

УДК 621.391:681.3.06.

РАСПОЗНАВАНИЕ ФОНЕМ С ИСПОЛЬЗОВАНИЕМ АПРИОРНОЙ
ИНФОРМАЦИИ

В.М. Величко

В данной статье описываются эксперименты по распознаванию фонем в отдельно взятых словах. Словарь был подобран так, чтобы все фонемы русского языка встречались в возможно большем числе различных звуко сочетания (этот же словарь использовался в [1]).

Для описания речевого сигнала была принята система признаков, аналогичная использованной в предыдущих экспериментах по распознаванию речи [1,2]. Акустический сигнал пропускался через систему из 4 цифровых октавных фильтров с центральными частотами 900, 1800, 3600 и 7200Гц. Параметры применяемых фильтров были выбраны в работе [3]. Слово разбивалось на сегменты фиксированной длительности 14 мсек. Длительность сегмента имеет тот же порядок, что и максимально ожидаемый период основного тона. В каждом сегменте подсчитывалась энергия сигнала в общей полосе E_0 и на выходе каждого из 4 фильтров E_i ($i=1, \dots, 4$). В качестве параметров, характеризующих сегмент, брались величины $x_0 = \ln E_0$, $x_i = \ln \frac{E_i}{E_0}$, т.е. сегмент описывался точкой \bar{X} в 5-мерном пространстве. Слово описывалось как последовательность сегментов $\bar{X}^{(k)}$.

Сравнение сегмента $\bar{X}^{(k)}$ с сегментом $\bar{X}^{(l)}$ проводилось с помощью меры сходства

$$a = \frac{\alpha^2}{\alpha^2 + \rho^2 \kappa \epsilon}$$

где $\alpha^2 = 2$

$$\rho^2 = \beta^2 + \sum_{i=1}^4 (x_i^{(k)} - x_i^{(l)})^2$$

$$\rho_0 = \begin{cases} |x_0^{(k)} - x_0^{(l)}| - 1 & \text{при } |x_0^{(k)} - x_0^{(l)}| > 1 \\ 0 & \text{при } |x_0^{(k)} - x_0^{(l)}| \leq 1 \end{cases}$$

Целью работы являлось получение фонетической транскрипции распознаваемого слова. Для достижения этой цели задача распознавания была сформулирована как известная задача об оптимальном размещении точек на отрезке (см., например, [4,5], там же приведена библиография). В работе [6] сведение к задаче размещения используется для членения речевого сигнала на квазистационарные участки.

В данной работе результатом явилось распознавание фонем без предварительного членения распознаваемого слова на фонемы. Принципиальная возможность обойти трудный вопрос о членении на фонемы показана в работах [1,7].

Для выделения эталонов фонем параметры всех слов обучающей последовательности были напечатаны в виде графиков. Затем вручную были выбраны участки слов, соответствующие отдельным фонемам в различных звуко сочетаниях. Эталоном фонемы считалась выбранная последовательность сегментов $\bar{X}^{(k)}$. Всего было выбрано 259 эталонов фонем, на каждую из 50 фонем русского языка приходилось от 1 до 8 эталонов из различных звуко сочетаний.

Оптимизируемый функционал учитывает, во-первых, величину меры сходства между участками распознаваемого слова и эталонными соответствующих этим участкам фонем и, во-вторых, вероятности сочетаний соседних фонем.

Мера сходства A_n между участком слова и эталоном фонемы определялась как ненормированная длина максимального пути на матрице сходства $\{a_{kl}\}$, элементами которой являются меры сходства a_{kl} между сегментами участка слова $\bar{X}^{(k)}$ и сегментами эталона фонемы $\bar{X}^{(l)}$. Для нахождения длины максимального пути используется метод динамического программирования [2,8]. Изменение выбранных для описания сигнала параметров может быть весьма значительным на длительности фонемы и зависит как от самой фонемы (например, для взрывных, включающих резко отличающиеся части - смычку, взрыв и придыхание), так и от звуко сочетания, в которое входит фонема (например, для гласных). Поэтому существенным в принятом определении меры сходства A_n является

ся использование представления о "траектории" фонемы, т.е. учет изменения параметров фонемы во времени [2, 9, 10].

Для улучшения качества распознавания учитывается лингвистическая информация о частоте встречаемости сочетаний двух фонем (диад) [11]. В настоящее время получены достаточно полные и надежные данные о двухфонемных сочетаниях русской речи [12]. Однако в данной работе была применена статистика конкретного словаря, которая ввиду специфических требований, предъявленных к словарю, значительно отличается от обычной.

При распознавании слова сначала производилось предварительное определение возможных границ фонем. Общее число границ было взято приблизительно в 3 раза меньше, чем число сегментов в слове, но значительно превышало число фонем. Границы ставились в точках с максимальными $\rho_{i-1, i}^A$ между соседними сегментами, причем каждый интервал между ближайшими границами содержал не менее двух сегментов.

Затем проводилось построение матрицы сходства $f_m(i, \ell)$. Каждый участок слова, начинающийся с i -го интервала и имеющий длину ℓ интервалов, сравнивался с эталонами фонем по сформулированному выше критерию сходства для определения меры сходства $A_{i, m}$. Все $A_{i, m}$ упорядочивались по величине, и 4 различные фонемы, имеющие максимальные меры сходства с рассматриваемым участком слова, составляли элементы "матрицы сходства"

$f_m(i, \ell), m = 1, \dots, 4$. Для сокращения довольно громоздкой процедуры вычисления $f_m(i, \ell)$ применялся предварительный отбор эталонов фонем, у которых средние значения параметров на длине эталона близки к соответствующим средним величинам на участке слова. Таким образом, каждый участок слова мог быть распознан как одна из 4-х фонем, на которые он больше всего похож. Оптимальный подбор сочетаний участков производился в соответствии с функционалом, задаваемым рекуррентно с помощью функционального уравнения:

$$F = \max_n [L_n \rho(t, n) + F_n(t, N)] \quad (I)$$

$$F_m(i, k) = \max_{\ell, n} [f_m(i, \ell) + F_n(i + \ell, k - 1) + L_n \rho(m, n)].$$

$$1 \leq \ell \leq \max\{1, \min[6, \ell_i \max(24), M - i + 1]\}$$

$$N - k + 1 \leq i \leq M - k + 1$$

$$1 \leq m, n \leq 4$$

$$1 \leq k \leq N.$$

Начальные условия: $F_n(M+1, 0) = 0$,
при этом $L_n \rho(m, n) = L_n \rho(m, t)$.
Здесь

$F_m(i, k)$ - значение функционала при оптимальном членении на k фонем участка слова от i -го интервала до конца слова,

i - номер начального интервала участка слова,
 k - число фонем, на которое делится участок слова,

ℓ - число интервалов в участке слова,

M - общее число интервалов в слове,

N - число фонем в слове ($N < M$),

$\ell_i \max(24)$ - максимальное число интервалов, начиная с i -го, с общим количеством сегментов, не превышающим 24. $\ell_i \max(24)$ вводит ограничение на максимальную длину фонемы ввиду того, что длительность фонемы больше, чем $24 \times 14 = 336$ мс, маловероятна.

$L_n \rho(m, n)$ - логарифм вероятности следования фонемы n за фонемой m .

\dagger означает начало или конец слова.

Функционал (I) учитывает как сходство участков слова с эталонами фонем, так и вероятности сочетания фонем. Величина функционала для оптимального распознавания была найдена методом динамического программирования с использованием "матрицы сходства" $\{f_m(i, \ell)\}$ и таблиц сочетаний диад путем последовательного увеличения k от 1 до N и перебора по n, ℓ, m, i . Затем восстанавливались оптимальное размещение границ фонем и соответствующие этому размещению фонемы. Число фонем N в слове задавалось переменным, и для каждого N определялся результат распознавания.

Эксперимент по проверке предложенного алгоритма распознавания проводился на универсальной ЭВМ БЭСМ-6. Две последовательности слов, произнесенные одним диктором (мужчина), записывались на магнитофон в тихой комнате с обычной акустикой. Обе последовательности с помощью 9-разрядного аналого-цифрового преобразователя с частотой квантования 20 кГц были введены в ЭВМ

БЭСМ-6 и записаны на магнитную ленту. При этом автоматически определялись границы слова. Затем производилась цифровая фильтрация, выделение логарифмических параметров и запись их на магнитную ленту БЭСМ-6. После этого первая последовательность выводилась на графики для ручного выделения эталонов фонем. Вторая последовательность предъявлялась на распознавание в качестве контрольной. После распознавания печатались варианты фонетической транскрипции слова при разном числе фонем N . Приведем пример печати для слова "лопасть" (диктор произнес это слово как "ЛОПАС'Т'"):

- ЛОП А -
- ЛОПАС'Т'
- - ЛОПАС'Т'
- - - ЛОПАС'Т'
- - ЛОПАС'Т' -
- - - ЛЫЛАС'Т' -

Отобранные вручную лучшие варианты распознавания для каждого слова приведены в таблице I. В таблице указаны предъявленные слова, фонетическая транскрипция лучшего варианта распознавания (смычка опущена), число правильно распознанных фонем в слове N_p , число недостающих фонем N_n и число лишних фонем N_l . При записи некоторых слов в память БЭСМ-6 отсутствовало начало слова (например, в слове "взрыв"). Это не считалось ошибкой распознавания. Кроме того, не считалась ошибкой замена ударной гласной на безударную и наоборот. Общая надежность распознавания подсчитывалась по формуле

$$G = \frac{N_p}{N_p + 0,5(N_n + N_l)}$$

При ручном выборе наилучшего варианта распознавания надежность получилась равной 83,6% на материале в 374 фонемы. Работа по автоматическому выбору наилучшего варианта пока не закончена, но было подсчитано число ошибок при выборе варианта, равноудаленного от начала и конца напечатанного списка. (В приведенном примере это четвертое сверху слово "---- ЛОПАС'Т'"). При этом надежность составила 77,5%. Среднее время распознавания одного слова ~ 5 сек.

Для сравнения автоматического распознавания фонем с человеческим распознаванием при одинаковой априорной информации был про-

веден следующий эксперимент.^{*)} На основе статистики дикта для русской речи [12] с помощью таблицы случайных чисел был составлен словарь из 47 бессмысленных слов. Словарь был записан на магнитофон в тех же условиях и тем же диктором, что и при автоматическом распознавании, и предъявлен 5 аудиторам. Разрешалось многократное проигрывание записи по просьбе аудиторов. Надежность распознавания, подсчитанная по формуле (4), колебалась для разных аудиторов от 86,3% до 95% и в среднем составила 90,1%. Нижний порог человеческого распознавания в этом эксперименте довольно близок к полученной надежности автомата, но следует отметить, во-первых, использование различной статистики в этих экспериментах и, во-вторых, непроизвольное использование человеком дополнительной априорной информации, что затрудняет сравнение результатов.

В заключение автор выражает благодарность Н.Г.Загоруико за руководство работой и полезные советы, Л.С.Юдиной за составление и транскрибирование словаря, Э.Х.Гимади и В.Д.Гусеву за обсуждение алгоритма.

^{*)} Идея эксперимента и методика его проведения предложены Н.Г. Загоруико.

Таблица I

Слово	Транскрипция	N _п	N _д	N _ж
помпа	пoмпa	5	0	0
свист	св'ист	5	0	0
лопасть	лoпac'т'	6	0	0
жемчужина	жeмчyжинa	8	1	0
фельдфебель	ф'ьл'тф'эб'ьл'	9	0	0
дифференциация	г'иф'эр'энциация	12	2	1
смех	см'эх	4	0	0
специфика	п'ицип'ика	7	1	1
устройство	урoстa	6	4	0
взрыв	гр'ив	3	1	1
гегемон	гeг'эмон	6	1	1
груздь	грyс'т'	5	0	0
гимн	г'имн	4	0	0
кизил	п'из'ил	4	1	1
медянка	гн'ид'янкa	6	1	2
лемех	л'эм'эх	5	0	0
лыжи	л'ижи	4	0	0
двухвостка	двoпoпa	4	6	3
скипетр	к'ип'этр	6	0	0
хищник	х'ищ'ник	6	0	0
жили	жэ	1	3	1
чешуя	чэшy/a	6	0	0
шихта	сyxтa	4	1	1
чечевича	чичэда	5	3	1
резерв	к'из'эф	4	2	1
бегемотики	б'иг'имот'ик'э	9	1	1
экив	п'эф	1	2	2
физик	ф'из'ик	5	0	0

Продолжение таблицы I

Слово	Транскрипция	N _п	N _д	N _ж
шавель	ш'ав'э	4	1	0
шебень	ш'эб'имн'	5	0	1
ходьба	ф'эд'ба	3	2	2
шуплый	ш'уплyj	6	0	0
буддизм	бyд'изм	6	0	0
жуелица	жyэцa	6	2	0
бубен	бyб'ьн	5	0	0
дуги	дyг'ьр'	4	0	1
бурнус	бyс	3	3	0
безденежье	б'ьэд'эн'ьж'э	10	0	0
кимограф	х'имoгрa	6	2	1
гнездовье	гн'издoн	6	3	1
хвоц	лoй'	2	1	1
хилость	х'илac'т'	6	0	0
пикет	п'ик'э	4	1	0
химик	х'им'ьк	5	0	0
питомник	п'ит'ьмн'ьк	6	2	2
грязища	б'ир'эв'ищ'а	6	1	2
пигмей	п'игм'як	4	2	2
сдоба	здoпa	4	1	1
пичуга	п'ичyпa	5	1	1
пелица	п'ив'ицa	6	0	0
воздух	лoздyx	5	1	1
дверь	д'э	1	3	1
гнев	гн'эм'э	3	1	2
Щепи	ш'эцин	5	0	0
бирюза	б'ьр'ьжнa	4	2	3

Продолжение таблицы I

Слово	Транскрипция	№	№	№
иквал	иквал	3	I	0
кунальда	кунальда	5	2	I
кефаль	кефаль	4	I	0
гистохимия	дв'исто/эм'квл	7	3	5
хихикальце	х'ик'идальн	5	4	4

304 70 49

Л и т е р а т у р а

1. Загоруйко Н.Г., Величко В.М., Волошин Г.Я., Гусев В.Д., Ёлкина В.Н., Бахмутова И.В., Хайретдинова А.Г., Юдина Л.С. Эксперименты по автоматическому распознаванию речевых сигналов. Труды АРСО-ГУ. Киев-Канев, 1968.
2. Величко В.М., Загоруйко Н.Г. Автоматическое распознавание ограниченного набора устных команд. "Вычислительные системы", Новосибирск, вып. 56, 1969.
3. Курилов Б.М., Гаврилко Б.П. Сжатие описания сигнала и членение речи на фонемы. "Вычислительные системы", Новосибирск, вып. 56, 1969.
4. Дементьев В.Т. Об одной задаче оптимального размещения точек на отрезке. "Дискретный анализ", Новосибирск, вып. 4, 1965.
5. Гимади (Гимадутдинов) Э.Х. Выбор оптимальных шкал в одном классе задач типа размещения, унификации и стандартизации. "Управляемые системы", Новосибирск, вып. 6, 1970.
6. Винцок Т.К. Оптимальное разбиение последовательности элементов на подпоследовательности. "Кибернетика", Киев, вып. 4, 1970.
7. Загоруйко Н.Г. Алгоритм распознавания фонем по последовательности сегментных решений. Труды АРСО-ГУ, Киев-Канев, 1968.
8. Беллман Р., Калаба Р. Динамическое программирование и современная теория управления. "Наука", Москва, 1969.
9. Бондарко Л.В., Загоруйко Н.Г., Кожевников В.А., Молчанов А.П., Чистович Л.А. Модель восприятия речи человеком. "Наука", Новосибирск, 1968.
10. Голубцов С.В. Труды АРСО-ГУ. Киев-Канев, 1968.
11. Волошин Г.Я. Об использовании языковой избыточности для повышения надежности автоматического распознавания речевых сигналов. "Вычислительные системы", Новосибирск, вып. 28, 1967.
12. Ёлкина В.Н., Юдина Л.С., Хайретдинова А.Г. Статистика двух- и трехфонемных сочетаний русской речи. "Вычислительные системы", Новосибирск, вып. 37, 1969.

Поступила в редакцию
20.12.1970 г.