

Kompilator języka strukturalnego

David Korenchuk

August, 23, 2023

Spis treści

1	Wprowadzenie	3
2	Teoria	4
2.1	Języki formalne	4
2.2	Klasyfikacja gramatyczna	4
3	Analiza leksykalna	5
3.1	Wyrażenia regularne	5
3.2	Flex	6
4	Analiza syntaksyczna	7
4.1	Definicja	7
4.2	Znane problemy	7
4.3	Implementacja AST	8
4.4	Implementacja analizatora składniowego	8
5	Analiza semantyczna	9
5.1	Analiza nieużywanych zmiennych	9
5.2	Analiza typów	10
5.2.1	Matematyczny opis systemu typów	10
6	Generacja warstwy pośredniej	12
7	Interpreter	13
8	Annex: Gramatyka w BNF	14

1 Wprowadzenie

W dniu dzisiejszym istnieje wiele języków programowania. Przyczyna na istnienie narzędzia takiego rodzaju jest taka, że człowiek myśli i mówi zdaniami. Aby móc przeprowadzić ludzkie zdania do formy, zrozumiałej do komputera, niezbędnie te zdania muszą być przeprowadzone do zestawu dużo prostszych zdań, niż w ludzkich językach.

...

2 Teoria

2.1 Języki formalne

Według teorii automatów, automat – jest to jednostka wykonawcza. Jednostki te, zależnie od swojej struktury i tego, jaki **język formalny** oni mogą obrobić, dzielą się na klasy.

Klasy te opisane są **hierarchią Chomsky’ego**. Mówi ona o tym, że języki formalne dzielą się na 4 typy:

- Typ 3 – języki regularne
- Typ 2 – języki bezkontekstowe
- Typ 1 – języki kontekstowe
- Typ 0 – języki rekurencyjnie przeliczalne

Jako przykład języka typu 3 według hierarchii Chomsky’ego można podać wyrażenia regularne. Język ten opisuje się automatem skończonym deterministycznym (DFA). Bardziej szczegółowo wyrażenia regularne będą rozpatrzone w opisanu analizy leksykalnej.

2.2 Klasyfikacja gramatyczna

Niniejszy język nie może być odniesiony do żadnej z klas hierarchii Chomsky’ego, chociaż jest on językiem regularnym. Tak jest dlatego, że można napisać gramatycznie poprawny kod, który jednak prowadzi do błędów kontekstowych i logicznych. Naprzykład

```
void f() {  
    return argument + 1;  
}
```

Kolejną z przyczyn niemożliwości odniesienia naszego języka do jednej z klas hierarchii Chomsky’ego jest niejednoznaczność konstrukcji językowych. Przykład niżej pokazuje, że nie można jednoznacznie stwierdzić, czy `data * d` jest deklaracją zmiennej albo operatorem mnożenia dwóch zmiennych. Aby móc poprawnie prowadzić analizę składniową, musimy zadbać o rozróżnienie kontekstu.

```
void f() {  
    data *d;  
}
```

3 Analiza leksykalna

Jednym ze sposobów na sprowadzanie kodu źródłowego do postaci listy tokenów jest narzędzie flex. Przyjmuje ono zestaw reguł w postaci wyrażeń regularnych, według których działa rozbięcie tekstu wejściowego. Można jednak ominąć lex i zaimplementować lexer ręcznie, ale ta praca nie skupia się na tym.

3.1 Wyrażenia regularne

Wyrażenie regularne – łańcuch znaków, zawierający pewne polecenia do wyszukiwania tekstu.

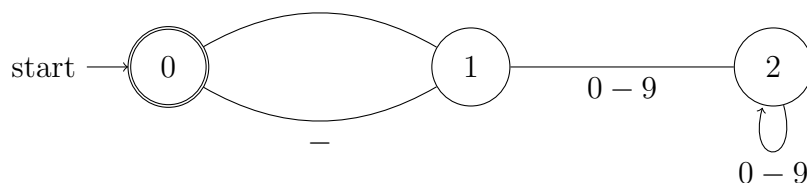
Mówimy, że wyrażenie regularne określone nad alfabetem Σ , jeżeli zachodzą następujące warunki:

- \emptyset – wyrażenie regularne, reprezentujące pusty zbiór.
- ϵ – wyrażenie regularne, reprezentujące pusty łańcuch.
- $\forall a \in \Sigma$, a reprezentuje jeden znak.
- Warunek indukcyjny: jeżeli R_1, R_2 – wyrażenia regularne, $(R_1 R_2)$ stanowi konkatencję R_1 i R_2 .
- Warunek indukcyjny: jeżeli R – wyrażenie regularne, R^* stanowi domknięcie Kleene'ego.

W rzeczywistości, takich zasad może być więcej.

Zazwyczaj wyrażenie regularne jest realizowane za pomocą DFA (Deterministic finite automaton, Deterministyczny automat skończony). Lex sprowadza podany zbiór zasad do takiego automatu.

Podamy przykład automatu dla wyrażenia $-?[0-9]^+$



Aby odśledzić wykonane kroki, można wypełnić tabelę przejść pomiędzy stanami. Podamy przykład dla łańcucha -22

Bieżący stan	Akcja
0	zaakceptować -
0, 1	zaakceptować 2
0, 1, 2	zaakceptować 2

3.2 Flex

Flex jest narzędziem projektu GNU. Pozwala ono w wygodny sposób podać listę reguł dla analizy leksykalnej (ang. Scanning). Flex jest mocno powiązany z językiem C, dlatego program w flex'u korzysta z konstrukcji języka C. Pokażemy przykład użycia flex'u

Listing 1: Przykład użycia flex

```
%{
#include "portrzebny-do-analizy-plik.h"

/* Kod w języku C. */
%}

/* Opcje flex */
%option noyywrap nounput noinput
%option yylineno

%% /* Reguły w postaci wyrażen regularnych. */

/*****
/* Wzorzec | Akcja przy znalezieniu takiego wzorcu */
*****/
-?[0-9]+ LEX_CONSUME_WORD(TOK_INTEGRAL_LITERAL)
-?[0-9]+\.[0-9]+ LEX_CONSUME_WORD(TOK_FLOATING_POINT_LITERAL)
\"([^\\"\\]*(\\.[^\\"\\]*)*)\" LEX_CONSUME_WORD(TOK_STRING_LITERAL)
\'.\\' LEX_CONSUME_WORD(TOK_CHAR_LITERAL)

. { /* Znaleziony niewiadomy znak.
    Zglosic blad.
    */ }

%%
```

Zauważmy, że flex próbuje szukać wzorców w tekście dokładnie w takiej kolejności, która jest podana w jego kodzie. Dlatego często robią ostatnią regułę z wyrażeniem regularnym ".", który akceptuje dowolny znak, i umieszczają tam komunikat o błędzie.

W naszym przypadku, lex generuje kod, który gromadzi wszystkie znalezione lexemy do tablicy.

4 Analiza syntaktyczna

4.1 Definicja

Mając listę składników elementarnych wejściowego programu, jesteśmy w stanie przejść do następnego etapu kompilacji – analizy składniowej. Jest to proces generacji struktury drzewiastej, a mianowicie AST (Abstract Syntax Tree).

AST może być stworzony po zdefiniowaniu gramatyki regularnej danego języka. Stosuje się do tego notacja BNF (Backus–Naur form). Pełny opis gramatyki pokazany jest w końcu pracy. Pokażemy tylko kilka przykładów:

$$\begin{aligned} \langle program \rangle & ::= (\langle function-decl \rangle \mid \langle structure-decl \rangle)^* \\ \langle var-decl \rangle & ::= \langle type \rangle (*)^* \langle id \rangle = \langle logical-or-stmt \rangle ; \\ \langle stmt \rangle & ::= \langle block-stmt \rangle \\ & \quad \mid \langle selection-stmt \rangle \\ & \quad \mid \langle iteration-stmt \rangle \\ & \quad \mid \langle jump-stmt \rangle \\ & \quad \mid \langle decl \rangle \\ & \quad \mid \langle expr \rangle \\ & \quad \mid \langle assignment-stmt \rangle \\ & \quad \mid \langle primary-stmt \rangle \end{aligned}$$

4.2 Znane problemy

Projektując gramatykę, należy wziąć pod uwagę problem rekurencji lewej (Left recursion). Są produkcje gramatyczne, nie pozwalające kodu, które je implementuje przejść do następnego terminalu, stosując tą samą produkcję, co prowadzi do rekurencji nieskończonej.

Rekurencja lewa może wyglądać następująco:

$$\langle factor \rangle ::= \langle factor \rangle '+' \langle term \rangle$$

Kod, wykonujący tą regułę będzie miał postać:

Listing 2: Rekurencja lewa

```
void factor() {
    factor(); // Rekurencja bez zadnego warunku wyjścia
    consume('+');
    term();
}
```

4.3 Implementacja AST

Zaimplementowany AST składa się ze struktury `ast_node`. Jest to główny typ węzła, zawierający niektóre zbędne informacje dla każdego typu węzła AST, i przechowujący konkretny węzeł jako wskaźnik.

Listing 3: Główny węzeł AST

```
struct ast_node {
    enum ast_type  type;    /* Rozrozniamy typ wedlug tej flagi */
    void          *ast;    /* ast_num, ast_for, ast_while, et cetera */
    uint16_t      line_no;
    uint16_t      col_no;
};
```

Konkretne węzły definiujemy w następujący sposób:

Listing 4: Konkretny węzeł AST

```
struct ast_num {
    int32_t value;
};
```

Taki AST stanowi strukturę drzewiastą, mającą wszystkie zalety i wady drzew jako struktur danych. Mając takie drzewo, jesteśmy w stanie prowadzić zwykłe przeszukiwanie w głąb i wszerz. Dokładnie w ten sposób działa każda z przedstawionych niżej analiz semantycznych oraz generacja kodu pośredniego.

Algorithm 1 Przeszukiwanie AST

```
1: procedure DFS(AST)
2:   for each child node Child of AST do
3:     DFS(Child)
4:   end for
5: end procedure
```

4.4 Implementacja analizatora składniowego

W danym przypadku, analizator składniowy jest napisany ręcznie, chociaż są narzędzia od projektu GNU, takie jak GNU Bison i UNIX'owe, takie jak YACC. Niniejszy analizator jest napisany bez pomocy tych programów, aby jawnie pokazać, jak się przekładają produkcje BNF na język C.

Aby poradzić sobie z zadaniem pisania takiego analizatora, możemy zauważyć, że zadanie to sprowadza się do implementacji każdej produkcji gramatycznej osobno.

5 Analiza semantyczna

Aby zapewnić poprawność napisanego kodu, stosuje się wiele rodzajów analiz. Niniejszy kompilator dysponuje trzema:

- Analiza nieużytych zmiennych, oraz zmiennych, które są zdefiniowane, ale nie zostały użyte
- Analiza poprawności typów
- Analiza prawidłowego użycia funkcji

5.1 Analiza nieużywanych zmiennych

Podamy przykłady kodu prowadzącego do odpowiednich ostrzeżeń

Listing 5: Nieużywana zmienna

```
void f() {
    int argument = 0; // Warning: unused variable 'argument'
}
```

Listing 6: Nieużywany parametr

```
void f(int argument) {} // Warning: unused variable 'argument'
```

Listing 7: Nieodczytana zmienna

```
void f() {
    int argument = 0;
    ++argument; // Warning: variable 'argument' written, but never read
}
```

Rzecz polega na przejściu drzewa syntaksycznego i zwiększania liczników `read_uses` i `write_uses` dla każdego węzła typu `ast_sym`.

Algorytm operuje na blokach kodu, zawartego w `{ ... }`. Po przejściu każdego bloku (w tym rekurencyjnie), analiza jest wykonana w następujący sposób:

Algorithm 2 Wyszukiwanie nieużywanych zmiennych

```
1: procedure ANALYZE(AST)
2:   Set ← all declarations at current scope depth
3:   for each collected declaration Use in Set do
4:     if Use is not a function & Use.ReadUses is 0 then
5:       Emit warning
6:     end if
7:   end for
8: end procedure
```

Do analizy nieużywanych funkcji stosuje się ten sam algorytm. Jedyne, co jest wtedy zmienione – sprawdzenie, czy nazwa rozpatrywanej funkcji nie jest **main**. Funkcja **main** jest wywołana automatycznie.

5.2 Analiza typów

Podczas analizy typów należy sprawdzić, czy:

- Oba operandy wyrażenia binarnego mają ten sam typ
- Faktycznie zwracana z funkcji wartość zgadza się z jej deklaracją
- Prawidłowa ilość i typy argumentów były przekazane do wyrażenia wywołania funkcji
- Rozmiar zadeklarowanej tablicy jest dopuszczalny
- Indeks tablicy nie wykracza poza jej rozmiar (w przypadku stałego indeksu)

Analiza typów w naszym przypadku jest dość prosta i tak samo polega na przejściu drzewa syntaktycznego.

5.2.1 Matematyczny opis systemu typów

$$\frac{\Gamma \vdash \diamond}{\Gamma \vdash \text{true} : \text{bool}}$$

$$\frac{\Gamma \vdash \diamond}{\Gamma \vdash \text{false} : \text{bool}}$$

$$\frac{\Gamma \vdash \diamond}{\Gamma \vdash n : \text{int}}$$

$$\frac{\Gamma \vdash \diamond}{\Gamma \vdash c : \text{char}}$$

$$\frac{\Gamma \vdash \diamond}{\Gamma \vdash x : \text{float}}$$

Oznaczmy tutaj dla $\mathbb{N}, \mathbb{R} : \oplus \in \{+, -, *, /, <, >, \leq, \geq, ==, \neq, ||, \&\&\}$, wtedy

$$\frac{\Gamma \vdash e_1 : \text{float} \quad \Gamma \vdash e_2 : \text{float}}{\Gamma \vdash e_1 \oplus e_2 : \text{float}}$$

Dodamy do \oplus operacje tylko dla $\mathbb{N} : \oplus \cup \{||, \&, ^, <<, >>, \% \}$, wtedy

$$\frac{\Gamma \vdash e_1 : \text{int} \quad \Gamma \vdash e_2 : \text{int}}{\Gamma \vdash e_1 \oplus e_2 : \text{int}}$$

$$\frac{\Gamma \vdash \text{condition} : \text{bool} \quad \Gamma \vdash e_1 : \tau \quad \Gamma \vdash e_2 : \tau}{\Gamma \vdash \text{if (condition) } \{ e_1 \} \text{ else } \{ e_2 \} : \tau}$$

$$\frac{\Gamma \vdash \text{condition} : \text{bool} \quad \Gamma \vdash e : \tau}{\Gamma \vdash \text{while (condition) } \{ e \} : \tau}$$

$$\frac{\Gamma \vdash \text{condition} : \text{bool} \quad \Gamma \vdash e : \tau}{\Gamma \vdash \text{do } \{ e \} \text{ while (condition) : } \tau}$$

$$\frac{\Gamma(x) = \tau}{\Gamma \vdash x : \tau}$$

6 Generacja warstwy pośredniej

7 Interpreter

8 Annex: Gramatyka w BNF

$\langle \text{program} \rangle$	$::= (\langle \text{function-decl} \rangle \mid \langle \text{structure-decl} \rangle)^*$
$\langle \text{structure-decl} \rangle$	$::= \text{struct } \{ \langle \text{structure-decl-list} \rangle \}$
$\langle \text{structure-decl-list} \rangle$	$::= (\langle \text{decl-without-initialiser} \rangle ;$ $\mid \langle \text{structure-decl} \rangle ;)^*$
$\langle \text{function-decl} \rangle$	$::= \langle \text{ret-type} \rangle \langle \text{id} \rangle (\langle \text{parameter-list-opt} \rangle) \{ \langle \text{stmt} \rangle^* \}$
$\langle \text{ret-type} \rangle$	$::= \langle \text{type} \rangle$ $\mid \langle \text{void-type} \rangle$
$\langle \text{type} \rangle$	$::= \text{int}$ $\mid \text{float}$ $\mid \text{char}$ $\mid \text{string}$ $\mid \text{boolean}$
$\langle \text{void-type} \rangle$	$::= \text{void}$
$\langle \text{constant} \rangle$	$::= \langle \text{integral-literal} \rangle$ $\mid \langle \text{floating-literal} \rangle$ $\mid \langle \text{string-literal} \rangle$ $\mid \langle \text{char-literal} \rangle$ $\mid \langle \text{boolean-literal} \rangle$
$\langle \text{integral-literal} \rangle$	$::= \langle \text{digit} \rangle^*$
$\langle \text{floating-literal} \rangle$	$::= \langle \text{digit} \rangle^* . \langle \text{digit} \rangle^*$
$\langle \text{string-literal} \rangle$	$::= \text{' ' (x00000000-x0010FFFF) * ' '}$
$\langle \text{char-literal} \rangle$	$::= \text{'ASCII(0)-ASCII(127)'}$
$\langle \text{boolean-literal} \rangle$	$::= \text{true}$ $\mid \text{false}$
$\langle \text{alpha} \rangle$	$::= \text{a} \mid \text{b} \mid \dots \mid \text{z} \mid _$
$\langle \text{digit} \rangle$	$::= 0 \mid 1 \mid \dots \mid 9$
$\langle \text{id} \rangle$	$::= \langle \text{alpha} \rangle (\langle \text{alpha} \rangle \mid \langle \text{digit} \rangle)^*$
$\langle \text{array-decl} \rangle$	$::= \langle \text{type} \rangle (*)^* \langle \text{id} \rangle [\langle \text{integral-literal} \rangle]$
$\langle \text{var-decl} \rangle$	$::= \langle \text{type} \rangle (*)^* \langle \text{id} \rangle = \langle \text{logical-or-stmt} \rangle ;$
$\langle \text{structure-var-decl} \rangle$	$::= \langle \text{id} \rangle (*)^* \langle \text{id} \rangle$

$\langle decl \rangle$	$::=$	$\langle var-decl \rangle$ $ $ $\langle array-decl \rangle$ $ $ $\langle structure-var-decl \rangle$
$\langle decl-without-initialiser \rangle$	$::=$	$\langle type \rangle (*) * \langle id \rangle$ $ $ $\langle array-decl \rangle$ $ $ $\langle structure-var-decl \rangle$
$\langle parameter-list \rangle$	$::=$	$\langle decl-without-initialiser \rangle , \langle parameter-list \rangle$ $ $ $\langle decl-without-initialiser \rangle$
$\langle parameter-list-opt \rangle$	$::=$	$\langle parameter-list \rangle \mid \epsilon$
$\langle stmt \rangle$	$::=$	$\langle block-stmt \rangle$ $ $ $\langle selection-stmt \rangle$ $ $ $\langle iteration-stmt \rangle$ $ $ $\langle jump-stmt \rangle$ $ $ $\langle decl \rangle$ $ $ $\langle expr \rangle$ $ $ $\langle assignment-stmt \rangle$ $ $ $\langle primary-stmt \rangle$
$\langle member-access-stmt \rangle$	$::=$	$\langle id \rangle . \langle member-access-stmt \rangle$ $ $ $\langle id \rangle . \langle id \rangle$
$\langle iteration-stmt \rangle$	$::=$	$\langle stmt \rangle$ $ $ $\text{break};$ $ $ $\text{continue};$
$\langle block-stmt \rangle$	$::=$	$\{ \langle stmt \rangle * \}$
$\langle iteration-block-stmt \rangle$	$::=$	$\{ \langle iteration-stmt \rangle * \}$
$\langle selection-stmt \rangle$	$::=$	$\text{if } (\langle expr \rangle) \langle block-stmt \rangle$ $ $ $\text{if } (\langle expr \rangle) \langle block-stmt \rangle \text{ else } \langle block-stmt \rangle$
$\langle iteration-stmt \rangle$	$::=$	$\text{for } (\langle expr-opt \rangle ; \langle expr-opt \rangle ; \langle expr-opt \rangle) \langle iteration-block-stmt \rangle$ $ $ $\text{while } (\langle expr \rangle) \langle iteration-block-stmt \rangle$ $ $ $\text{do } \langle iteration-block-stmt \rangle \text{ while } (\langle expr \rangle) ;$
$\langle jump-stmt \rangle$	$::=$	$\text{return } \langle expr \rangle ? ;$
$\langle assignment-op \rangle$	$::=$	$=$ $ $ $*$ $ $ $/$ $ $ $\%$ $ $ $+$ $ $ $-$ $ $ $<<=$ $ $ $>>=$

	$\&=$
	$ =$
	$\wedge=$
$\langle expr \rangle$	$::= \langle assignment-stmt \rangle$ $\langle var-decl \rangle$
$\langle expr-opt \rangle$	$::= \langle expr \rangle \mid \epsilon$
$\langle assignment-stmt \rangle$	$::= \langle logical-or-stmt \rangle$ $\langle logical-or-stmt \rangle \langle assignment-op \rangle \langle assignment-stmt \rangle$
$\langle logical-or-stmt \rangle$	$::= \langle logical-and-stmt \rangle$ $\langle logical-and-stmt \rangle \parallel \langle logical-or-stmt \rangle$
$\langle logical-and-stmt \rangle$	$::= \langle inclusive-or-stmt \rangle$ $\langle inclusive-or-stmt \rangle \&\& \langle logical-and-stmt \rangle$
$\langle inclusive-or-stmt \rangle$	$::= \langle exclusive-or-stmt \rangle$ $\langle exclusive-or-stmt \rangle \mid \langle inclusive-or-stmt \rangle$
$\langle exclusive-or-stmt \rangle$	$::= \langle and-stmt \rangle$ $\langle and-stmt \rangle \wedge \langle exclusive-or-stmt \rangle$
$\langle and-stmt \rangle$	$::= \langle equality-stmt \rangle$ $\langle equality-stmt \rangle \& \langle and-stmt \rangle$
$\langle equality-stmt \rangle$	$::= \langle relational-stmt \rangle$ $\langle relational-stmt \rangle == \langle equality-stmt \rangle$ $\langle relational-stmt \rangle != \langle equality-stmt \rangle$
$\langle relational-stmt \rangle$	$::= \langle shift-stmt \rangle$ $\langle shift-stmt \rangle > \langle relational-stmt \rangle$ $\langle shift-stmt \rangle < \langle relational-stmt \rangle$ $\langle shift-stmt \rangle >= \langle relational-stmt \rangle$ $\langle shift-stmt \rangle <= \langle relational-stmt \rangle$
$\langle shift-stmt \rangle$	$::= \langle additive-stmt \rangle$ $\langle additive-stmt \rangle << \langle shift-stmt \rangle$ $\langle additive-stmt \rangle >> \langle shift-stmt \rangle$
$\langle additive-stmt \rangle$	$::= \langle multiplicative-stmt \rangle$ $\langle multiplicative-stmt \rangle + \langle additive-stmt \rangle$ $\langle multiplicative-stmt \rangle - \langle additive-stmt \rangle$
$\langle multiplicative-stmt \rangle$	$::= \langle prefix-unary-stmt \rangle$ $\langle prefix-unary-stmt \rangle * \langle multiplicative-stmt \rangle$ $\langle prefix-unary-stmt \rangle / \langle multiplicative-stmt \rangle$ $\langle prefix-unary-stmt \rangle \% \langle multiplicative-stmt \rangle$

$\langle \text{prefix-unary-stmt} \rangle$	$::=$	$\langle \text{postfix-unary-stmt} \rangle$ $++ \langle \text{postfix-unary-stmt} \rangle$ $-- \langle \text{postfix-unary-stmt} \rangle$ $* \langle \text{postfix-unary-stmt} \rangle$ $\& \langle \text{postfix-unary-stmt} \rangle$ $! \langle \text{postfix-unary-stmt} \rangle$
$\langle \text{postfix-unary-stmt} \rangle$	$::=$	$\langle \text{primary-stmt} \rangle$ $\langle \text{primary-stmt} \rangle ++$ $\langle \text{primary-stmt} \rangle --$
$\langle \text{primary-stmt} \rangle$	$::=$	$\langle \text{constant} \rangle$ $\langle \text{symbol-stmt} \rangle$ $(\langle \text{logical-or-stmt} \rangle)$
$\langle \text{symbol-stmt} \rangle$	$::=$	$\langle \text{function-call-stmt} \rangle$ $\langle \text{array-access-stmt} \rangle$ $\langle \text{member-access-stmt} \rangle$ $\langle \text{id} \rangle$
$\langle \text{array-access-stmt} \rangle$	$::=$	$\langle \text{id} \rangle ([\langle \text{expr} \rangle])^*$
$\langle \text{function-call-arg-list} \rangle$	$::=$	$\langle \text{logical-or-stmt} \rangle , \langle \text{function-call-arg-list} \rangle$ $\langle \text{logical-or-stmt} \rangle$
$\langle \text{function-call-arg-list-opt} \rangle$	$::=$	$\langle \text{function-call-arg-list} \rangle \mid \epsilon$
$\langle \text{function-call-expr} \rangle$	$::=$	$\langle \text{id} \rangle (\langle \text{function-call-arg-list-opt} \rangle)$

Literatura

- [1] <https://www.bates.edu/biology/files/2010/06/How-to-Write-Guide-v10-2014.pdf>
- [2] <http://lucacardelli.name/papers/typesystems.pdf>