

МАТЕМАТИЧЕСКОЕ ОБЕСПЕЧЕНИЕ ПОДСИСТЕМЫ СИНТЕЗА ТЕКСТА В САПР

И.Л. Коробова

*Кафедра «Системы автоматизированного проектирования»,
ГОУ ВПО «ТГТУ»; ira.sapr.tstu@mail.ru*

Представлена членом редколлегии профессором В.И. Коноваловым

Ключевые слова и фразы: грамматика языка; диалоговая система; компьютерная лингвистика; синтаксис; семантика; фрейм; шаблонизация.

Аннотация: Рассмотрена методика разработки систем синтеза текста на естественном языке. Приведено описание математического обеспечения, в том числе методов разработки шаблонов, генерации текста, синтаксического анализа.

При разработке подсистемы синтеза текста [1, 2, 5] мы ставили задачу формализовать процесс шаблонизации [3], сделать его понятным и применимым в различных процедурах САПР. В связи с этим, наиболее трудоемким оказался процесс описания математического обеспечения.

Математическое обеспечение подсистемы автоматизированного синтеза текста на основе технологии шаблонизации состоит из множества методов, образующих три группы: методы разработки шаблонов, методы генерации текста, методы обеспечения синтаксического анализа и корректировки.

I. Методы разработки шаблонов включают:

- разработку шаблонов последовательности и связей функций, которые отвечают за правильное следование функций и правильные связи между ними;
- разработку шаблонов функций, которые предусматривают возможность добавления или удаления элемента функции, изменения приоритета элемента, его редактирование;
- разработку шаблонов элементов функций, которые предусматривают возможность добавления или удаления подэлемента, изменения приоритета подэлемента (используется сдвиг набора подэлементов), изменение значения (если это глагол-действие), изменение падежа (если это имя существительное или прилагательное).

II. Методы генерации текста включают:

- генерацию текста-шаблона, в котором отсутствуют индивидуальные данные. Производится перебор всех функций набора, начиная с первой. Для каждой функции производится перебор всех элементов в зависимости от приоритета, начиная с наименьшего. Для каждого элемента производится перебор всех подэлементов в зависимости от приоритета, начиная с наименьшего. Для каждого подэлемента производится вывод в текст-шаблон значения подэлемента;
- генерацию текста, в котором присутствуют индивидуальные данные. Подобен предыдущему методу, но в значения подэлементов, характеризующихся

индивидуальными данными, подставляются введенные ранее индивидуальные данные.

III. Методы обеспечения синтаксического анализа и корректировки.

Для их реализации в COM-сервер MS WordDocument передается текст, и в диалоговом режиме производится проверка правописания, затем производится ручная корректировка текста – в выходном тексте заменяются выбранные пользователем фрагменты.

Рассмотрим подробно алгоритм блока формирования шаблона последовательности функций для выбранного сюжета (рис. 1).

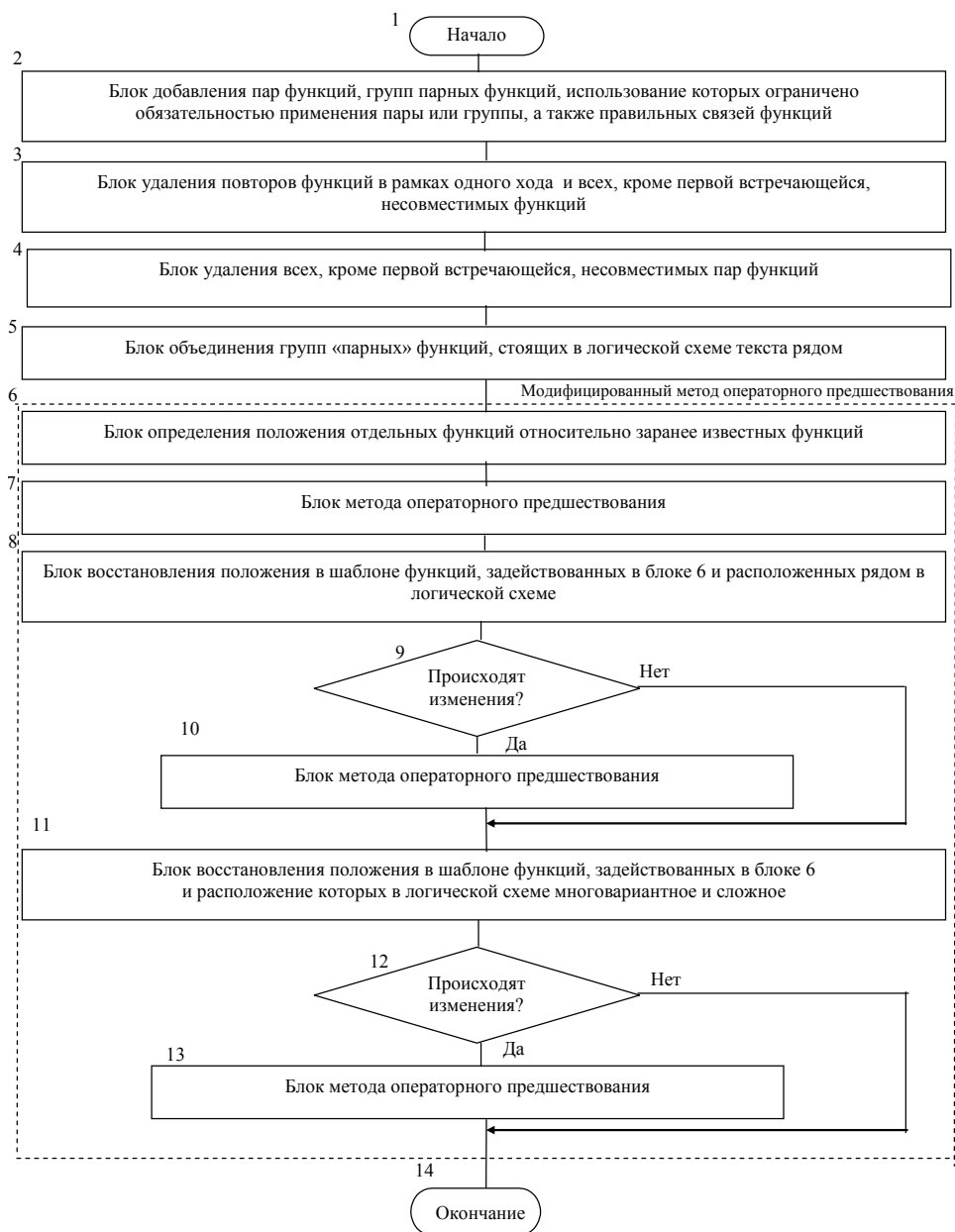


Рис. 1. Блок формирования шаблона последовательности функций для выбранного сюжета

Блок определения наличия в массиве функций, необходимых для существования текста, основан на правиле:

$$\begin{aligned}
 &\text{ЕСЛИ } f_j \triangleleft f_i \text{ ТО count}[j] := \text{count}[j] + 1; \\
 &\text{ЕСЛИ count}[j] = n \\
 &\quad \text{ТО количество несовпадений} := \text{количество несовпадений} + 1; \\
 &\text{ЕСЛИ количество несовпадений} > 0, \text{ ТО результат} := \text{ложь} \\
 &\text{ИНАЧЕ результат} := \text{истина};
 \end{aligned} \tag{1}$$

где n – количество функций в наборе; $\text{count}[j]$ – счетчик несовпадений f_j функции, необходимой для существования текста ($j \in [1, m]$), с f_i функциями набора ($i \in [1, n]$), m – количество функций, необходимых для существования текста; количество несовпадений – количество функций из m , которые не встретились в тексте; результат – наличие в наборе m функций, необходимых для существования текста.

Блок определения сюжета текста основан на правиле:

$$\begin{aligned}
 &\text{ЕСЛИ } f_j = f_i \text{ ТО stroka}[j] := 1; \\
 &\text{Сюжет} := \sum_{j=1}^m (\text{stroka}[j] \cdot 2^j) + 1
 \end{aligned} \tag{2}$$

где $\text{stroka}[j]$ – строка (изначально равна '000...00') со строковым номером сюжета в двоичном виде ($j \in [1, m]$); сюжет – номер сюжета в десятичном виде (сюжет $\in [1; 2^m]$).

Блок добавления пар функций, групп парных функций, правильных связей функций основан на правиле:

$$\begin{aligned}
 &\text{ЕСЛИ } f_j^l = f_i \text{ ТО } [\text{ЕСЛИ НЕ } f_r \triangleleft f_k^l \text{ ТО count}[k] := \text{count}[k] + 1] \\
 &\text{ЕСЛИ count}[k] = n \text{ ТО } [\text{массив}[n] := f_k^l; n := n + 1]
 \end{aligned} \tag{3}$$

где $\text{массив}[n]$ – ячейка набора функций.

Добавление правильных связей функций выполняется блоком добавления правильных связей функций. Используется утверждение, что не для всех функций набора необходимы взаимные связи.

Блок удаления повторов функций в рамках одного хода и всех, кроме первой встречающейся, несовместимых функций основан на правиле:

$$\text{ЕСЛИ } f_j = f_i \text{ ТО } [\text{ЕСЛИ } f_r = f_j \text{ ТО } (\text{массив}[p] := \text{массив}[p + 1]; n := n - 1)] \tag{4}$$

Блок удаления всех, кроме первой встречающейся, несовместимых функций основан на правиле:

$$\text{ЕСЛИ } f_j^l = f_i \text{ ТО } [\text{ЕСЛИ } f_r = f_k^l \text{ ТО } (\text{массив}[p] := \text{массив}[p + 1]; n := n - 1)] \tag{5}$$

где l – номер группы несовместимых функций ($l \in [1, m]$); s – количество функций в l -й группе ($j \in [1, s]$, $k \in [1, s]$, $j \neq k$); $\text{массив}[p]$ – ячейка набора функций.

Блок объединения групп «парных» функций, стоящих в логической схеме текста рядом, основан на правиле:

$$\begin{aligned}
 &\text{ЕСЛИ } f_i^l = f_j \text{ ТО } [\text{ЕСЛИ } ((f_r = f_k^l) \text{ И } (k \triangleleft r)) \text{ ТО} \\
 &\quad \text{ЕСЛИ } r > i \text{ ТО } (\text{переменная} := \text{массив}[r]; \\
 &\quad \text{массив}[p] := \text{массив}[p - 1]; \\
 &\quad \text{массив}[i + k + l] := \text{переменная}) \\
 &\quad \text{ИНАЧЕ } (\text{переменная} := \text{массив}[r]; \\
 &\quad \text{массив}[p] := \text{массив}[p + 1]; \\
 &\quad \text{массив}[i + k - 1] := \text{переменная})
 \end{aligned} \tag{6}$$

Блок определения положения отдельных функций относительно заранее известных функций основан на правиле:

$$\boxed{\text{ЕСЛИ } f_i = f_j \text{ ТО } [\text{ЕСЛИ } i < m \text{ ТО символ } := ' < ' \text{ ИНАЧЕ символ } := ' > ']} \quad (7)$$

где символ – ячейка матрицы, схожей с матрицей операторного предшествования (изначально равен ' ').

Блок коррекции порядка следования функций основан на методе операторного предшествования (между функциями f_i и f_j существует отношение '<', если $i < j$; отношение '>', если $i > j$; 'n' (несовместимость), если f_i и f_j не могут располагаться рядом в наборе функций).

Модифицированный метод операторного предшествования основан на правиле:

$$\boxed{\begin{array}{l} \text{ЕСЛИ } ((i > j \text{ И матрица}[i, j] = ' < ') \text{ ИЛИ } (i < j \text{ И матрица}[i, j] = ' > ')) \text{ ТО} \\ [\text{ЕСЛИ } i > j \text{ ТО } (k1 := j; k2 := i) \\ \text{ИНАЧЕ } (k1 := i; k2 := j) \\ \text{переменная } := \text{массив}[k1]; \\ \text{массив}[p] := \text{массив}[p - 1]; \\ \text{массив}[k2] := \text{переменная}] \\ \text{ИНАЧЕ переход на следующий шаг} \end{array}} \quad (8)$$

Блок восстановления положения в шаблоне функций, задействованных в блоке определения положения отдельных функций относительно заранее известных функций, расположенных рядом в логической схеме, основан на правиле:

$$\boxed{\begin{array}{l} \text{ЕСЛИ } f_i = f_j \text{ ТО } [\text{ЕСЛИ } ((i > m \text{ И символ } = ' < ') \\ \text{ИЛИ } (i < m \text{ И символ } = ' > ')) \text{ ТО} \\ (a := \text{массив}[i]; \\ \text{массив}[i] := \text{массив}[m]; \\ \text{массив}[m] := a)] \end{array}} \quad (9)$$

Математические методы и алгоритмы синтеза текста решают задачу выбора выходного параметра системы высказываний (1) – (9) на основе правила Modus Ponens [4].

В настоящее время подсистема синтеза текста работает для формирования системы объяснений в процессе принятия решения в условиях нечеткой экспертной информации [4], а также, для создания технического задания на разработку САПР и ее составных частей [5].

Список литературы

1. Коробова, И.Л. Автоматизированная система синтеза текста на основе технологии шаблонизации [Электронный ресурс] / И.Л. Коробова // Мат. межрегионал. науч.-практ. конф. «Информатизация системы образования Тамбовского региона». – Режим доступа : <http://club-edu.tambov.ru/main/news/index.php?r=konf1&f=t12>. – Загл. с экрана.
2. Коробова, И.Л. Информационное обеспечение подсистемы синтеза текста при автоматизированном проектировании технологических объектов / И.Л. Коробова, И.А. Дьяков // Теплофизика в энергосбережении и управлении качеством : мат. Шестой междунар. теплофиз. шк. / Тамб. гос. техн. ун-т. – Тамбов, 2007. – Ч. 2. – С. 27–31.
3. Информатика : энциклопед. слов. для начинающих / сост. Д.А. Поспелов. – М. : Педагогика-Пресс, 1994 – 352 с.

4. Коробова, И.Л. Анализ знаний в экспертной системе нечеткого принятия решений / И.Л. Коробова // Вест. Тамб. гос. техн. ун-та. – 2005. – Т. 11, № 4. – С. 873–881.

5. Коробова, И.Л. Подсистема синтеза текста в САПР / И.Л. Коробова, Н.В. Майстренко // Вест. Тамб. гос. техн. ун-та. – 2009. – Т. 15, № 1. – С. 49–55.

Mathematical Software for Subsystem of CAD Text Synthesis

I.L. Korobova

*Department "Computer Aided Design Systems", TSTU;
ira.sapr.tstu@mail.ru*

Key words and phrases: dialogue system; computer linguistics; grammar of the language; frame; syntax; semantics; standardization.

Abstract: The paper studies the method of developing systems of text synthesis in the natural language. The description of mathematical software, including the techniques for developing templates, text generation and syntactic analysis is given.

Matematische Versorgung des Subsystemes der Synthese des Textes in SAPR

Zusammenfassung: Es wird die Methodik der Erarbeitung der Systeme der Synthese des Textes auf der Natursprache betrachtet. Es ist die Beschreibung der matematischen Versorgung, darunter der Methoden der Erarbeitung der Schablonen, der Textgeneration, der syntaktischen Analyse, angeführt.

Le logiciel du sous-système de la synthèse du texte dans CAO

Résumé: Est examinée la méthode de l'élaboration des systèmes de la synthèse du texte en langue naturelle. Est donnée la description du logiciel y compris des méthodes de l'élaborations des clichés, de la génération des textes, de l'analyse syntaxique.

Автор: *Коробова Ирина Львовна* – кандидат технических наук, доцент кафедры «Системы автоматизированного проектирования», ГОУ ВПО «ТГТУ».

Рецензент: *Подольский Владимир Ефимович* – доктор технических наук, профессор, заведующий кафедрой «Системы автоматизированного проектирования», проректор по информатизации, ГОУ ВПО «ТГТУ».
