

Амиргалиев Е.Н (доктор технических наук, профессор),

Калимолдаев М.Н (доктор физико-математических наук, профессор),

Мусабаев Р.Р.(кандидат технических наук)

Университет имени Сулеймана Демиреля,

Институт проблем информатики и управления КН МОН РК. Казахстан, Алматы

О некоторых методах синтеза интонации в системах синтеза речи

Введение. Важной составной частью систем синтеза речевого сигнала по тексту является модуль синтеза интонации. Данный модуль предназначен для генерации интонационного контура и его последующего наложения на синтезируемый сигнал. Естественность синтезированного сигнала в значительной степени определяется качеством задания интонационного контура. В процессе синтеза речи важно осуществлять плавное изменение параметров речевого сигнала. В противном случае синтезированная речь будет обладать неестественным звучанием. Таким образом, при построении систем синтеза и распознавания речи актуальной является задача моделирования плавных речевых интонационных процессов. В данной статье дается описание метода синтеза интонационной составляющей речевого сигнала на основе сплайнов – математически рассчитанных кривых, плавно соединяющих отдельные опорные точки интонационного контура.

Известны классические работы ряда зарубежных учёных: Г. Фанта [2], Дж. Фланагана [3], С. Фуруи [4], П. Тэйлора [5], Х. Хуанга [6]. Подобные вопросы также изучаются в работах белорусских и российских учёных: Б. М. Лобанова [1], М. А. Сапожкова [7] и др.

Постановка задачи. Для синтеза речевого сигнала по компилятивному принципу необходимо предварительно получить формализованное описание его фонетических и интонационных свойств. В рамках данного описания для всех фонем необходимо указать интонационные характеристики. В их число входит и множество опорных точек параметрических кривых. При этом параметры соседних фонем должны быть плавно согласованы. Таким образом, в качестве задачи ставится разработка специализированного языка, с помощью которого будет производиться предварительное описание фонетических и интонационных свойств синтезируемого речевого сигнала. Также необходимо осуществить алгоритмизацию процесса расчета гладких параметрических кривых, с помощью которых будет задаваться динамика изменения регулируемых параметров.

Предложенное решение. На рисунках с 1 по 3 показаны основные этапы синтеза речевого сигнала по компилятивному принципу с применением гладких параметрических кривых заданных ограниченным множеством опорных точек. На рисунке 4 показан результат синтеза.

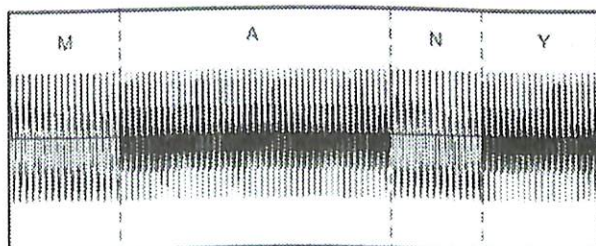


Рисунок 1 – Исходный речевой сигнал после согласования и конкатенации

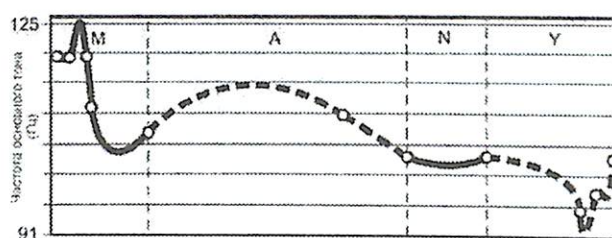


Рисунок 2 – Наложение контура частоты основного тона

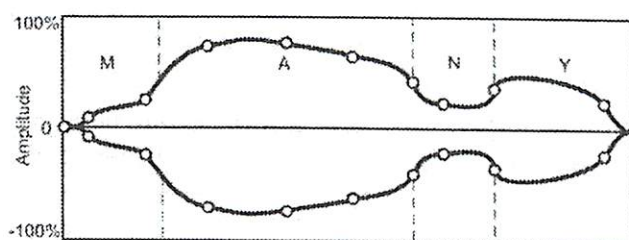


Рисунок 3 – Наложение амплитудных огибающих

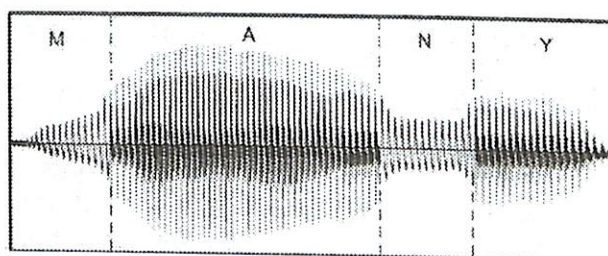


Рисунок 4 – Результат синтеза

Для достижения качественного синтеза важно плавно регулировать следующие параметры речевого сигнала:

1. Контур частоты основного тона – это главная интонационная составляющая речи (рисунок 2).

2. Амплитудные огибающие, основным назначением которых является динамическое регулирование амплитудного уровня сигнала (рисунок 4). Совместное увеличение амплитуды и частоты сигнала приводит к увеличению его громкости.

При компилятивном синтезе [4] на основе базовых фрагментов речи методом различных алгоритмических манипуляций звуковому сигналу придают необходимую форму. Заданная форма речевого сигнала может зависеть от множества различных факторов: от языка, индивидуальных особенностей голоса, синтезируемого текста, требуемой интонации, скорости и громкости произношения и т. д.

Заранее подготовленный, нормализованный по длительности фонем, общему уровню амплитуд и плавно соединённый из различных фрагментов речевой сигнал подаётся на вход системы регулирования параметров (рисунок 1). В зависимости от требуемых интонационных характеристик формируется контур частоты основного тона и накладывается на исходный речевой сигнал (рисунок 2). Затем на сигнал накладываются амплитудные огибающие (рисунок 3).

Для задания кривой выделяется ограниченное множество опорных точек. Выбирается их оптимальное расположение так чтобы наилучшим образом аппроксимировать исходную функцию контролируемого параметра. Изначально в качестве опорных точек выбираются экстремумы аппроксимируемой функции. В рамках решения задачи выбора оптимального расположения опорных точек требовалось оценить значения их координат в заданных диапазонах поисковым методом в смысле минимума критерия (суммы квадратов невязок). Критерий имеет следующий вид (1):

$$K = \sqrt{\frac{1}{N} \sum_{j=1}^N \left(\frac{Y_j^* - Y_j}{Y_j^*} \right)^2}, \quad (1)$$

где Y_j^* , Y_j - соответственно значения аппроксимируемой функции и полученные значения при расчете кривой; N – количество выборок.

Ниже приводится алгоритм вычисления произвольной точки гладкой параметрической кривой.

Входные данные:

- A – множество опорных точек заданных своими координатами $(X; Y)$
- Ax – компонента X элемента множества A
- Ay – компонента Y элемента множества A
- T – задаёт положение вычисляемой точки на кривой, $t \in [0, 1]$

Выходные данные:

- X – координата вычисленной точки по оси X
- Y – координата вычисленной точки по оси Y

Нотация:

- $f(g) = g^3 - g$
- i, j – переменные для счётчиков циклов
- Num – количество элементов множества A
- dT – значение приращения для T на каждый элемент множества A
- dX – значение приращения по оси X
- Px, Py, Wx, Wy, D – множества с количеством элементов равным Num
- div – операция целочисленного деления

1. Инициализация: $Num = \text{Длина}(A) - 1$,

2. Цикл для каждого $i = 1..Num-1$

$$D_i = 4$$

$$W_x = 6 \cdot ((Ax_{i+1} - Ax_i) - (Ax_i - Ax_{i-1})); W_y = 6 \cdot ((Ay_{i+1} - Ay_i) - (Ay_i - Ay_{i-1}))$$

Конец цикла

$$3. Px_0 = 0, Py_0 = 0, Px_{Num} = 0, Py_{Num} = 0$$

4. Цикл для каждого $i = 1..Num-2$

$$Wy_{i+1} = Wy_{i+1} - Wy_i \cdot 0.25; Wx_{i+1} = Wx_{i+1} - Wx_i \cdot 0.25$$

$$D_{i+1} = D_{i+1} - 0.25$$

Конец цикла

5. Цикл для каждого $i = Num-1..1$

$$Px_i = \frac{Wx_i - Px_{i+1}}{D_i}; Py_i = \frac{Wy_i - Py_{i+1}}{D_i}$$

Конец цикла

$$6. X = Ax_0; Y = Ay_0$$

$$7. dX = Ax_{Num} - Ax_0$$

8. Если $dX > 0$ тогда

Начало

$$dT = \frac{1}{Num}$$

Цикл для каждого $i = 0..Num-1$

Если $(dT \cdot i \leq T)$ и $(dT \cdot (i+1) \geq T)$ тогда

Прерывание цикла

$$T = (T - (T \text{ div } dT) \cdot dT) \cdot Num$$

$$X = T \cdot Ax_{i+1} + (1-T) \cdot Ax_i + \frac{f(T) \cdot Px_{i+1} + f(1-T) \cdot Px_i}{6}$$

$$Y = T \cdot Ay_{i+1} + (1-T) \cdot Ay_i + \frac{f(T) \cdot Py_{i+1} + f(1-T) \cdot Py_i}{6}$$

Конец.

Для решения задачи синтеза речевого сигнала [4] разработан унифицированный язык фонетического представления (Unified Phonetic Language - UPL). Фактически данный язык является расширенной фонетической транскрипцией. На языке UPL описываются требуемые характеристики речевого сигнала, на основе которых компилятор выбирает наиболее подходящие элементы компиляции и осуществляет последующую генерацию речевого сигнала.

Синтезируемый сигнал на унифицированном языке фонетического представления описывается в следующем виде:

Фонема1Ударение1(Длительность1;Аллофон1;[ЧОТ1];{Амплитуда1})

Фонема2Ударение2(Длительность2;Аллофон2;[ЧОТ2];{Амплитуда2})

...

ФонемаNУдарениеN(ДлительностьN;АллофонN;[ЧОТN];{АмплитудаN}),

где *ФонемаN* – мнемоническое обозначение фонемы, *УдарениеN* – признак ударения, *ДлительностьN* – значение длительности звучания фонемы в миллисекундах, *АллофонN* – номер аллофонной реализации фонемы, *ЧОТN* и *АмплитудаN* – соответственно контура частоты основного тона и амплитуды, которые задаются следующим образом: $X_1, Y_1; X_2, Y_2; \dots; X_m, Y_m$, где X_m – относительная координата m -ой опорной точки на отрезке от начала (0) и до конца (1) звучания фонемы $X_m \in [0;1]$, Y_m – значение частоты основного тона ($X_m \in [0;+\infty]$) либо амплитуды ($X_m \in [0;1]$). Приведём пример описания слова «тапу» на унифицированном языке фонетического представления:

PAU(160;1)

M(43;1199:[0,98;0.534,99.46;1,101.28];{0,0;0.6,0.1;1,0.2})

EH1(82;1:[0,101.28;0.5,106.1,103.46];{0,0.2;0.5,0.21;1,0.2})

N(78;1:[0,102.92;0.452,104.38;1,106.74];{0,0.2;0.5,0.1;1,0.2})

IY1(127;1184:[0,106.74;0.473,108.81;1,102.03];{0,0.2;0.5,0.21;1,0})

PAU(160;1).

Выводы. Произведено сравнение предложенного в данной работе метода с методом линейной интерполяции, который используется в большинстве существующих систем синтеза речи [ссылка на Лобанова]. Оценка производилась по критерию минимума суммы квадратов невязок по аналогии с формулой (1) между расчетными значениями по двум методам и натуральным эталонным контуром. В результате для метода линейной интерполяции критерий в среднем равен 0.25, в то время как для предложенного метода значение критерия составляет в среднем 0.07. Таким образом, параметрическое описание синтезируемого речевого сигнала на языке UPL в совокупности с методом аппроксимации его интонационной составляющей сплайнами позволяют добиться улучшения качества аппроксимации по сравнению с существующими методами.

Разработанный язык UPL позволяет задавать и описывать разнообразие фонетических и интонационных форм устной речи. Все исходные данные описываются с помощью унифицированного языкового представления, что позволяет осуществлять гибкое межсистемное взаимодействие и на качественном уровне решать задачу синтеза речевого сигнала.

Предложенный подход может также использоваться в системах распознавания речи и идентификации диктора по особенностям его интонации.

Литература

1. Амиргалиев Е.Н., Мусабаев Р.Р. Информационные технологии создания систем синтеза речи// Вестник КазНПУ имени Абая. -2007.-№4(20). -С.26-34

2. Амиргалиев Е.Н., Мусабаев Р.Р. Некоторые направления и задачи обработки лингвистических данных// Вестник КазНТУ им. К.И.Сатпаева. -2007. - №6(63). – С.182-187.
3. Амиргалиев Е.Н., Мусабаев Р.Р. The algorithms of phonemes classification in field of complete speech synthesis systems realization// Труды международной конференции «Проблемы кибернетики и информатики, РСІ2008», Баку, 2008.

Resume

An important component part of the most text-to-speech synthesis systems is the intonation synthesis module. This module is intended for generation of an intonation contour and its subsequent applying to a synthesized signal. Naturalness of the synthesized signal is substantially determined by intonation contour definition quality.

Түйін

Бұл мақалада сөйлеуді синтездеу жүйелерінде интонацияны синтездеу әдістері ұсынылған. Ұсынылған әдісте интонацияның контурын анықтау арқылы сөйлемдік сигналға енгізу модулі қарастырылған, әрине сөйлеуді синтездеудің сапасы анықталған контурға тәуелді.

Özet

Text-to-speech sentez sistemlerinin önemli bir bileşeni parçası tonlama sentez modülüdür. Bu modül bir tonlama kontur üretimi için tasarlanmıştır ve bir sonraki sentezlenmiş bir sinyali için geçerli. Sentezlenmiş sinyal Doğallık tonlama kontur tanımı kalitesini önemli ölçüde tarafından belirlenir.