

УДК 519.876.5

С.А. Панов

## Формальный язык описания структуры документов и его интерпретация в формат метода компонентных цепей

Статья посвящена алфавиту и правилам формального языка  $L(G)$ , используемого для описания структуры документов и интерпретации символов этого языка в компьютерные модели, создаваемые с помощью метода компонентных цепей. Разработанные компьютерные модели будут использоваться в интерактивном генераторе отчетных форм, представляющем собой систему автоматизированного документирования.

**Ключевые слова:** автоматизация, формальный язык, документирование, автоматизация документирования, метод компонентных цепей, компьютерное моделирование.

С целью повышения скорости создания документов и точности содержащихся в них данных активно разрабатываются и повсеместно используются системы автоматизированного документирования. К сожалению, существующие к настоящему моменту программные средства автоматизированного документирования не позволяют описывать формальную структуру документов и соответственно, направлены на создание документов только определенного формата и только для определенной области применения (промышленность, экономика, образование и т.д.). Разработка универсального формального языка позволит создать систему автоматизированного формирования документов любого состава и назначения.

В текущий момент в ТУСУРе разрабатывается собственная универсальная система автоматизированного документирования, называемая «Интерактивный генератор отчетных форм» (ИГОФ). ИГОФ предоставляет широкие возможности интерактивной работы над документом [1]. В данной работе под документом понимается набор структурных элементов (фрагментов), использующийся для хранения и передачи информации. Фрагменты могут быть двух видов: элементарные фрагменты представляют простейшие неделимые элементы, а составные фрагменты содержат в себе элементарные и другие фрагменты.

**Формальный язык описания структуры документов.** Модель документа в ИГОФ создается с помощью метода компонентных цепей [2], который идеально подходит для моделирования информационных процессов и систем. Сущность данного метода заключается в том, что объект (в данном случае это документ) представляет собой цепь из специальных взаимосвязанных между собой компонентов. Каждый компонент в такой цепи имеет входные и выходные узлы, с помощью которых выполняется соединение компонентов между собой. На входной узел поступают данные, которые обрабатываются и отправляются на другие компоненты с помощью выходного узла.

Задачу построения компонентов ИГОФ существенно упрощает использование формального языка описания структуры документов, каждое слово в котором отражает определенный фрагмент документа.

Формальный язык  $L(G)$  задается формальной грамматикой  $G$ :

$$G = (N, T, S, P),$$

где  $N$  – набор (алфавит) нетерминальных символов (нетерминалов):

$$N = \{\text{Документ, Текст, Таблица, Рисунок, МатВыражение, ИсЛит, Число, Слово, Цифра, ЗнакПрепинания, ЗнакОкончания, Буква, АрифмOp}\},$$

•  $T$  – набор (алфавит) терминальных символов (терминалов):

$$T = \{0,1,2,3,4,5,6,7,8,9\} \cup \{\dots, \text{”}, \text{”}, \text{”}, \text{”}\} \cup \{!,?,?!,!!!,\dots\} \cup \{A,B,B,G,D,E,\dots,Я,a,b,в,z,d,e,\dots,я,A,B,C,D,E,\dots,Z,a,b,c,d,e,\dots,z\} \cup \{+,-,*,/\} \cup \{\square\},$$

где  $\square$  – остальные символы ASCII или Unicode;  $S$  – стартовый (начальный) нетерминал (источник):

$$S = \text{Документ},$$

$P$  – конечный набор (множество) правил формальной грамматики.

Таким образом, нетерминальные символы соответствуют составным, а терминальные – элементарным фрагментам документа.

Основным обязательным элементом формальной грамматики является набор специальных правил, по которым каждый нетерминальный символ может быть представлен в виде совокупности терминальных символов. Такие правила имеют вид: «левая часть»  $\rightarrow$  «правая часть», где «левая часть» – это нетерминал, а «правая часть» – терминал, нетерминал или их совокупность. Так как в разрабатываемом формальном языке  $L(G)$  имеется 13 нетерминалов, соответственно имеется 13 правил формальной грамматики:

1. Символ **<Документ>** является начальным и служит для описания всего документа в целом:  
 $\langle \text{Документ} \rangle \rightarrow \langle \text{Текст} \rangle \mid \langle \text{Текст} \rangle \langle \text{Таблица} \rangle \mid \langle \text{Текст} \rangle \langle \text{Рисунок} \rangle \mid \langle \text{Текст} \rangle \langle \text{Документ} \rangle$

2. Символ **<Текст>** представляет собой текст, который включает в себя набор всех символов из таблицы ASCII (или Unicode) и поэтому синтаксическому анализу не подлежит. В самом простом случае – тексте, где нет цифровых обозначений, рисунков, формул и ссылок на литературные источники, символ **<Текст>** может быть представлен в виде

$$\langle \text{Текст} \rangle \rightarrow \langle \text{Предложение} \rangle \mid \langle \text{Предложение} \rangle \langle \text{Текст} \rangle$$

$$\langle \text{Предложение} \rangle \rightarrow \langle \text{Слово} \rangle \langle \text{КонецПредложения} \rangle$$

$$\langle \text{КонецПредложения} \rangle \rightarrow \langle \text{Слово} \rangle \langle \text{ЗнакОкончания} \rangle \mid \langle \text{Слово} \rangle \langle \text{КонецПредложения} \rangle \mid \langle \text{ЗнакПрепинания} \rangle \langle \text{КонецПредложения} \rangle \mid \langle \text{ЗнакОкончания} \rangle$$

$$\langle \text{Слово} \rangle \rightarrow \langle \text{Буква} \rangle \mid \langle \text{Слово} \rangle \langle \text{Буква} \rangle$$

3. Символ **<Таблица>** служит для обозначения в документе таблиц различной формы:

$$\langle \text{Таблица} \rangle \rightarrow \langle \text{Строка} \rangle \mid \langle \text{Столбец} \rangle \mid \langle \text{Строка} \rangle \langle \text{Столбец} \rangle \mid \langle \text{Таблица} \rangle \langle \text{Строка} \rangle \mid \langle \text{Таблица} \rangle \langle \text{Столбец} \rangle$$

$$\langle \text{Строка} \rangle \rightarrow \langle \text{Ячейка} \rangle \mid \langle \text{Строка} \rangle \langle \text{Ячейка} \rangle$$

$$\langle \text{Столбец} \rangle \rightarrow \langle \text{Ячейка} \rangle \mid \langle \text{Столбец} \rangle \langle \text{Ячейка} \rangle$$

$$\langle \text{Ячейка} \rangle \rightarrow \langle \text{Текст} \rangle \mid \langle \text{МатВыражение} \rangle \mid \langle \text{Формула} \rangle \mid \langle \text{Рисунок} \rangle \mid \langle \text{Число} \rangle$$

4. Символ **<Рисунок>** служит для обозначения рисунков с подписью:

$$\langle \text{Рисунок} \rangle \rightarrow \langle \text{Изображение} \rangle \mid \langle \text{Изображение} \rangle \langle \text{Описание} \rangle$$

Символ **<Изображение>** является терминальным и представляет собой изображение в формате JPEG, PNG, GIF, TIFF.

$$\langle \text{Описание} \rangle \rightarrow \langle \text{НачОп} \rangle \langle \text{Слово} \rangle \mid \langle \text{НачОп} \rangle \langle \text{Слово} \rangle \langle \text{КонОп} \rangle$$

$$\langle \text{НачОп} \rangle \rightarrow \text{Рис.} \langle \text{Число} \rangle$$

$$\langle \text{КонОп} \rangle \rightarrow \langle \text{Слово} \rangle \mid \langle \text{Слово} \rangle \langle \text{КонОп} \rangle \mid \langle \text{ЗнакПрепинания} \rangle \langle \text{КонОп} \rangle \mid \langle \text{ЗнакПрепинания} \rangle \langle \text{Слово} \rangle$$

5. Символ **<МатВыражение>** служит для обозначения математических выражений и формул [3].

6. Символ **<ИстЛит>** служит для обозначения в документе номеров источников использованной литературы:

$$\langle \text{ИстЛит} \rangle \rightarrow \langle \text{НачалоИстЛит} \rangle \langle \text{НомерИсточника} \rangle \langle \text{КонецИстЛит} \rangle$$

$$\langle \text{НачалоИстЛит} \rangle \rightarrow [$$

$$\langle \text{КонецИстЛит} \rangle \rightarrow ]$$

$$\langle \text{НомерИсточника} \rangle \rightarrow 1 \mid 2 \mid 3 \mid \dots \mid N,$$

где  $N$  – целое положительное число.

Ниже представлены символы для обозначения чисел, слов, цифр, знаков препинания, знаков окончания предложения, букв и арифметических операций:

7. **<Число>**:

$$\langle \text{Число} \rangle \rightarrow [-] \langle \text{Цифра} \rangle \mid [-] \langle \text{Число} \rangle \langle \text{Цифра} \rangle$$

8. **<Слово>**:

$$\langle \text{Слово} \rangle \rightarrow \langle \text{Буква} \rangle \mid \langle \text{Слово} \rangle \langle \text{Буква} \rangle$$

9. **<Цифра>**:

$$\langle \text{Цифра} \rangle \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

10. **<ЗнакПрепинания>**:

$$\langle \text{ЗнакПрепинания} \rangle \rightarrow \dots \mid , \mid : \mid -$$

11. **<ЗнакОкончания>**:

$$\langle \text{ЗнакОкончания} \rangle \rightarrow . \mid ! \mid ? \mid ?! \mid !!! \mid \dots$$

12. **<Буква>**:

$$\langle \text{Буква} \rangle \rightarrow A \mid B \mid V \mid \Gamma \mid D \mid E \mid \dots \mid Я \mid a \mid \bar{b} \mid \bar{v} \mid \bar{z} \mid \bar{d} \mid e \mid \dots \mid я \mid A \mid B \mid C \mid D \mid E \mid \dots \mid Z \mid a \mid b \mid c \mid d \mid e \mid \dots \mid z$$

13. **<АрифмОп>**:

$$\langle \text{АрифмОп} \rangle \rightarrow + \mid - \mid * \mid /$$

Разработанная грамматика содержит всю полноту описания фрагментов документов различного назначения (технических, экономических, учебных, научных и т.д.).

С целью разработки компьютерных имитационных моделей фрагментов документа и последующей реализации этих моделей в виде компонентов ИГОФ необходимо выполнить интерпретацию между фрагментами документа, основными нетерминальными символами в формальном языке  $L(G)$  и компонентами ИГОФ.

Интерпретация символов формального языка  $L(G)$  в компьютерные имитационные модели. Интерпретация осуществляется между:

- 1) фрагментами документа;
- 2) основными нетерминальными символами формального языка  $L(G)$ ;
- 3) компьютерными имитационными моделями фрагментов документа, разработанными с помощью метода компонентных цепей и реализованными в виде компонентов ИГОФ.

На текущий момент спроектировано шесть основных компонентов ИГОФ, соответствующих нетерминальным символам формального языка  $L(G)$ : документ, источник текста, формирователь таблицы / расширенная таблица, схема / диаграмма / график / рисунок, источник формул и числовой интерпретатор (таблица).

Компонентный базис ИГОФ

№ п/п	Фрагмент документа	Нетерминальный символ в языке формальном языке $L(G)$	Название компонента ИГОФ	Графическое представление в ИГОФ	Краткое описание
1	Документ	<Документ>	«Документ»		Служит для сбора информации с других компонентов и формирования итогового документа
2	Текст	<Текст>	«Источник текста»		Служит для вставки в документ текста
3	Таблица	<Таблица>	«Формирователь таблицы» / «Расширенная таблица»		Служит для вставки в документ таблиц [4]
4	Рисунок	<Рисунок>	«Схема», «Диаграмма», «График», «Рисунок»		Служит для вставки в документ рисунков
5	Математическое выражение	<МатВыражение>	«Источник формул»		Служит для вставки в документ формул
6	Число	<Число>	«Числовой интерпретатор»		Служит для задания числовых констант или переменных

В процессе разработки ИГОФ могут быть получены новые компоненты для работы с фрагментами документов, поэтому таблица будет постепенно пополняться все новыми и новыми элементами.

**Заключение.** Предложенный формальный язык  $L(G)$  разработан с целью упрощения процесса создания компьютерных моделей фрагментов документов и является универсальным языком описания структуры документов различного назначения. Формальный язык  $L(G)$  может быть встроен в проектируемые и уже используемые системы автоматизированного документирования, что позволит убрать ограничение таких систем на область их применения. Полнота языка подтверждается наличием всех обязательных элементов: алфавита терминальных и нетерминальных символов, грамматики и правил построения слов.

Исследование выполнено при финансовой поддержке гранта РФФИ 11-07-00384 «Метод многоуровневого моделирования алгоритмов управления технологическими процессами в сложных системах».

#### *Литература*

1. Ганджа Т.В. Задачи и архитектура подсистемы документирования исследований в среде многоуровневого моделирования МАРС / Т.В. Ганджа, С.А. Панов // Доклады Том. гос. ун-та систем управления и радиоэлектроники. – 2011. – № 2(24), ч. 2. – С. 334–338.
2. Дмитриев В.М. Автоматизация моделирования промышленных роботов // В.М. Дмитриев, Л.А. Арайс, А.В. Шутенков. – М.: Машиностроение, 1995. – 304 с.
3. Дмитриев В.М. Алгоритм формирования и вычисления математических выражений методом компонентных цепей / В.М. Дмитриев, Т.В. Ганджа // Математические машины и системы. – 2010. – № 3. – С. 9–21.
4. Ганджа Т.В. Генератор табличных форм как компонент системы автоматизированного документирования / Т.В. Ганджа, С.А. Панов // Наука. Технологии. Инновации: матер. Всерос. науч. конф. молодых ученых, Новосибирск, 29 ноября – 2 декабря 2012 г.: в 7 ч. – Новосибирск: Изд-во НГТУ, 2012. – Ч. 3. – С. 88–91.

---

#### **Панов Сергей Аркадьевич**

Аспирант каф. моделирования и системного анализа ТУСУРа

Тел.: 8 (382-2) 41-39-15

Эл. почта: spytech@ieee.org

Panov S.A.

#### **A formal language for describing the structure of technical documents and interpretation of it to the format method of component circuits**

The paper describes a formal language that contains terminal and non-terminal symbols that reflect the structure of the technical documentation and interpretation of these symbols in the corresponding computer model presented in the format of a method of component circuits. The development of these models will create a system of automatic formation of documents which is an important task in the field of automation and process control industries.

**Keywords:** automation, formal language, documentation, document automation, technical paper, the method of component circuits, computer simulation.

---