo — мемоизатор. eta — расширение, которое делает функцию из того, что следует далее.

Driver.lama — драйвер. Компонент компилятора, с которого начинается его работа. Разбирает опции командной строки, что-то запускает...

Пример простого парсера:

```
fun eof (stream) {
  case stream of
    {} -> Ok ({}, stream)
  | _ -> Fail ("EOF expected")
  esac
}
```

Комбинатор — это функция, которая получает действие и возвращает измененное действие.

Пример парсер-комбинатора, распознающего символ:

```
fun token (t) {
  fun (stream) {
    case stream of
      t0 : stream0 -> if compare (t0, t) == 0
                      then Ok (t, stream0)
                      else Fail (t.string ++ " expected")
      fi
    | _ -> Fail (t.string ++ " expected")   поток пустой
    esac
  }
}
```

В данном случае действием является распознающая символ функция, которая определяется самим комбинатором и изменяется так, чтобы распознавать определенный символ.

Пример монадического парсер-комбинатора:

```
fun seq (p, q) {    p - парсер, q - комбинатор
  fun (stream) {
    case p (stream) of
      failure@Fail (_) -> failure
    | Ok (result, stream0) -> q (result) (stream0)
    esac
  }
}
```

```

}
}

```

Монада определяет следующее вычисление на основе результата предыдущего.

Монадический парсер-комбинатор реализует схему L-атрибутной грамматики.

В модуле Ostar улучшенная реализация парсеров.

Семантические терминалы называются semantic action

## 4 Четвертая домашка

Написать оптимальный компилятор невозможно

Расширим язык операторов структурным потоком управления (structural control flow).

Называется структурным, так как нет перехода по метке (оператор goto).

$\mathcal{S} \quad + = \quad \mathbf{if} \ \mathcal{E} \ \mathbf{then} \ \mathcal{S} \ \mathbf{else} \ \mathcal{S}$   
 $\mathbf{while} \ \mathcal{E} \ \mathbf{do} \ \mathcal{S}$

Добавляем два оператора. Они записаны в абстрактном синтаксисе, но в символической форме. В виде дерева их можно записать так:

В конкретном синтаксисе обязательно присутствуют закрывающие конструкции (fi, od).

Операционная семантика для новых конструкций состоит из 4 правил:

$$\begin{array}{c}
 \frac{\sigma \xRightarrow{e}_{\mathcal{E}} n \neq 0 \quad \langle \sigma, w \rangle \xRightarrow{S_1}_{\mathcal{S}} c'}{\langle \sigma, w \rangle \xRightarrow{\mathbf{if} \ e \ \mathbf{then} \ S_1 \ \mathbf{else} \ S_2}_{\mathcal{S}} c'} \quad [\text{IF-TRUE}] \\
 \\
 \frac{\sigma \xRightarrow{e}_{\mathcal{E}} 0 \quad \langle \sigma, w \rangle \xRightarrow{S_2}_{\mathcal{S}} c'}{\langle \sigma, w \rangle \xRightarrow{\mathbf{if} \ e \ \mathbf{then} \ S_1 \ \mathbf{else} \ S_2}_{\mathcal{S}} c'} \quad [\text{IF-FALSE}] \\
 \\
 \frac{\sigma \xRightarrow{e}_{\mathcal{E}} n \neq 0 \quad \langle \sigma, w \rangle \xRightarrow{S}_{\mathcal{S}} c' \quad c' \xRightarrow{\mathbf{while} \ e \ \mathbf{do} \ S}_{\mathcal{S}} c''}{\langle \sigma, w \rangle \xRightarrow{\mathbf{while} \ e \ \mathbf{do} \ S}_{\mathcal{S}} c''} \quad [\text{WHILE-TRUE}] \\
 \\
 \frac{\sigma \xRightarrow{e}_{\mathcal{E}} 0}{\langle \sigma, w \rangle \xRightarrow{\mathbf{while} \ e \ \mathbf{do} \ S}_{\mathcal{S}} \langle \sigma, w \rangle} \quad [\text{WHILE-FALSE}]
 \end{array}$$

Семантика есть описание поведения программы.

Добавим синтаксическое расширение. Или синтаксический сахар. Это расширение которое существует в конкретном синтаксисе, но не в абстрактном, в абстрактный они проецируются.

```

if e1 then s1  ->  if e1 then s1
elif e2 then s2      else if e2 then s2
else s3              else sk

```

```

fi
fi
fi

```

Достаточно добавить изменения в синтаксический анализатор (парсер), чтобы все заработало.

```

for s1, e, s2 do s3 od -> s1;
    while e do
        s3; s2
    od

```

Цикл с постусловием:

```

do s while e od -> s;
    while e do
        s
    od

```

Реализовывать его при помощи синтаксического расширения ужасно, так как получаем экспоненциальную сложность в преобразовании

```

do
    do
        while od
    while od ->
do
    while do
        do
            while od
        while od
    od

```

Если при синтаксическом расширении получается нелинейное преобразование, то лучше его не реализовывать. Поэтому будем добавлять эту конструкцию на уровне абстрактного синтаксиса. Дз определить семантику.

## 5 Пятая домашка

Задача: расширить компиляцию в стековую машину и X86.

На данный момент язык стековой машины не достаточно выразителен, чтобы отобразить циклы. Поэтому расширим ее, то есть добавим новые конструкции и опишем семантику для них. множество меток:

$$\mathcal{L} = \{l_1, l_2, \dots\}$$

$$\begin{aligned} \mathcal{I} \quad + = \quad & \text{LABEL } \mathcal{L} \\ & \text{JMP } \mathcal{L} \\ & \text{CJMP}_x \mathcal{L}, \text{ where } x \in \{\text{nz}, \text{z}\} \end{aligned}$$

JMP — инструкция безусловного перехода к метке, CJMP — инструкция перехода по 0 или 1 (значение берется из стека).

В описании семантики будем использовать окружение (environment)  $\Gamma$ . В нашем случае окружением является программа  $P$  (последовательность инструкций),

которую мы интерпретируем.  $P[l]$  — подпрограмма  $P$  которая начинается с метки  $l$ . В последовательности инструкций ..., LABEL  $l$  ... правая часть и есть  $P[l]$ .

В семантических правилах будем исходную программу обозначать  $P$ , а  $p$  — хвост, до которого мы добрались.

$$\frac{p \vdash \langle \epsilon, \langle \Lambda, \langle i, \epsilon \rangle \rangle \rangle \xRightarrow{P} \mathcal{SM} \langle s, \langle \sigma, \omega \rangle \rangle}{\llbracket p \rrbracket_{\mathcal{SM}} i = \mathbf{out} \ \omega}$$

$$\begin{array}{l} \frac{P \vdash c \xRightarrow{P} \mathcal{SM} c'}{P \vdash c \xRightarrow{[\text{LABEL } l]P} \mathcal{SM} c'} \quad [\text{LABEL}_{SM}] \\ \frac{P \vdash c \xRightarrow{P[l]} \mathcal{SM} c'}{P \vdash c \xRightarrow{[\text{JMP } l]P} \mathcal{SM} c'} \quad [\text{JMP}_{SM}] \\ \frac{z \neq 0, \quad P \vdash \langle s, \theta \rangle \xRightarrow{P[l]} \mathcal{SM} c'}{P \vdash \langle zs, \theta \rangle \xRightarrow{[\text{CJMP}_{nz} l]P} \mathcal{SM} c'} \quad [\text{CJMP}_{nz}^{+} SM] \\ \frac{z = 0, \quad P \vdash \langle s, \theta \rangle \xRightarrow{P} \mathcal{SM} c'}{P \vdash \langle zs, \theta \rangle \xRightarrow{[\text{CJMP}_{nz} l]P} \mathcal{SM} c'} \quad [\text{CJMP}_{nz}^{-} SM] \\ \frac{z = 0, \quad P \vdash \langle s, \theta \rangle \xRightarrow{P[l]} \mathcal{SM} c'}{P \vdash \langle zs, \theta \rangle \xRightarrow{[\text{CJMP}_z l]P} \mathcal{SM} c'} \quad [\text{CJMP}_z^{+} SM] \\ \frac{z \neq 0, \quad P \vdash \langle s, \theta \rangle \xRightarrow{P} \mathcal{SM} c'}{P \vdash \langle zs, \theta \rangle \xRightarrow{[\text{CJMP}_z l]P} \mathcal{SM} c'} \quad [\text{CJMP}_z^{-} SM] \end{array}$$

Метка ничего не делает, из прошлой конфигурации выполняем хвост. Безусловный переход изменяет только поток управления.

$\llbracket c \rrbracket - - - SMc., .CJMP_z l_{s2}$

$$\begin{array}{ll} \llbracket c \rrbracket - \text{код SM для } c & \\ \text{if } c & \text{CJMP}_z^{+} l_{s2} \\ \text{then } s_1 & \llbracket s_1 \rrbracket \\ \text{else } s_2 & \text{JMP } l_{\text{eof}} \\ \text{fi} & \text{LABEL } l_{s2} \\ & \llbracket s_2 \rrbracket \\ & \text{LABEL } l_{\text{eof}} \end{array} \quad (1)$$

Проблема в такой генерации возникает, когда встречается множество вложенных условий. Тогда переход для выхода будет совершаться многократно (линейно).

Процессоры общего назначения (ARM, X86) устроены так, что используют предсказание переходов. Условные переходы выполняются заранее, поэтому имеют нулевую вычислительную стоимость (процессор заранее берет инструкции для выполнения из нужного места). Для условных переходов выполняется угадывание,

которое почти всегда верно, но стоимость ошибки выше, чем переход.

$$\begin{array}{rcl}
 & \text{LABEL } l_c & \text{JMP } l_c \\
 \text{while } c & \llbracket c \rrbracket & \text{LABEL } l_s \\
 \text{do } s & \rightarrow \text{CJMP}_z^+ l_{\text{eof}} & \llbracket s \rrbracket \\
 \text{od} & \llbracket s \rrbracket & \text{LABEL } l_c \\
 & \text{JMP } l_c & \llbracket c \rrbracket \\
 & \text{LABEL } l_{\text{eof}} & \text{CJMP}_{nz}^+ l_s
 \end{array} \tag{2}$$

Вторая генерация более правильная, так как в теле цикла выполняется только одна команда перехода, хотя такая разница незначительная для современных процессоров.

Хвост программы по метке сохраняется в окружении.

Метка переходит в метку. Безусловный переход превращается в безусловный. Условный переход использует СМ, использует стек, в x86 нужно смотреть в символический стек. Значение с символического стека нужно переместить в регистр флагов, так как инструкция условного перехода использует их (см. сравнивать с нулем).

Рассинхронизация символического стека не происходит, так как в момент перехода имеем пустой стек.