

## КОНСТРУКТИВНО-ПРОДУКЦИОННАЯ МОДЕЛЬ ГРАФОВОГО ПРЕДСТАВЛЕНИЯ ТЕКСТА

*В.І. Шинкаренко, Е.С. Куропятник*

У статті розглянута графова модель тексту, що дозволяє прискорити обробку інформації. Данна модель дозволяє виявляти однакові фрагменти в документах зі зміною порядку слідування речень та інших частин. Використання конструктивно-продукційних структур для формалізації даної моделі є перспективним підходом для подальшої автоматизації процесу роботи з моделлю і відповідно з текстом.

Ключові слова: графова модель тексту, стискання графа, конструктивно-продукційна структура, співставлення текстів.

В статье рассмотрена графовая модель текста, позволяющая ускорить обработку информации. Данная модель позволяет выявлять одинаковые фрагменты в документах с изменением порядка следования предложений и других частей. Использование конструктивно-продукционных структур для формализации данной модели является перспективным подходом для дальнейшей автоматизации процесса работы с моделью и соответственно с текстом.

Ключевые слова: графовая модель текста, сжатие графа, конструктивно-продукционная структура, сопоставление текстов.

The article describes the graph model of the text, allowing speeds up processing. This model allows us to identify the same fragments in the documents with the change in the order of sentences and other parts. Using constructive-synthesizing structure to formalize this model is a promising approach to further automate the process of working with the model and the text accordingly.

Key words: graph model of the text, graph compression, constructive-synthesizing structure, text comparison.

### Введение

Исследования и моделирование в различных отраслях науки и общественного производства, а также образовательные процессы неотрывно связаны с обработкой уже имеющихся научных изданий, разработкой и документированием новых методик, разработкой текстов компьютерных программ и тому подобное. Скорость и интенсивность развития сфер науки, образования и производства, средств их информационной поддержки порождает проблему обработки текстов, которая может быть разделена на несколько составляющих. На сегодня наиболее актуальными являются поиск, сравнение, проверка правописания, выделение знаний. Разработка и модернизация средств автоматизации для решения данных вопросов позволяет улучшить временные показатели решения данных задач.

Существует множество методик для предварительной обработки текстов с целью упрощения их дальнейшей обработки. Среди них приведение символов к нижнему регистру и единой кодировке, удаление стопслов и замена слов их грамматическими основами [1, 2], удаление или игнорирование пунктуационных знаков [3], игнорирование пробелов, пустых строк, комментариев в текстах программ [4].

Особое внимание уделяется обработке текстов с целью выявления плагиата. Для этого разработано множество алгоритмов, моделей и программ на их основе [4]. Среди них графовое представление семантики текста [3, 5], токенизации [6], морфологический анализ и синтаксический разбор [4], латентный семантический анализ [6, 7] и др.

Представление текста в виде графовой модели [8] позволяет ускорить обработку текстовой информации. Использование КПС для формализации данной модели является перспективным подходом для дальнейшей автоматизации процесса работы с моделью и соответственно текстом.

Формализация средствами конструктивно-продукционных структур (КПС) позволяет описать не только структуру объектов, но и их свойства, определить допустимые операции над ними, алгоритмические элементы построения конструкций на их основе [9]. Ранее была рассмотрена обобщенная КПС [10]. Данный подход является универсальным для решения многих задач и применяется в разноплановых работах, посвященных адаптации алгоритмов сжатия [11], формированию структур данных на логическом уровне [12], обобщенных грамматик [13].

Особенностью КПС является возможность формирования конструкций различной природы на основе операций подстановки, связывания, вывода и других вспомогательных операций, а также использования атрибутов [10].

### Обобщенная конструктивно-продукционная структура

Обобщённая КПС определяется носителем  $M$ , сигнатурой  $\Sigma$  и аксиоматикой  $\Lambda$   $C = \langle M, \Sigma, \Lambda \rangle$ .

Носитель включает в себя конструктивные элементы с набором атрибутов. Сигнатура состоит из множества имен операций и отношений. Аксиоматика включает множества определений, аксиом, правил, свойств, правил и пр. Носитель и сигнатура определяются аксиоматикой.

Преобразования КПС позволяют применять их для различных предметных областей. Специализация КПС подразумевает определение атрибутов носителя, семантическую составляющую, а также определения конечного конкретного набора операций, их атрибутов. В ходе интерпретации каждой операции сигнатуры становится в соответствие алгоритм ее выполнения, являющийся элементом носителя некоторой алгоритмической

структурой [14]. Реализации КПС заключается в построении конструкций из элементов носителя КПС путем выполнения алгоритмов, связанных с операциями сигнатуры.

Подход на основе КПС может быть использован для формализации понятий текста, его составляющих. Особенности определения понятий текста, предложения, слова определяют носитель КПС и суть операций по ее преобразованию. Так, например, если слово рассматривать с точки зрения естественного языка, то оно не может включать символы-цифры, а если с точки зрения языков программирования, слово – лексема, которая может содержать цифры и иные символы (имена переменных, функций, классов и т.п.). Таким образом, определение подобных понятий может как расширять, так и сужать область определения элементов носителя, конструкций, операций над ними.

## Формализованная спецификация текста средствами КПС

Формализованная спецификация текста и его составляющих дает возможность автоматизировать разработку и обновление программ для решения поставленных задач.

Введем некоторые ограничения и дополнения аксиоматики и обозначений [10] в соответствие с рассматриваемой предметной областью.

С данной целью определим КПС и выполним ее специализацию:

$$C = \langle M, \Sigma, \Lambda \rangle \xrightarrow{S} C_T = \langle M_T, \Sigma_T, \Lambda_T \rangle, \quad (1)$$

где  $M_T$  – носитель, включающий все символы электронного представления текста и конструкции, построенные на них,  $\Sigma_T$  – операции и отношения на элементах  $M_T$ ,  $\Lambda_T$  – аксиоматика, определяющая  $M_T$  и  $\Sigma_T$ .

**Частичная аксиоматика носителя.** Носитель включает множества терминалов и нетерминалов.  $M_T \supset T_T \cup N_T$ ,  $T_T = T_K \cup T_L \cup T_D \cup T_N \cup T_{NP}$ ,  $T_K, T_L$  – множество кириллических и латинских символов соответственно, которые входят в текст;  $T_D$  – множества символов-разделителей для отображения знаков пунктуации и пробела;  $T_N$  – множество цифр;  $T_{NP}$  – множество непечатных символов, не включая пробел;  $N_T$  – множество нетерминалов (вспомогательных символов).

Элементы носителя  $\bar{w}_i m_i$  имеют набор атрибутов  $\bar{w}_i = \langle type, code \rangle$ . Принадлежность конкретного атрибута элементу носителя будем обозначать  $w_i \leftarrow m$ . Тип ( $type \leftarrow m$ ) принимает значение терминал/нетерминал ( $ter / nter$ ). Атрибут  $code$  указывает на машинное представление символа-терминала. В качестве идентифицирующего атрибута выступает код.

**Частичная аксиоматика отношений и операций.** Сигнтура  $\Sigma_T$  состоит из множества операций  $\Sigma_T = \langle \Xi, \Theta, \Phi, \{\rightarrow\} \rangle \cup \Psi$ ,  $\Xi = \{\cdot\}$ , где  $\Xi$  – множество операций связывания,  $\Theta$  – множество операций вывода,  $\Phi = \{=, :=\}$  – множество операций над атрибутами;  $\Psi$  – множество правил продукций вида  $\psi_i : s_i, g_i \rightarrow, i$  – номер правила,  $s$  – последовательность операций подстановки,  $g$  – последовательность операций над атрибутами, « $\rightarrow$ » – отношение подстановки.  $\Theta = \{\Rightarrow, \mid\Rightarrow, \parallel\Rightarrow\}$  – операции подстановки, частичного и полного вывода.

Операция конкатенации  $\cdot(m_i, m_j)$  связывает два элемента носителя. Результатом выполнения операции является форма  $_{w_i} l$ . Форма может быть построена таким образом:

- $_{w_i} l = _{w_0} \cdot (_{w_j} m_j, _{w_k} m_k)$  для  $\forall _{w_i} m_i \in M_T$ ;
- $_{w_i} l = _{w_j} m_j$ , если  $l = _{w_0} \cdot (\varepsilon, _{w_j} m_j) = _{w_0} \cdot (_{w_j} m_j, \varepsilon)$ , где  $\varepsilon$  – пустой элемент;
- $_{w_i} l = _{w_0} \cdot (_{w_1} l_1, _{w_2} l_2)$  для  $\forall _{w_i} l_i = _{w_j} m_j \mid (_{w_j} m_j, _{w_k} m_k) \mid \cdot (_{w_j} l_j, _{w_k} m_k)$ ,  $\forall _{w_i} m_i \in M_T$ .

Множество значений атрибутов формы определяется совокупностью значений атрибутов ее составляющих.

Сентенциальная форма – форма, полученная в результате вывода из аксиомы (начального нетерминального символа, принадлежащего носителю) согласно правилам вывода конкретизированной КПС.

Формы, в которых отсутствуют нетерминальные элементы – конструкции. Конструкции имеют такие атрибуты: тип языковой конструкции, набор кодов. Атрибут тип ( $type.l$ ) может принимать такие значения: к-слово, к-предложение, к-абзац, к-текст.

Отношение подстановки – отношение с атрибутами  $_{w_i} l_i \rightarrow _{w_j} l_j$ , где  $l_i, l_j$  – сентенциальные формы [11].

Для заданной формы  $w_l l =_{w_0} \cdot(w_1 l_1, w_2 l_2, \dots, w_h l_h, \dots, w_k l_k)$  и доступного отношения подстановки  $w_p \rightarrow (w_h l_h, w_q l_q)$  такого, что  $w_h l_h$  – подформа  $w_l l$  результатом трехместной операции подстановки  $w_l^* l^* =_{w_p} \Rightarrow (w_h l_h, w_q l_q, w_l l)$  будет форма  $w_l^* l^* =_{w_0} \cdot(w_1 l_1, w_2 l_2, \dots, w_q l_q, \dots, w_k l_k)$ .

Двухместная операция частичного вывода  $w_l^* l^* =_{w_p} \Rightarrow (\Psi, w_l l)$  ( $\Rightarrow \in \Theta$ ) заключается в:

- выборе одного из доступных правил подстановки  $\psi_r : \langle s_r, g_r \rangle$  с отношениями подстановки  $s_r$ ;
- выполнении на его основе операций подстановки;
- выполнении операций над атрибутами  $g_r$  в соответствующей последовательности.

Операция полного вывода или просто вывода ( $\Rightarrow \in \Theta$ ) заключается в пошаговом преобразовании форм, начиная с начального нетерминала и заканчивая конструкцией, удовлетворяющей условию окончания вывода, что подразумевает циклическое выполнение операций частичного вывода. Операция двухместная  $w_l^* l^* = \Rightarrow (\Psi, w_l l)$ .

Операция присваивания  $\mathbf{:=}(a, b)$  копирует значение операнда  $b$  в  $a$ .

Операция  $= (w_i, w_j)$  выполняет сравнения атрибутов. Результатом операции сравнения является значение «истина», если  $w_i = w_j$ , иначе – «ложь».

Элемент принадлежит форме  $m \in l$ , если  $\exists l_1, l_2, l_3 : l = \cdot(l_1, l_2) \& ((l_1 = \cdot(l_3, m) | l_1 = \cdot(m, l_3)) | l_2 = \cdot(l_3, m))$ , где  $l_1, l_2, l_3$  – формы, образованные с помощью операции конкатенации.

**Расширение аксиоматики носителя.** Конструкция является к-словом  $type \dashv l = cw$ , если  $\forall m_i \in l : m_i \in T_K \cup T_L \cup T_N$ .

Конструкция является к-предложением  $type \dashv l = cs$ , если  $l = \cdot(l_1, m)$ ,  $m \in T_{SD}$ ,  $T_{SD} = \{"!", "?", ".", "..."\}$ ,  $T_{SD} \subset T_D$ .

Конструкция является к-абзацем  $type \dashv l = cp$ , если  $l = \cdot(l_1, l_2) \& l_2 = \cdot(m_1, m_2)$ ,  $code \dashv m_1 = 13$ ,  $code \dashv m_2 = 10$  (переход на новую строку, возврат каретки).

К-абзац может включать в себя несколько предложений, к-предложение – несколько абзацев.

Конструкция является к-текстом  $type \dashv l = ct$ , если  $l = \cdot(l_1, l_2) \& type \dashv l_1 = cp \& (type \dashv l_2 = cp \mid l_2 = \varepsilon)$  и  $l_1, l_2$  имеют смысловую связь.

**Интерпретация КПС текста.** Частично интерпретируем структуру  $C_T$  (1) с помощью алгоритмической структуры  $C_A$ :

$$\langle C_T, C_{A,T} = \langle M_{A,T}, V_{A,T}, \Sigma_{A,T}, \Lambda_{A,T} \rangle \rangle_I \mapsto \langle {}_A C_T, {}_A C_T = \langle M_1, \Sigma_1, \Lambda_1 \rangle \rangle, \quad (2)$$

где  $\Lambda_1 \supset \Lambda_T$ ,  $V_{A,T} = \{A_i^0 | X_i\}$  – множество базовых алгоритмов [10, 13],  $X_i, Y_i$  – множества определений и значений алгоритма  $A_i^0 | X_i$ .  $\Lambda_{A,T} = \{M_{A,T} = \bigcup_{A_i^0 \in V_{A,T}} (X(A_i^0) \cup Y(A_i^0)) \cup \Omega(C_T)\}$  – неоднородный носитель,  $\Omega(C_T)$  – множество языковых конструкций, которые удовлетворяют  $C_T$ ;  $\Lambda_1 = \{(A_3 |_{l_1, l_2}^l \leftarrow "!" \Rightarrow"); (A_4 |_{l_h, l_q, f_i}^{f_i} \leftarrow "!" \Rightarrow"); (A_5 |_{f_i, \Psi}^{f_j} \leftarrow "!" \Rightarrow"); (A_6 |_{\sigma, \Psi}^{\bar{\Omega}} \leftarrow "!" \Rightarrow"); (A_7 |_{a,b}^a \leftarrow "!" \Rightarrow"); (A_8 |_{a,b}^c \leftarrow "!" \Rightarrow)\}$ .

Структура  ${}_A C_T$  включает алгоритмы выполнения операций:

- $A_1^0 |_{A_i, A_j}^{A_i \cdot A_j}$  – композиция алгоритмов,  $A_i \cdot A_j$  – последовательное выполнение алгоритма  $A_j$  после алгоритма  $A_i$ ;
- $A_2^0 |_Z^{A_1}$  – условное выполнение: алгоритм  $A_i$  выполняется, если условие  $Z$  истинно;
- $A_3 |_{l_1, l_2}^l$  – конкатенация,  $l, l_1, l_2$  – формы;

- $A_4 |_{l_h, l_q, f_i}^{f_i}$  – подстановка,  $l_h, l_q, f_i$  – формы;
- $A_5 |_{f_i, \Psi}^{f_j}$ ,  $A_6 |_{\sigma, \Psi}^{\bar{\Omega}}$  – частичный и полный вывод, где  $f_i, f_j$  – формы,  $\sigma$  – аксиома,  $\bar{\Omega}$  – множество сформированных конструкций;
- $A_7 |_{a, b}^a$  – присвоение операнду  $a$  значение операнда  $b$ ;
- $A_8 |_{a, b}^c$  – сравнение атрибутов  $a$  и  $b$ , если  $a = b$ , то  $c = true$ , в противном случае  $c = false$ .

**Конкретизация КПС текста.** Выполним конкретизацию структуры  $C_T$ :

$$C_T \hookrightarrow {}_K C_T = \langle M_2, \Sigma_2, \Lambda_2 \rangle,$$

где  $\Lambda_2 \supset \Lambda_1$ ,  $\Lambda_2 \supset \{M_2 = T_T \cup N; N = \{w_i\} = \{\alpha, \gamma, \delta, \kappa, \pi, \rho, \sigma, \tau\}\}$ . Атрибутом нетерминала ( $w_i$ ) является тип конструкции ( $kind \downarrow w_i$ ), для построения которой он используется.

Далее рассмотрим правила для построения к-текста, а также других конструкций, которые могут в него входить.

Правило  $s_1$  позволяет начать выполнение построения текста, определив значение соответствующего атрибута как текст:

$$s_1 = \langle \sigma \rightarrow \tau \rangle, g_1 = \langle kind \downarrow \tau := ct \rangle.$$

Правило  $s_2 - s_4$  позволяют определить составляющие текста – абзацы (одного или много), определив значение соответствующего атрибута как абзац

$$s_2 = \langle \tau \rightarrow \alpha \pi \rangle, g_2 = \langle kind \downarrow \alpha := cp \rangle, s_3 = \langle \pi \rightarrow \alpha \pi \rangle, s_4 = \langle \pi \rightarrow m_1 m_2 \rangle,$$

где  $m_1, m_2$  – символы возврата каретки и перехода на новую строку.

Правила  $s_2 - s_4$  позволяют определить составляющие абзаца – предложение (много или одно), определив значение соответствующего атрибута как предложение

$$s_5 = \langle \alpha \rightarrow \rho \gamma \alpha \rangle, s_6 = \langle \alpha \rightarrow \rho \gamma \rangle, g_6 = \langle kind \downarrow \rho := cs \rangle, s_7 = \langle \gamma \rightarrow m_{end} \rangle,$$

где  $m_{end} \in T_{SD}$  – символ признак окончания предложения.

Правила  $s_8 - s_9$  позволяют определить составляющие предложения: одно или более слов

$$s_8 = \langle \rho \rightarrow \kappa \rangle, g_8 = \langle kind \downarrow \kappa := cw \rangle, s_9 = \langle \rho \rightarrow \kappa \delta \rho \rangle.$$

Правила  $s_{10} - s_{11}$  позволяют определить разделитель между словами

$$s_{10} = \langle \delta \rightarrow m_{sep} \rangle, s_{11} = \langle \delta \rightarrow m_{sep} \delta \rangle,$$

где  $m_{sep} \in T_{WD}, T_{WD} = T_D \setminus T_{SD}$ .

Правила  $s_{12} - s_{13}$  позволяют построить слово из одной и более букв:

$$s_{12} = \langle \kappa \rightarrow m_c \kappa \rangle, s_{13} = \langle \kappa \rightarrow m_c \rangle,$$

где  $m_c \in T_K \cup T_L \cup T_N$ .

Тут под записью типа  $\kappa \rightarrow m_c, m_c \in T_K \cup T_L \cup T_N$  следует понимать множество альтернативных правил, что можно также записать в виде  $\kappa \rightarrow a | b | c | \dots$ .

Операции над атрибутами выполняются после операции частичного вывода.

**Реализация КПС текста.** Реализация структуры (2) заключается в формировании языковых конструкций из элементов ее носителя путем выполнения алгоритмов, связанных с операциями сигнатуры, по правилам аксиоматики:

$${}_A C_T \underset{R}{\mapsto} \overline{\Omega}({}_A C_T),$$

где  $\overline{\Omega}({}_A C_T) \subset \Omega({}_A C_T)$ .

Рассмотрим пример построения конструкции «Колпак под колпаком». Далее показан вывод текста, состоящего из одного абзаца и одного предложения, заканчивающегося точкой:

$$\sigma \overset{1}{\Rightarrow} \tau \overset{2}{\Rightarrow} \alpha \pi \overset{4}{\Rightarrow} \alpha' 13' 10' \overset{6}{\Rightarrow} \rho \gamma' 13' 10' \overset{7}{\Rightarrow} \rho' . 13' 10'$$

Определим количество слов и структуру предложения:

$$\rho' . 13' 10' \overset{9}{\Rightarrow} \kappa \delta \rho' . 13' 10' \overset{9}{\Rightarrow} \kappa \delta \kappa \delta \rho' . 13' 10' \overset{8}{\Rightarrow} \kappa \delta \kappa \delta \kappa' . 13' 10' \overset{10}{\Rightarrow} \kappa' ' \kappa' ' \kappa . 13' 10' .$$

Далее используя правила 12, 13 получаем фразу «Колпак'32'под'32'колпаком.'13''10'». Здесь в кавычки взяты символы пробела и перехода к новому абзацу.

Для обработки конструкций рассмотрим задачи синтеза и анализа графовой модели текста [8].

## Синтез графа конструктивно-продукционными структурами

Пусть есть множество языковых конструкций  $\overline{\Omega}({}_A C_T)$ , порожденное структурой  ${}_A C_T$  (2). Задача состоит в определенные структуры  $C_g$ , порождающей множество конструкций-графов  $\overline{\Omega}(C_g)$  такое, что существует биективное отображение  $f : \overline{\Omega}({}_A C_T) \rightarrow \overline{\Omega}_g(C_g)$ .

Для решения поставленной задачи определим структуру и специализируем ее соответствующим образом:

$$C = \langle M, \Sigma, \Lambda \rangle \underset{S}{\mapsto} C_g = \langle M_g, \Sigma_g, \Lambda_g \rangle,$$

где  $M_g$  – расширяемый носитель, включающий множества конструкций-графов, языковых конструкций и их элементов,  $\Sigma_g$  – множество операций и отношений на элементах  $M_g$ ,  $\Lambda_g$  – аксиоматика.

**Частичная аксиоматика носителя.** Носитель включает множества терминальных и нетерминальных элементов  $M_g = T_g \cup N_g$ . Терминалами являются языковые конструкции, построенные КПС (2) и их составляющие ( $T_T$ ), а также конструкции графов и их составляющих  $T_g = \overline{\Omega} \cup \Omega_g \cup T_T \cup V \cup E$ ,  $\Omega_g$  – множество конструкций-графов,  $V$ ,  $E$  – множества вершин и дуг с их атрибутами.

Вершина имеет атрибуты  $\overline{W}_v = \langle id, content, tokens \rangle$ ,  $id$  – идентификатор, принимает целочисленные значения,  $content$  – часть текстовой конструкции,  $tokens$  – список, содержащий признаки начала языковых конструкций. Атрибуты дуги  $\overline{w}_e = \langle id, routes, start, end \rangle$ ,  $id$  – идентификатор, принимает целочисленные значения,  $routes$  – множество номеров путей, в которые входит дуга (указывает на порядок обхода графа),  $start$ ,  $end$  – вершин, которые являются началом и концом дуги.

Нагруженный граф будем обозначать  $\overline{w}_g G = \langle V, E \rangle$ ,  $V = \{ \overline{w}_v v_i \}$ ,  $E = \{ \overline{w}_e e_j \}$  – множества вершин и дуг, нагруженных атрибутами. Каждое множество содержит пустой элемент.

Граф имеет такие атрибуты  $\overline{w}_g = \langle start\_v, last\_v, current\_v, amount\_l \rangle$ , где  $start\_v$  – стартовая вершина графа,  $last\_v$  – последняя добавленная вершина,  $current\_v$  – текущая вершина при формировании графа,  $amount\_l$  – количество циклов, в которые входит стартовая вершина.

**Частичная аксиоматика операций.** Рассмотрим сигнатуру  $\Sigma_g$ :

$$\Sigma_g = \langle \Xi_g, \Theta_g, \Phi, \{\rightarrow\} \rangle \cup \Psi_g,$$

где  $\Xi_g = \{ ;, \equiv, \doteq, \tilde{\cup}, \tilde{\cup} \} – множество операций преобразования и связывания, \Theta_g = \{ \Rightarrow, | \Rightarrow, || \tilde{\Rightarrow} \} – множество операций вывода, \Phi_g = \{ \div, \doteq, \#, + \} – множество операций над атрибутами, \Psi_g – множество правил продукции вида \psi_i : s_i, g_i, i – номер правила, s – последовательность операций подстановки, g – последовательность операций над атрибутами, \ll \rightarrow \gg – отношение подстановки.$

Операция  $e \doteq (v_1, v_2, G)$  заключается в определении дуги  $e$ , соединяющей вершины  $v_1$  и  $v_2$  графа  $G$ . Если дуга со значениями соответствующих атрибутов отсутствует, то возвращается пустой элемент множества дуг графа  $G$ , его идентификатор равен нулю.

Операция  $v \doteq (x, V)$  заключается в нахождении вершины  $v$  с атрибутом веса, равным  $x$ , из множества вершин  $V$ . Если вершина со значениями соответствующих атрибутов отсутствует, то возвращается пустой элемент множества вершин  $V$ , его идентификатор равен нулю.

Операция  $\div(c, n, L)$  заключается в выполнении  $n$  операций из списка  $L$ , если  $c = true$ .

Операция вычисления мощности множества  $\#Q$  определяет число, равное количеству элементов в  $Q$ .

Операция сложения двух чисел  $+(a, b)$  предполагает нахождение третьего числа, являющегося их суммой.

Операция объединения графов  $\overline{w}_g G = \widetilde{\cup}(\overline{w}_1 G_1, \overline{w}_2 G_2)$  предполагает формирование нового графа  $\overline{w}_g G$ , включающего объединенные множества вершин и дуг исходных графов  $\overline{w}_g G = \langle V, E \rangle$ ,  $V = V_1 \cup V_2$ ,  $E = E_1 \cup E_2$ ,  $\overline{w}_1 G_1 = \langle V_1, E_1 \rangle$ ,  $\overline{w}_2 G_2 = \langle V_2, E_2 \rangle$ , при этом  $\cup$  – традиционная операция объединения множеств.

Отношение подстановки имеет вид

$$\psi_i = \langle s_i, g_i \rangle, \quad s_i = \langle \bar{s}_i, \tilde{s}_i \rangle, \quad g_i = \langle \bar{g}_i, \tilde{g}_i \rangle,$$

где  $\bar{s}_i$ ,  $\tilde{s}_i$  – отношение подстановки для распознавания языковой конструкции и построения конструкции графа соответственно,  $\bar{g}_i$ ,  $\tilde{g}_i$  – операции над атрибутами языковой конструкции и графа, его вершин и дуг соответственно. В случае если операции над атрибутами не выполняются, отношение подстановки имеет вид  $\psi = \langle s, \epsilon \rangle$ .

Операция полного вывода  $\|\tilde{\rightarrow}\| (\Psi, {}_{w_l} l)$  состоит из:

- определении входной конструкции, набора отношений подстановки и доступных из них;
- выполнении операции частичного вывода, пока языковая конструкция  $\omega \in \overline{\Omega}$  полностью не распознана, то есть для каждого ее элемента  $\omega_i$  в графе нет соответствующего элемента-вершины или форма графа содержит нетерминалы.

Результатом операции вывода является конструкция-граф.

**Конкретизация КПС графа.** Выполним конкретизацию структуры  $C_g$ :

$$C_g \mapsto {}_K C_g = \langle M_3, \Sigma_3, \Lambda_3 \rangle,$$

где  $\Lambda_3 \supset \Lambda_g$ ,  $\Lambda_3 \supset \{M_3 \supset M_g, c \in T_T, N_g = \{\alpha, \delta\}, T_g \supset \{G, G^*, G^{**}\}, G = \langle V, E \rangle, V = \{v\}, E = \emptyset, G^* = \langle V^*, G^* \rangle, V^* = \{v_1^*, v_2^*\}, E^* = \{e_1^*\}, G^{**} = \langle V^{**}, G^{**} \rangle, V^{**} = \{v_1^{**}, v_2^{**}\}, E^{**} = \{e_1^{**}\}\}$ .

Для распознавания языковых конструкций определим такие правила:

$$\bar{s}_1 = \langle \sigma_{d_1} \rightarrow c \sigma \rangle, \quad \bar{g}_1 = \langle \div(\text{code} \dashv c \neq EOF, 1, d_1 := true) \rangle, \quad \bar{s}_2 = \langle \sigma_{d_2} \rightarrow c \rangle, \quad \bar{g}_2 = \langle \div(\text{code} \dashv c = EOF, 1, d_2 := true) \rangle,$$

где  $c$  – символ текста, кроме  $EOF$  – признак конца текста в его электронном представлении.

Следующие правила описывают добавление первой вершины в пустой граф:

$$\tilde{s}_1 = \langle \sigma \tilde{\rightarrow} G \alpha \rangle, \quad \tilde{g}_1 = \langle id \dashv v := \#V, content \dashv v := c, tokens := \langle cw, cs, cp, ct \rangle \rangle,$$

$$start\_v \dashv G := v, current\_v \dashv G := v, last\_v \dashv G := v, amount\_l \dashv G := 0 \rangle.$$

Правило  $\tilde{s}_2$  позволяет добавить к графу новую вершину и дугу, связывающую новую вершину с текущей в графе

$$\tilde{s}_2 = \langle G \alpha \underset{d_1}{\tilde{\rightarrow}} \widetilde{\cup}(G, G^*) \alpha \rightarrow G \alpha \rangle, \quad \tilde{g}_2 = \langle v_1 \doteq (c, V), e_1 \doteq (current\_v \dashv G, v, G) \rangle,$$

$$\begin{aligned} & \div (id \downarrow v_1 = 0 \& id \downarrow e_1 = 0, 14, d_1^* := true, id \downarrow v_1^* := id \downarrow current\_v \downarrow G, content \downarrow v_1^* := content \downarrow current\_v \downarrow G, \\ & id \downarrow v_2^* := \#V + 1, content \downarrow v_2^* := c, id \downarrow e_1^* := \#E + 1, start \downarrow e_1^* := v_1^*, end \downarrow e_1^* := v_2^*, routes \downarrow e_1^* := \{amount\_l \downarrow G\}, \\ & \div (content \downarrow current\_v \downarrow G = x \& x \in T_{WD}, tokens \downarrow v_2^* := \langle cw \rangle), \\ & \div (content \downarrow current\_v \downarrow G = x \& x \in T_{SD}, tokens \downarrow v_2^* := \langle cs \rangle), \\ & \div (content \downarrow current\_v \downarrow G = '10', tokens \downarrow v_2^* := \langle cp \rangle), last\_v \downarrow G := v_2, current\_v \downarrow G := v_2 \rangle). \end{aligned}$$

Правило  $\tilde{s}_3$  позволяет добавить к графу новую дугу, связывающую текущую вершину со стартовой.

$$\begin{aligned} \tilde{s}_3 = & \left\langle \alpha \xrightarrow{d_2^*} \tilde{\cup}(G, G^{**})\alpha, \tilde{\cup}(G, G^{**})\alpha \rightarrow G\alpha \right\rangle, \tilde{g}_3 = \langle v_1 := (c, V), e_1 := (current\_v \downarrow G, v_1, G), \right. \\ & \div (v_1 = start\_v \downarrow G \& id \downarrow e_1 = 0, 14, d_2^* := true, id \downarrow v_1^{**} := id \downarrow current\_v \downarrow G, content \downarrow v_1^{**} := \\ & := content \downarrow current\_v \downarrow G, id \downarrow v_2^{**} := id \downarrow start\_v \downarrow G, content \downarrow v_2^{**} := c, id \downarrow e_1^{**} := \#E + 1, start \downarrow e_1^{**} := \\ & := v_1^{**}, end \downarrow e_1^{**} := v_2^{**}, routes \downarrow e_1^{**} := \{amount\_l \downarrow G\}, amount\_l \downarrow G := amount\_l \downarrow G + 1, \\ & \div (content \downarrow current\_v \downarrow G = x \& x \in T_{WD}, 1, tokens \downarrow v_2^{**} := \cdot(tokens \downarrow start\_v \downarrow G, \langle cw \rangle)), \\ & \div (content \downarrow current\_v \downarrow G = x \& x \in T_{SD}, 1, tokens \downarrow v_2^{**} := \cdot(tokens \downarrow start\_v \downarrow G, \langle cs \rangle)), \\ & \div (content \downarrow current\_v \downarrow G = '10', 1, tokens \downarrow v_2^{**} := \cdot(tokens \downarrow start\_v \downarrow G, \langle cp \rangle)), current\_v \downarrow G = start\_v \downarrow G \rangle. \end{aligned}$$

Правило  $\tilde{s}_4$  позволяет изменить нагрузку имеющейся дуги:

$$\begin{aligned} \tilde{s}_4 = & \left\langle G\alpha \xrightarrow{d_3^*} G\alpha \right\rangle, \tilde{g}_4 = \langle v_1 := (c, V \downarrow G), e_1 := (current\_v \downarrow G, v_1, G), \div (id \downarrow v_1 \neq 0 \& id \downarrow e_1 \neq 0, 3, d_3^* := true, \right. \\ & \div (start\_v \downarrow G \neq v_1, 5, routes \downarrow e_1 := routes \downarrow e_1 \cup \{amount\_l \downarrow G\}), \div (content \downarrow current\_v \downarrow G = x \& x \in T_{WD}, 1, \\ & tokens \downarrow v_1 := \cdot(tokens \downarrow v_1, \langle cw \rangle)), \div (content \downarrow current\_v \downarrow G = x \& x \in T_{SD}, 1, tokens \downarrow v_1 := \cdot(tokens \downarrow v_1, \langle cs \rangle)), \\ & \div (content \downarrow current\_v \downarrow G = '10', 1, tokens \downarrow v_1 := \cdot(tokens \downarrow v_1, \langle cp \rangle)), current \downarrow v := v_1), \div (start\_v \downarrow G = v, 6, \\ & routes \downarrow e_1 := routes \downarrow e_1 \cup \{amount\_l \downarrow G\}), \div (content \downarrow current\_v \downarrow G = x \& x \in T_{WD}, 1, tokens \downarrow v_1 := \cdot(tokens \downarrow v_1, \langle cw \rangle)), \\ & \div (content \downarrow current\_v \downarrow G = x \& x \in T_{SD}, 1, tokens \downarrow v_1 := \tilde{\cup}(tokens \downarrow v_1, \langle cs \rangle)), \div (content \downarrow current\_v \downarrow G = '10', \\ & tokens \downarrow v_1 := \cdot(tokens \downarrow v_1, \langle cp \rangle)), current \downarrow v := start\_v \downarrow G, amount\_l \downarrow G := amount\_l \downarrow G + 1) \rangle. \end{aligned}$$

Следующее правило позволяют завершить процесс построения конструкции-графа:

$$\tilde{s}_5 = \langle \alpha \xrightarrow{\sim} \varepsilon \rangle.$$

**Интерпретация КПС графа.** Интерпретируем графовую структуру:

$$\langle C_g, C_{A,G} = \langle M_A, V_A, \Sigma_A, \Lambda_A \rangle \rangle_I \mapsto \langle {}_A C_g, {}_A C_g = \langle M_g, \Sigma_g, \Lambda_g \rangle \rangle,$$

где  $\Lambda_4 \supset \Lambda_3$ ,  $V_A = \{A_i^0 | X_i\}$  – множество базовых алгоритмов [10, 13],  $X_i, Y_i$  – множества определений та значений алгоритма  $A_i^0 | X_i$ .  $\Lambda_A = \{M_A = \bigcup_{A_i^0 \in V_A} (X(A_i^0) \cup Y(A_i^0)) \cup \Omega(C_l) \cup \Omega(C_g)\}$  – неоднородный носитель,

$\Omega(C_l), \Omega(C_g)$  – множество языковых и графовых конструкций, которые удовлетворяют  $C_l$  и  $C_g$  соответственно;  $\Lambda_4 \supset \{(A_6 \leftarrow "||\tilde{\Rightarrow}"); (A_9 \leftarrow ":="); (A_{10} \leftarrow ":="); (A_{11} \leftarrow "||"); (A_{12} \leftarrow "||\#"); (A_{13} \leftarrow "||+"), (A_{14} \leftarrow "||\tilde{U}"), (A_{15} \leftarrow "||U")\}$ .

Структура  ${}_A C_g$  включает алгоритмы, реализующие такие операции:

- $A_1^0, A_2^0, A_3 - A_5, A_7 - A_8$  – аналогичны алгоритмам структуры  $C_A$ ;
- $A_6 |_{\sigma, \Psi_g}^{\bar{\Omega}}$  – полный вывод, где  $\sigma$  – аксиома,  $\bar{\Omega}$  – множество сформированных конструкций;
- $A_9 |_{v_1, v_2, G}^e$  – определение дуги  $e$ , соединяющей две заданные вершины  $v_1, v_2$  в графе  $G$ ;
- $A_{10} |_{x, V}^v$  – нахождение вершины  $v$  с заданным значением атрибута веса  $x$  в множестве вершин  $V$ ;
- $A_{11} |_{c, n, L}^L$  – выполнение  $n$  операций из списка  $L$ , если  $c = true$ ,  $L$  – список из  $n$  операций;
- $A_{12} |_Q^x$  – вычисление мощности  $x$  множества  $Q$ ;
- $A_{13} |_{a, b}^c$  – сложения двух чисел  $a, b$ ,  $c$  – результат;
- $A_{14} |_{G_1, G_2}^G$  – объединение графов  $G_1, G_2$  в  $G$ ;
- $A_{15} |_{Q_1, Q_2}^Q$  – объединение множеств  $Q_1, Q_2$  в  $Q$ .

**Реализация КПС графа.** Реализация структуры  ${}_A C_g$  заключается в формировании графовых конструкций, которые имеют однозначное соответствие конструкциям  $\bar{\Omega}({}_A C_T)$  путем выполнения алгоритмов, связанных с операциями сигнатуры, по правилам аксиоматики:

$${}_A C_g \ R \mapsto \bar{\Omega}({}_A C_g),$$

где  $\bar{\Omega}({}_A C_g) \subset \Omega({}_A C_g)$ .

## Сжатие графа

Данная графовая модель может быть использована для сравнения текстов, поиска подстроки в строке. Для ускорения процесса сравнения можно применить сжатия графа. Для этого рассмотрим такую структуру:

$$C = \langle M, \Sigma, \Lambda \rangle_S \mapsto C_C = \langle M_C, \Sigma_C, \Lambda_C \rangle,$$

где  $\Lambda_C = \{M_C = V \cup E \cup \Omega_g, \Sigma_C = \{\{, :=\}, \{\rightarrow, \Rightarrow, ||\tilde{\Rightarrow}\} \cup \Psi_C\}, \Psi_C = \{\psi_C = \langle s_i, g_i \rangle\}\}$ .

### Частичная аксиоматика операций.

Отношение подстановки и операции подстановки, частичного и полного вывода, конкатенации, присваивания аналогичны одноименным операциям КПС графа.

**Интерпретация КПС сжатия графа.** Алгоритмы операций, использованы в рассматриваемой структуре, аналогичны алгоритмам одноимённых операций структуры для построения графа.

**Конкретизация КПС сжатия графа.** Выполним конкретизацию структуры для сжатия графа:

$$C_C \ K \mapsto C_{CK} = \langle M_{CK}, \Sigma_{CK}, \Lambda_{CK} \rangle,$$

где  $M_{CK} = M_C \cup N_{CK}, N_{CK} = \{\sigma\}$ .

Рассмотрим правила, позволяющие сжать граф.

$$s_1 = \langle \sigma \rightarrow G\sigma \rangle, \quad s_2 = \langle G\sigma_d \rightarrow G\sigma \rangle,$$

$g_1 = g_2 = \langle e_1 := (v_i, v_j, G), e_2 := (v_j, v_k, G), \div (routes \leftarrow e_1 = routes \leftarrow e_2, 4, d := true, content \leftarrow v_i := \cdot (content \leftarrow v_i, content \leftarrow v_j),$

$$tokens \leftarrow v_i := \cdot (tokens \leftarrow v_i, tokens \leftarrow v_j), end \leftarrow e_1 := v_k \rangle \rangle, \quad s_3 = \langle \sigma \rightarrow \varepsilon \rangle.$$

**Реализация КПС сжатия графа.** Реализация интерпретируемой структуры  ${}_A C_C$  заключается в формировании графовых конструкций, которые имеют однозначное соответствие конструкциям  $\overline{\Omega}({}_A C_T)$  и  $\overline{\Omega}({}_A C_g)$  путем выполнения алгоритмов, связанных с операциями сигнатуры, по правилам аксиоматики:

$${}_A C_C \ R \mapsto \overline{\Omega}({}_A C_C),$$

где  $\overline{\Omega}({}_A C_C) \subset \Omega({}_A C_C)$ .

## Выводы

Проблема поиска информации, ее сравнении требует автоматизированных средств решения. С данной целью целесообразной является разработка средств формализации для представления текстов.

Формальные грамматические и алгоритмические структуры использованы как средство формализации текста в виде конструкций с атрибутами. Конструктивно-продукционный подход позволяет определить порядок конструирования и обработки текста, представить его в виде графовых структур с целью снижения алгоритмической сложности.

В рамках рассмотренного атрибутивного подхода могут быть решены задачи по таким направлениям, как семантический поиск и сопоставление текстов для выявления заимствований, что является ключевыми моментами в решении проблем пластиата.

1. Efstatios Stamatatos Plagiarism Detection Based on Structural Information – CIKM’11, October 24–28, 2011, Glasgow, Scotland, UK.
2. Leilei Kong, Zhimao Lu, Haoliang Qi and Zhongyuan Han. Detecting High Obfuscation Plagiarism: Exploring Multi-Features Fusion via Machine Learning // International Journal of u-and e-Service, Science and Technology – 2014. – Vol. 7, N 4. – P. 385–396.
3. Ahmed Hamza Osman, Naomie Salim, Mohammed Salem Binwahlan. Plagiarism Detection Using Graph-Based Representation // JOURNAL OF COMPUTING. – 2010. – Vol. 2, N 4. – P. 36 – 41.
4. Bin-Habtoor A.S., Zaher M.A. A Survey on Plagiarism Detection Systems [Text] // International Journal of Computer Theory and Engineering. – 2012. – Vol. 4, N 2. – P. 185–188.
5. Ahmed Hamza Osman, Naomie Salim, Mohammed Salem Binwahlan, Hanza Hentably, Albaraa M. Ali. Conceptual Similarity and graph-based method for plagiarism detection // Journal of Theoretical and Applied Information Technology. – 31st October 2011. – Vol. 32, N 2. – P. 135–145.
6. Mozgovoy M. Maxim Mozgovoy Desktop Tools for Offline Plagiarism Detection in Computer Programs // Informatics in Education. – 2006. – Vol. 5, N 1. – P. 97–112.
7. Mozgovoy M., Kakkonen T., Cosma G. Automatic Student Plagiarism Detection: Future Perspectives // Journal of Educational Computing Research. – 2010. – Vol. 43(4). – P. 507–527.
8. Шинкаренко В.І., Куропятник О.С. Система контролю плагіату в студентських роботах // Восточно-Европейский журнал передовых технологий. – Харьков: Технологический центр, 2012. – № 4/2 (58). – С. 32–36.
9. Ільман В.М., Скалоуб В.В., Шинкаренко В.І. Формальні структури та їх застосування: монографія. – Д.: Вид-во Дніпропетр. нац. ун-ту залізн. трансп. ім. акад. В. Лазаряна, 2009. – 205 с.
10. Шинкаренко В.І., Ільман В.М. Конструктивно-продукционные структуры и их грамматические интерпретации. I. Обобщенная формальная конструктивно-продукционная структура // Кибернетика и системный анализ. – Киев, 2014. – Т. 50, № 5 – С. 8–16.
11. Шинкаренко В.І. , Васецкая Т.Н. Моделирования процесса адаптации алгоритмов сжатия средствами конструктивно-продукционных структур // Кибернетика и системный анализ. – 2015. – Т. 51. – № 6. – С. 19–34.
12. Шинкаренко В.І., Ільман В.М., Забула Г.В. Конструктивно-продукционная модель структур данных на логическом уровне // Проблемы программирования. – 2014. – № 2–3. – С. 10–16.
13. Шинкаренко В.І., Ільман В.М. Конструктивно-продукционные структуры и их грамматические интерпретации. II. Уточняющие преобразования // Кибернетика и системный анализ. – 2014. – № 6. – С. 15–28.
14. Шинкаренко В.І., Ільман В.М., Скалоуб В.В. Структурные модели алгоритмов в задачах прикладного программирования. Часть I. Формальные алгоритмические структуры // Кибернетика и системный анализ. – 2009. – № 3. – С. 3–14.

## References

1. EFSTATHIOS, STAMATATOS (2011) Plagiarism Detection Based on Structural Information in CIKM’11. Glasgow, Scotland, UK October 24–28, 2011.
2. LEILEI KONG & ZHIMAO LU & HAOLIANG QI & ZHONGYUAN HAN (2014) Detecting High Obfuscation Plagiarism: Exploring Multi-Features Fusion via Machine Learning. International Journal of u-and e-Service, Science and Technology. 7 (4). P. 385–396.
3. AHMED HAMZA OSMAN & NAOMIE SALIM & MOHAMMED SALEM BINWAHLAN (2010) Plagiarism Detection Using Graph-Based Representation Journal of Computing. 2 (4). P. 36 – 41.
4. A. S. BIN-HABTOOR & M. A. ZAHER (2012) A Survey on Plagiarism Detection Systems. International Journal of Computer Theory and Engineering. 4 (2). P. 185 – 188.
5. AHMED HAMZA OSMAN & NAOMIE SALIM & MOHAMMED SALEM BINWAHLAN & HANZA HENTABLY & ALBARAA M. ALI (2011) Conceptual Similarity and graph-based method for plagiarism detection. Journal of Theoretical and Applied Information Technology. 32 (2). P. 135–145.
6. MOZGOVOY, M. (2006) Desktop Tools for Offline Plagiarism Detection in Computer Programs. Informatics in Education. 5 (1). pp. 97–112
7. M. MOZGOVOY & T. KAKKONEN & G. COSMA (2010) Automatic Student Plagiarism Detection: Future Perspectives. Journal of Educational Computing Research. 43 (4). P. 507-527.

8. V. SHYNKARENKO & E. KUROPYATNICK (2012) Monitoring system of plagiarism in student works. East Europe Journal of Enterprise Technologies. 4/2 (58). p. 32 – 36.
9. V. ILMAN, V. SKALOZUB & V. SHYNKARENKO Formal structures and their applications: monograph. Dnipropetrovsk: DNURT named after academician V. Lazaryan.
10. SHYNKARENKO V. & ILMAN V. (2014) Constructive-synthesizing structures and their grammatical interpretations. I. Generalized formal constructive-synthesizing structure. Cybernetics and Systems Analysis. 50 (5). P. 8–16.
11. SHYNKARENKO V. & VASETSKA T. (2015) Modeling the Adaptation of Compression Algorithms by Means of Constructive-Synthesizing Structures. Cybernetics and Systems Analysis. 51(6). P. 19 – 34.
12. SHYNKARENKO V. I., ILMAN V. M., ZABULA H. V. (2014) Logical view for construction-synthesis model of data structeres. Problems in programming. 2-3 Special issue.
13. SHYNKARENKO V. & ILMAN V (2014). Constructive-production structures and their grammatical interpretations. II. Clarifying conversions. Cybernetics and Systems Analysis. 6. P. 15 – 28.
14. SHYNKARENKO V. & ILMAN V., SKALOZUB V. (2009) Structural models of algorithms in problems of applied programming. i. formal algorithmic structures. Cybernetics and Systems Analysis. 3. P. 3 – 14.

### ***Об авторах:***

*Шинкаренко Виктор Иванович,*  
доктор технических наук, профессор,  
заведующий кафедрой Компьютерные информационные технологии.  
Количество научных публикаций в украинских изданиях – 200.  
Индекс Гирша – 5.  
<http://orcid.org/0000-0001-8738-7225>,

*Куропятник Елена Сергеевна,*  
аспирант кафедры Компьютерные информационные технологии.  
Количество научных публикаций в украинских изданиях – 14.  
Индекс Гирша – 1.  
<http://orcid.org/0000-0003-2286-884X>.

### ***Место работы авторов:***

Днепропетровский национальный университет  
железнодорожного транспорта имени академика В. Лазаряна.  
49010, Днепропетровск, ул. Лазаряна, 2, ДНУЖТ, кафедра КИТ.  
Тел.: (056) 373 1535.  
E-mail: shinkarenko\_v@ua.fm, elenadiit@rambler.ru