

Modul universal pentru sinteza mesajelor cu voce multilingve

Sergiu TINCOVAN, Iurie SOROCEANU, Vitalie SECRIERU, Eugeniu MUNTEANU, Valerian DOROGAN
Technical University of Moldova

s_tincovan@mail.ru, iuis@mail.md, primcast@mail.com, eugeniumunteanu@gmail.com, dorogan_lme@yahoo.com

Abstract — Most voice messaging devices and systems are provided for a single language, and changing messages in to another language represents a technical difficulties. Elaborated method supports two and more languages at the same time, as will as gives possibility of completing and correcting the dictionary with reasonable expenses. The method of synthesizing voice messages is based on the original procedure for using standard methods.

Index Terms — Microcontroller, delta-modulation, spectral analysis, phonemes, frequency formats, vocal sound, consonant sound.

I. INTRODUCERE

În tehnică tot mai frecvent necesitatea de a transforma mesajele de text sau codul mesajului de eroare în expresii cu voce. Ca exemplu pot servi sistemele de diagnosticare, robot-secretar al serviciilor de dispecerat, sistemul de bord în transportului, telefonie, bioacustica, medicină, sistem de educație, etc.

Majoritatea sistemelor de tipul TTS (Text To Speech, sau text în voce) sunt prevăzute pentru menținerea unei singure limbi, ce nu întotdeauna poate asigura schimbarea operativă a limbii mesajelor în sistem. Schimbarea limbii în dispozitivele (sistemele) se sintetizează a mesajelor cu voce necesită schimbarea chip-ului specializat sau reprogramarea circuitelor de memorie cu altă bază de date a altei limbi. Altă problemă asemenea sistem este imposibilitatea de menținere în paralel a 2 și mai multe limbi, unde e posibilă complectarea dicționarului pe parcurs.

În lucrarea dată sunt analizate modalitățile de realizare a modulelor de sinteză a mesajelor cu voce de uz universal, ce permite soluționat majoritatea problemelor menționate.

II. Formularea problemei

Pentru mărirea eficienței proceselor este necesar de soluționat următoarele sarcini:

1) De analizat metodele de sinteză a mesajelor la nivel de algoritmi și mijloace tehnice.

2) De obținut o variantă optimală a structurii și algoritmului de funcționare pentru modulul de sinteză, ce permite de modificat operativ limba și dicționarul ei.

3) De obținut structura și algoritmul de funcționare optimal a modulului de sinteză a sunetelor și mesajelor cu voce, unde dicționarul (vocabularul) de sunete și foneme este limitat după volum.

Pe parcursul soluționării este necesar de ținut cont de repartizarea funcțiilor executate între partea HARD și SOFT a modulului de sinteză mesajelor cu voce, specificul de preparare a bancului de foneme și sunete cu ajutorul mijloacelor SOFT specializate al PC aparte.

III. Căile de soluționare și implementare

Variantele posibile de realizare a sintezei sunetelor speciale și vocii pot utiliza următoarele tehnologii:

1) Metoda directă de codare-restabilire a semnalelor sunetelor și vocii;

2) Modelarea digitală a tractului vocii;

3) Sinteza analogică a formantelor de frecvență;

4) Sinteza digitală a fonemelor.

Pentru metoda directă codare-restabilire a semnalelor sunetelor și vocii codarea semnalelor vocii se efectuează cu ajutorul PCM (Pulse Code Modulation) sau delta-modulare [1, 2].

Informația despre semnal vine în formă succesiunii de eșanțonări cu frecvența de discretizare f_{discr} , Conform teoremei Kotelnicov este necesar de respectat $2 \cdot f_{max} \leq f_{discr}$, unde f_{max} este frecvența maximă a spectrului semnalului vocii. Această relație este valabilă pentru filtrare ideală a funcției restabilite, pentru condițiile reale este acceptată condiția $f_{discr} = 45 \cdot f_{max}$.

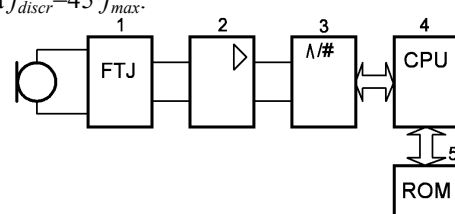


Fig. 1. Convertirea semnalului cu ajutorul PCM.

Pentru intervalul de frecvențe, unde valoarea maximă nu depășește 8-10 KHz în tractul preliminar de prelucrare a semnalului se include filtru trece jos cu scopul de a elimina frecvențele mai înalte $f_{discr}/2$ pentru intrarea CAD. Pentru imprimarea și redarea vocii este necesar să fie asigurată compatibilitatea CAD și DAC după binaritate și frecvența discretizării. Rata de transfer a datelor constituie circa 96 Kbit/s pentru o secundă de sonorizare, care permite de simplificat partea HARD a echipamentului și asigură calitate suficientă a sintezei sunetului.

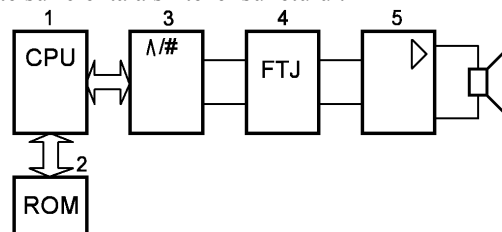


Fig. 2. Schema de structură a sintezării vocii cu ajutorul PCM.

În calitate de alternativă a metodei directe de codare-restabilire a semnalelor sunetelor și vocii poate servi delta-modularea, care permite de redus rata de transfer, care se bazează pe variația relativă a amplitudinii și nu pe valorile

absolute. Semnalul vocii ca și pentru PCM este prelucrat preliminar, care apoi este supus delta-modulării.

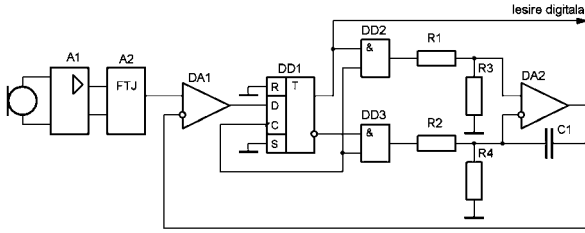


Fig. 3. Schema funcțională a delta-modulatorului.

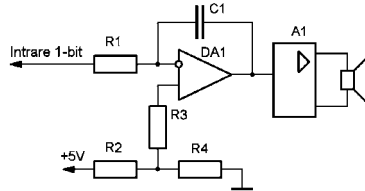


Fig. 4. Schema funcțională a delta-demodulatorului.

Sinteza analogică a formantelor de frecvență se bazează pe cunoașterea detaliată a fonemelor și descompunerii fonetice a mesajelor, la temelia cărora stau două noțiuni: lingvistică – foneme, și acustică – formantă.

Sub noțiunea de formantă se subînțeleg frecvențele de rezonanță (polurile funcției de transfer) a sistemului acustic vocal. Parametrii formantelor (frecvența, banda de trecere și amplitudinea) sunt determinate de parametrii tractului vocal. Cel mai important parametru, frecvența este strâns legat de configurația geometrică a tractului vocal, unde în procesul vorbirii odată cu schimbarea parametrilor configurației se modifică și valorile frecvențelor formantelor (fig. 5).

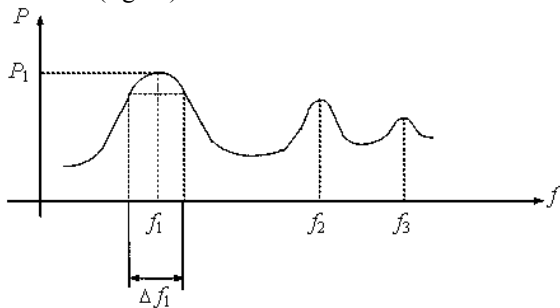


Fig. 5. Exemplu al spectrului vocii.

Pentru sinteza vocii sunt suficiente 2-4 frecvențe a formantelor, unde prima formantă are frecvența 200 Hz (prima formantă a vocii bărbătești) până la 2000 Hz (a treia formantă a vocii feminine). În procesul vorbirii frecvențele formantelor și obertoanelor sunt prezente simultan și se deplasează pe axa spectrului, ce corespund particularităților cuvântului pronunțat. Deaceia mesajele vocale sau sunetele lumii animale se percep nu ca o singură frecvență, ci o mulțime de armonici, care se formează la filtrarea impulsurilor, formate la ieșirea tractului vocal (fig. 6).

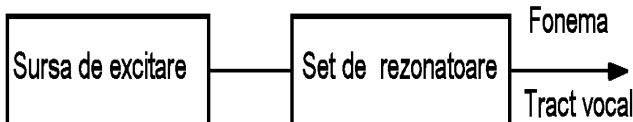


Fig. 6. Modelul tractului vocal.

În proces de vorbire tractul vocal funcționează în calitate de un set de rezonatoare ce filtrează spectrul semnalului de excitație, ca rezultat apare un tablou pe caracteristica spectrală, în care se conțin un șir de maxime (ele și reprezintă rezonanțele formantelor). Ca exemplu este format un tabel pentru vocalele „o”, „a” și „i”.

Tabelul 1.

Fonema	Frecvențele formantelor		
	F ₁	F ₂	F ₃
o	275	850	2400
i	250	2300	3000
a	575	900	2450

Prin generarea simultană a frecvențelor F₁, F₂, și F₃ conform tabelului 1 putem obține sunetele vocale, schema de structură a sintetizatorului este dată mai jos (fig. 7.)

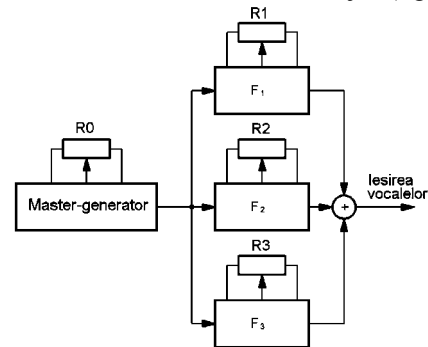


Fig. 7. Schema de structură a sintezei sunetelor vocale.

Schema sintetizatorului de sunete vocale conține master-generator al tonului principal și filtre trece bandă reglabile, ce pot fi acordate pe frecvențele formantelor (cu rezistențele R₁-R₃) și sumator, ce mixează semnalele de la ieșirea filtrelor. Corespunzător spectrograma la ieșire va conține trei frecvențe a formantelor identice acelor sunete vocale ce sunt pronunțate pe viu.

Este mult mai dificil de sintetizat sunetele consonante, ele se formează când în cavitatea bucală se formează obstacole pentru aerul expirat, așa tipuri de consonante sunt:

- explozive (sunetele „p”, „t”, „k”)
- fricative (sunetele „s”, „f”, „h”)
- nazale (sunetele „n”, „m”)
- affricata (sunete „ci”, „t-ș-ci”, „ț”, „ts”)

Schema de structură unui asemenea modul permite sinteza a sunetelor vocale și consonante, din care se poate de format sunete speciale și mesaje (fig. 8), unde sunetele sunt sintetizate în ordinea prescrisă de vocabular (dicționar) [3].

Sinteza digitală a fonemelor se bazează pe generarea fonemelor și compilarea ulterioară din ele a sunetelor speciale, cuvinte, propoziții și fraze. Realizarea ei combină metode compacte de prelucrare digitală și flexibilitatea de gestionare cu parametrii principali a vocii, ce este specific modelelor de formante (fig. 9). Procesul codării vocabularului necesar este substituit cu compilarea mesajelor arbitrare dintr-un set de elemente a vocii, pregătit preliminar [4].

Sintetizatorul de foneme conține 3 nivele de prelucrare, unde la primul nivel se efectuează translatarea simbolurilor orfografice în codurile fonemelor, la nivelul doi se calculează setul parametrilor acustici, care servește pentru gestionarea nivelului trei – formarea semnalului vocii.

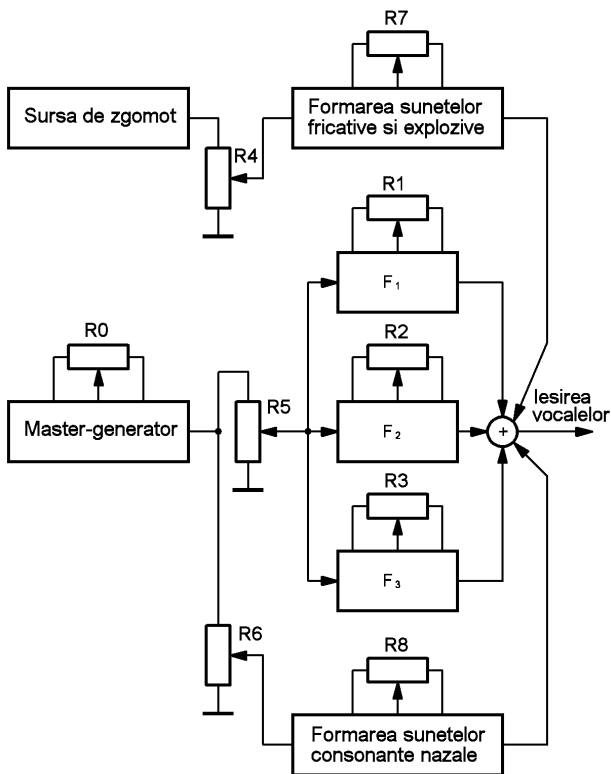


Fig. 8. Schema de structură a sintetizatorului sunetelor vocale și consoane.

Pentru nivelul doi de prelucrare se calculează setul parametrilor de gestionare, care reprezintă frecvențele formantelor (F_1, F_2, F_3), banda de frecvențe ($\Delta F_1, \Delta F_2, \Delta F_3$), frecvența tonului principal și amplitudinea de vocalizare. Setul parametrilor regenerează cu interval de 6,4 ms, ce asigură de urmărit cele mai rapide variații între foneme, rata de transfer de la nivelul doi la nivelul trei constituie circa 45 Kbit/s. Sintezarea propriu zisă se efectuează la nivelul trei, unde semnalele de excitație (armonic și de zgomot) se filtrează cu ajutorul filtrelor de rezonanță a formantelor a tractului vocii.

Frecvența de eșanționare constituie 10 kHz, ce permite de redat componentele de frecvență până la 5 kHz. Ca exemplu poate servi procesorul DSP de tipul TMS32010 a firmei „Texas Instruments” (SUA)

Codarea vocii cu ajutorul coeficienților liniari de predicție (CPL) se bazează pe teoria analizei statistice a seriilor de timp. Seria de timp reprezintă o succesivitate de observații (referințe) aranjate în timp. Esența acestei metode este următoare. Admitem succesivitatea seriei a semnalului vocal x_1, x_2, \dots, x_t . Pentru această mulțime se calculează „valoare medie”. Pentru un interval prestabilit (de exemplu 1020 ms) proprietățile statistice a semnalului vocal sunt neschimbate. Acest interval se codează cu un set de coeficienții a_s , care minimizează eroarea medie patrată a predicției, cu alte cuvinte reduce la minim eroarea între seria inițială și „netezită”. Calcularea coeficienților constituie o procedură dificilă (se calculează ecuațiile diferențiale cu metoda patratelor minime). Principiul analizei CPL și codării sunt date pe fig. 10.

Ca rezultat analiza CPL a vocii codate PCM, alcătuită din eșanționări cu frecvența 1020kHz, se convertează în serie de vectori a parametrilor cu frecvența 50100Hz, unde

compresia descrierii vocii se comprimă de 50-100 ori cu pierdere neesențială a calității.

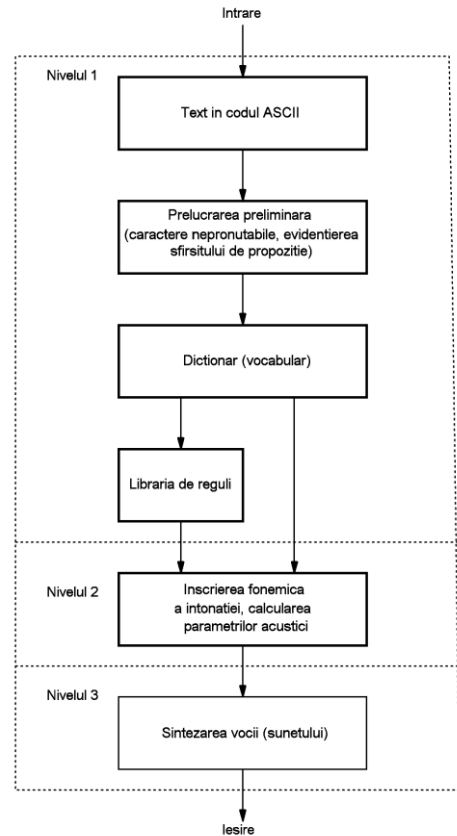


Fig. 9. Structura sintetizatorului digital de foneme.

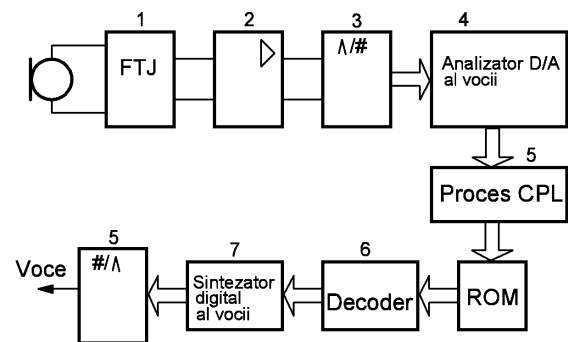


Fig. 10. Schema de structură a codării/decodării CPL.

Metoda dată de sinteză a vocii combină în sine performanțele PCM și sintezei cu formante. În CPL sinteză a vocii se utilizează filtre liniare recursive de imitare a tractului vocal (fig. 11).

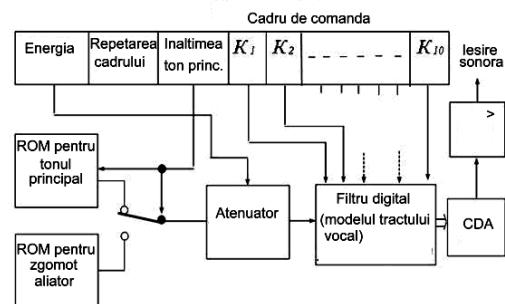


Fig. 11. Structura sintetizatorului CPL.

Ultimele 10 elemente a cadrului CPL corespund coeficienților filtrului digital pentru sinteza vocii. Pentru implementare sunt utilizate procesoare specializate sau procesoare programabile, deoarece viteza de lucru a procesoarelor tradiționale este insuficientă pentru asemenea scop.

Metoda dată este destul de complicată în realizare, deoarece ea cere echipament de implementare foarte rapid. Însă această metodă se consideră de perspectivă, deoarece el se bazează pe procedee perfecte de formare a vocii. Posibilitatea de gestionat cu parametrii modelului permite de acordat sunetele și cuvintele după nivelul energiei, tempou și tonalitate. Acest lucru permite de sintetizat mesaje mai complicate pe baza setului de elemente cu ajutorul regulilor.

Pentru implementarea modulului de sinteză din variantele menționate pentru cerințele de menținere dicționarului nelimitat cele mai reușite metode este analiza formantelor (poate sinteza orice sunet a limbilor existente) și metoda fonemică (poate sinteza mesaje în mai multe limbi). Alt argument în folosul sintezei fonemice este cea mai joasă rată transfer pentru canalele de legătură (tabelul 2.), unde se păstrează calitatea sunetului sintetizat.

Tabelul 2.

Metoda de codare a vocii	Cheltuieli informaționale, bit/s
PCM	$(40-100) \cdot 10^3$
Delta-modulare	$(20-50) \cdot 10^3$
Delta-modulare adaptată	$(10-25) \cdot 10^3$
Analiza formantelor	$(2-4) \cdot 10^3$
Metodă fonemică	50-100

IV. CONCLUZII

Pentru implementarea modulului de sinteză a mesajelor cu voce este necesar de utilizat algoritmul din fig. 9, unde la nivelul 3 de recurs la redarea samplerelor fonemelor, ce sunt extrase dintr-un banc foneme, unde tipul și

succesiunea este definită de către program la nivelul 2 (fig. 9.).

1) Pentru sinteza sunetelor și mesajelor este suficient să fie la dispoziție samplere a circa 300-600 silabe în diferite combinații, pentru obținerea parametrilor mai performanți este necesar de mărit numărul samplerelor a variantelor de silabe.

2) Calitatea sunetului va fi determinată de binaritatea și frecvența de eșanționare a semnalului vocii-donor, pentru telefonie, sisteme de bord, servicii de dispecerat, etc. este suficient de utilizat binaritatea discretizării 8 biți și frecvența de eșanționare 22500 Hz.

3) Pentru convertirea codului în semnal acustic e suficient de recurs la metoda PCM sau PWM, alte modalități sunt mai costisitoare în resurse HARD

Ca un neajuns a variantei propuse este necesitatea de a obține semnal a vocii-donor, unde dicatorul citește un text special, ce conține maximul combinațiilor de silabe, apoi imprimarea este prelucrată minuțios. Rezultatul prelucrării imprimării va fi acel banc de foneme, ce se obține prin decuparea silabelor din fonograma textului citit.

Pentru înlăturarea acestor neajunsuri este necesar de inclus elemente de intelect artificial ce va constitui direcția de cercetări ulterioare.

REFERINȚE

1. Г. Нуссбаумер. Быстрое преобразование Фурье и алгоритмы вычисления сверток.— М.: Радио и связь, 1985.

2. А. В. Фролов, Г. В. Фролов. Мультимедиа для Windows. Библиотека системного программиста, т. 15 М: Диалог-МИФИ, 1994

3. Л. Захаров. Проблемы создания аллофонной базы автоматического синтеза речи (<http://art.bdk.com.ru/govor/rasp.htm>).

4. Black A.W., Taylor P. Automatically clustering similar units for unit selection in speech synthesis // In Proceedings of Eurospeech 97. Rhodes, Greece, 1997. Vol.2, pp. 601-604.