

УДК 655.25+801.4

АВТОМАТИЧЕСКОЕ ОПРЕДЕЛЕНИЕ МЕСТА ПЕРЕНОСА
В ПРЕДЛОЖЕНИИ

Л.С.Юдина, А.С.Нудельман

В связи с дальнейшим увеличением потока информации в последнее время значительно возрос интерес специалистов к автоматизации редакционно-издательской деятельности (см., например, [1,2]).

Предлагаемая работа также выполнена в плане проблемы совершенствования редакционно-издательской деятельности на базе современной вычислительной техники и является составной частью проекта ПАРИС [3].

Проект ПАРИС (Полный Автоматизированный Редакционно-Издательский Сервис) представляет собой попытку разработать общие принципы и конкретные пути использования ЭВМ для ускорения и повышения качества редакционно-издательской подготовки рукописи к печати ротапринтным, фотонаборным и другими способами [3].

В данной статье предлагается алгоритм автоматического принятия решения о переносе на другую строку части слова в соответствии с правилами русской орфографии. Потребность в этом возникает в процессе печати при автоматизированном выравнивании строки.

Из известных авторам работ, предшествующих данной, следует назвать [1]. Отличие их - прежде всего в назначении: алгоритм, описанный в [1], разработан для систем программирования набора, данный алгоритм ориентирован на ротапринтный способ издания. Кроме того, алгоритм [1] определяет места переноса только внутри слова, предлагаемый в данной работе - как в слове, так и в предложении в целом. Можно назвать также некоторые частности, которые в конечном счете отражаются на результатах работы: по-разному организован массив приставок и метод проверки на него; различаются наборы формализуемых правил, составляющих основу алгоритма.

Основное содержание работы составляет обсуждение общих принципов переноса и возможностей их алгоритмизации, а также описание самого алгоритма и экспериментов по автоматизации переносов в реальном тексте.

I. Общие принципы переносов

Задача автоматического определения места переноса должна включать два случая:

- а) перенос внутри слова;
- б) перенос между словами.

Авторами разработаны алгоритм и программа, осуществляющие автоматическое определение места переноса как внутри слова, так и между словами. Алгоритм реализован на ЭВМ "Минск-32". Следует отметить, что при написании программы преследовалась основная цель – добиться высокого качества переносов. По времени работы программа не оптимизировалась.

I.I. Перенос внутри слова. Алгоритм автоматизации переносов слов основан на формализации существующих правил переносов [4,5]. Известно, что в основе действующих в настоящее время правил переносов лежат как фонетический, так и морфологический принципы [6]. Перенос внутри слова осуществляется, главным образом, по слогам (фонетический принцип). Однако это правило ограничивается запрещением переносить слоги, если они состоят из одного звука (буквы), и разрывать морфемы (морфологический принцип).

История развития языка свидетельствует о том, что стихийно появившиеся переносы усложнялись постепенно: ранние писцы при необходимости переносили любую оставшуюся часть текста. Впоследствии стали переносить по слогам, еще позже – с учетом морфемного состава слов. Реформа 1917–1918 гг. свела правила переносов к переносам по слогам (с небольшими ограничениями). Однако окльязыки стихийно стали появляться и ломорфемные переносы. Считается, что существующие к настоящему времени правила переносов учитывают оба этих момента – слогоделение и морфемный состав слов [6]. Дело осложняется еще тем, что толкование самого слогоделения не единобразно.

В результате издательская практика содержит многочисленные отступления от правил в виде повсеместно встречающихся переносов

типа: ли-менование, ох-ранять, сог-лашения, ун-решение, на-бо-
лее, сп-равдана, энерго-тэя. Подтверждаясь практикой печати
нестрогость действующих правил переносов усложняет их формали-
зацию.

Следует отметить, что авторы старались максимально придерживаться существующих правил переносов; пересмотр их в задачи работы ни в коей мере не входил.

К настоящему времени считается установленным, что наименьши-
ми произносительными единицами, на которые распадается речь, явля-
ются слоги, и слова следует переносить по слогам. Известно также,
что наиболее распространенный, "классический" тип русского слога —
сочетание согласного (С) с гласным (Г) — (СГ).

Возьмем эти предпосылки за основу алгоритмизации. Формально,
без учета морфемной структуры, слово можно рассматривать как це-
почку из элементов С и Г в различных (но фиксированных) комбина-
циях.

В задаче формализации переносов, очевидно, было бы оптималь-
ным оперировать трех- и четырехбуквенными сочетаниями, однако про-
блема сочетаемости в русском языке изучена еще недостаточно, и мы
вынуждены ограничиться билаграммами с привлечением по мере надобно-
сти информации о предшествующих и последующих буквах.

На уровне билагрэмм сочетания внутри слова представлены типа —
ми: ГГ, СГ, ГС, СС, _Г, _С, Г_, С_, где символ _ обозначает про-
бел — начало или конец слова. Относительно сочетаний, содержащих
его, принимается решение на основании правила: одна буква не остав-
ляется на строке и не переносится, поскольку она ассоциируется
в первую очередь не со словом, а с соплем, предлогом, междометием.

В числе оставшихся сочетаний группа ГГ самая малочисленная;
сочетания ГГ для русского языка не характерны, они встречаются,
главным образом, в заимствованных словах. В сочетаниях гласных ал-
горитмом предусмотрен перенос как между ними, так и после послед-
ней гласной, за исключением двух случаев, когда после последней
гласной идет сочетание: а) согласный илис (ъ, ъ); два согласных,
но не -ск-, -ст-, -гр-.

В сочетаниях СГ согласный всегда отходит к последующему гла-
сному, и перенос между ними запрещен. Если при этом иметь в виду,
что более половины общего числа русских слогов^{*)} состоят из слоги

^{*)} Данные получены для русской речи, но с известными допущениями
их можно распространить на печатный текст.

така СГ [7], то, очевидно, можно заранее сократить свыше 50% правильных переносов.

Сочетания ИС делятся переносом, кроме случаев, когда после согласного идет согласный или (ъ, ѿ, и).

В сочетаниях СС - тоже немногочисленной группы сочетаний в русском языке - перенос предусмотрен между ними. Следовательно, в сочетаниях из трех и более согласных перенос может быть поставлен после каждой из них, кроме последней.

По-видимому, четкую грань между фонетическим и морфологическим принципами провести нельзя (это, собственно, отражено в самой истории вопроса и в формулировках правил переносов). При разработке алгоритма также шли не место случая, требующие обоих подходов одновременно. Так, в алгоритме не делается переносом сочетания согласных -ст-, -ск-, -гр- - наиболее частые двухбуквенные сочетания согласных в русском языке [7].

Есть основания полагать, что они, как правило, входят в состав одних морфем; пример тому - -гр-, -ск-, -ст- в наиболее продуктивных морфемах (во вторых основах сложных слов): -граф, -грам, -скол, -стат, а также -ск-, -ст- в качестве одних из самых продуктивных словообразовательных суффиксов.

Кроме описанного, в алгоритме учтены также следующие правила:

- а) не оставляется на строке и не переносится часть слова, состоящая из одних согласных;
- б) буквы (ъ, ѿ, и) не отделяются от предшествующей буквы;
- в) две одинаковые согласные делются переносом, за исключением двойных согласных корня, идущих после приставки;
- г) часть слова, начинавшаяся с (и), не переносится.

Наибольшие трудности вызывает, естественно, реализация морфологического принципа. Конкретно это выражало необходимость учитывать морфемную структуру слова, а именно:

- не разбивать переносом односложную приставку;
- не отделять от корня начальную часть, не составляющую слова;
- в сложных словах не отделять переносом начальную часть второй основы, не составляющую слова;
- не разбивать переносом односложную часть сложносокращенного слова.

Это и определяют порядок работы алгоритма: в память машинных изведок "словарь приставок" - наиболее употребительные пристав-

ки и наиболее частые элементы сложных слов [9]. Назовем их все квазиприставками (перечень их приводится в приложении I). Введение в память машины этого более или менее обширного "списка приставок" соответственно снижает процент ошибок, т.е. повышает качество машинных переносов.

Сначала слово анализируется на наличие квазиприставки, после которой ставится перенос; затем описанным выше способом происходит деление на слоги как квазиприставки, если в ней более одного слога, так и оставшейся части слова.

В словах с приставками, не помещенными в память машины, перенос осуществляется по описанному выше фонетическим правилам.

Требования морфологии при этом могут быть либо учтены (автоматически), либо не выполнены.

К случаям, не поддавшимся формализации по общим принципам и явившимся собой примеры неправильных "машинных" переносов, относятся прежде всего следующие:

- сложные слова, поскольку все разнообразие частных, иногда единичных, случаев написания известны в память машины невозможно, например: пицер-кусовой, пенико-ниматель, соломот-рас^{*)}; в их числе следует особо отметить заимствование слова и термины, состоящие, как известно, в русской языке особую языковую подсистему [10], например, ай-с-берг, бран-д-муэр, фенол-ф-талени;

- одноковая реализация на письме неодинаковых по морфемному составу частей слова, для различения которых простых формальных признаков найти пока не удается: двухстка, но двухвершинный, необитаемый - неологизм; не поддаются такие формальному различию случаям, когда одно и то же буквосочетание в одних словах является приставкой, в других входит в состав корня: ретресс - резной; в последнем случае на наиболее частые буквосочетания, в ходящие в состав корня, производится специальная проверка с целью "отмены действия приставки", и слово переносится по общему принципу: сол-и-це (но со-адить).

1.2. Перенос между словами. Как известно [4,8], не любое слово и не все, что написано через пробел или де-

^{*)} Здесь указаны места только неправильных переносов.

фис, может быть разделено переносом. Наиболее распространеными примерами этого положения являются:

1. Аббревиатуры, написанные:

- а) прописными буквами (ИК ВИКСМ),
- б) частью прописными, частью строчными (КЗоТ),
- в) прописными с цифрами, дефисом (ТУ-104).

2. Наречия – грамматические окончания, соединенные дефисом с цифрами (2-ю).

3. Единицы измерения – сокращенные обозначения мер, пишущиеся вместе с цифрами, указывающими число измеряемых единиц (10 кг, 20 см).

4. Общепринятые графические сокращения (т.е., о-во, и т.п.).

5. Сокращенные слова в сочетаниях с именами собственными (г.Пушкин, пл.Свердлова).

6. Цифры или буквы со скобкой или точкой (при перечислении) в сочетании со следующими за ними словами (3. При наборе..., б) сокращенные выражения...).

7. Знаки и обозначения (§, №, %, С, °) в сочетаниях с цифрами.

8. Инициалы, фамилия и инициалы.

9. Пунктуационные знаки.

10. Открывающие скобки и отыравающие кавычки, а также знак сноски.

Подробнее формальные правила переносов между словами (через пробел) приведены при описании соответствующего алгоритма (п.2).

2. Алгоритм автоматического анализа текста

2.1. Исходные определения. Алфавит – полный набор различных символов, из которых составлен текст. В нашем случае это: строчные и прописные буквы русского алфавита (кроме "ё" и "Ё"); арабские цифры; знаки препинания (., : ! ?); скобки ([, (,),]); кавычки (""); знаки параграфа, номера, процента (§, №, %); наклонная черта (/); минус (-); пробел (—); пробел индексный (—); пробел формульный (Φ).

Формальный текст (или просто текст) есть последовательность символов алфавита, составленная в соответствии с общепринятыми правилами манипуляции. Все особенности оформления входного текста перечислены ниже:

1) Не допускается печать слов "разрядку" (разрядку можно обозначать кодом).

2) Наклонная черта не используется в качестве скобки.

3) Пробел эквивалентен по употреблению "пропуску символа" в машинописном тексте.

4) Пробел индексный ставится вместо любого символа, пропущенного над строкой (одущевленного под строку) на один интервал в машинописном тексте.

5) Пробел формульный резервирует место в строке, равное ширине символа, для вписывания формулы.

6) Минус, стоящий между пробелами, обозначает математический символ во вписанной формуле, на котором возможен перенос.

7) Минус не используется в качестве знака переноса.

Отсутствие в формальном тексте знака переноса говорит о том, что формальный текст всегда выполняется в виде одной строки.

ПРИМЕР. Машинописный текст:

"Среди вершин $\Gamma(x_1)$ множества $\Gamma(x_1)$ выбираем одну из вершин x_2 , для которой
 $\lambda'(x_1) - \lambda'(x_2) = \lambda(x_1, x_2) + (21 + 1),$
и обозначаем выбранную вершину через x_2 ".

в формальном варианте выглядит так:

"Среди вершин $\Gamma(x_1)$ множества $\Gamma(x_1)$ выбираем одну из вершин x_2 , для которой
 $\lambda(x_1, x_2) + (21 + 1),$
и обозначаем выбранную вершину через x_2 ".

Любая буква русского алфавита является либо гласной, либо согласной, либо особой (и, ю, ё).

Априори выбранную последовательность букв, являющуюся начлом некоторого слова, будем называть квазиприставкой. Каждая квазиприставка либо является приставкой, либо не является таковой. Например, из выбранных в настоящей работе квазиприставок - "под" и "за" - приставки, а "жизн" и "авка" - не приставки.

Каждой квазиприставке Q априори ставится в соответствие множество $Z(Q)$ запрещающих (выделение этой квазиприставки) буквосочетаний (это множество может быть и пустым; все пустые множества запрещающих буквосочетаний описаны в приложении I). В настоящей работе выбраны, например, $Z(\text{под}) = \emptyset$, $Z(\text{за}) = \{\text{ч}, \text{ри}\}$, $Z(\text{жизн}) =$

$= \{ \text{ни} \}$, $Z(\text{ни}) = \emptyset$. Отсюда следует, что последовательность букв "за" в слове "забег" является квазиприставкой, а в слове "зарыла" не является таковой.

В работе тоже алгоритм фиксирует следующие множества:

а) множество нерасчлененных буквосочетаний $U = \{ \text{гр}, \text{ск}, \text{от} \}$;

б) множество сокращений первого рода

$S1 = \{ \text{г.}, \text{гор.}, \text{гр.}, \text{ни.}, \text{о.}, \text{пос.}, \text{с.}, \text{т.}, \text{тт.}, \text{тоб.} \}$;

в) множество сокращений второго рода $S2 = \{ \text{др.}, \text{пр.}, \text{т.д.}, \text{т.п.} \}$.

В дальнейшем символы формального текста будем обозначать через индексированные a_1, a_2, \dots, a_n . Предполагается, что символы любого текста пронумерованы в естественном порядке (номер первого символа больше или равен 1). Условимся номер символа в тексте отождествлять с индексом. Таким образом, любой текст T длины n можно представить как последовательность $T = a_1 a_{i+1} \dots a_{i+n-1}$.

Текст, составленный только из букв, будем называть словом.

Пусть $T = a_1 a_2 \dots a_n$ — текст. Символ a_i будем называть переносимым в данном тексте, если часть текста $a_1 a_{i+1} \dots a_n$ можно перенести на следующую строку. Если такого переноса сделать нельзя, то a_i считается непереносимым.

2.2. Алгоритм выделения квазиприставки.

Входная информация:

а) $w = a_1 a_2 \dots a_n$ — слово длины n ($n > 0$);

б) Q_1, Q_2, \dots, Q_m — перечень всех квазиприставок такой, что $q_{i+1} \leq q_i$, где q_i — длина квазиприставки Q_i ;

в) $Z(Q_i)$, $1 \leq i \leq m$, — множество запрещенных буквосочетаний.

Выходная информация: целое число q ($q \geq 0$), при этом $q > 0$ означает, что последовательность $a_1 a_2 \dots a_q$ является квазиприставкой максимальной длины в слове w ; $q = 0$ означает, что слово w не содержит квазиприставки.

Общая идея алгоритма: перебором в множестве квазиприставок ведется поиск квазиприставки (максимальной длины), совпадающей с началом слова w и такой, что в слове w после этой квазиприставки нет запрещающего ее буквосочетания.

Работа алгоритма состоит в следующем:

Шаг I. $i := 1$.

Шаг 2. Если $i > m$, то $q := 0$ и шаг II.

Шаг 3. Если $q_1 > n$, то шаг 10.

Шаг 4. Пусть $Q_{q_1} = b_1 b_2 \dots b_k$, тогда если $\forall j (1 \leq j \leq k \rightarrow b_j = a_j)$, то шаг 5, иначе шаг 10.

Шаг 5. Если в $Z(Q_{q_1})$ имеется буквосочетание $c_1 c_2 \dots c_p$ такое, что $\forall j (1 \leq j \leq p \rightarrow c_j = a_{k+j})$, то шаг 10.

Шаг 6. Если $q_1 > n - 3$, то $q := 0$ и шаг II.

Шаг 7. Если среди букв $a_{q_1+1}, a_{q_1+2}, \dots, a_n$ нет ни одной гласной, то $q := 0$ и шаг II.

Шаг 8. Если a_{q_1+1} — особая буква, то $q := q_1 + 1$ и шаг II.

Шаг 9. Если Q_{q_1} — приставка, a_{q_1} — согласная и a_{q_1+1} — гласная буква, то $q := q_1 - 1$ и шаг II, иначе $q := q_1$ и шаг II.

Шаг 10. $i := i + 1$ и шаг 2.

Шаг II. Конец.

2.3. Алгоритм анализа букв и.

Входная информация:

а) $W = a_1 a_2 \dots a_n$ — слово длины n ($n \geq 4$);

б) U — множество нерасчленяемых буквосочетаний;

в) m — номер символа в W , $3 \leq m \leq n$.

Выходная информация: целое число q ($q = 0$ или I), при этом $q = 0$ означает, что буква a_m непереносима в слове W ; $q = I$ означает, что a_m переносима в W .

Общая идея алгоритма: решение о переносимости или непереносимости буквы a_m в слове W принимается на основе анализа части слова $a_{m-2} a_{m-1} a_m a_{m+1}$ в терминах гласная, согласная и особая (буквы), при этом учитывается расположение в этой части слова нерасчленяемого буквосочетания (если такое есть).

Работа алгоритма заключается в следующем:

Шаг 1. Если a_m — особая буква, то $q := 0$ и шаг 8.

Шаг 2. Если среди a_1, a_2, \dots, a_{m-1} нет гласных или среди a_m, a_{m+1}, \dots, a_n нет гласных, то $q := 0$ и шаг 8.

Шаг 3. Если a_m — гласная, тогда, если a_{m-1} — согласная, то $q := 0$ и шаг 8, иначе $q := I$ и шаг 8.

Шаг 4. Если $a_{m+1} = a_m$, то $q := 0$ и шаг 8.

Шаг 5. Если a_{m-1} — особая буква, то $q := I$ и шаг 8.

Шаг 6. Если a_{m-1} — гласная, тогда, если a_{m+1} — особая или a_{m+1} — согласная и буквосочетание $a_m a_{m+1}$ не содержится в U , то $q := 0$ и шаг 8, иначе $q := I$ и шаг 8.

Наг 7. Если $a_{n-2} = a_{n-1}$, или буквосочетание $a_{n-1}a_n$ содержится в U , то $q := 0$ и наг 8, иначе $q := 1$ и наг 8.

Наг 8. Конец.

2.4. Алгоритм анализа слова.

Входная информация:

- $T = a_1 a_2 \dots a_n$ – текст длины n ($n \geq 4$);
- $W = a_k a_{k+1} \dots a_{k+n-1}$ – слово длины m в тексте T ($m \geq 4$);
- q_0 – номер символа в тексте T ;
- г) АВК (W' , q) – алгоритм выделения квазиприставки (W' – входное слово, q – выходное число);

- д) АБ (W' , i , q) – алгоритм анализа буквы (W' – входное слово, i – номер символа в W' , q – выходное число).

Выходная информация: целые числа q^- и q^+ ($q^- \geq 0$, $q^+ \leq n+1$), при этом $q^- = 0$ ($q^+ = n+1$) означает, что не существует номера i такого, что $k < i \leq k+m-1$, $i \leq q_0$ ($i \geq q_0$), и a_i переносим в слове W (и, следовательно, в тексте T); $q^- > 0$ и $q^- \neq q_0$ ($q^- \leq n+1$ и $q^+ \neq q_0$) означает, что символ a_{q^-} (a_{q^+}) из слова $a_{k+1} a_{k+2} \dots a_{k+m-1}$ является ближайшим слова (справа) к a_{q_0} символом, переносимым в W (в T).

Если $k < q_0 \leq k+m-1$ и символ a_{q_0} из слова W переносим в W , то $q^- = q^+ = q_0$.

Основная идея алгоритма: слово W рассматривается как состоящее из квазиприставок и остатка, т.е. $W = K_1 K_2 \dots K_1 R$, где K_1 – квазиприставка в W , K_2 – квазиприставка в W без K_1 , и т.д., а R – часть слова W , в которой отсутствует квазиприставка; буква, стоящая после неоднобуквенной квазиприставки, переносима в W , вторая и последняя буквы в K_1 (и R) запереносимы в W , переносимость в W всех остальных букв слова W совпадает с переносимостью этих букв в соответствующей квазиприставке K_1 (остатке R).

Будем обозначать через МИ множество номеров таких, что если $i \in \text{МИ}$, то символ a_i переносим в W (в T).

Работа алгоритма состоит в следующем:

Наг 1. $t := k$ и МИ := \emptyset .

Наг 2. Слово $W_1 := b_1 b_2 \dots b_r$ ($b_i = a_{t+i-1}$, $r = n+k-t$).

Наг 3. АВК (W_1 , p).

Наг 4. Если $p = 0$, то наг 9.

Наг 5. Слово $W_2 := b_1 b_2 \dots b_p$.

Наг 6. Для всех j , $3 \leq j < p$, АБ (W_2, j, h), и если $h = 1$, то МИ := МИ $\cup \{ t + j - 1 \}$.

Шаг 7. Если $p > 1$, то $MII := MII \cup \{t + p\}$.

Шаг 8. $t := t + p$ и шаг 2.

Шаг 9. Для всех j , $3 \leq j \leq r$, ААБ (w_i, j, h) , и если $h = I$, то $MII := MII \cup \{t + j - 1\}$.

Шаг 10. Из анализа множества MII устанавливаются значения q^- и q^+ .

Шаг II. Конец.

2.5. Алгоритм анализа текста.

Входная информация:

а) $T = a_1 a_2 \dots a_n$ – текст длины n ($n \geq 6$ и символы a_1 и a_n – пробелы);

б) q_0 – номер символа в тексте T ;

в) ААС (w^*, q) – алгоритм анализа слова (w^* – входное слово, q – входной номер).

Выходная информация: целые числа q^- и q^+ ($q^- \geq 0$, $q^+ \leq n+1$), при этом $q^- = 0$ ($q^+ = n+1$) означает, что не существует номера i такого, что $4 < i < n$, $i \leq q_0$ ($i \geq q_0$), и a_i переносим в тексте T ; $q^- > 0$ и $q^- \neq q_0$ ($q^+ < n+1$ и $q^+ \neq q_0$) означает, что символ $a_{q^-}(a_{q^+})$ из текста T является ближайшим слева (справа) к a_{q_0} символом, переносимым в T .

Если $4 < q_0 < n$ и символ a_{q_0} переносим в T , то $q^- = q^+ = q_0$.

Общая идея алгоритма: могут быть переносимы в T только:

а) внутренняя буква некоторого слова из T , не являющегося единицей измерения при числе или аббревиатурой,

б) символ, стоящий после минуса,

в) символ, стоящий после пробела.

В случае "а" решение о переносимости принимает алгоритм анализа слова; случаи "б" и "в" рассматриваются данным алгоритмом непосредственно, при этом в обоих случаях анализируется контекст, содержащий исследуемый символ.

Через MII будем обозначать множество номеров таких, что если $i \in MII$, то символ a_i возможен для переноса.

Слово $w = a_1 a_{i+1} \dots a_{i+m-1}$ из текста T будем называть максимальным, если в T a_{i-1} не буква и a_{i+m} не буква.

Работа алгоритма состоит в следующем:

Шаг 1. $MII := \emptyset$, $i := 4$.

Шаг 2. Если $i = n-1$, то шаг 28.

Шаг 3. Если a_i – буква, то шаг 8.

Шаг 4. Если a_{i_1} - минус, то шаг 10.

Шаг 5. Если a_{i_1} - пробел, то шаг 14.

Шаг 6. $i:=i+1$ и шаг 2.

Шаг 7. $M:=M \cup \{ i+1 \}$ и шаг 6.

Шаг 8. Если $a_{i_{i+1}}$ - буква, то шаг 6.

Шаг 9. Пусть $w = a_{i_1}a_{i_1+1} \dots a_{i_1+n-1}$ - максимальное слово, тогда если $n \geq 4$, количество гласных в w больше единицы и среди $a_{i_1+2}, a_{i_1+3}, \dots, a_{i_1+n-1}$ нет ни одной прописной буквы, то AAC (w, q_0), иначе шаг 6.

Шаг 10. Если a_{i_1-1} и a_{i_1+1} - цифры, либо a_{i_1-1} и a_{i_1+1} - пробелы формульные, то шаг 7.

Шаг 11. Если a_{i_1-1} не буква или a_{i_1+1} не буква, то шаг 6.

Шаг 12. Пусть w_1 - ближайшее слева от a_{i_1} максимальное слово, а w_2 - ближайшее справа от a_{i_1} максимальное слово. Если w_1 или w_2 содержит только одну букву, или w_1 или w_2 не имеют гласных, или первая буква w_2 прописная, то шаг 6.

Шаг 13. Если w_1 содержит прошисную не первую букву и w_2 содержит прошисную не первую букву, то шаг 6, иначе шаг 7.

Шаг 14. Если a_{i_1+1} - пробел или %, или a_{i_1-1} - пробел или \$, или %, то шаг 6.

Шаг 15. Если a_{i_1+1} - минус, тогда, если a_{i_1+2} - цифра, то шаг 7; если a_{i_1+2} - пробел, то шаг 6.

Шаг 16. Если a_{i_1-1} - скобка и a_{i_1-3} - пробел, то шаг 6.

Шаг 17. Если a_{i_1-1} не цифра и не пробел формульный, то шаг 22.

Шаг 18. Если a_{i_1+1} не пробел формульный, то шаг 20.

Шаг 19. Пусть $w' = a_{i_1+1}a_{i_1+2} \dots a_{i_1+n-1}$ - последовательность формульных пробелов такая, что a_{i_1+n} не формульный пробел. Тогда, если a_{i_1+n} - минус, то шаг 7, иначе, если $n > 3$, то шаг 7, иначе шаг 6.

Шаг 20. Если a_{i_1+1} не буква, то шаг 7.

Шаг 21. Пусть w_2 - ближайшее справа от a_{i_1} максимальное слово, тогда если в w_2 не более одной гласной, то шаг 6, иначе шаг 7.

Шаг 22. Если в тексте Т слева от a_{i_1} нет ни одной буквы, то шаг 7.

Шаг 23. Если a_{i_1+1} не буква, то шаг 7.

Шаг 24. Если существует такое n , что последовательность $a_{i_1-n}a_{i_1-n+1} \dots a_{i_1-1}$ есть сокращение первого рода и a_{i_1+1} - прописная буква, то шаг 6.

Шаг 25. Если существует такое i^+ , что последовательность $a_{i+1} a_{i+2} \dots a_{i+n}$ есть сокращение второго рода, то шаг 6.

Шаг 26. Если a_{i-1} не буква, то шаг 7.

Шаг 27. Пусть w_1 - ближайшее слева от a_1 максимальное слово, а w_2 - ближайшее справа от a_1 максимальное слово. Если w_1 и w_2 начинаются с прописных букв и w_2 состоит только из одной буквы, то шаг 6, иначе шаг 13.

Шаг 28. Из анализа множества МП устанавливаются значения c^- и q^+ .

Шаг 29. Конец.

3. Эксперименты и их результаты

Алгоритм анализа текста был реализован в виде программы ЮСТРУ, оформленной на языке ФОРТРАН для ЭВМ "Минск-32". Эта программа предназначена, в частности, для использования в надпрограмме, формирующей строки фиксированной длины, и определяет в тексте места возможных переносов.

С целью проверки качества переносов, реализуемых программой ЮСТРУ, было проведено два эксперимента.

3.1. Эксперимент "слово". Цель эксперимента состояла в выяснении точности определяемых программой ЮСТРУ переносов внутри слов. Экспериментальный текст (3000 символов) был составлен из слов русского языка, выбранных случайным образом из словаря [11]. Экспериментальный текст представлен в приложении 2.

3.2. Эксперимент "текст". Цель эксперимента состояла в выяснении корректности определяемых программой ЮСТРУ переносов внутри произвольного (формального) текста. Экспериментальный текст (3000 символов) был взят из [12] со случайно выбранной страницы. Экспериментальный текст представлен в приложении 2.

3.3. Распечатка экспериментальных текстов. Из-за ограниченности набора символов УПЧ ЭВМ "Минск-32" при распечатке (формального) текста была использована следующая символика:

1. Прописная буква русского алфавита обозначает соответствующую строчную букву русского алфавита.

2. Прописная буква латинского алфавита обозначает аналогичную букву русского алфавита.

3. Вписанная буква обозначает прописную букву.

4. "Стрелка" - вопросительный знак.

5. "Равно" - пробел индексный.

6. "Больше" - прописную букву Б.

7. "Звездочка" - символ номера.

3.4. Оценка результатов экспериментов. При оценке качества "машинных" переносов производилось сравнение переносов, разрешенных программой ИСТР®, с переносами, допустимыми с точки зрения человека. Количественно такое сравнение выражается двумя числами: коэффициентом полноты и коэффициентом ошибки.

Пусть в тексте Т человек-специалист выделяет точно n переносимых (в Т) символов; пусть в этом же тексте программа (в частности, ИСТР®) относит к переносимым m_0 символов, причём m_0 из них являются переносимыми и с точки зрения человека (понятно, что $m_0 \leq n$ и $m_0 \leq n$). Тогда коэффициент полноты

$$P = \frac{m_0}{n},$$

а коэффициент ошибки

$$R = \frac{n - m_0}{n}.$$

Следует отметить, что коэффициенты полноты и ошибки характеризуют различные стороны автоматического переноса и поэтому являются достаточно независимыми. Большой коэффициент полноты указывает на большее использование "ресурсов" переноса, имеющихся в русской грамматике. При меньшем коэффициенте ошибки повышается грамотность разбивки текста на строки.

Результат эксперимента "слово": $P_1 = 0.963$, $R_1 = 0.053$ (при $n_1 = 750$).

Результат эксперимента "текст": $P_2 = 0.955$, $R_2 = 0.005$ (при $n_2 = 940$).

3.5. Выводы. Полнота и ошибочность переносов, выявляемых программой ИСТР®, находится, по-видимому (см. I.I), в пределах требований существующих правил переносов. Для более надежной оценки качества машинных переносов необходимы экспериментальные тексты различного характера (технические, общественно-политические, художественные) объемом порядка 10^6 символов.

Анализ автоматических переносов, полученных в экспериментах, позволяет предположить, что их качество в основе своей удовлетворяет требованиям практики печати. В то же время преимущество автоматического определения места переноса в слове в свете проблемы программирования набора несомненно.

Авторы выражают благодарность В.Г.Косареву за постановку задачи.

Л и т е р а т у р а

1. БЕРЛИН А.С. Системы программирования набора. М., "Книга", 1971.
2. ВОРОНОВА Т.П., ВОСКРЕСЕНСКИЙ М.И., КОСТИН Б.А. Автоматизированная система набора, правки и верстки. -"Полиграфия", 1973, № 2.
3. КОСАРЕВ В.Г. Проект ПАРИС - полный автоматизированный редакционно-издательский сервис. Отчет ИМ СО АН СССР, Новосибирск, 1974.
4. БЫЛИНСКИЙ К.И., НИКОЛЬСКИЙ Н.Н. Справочник по орфографии и пунктуации для работников печати. М., "Просвещение", 1970.
5. Правила русской орфографии и пунктуации. М., изд-во "Искусство", 1962.
6. ИВАНОВА В.Ф. Современный русский язык. Графика и орфография. М., "Учпедгиз", 1966.
7. ЮНИНА Л.С. Автоматизация транскрипции и некоторые количественные характеристики русской речи. Автореф. дис. на соиск.учен. степени канд. филологич. наук. Л., 1974, 23 с. (Ленингр. гос. ун-т).
8. РОЗЕНТАЛЬ Д.Э. Справочник по правописанию и литературной правке. М., "Книга", 1971.
9. БУКЧИНА Б.З., КАЛАУШКАЯ Л.П. Сложные слова, М., "Наука", 1964.
10. ПАНК М.В. Русская фонетика. М., "Просвещение", 1967.
11. Орфографический словарь русского языка. М., изд-во "Советская энциклопедия", 1970.
12. Проблемы структурной лингвистики. М., изд-во АН СССР, 1963.

Поступила в ред.-изд.отд.
3 августа 1977 года

ПРИЛОЖЕНИЕ I

**Вспомогательные массивы, используемые для
автоматизации переносов**

19.11.1975

III
001 011 КРИСТАЛЛО, ФОРМАЛЬНО,)
001 012 РЕНТГЕНБС;
001 013 ВОЗДУХО, ГРУБЬКО, ЧАРГИТО-МЕТАЛО,)
001 014 ОВРАТНО, ОСЦИЛЛО, ПРЕГАНЕ, ПРОТИВ,)
001 015 СЕЛЬСКО, СЧЕГАНО, УСЛОВНО, ЧЕТЫРЕХ,)
001 016 ЭЛЕКТРО,;)
001 017 БЫСТРО-АДГОНО-ВЗДЫМН-ВНУТРН-ВЫС-КО-
ГЕРБР, АДМИНО-АВЕРН-НЕЛЕЗД-ЗЛЯТА-
КАРАНО-КЛЕТВ-КОНТВ, КРАСНО-КРЕЧЕ-
КРУГЛНО, КРУПНО, МАШИНО, НЕВНО, ОРГАНО,
ПЛОСКО-ПЛЕВДО-ПРОЧНО-РЕЛЬСО-СИЛЬНО-
СЛОЖНО-СВЕАН-СТЕХДО-СТЕРЕО-ГОЛГО-
ТРУДНО-ТОЧЕЛО-УЛЬТРН-ЧАСТИНО-ЩИРКО-
ЭКСТРН, ЭНЕРГО,)
001 018 БЛАГО, ВИБРО, ВЛАГО, РЕЛМО, ГЕКТО, ГУРРО,
ГИАРО-ГИГЕР, ГОРНОГРУСТО-ДИНОН-УРЗН-
ЗААНЕ, ЗЛУКО, ЗЕРНО, ЗЛАТО, КАЗНО, КАННЕ,
КВАЗИ, КЛЮНО, КОНТР, КОСНО, КРОВО, КРУТО,
ЛЕРКО, МАКРО, МЕЛКО, МИКРО, МИЛГО, НИФЕ,
НИЗКО-ОШГО; ПЕРВО-ПЛОДО-ПРЯМО-ПУСТО-
РАВНО-РАДНО-РЭЗКО-СВЕРЕ, СВЕТО, СУЕРХ,
СЛАБО, СЛОВО, СТАРД, СТЕНО, СТЕТО, СТИХО,
ТРЛДО, ТЕЧНО, ТОНКО, ТРАНС, ТРУДО,
ТУРБО-ХЛЕБО, ХРОНО-ХРОНО-ЦИКЛО-ЧИРНО,
ЧИСТО, ЧИФКО, ВЕНЗО, ГЕКСО, НЕТРО, ПРАВО,
АВИА, АВТОГРБ, АНТИ-АРХИ; АВРО; ВСИО,
ВРДО, ВЫЧЕ, ГАЗО, ГИДР, ГИРО, ГОМО-ДВУХ,
ДИКО, ЖИЯЮ, ЗУБО, ЕРНО-ЛЕГО, ЛЕСО, ЛИТБ-
ЛЬНО, МАЛО, МЕТА, МОНО, МОТО; МИСО,
МИРЕ-НОДО, НССО, ОВЧЕ-ОДНО, ОРТО-ОФФО,
ПЕРЕ-ПИНО, ПОЛУ-ПРЕА-ПУТЕ, РУКО, САНО,
ТЕЛЕ, ТРЕХ, УГЛЕ, УЗКО, ФИТО, ФОНД, ГЕСТ,
ЦИТО, ЧУЧЕ, ЭНАД,)
001 034 БИО, ВНЕ, ВСЕ, ГЕО, ДВУ, ЕХЕ, ЗЛД; ЗОО; ИЗО, ИНО,
ИНЕ-НАА-ОТО; ПОД. ПРЕ-ПРИ-ПРО-ЭКС. ЭПН,)
001 035 ВЫДБ; АН-ЗА-НА, НЕ-ОБ, ОТ. ПО-РЕ; ГР-
001 036 0,++++++
001 037 +
001 038 +
001 039 РЕНТГЕНБС;
001 040 ЭЛЕКТРО-ИН;
001 041 МИНЕ-ИН;
001 042 ОСТРО-ИН;
001 043 АВТО-РС;
001 044 АНТИ-ЧН;
001 045 АРХИ-ВН;
001 046 АВРО-ВН;
001 047 АЛТО-ВС;
001 048 МОТО-РК, РН;
001 049 ПЕРЕДН-ЧН;
001 050 ЗЛО-СТ;
001 051 ОТВ-ИА, ПЛ, ПТ;
001 052 ВИ-ИП;
001 053 АЕ-ВС, ВЧ, ГТ, ЛВ, НН, НТ, Нб, РГ, РВ, РЗ, РН, РН,
РР, СН, СР, ТС;
001 054 АО-ВЛ, ВО, ГИ, КТ, ЛВ, АГ, АЧ
, МР, НК, НН, ЕК, ХЛ, ХН, ЧЕ;
001 055 ЗА-ПЧ, РП;
001 056 НЕ-ВР, ВН-КТ, НК, НЦ, НЧ, НЦ, РВ, РП;
001 057 ПО-ЗИ, ЧН;
001 058 РЕ-ДАК, ЗЛ, ЛВ, МТ;
001 059 СО-ВС, ВХ, ЛВ, ЛН, НЛ, НН, РТ, ТН, ХН, ШН;
001 060 ++++++
001 061 ++++++
001 062 ++++++
001 063 +
IV

MS 001 064 FP,CK,CT;+++++
001 065 ++
MPS 001 066 1-39,2-18,3-00,4-00,
001 067 P..PP..PP..MM..0..00C..C..T..TT..T08-?
001 068 A.P.,P.,T.A.,T.M.;
001 069 +++++++
001 070 +++++++
001 071 +++++++
001 072 ++
MB 001 073 0000010n203004005006007008009
001 074 90+40=5n/10,11,70w=0M10+20(21)
001 075 +0E80=1n/20l21]30w=0w=0u=0x
001 076 =08101+14-15-16-17-18+1E-1E-13
001 077 +1u/18-1k-1l-1m-1n+1o-1p-1c
001 078 -+7+1y-1d-1x-1z+1q-1w-1s+1b/1b
001 079 +1g+1e+1g-0A-0e-0f+0M/pB-0L-0H
001 080 +0R-0P-0C-0Y-0W-0B-03/1-0K-09
001 081 -01/0b-0p-0o-0q-0n60%-04221-0a
001 082 /0-10! +++++++
001 083 +++++++
001 084 +++++++
001 085 +++++++

ПРИЛОЖЕНИЕ 2

**Примеры автоматического находления документных
мест переносов**

ЭКСПЕРИМЕНТ "TEGCT"

ЗАСТЕРЖЕНІ "СІДІВО"