

УДК 62.506.2

*Ю. П. ШАБАНОВ-КУШНАРЕНКО*, д-р техн. наук,  
*Е. А. СОЛОВЬЕВА*

**БИОНИЧЕСКИЕ АСПЕКТЫ МОДЕЛИРОВАНИЯ  
РЕЧЕВОГО ПОВЕДЕНИЯ ЧЕЛОВЕКА**

Математические модели различных способностей человека широко применяются в технике, которая, развиваясь, в свою очередь выдвигает перед бионикой новые задачи. В связи с созданием вычислительных машин, способных частично заменять человека в сфере интеллектуальной деятельности, появился ряд требующих разрешения проблем, в частности проблемы общения человека с машиной и машины между собой.

Несмотря на совершенствование входных и выходных устройств ЭЦВМ, а также методов автоматизации программирования, процесс общения до сих пор недостаточно эффективен. Человеку удобнее выражать свои мысли в языковой форме, поэтому его непосредственное общение с машиной на естественном человеческом языке или на языке, близком к человеческому, было бы наиболее целесообразным. Таким образом, возникает необходимость создания транслятора с естественного языка на машинный, т. е. действующей модели человеческого языка. В сущности, не имеет значения, происходит ли общение человека с машиной в виде устной или письменной речи. Важно другое — умение машины владеть языком, понимать его смысловую и грамматическую структуру. В данной работе не затрагивается вопрос о том, насколько такой язык общения будет отличаться от человеческого (видимо, он станет приближаться к естественному языку по мере совершенствования).

Известно, что современные искусственные языки плохо приспособлены для описания творческих функций человека, поэтому при использовании машин для автоматизации интеллектуальной деятельности также следует обучать их естественному языку. Способность машины перерабатывать словесную информацию особенно необходима при решении таких важных прикладных задач, как машинный перевод, автоматическое реферирование и аннотирование текстов, обработка данных в системах материально-технического снабжения и др. Модель естественного языка необходима в деле разработки и внедрения автоматизированных систем управления.

Итак, практика выдвинула задачу создания действующей модели человеческого языка, которая по мере совершенствования будет все более ему соответствовать. Пока полная модель не создана, большое значение для теории и практики представляют частные математические модели, описывающие отдельные компоненты речевого общения.

Не вызывает сомнения, что «язык есть не нечто, постороннее по отношению к человеку, что можно изучать лишь как некоторый «памятник» эпохи, направления или художественного творчества отдельных людей, а часть самого человека в такой мере, в какой частью человека является его способность ходить на двух ногах в вертикальном положении, создавать орудия труда, мыслить понятиями и пр.» [1, с. 19]. Поэтому для моделирования языка (в частности, русского) требуется прежде всего рассмотрение его как «части самого человека», что является одной из первоочередных задач бионики.

Речевое поведение можно изучать с помощью методов биофизиологии, нейрохирургии и других наук, исследуя процессы, происходящие в мозгу человека при речевом общении. Подобные исследования важны, однако проще изучать и моделировать язык с помощью макроподхода — кибернетического метода «черного ящика». Информационно-функциональный характер кибернетического

моделирования позволяет без анализа структуры человеческого языка достичь хороших результатов при исследовании способности человека перерабатывать словесную информацию.

В данной работе на материале русского языка рассматриваются результаты моделирования словесного поведения человека, которое является неотъемлемой частью его речевого поведения. Под словесным поведением понимаем не обычное речевое общение, переработку человеком (как преобразователем информации, которому применимы понятия входных и выходных сигналов) отдельных словоформ. Одним из этапов моделирования словесного поведения является создание алгоритмов автоматического морфологического анализа, которые позволяют получать морфологическую информацию о каждой словоформе.

Психологические эксперименты показывают, что результаты решения каким-либо человеком даже простейших морфологических или субморфологических (связанных с частями слов, например [2]) задач могут несколько отличаться от результатов других людей, что объясняется прежде всего различием словарного запаса неодинаковым знанием грамматики испытуемыми. В принципе можно создать столько моделей какой-нибудь способности человека обрабатывать словесную информацию, а также моделей речи человека, сколько существует их прототипов — людей, говорящих на русском языке. Математическое описание речи отдельного индивидуума представляет известный интерес, однако особую практическую ценность имеет изучение и моделирование поведения идеально грамотного человека (точно руководствующегося правилами грамматики [3]), т. е. создание моделей языка.

Во многих работах, посвященных моделированию языка, в явном или неявном виде применяются морфологические анализ и синтез, которые обычно осуществляются при помощи словаря основ или словоформ, а также вручную (частично или полностью) переданием информации в машину. Множество входных слов велико и часто ограничивается какой-либо узкой областью, например математическими текстами. Все это объясняется не только трудностями формализации языка, но и специфической задачей машинного перевода, на решение которой в основном ориентируются исследователи. Общепринятый путь анализа, при котором сначала осуществляются словарный поиск и извлечение информации основы, иногда стараются объединить с подходом к языку как к цепи морфем [4], использующим максимум информации флекций (и некоторых других элементов) до обращения к словарю основ.

Предлагаемые нами алгоритмы отличаются прежде всего тем, что множество входных слов достаточно велико (за основу берется словарь русского языка [5] объемом в 104 000 слов с учетом всевозможных синтетических форм этих слов), а используемые словари относительно малы по сравнению с входным множеством. При получении искоемых грамматических признаков исследуется минимально

необходимое для этой цели количество букв слова (обычно по ним), независимо от его членения на морфемы.

В единичных случаях такой информации о словоформе окажется недостаточно для точного решения. Если же оно необходимо для анализа слов, неточно обрабатываемых алгоритмом, можно использовать предлагаемые ниже методы (отдельно или в комбинации друг с другом).

1. *Метод словарей.* В качестве формального признака используется все слово. При этом появится словарь исключений, объем которого должен быть намного меньше объема входного множества.

2. *Метод обращения.* В алгоритме (программе) должно быть предусмотрено обращение к человеку или к алгоритму (если таковой имеется) в целях получения признака, необходимого для точного решения.

3. *Метод дополнения.* К входному слову приписываются одно или несколько дополнительных признаков, позволяющих решить задачу точно.

4. *Метод ограничения.* Входное множество слов ограничивается таким образом, чтобы модель функционировала безошибочно.

Хотя метод 1 является самым трудоемким, его использование наиболее целесообразно, если объем словаря исключений невелик. В противном случае наиболее приемлемым окажется метод 2. Методы 3, снижающий автоматизированность модели, и метод 4 желательнее не применять, так как они несколько уменьшают ценность алгоритмов.

В предлагаемых здесь моделях получение высокой вероятности достоверных предсказаний основано прежде всего на применении метода 1.

При построении моделей используем принцип, заключающийся в том, что словоформа, изолированная от контекста, обладает всеми грамматическими значениями некоторой грамматической категории, которыми эта словоформа может характеризоваться в различных контекстах. Это означает, что допускаются неопределенные ответы машины, которые считаются точными, если соответствуют предложенному выше принципу. Например, для глагольной формы *вели* ответ «изъявительное наклонение и повелительное наклонение» будет верным, так как *вели* может использоваться в контекстах в грамматических значениях как изъявительного, так и повелительного наклонения. Подавляющее большинство словоформ характеризуется одним грамматическим значением любой грамматической категории за исключением слов-омографов, для которых выходной сигнал будет состоять из конъюнктивных признаков. Предложенный принцип подтверждается многочисленными психологическими экспериментами.

Плодотворной является идея об использовании взаимосвязи и взаимозависимости между различными грамматическими категориями и значениями. Исследование таких зависимостей полезно

например, при составлении алгоритмов анализа многих грамматических категорий для минимизации таких алгоритмов.

В целях использования глубинных структур языка мы неограниченно расширяем входное множество всевозможными словами, псевдословами и их формами, образованными на основании правил грамматики русского языка. При этом подразумевается, что существует распознающая процедура, способная отличать слова и псевдослова от произвольных цепочек букв и позволяющая таким образом формировать входное множество. Алгоритмы, оперирующие словами из неограниченного множества, не содержат исключений, довольно просты и дают достоверные предсказания с высокой вероятностью. Для получения точного решения при неограниченном множестве входных слов можно использовать методы 2—4.

Полученные алгоритмы позволяют сделать заключение о больших возможностях формального описания морфологии русского языка.

Математические модели способности идеального грамотного человека решать задачи морфологической классификации получены в виде алгоритмов, реализованных на ЭЦВМ. Решена задача морфологической классификации глаголов русского языка на ряды (инфинитив, личные формы, причастия, деепричастия), т. е. анализа категории репрезентации. Получены модели определения склонения, времени, лица, числа, рода, спряжения, анализа суженной парадигмы, повелительного склонения (собственно глаголов), проанализированы атрибутивные формы (причастия и деепричастия). Составлены алгоритмы определения глагольных признаков (времени и залога) причастий, а также разделения причастий на полные и краткие. Рассматриваются задачи об определении признаков вида и залога глагольных форм, задачи синтеза глаголов и субморфологической классификации.

Приведем записанные на алгоритмическом языке АЛГОЛ (для транслятора АЛГОЛ—ЦЭМИ на ЭЦВМ «Урал-14») отлаженные программы алгоритмов анализа категорий времени и числа глаголов русского языка в личной форме. Входным сигналом этих алгоритмов может быть произвольный глагол или псевдоглагол русского языка; входной алфавит — русский, расширенный знаком «—», который может применяться в глаголах повелительного склонения (на *-ка*) и для выражения редупликации (или удвоения слов). Выходными сигналами для первого алгоритма служат признаки времени: не прошедшее и прошедшее, а также сигнал отсутствия времени (для глаголов повелительного склонения). Множество выходных сигналов второго алгоритма объединяет два грамматических значения категории числа — единственное и множественное.

Количество входных слов  $N$  для каждого просчета на машине может быть любым, для изменения этого количества в программе достаточно заменить одну перфокарту. В данном случае  $N = 100$  (для определения вероятностей достоверного предсказания

моделей выбрано  $N=1000$ ). Максимальная длина слова в предлагаемых алгоритмах равна 80 символам.

Эталонный язык в трансляторе АЛГОЛ — ЦЭМИ расширен введением текстовых величин и действий над ними. Значения текстовых величин являются последовательности литер (пустой строка, буква, цифра, знак операции и т. п.). Мы пользовались операциями соединения текстовых значений (ТЗ) и выделения частей значения текстового выражения, осуществляемого с помощью выделителя (**from** < индексное выражение > **thru** < индексное выражение >), следующего за текстовым выражением, которое заключено в круглые скобки. Операция соединения ТЗ представляется знаком | (вертикальная черта); результат ее состоит из литер первого ТЗ, за которыми следуют литеры второго ТЗ. Действие выделителя состоит в выделении подпоследовательности литер значения предшествующего ему текстового выражения, начинающейся с литеры, номер которой равен значению первого индексного выражения, и кончающейся литерой, номер которой равен значению второго индексного выражения включительно.

Программа определения времени глагола имеет вид

**Begin integer N;**

**N := 100;**

**begin text array A[1:N], B[1:80];**

**integer i, j, L;**

**text C;**

**for i := 1 step 1 until N do**

**A[i] := intext (ПК, 80);**

**for i := 1 step 1 until N do**

**begin for j := 1 step 1 until 80 do**

**begin B[j] := (A[i] from j thru j);**

**if B[j] = ( ) then**

**begin L := j - 1; go to M1; end;**

**end;**

**M1: if (B[L - 1] | B[L] = (ся) ∨ (B[L - 1] | B[L] = (сь)**

**then L := L - 2;**

**go to if (B[L] = (y) ∨ (B[L] = (ю) ∨**

**∨ (B[L] = (м) ∨ (B[L] = (т)**

**then HB else if (B[L] = (н) ∨**

**∨ (B[L - 2] | B[L - 1] | B[L] = (ка)**

**then H;**

**go to if (B[L] = (н) ∨ (B[L] = (о) then**

**(if B[L - 1] = (л) then ПВ else H);**

**go to if B[L] = (ь) then**

**(if B[L - 1] = (ш) then HB else H);**

**go to if B[L] = (е) then**

**(if (B[L - 2] | B[L - 1] = (ер) ∨**

**∨ (B[L - 2] | B[L - 1] = (ит)**

**then HB else H) else ПВ;**

```

ЛВ:      outtext (АЦПУ, А [i] | 'непрошедшее время');
         go to M2;
ЛВ:      outtext (АЦПУ, А [i] | 'прошедшее время');
         go to M2;
Н:       outtext (АЦПУ, А [i] | 'нет времени');
М2:      newline (АЦПУ);
end;
end;
end.

```

Программа определения категории числа глаголов  
получена в виде

```

Begin integer N;
N: = 100;
begin text array A [1: N], Б [1: 80];
integer i, j, L;
text C;
for i: = 1 step 1 until N do
  A [i]: = intext (ПК, 80);
for i: = 1 step 1 until N do
begin for j: = 1 step 1 until 80 do
begin Б [j]: = (А [i]) from j thru j;
if Б [j] = ( ) then
begin L: = j - 1; go to M1; end;
M1:  if Б [L - 2] | Б [L - 1] | Б [L] = ('ка') then L: = L - 3;
if (Б [L - 1] | Б [L] = ('ся') ∨
∨ (Б [L - 1] | Б [L] = ('сь')) then L: = L - 2;
C: = Б [L - 1] | Б [L];
if (C = ('ат')) ∨ (C = ('ят')) ∨ (C = ('ут')) ∨ (C = ('ют')) ∨
∨ (C = ('ем')) ∨ (C = ('им')) ∨ (C = ('е')) then go to M;
if Б [L] = ('и') then L: = L - 1 else go to E;
if Б [L] = ('л') then go to M;
E:   outtext (АЦПУ, А [i] | 'единственное число');
      go to M2;
M:   outtext (АЦПУ, А [i] | 'множественное число');
M2:  newline (АЦПУ);
end;
end;
end.

```

Структура программ определения лица, рода и некоторых других грамматических категорий аналогична.

Ставя своей целью автоматическое получение грамматической информации о слове, в ряде случаев можно автоматически получать частичную смысловую информацию. Это обусловлено тем, что, определяя несинтаксические грамматические значения слов (например, значения времени), «отражающие различные смысловые абстракции» [3, с. 303], мы тем самым частично формализуем смысл. Решение задач морфологического анализа и синтеза позволит

определить для каждого слова его парадигму, т. е. класс слов с одним и тем же лексическим значением.

Полученные алгоритмы используются при дальнейших исследованиях (анализ пар слов, фраз и т. д.), но представляют и самостоятельный интерес как математические описания психических функций человека, а также могут найти применение, например в различных системах автоматического анализа текстов.

Процессы переработки человеком словесной информации требуются изучать не только при решении морфологических задач. Грамматика русского языка (представляющая собой обобщенный опыт владения языком многими поколениями людей и являющаяся эффективным средством при обучении человека языку, закреплении языковых норм и т. д.) предназначена для человека, и ее правила в том виде, в каком они записаны, обычно не пригодны для обучения машины. Бывает, что правила грамматики оказываются недостаточными для формализации языка. В таких случаях приходится проводить психологические эксперименты и изучать языковые структуры для выявления закономерностей, еще не описанных в языке. Исследование речевого поведения человека позволяет глубже проникать в структуру естественного языка и создавать «машинные» правила и модели.

#### ЛИТЕРАТУРА

1. Звегинцев В. А. Теоретическая и прикладная лингвистика. М. «Просвещение», 1968. 335 с.
2. Соловьева Е. А. Математическое описание способности человека анализировать правильность переноса слов. — В сб.: Проблемы бионики. Вып. 8. Харьков, 1972, с. 61—67.
3. Грамматика современного русского литературного языка. Отв. ред. Н. Ю. Шведова, М., «Наука», 1970. 767 с.
4. Севбо И. П. Пивоварова Е. П. Об алгоритме независимого флективного анализа русского текста. — В кн.: Прикладная лингвистика и машинный перевод. [Сб. статей]. Киев, 1962, с. 66—71.
5. Орфографический словарь русского языка. Изд. 11-е. Под ред. С. Г. Бархударова и др. М., «Сов. энциклопедия», 1971. 520 с.