

УДК 534.78:681.142:62-501.2:513.62

СЖАТИЕ ОПИСАНИЯ СИГНАЛА И ЧЛЕНЕНИЕ РЕЧИ
НА ФОНЕМЫ

Б.М. Курилов, Б.Л. Гаврилко

При проведении настоящей работы авторы руководствовались, в основном, представлениями о модели восприятия речи человеком, изложенными в [1,2]. Основным предположением в работе было то, что в органе слуха с помощью основной мембраны и волосковых клеток кортиева органа происходит разложение входного сигнала $f(t)$ на два сигнала: сигнал мгновенной амплитуды $\alpha(t)$ и сигнал мгновенной частоты $\rho(t)$.

$$f(t) = \alpha(t) \sin \left[\int_0^t \rho(t) dt \right]. \quad (1)$$

Такое разложение квазигармонического сигнала возможно без потери информации лишь в том случае, если на спектр сигнала наложены некоторые ограничения. Рассмотрим эти ограничения.

Заметим, что значения мгновенной амплитуды и мгновенной частоты могут быть измерены не более одного раза за период функции $f(t)$. То есть ограничения, накладываемые на спектр функции $f(t)$, необходимо рассматривать при условии, что $\alpha(t)$ и $\rho(t)$ определяются по дискретным значениям этой функции один раз за период. В работе [3] на основании теорем Котельникова показано, что в этом случае для функции $f(t)$, имеющей спектр частот от ω_1 до ω_2 , должно выполняться неравенство $\omega_2 \leq 2\omega_1$. То есть для представления функции $f(t)$ без потери информации в виде двух функций — мгновенной амплитуды

$\alpha(t)$ и мгновенной частоты $\rho(t)$ – ширина спектра исходной функции должна быть не более одной октавы. Если ширина спектра исходного сигнала более одной октавы, его необходимо предварительно пропустить через набор октавных фильтров $\Phi(j)$. Но для того, чтобы в сигнале полностью отсутствовали частотные составляющие ниже ω_1 и выше ω_2 , амплитудно-частотные характеристики октавных фильтров должны быть, вообще говоря, П-образными. Однако конечность длительности сигнала накладывает ограничения на добротность фильтров (или на их постоянную времени τ). Известно, что для того, чтобы сигнал полностью раскачал фильтр, постоянная времени фильтра τ должна быть, по крайней мере, не больше длительности сигнала. Квазигармонический сигнал минимальной длительности, который может присутствовать в j -м фильтре, очевидно, равен длительности одного периода верхней частоты этого фильтра $f_2[j]$

$$t_{\min}[j] = \frac{1}{f_2[j]} = \frac{1}{2f_1[j]} . \quad (2)$$

Тогда, учитывая (2), запишем:

$$\tau[j] \leq \frac{1}{2f_1[j]} ,$$

или, немного усиливая это требование, что обычно принято делать в теории расчета фильтров, получим:

$$\tau[j] \leq \frac{1}{ef_1[j]} . \quad (3)$$

В силу неравенства (3) и того, что для октавных фильтров

$$\Delta f[j] = f_2[j] , \quad (4)$$

получим следующее соотношение:

$$\tau[j] \times \Delta f[j] \leq \frac{1}{e}$$

или

$$\max(\tau \cdot \Delta f) = \frac{1}{e} . \quad (5)$$

Из (5) следует, что максимально допустимая величина произведения $\tau \cdot \Delta f$ для всех октавных фильтров постоянна и равна $1/e$.

Для одиночного октавного резонансного контура имеем:

$$\tau = \frac{Q}{\pi f_0} ; \quad f_0 = \sqrt{f_1 \cdot f_2} = f_1 \sqrt{2} ,$$

где Q – добротность контура.

Тогда, учитывая (5), запишем:

$$\max Q = \pi f_0 \tau = \pi \sqrt{2} \cdot \tau \cdot \Delta f = \frac{\pi \sqrt{2}}{e} \approx 1,62 .$$

Этот вывод хорошо согласуется с измерениями эквивалентной добротности основной мембраны слухового анализатора (по данным Бекети, эквивалентная добротность порядка 1,5 + 2).

Далее, зная максимально допустимую величину добротности фильтров, максимизируем их частотно-селективные свойства. Под последними будем понимать степень подавления сигналов вне полосы пропускания.

В работе [4] теоретически показано, что максимальными частотно-селективными свойствами при заданной и постоянной для всех фильтров величине обобщенной расстройки α обладает полосовой фильтр, состоящий из цепочки слабо связанных одиночных контуров и имеющий частотную характеристику вида:

$$\sigma'' = [1 + \alpha^2]^{1/2} , \quad (6)$$

где $\alpha = 2Q \frac{\Delta f}{f_0}$ – обобщенная расстройка,

n – число одиночных контуров в полосовом фильтре.

Таким образом, исходный сигнал $f(t)$ можно представить в виде:

$$f(t) = \sum_{j=1}^n a_j(t) \sin \left[\int_0^t p_j(t) dt \right] , \quad (7)$$

где $a_j(t)$ и $p_j(t)$ – мгновенная амплитуда и мгновенная частота на выходе полосового октавного фильтра, имеющего частотную характеристику (6) и временную характеристику, удовлетворяющую выражению (5); n – число полосовых октавных фильтров на заданном диапазоне частот исходного сигнала.

Следует отметить, что представление исходного сигнала в виде (7) является полным в силу второй и четвертой теорем Котельникова и минимальным в том смысле, что уменьшение числа полос системы фильтров или частоты квантования модулирующих функций в общем случае приведет к потере информации.

Однако, при таком представлении исходного сигнала сматия

описания этого сигнала не происходит из-за того, что квантование модулирующих функций ведется из расчета максимально возможной частоты этих функций независимо от их истинной мгновенной частоты. Для того, чтобы квантование модулирующих функций велось в зависимости от их мгновенной частоты, преобразуем выражение (7).

Прежде всего заметим, что функция $a_j(t)$ в (7) представляется в виде временной последовательности:

$$a_j(t) = a_j(t_0); a_j(t_1); \dots; a_j(t_i); \dots,$$

которую можно записать как

$$\begin{aligned} a_j(t) &= a_j(t_0); a_j(t_1) + \frac{\Delta a_j(t_1)}{\Delta t_1} \cdot \Delta t_1; \dots \\ &\dots; a_j(t_{i-1}) + \frac{\Delta a_j(t_i)}{\Delta t_i} \cdot \Delta t_i; \dots \end{aligned} \quad (8)$$

Тогда без потери информации о модулирующей функции $a_j(t)$ можно брать лишь такие дискретные отсчеты этой функции на временной оси, для которых

$$\frac{| \Delta a_j(t_i) |}{\Delta t_i} > 0$$

или, задавшись некоторым уровнем точности ε представления функции $a_j(t)$, такие отсчеты, для которых

$$\frac{| \Delta a_j(t_i) |}{\Delta t_i} > \varepsilon.$$

Тогда

$$\begin{aligned} a_j(t) &= a_j(t_0), t_0; a_j(t_k), t_k; \dots \\ &\dots; a_j(t_\ell), t_\ell; \dots = A_j(t). \end{aligned}$$

Проводя аналогичное преобразование для $R_j(t)$, запишем (7) в виде:

$$f(t) = \sum_{j=1}^m A_j(t) \sin \left[\int_0^t R_j(t) dt \right]. \quad (9)$$

Такое представление исходной функции $f(t)$ по-прежнему является полным, а сжатие описания происходит без потери информации (или с допустимыми потерями ε) и зависит от свойств модулирующих функций. Это сжатие будет особенно эффективным, если модулирующие функции сигнала имеют постоянные значения на больших участках.

Далее, для сокращения динамического диапазона представления функций $A_j(t)$ и $R_j(t)$ проведем их логарифмирование,

сместив каждую из функций на $+\lambda$.

$$\ln[A_j(t)+1] = A_j^*(t); \ln[R_j(t)+1] = R_j^*(t).$$

Тогда выражение (9) можем представить в виде:

$$f(t) = \sum_{j=1}^m [e^{A_j^*(t)-1}] \sin \left\{ \int_0^t [e^{R_j^*(t)-1}] dt \right\}. \quad (10)$$

И, наконец, для представления исходной функции $f(t)$ на временной оси в виде сегментов, содержащих или не содержащих сигнал хотя бы в одной из частотных полос, введем функцию сегментации $\Psi(A_j^*, t)$ по следующему закону:

$$\Psi(A_j^*, t) = \begin{cases} 1, & \text{если } A_j^*(t) = \lambda, \\ 0, & \text{если } A_j^*(t) \neq \lambda. \end{cases} \quad (II)$$

Здесь λ – некоторое пороговое значение. Тогда, разбивая исходный сигнал $f(t)$ по временной оси на сегменты, находящиеся между единичными значениями функции $\Psi(A_j^*, t)$, и вводя обозначения

$$B_\varphi = \sum_{j=1}^m [e^{A_j^*(t)-1}] \sin \left\{ \int_0^t [e^{R_j^*(t)-1}] dt \right\} \Big|_{t=t_1}^{t=t_\varphi} \quad (\varphi = 1, 2, 3, \dots),$$

преобразуем (10) к виду:

$$f(t) = B_1, B_2, B_3, \dots, B_\varphi, \dots \quad (12)$$

Здесь следует отметить, что положение сегментов B_φ на временной оси будет зависеть как от свойств самого сигнала, так и от значения параметров системы $f_1[j], f_2[j], \tau[j], G[j]$, t и λ , причем между параметрами $f_1[j], f_2[j], \tau[j], G[j]$ и t получена взаимосвязь исходя из критерия максимума частотно-селективных свойств системы и минимума потери информации об исходной функции $f(t)$. Таким образом, свободными параметрами системы остались $f_1[1]$ и λ .

Цель настоящей работы заключается в том, чтобы, варируя свободные параметры системы описания речевого сигнала, получить сегментные отрезки B_φ , максимально приближающиеся к фонемам.

Для достижения этой цели был проведен эксперимент с речевыми сигналами на ЭВМ "БЭСМ-6".

Таблица I

Речевой сигнал вводился в ЭВМ через 9-разрядный аналого-цифровой преобразователь с частотой квантования 20 кГц и затем пропускался через гребенку цифровых математических фильтров [5]. На выходе фильтров вычислялись значения функций $A_j^*(t)$ и $R_j^*(t)$ и выдавались на печать в виде графиков. По полученным на ЭВМ "БЭСМ-6" значениям функций вручную вычислялось значение функции сегментации $\Psi(A_j^*, t)$. Функция $R_j^*(t)$ для членения не использовалась.

В качестве речевого материала было взято 23 русских слова в исполнении трех дикторов (двою мужчин и одна женщина), десять английских слов для одного диктора (мужской голос) и десять японских слов для одного диктора (женский голос). Общее число фонем составляло 410. В эксперименте варьировались не только значения свободных параметров системы f_1 , [I] и λ , но и значения всех остальных параметров. Из соображений не принципиального характера были взяты математические фильтры, аппроксимированные по Баттерворту [5]. Частотный диапазон речи разбивался на полосы от 3 до 15. Временные характеристики фильтров варьировались для значения произведения $\tau \cdot \Delta f$ от 1/2 до 1/10; значение λ - во всем динамическом диапазоне. Результаты членения речевого потока на фонемы оценивались методом экспертной оценки (невыделение фонемы "Р" не считалось ошибкой). Наилучший результат членения 99,2% (четыре ошибки) был получен при значениях параметров (см. табл. I): f_1 [I], равном 5 гц, λ , равном 1,2 от среднего уровня шума на участке "молчания", и остальных параметров системы, близких к расчетным.

В заключение следует отметить, что экспертная оценка результатов членения не является, вообще говоря, объективной; действительно, объективной оценкой качества членения может быть надежность распознавания вычисленных фонем.

| № фильтра (j) | Границные частоты f_1 и f_2 | Полоса пропускания Δf | Постоянная времени τ |
|---------------------|------------------------------------|-------------------------------------|------------------------------|
| 1. | 5 + 150 гц | 5 гц | 7 мсек |
| 2. | 150 + 300 гц | 150 гц | 3,5 мсек |
| 3. | 300 + 600 гц | 300 гц | 1,5 мсек |
| 4. | 600 + 1200 гц | 600 гц | 0,87 мсек |
| 5 | 1200 + 2400 гц | 1200 гц | 0,43 мсек |
| 6 | 2400 + 4800 гц | 2400 гц | 0,22 мсек |
| 7 | 4800 + 9600 гц | 4800 гц | 0,11 мсек |

ЛИТЕРАТУРА

1. Л.В. БОНДАРКО и др. Модель восприятия речи человеком. Изд-во "Наука", Новосибирск, 1968 г.
2. Распознавание слуховых образов. Под ред. Н.Г. ЗАГОРУЙКО и Г.Я. Волошина. Изд-во "Наука", Новосибирск, 1966.
3. И.Т. ТУРБОВИЧ, О.А. ПЕТРОВ. Об одном методе полного описания одномерных образов совокупностью простых функций (применительно к речевым сигналам). - В сб. "Опознавание образов". Теория передачи информации. М., Изд-во "Наука", 1965, стр. 25.
4. В.С. ПЕРЕВЕРЗЕВ-ОРЛОВ. Об одном способе улучшения качества радиоприема. - В сб. "Опознавание образов", Теория передачи информации. Изд-во "Наука", М., 1965, стр. 129.
5. В.С. ЛОЗОВСКИЙ и др. Отчет по теме "Ф-7Г" НИО "Факел", Новосибирск, 1968.

Поступила в редакцию
10.1.1969г.