

Tytuł

Patryk Wałach

January 2022

Contents

1	[Wstęp]	2
2	Ogólne wprowadzenie	2
2.1	Analiza leksykalna	2
2.2	Parser	2
2.3	System typów	3
2.4	System typów ML	3
3	Inferencja typów w teorii	4
3.1	System typów Hindley–Milner	4
3.2	Algorytm W	4
4	Rescript/ReasonML jako języki realizujące podobne zadania	4
5	Założenia i priorytety opracowanej aplikacji	5
5.1	Opis formalny składni języka	9
6	Narzędzia	9
6.1	Język python	9
6.2	Parsowanie i tokenizowanie przy użyciu biblioteki sly	9
6.3	Środowisko nodejs do uruchomienia skompilowanego kodu	9
7	Implementacja	9
7.1	lexer	9
7.2	parser	9
7.3	inferencja typów	9
7.4	kompilacja	9
8	Opis działania	9
8.1	Co działa	9
8.2	Uwagi co do obsługi błędów	9
9	[Podsumowanie]	9

1 [Wstęp]

2 Ogólne wprowadzenie

2.1 Analiza leksykalna

Analiza leksykalna w informatyce jest to proces rozbijania program źródłowych na jednostki logiczne (zwane leksemami) złożone z jednego lub więcej znaków, które łącznie mają jakieś znaczenie[1]. Przykładami leksemów mogą być słowa kluczowe (np. `while`), identyfikator lub liczba składająca się z cyfr. Rozdzielaniem programu źródłowego, na leksemy, zajmuje się lekser.

Token jest strukturą reprezentującą leksem i wprost go kategoryzującą[3], co ułatwia późniejszą pracę parserowi. Tokeny kategoryzuje się na komputerowy odpowiednik tego, co lingwiści określiliby mianem części mowy. Biorąc jako przykład poniższy kod w języku C:

```
1  x = a + b * 2;
```

analiza leksykalna, zwraca tokeny:

```
1  [(identifier, x), (operator, =), (identifier, a), (operator, +),
    (identifier, b), (operator, *), (literal, 2), (separator, ;)]
```

Dwoma ważnymi przypadkami są znaki białe i komentarze. One również muszą być uwzględnione w gramatyce i przeanalizowane przez lexer, lecz mogą być odrzucone (nie produkować żadnych tokenów) i traktowane jako spełniające mało znaczące zadanie, rozdzielania dwóch tokenów (np. w `if x` zamiast `ifx`).

2.2 Parser

Analizator składniowy, parser – program dokonujący analizy składniowej danych wejściowych w celu określenia ich struktury gramatycznej w związku z określoną gramatyką formalną. Analizator składniowy umożliwia przetworzenie tekstu czytelnego dla człowieka w strukturę danych przydatną dla oprogramowania komputera. Wynikiem analizy składni, dokonywanej przez parser, najczęściej jest drzewo składniowe nazywane czasami drzewem wyprowadzenia[2].

Zadanie parsera sprowadza się do sprawdzenia czy i jak dane wejściowe mogą zostać wyprowadzone z symbolu startowego. To zadanie można zrealizować na dwa sposoby:

- Analiza zstępująca (ang. top-down parsing) to strategia znajdowania powiązań między danymi przez stawianie hipotez dotyczących drzewa rozbioru składniowego i sprawdzanie, czy zależności między danymi są zgodne z tymi hipotezami.
- Analiza wstępująca (ang. bottom-up parsing) – ogólna metoda analizy składniowej, w której zaczyna się od słowa wejściowego i próbuje się zredukować je do symbolu startowego. Drzewo wyprowadzenia jest konstruowane od liści do korzenia (stąd nazwa). W każdym momencie w trakcie tego procesu mamy formę zdaniową, która zawiera segment, powstały w ostatnim kroku wyprowadzenia. Segment ten (nazywany uchwytym (ang. handle) jest prawą stroną produkcji i powinien zostać w tym kroku zredukowany do jej lewej strony, w wyniku czego powstanie poprzednia forma zdaniowa z wyprowadzenia. Główna trudność w analizie wstępującej polega właśnie na odpowiednim znajdowaniu uchwytów. Analiza wstępująca może przebiegać w określonym kierunku (np. od lewej do prawej), lub w sposób bezkierunkowy, wtedy analizowane jest całe słowo naraz. Jednym z bardziej znanych przedstawicieli metody bezkierunkowej jest algorytm CYK. Do metod kierunkowych zalicza się między innymi parsery shift-reduce czyli LR, LALR, SLR, BC, pierwszeństwa.

2.3 System typów

System typów – system klasyfikacji wyrażeń w zależności od rodzajów wartości, jakie one generują[1]. Każdej obliczonej wartości przypisywany jest pewien typ, który jednoznacznie definiuje, jakie operacje można na niej wykonać. Śledząc przepływ wartości, system typów stara się udowodnić, że w programie występuje poprawne typowanie, tzn. nie dochodzi do sytuacji, w której na wartości określonego typu próbujemy wykonać niedozwoloną operację.

2.4 System typów ML

– silny system typów stosowany w językach rodziny ML (Ocaml, Standard ML) oparty na inferencji.

Podstawowy system typów jest następujący: istnieją typy proste, takie jak string, int, bool, unit (typ pusty) itd. Z dowolnych typów można też generować typy złożone – przez krotki (typ1 * typ2, typ1 * typ2 * typ3 itd.), konstruktory typów (typ list, typ tree itd.) i funkcje (typ1 → typ2).

System próbuje nadać typy każdemu wyrażeniu języka, i nie licząc kilku rzadkich przypadków, udaje mu się to całkiem dobrze.

Generalnie system taki wyklucza polimorfizm (nie licząc typów polimorficznych), jednak w Standard ML stworzono specjalne reguły umożliwiające polimorfizm dla wyrażeń arytmetycznych.

System typów ML jest interesujący z teoretycznego punktu widzenia – wiele problemów ma bardzo wysoką złożoność, jednak w praktyce inferencja zachodzi bardzo szybko – typy, które są rzeczywiście używane, są zwykle bardzo proste – rzadko używa się funkcji rzędów wyższych niż trzeci-czwarty, oraz liczby argumentów większej niż kilkanaście.

W rzeczywistych implementacjach dochodzą do tego bardziej złożone problemy typizacji obiektów, modułów itd.

3 Inferencja typów w teorii

3.1 System typów Hindley–Milner

3.2 Algorytm W

4 Rescript/ReasonML jako języki realizujące podobne zadania

5 Założenia i priorytety opracowanej aplikacji

Tworząc aplikację, chciałem, by język posiadał podstawowe typy danych (liczby, stringi, wartość logiczna), kilka typów generycznych (funkcje, tablice), typ ‘Option’, oraz możliwość tworzenia własnych typów. Dodatkowo nie powinno być potrzeby podawania typów zmiennych w większości przypadków, kompilator sam powinien wykrywać typy zmiennych na podstawie ich użycia.

Język poza zmiennymi, potrzebuje możliwości wykonywania operacji na danych, dlatego ważne było dla mnie, by zaimplementować operatory binarne, oraz unarne. Operatory te miały też spełniać ważną rolę w trakcie inferencji typów. W języku javascript operator ‘+’ może być wykorzystywany do dodawania liczb jak i konkatencji stringów, ważne więc było by stworzyć dwa oddzielne operatory.

Prymitywne typy danych:

- string

```
1 greeting = 'Hello world!'
```

examples/string.uwu

Do konkatencji stringów służy operator ++

- Wartość logiczna ma typ `Bool` i wartość `True` lub `False`. Powiązane operacje:
 - `<>`, równość pomiędzy dwiema liczbami
 - `>`, `<`
 - `!=` równość
- liczby Powiązane operacje: `+`, `-`, `*`, `/`, `*`, `%`, `//`

Chciałem również by funkcje wieloargumentowe kompilowane były jako funkcje jednoargumentowe zwracające kolejne funkcje, co pozwala na wywoływanie funkcji bez wszystkich argumentów w celu zwrócenia funkcji przyjmującej resztę argumentów tzw. currying.

```
1 def classy_greeting(first_name, last_name) do
2   "The name's " ++ last_name ++ ', ' ++ first_name ++ ' ' ++
    last_name
3 end
4
5 def compose1(f, g, a) do
6   f(g(a))
7 end
8
9 def compose2(f, g, a, b) do
10  f(g(a, b))
11 end
12
13 yell_greetings = compose2(to_upper, classy_greeting)
14 yell_greetings('James', 'Bond') # "THE NAME'S BOND, JAMES BOND"
15
16 compose1(compose1(abs, add(1)), multiply(2))(-4) # 7
```

examples/currying.uwu

examples/currying.uwu.js

Jednym z ważniejszych elementów każdego języka jest możliwość wykonywania różnego zbioru instrukcji, warunkowo. W tym celu planowałem zaimplementowanie instrukcji 'if', oraz 'case'. Instrukcja 'case' wykonywać ma dopasowanie do wzorca (tzw. pattern-matching), wykonywać, odpowiedni zbiór instrukcji zależnie od wprowadzonych danych. Kompilator, powinien ostrzegać, jeżeli ścieżka dla jednego z typów danych nie została zaimplementowana.

- Przykład - funkcja łącząca dwie posortowane tablice

```
1 def merge<A>(a, b) do
2   merge2: Callable<Array<A>, Callable<Array<A>, Array<A>>> =
3     'merge'
4
5   case Tuple(get_head(a), get_head(b)) of
6     Tuple(Empty(), _) do b end
7     Tuple(_, Empty()) do a end
8     Tuple(Head(head_a, rest_a), Head(head_b, rest_b)) do
9       if head_a < head_b then
10        [head_a] | merge2(rest_a, b)
11      else
12        [head_b] | merge2(a, rest_b)
13      end
14    end
15  end
16 end
```

examples/merge.uwu

examples/merge.uwu.js

Kolejnym dość ważnym elementem języka jest brak wyrażenia ‘return’, które jest wykorzystywane do zwrócenia wartości z funkcji. Zamiast tego każdy blok instrukcji powinien zwracać ostatnie wyrażenie. Pozwoli to na łatwiejsze inicjowanie zmiennych, w przypadku gdy inicjalizacja wymaga więcej niż jednej linii kodu.

- Przykład

```
1 message = if is_morning then
2   'Good morning!'
3 else
4   'Hello!'
5 end
6
7
8 result = do
9   arr1 = [1, 2, 3]
10  arr2 = map(arr1, add_one)
11
12  filter(arr2, def is_even(x) do
13    x % 2 == 0
14  end)
15 end
```

examples/return.uwu

examples/return.uwu.js

5.1 Opis formalny składni języka

6 Narzędzia

6.1 Język python

6.2 Parsowanie i tokenizowanie przy użyciu biblioteki sly

6.3 Środowisko nodejs do uruchomienia skompilowanego kodu

7 Implementacja

7.1 lexer

7.2 parser

7.3 inferencja typów

7.4 kompilacja

8 Opis działania

8.1 Co działa

8.2 Uwagi co do obsługi błędów

9 [Podsumowanie]

10 [spisy – rysunków, tabel, listingów itp.]

References

- [1] J.E. Hopcroft, R. Motwani, J.D. Ullman: Wprowadzenie do Teorii automatów, języków i obliczeń. Warszawa: Wydawnictwo Naukowe PWN, 2005. ISBN 83-01-14502-1.
- [2] Alfred V. Aho, Ravi Sethi, Jeffrey D. Ullman: Kompilatory. Wydawnictwo Naukowe PWN SA, 2002. ISBN 83-204-2656-1.
- [3] Aho, A.V., Sethi, R. and Ullman, J.D. (1986) " Compilers: principles, techniques, and tools." Addison-Wesley Longman Publishing Co., Inc. Boston, MA, USA.
- [4] Sikkels, Klaas, 1954- (1997). Parsing schemata : a framework for specification and analysis of parsing algorithms. Berlin: Springer. ISBN 9783642605413. OCLC 606012644.