

*В. А. ЛОВИЦКИЙ*, канд. техн. наук

## **СТРУКТУРНЫЙ ПОДХОД К РЕШЕНИЮ МОРФОЛОГИЧЕСКИХ ЗАДАЧ**

Построение искусственных систем, понимающих естественный язык, связано с реализацией процессов грамматической обработки отдельных слов языка. Понимание слов и их значений является первым из трех условий, выполнение которых обеспечивает понимание естественноречевого высказывания [1]. Условно можно выделить два подхода к решению морфологических задач (МЗ): алгоритмический [2] и логический [3]. В первом случае МЗ решаются путем составления алгоритмов, с привлечением эвристической информации, во втором — описание процесса решения МЗ осуществляется на языке обобщенных логических функций. Основной недостаток этих подходов состоит в том, что описание процессов решения МЗ связано с конкретным естественным языком и отражает его особенности. Этот недостаток привел, напри-

мер, к необходимости разрабатывать различные системы для морфологической классификации глагольных форм [4] и име прилагательных русского языка [5]. Более того, система для морфологической классификации глагольных форм русского языка не может быть использована для решения таких же МЗ в английском языке без существенного изменения программного обеспечения системы, отражающего особенности естественного языка.

В то же время мозг ребенка легко усваивает любой язык, выделяя закономерности его построения на основании анализ элементов «речевой» среды. Этот общеизвестный факт позволяет утверждать, что мозг ребенка располагает определенными механизмами и структурами, обеспечивающими усвоение особенностей любого естественного языка. Поскольку закономерности языка выявляются в результате анализа относительно небольшого множества конкретных примеров, логично предположить, что человеческий мозг использует индуктивные механизмы формирования понятий [6]. В самом деле, с помощью индуктивного метода легко определить, например, понятие «существительное среднего рода единственного числа именительного падежа», как слова, оканчивающиеся на «О», «Е» или «МЯ». Но с помощью данного понятия можно решать только метазадачи, позволяющие определить принадлежность любого входного слова заданному классу. А следовательно, решение задач преобразования слов (именно к этому классу относится большинство МЗ) требует подхода, отличного от индуктивного формирования понятий функциональным путем. Назовем его структурным методом формирования индуктивных понятий. Суть этого метода состоит в том, что в результате анализа небольшого числа слов, заданного конечным множеством  $X$ , строится не функция, включающая в себя существенные значения признаков всех слов множества  $X$  [6], а структура, которая содержит все слова множества  $X$  и отражает закономерности объединяющую слова  $x_i \in X$  в одно множество.

Продемонстрируем работу алгоритма формирования индуктивных понятий структурным методом на конкретном примере. Структуру будем строить в виде графа, вершины которого представлены буквами слова, а дуги указывают на связь букв в слове. Будем считать, что на множестве  $X$  задано разбиение его на подмножества  $X_i$ , каждое из которых характеризуется некоторым именем понятия  $CNP(<X_i>)$  [6]. Для рассматриваемого примера  $\bigcup_i X_i = X$ , а  $\bigcap_i X_i = \emptyset$ . Итак, пусть  $X = \{\text{торчать, трещать, бить, торчит, торчат, торчал, торчала, торчало, торчали}\}$ , причём  $X_1 = \{\text{торчать, трещать, бить}\}$ ,  $X_2 = \{\text{торчит}\}$ ,  $X_3 = \{\text{торчат}\}$ ,  $X_4 = \{\text{торчал}\}$ ,  $X_5 = \{\text{торчала}\}$ ,  $X_6 = \{\text{торчало}\}$ ,  $X_7 = \{\text{торчали}\}$ , а  $CNP(<X_1>) = \{\langle \text{что делать} \rangle\}$ ,  $CNP(<X_2>) = \{\langle \text{что делает} \rangle, \langle \text{муж. род} \rangle, \langle \text{ед. число} \rangle\}$ ,  $CNP(<X_3>) = \{\langle \text{что делают} \rangle, \langle \text{множ. число} \rangle\}$ ,  $CNP(<X_4>) = \{\langle \text{что делал} \rangle, \langle \text{муж. род} \rangle, \langle \text{ед. число} \rangle\}$ ,  $CNP(<X_5>) = \{\langle \text{что делала} \rangle, \langle \text{жен. род} \rangle, \langle \text{ед. число} \rangle\}$ ,  $CNP(<X_6>) = \{\langle \text{что делало} \rangle, \langle \text{ср. род} \rangle, \langle \text{ед. число} \rangle\}$ ,  $CNP(<X_7>) = \{\langle \text{что делали} \rangle, \langle \text{множ. число} \rangle, \langle \text{ср. род} \rangle\}$ .

$= \{ \langle \text{что делала} \rangle, \langle \text{жен. род} \rangle, \langle \text{ед. число} \rangle \}$ , CNP ( $\langle X_6 \rangle$ ) =  
 $= \{ \langle \text{что делало} \rangle, \langle \text{ср. род} \rangle, \langle \text{ед. число} \rangle \}$ , CNP ( $\langle X_7 \rangle$ ) =  
 $= \{ \langle \text{что делали} \rangle, \langle \text{множ. число} \rangle \}$ . Приведенные характеристики  $X_i$  ни в коей мере не претендуют на полноту и указаны только с целью демонстрации работы алгоритма. Идея работы алгоритма основана на построении древовидной структуры (ДС) [7]. Если ДС строится на основании элементов множества  $X$ , то будем ее называть прямой ДС (ПДС). Если же в качестве элементов множества взять обращенные слова, то получим  $X^R = \{ \text{ТАЧРОТ, БТАЩЕРТ, БТАЖЕБ, ТИЧРОТ, ТАЧРОТ, ПАЧРОТ, АЛАЧРОТ, ОЛАЧРОТ, ИЛАЧРОТ} \}$ , а ДС, построенную на основании этих слов, назовем обратной ДС (ОДС). Совместив эти две структуры соответствующим образом, получим ТВ-структуру (сокращение от английского There and Back: «туда и обратно»).

Алгоритм построения ТВ-структуры начинается с первого обращенного слова множества  $X_1$  и продолжается до тех пор, пока не встретится «незнакомая» буква. Если этой «незнакомой» буквой является первая буква обращенного слова, то она заносится в ОДС, исключается из слова и алгоритм, используя остаток слова, переходит к построению ПДС. Так, для нашего примера построение ТВ-структуры началось со слова «торчать». Вначале в ОДС был занесен «Ь», затем в ПДС — «торчат» и наконец элемент ПДС с «Т» и элемент ОДС с «Ь» были соединены дугой. Для слова «трещать» фрагмент обращенного слова «та» читается в ОДС, фрагмент слова «трещ» заносится в ПДС, элемент ПДС с «Щ» и элемент ОДС с «А» связываются между собой дугой. Таким же образом заносятся все слова множеств  $X_i$ . Иначе говоря, в ОДС заносятся только «незнакомые» последние буквы (у соответствующих им элементов (см. рисунок) заштрихована нижняя часть кружка), а в ПДС — фрагмент слова, который остался или после занесения последней буквы в ОДС, или после чтения в ОДС последовательности букв. Для рассматриваемого примера, пользуясь данным алгоритмом, получена ТВ-структура, представленная на рисунке. Вершины, соответствующие начальным буквам, заштрихованы вверх. При построении ТВ-структуры должно выполняться требование, согласно которому эти вершины не должны иметь входных связей, а вершины, отвечающие последним буквам слова — выходных связей. Названное достигается тем, что в ОДС читаются слова без начальной буквы. В процессе построения ТВ-структуры алгоритм ставит в соответствие каждому элементу указатели прямого и обратного хода, позволяющие читать слова в прямом и обратном направлениях. В ТВ-структуре, приведенной на рисунке, ориентация связей соответствует примерам рассматриваемых ниже задач.

После занесения каждого слова в ТВ-структуру, алгоритм связывает последнюю букву слова с областью, в которой хранятся

коды слов, оканчивающихся на эту букву, и их характеристик т. е.  $CNP(< X_i >)$ .

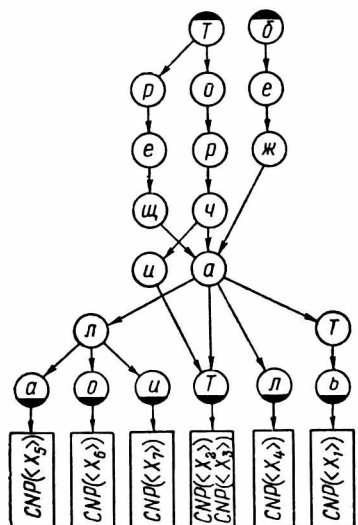
Если построить систему, память которой представлена Т структурой, с ее помощью можно решать следующие задачи.

1. Необходимо охарактеризовать новое слово «бежала». Слово  $x_i$  называется новым, если  $x_i \notin X$ . Легко видеть, что в результате решения этой задачи получим  $CNP(< X_5 >)$ .

Для глаголов, имен существительных и т. д. должны строиться самостоятельные ТВ-структуры. Дано слово «бежала». Его необходимо так преобразовать, чтобы его характеристика, заданная  $CNP(< X_4 >)$ .

При решении задач данного типа система будет совершать «человеческие» ошибки. Так, преобразуя слово «бежала» в соответствии с  $CNP(< X_3 >)$ , система выдаст слово «бежат».

Как видно из рисунка, слова, оканчивающиеся буквой «Т», характеризуются  $CNP(< X_2 >)$  и  $CNP(< X_3 >)$ . В этом случае для установления взаимно-однозначного соответствия между классами слов и характеристиками берется соответствующее число последних букв.



Коды слов и их характеристики

Так, для нашего примера  $CNP(< X_5 >)$  будет характеризовать класс слов, оканчивающихся на «ИТ», а  $CNP(< X_3 >)$  — оканчивающихся на «АТ».

Характеристики слов связываются с начальными буквами. Они имеют более сложный характер и связаны с определением значений слов. В этом читатель сможет легко убедиться, построив самостоятельно ТВ-структуру для  $X = \{ХОДИТЬ, ПРИХОДИТЬ, УХОДИТЬ, ЗАХОДИТЬ, ПОДХОДИТЬ\}$ .

ТВ-структура, совместно с пирамидальной и «И/ИЛИ» структурами [6] представляет собой память диалоговой системы, позволяющей естественный язык (система ДЕСТА) и реализована в Харьковском институте радиоэлектроники.

Список литературы: 1. Лурия А. Р. Основные проблемы нейролингвистики. М., Изд-во Моск. ун-та, 1975. 253 с. 2. Пиотровский Р. Г. Текст, машина, человек. Л., Наука, 1975. 326 с. 3. Шабанов-Кушнаренко Ю. П. Применение метода нуль-органа в лингвистике. — В кн.: Проблемы бионики. Харьков, 1978. Вып. 21, с. 109 — 112. 4. Соловьева Е. А. К вопросу о построении общего алгоритма морфологической классификации глагольных форм русского языка. — В кн.: Проблемы бионики. Харьков, 1975. Вып. 15, с. 143 — 151. 5. Морфологическая классификация имен прилагательных русского языка. М. Ф. Бондаренко, Э. М. Бузницкая, Ю. В. Лопухин, Н. К. Свиридов.