

УДК 681.14:621.01

ОБ ЭФФЕКТИВНОСТИ АВТОМАТИЧЕСКОГО КОДИРОВАНИЯ
И ИСПРАВЛЕНИЯ ОШИБОК ПРИ ПОДГОТОВКЕ ДАННЫХ

В.В.Хабаров, Ю.Г.Косарев

Одно из серьёзных препятствий на пути широкого применения ЭВМ для автоматизации многих процессов - сложность подготовки информации. Как правило, на пользователя ЭВМ возлагается обязанность представлять информацию в виде, удобном для машины. Пользователь должен привычные для него термины естественного языка (точнее, его модификации, учитывающей специфику данной области) заменить специальными (обычно цифровыми) кодами.

Процесс такого кодирования заметно затрудняется как необходимостью обращаться к перекодировочным словарям, так и тем, что от пользователя требуется гораздо большая скрупулезность, чем та, к которой он привык при составлении документов, рассчитанных на человека. В то же время информация после подобного кодирования теряет свою наглядность, что затрудняет обнаружение ошибок "на глаз", как это обычно имеет место с текстами на естественном языке. Поэтому кодирование приходится повторять дважды, либо делать для контроля декодирование. Возможности автоматического обнаружения ошибок, а тем более их исправления, в данном случае ограничены из-за сложности, а нередко и невозможности построения формальных методов для обнаружения некоторых типичных видов ошибок (например, замена одного кода другим).

Важно также отметить, что этот трудоемкий и малопривлекательный процесс подготовки данных обычно приходится выполнять высококвалифицированным специалистам.

Естественно попытаться облегчить труд этих специалистов переложив его хотя бы частично на ЭВМ, т.е. предоставить чело-

веку возможность описывать данные в привычных для него терминах естественного языка, а остальную часть работы: кодирование, обнаружение и исправление ошибок поручить машине.

Из общих соображений вполне понятно, что при словесном описании данных процесс кодирования должен ускориться, так как отпадает необходимость пользоваться перекодировочными словарями, становитсяunnecessary повторное контрольное кодирование. Кроме того, что весьма важно, из-за большей наглядности информации можно ожидать уменьшения числа формально не обнаруживаемых ошибок, а увеличение избыточности позволяет надеяться на увеличение доли автоматически исправляемых ошибок.

Наряду с этим может заметно увеличиться объём записи, что повлечет за собой рост времени собственно записи у кодировщика и времени перфорации.

Таков общеизвестный качественный характер ожидаемых изменений при переходе к естественному языку (словесному кодированию).

Вполне понятно, что вывод о практической целесообразности перехода к словесному кодированию можно сделать лишь на основании количественных оценок и при непременном условии - создании простых и эффективных программных средств для реализации автоматического кодирования, обнаружения и исправления ошибок.

Далее описывается эксперимент, который позволяет получить требуемые количественные оценки для одной из областей применения - автоматизации проектирования маршрутной технологии изготовления деталей машин, - и предлагается методика автоматического кодирования и исправления ошибок при подготовке данных.

При этом за основу взята одна из реальных систем автоматизации маршрутной технологии, применяемых при проектировании промышленных предприятий.

Для проектирования маршрутной технологии, т.е. состава и размещения оборудования, наиболее опасны ошибки, которые могут привести к неправильному выбору вида оборудования (например, типа станка) либо его марки (т.е. параметров станка). Иными словами, опасны искажения при кодировании и большие изменения значений числовых параметров. Так как последние

можно относительно просто обнаружить, сравнивая их с габаритами детали, то основное внимание в данной работе уделяется процессу кодирования.

1. Методика автоматического кодирования и исправления ошибок. Решалась следующая задача. Задан словарь, каждому слову которого поставлен в соответствие код. Требуется найти для каждого слова исходного текста соответствующее слово из словаря и заменить его кодом.

При отсутствии искажений задача не вызывает каких-либо затруднений. При искажениях нужно найти слово, наиболее близкое к данному, т.е. возникает типичная задача распознавания образов [1]. Методы решения подобных задач обычно опираются на трудоемкую процедуру сравнения данного слова со всеми эталонными словами (словами из словаря), поэтому применялся более простой метод, который теоретически не всегда приводит к успеху, но зато весьма прост в реализации.

Суть этого метода [2] заключается в следующем. Анализируется слово буква за буквой до тех пор, пока не встретится буква, которой данное слово отличается от всех других слов (рис.1). После этого слову присваивается код, а затем проверяется окончание, при совпадении осуществляется переход к следующему слову.

1. Полоса
2. Проволока
3. Пружина
4. Груток

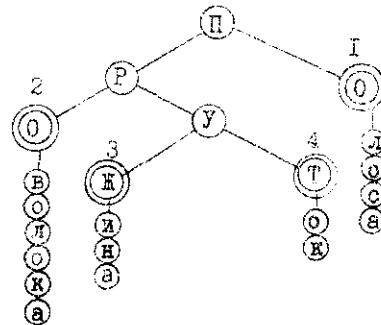


Рис. 1

Буквы, которые отличают данное слово от всех других слов при анализе с начала слова (I), определяют первую базу, при ана-

лизе с конца слова (II) – вторую базу. (На рис. I первые базы подчеркнуты снизу, вторые – сверху.)

Если при анализе слова окончания не совпадают, то слово заносится в список исправлений.

При ошибке в базе анализ слова повторяется по другой базе. Если ошибка при определении этой базы не обнаружена, то слову присваивается код и оно заносится в список исправлений. При ошибке и во второй базе слово заносится в список неисправлений ошибок. Оба списка печатаются в удобной для пользователя форме. Каждой ошибке присвоено имя. В списке исправлений каждое слово с ошибкой сопровождается тем словом, на которое оно исправлено, т.е. пользователю достаточно для контроля правильности исправлений бегло просмотреть этот список.

В качестве языка программирования применяется аппарат формальных грамматик, аналогичный тому, который был предложен И.В.Вельбицким [3].

Важная особенность предлагаемой методики – простота ее использования. От пользователя требуется лишь словарь с указанием кодов. Составление соответствующих программных средств (*R* – таблиц) для кодирования, обнаружения и исправления ошибок выполняется автоматически.

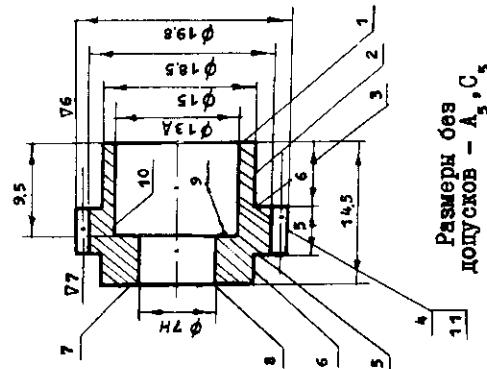
2. Экспримент.

2.1. Структура данных. Информация о детали обычно представлена в виде двух разделов: описания общих параметров детали и описания параметров ее элементов (рис.2). Технология при кодировании использует словарь, который содержит наименования материалов, виды заготовок, названия элементов поверхностей, операции обработки, обозначения параметров.

Описание детали и каждого ее параметра состоит из нескольких уровней записи: символа (μ_0), слова (μ_1), предложения (μ_3), фразы (μ_4). При цифровом кодировании код соответствует предложению, группа кодов – фразе, т.е. уровни μ_0 и μ_1 отсутствуют.

Фразы состоят из последовательности предложений. Первое – определяет название детали или ее элемента, последующие – наименования параметров и их значения. Предложение может состоять из одного слова. Словами языка служат слова из словаря, их со-

$m = 0,3$;
 $z = 64$;
 $cm = 7 \text{ мм}$;
 $d\vartheta = 19,2$.



III

№	черт.	издел.	матер.	к	q
8410017	I	0,6	4	I4	14
заготов.	φ, В	L	H	S	сталь эи-474, К4, Q14
02	19,8	14,5	0	0	штурок, Д19,8, L14,5,
					г5, ч5
№ зод параметры					
1	01	15	I3	I	55 горец, Д15, L1
2	02	15	6	0	55 цилиндр, Д15, L6
3	26	19,8	15	2,4	55 уступ, Д19,8, L2,4
4	02	19,8	5	0	56 цилиндр, Д19,8, L5, ч6
5	26	19,8	15	2,4	55 уступ, Д19,8, L2,4
6	02	15	3,5	0	55 цилиндр, Д15, L3,5
7	01	15	7	4	55 горец, Д15, L4
8	09	7	5	0	27 отверстие, Д7, L5, ч2, ч7
9	10	13	7	3	55 дно отверстия, Д7, L3
10	09	13	9,5	0	27 отверстие, Д7, L9,5, ч7
11	42	19,8	0,3	64	27 зубья шлицевые, Д19,8, М0,3, ч64, ч5, ч7

Рис.2. Зубчатое колесо. Примеры кодирования

кращения, буквы, обозначающие параметры, а также числовые значения параметров. В качестве символов используются все русские буквы, некоторые латинские, а также цифры и разделители.

2.2. Описание эксперимента. За основу словесного кодирования был взят полный вариант языка, который представлен примером на рис.2. Кроме того, был рассмотрен вариант языка (II), в котором допускались сокращения слов.

В эксперименте на этапе кодирования информации участвовали две равноценные по квалификации группы технологов, прошедших предварительное обучение. Квалификационная структура каждой группы отражала реальный состав групп проектного производства. Кодировались две партии чертежей (партии A и B) по 800 чертежей в каждом (табл.I). Чертежи подбирались с учетом практической разноценностии деталей по их технологической сложности.

Данные перфорировались группой операторов на телеграфных аппаратах в цифровом коде ЭВМ "Минск-22" и в коде "М-2". Перфолента проверялась методом дубль-перфорации с последующим сравнением на контрольно-считывающем устройстве.

В ходе эксперимента применялась следующая методика обнаружения и исправления ошибок:

- при цифровом кодировании (вариант I) вся информация повторно кодировалась (без записи), и результат сравнивался с предыдущим (обнаруженные в тексте ошибки тут же исправлялись);
- при словесном кодировании (вариант II и III) непосредственно проверялось соответствие между элементами чертежа и их записью (ошибки также исправлялись сразу же после обнаружения);
- при перфорации (для всех вариантов) применялась повторная перфорация, сверка и печать исправлений.

Общая продолжительность эксперимента с учетом времени на обучение составила 7 месяцев.

Для оценки вариантов кодирования измерялись следующие параметры:

- трудоемкость подготовки данных по видам работ;
- объем информации;
- количество, вид и тип ошибок.

Рассматривались три вида ошибок: пропуск (x), дополнение (y) и замена (z) записи, которые свойственны каждому уровню записи. Эти ошибки, в зависимости от возможности их обнаружения и исправления, классифицировались по следующим типам:

III

Таблица I

Обработка двух партий первичных документов А и В ($A=B=800$)

Этапы подготовки данных	Исполнители	Коли-чество членов группы	Средний оклад, руб.	Количество обрабатываемых исполнителем языка		
				I	II	III
Этап I Заполнение первичных документов	1 группа Инженер-технологи	4	140	1,4	A/2	B
Этап 2 Перфорация данных	2 группа Инженер-технологи	4	140	1,4	B/2	-
Всего первичных документов				800	800	800
III2				200	200	200

А – автоматически исправляемые;

Б – автоматически обнаруживаемые, но не исправляемые;

С – семантические (синтаксически не обнаруживаемые).

Методика обнаружения и исправления ошибок описана выше (п. I).
 З. Обсуждение результатов эксперимента.

З.1. Сравнение цифрового и словесного кодирования. Переход от цифрового кодирования (вариант I) к словесному (вариант II) привел (табл. 2-4) к следующему:

- трудозатраты инженеров (этап I) сократились примерно на 40%;

- трудозатраты перфорации (этап II) возросли примерно вдвое;

- общие трудозатраты уменьшились мало (примерно на 10%).

Однако если учесть качество труда, то выигрыш получится вполне ощутимым. Даже при оценке качества труда по различию в зарплате выигрыш в стоимости получается около 20%. Если учесть различие во времени подготовки инженера (5 лет) и перфораторщика (0,5 года), то выигрыш составит около 35% (из 40% возможных).

Указанное сокращение трудозатрат на кодирование (на 40%) произошло из-за применения кодов, непосредственно отражающих семантику терминов, которые они обозначают (эффект наглядности).

При этом около 3/8 этого выигрыша прометекает от упрощения самого процесса кодирования, так как отпада необходимость обращения к словарю, хотя и увеличился (в 2,7) объем записи. Остальные 5/8 выигрыша приходятся на контроль. Здесь две основные причины: исключение обращений к словарю и уменьшение числа исправлений в тексте.

Число ошибок при кодировании уменьшилось примерно в 2,7 раза, что особенно важно, более чем на порядок сократилось число ошибок, обнаружение которых формальными методами особенно сложно. Более трети всех обнаруженных ошибок исправляются с помощью описанной выше методики.

Число ошибок перфорации возросло примерно в 2,1 раза при росте объема данных в 2,7. Непропорциональность в увеличении числа ошибок подтверждает известный факт, что осмысленный текст печатать проще. Все ошибки перфорации, что весьма важно, ока-

Таблица 2

Среднестатистическое описание I первичный документ

Блоки данных	Этапы подготовки данных	Объем знаков	Трудоемкость обработ- ки, чел/мин	Распредел. труда по этапам	Производст. труда в сравнении с вариантом I	Количество обнаружен. и исправлен. ошибок на ЭВМ	Распредел. по этапам, %
I	Этап 1	15,3	12,8	28,1	78	100	100
	Этап 2	2,1	5,7	7,8	22	100	0,146
	Итого		17,4	18,5	35,9	100	100
II	Этап 1	11,2	5,7	16,9	52	166	251
	Этап 2	380	4,9	10,7	15,6	48	50
	Итого		16,1	16,4	32,5	100	160
III	Этап 1	10,8	6,6	17,4	53	161	223
	Этап 2	270	4,0	11,3	15,3	47	51
	Итого		14,8	17,9	32,7	100	260
							0,100
							24
							76
							100
							83
							17

III4

Таблица 3

Стоимость обработки документов

Вариант языка	Этапы подготовки данных	Стоимость обработки I первичного документа, руб.		Стоимость в сравнении с вариантом I, %
		подготовка	контроль и исправлен. ошиб.	
I	Этап 1	0,214	0,205	0,419
	Этап 2	0,017	0,050	0,067
	Итого	0,231	0,255	0,486
II	Этап 1	0,157	0,091	100,0
	Этап 2	0,039	0,036	100,0
	Итого	0,196	0,187	100,0
III	Этап 1	0,151	0,106	61,3
	Этап 2	0,032	0,102	0,134
	Итого	0,183	0,208	0,391

III5

Министерство спорта по спортивным зонам

зались формально обнаруживаемыми. Около 80% из них исправляются с помощью простейшей методики, описанной выше. Эффект "осмысленности" текста в данном эксперименте, по-видимому, проявился не полностью, так как перфораторщицы, участвующие в эксперименте, обладали в основном опытом работы с числовыми данными.

3.2. О возможности отказа от ручного контроля. Как пока-
зал эксперимент, на 400 деталей встретилось всего 8 ошибок, ко-
торые формально не обнаруживаются методами, описанными выше.
В семи случаях кодировщик неправильно назвал вид поверхности и
в одном пропустил фразу. Допустимо ли такое число ошибок — воп-
рос дискуссионный. Однако одно из серьёзных последствий этих
ошибок — неправильный выбор парка оборудования — может быть в
значительной мере устранено при анализе малозагруженного обо-
рудования на следующих этапах проектирования.

В целом вопрос об отказе от ручного контроля, который при словесном кодировании занимает около 1/3 времени кодировщика, требует дополнительного исследования.

Все ошибки перфорации оказались обнаруживаемыми и в подавляющем большинстве исправимыми. Специальное исследование ошибок перфорации [4] также подтверждает данный вывод. Эти результаты дают основание сделать вывод о возможности для определенных групп перфораторщиц отказаться от повторной перфорации. В согласии с данными эксперимента (табл. 2,3), это более чем втрое сократит трудозатраты на перфорацию при словесном кодировании, и они станут меньше (примерно в 1,6 раза), чем при цифровом кодировании.

3.3. Сравнение вариантов словесного кодирования. Полная форма словесного кодирования дает увеличение количества символов в 2,7 раза. Возникает вопрос, нельзя ли уменьшить это соотношение путем перехода к сокращенным словам, применить такие сокращения, которые не привели бы к потере наглядности? Вопрос однозначно не решается, поскольку при сокращении объема уменьшается и избыточность, т.е. возникает задача на оптимум, основным критерием которой является возможность автоматического исправления ошибок.

Для рассмотренного сокращенного варианта языка (III) (табл.2) эксперимент не выявил каких-либо существенных преимуществ перед вариантом (II) по производительности. В то же время

число ошибок в этом варианте больше, чем у варианта II. Особен-
но плохо то, что из-за уменьшения избыточности произошло уве-
личение числа неисправляемых ошибок [4], а это не позволяет при-
менить для этого варианта режим автоматического исправления
ошибок и отказаться от ручного контроля.

Таким образом, введение сокращенного языка требует допол-
нительных исследований. Поэтому в настоящее время для практи-
ческих целей обоснованным является применение полного варианта
языка.

Предлагаемая методика подготовки нечисловой информации
для ЭВМ даже в простейшем ее виде, благодаря применению привыч-
ного для пользователя языка и переложению на машину процессов
кодирования, обнаружения и исправления ошибок, сокращает при-
мерно вдвое трудозатраты по сравнению с ручным кодированием и
контролем.

Л и т е р а т у р а

1. ЗАГОРУЙКО Н.Г. Методы распознавания и их применение. М.,
"Сов. Радио", 1972.
2. КОСАРЕВ Ю.Г., КОНСТАНТИНОВ В.И., НУРИЕВ Р.М. Метод эф-
фективной обработки текстовой информации (7R-грамматики). От-
чет Института математики СО АН СССР, 1974.
3. ВЕЛЬБИЦКИЙ И.В. Метаязык R-грамматик. -"Кибернетика",
1973, № 3, с. 47-63.
4. ПИЛИПОВИЧ Ю.В., КОСАРЕВ Ю.Г. Статистика ошибок перфо-
рации. Отчет института математики СО АН СССР, 1974.

Поступила в ред.-изд. отд.
19 декабря 1974 года