

自然言語処理 parser 試作に関する報告書

島田 歩[†]

[†]ayumu.shimada@gmail.com

1 始めに

本実験報告書は以下の章立てで構成されている:

- 第 1 章 始めに
- 第 2 章 実験の目的
- 第 3 章 自然言語処理の概要
- 第 4 章 実験の結果
- 第 5 章 考察
- 第 6 章 まとめ

また、巻末には、参考文献リストと付録および手書きの図がある。

特に、第 5 章の考察では、文の解釈を一意に特定できない問題と非文が生成される問題について述べている。

2 実験の目的

自然言語処理において、基本技術の一つになっている構文解析技法について、その原理を学ぶ¹⁾。また、自然言語処理における構文解析の困難さについても実際に体験する。

3 自然言語処理の概要

自然言語処理 (natural language processing: NLP) とは、人間が書いたり話したりした文章をコンピュータに処理させることである¹⁾。これは、文章の中の単語の品詞を決定するための形態素解析、単語間のつながりを調べる構文解析、単語や文の意味を解析する意味解析および必要な応答文を生成する等の処理を含む¹⁾。

また、Chomsky によれば、言語は、形式言語を用いて階層化できる (図 1)。なお、正規文法は有限オートマトンで、文脈自由文法 (context-free grammar: CFG) はプッシュダウンオートマトンで、それぞれ認識できる。

3.1 有限オートマトン

有限オートマトンは入力記号列と現在の状態に応じて次の状態へ遷移する。入力記号列を入力し終えたときに

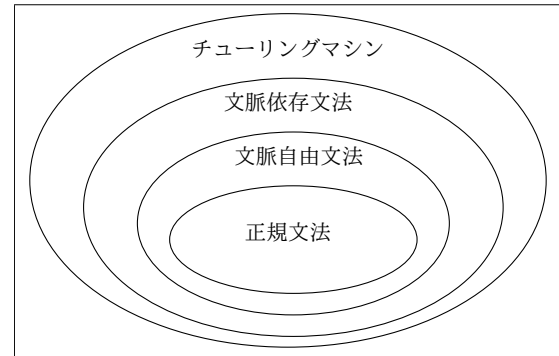


図 1 Chomsky の言語階層¹⁾

受理状態にあれば、その入力記号列は受理可能である。

3.2 プッシュダウンオートマトン

有限オートマトンに補助記憶装置としてプッシュダウンスタックを取り付けたものがプッシュダウンオートマトンである。入力記号列を入力し終えたときに受理状態にあり、かつ、スタックが空であれば、その入力記号列は受理可能である。

3.3 あいまい性

文脈自由文法では、構文木が複数得られることがある。この場合、その文法はあいまいであると言う。

4 実験の結果

4.1 ネットワーク文法の実験

4.1.1 手による解析

表 1 の文法規則と表 2 の単語辞書を用いて “I saw a girl with a telescope .” を形態素解析した結果を図 2 に示す。

4.1.2 プログラムによる解析

文法規則と単語辞書に表 1, 表 2 および表 3 を用い、作成したプログラムで “I saw a girl with a telescope .” と “the child runs quickly to the large house .” の 2 文を形態素解析した結果を示す。

4.2 文脈自由文法の実験

4.2.1 手による解析

表 4 の文法規則と表 2 の単語辞書を用いて top-down 法および CYK 法により構文解析した結果を、それぞ

¹⁾2017 年 12 月 1 日 仮提出; 2017 年 12 月 6 日 受理;
2017 年 12 月 12 日 更新

	I	saw	a	girl	with	a	telescope	.	
start	NOUN	VERB	DET	NOUN	PREP	DET	NOUN		end
	the	child	runs	quickly	to	the	large	house	.
start	DET	NOUN	VERB	ADV		PREP	DET	ADJ	NOUN
									end

図 3 プログラムによる形態素解析結果

表 1 ネットワーク文法の文法規則¹⁾

start	→	DET (冠詞)
start	→	NOUN (名詞)
DET (冠詞)	→	NOUN (名詞)
DET (冠詞)	→	ADJ (形容詞)
DET (冠詞)	→	NOUN (名詞)
ADJ (形容詞)	→	NOUN (名詞)
NOUN (名詞)	→	PREP (前置詞)
PREP (前置詞)	→	NOUN (名詞)
PREP (前置詞)	→	ADJ (形容詞)
PREP (前置詞)	→	DET (冠詞)
NOUN (名詞)	→	VERB (動詞)
VERB (動詞)	→	ADV (副詞)
VERB (動詞)	→	DET (冠詞)
ADV (副詞)	→	PREP (前置詞)
NOUN (名詞)	→	ADV (副詞)
ADV (副詞)	→	NOUN (名詞)
VERB (動詞)	→	end
NOUN (名詞)	→	end

表 3 “the child runs quickly to the large house.” の単語辞書

the	→	DET (冠詞)
child	→	NOUN (名詞)
runs	→	VERB (動詞)
quickly	→	ADV (副詞)
to	→	PREP (前置詞)
large	→	ADJ (形容詞)
house	→	NOUN (名詞)
.	→	end

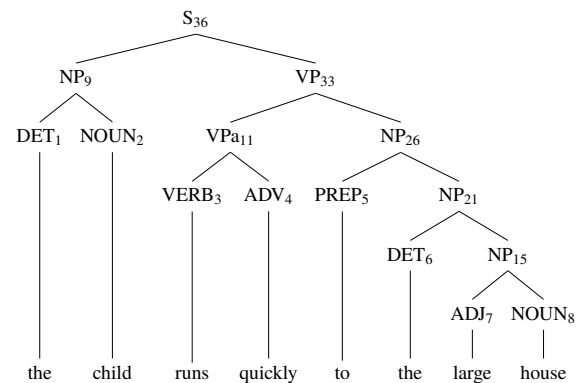


図 6 “the child runs quickly to the large house” の構文木

表 2 “I saw a girl with a telescope.” の単語辞書

I	→	NOUN (名詞)
saw	→	VERB (動詞)
a	→	DET (冠詞)
girl	→	NOUN (名詞)
with	→	PREP (前置詞)
telescope	→	NOUN (名詞)
.	→	end

れ図 4, 図 5 に示す。

4.2.2 プログラムによる解析

表 5 の文法規則と表 3 の単語辞書を用い、作成したプログラムで “the child runs quickly to the large house” を CYK 法により構文解析した結果得られた構文木と S 式をそれぞれ図 6 と図 7 に示す。

さらに、表 4 の文法規則と表 2 の単語辞書を用い、“I saw a girl with a telescope” を構文解析した結果得られた構文木と S 式をそれぞれ図 8a と図 7 に示す。

5 考察

5.1 “I saw a girl with a telescope.” の 2 つの解釈

“I saw a girl with a telescope.” には大きくわけて、以下の 2 つの意味（解釈）がある：

- 私は望遠鏡を持った少女に会った。（図 8a）
- 私は望遠鏡で少女を見た。（図 8b）

今回の実験の結果得られた解釈は前者だが、表 4 の文法規則に $VP \rightarrow VP + PP$ を追加することにより、後者の解釈となる構文木（図 8b）と S 式（図 9）が得られる。

表 4 “I saw a girl with a telescope” の文法規則

S (文)	→	NP (名詞句)	+	VP (動詞句)
NP (名詞句)	→	NOUN (名詞)		
NP (名詞句)	→	DET (冠詞)	+	NOUN (名詞)
VP (動詞句)	→	VERB (動詞)		
VP (動詞句)	→	VERB (動詞)	+	NP (名詞句)
VP (動詞句)	→	VERB (動詞)	+	SS
SS	→	NP (名詞句)	+	PP (前置詞句)
PP (前置詞句)	→	PREP (前置詞)	+	NP (名詞句)

表 5 “the child runs quickly to the large house” の文法規則

S (文)	→	NP (名詞句)	+	VP (動詞句)
NP (名詞句)	→	DET (冠詞)	+	NOUN (名詞)
NP (名詞句)	→	DET (冠詞)	+	NP (名詞句)
NP (名詞句)	→	PREP (前置詞)	+	NP (名詞句)
NP (名詞句)	→	ADJ (形容詞)	+	NOUN (名詞)
VP (動詞句)	→	VERB (動詞)	+	NP (名詞句)
VP (動詞句)	→	VPa (動詞句)	+	NP (名詞句)
VP (動詞句)	→	VERB (動詞)	+	ADV (副詞)
VP (動詞句)	→	VERB (動詞)		
VPa (動詞句)	→	VERB (動詞)	+	ADV (副詞)
VPa (動詞句)	→	VERB (動詞)	+	NP (名詞句)

(S36(NP9(DET1 "the")(NOUN2 "child"))(VP33(VPa11(VERB3 "runs")(ADV4 "quickly"))(NP26(PREP5 "to")(NP21(DET6 "the")(NP15(ADJ7 "large")(NOUN8 "house"))))))))

(S28(NP1(NOUN1 "I"))(VP27(VERB2 "saw")(SS25(NP10(DET3 "a")(NOUN4 "girl"))(PP18(PREP5 "with")(NP13(DET6 "a")(NOUN7 "telescope"))))))))

図 7 プログラムによる構文解析結果（印刷の関係上，出力は改行している。）

なお，複数の S 式を得るためには，合致する文法規則を見つけるとそこで中止して次のセルに移動するような実装では不可能で，必ず全ての文法規則についてチェックを行い，合致する文法規則を全て記憶しておく必要がある。本プログラムはそれを考慮した実装になっているため，複数の S 式が出力可能である。

しかしながら，複数の S 式が出力されるということは文の解釈を一意に特定できないということであり，この問題の解決は困難であると考えられる。

5.2 “Time flies like an arrow.” の解釈

表 6 の文法規則と表 7 の単語辞書を用い “Time flies like an arrow” を構文解析した結果，4 つの構文木（図 10）と S 式（図 11）が得られた。これらの解釈はそれぞれ以下である：

- 時は矢のように過ぎる。⇒ 光陰矢のごとし。（図 10a）
- time fly（トキバエ？）たちは矢を好む。（図 10b）
- “矢のようなハエ”たちの速度を測定せよ。（図 10c）
- 矢（の速度を測定するとき）のように，ハエたちの速度を測定せよ。（図 10d）

さらに，“Time” が雑誌の名前である可能性も考慮すると，さらに別の解釈も可能である。

実際には，「光陰矢のごとし」以外の候補はいわゆる非文とみなせるが，候補の絞り込みには別の方法を用いる必要性が示唆される。

5.3 ネットワーク文法と文脈自由文法の能力の比較

第 3 章で触れたように，形式文法は階層構造をなしており，文脈自由文法はネットワーク文法を含む。今回作

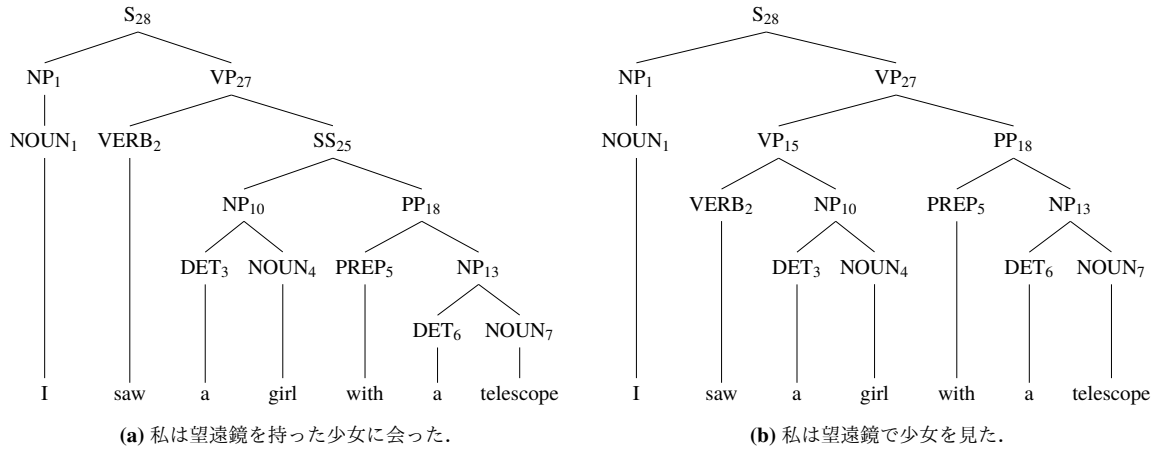


図 8 “I saw a girl with a telescope” の構文木

(S28(NP1(NOUN1 "I"))(VP27(VP15(VERB2 "saw")(NP10(DET3 "a")(NOUN4 "girl")))(PP18(PREP5 "with")(NP13(DET6 "a")(NOUN7 "telescope")))))

図 9 私は望遠鏡で少女を見た（図 8b）という解釈になる S 式（印刷の関係上，出力は改行している。）

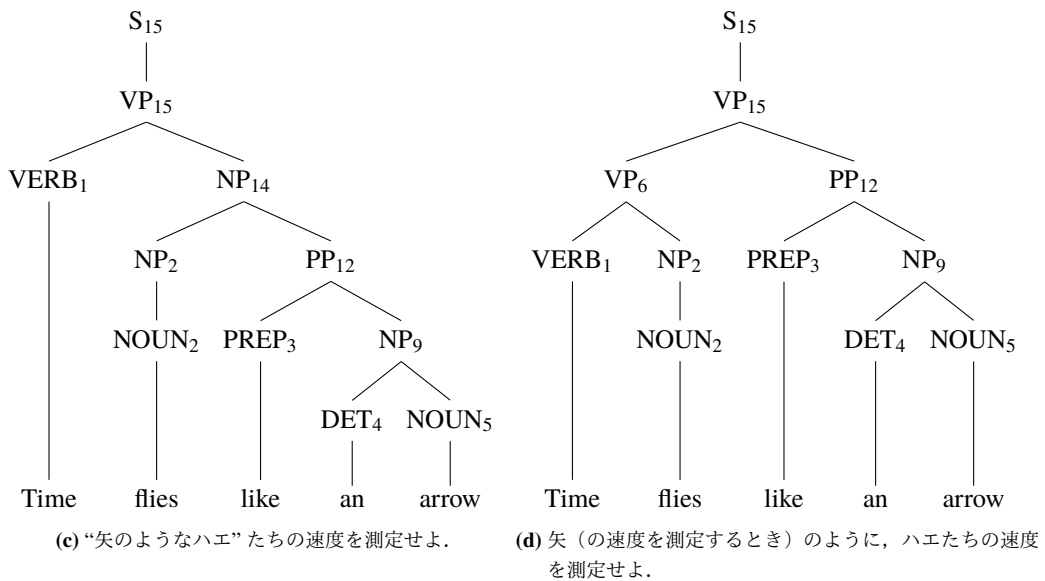
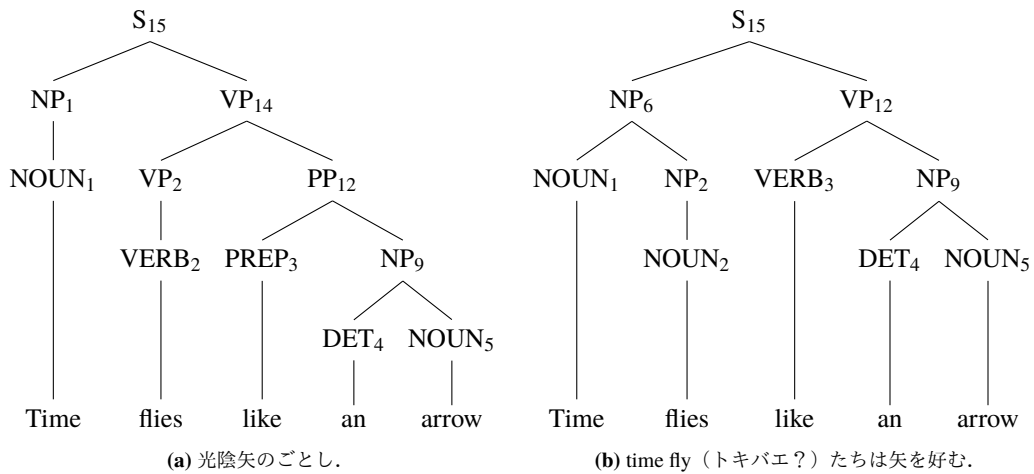


図 10 “Time flies like an arrow” の構文木

```

(S15(NP1(NOUN1 "Time"))(VP14(VP2(VERB2 "flies"))
  (PP12(PREP3 "like")(NP9(DET4 "an")(NOUN5 "arrow")))))

(S15(NP6(NOUN1 "Time")(NP2(NOUN2 "flies"))
  (VP12(VERB3 "like")(NP9(DET4 "an")(NOUN5 "arrow"))))

(S15(VP15(VERB1 "Time")(NP14(NP2(NOUN2 "flies"))
  (PP12(PREP3 "like")(NP9(DET4 "an")(NOUN5 "arrow")))))

(S15(VP15(VP6(VERB1 "Time")(NP2(NOUN2 "flies"))
  (PP12(PREP3 "like")(NP9(DET4 "an")(NOUN5 "arrow")))))

```

図 11 “Time flies like an arrow” の S 式（印刷の関係上，出力は改行している。）

表 6 “Time flies like an arrow” の文法規則

S	→	NP	+	VP
S	→	VP		
NP	→	NOUN		
NP	→	DET	+	NOUN
NP	→	NOUN	+	NP
NP	→	NP	+	PP
VP	→	VERB		
VP	→	VERB	+	NP
VP	→	VP	+	PP
PP	→	PREP	+	NP

表 7 “Time flies like an arrow” の単語辞書

Time	→	NOUN
Time	→	VERB
flies	→	NOUN
flies	→	VERB
like	→	PREP
like	→	VERB
an	→	DET
arrow	→	NOUN

成する際も，形態素解析の部分にはネットワーク文法のプログラムを流用したことから，これが確認できる。

ネットワーク文法においても，遷移規則（文法規則）によって簡易的に構文解析ができたが，構文木（S 式）まで得るためには文脈自由文法を用いる必要がある。

しかしながら，文脈自由文法も，今回の考察でも分かったように，文法規則を作成することそのものも困難であるし，複数の構文木が得られてあいまいさが残る問題もある。

5.4 実験の感想

C 言語での実装は比較的大変であったが，オートマトンと構文解析について理解が深まり，本実験の目的は達せられたのではないと思う。

6 まとめ

構文解析そのものはプログラムにより実現できたが，構文木が複数得られることから，文の解釈を一意に特定することや非文の生成を防ぐことの困難さが示唆された。

参考文献

- 1) 実験手順書.

成したプログラムにおいても，CYK 法による parser を作

付録 A ネットワーク文法のプログラム

ソースコード 1 Makefile

```

1 CC = gcc
2 CFLAGS = -Wall
3
4 all: main
5
6 main: main.o dict_builder.o syntax_builder.o splitter.o post.o network.o output.o
7     $(CC) $(CFLAGS) -o $@ $^
8
9 %.o: %.c global.h
10    $(CC) $(CFLAGS) -c $<
11
12 clean:
13     $(RM) *.o
14
15 run: all
16     cat data.txt | ./main dict.csv syntax.csv > result.txt

```

ソースコード 2 global.h

```

1 #ifndef GLOBAL_H
2 #define GLOBAL_H
3
4 #include <stdio.h>
5 #include <stdlib.h>
6 #include <string.h>
7
8 #define bool int
9 #define true 1
10 #define false 0
11
12 typedef struct {
13     char *word;
14     char *type;
15 } WORD;
16
17 typedef struct {
18     char *current;
19     char *next;
20 } SYNTAX;
21
22 extern WORD START;
23 extern WORD DUMMY;
24 extern SYNTAX DUMMY_SYNTAX;
25
26 bool type_analyzer(WORD *words, const WORD *dict);
27 bool syntax_analyzer(const WORD *words, const SYNTAX *syntax);
28 WORD* sentence_splitter(char *sentence);
29 WORD* build_dict(const char *srcfile);
30 SYNTAX* build_syntax(const char *srcfile);
31 void print_words(const WORD *words, const size_t buf_size);
32
33 #endif

```

ソースコード 3 main.c

```

1 #include "global.h"
2
3 WORD START;
4 WORD DUMMY;
5 SYNTAX DUMMY_SYNTAX;
6

```

```

7
8 int main(int argc, char **argv) {
9     const size_t DEFAULT_LINE_LENGTH = 1024;
10
11     size_t max_line_length;
12     if (argc > 3) {
13         max_line_length = strtoul(argv[3], NULL, 10);
14     } else {
15         max_line_length = DEFAULT_LINE_LENGTH;
16     }
17     fprintf(stderr, "[INFO]: up to %ld characters are allowed for each sentence\n",
18         max_line_length - 1);
19
20     fprintf(stderr, "[INFO]: each sentence must be separated by newline character '\n'");
21     fprintf(stderr, "[INFO]: each word must be separated by whitespace characters\n");
22
23     // init global variables
24     DUMMY.word = "";
25     DUMMY.type = "?";
26     DUMMY_SYNTAX.current = "";
27     DUMMY_SYNTAX.next = "";
28     START.word = "_";
29     START.type = "start";
30
31     WORD *dict = build_dict(argv[1]);
32     SYNTAX *syntax = build_syntax(argv[2]);
33
34     WORD *words = NULL;
35     char *sentence = malloc(sizeof(char) * max_line_length);
36
37     while (fgets(sentence, max_line_length, stdin) != NULL) {
38         fprintf(stderr, "[INFO]: Parsing started.\n");
39
40         words = sentence_splitter(sentence);
41         if (words == &DUMMY) {
42             fprintf(stderr, "\x1b[31m[WARN]: Parsing stopped.\x1b[39m\n");
43             fprintf(stderr, "\x1b[31m[WARN]: There is nothing to process. (maybe empty line?)\x1b[39m\n");
44             fprintf(stderr, "[INFO]: continue...\n");
45             words = NULL;
46             continue;
47         }
48
49         bool typeOK = type_analyzer(words, dict);
50         print_words(words, max_line_length);
51
52         if (typeOK == true) {
53             if (syntax_analyzer(words, syntax) == true) {
54                 fprintf(stderr, "[INFO]: \x1b[42mSentence OK\x1b[49m\x1b[32m (sentence is acceptable.)\x1b[39m\n");
55             } else {
56                 fprintf(stderr, "\x1b[31m[WARN]: \x1b[39m\x1b[41mSentence NG\x1b[49m\x1b[31m (sentence is not acceptable.)\x1b[39m\n");
57             }
58             fprintf(stderr, "[INFO]: Parsing finished.\n");
59         } else {
60             fprintf(stderr, "\x1b[31m[WARN]: can't start syntax checking because of unknown words.\x1b[39m\n");
61             fprintf(stderr, "\x1b[31m[WARN]: Parsing aborted.\x1b[39m\n");
62         }
63
64         free(words);
65     }
66
67     fprintf(stderr, "[INFO]: exiting...\n");

```

```

67     fflush(stderr);
68
69     /*
70     // no need?
71     free(sentence);
72     free(syntax[0].current);
73     free(syntax);
74     free(dict[0].word);
75     free(dict);
76     */
77
78     fprintf(stderr, "_Goodbye.\n");
79
80     return 0;
81 }

```

ソースコード 4 dict_builder.c

```

1  #include "global.h"
2
3
4  WORD* build_dict(const char *srcfile) {
5      const int FGETS_MARGIN = 2; // do NOT decrease to < 2 (do not change)
6
7      fprintf(stderr, "[INFO]:_building_dictionary_database...");
8      fflush(stderr);
9
10     FILE *fp;
11     if ((fp = fopen(srcfile, "r")) == NULL) {
12         fprintf(stderr, "\n[ERROR:]_Failed_to_open_dictionary_file.\n");
13         exit(EXIT_FAILURE);
14     }
15
16     // get file size
17     fseek(fp, 0, SEEK_END);
18     int size = ftell(fp) + FGETS_MARGIN;
19     fseek(fp, 0, SEEK_SET);
20
21     // checking csv and copying to memory
22     char *bin = malloc(sizeof(char) * size);
23     int count = 0;
24     char *ptr;
25     int remain;
26     int entry = 0;
27     char c;
28     for (ptr = bin, remain = size; fgets(ptr, remain, fp) != NULL; /**/) {
29         if (sscanf(ptr, "%[^,],%s%c\n", &c, &count) == 1) { // csv format checking
30             entry++;
31             // count = strlen(ptr);
32             ptr[count - 1] = '\0'; // %c (CR or LF (or space?)) -> '\0'
33             ptr += count;
34             remain -= count;
35         } else {
36             // do nothing
37         }
38     }
39
40     WORD *dict = malloc(sizeof(WORD) * (entry + 1));
41     int n[4] = {0};
42     int i = 0;
43     for (ptr = bin; sscanf(ptr, "%n%[^,],%n,%n%s\n", &n[0], &n[1], &n[2], &n[3]) !=
44         EOF; ptr += n[3] + 1) {
45         // fprintf(stderr, "%d, %d, %d, %d\n", n[0], n[1], n[2], n[3]);
46         dict[i].word = ptr + n[0];
47         dict[i].type = ptr + n[2];
48         i++;

```



```

48     ptr[n[1]] = '\0';
49     ptr[n[3]] = '\0';
50 }
51 dict[i] = DUMMY;
52
53 fclose(fp);
54 fprintf(stderr, "_done.\n");
55 fprintf(stderr, "[INFO]:_Dictionary_database_has_%d_entries.\n", i);
56
57 /*
58 // write out to dict.bin (for debug)
59 if ((fp = fopen("dict.bin", "wb")) == NULL) {
60     fprintf(stderr, "\n[ERROR:] Failed to write dictionary file.\n");
61     exit(EXIT_FAILURE);
62 }
63 fwrite(dict, sizeof(WORD), entry + 1, fp);
64 fwrite(bin, sizeof(char), size, fp);
65 fclose(fp);
66 */
67
68 return dict;
69 }

```

ソースコード 5 syntax_builder.c

```

1  #include "global.h"
2
3
4  SYNTAX* build_syntax(const char *srcfile) {
5      const int FGETS_MARGIN = 2; // do NOT decrease to < 2 (do not change)
6
7      fprintf(stderr, "[INFO]:_building_syntax_database...");
8      fflush(stderr);
9
10     FILE *fp;
11     if ((fp = fopen(srcfile, "r")) == NULL) {
12         fprintf(stderr, "\n[ERROR]:_Failed_to_open_syntax_file.\n");
13         exit(EXIT_FAILURE);
14     }
15
16     // get file size
17     fseek(fp, 0, SEEK_END);
18     int size = ftell(fp) + FGETS_MARGIN;
19     fseek(fp, 0, SEEK_SET);
20
21     // checking csv and copying to memory
22     char *bin = malloc(sizeof(char) * size);
23     int count = 0;
24     char *ptr;
25     int remain;
26     int entry = 0;
27     char c;
28     for (ptr = bin, remain = size; fgetc(ptr, remain, fp) != NULL; /**/) {
29         if (sscanf(ptr, "%[^,]%s%c\n", &c, &count) == 1) { // csv format checking
30             entry++;
31             // count = strlen(ptr);
32             ptr[count - 1] = '\0'; // %c (CR or LF (or space?)) -> '\0'
33             ptr += count;
34             remain -= count;
35         } else {
36             ; // do nothing
37         }
38     }
39
40     SYNTAX *syntax = malloc(sizeof(SYNTAX) * (entry + 1));
41     int n[4] = {0};

```

```

42  int i = 0;
43  for (ptr = bin; sscanf(ptr, "%n%*[^,]%n,%n%*s%n", &n[0], &n[1], &n[2], &n[3]) !=
      EOF; ptr += n[3] + 1) {
44      // fprintf(stderr, "%d, %d, %d, %d\n", n[0], n[1], n[2], n[3]);
45      syntax[i].current = ptr + n[0];
46      syntax[i].next    = ptr + n[2];
47      i++;
48      ptr[n[1]] = '\0';
49      ptr[n[3]] = '\0';
50  }
51  syntax[i] = DUMMY_SYNTAX;
52
53  fclose(fp);
54  fprintf(stderr, "\done.\n");
55  fprintf(stderr, "[INFO]: Syntax database has %d entries.\n", i);
56
57  /*
58  // write out to syntax.bin (for debug)
59  if ((fp = fopen("syntax.bin", "wb")) == NULL) {
60      fprintf(stderr, "\n[ERROR:] Failed to write syntax file.\n");
61      exit(EXIT_FAILURE);
62  }
63  fwrite(syntax, sizeof(SYNTAX), entry + 1, fp);
64  fwrite(bin, sizeof(char), size, fp);
65  fclose(fp);
66  */
67
68  return syntax;
69 }

```

ソースコード 6 splitter.c

```

1  #include "global.h"
2
3
4  WORD* sentence_splitter(char *sentence) {
5      char *ptr;
6      int n[2] = {0};
7
8      // check wordage
9      int wordage = 0;
10     for (ptr = sentence; sscanf(ptr, "%*s%n", &n[1]) != EOF; ptr += n[1]) {
11         wordage++;
12     }
13     if (wordage == 0) {
14         return &DUMMY;
15     }
16
17     // memory allocation
18     WORD *words = malloc(sizeof(WORD) * (wordage + 2));
19     // wordage + 2?: last word will be DUMMY; first word will be START.
20
21     // word separation
22     int i = 1; // words[0] (i = 0) is START
23     n[1] = 0;
24     for (ptr = sentence; sscanf(ptr, "%n%*s%n", &n[0], &n[1]) != EOF; ptr += n[1] +
        1) {
25         words[i].word = ptr + n[0];
26         i++;
27         if (ptr[n[1]] == '\0') { // if reach the end of sentence
28             break;
29         } else {
30             ptr[n[1]] = '\0';
31         }
32     }
33 }

```

```

34 // head and tail editing
35 words[0] = START;
36 words[i] = DUMMY;
37
38 return words;
39 }

```

ソースコード 7 post.c

```

1 #include "global.h"
2
3
4 bool type_analyzer(WORD *words, const WORD *dict) {
5     fprintf(stderr, "[INFO]: Part-of-Speech tagging (POST) has been executed.\n");
6
7     int i, j;
8     bool typeFound = false;
9     bool noUnknown = true;
10    for (i = 1; words[i].word != DUMMY.word; i++) { // words[0] is START
11        fprintf(stderr, "[DEBUG]: %s => ", words[i].word); fflush(stderr);
12        for (j = 0; dict[j].word != DUMMY.word; j++) {
13            if (strcmp(words[i].word, dict[j].word) == 0) {
14                fprintf(stderr, "%s\n", dict[j].type); // [DEBUG]:
15                words[i].type = dict[j].type;
16                typeFound = true;
17                break;
18            }
19        }
20
21        if (typeFound == true) {
22            typeFound = false;
23        } else {
24            fprintf(stderr, "\x1b[41m?\x1b[49m\n"); // [DEBUG]:
25            fprintf(stderr, "\x1b[31m[ERROR]: Unknown word \"%s\" detected.\x1b[39m\n",
26                words[i].word);
27            noUnknown = false;
28            words[i].type = DUMMY.type;
29        }
30    }
31    fprintf(stderr, "[INFO]: POST finished.\n");
32    return noUnknown;
33 }

```

ソースコード 8 network.c

```

1 #include "global.h"
2
3
4 bool syntax_analyzer(const WORD *words, const SYNTAX *syntax) {
5     fprintf(stderr, "[INFO]: Syntax checking started.\n");
6
7     int i, j;
8     bool transOK = false;
9     for (i = 0; words[i+1].type != DUMMY.type; i++) {
10        fprintf(stderr, "[DEBUG]: %s (%s) -> %s (%s) => ", words[i].type, words[i].word,
11            words[i+1].type, words[i+1].word); fflush(stderr);
12        for (j = 0; syntax[j].current != DUMMY_SYNTAX.current; j++) {
13            if (strcmp(words[i].type, syntax[j].current) == 0) {
14                if (strcmp(words[i+1].type, syntax[j].next) == 0) {
15                    fprintf(stderr, "\x1b[42mOK\x1b[49m\n"); // [DEBUG]:
16                    transOK = true;
17                    break;
18                }
19            }
20        }
21    }
22 }

```

```

20
21     if (transOK == true) {
22         transOK = false;
23     } else {
24         fprintf(stderr, "\x1b[41mNG\x1b[49m\n"); // [DEBUG]:
25         fprintf(stderr, "\x1b[31m[WARN]:\x1b[39mSyntax\x1b[39merror\x1b[39mdetected.\x1b[39m\n");
26         return false;
27     }
28 }
29
30 fprintf(stderr, "[INFO]:\x1b[39mSyntax\x1b[39mchecking\x1b[39mfinished.\x1b[39m\n");
31 return true;
32 }

```

ソースコード 9 output.c

```

1  #include "global.h"
2
3
4  void print_words(const WORD *words, const size_t buf_size) {
5      const int WORD_SEPARATION = 2;
6      const int LINES = 3;
7
8      int i;
9
10     // memory allocation
11     char *line[LINES];
12     line[0] = malloc(sizeof(char) * buf_size * LINES); // get heap area
13
14     for (i = 1; i < LINES; i++) {
15         line[i] = line[0] + i * buf_size;
16     }
17
18     // initializing counters and lines
19     int count[LINES];
20     for (i = 0; i < LINES; i++) {
21         line[i][0] = '\0';
22         count[i] = 0;
23     }
24
25     // padding
26     int max = 0;
27     int j, k;
28     int remain = buf_size;
29     for (i = 0; words[i].word != DUMMY.word; i++) {
30         count[0] += snprintf(line[0] + count[0], remain, "%s", words[i].word);
31         count[1] += snprintf(line[1] + count[1], remain, "|");
32         count[2] += snprintf(line[2] + count[2], remain, "%s", words[i].type);
33
34         // getting the longest elements
35         for (j = 0; j < LINES; j++) {
36             if (count[j] > buf_size - 1) { // minus 1?: last element must be terminated
37                 // with '\0'
38                 // Because snprintf doesn't return the actual number of written characters,
39                 // count[] may exceeds buf_size.
40                 count[j] = buf_size - 1;
41             }
42             if (max < count[j]) {
43                 max = count[j];
44             }
45         }
46
47         // padding with spaces
48         for (j = 0; j < LINES; j++) {
49             while (count[j] < max) {
50                 line[j][count[j]++] = ' ';
51             }
52         }
53     }
54 }

```

```

50     }
51 }
52 // now, count[] == max
53
54 remain = buf_size - max;
55
56 if (remain > WORD_SEPARATION) {
57     for (j = 0; j < LINES; j++) {
58         for (k = 0; k < WORD_SEPARATION; k++) {
59             line[j][count[j]++] = '␣';
60         }
61     }
62     remain -= WORD_SEPARATION;
63 } else {
64     break;
65 }
66 }
67
68 // terminating and printing each line
69 for (i = 0; i < LINES; i++) {
70     if (buf_size - count[i] > WORD_SEPARATION) {
71         line[i][count[i] - WORD_SEPARATION] = '\0'; // remove last padding
72     } else {
73         line[i][count[i]] = '\0';
74     }
75     printf("%s\n", line[i]);
76 }
77
78 free(line[0]);
79 return;
80 }

```

A.1 一文の長さの制限変更について

初期状態では 1 文の長さは 1023 文字に制限されているが、再コンパイルの必要なく、第 3 引数に n を渡すことで $n - 1$ 文字までの制限に変更することができる。

これは文脈自由文法のプログラムにおいても同様である。

```

user@CYNTHIA: ~/jikken3/c1
user@CYNTHIA:~$ cd ~/jikken3/c1/
user@CYNTHIA:~/jikken3/c1$ cat data.txt | ./network dict.csv syntax.csv 1024
[INFO]: up to 1023 characters are allowed for each sentence
[INFO]: each sentence must be separated by newline character '\n'
[INFO]: each word must be separated by whitespace characters
[INFO]: building dictionary database... done.
[INFO]: Dictionary database has 15 entries.
[INFO]: building syntax database... done.
[INFO]: Syntax database has 17 entries.
[INFO]: Parsing started.
[INFO]: Part-of-Speech tagging (POST) has been executed.
[DEBUG]: the => DET
[DEBUG]: child => NOUN
[DEBUG]: runs => VERB
[DEBUG]: quickly => ADV
[DEBUG]: to => PREP
[DEBUG]: the => DET
[DEBUG]: large => ADJ
[DEBUG]: house => NOUN
[DEBUG]: . => end
[INFO]: POST finished.
      the  child  runs  quickly  to    the  large  house  .
      |    |    |    |    |    |    |    |    |
start DET  NOUN  VERB  ADV    PREP  DET  ADJ   NOUN  end
[INFO]: Syntax cheking started.
[DEBUG]: start ( ) -> DET (the) => OK
[DEBUG]: DET (the) -> NOUN (child) => OK
[DEBUG]: NOUN (child) -> VERB (runs) => OK
[DEBUG]: VERB (runs) -> ADV (quickly) => OK
[DEBUG]: ADV (quickly) -> PREP (to) => OK
[DEBUG]: PREP (to) -> DET (the) => OK
[DEBUG]: DET (the) -> ADJ (large) => OK
[DEBUG]: ADJ (large) -> NOUN (house) => OK
[DEBUG]: NOUN (house) -> end (.) => OK
[INFO]: Syntax cheking finished.
[INFO]: Sentence OK (sentence is acceptable.)
[INFO]: Parsing finished.
[INFO]: Parsing started.
[INFO]: Part-of-Speech tagging (POST) has been executed.
[DEBUG]: I => NOUN
[DEBUG]: saw => VERB
[DEBUG]: a => DET
[DEBUG]: girl => NOUN
[DEBUG]: with => PREP
[DEBUG]: a => DET
[DEBUG]: telescope => NOUN
[DEBUG]: . => end
[INFO]: POST finished.
      I    saw  a    girl  with  a    telescope  .
      |    |    |    |    |    |    |
start NOUN  VERB  DET  NOUN  PREP  DET  NOUN      end
[INFO]: Syntax cheking started.
[DEBUG]: start ( ) -> NOUN (I) => OK
[DEBUG]: NOUN (I) -> VERB (saw) => OK
[DEBUG]: VERB (saw) -> DET (a) => OK
[DEBUG]: DET (a) -> NOUN (girl) => OK
[DEBUG]: NOUN (girl) -> PREP (with) => OK
[DEBUG]: PREP (with) -> DET (a) => OK
[DEBUG]: DET (a) -> NOUN (telescope) => OK
[DEBUG]: NOUN (telescope) -> end (.) => OK
[INFO]: Syntax cheking finished.
[INFO]: Sentence OK (sentence is acceptable.)
[INFO]: Parsing finished.
[INFO]: exiting... Goodbye.
user@CYNTHIA:~/jikken3/c1$

```

図 12 ネットワーク文法のプログラムの動作風景

付録 B 文脈自由文法のプログラム

ソースコード 10 Makefile

```

1 CC = gcc
2 CFLAGS = -Wall
3
4 all: main
5
6 main: main.o dict_builder.o syntax_builder.o splitter.o post.o cyk.o output.o
7     $(CC) $(CFLAGS) -o $@ $^
8
9 %.o: %.c global.h
10     $(CC) $(CFLAGS) -c $<
11
12 clean:
13     $(RM) *.o
14
15 run: all
16     cat child.txt | ./main dict.csv syntax_child.csv > result.txt
17     cat telescope.txt | ./main dict.csv syntax_telescope.csv >> result.txt
18     cat flies.txt | ./main dict.csv syntax_flies.csv >> result.txt

```

ソースコード 11 global.h

```

1 #ifndef GLOBAL_H
2 #define GLOBAL_H
3
4 #include <stdio.h>
5 #include <stdlib.h>
6 #include <string.h>
7
8 #define bool int
9 #define true 1
10 #define false 0
11
12 #define DEFAULT_LENGTH_OF_ARRAY 64
13
14 typedef struct {
15     char *word;
16     char **type;
17     bool type_is_new_array;
18 } WORD;
19
20 typedef struct { // ex.) NP -> DET + NOUN
21     char *lhs; // left-hand side (NP)
22     char *rhs1; // right-hand side (DET)
23     char *rhs2; // right-hand side (NOUN)
24 } SYNTAX;
25
26 typedef struct type {
27     int index;
28     char *type;
29     struct cell *cell;
30     struct type *prev_left_type;
31     struct type *prev_right_type;
32 } TYPE;
33
34 typedef struct cell {
35     int index;
36     char *word;
37     TYPE *type;
38     bool type_is_new_array;
39 } CELL;
40

```

```

41 extern WORD DUMMY;
42 extern SYNTAX DUMMY_SYNTAX;
43 extern CELL DUMMY_CELL;
44 extern int WORDAGE; // the number of words
45 extern int CELL_LENGTH; // the number of cells
46
47 bool type_analyzer(WORD *words, const WORD *dict);
48 CELL* syntax_analyzer(const WORD *words, const SYNTAX *syntax);
49 WORD* sentence_splitter(char *sentence);
50 WORD* build_dict(const char *srcfile);
51 SYNTAX* build_syntax(const char *srcfile);
52 void print_words(const WORD *words, const size_t buf_size, const int *type_num);
53 void print_sexps(const TYPE *type, int *type_num, int indent);
54
55 #endif

```

ソースコード 12 main.c

```

1  #include "global.h"
2
3  // global variables
4  WORD DUMMY;
5  SYNTAX DUMMY_SYNTAX;
6  CELL DUMMY_CELL;
7
8
9  int main(int argc, char **argv) {
10     size_t max_line_length;
11     if (argc > 3) {
12         max_line_length = strtoul(argv[3], NULL, 10);
13     } else {
14         max_line_length = 1024;
15     }
16     fprintf(stderr, "[INFO]: up to %ld characters are allowed for each sentence\n",
17             max_line_length - 1);
18     fprintf(stderr, "[INFO]: each sentence must be separated by newline character '\\n\\n');
19     fprintf(stderr, "[INFO]: each word must be separated by whitespace characters\n");
20
21     // init global variables
22     int i;
23     DUMMY.word = "";
24     DUMMY.type = (char *)"?";
25     DUMMY.type_is_new_array = false;
26     DUMMY_SYNTAX.lhs = "";
27     DUMMY_SYNTAX.rhs1 = "";
28     DUMMY_SYNTAX.rhs2 = "";
29
30     TYPE dummy_type;
31     DUMMY_CELL.index = -1;
32     DUMMY_CELL.word = DUMMY.word;
33     DUMMY_CELL.type = &dummy_type;
34     DUMMY_CELL.type_is_new_array = false;
35     dummy_type.index = -1;
36     dummy_type.type = DUMMY.type[0];
37     dummy_type.cell = &DUMMY_CELL;
38     dummy_type.prev_left_type = &dummy_type;
39     dummy_type.prev_right_type = &dummy_type;
40
41     WORD *dict = build_dict(argv[1]);
42     SYNTAX *syntax = build_syntax(argv[2]);
43
44     WORD *words = NULL;
45     char *sentence = malloc(sizeof(char) * max_line_length);
46

```



```

47 while (fgets(sentence, max_line_length, stdin) != NULL) {
48     fprintf(stderr, "[INFO]: Parsing started.\n");
49     words = sentence_splitter(sentence);
50     if (words == &DUMMY) {
51         fprintf(stderr, "\x1b[31m[WARN]: Parsing stopped.\x1b[39m\n");
52         fprintf(stderr, "\x1b[31m[WARN]: There is nothing to process. (maybe empty line?)\x1b[39m\n");
53         fprintf(stderr, "[INFO]: continue...\n");
54         words = NULL;
55         continue;
56     }
57
58     bool typeOK = type_analyzer(words, dict);
59     // print_words(words, max_line_length, NULL);
60
61     CELL *cells = NULL;
62     int *type_num = malloc(sizeof(int) * WORDAGE);
63     bool sentenceOK = false;
64     if (typeOK == true) {
65         cells = syntax_analyzer(words, syntax);
66         if (cells != NULL) {
67             for (i = 0; cells[CELL_LENGTH - 1].type[i].type != DUMMY_CELL.type[0].type; i++) {
68                 if (strcmp(cells[CELL_LENGTH - 1].type[i].type, "S") == 0) {
69                     sentenceOK = true;
70                     print_sexps(&(cells[CELL_LENGTH - 1].type[i]), type_num, 0);
71                     // printf("\n");
72                     // print_words(words, max_line_length, type_num);
73                     // printf("\n");
74                 }
75             }
76             if (sentenceOK == true) {
77                 fprintf(stderr, "[INFO]: \x1b[42mSentence OK\x1b[49m\x1b[32m (sentence is acceptable.)\x1b[39m\n");
78             } else {
79                 fprintf(stderr, "\x1b[31m[WARN]: \x1b[39m\x1b[41mSentence NG\x1b[49m\x1b[31m (sentence is not acceptable.)\x1b[39m\n");
80             }
81             fprintf(stderr, "[INFO]: Parsing finished.\n");
82         } else {
83             fprintf(stderr, "\x1b[31m[WARN]: Parsing aborted.\x1b[39m\n");
84         }
85     } else {
86         fprintf(stderr, "\x1b[31m[WARN]: can't start syntax checking because of unknown words.\x1b[39m\n");
87         fprintf(stderr, "\x1b[31m[WARN]: Parsing aborted.\x1b[39m\n");
88     }
89
90     for (i = 0; cells[i].word != DUMMY_CELL.word; i++) {
91         if (cells[i].type_is_new_array == true) {
92             free(cells[i].type);
93         }
94     }
95     free(cells);
96     free(type_num);
97     for (i = 0; words[i].word != DUMMY.word; i++) {
98         if (words[i].type_is_new_array == true) {
99             free(words[i].type);
100         }
101     }
102     free(words);
103 }
104
105 fprintf(stderr, "[INFO]: exiting...");
106 fflush(stderr);
107

```

```

108  /*
109  // no need?
110  free(sentence);
111  free(syntax[0].lhs);
112  free(syntax);
113  free(dict[0].word);
114  free(dict[0].type);
115  free(dict);
116  */
117
118  fprintf(stderr, "\Goodbye.\n");
119
120  return 0;
121  }

```

ソースコード 13 dict_builder.c

```

1  #include "global.h"
2
3
4  WORD* build_dict(const char *srcfile) {
5      const int FGETS_MARGIN = 2; // do NOT decrease to < 2 (do not change)
6
7      fprintf(stderr, "[INFO]:_building_dictionary_database...");
8      fflush(stderr);
9
10     FILE *fp;
11     if ((fp = fopen(srcfile, "r")) == NULL) {
12         fprintf(stderr, "\n[ERROR:]_Failed_to_open_dictionary_file.\n");
13         exit(EXIT_FAILURE);
14     }
15
16     // get file size
17     fseek(fp, 0, SEEK_END);
18     int size = ftell(fp) + FGETS_MARGIN;
19     fseek(fp, 0, SEEK_SET);
20
21     // checking csv and copying to memory
22     char *bin = malloc(sizeof(char) * size);
23     int count = 0;
24     char *ptr;
25     int remain;
26     int entry = 0;
27     char c;
28     for (ptr = bin, remain = size; fgetc(ptr, remain, fp) != NULL; /**/) {
29         if (sscanf(ptr, "%[^,],%s%c\n", &c, &count) == 1) { // csv format checking
30             entry++;
31             // count = strlen(ptr);
32             ptr[count - 1] = '\0'; // %c (CR or LF (or space?)) -> '\0'
33             ptr += count;
34             remain -= count;
35         } else {
36             // do nothing
37         }
38     }
39
40     WORD *dict = malloc(sizeof(WORD) * (entry + 1));
41     char **type = malloc(sizeof(char *) * entry);
42     int n[4] = {0};
43     int i = 0;
44     for (ptr = bin; sscanf(ptr, "%n%n%s%n", &n[0], &n[1], &n[2], &n[3]) !=
45         EOF; ptr += n[3] + 1) {
46         // fprintf(stderr, "%d, %d, %d, %d\n", n[0], n[1], n[2], n[3]);
47         dict[i].type = &type[i];
48         dict[i].word = ptr + n[0];
49         dict[i].type[0] = ptr + n[2];

```

```

49     i++;
50     ptr[n[1]] = '\0';
51     ptr[n[3]] = '\0';
52 }
53 dict[i] = DUMMY;
54 dict[0].type_is_new_array = true;
55
56 fclose(fp);
57 fprintf(stderr, "\done.\n");
58 fprintf(stderr, "[INFO]:_Dictionary_database_has_%d_entries.\n", i);
59
60 /*
61 // write out to dict.bin (for debug)
62 if ((fp = fopen("dict.bin", "wb")) == NULL) {
63     fprintf(stderr, "\n[ERROR:] Failed to write dictionary file.\n");
64     exit(EXIT_FAILURE);
65 }
66 fwrite(dict, sizeof(WORD), entry + 1, fp);
67 fwrite(bin, sizeof(char), size, fp);
68 fclose(fp);
69 */
70
71 return dict;
72 }

```

ソースコード 14 syntax_builder.c

```

1  #include "global.h"
2
3
4  SYNTAX* build_syntax(const char *srcfile) {
5      const int FGETS_MARGIN = 2; // do NOT decrease to < 2 (do not change)
6
7      fprintf(stderr, "[INFO]:_building_phrase_structure_rules_database...");
8      fflush(stderr);
9
10     FILE *fp;
11     if ((fp = fopen(srcfile, "r")) == NULL) {
12         fprintf(stderr, "\n[ERROR]:_Failed_to_open_syntax_file.\n");
13         exit(EXIT_FAILURE);
14     }
15
16     // get file size
17     fseek(fp, 0, SEEK_END);
18     int size = ftell(fp) + FGETS_MARGIN;
19     fseek(fp, 0, SEEK_SET);
20
21     // checking csv and copying to memory
22     char *bin = malloc(sizeof(char) * size);
23     int count = 0;
24     char *ptr;
25     int remain;
26     int entry = 0;
27     char c;
28     for (ptr = bin, remain = size; fgets(ptr, remain, fp) != NULL; /**/) {
29         // csv format checking
30         if (sscanf(ptr, "%[^,],%s%c\n", &c, &count) == 1 || // a,b,
31             sscanf(ptr, "%[^,],%[^,],%s%c\n", &c, &count) == 1) { // a,b,c
32             entry++;
33             // count = strlen(ptr);
34             ptr[count - 1] = '\0'; // %c (CR or LF (or space?)) -> '\0'
35             ptr += count;
36             remain -= count;
37         } else {
38             ; // do nothing
39         }

```

```

40 }
41
42 SYNTAX *syntax = malloc(sizeof(SYNTAX) * (entry + 1));
43 int n[6] = {0};
44 int i = 0;
45 ptr = bin;
46 for (;;) { // infinite loop
47     if (sscanf(ptr, "%n%[^,]%n,%n%[^,]%n,%n%s%n", &n[0], &n[1], &n[2], &n[3], &n
        [4], &n[5]) != EOF) {
48         // fprintf(stderr, "%d, %d, %d, %d\n", n[0], n[1], n[2], n[3], n[4], n[5]);
49         syntax[i].lhs = ptr + n[0];
50         syntax[i].rhs1 = ptr + n[2];
51         syntax[i].rhs2 = ptr + n[4];
52         i++;
53         ptr[n[1]] = '\0';
54         ptr[n[3]] = '\0';
55         ptr[n[5]] = '\0';
56         ptr += n[5] + 1;
57     } else
58     if (sscanf(ptr, "%n%[^,]%n,%n%s%n", &n[0], &n[1], &n[2], &n[3]) != EOF) {
59         // fprintf(stderr, "%d, %d, %d, %d\n", n[0], n[1], n[2], n[3]);
60         syntax[i].lhs = ptr + n[0];
61         syntax[i].rhs1 = ptr + n[2];
62         syntax[i].rhs2 = DUMMY_SYNTAX.rhs2;
63         i++;
64         ptr[n[1]] = '\0';
65         ptr[n[3]] = '\0';
66         if (ptr[n[3] - 1] == ',') {
67             // remove comma (rhs1 may include comma at its tail)
68             ptr[n[3] - 1] = '\0';
69         }
70         ptr += n[3] + 1;
71     } else {
72         break;
73     }
74 }
75 syntax[i] = DUMMY_SYNTAX;
76
77 fclose(fp);
78 fprintf(stderr, "\done.\n");
79 fprintf(stderr, "[INFO]:_Syntax_database_has_%d_phrase_structure_rules.\n", i);
80
81 /*
82 // write out to syntax.bin (for debug)
83 if ((fp = fopen("syntax.bin", "wb")) == NULL) {
84     fprintf(stderr, "\n[ERROR:] Failed to write syntax file.\n");
85     exit(EXIT_FAILURE);
86 }
87 fwrite(syntax, sizeof(SYNTAX), entry + 1, fp);
88 fwrite(bin, sizeof(char), size, fp);
89 fclose(fp);
90 */
91
92 return syntax;
93 }

```

ソースコード 15 splitter.c

```

1 #include "global.h"
2
3 int WORDAGE;
4
5
6 WORD* sentence_splitter(char *sentence) {
7     char *ptr;
8     int n[2] = {0};

```

```

9
10 // check wordage
11 WORDAGE = 0;
12 for (ptr = sentence; sscanf(ptr, "%s%n", &n[1]) != EOF; ptr += n[1]) {
13     WORDAGE++;
14 }
15 if (WORDAGE == 0) {
16     return &DUMMY;
17 }
18
19 // memory allocation
20 WORD *words = malloc(sizeof(WORD) * (WORDAGE + 1));
21 // WORDAGE + 1?: last word will be DUMMY.
22
23 // word separation
24 int i = 0;
25 n[1] = 0;
26 for (ptr = sentence; sscanf(ptr, "%n%s%n", &n[0], &n[1]) != EOF; ptr += n[1] +
27     1) {
28     words[i].word = ptr + n[0];
29     i++;
30     if (ptr[n[1]] == '\0') { // if reach the end of sentence
31         break;
32     } else {
33         ptr[n[1]] = '\0';
34     }
35 }
36 // head and tail editing
37 words[i] = DUMMY;
38
39 return words;
40 }

```

ソースコード 16 post.c

```

1 #include "global.h"
2
3
4 bool type_analyzer(WORD *words, const WORD *dict) {
5     fprintf(stderr, "[INFO]: Part-of-Speech tagging (POST) has been executed.\n");
6
7     int i, j, k;
8     bool typeFound = false;
9     bool noUnknown = true;
10
11     int array_size = DEFAULT_LENGTH_OF_ARRAY;
12     char **type = malloc(sizeof(char *) * array_size);
13     int type_index = 0;
14     for (i = 0; words[i].word != DUMMY.word; i++) {
15         outer: // outer loop
16         words[i].type = &type[type_index];
17         if (type_index == 0) {
18             words[i].type_is_new_array = true;
19         }
20         k = 0;
21         // fprintf(stderr, "[DEBUG]: %s => ", words[i].word); fflush(stderr);
22         for (j = 0; dict[j].word != DUMMY.word; j++) {
23             if (strcmp(words[i].word, dict[j].word) == 0) {
24                 typeFound = true;
25                 if (type_index >= (array_size - 1)) { // -1?: last element is preserved for
26                     dummy element.
27                     array_size *= 2; // new array will be expanded.
28                     fprintf(stderr, "[DEBUG]: words[i].type is full. increasing array size to %d elements.\n", array_size);
29                     if (k == type_index) {

```

```

29         free(type);
30     }
31     type = malloc(sizeof(char *) * array_size);
32     type_index = 0;
33     goto outer; // continue from outer loop
34 }
35 words[i].type[k] = dict[j].type[0];
36 // fprintf(stderr, "%s\n", words[i].type[k]); // [DEBUG]:
37 k++;
38 type_index++;
39 }
40 }
41 words[i].type[k] = DUMMY.type[0];
42 k++;
43 type_index++;
44
45 if (typeFound == true) {
46     typeFound = false;
47 } else {
48     fprintf(stderr, "[DEBUG]: %s=>", words[i].word);
49     fprintf(stderr, "\x1b[41m?\x1b[49m\n"); // [DEBUG]:
50     fprintf(stderr, "\x1b[31m[ERROR]: Unknown word \"%s\" detected.\x1b[39m\n",
51             words[i].word);
52     noUnknown = false;
53 }
54 }
55 fprintf(stderr, "[INFO]: POST finished.\n");
56 return noUnknown;
57 }

```

ソースコード 17 cyk.c

```

1  #include "global.h"
2
3  // global variables
4  int CELL_LENGTH;
5
6
7  int getIndex(const int x, const int y) {
8  // this method for upper triangular matrix
9      int index = 0;
10     int column = 0;
11     int i, j;
12     for (i = 0, j = 0; i < WORDAGE && j < WORDAGE; /* */) {
13         if (i == x && j == y) {
14             return index;
15         }
16
17         if (j == WORDAGE - 1) {
18             i = 0;
19             j = ++column;
20         } else {
21             i++;
22             j++;
23         }
24         index++;
25     }
26
27     return -1;
28 }
29
30
31 bool setNewArrayIfNeeded(TYPE **type, int *type_index, int *array_size, const int *k
32 ) {
33     if (*type_index >= (*array_size - 1)) { // -1?: last element is preserved for

```

```

    dummy element.
33  *array_size *= 2; // new array will be expanded.
34  fprintf(stderr, "[DEBUG]:_cells[index].type_is_full._increasing_array_size_to_%d
    _elements.\n", *array_size);
35  if (*k == *type_index) {
36      free(*type);
37  }
38  *type = malloc(sizeof(TYPE) * *array_size);
39
40  int i;
41  for (i = 0; i < *array_size; i++) { // init type[] with DUMMY_CELL.type[0]
42      memcpy(&(*type)[i], &DUMMY_CELL.type[0], sizeof(TYPE));
43  }
44
45  *type_index = 0;
46  return true;
47  } else {
48      return false;
49  }
50 }
51
52
53 CELL* syntax_analyzer(const WORD *words, const SYNTAX *syntax) {
54     fprintf(stderr, "[INFO]:_CYK_parser_started.\n");
55
56     // memory allocation
57     CELL_LENGTH = (WORDAGE * (WORDAGE + 1)) / 2; // upper triangular matrix (n*(n
        +1)/2)
58     if (CELL_LENGTH == 0) {
59         fprintf(stderr, "\x1b[31m[WARN]:_CYK_parser_stopped_(word_not_found).\x1b[39m\n"
            );
60         return NULL;
61     }
62     CELL *cells = malloc(sizeof(CELL) * CELL_LENGTH);
63
64     int index;
65
66     // init cells
67     for (index = 0; index < CELL_LENGTH; index++) {
68         cells[index] = DUMMY_CELL; // for unused-member initializing
69         cells[index].index = index;
70     }
71
72     int i, j; // for loop
73     int k; // cells[index].type[k]
74     int p, q, a; // comparison cells[a] (x, y) = (p, q)
75     int u, v, b; // comparison cells[b] (x, y) = (u, v)
76     int l, m, n; // for loop
77     int column = 0;
78     index = 0;
79     int array_size = DEFAULT_LENGTH_OF_ARRAY;
80     TYPE *type = malloc(sizeof(TYPE) * array_size);
81     for (i = 0; i < array_size; i++) { // init type[] with DUMMY_CELL.type[0]
82         memcpy(&type[i], &DUMMY_CELL.type[0], sizeof(TYPE));
83     }
84     int type_index = 0;
85     for (i = 0, j = column; i < WORDAGE && j < WORDAGE; /* */) {
86         outer: // outer loop
87         cells[index].type = &type[type_index];
88         if (type_index == 0) {
89             cells[index].type_is_new_array = true;
90         }
91         k = 0;
92
93         // (cells[index].word <= words[index].word)
94         if (index < WORDAGE) {

```

```

95     cells[index].word = words[index].word;
96     for (i = 0; words[index].type[i] != DUMMY.type[0]; i++) {
97         if (setNewArrayIfNeeded(&type, &type_index, &array_size, &k) == true) {
98             goto outer; // continue from outer loop
99         }
100
101         cells[index].type[k].index = k;
102         cells[index].type[k].type = words[index].type[i];
103         cells[index].type[k].cell = &cells[index];
104         fprintf(stderr, "[DEBUG]: %s%d\ \"%s\"\n", cells[index].type[k].type, cells[
            index].index + 1, cells[index].word);
105         k++;
106         type_index++;
107     }
108 }
109
110 // ex.) S -> NP + VP; VP -> VP + NP
111 p = i;
112 v = j;
113 for (q = i, u = i + 1; q < j && u <= j; q++, u++) {
114     a = getIndex(p, q);
115     b = getIndex(u, v);
116     // fprintf(stderr, "[DEBUG]: %d: (%d, %d) = %d, (%d, %d) = %d\n", index+1, p
        +1, q+1, a+1, u+1, v+1, b+1);
117
118     for (l = 0; cells[a].type[l].type != DUMMY_CELL.type[0].type; l++) {
119         for (m = 0; cells[b].type[m].type != DUMMY_CELL.type[0].type; m++) {
120             for (n = 0; syntax[n].lhs != DUMMY_SYNTAX.lhs; n++) {
121                 if (strcmp(cells[a].type[l].type, syntax[n].rhs1) == 0 && strcmp(cells[b]
                    .type[m].type, syntax[n].rhs2) == 0) {
122                     if (setNewArrayIfNeeded(&type, &type_index, &array_size, &k) == true)
123                     {
124                         goto outer; // continue from outer loop
125                     }
126
127                     cells[index].type[k].index = k;
128                     cells[index].type[k].type = syntax[n].lhs;
129                     cells[index].type[k].cell = &cells[index];
130                     cells[index].type[k].prev_left_type = &cells[a].type[l];
131                     cells[index].type[k].prev_right_type = &cells[b].type[m];
132                     fprintf(stderr, "[DEBUG]: %s%d->%s%d+%s%d\n",
                        cells[index].type[k].type, cells[index].type[k].cell->index + 1,
133                     cells[index].type[k].prev_left_type->type, cells[index].type[k].
                        prev_left_type->cell->index + 1,
134                     cells[index].type[k].prev_right_type->type, cells[index].type[k].
                        prev_right_type->cell->index + 1
135                 );
136                     k++;
137                     type_index++;
138                 }
139             }
140         }
141     }
142 }
143
144 // self-referencing (ex.) VP -> VERB; S -> VP
145 for (l = 0; cells[index].type[l].type != DUMMY_CELL.type[0].type; l++) {
146     for (n = 0; (syntax[n].lhs != DUMMY_SYNTAX.lhs); n++) {
147         if (strcmp(cells[index].type[l].type, syntax[n].rhs1) == 0 && syntax[n].rhs2
            == DUMMY_SYNTAX.rhs2) {
148             if (setNewArrayIfNeeded(&type, &type_index, &array_size, &k) == true) {
149                 goto outer; // continue from outer loop
150             }
151
152             cells[index].type[k].index = k;
153             cells[index].type[k].type = syntax[n].lhs;

```



```

154     cells[index].type[k].cell = &cells[index];
155     cells[index].type[k].prev_left_type = &cells[index].type[1];
156     fprintf(stderr, "[DEBUG]:_s%d->_s%d\n",
157         cells[index].type[k].type, cells[index].type[k].cell->index + 1,
158         cells[index].type[k].prev_left_type->type, cells[index].type[k].
            prev_left_type->cell->index + 1
159     );
160     k++;
161     type_index++;
162 }
163 }
164 }
165
166 if (j == WORDAGE - 1) {
167     i = 0;
168     j = ++column;
169 } else {
170     i++;
171     j++;
172 }
173
174 if (k == 0) { // when k == 0, the cell is empty.
175     if (type_index == 0) { // type_index == 0 means the head of type[].
176         type_index++;
177     } else {
178         cells[index].type--; // cells[index].type points last dummy element of
            previous cell (cells[index - 1])
179         // to save memory resource
180     } else {
181         type_index++; // make last element dummy
182     }
183
184     setNewArrayIfNeeded(&type, &type_index, &array_size, &k);
185     index++;
186 }
187
188 fprintf(stderr, "[INFO]:_CYK_parser_stopped.\n");
189 return cells;
190 }

```

ソースコード 18 output.c

```

1  #include "global.h"
2
3
4  void _print_sexps(const TYPE *type, int *word_num, int *type_num, int indent, int
    depth) {
5      if (type->type == DUMMY_CELL.type[0].type) {
6          return;
7      }
8
9      if (indent > 0) {
10         int i, j;
11         printf("\n");
12         for (i = 0; i < depth; i++) {
13             for (j = 0; j < indent; j++) {
14                 printf("_");
15             }
16         }
17     }
18     printf("(%s%d", type->type, (type->cell)->index + 1);
19
20     if ((type->cell)->word != DUMMY_CELL.word) {
21         if (type->prev_left_type == &DUMMY_CELL.type[0]) {
22             printf("_\"%s\"", (type->cell)->word);
23             if (type_num != NULL) {

```

```

24         //fprintf(stderr, "%d, ", type->index);
25         type_num[(*word_num)++] = type->index;
26     }
27 }
28 }
29
30 _print_sexps(type->prev_left_type, word_num, type_num, indent, depth + 1);
31 _print_sexps(type->prev_right_type, word_num, type_num, indent, depth + 1);
32 printf("\n");
33
34 return;
35 }
36
37
38 void print_sexps(const TYPE *type, int *type_num, int indent) {
39     int word_num = 0;
40     _print_sexps(type, &word_num, type_num, indent, 0);
41     printf("\n");
42     return;
43 }
44
45
46 // legacy
47 void print_words(const WORD *words, const size_t buf_size, const int *type_num) {
48     const int WORD_SEPARATION = 2;
49     const int LINES = 3;
50
51     int i;
52
53     // memory allocation
54     char *line[LINES];
55     line[0] = malloc(sizeof(char) * buf_size * LINES); // get heap area
56
57     for (i = 1; i < LINES; i++) {
58         line[i] = line[0] + i * buf_size;
59     }
60
61     // initializing counters and lines
62     int count[LINES];
63     for (i = 0; i < LINES; i++) {
64         line[i][0] = '\0';
65         count[i] = 0;
66     }
67
68     // padding
69     int max = 0;
70     int j, k;
71     int remain = buf_size;
72     for (i = 0; words[i].word != DUMMY.word; i++) {
73         count[0] += snprintf(line[0] + count[0], remain, "%s", words[i].word);
74         count[1] += snprintf(line[1] + count[1], remain, "|");
75         if (type_num == NULL) {
76             count[2] += snprintf(line[2] + count[2], remain, "%s", words[i].type[0]);
77         } else {
78             count[2] += snprintf(line[2] + count[2], remain, "%s", words[i].type[type_num[
79                 i]]);
80         }
81
82         // getting the longest elements
83         for (j = 0; j < LINES; j++) {
84             if (count[j] > buf_size - 1) { // minus 1?: last element must be terminated
85                 with '\0'
86                 // Because snprintf doesn't return the actual number of written characters,
87                 // count[] may exceeds buf_size.
88                 count[j] = buf_size - 1;
89             }
90         }
91     }
92 }

```

```

88     if (max < count[j]) {
89         max = count[j];
90     }
91 }
92
93 // padding with spaces
94 for (j = 0; j < LINES; j++) {
95     while (count[j] < max) {
96         line[j][count[j]++] = ' ';
97     }
98 }
99 // now, count[] == max
100
101 remain = buf_size - max;
102
103 if (remain > WORD_SEPARATION) {
104     for (j = 0; j < LINES; j++) {
105         for (k = 0; k < WORD_SEPARATION; k++) {
106             line[j][count[j]++] = ' ';
107         }
108     }
109     remain -= WORD_SEPARATION;
110 } else {
111     break;
112 }
113 }
114
115 // terminating and printing each line
116 // fprintf(stderr, "[INFO]:\n");
117 for (i = 0; i < LINES; i++) {
118     if (buf_size - count[i] > WORD_SEPARATION) {
119         line[i][count[i] - WORD_SEPARATION] = '\0'; // remove last padding
120     } else {
121         line[i][count[i]] = '\0';
122     }
123     fprintf(stdout, "%s\n", line[i]);
124 }
125
126 free(line[0]);
127 return;
128 }

```

```

File Edit View Search Terminal Help
[DEBUG]: VP14 -> VP2 + PP12
[DEBUG]: S14 -> VP14
[DEBUG]: NP15 -> NOUN1 + NP14
[DEBUG]: VP15 -> VERB1 + NP14
[DEBUG]: S15 -> NP1 + VP14
[DEBUG]: NP15 -> NP6 + PP12
[DEBUG]: S15 -> NP6 + VP12
[DEBUG]: VP15 -> VP6 + PP12
[DEBUG]: S15 -> VP15
[DEBUG]: S15 -> VP15
[INFO]: CYK parser stopped.
(S15(NP1(NOUN1 "time"))(VP14(VP2(VERB2 "flies"))(PP12(PREP3 "like")(NP9(
DET4 "an")(NOUN5 "arrow")))))

time flies like an arrow
|   |   |   |   |
NOUN VERB PREP DET NOUN

(S15(NP6(NOUN1 "time")(NP2(NOUN2 "flies")))(VP12(VERB3 "like")(NP9(
DET4 "an")(NOUN5 "arrow"))))

time flies like an arrow
|   |   |   |   |
NOUN NOUN VERB DET NOUN

(S15(VP15(VERB1 "time")(NP14(NP2(NOUN2 "flies"))(PP12(PREP3 "like")(NP9(
DET4 "an")(NOUN5 "arrow")))))

time flies like an arrow
|   |   |   |   |
VERB NOUN PREP DET NOUN

(S15(VP15(VP6(VERB1 "time")(NP2(NOUN2 "flies")))(PP12(PREP3 "like")(NP9(
DET4 "an")(NOUN5 "arrow")))))

time flies like an arrow
|   |   |   |   |
VERB NOUN PREP DET NOUN

[INFO]: Sentence OK (sentence is acceptable.)
[INFO]: Parsing finished.
[INFO]: exiting... Goodbye.
:~/jikken3/c1/cyk$

```

図 13 文脈自由文法のプログラムの動作風景

B.1 インデント機能の使用

S 式出力用の関数 `print_sexps` はインデント機能を有している。main 関数内で

```
print_sexps(&(cells[CELL_LENGTH - 1].type[i]), type_num, 0);
```

と呼び出している部分の最後の引数を例えば

```
print_sexps(&(cells[CELL_LENGTH - 1].type[i]), type_num, 4);
```

とすれば 4 文字でのインデントが有効になり、簡易的に木構造を確認できる (図 14)。

```
(S36
  (NP9
    (DET1 "the")
    (NOUN2 "child"))
  (VP33
    (VPa11
      (VERB3 "runs")
      (ADV4 "quickly"))
    (NP26
      (PREP5 "to")
      (NP21
        (DET6 "the")
        (NP15
          (ADJ7 "large")
          (NOUN8 "house"))))))))
```

図 14 インデント機能を有効にした場合の出力