



(12) CERERE DE BREVET DE INVENȚIE

(21) Nr. cerere: a 2014 00347

(22) Data de depozit: 07/05/2014

(41) Data publicării cererii:  
29/01/2016 BOPI nr. 1/2016

(71) Solicitant:

- BUZO ANDI,  
BD. GENERAL VASILE MILEA NR. 8,  
BL. B2, SC. 3, ET. 4, AP.57, SECTOR 6,  
BUCUREȘTI, B, RO;
- CUCU HORIA, ALEEA POLITEHNICII  
NR. 6, BL. 3, SC. 4, ET. 2, AP. 42,  
SECTOR 6, BUCUREȘTI, B, RO;
- PETRICĂ LUCIAN, ÎNTRAREA VÂSLEI  
NR. 1, BL. PM63, SC. 2, ET. 5, AP. 77,  
SECTOR 3, BUCUREȘTI, B, RO;
- BURILEANU DRAGOȘ,  
STR. JOHANNES KEPLER NR. 2 BL. 2,  
SC. 2, AP. 61, SECTOR 2, BUCUREȘTI, B,  
RO

(72) Inventatori:

- BUZO ANDI,  
BD. GENERAL VASILE MILEA NR. 8,  
BL. B2, SC. 3, ET. 4, AP.57, SECTOR 6,  
BUCUREȘTI, B, RO;
- CUCU HORIA, ALEEA POLITEHNICII  
NR. 6, BL. 3, SC. 4, ET. 2, AP. 42,  
SECTOR 6, BUCUREȘTI, B, RO;
- PETRICĂ LUCIAN, ÎNTRAREA VÂSLEI  
NR. 1, BL. PM63, SC. 2, ET. 5, AP. 77,  
SECTOR 3, BUCUREȘTI, B, RO;
- BURILEANU DRAGOȘ,  
STR. JOHANNES KEPLER NR. 2 BL. 2,  
SC. 2, AP. 61, SECTOR 2, BUCUREȘTI, B,  
RO

(74) Mandatar:

CABINET D.NICOLAESCU, STR.TURDA,  
NR.102, BL.30A, ET.7, AP.28, BUCUREȘTI

(54) METODĂ ȘI SISTEM PENTRU DIARIZARE ÎN TIMP REAL A SEMNALELOR AUDIO, UTILIZATE PENTRU RECUNOAȘTEREA AUTOMATĂ A VORBIRII ȘI A VORBITORULUI

(57) Rezumat:

Invenția se referă la o metodă și la un sistem pentru diarizarea în timp real a semnalelor vocale și este destinată a fi utilizată în domeniul sistemelor de procesare a semnalului audio pentru recunoașterea automată a vorbirii și identificarea automată a vorbitorului. Metoda conform invenției constă în citirea periodică a unui număr de vectori de caracteristici audio, segmentarea vectorilor, stocarea modelelor audio în istoria de diarizare, gestionarea periodică a istoriei și folosirea unei funcții de cost bazată pe vechimea modelului și ponderea sa în istoria de diarizare. Sistemul conform invenției constă dintr-o memorie tampon (303, 304) pentru vectori de caracteristici audio, o memorie (201) pentru segmente, o memorie (202) pentru grupuri de segmente, o memorie pentru modele (203), un serviciu (102) extern de segmentare și automate de control (301), care gestionează istoria de diarizare.

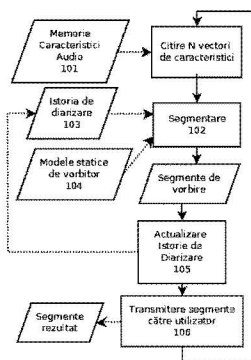


Fig. 1

Revendicări: 5  
Figuri: 3

Cu începere de la data publicării cererii de brevet, cererea asigură, în mod provizoriu, solicitantului, protecția conferită potrivit dispozițiilor art.32 din Legea nr.64/1991, cu excepția cazurilor în care cererea de brevet de invenție a fost respinsă, retrasă sau considerată ca fiind retrasă. Întinderea protecției conferite de cererea de brevet de invenție este determinată de revendicările conținute în cererea publicată în conformitate cu art.23 alin.(1) - (3).



## METODĂ ȘI SISTEM PENTRU DIARIZARE ÎN TIMP REAL A SEMNALELOR AUDIO, UTILIZATE PENTRU RECUNOAȘTEREA AUTOMATĂ A VORBIRII ȘI A VORBITORULUI

Invenția aparține domeniului sistemelor de procesare a semnalului audio pentru recunoașterea automată a vorbirii și identificare automată a vorbitorului.

Un sistem de recunoaștere a automată a vorbirii are ca scop transcrierea unui semnal audio, de cele mai multe ori înregistrarea unei conversații sau a unui monolog. Prin recunoaștere se obține un fișier text care conține cuvintele rostite. În cazul în care semnalul audio conține segmente de muzică, zgomot, efecte speciale, sau în cazul în care proprietățile semnalului audio se schimbă în timp, rezultatele procesului de recunoaștere sunt afectate în mod negativ, prin introducerea de erori. De exemplu, sistemul de recunoaștere va încerca transcrierea unui semnal muzical, rezultând text fără sens gramatical. Probleme similare sunt întâlnite și în cazul recunoașterii vorbitorului, care nu este posibilă atunci când înregistrarea audio conține muzică sau zgomot, sau când vorbirea înregistrată conține semnal vocal de la mai mulți vorbitori.

Diarizarea este procesul prin care semnalul audio este analizat și segmentat, astfel încât fiecare segment are proprietăți uniforme din punct de vedere audio și conține un singur fel de semnal, care poate fi liniște, muzică, vorbire sau alte tipuri de semnal audio. Recunoașterea se poate face apoi doar pe segmentele de vorbire.

Segmentele de vorbire identificate pot fi suplimentar procesate în cadrul procesului de diarizare pentru separarea vorbitorilor, folosind proprietăți audio cum ar fi frecvența fundamentală a semnalului vocal. În urma acestui pas, fiecare segment de vorbire aparține unui singur vorbitor și poate fi folosit pentru identificarea vorbitorului respectiv.

În cele ce urmează vom enunța terminologia folosită în descrierea prezentei invenții, cu mențiunea că se folosește terminologia consacrată specifică domeniului, în limba engleză:

- Frame – Fereastră audio – un număr de eșantioane ale semnalului audio digital, reprezentând în mod obișnuit un interval de timp fix, de ordinul milisecundelor (10-20ms)

- Speech Features – Vector de caracteristici audio - coeficienții cepstrali de frecvență (MFCC) și a alte măsuri ale semnalului din o fereastră audio, folosite pentru procesul de diarizare. Semnalul audio este reprezentat în procesul de diarizare ca o înșiruire de vectori de caracteristici consecutivi.
- Segment – Segment – un număr de vectori de caracteristici audio consecutivi, care au proprietăți similare
- Cluster – Grup – un număr de segmente consecutive, care au proprietăți similare, de exemplu aparțin aceluiași vorbitor
- Speaker model – Model audio - un model de mixtură Gaussiană (GMM) care aproximează caracteristicile vorbitorului. Un model de vorbitor poate fi antrenat folosind vectorii de caracteristici audio ai unui grup. Un GMM poate de asemenea să aproximeze clase mai largi de semnal audio, cum ar fi muzica, liniște, vorbire, voce de bărbat, voce de femeie, și altele.

Stadiul cunoscut al tehnicii, în ceea ce privește sistemele de diarizare automată, este analizat în [S. Galliano, G. Gravier, L. Chaubard, "The ester 2 evaluation campaign for the rich transcription of French radio broadcasts," In Proc. Interspeech, pp. 2583-2586, 2009]. Unul din sistemele cele mai performante este dezvoltat de laboratoarele LIUM [S. Meignier, T. Merlin, "LIUM SpkDiarization: An Open Source Toolkit For Diarization," in Proc. CMU SPUD Workshop, 2010.] Acesta procesează semnalul vocal în următoarele etape:

1. Extragerea vectorilor de caracteristici audio din semnalul audio
2. Segmentare
3. Agregarea segmentelor în grupuri
4. Antrenarea modelelor audio
5. Re-segmentarea Viterbi, folosind toate modelele identificate
6. Identificarea segmentelor care conțin vorbire
7. Identificarea bărbat/femeie, folosind modele pre-antrenate
8. Agregarea finală în grupuri a segmentelor care aparțin aceluiași vorbitor.

Așa cum se prezintă sistemele de diarizare cunoscute, dezavantajele lor sunt inabilitatea de a segmenta semnalul audio pe măsură ce acesta este primit de către

sistemul de diarizare. Sistemele cunoscute de diarizare nu pot determina caracteristicile segmentului (vorbitor, zgomot de fundal, etc) dacă nu au la dispoziție toate eșantioanele segmentului, astfel încât să poată antrena modelele audio necesare. Mai mult, sistemul de diarizare nu poate determina granița între două segmente dacă nu are la dispoziție toate eșantioanele pentru ambele segmente. În cazul ideal, diarizarea se face pe tot semnalul vocal, eliminând problemele enunțate, dar această metodă introduce intarzieri mari, inacceptabile pentru unele aplicații.

De exemplu, rezultatele diarizării sunt utile procesului de recunoaștere a vorbirii prin filtrarea vorbire/liniște sau informația despre vorbitor (pentru adaptarea modelului acustic), și este de dorit ca recunoașterea să aibă loc abia după terminarea diarizării. În multe aplicații, de exemplu cele care implică fluxuri audio foarte lungi, sau care necesită răspuns în timp real al sistemului de recunoașterea vorbirii, diarizarea pe tot semnalul vocal nu este o soluție viabilă și este necesară o soluție de diarizare în timp real, care să segmenteze semnalul audio pe măsură ce acesta este primit.

Invenția se referă la o metodă de diarizare în timp real a semnalelor vocale, care se realizează folosind, atât modele statice pre-antrenate de vorbitor, cât și o istorie de segmente și modele de vorbitor create dinamic, istoria fiind gestionată periodic prin actualizarea modelelor de vorbitor și prin eliminarea modelelor, folosind o funcție de cost bazată pe vechimea modelului și ponderea sa în istoria de diarizare, atunci când istoria de diarizare depășește o dimensiune prestabilită. Invenția se mai referă și la un sistem de diarizare în timp real, pentru implementarea metodei, care constă din module funcționale ce pot fi implementate ca programe software executabile pe un calculator sau circuite integrate, care mențin istoria de diarizare și comunică cu servicii externe de extragere a vectorilor de caracteristici audio, recunoaștere automată a vorbirii, actualizare a modelelor GMM.

Se prezintă în continuare, în detaliu, principiile și realizarea invenției, în legătură și cu figurile de la 1 la 3, care reprezintă:

Fig.1 prezintă metoda de diarizare propusă, ce presupune folosirea unei istorii de diarizare, conținând modele dinamice de vorbitor, și modele statice de vorbitor, pentru segmentarea unui flux de caracteristici de vorbire.

Fig. 2 prezintă metoda de gestiune a istoriei de diarizare, prin care segmente noi sunt adăugate și, atunci când istoria depășește o dimensiune dată, sunt

eliminate segmente și modele de vorbitor.

Fig. 3 prezintă un sistem de diarizare construit pe baza metodei propuse, constând din memorii pentru istoria de diarizare, și componente pentru gestiunea acestei istorii.

Invenția se referă la o metodă pentru separarea unui semnal audio în segmente omogene din punctul de vedere al proprietăților audio (diarizare), incluzând separarea segmentelor de vorbire de segmentele audio de liniște și, separarea segmentelor de vorbire în funcție de vorbitor. Metoda propusă are la bază faptul că segmentarea se face fără a aștepta primirea întregului fișier audio, iar segmentele rezultate dintr-o anumită porțiune a fluxului audio sunt calculate și livrate utilizatorului în timp real, cu o întârziere fixă, relativ la fluxul audio.

Metoda propusă pentru diarizare în timp real este prezentată în Figura 1. Metoda se bazează pe citirea, la fiecare T secunde, a unui număr N de vectori de caracteristici audio dintr-o memorie tampon 101. Vectorii sunt segmentați folosind un serviciu extern de segmentare 102. Modelele audio folosite pentru segmentare sunt stocate în istoria de diarizare 103 sau în memoria statică de modele pre-calculate 104. Istoria 103 acoperă ultimele S secunde de semnal audio, și conține, atât segmentele, cât și grupurile de segmente împreună cu modelele audio asociate grupurilor. Memoria statică 104 de modele pre-calculate conține modele pre-calculate pentru vorbitori considerați a fi importanți, a căror recunoaștere este esențială (persoane publice cunoscute, , VIP)

Segmentele rezultate din diarizarea de la pasul curent sunt adăugate la istoria de diarizare 103 și modelele vechi sunt eliminate din istoria 103 conform unei metode 105 de gestiune a istoriei de diarizare. Actualizarea modelelor de vorbitor din istoria 103, folosind informația audio de la pasul curent, este realizată de un serviciu extern. Vectorii audio corespunzători tuturor segmentelor, în afară de ultimul, sunt apoi transmise către utilizator. Vectorii audio din ultimul segment sunt păstrați și folosiți ca parte a următorului set de N vectori de caracteristici audio ce vor fi procesați.

Metoda propusă pentru gestiunea istoriei de diarizare este ilustrată în Figura 2. Istoria de diarizare este compusă din trei memorii distincte:

- memoria pentru segmente (MS) 201, ce conține caracteristicile de vorbire corespunzătoare segmentelor identificate anterior prin diarizare,
- memoria pentru grupuri de segmente (MG) 202, ce conține grupurile identificate anterior prin diarizare, și
- memoria pentru modele GMM (MM) 203, ce conține modelele de mixtură gaussiană calculate pentru fiecare grup de segmente în parte.

În urma diarizării unei ferestre audio, segmentele sunt adăugate istoriei de diarizare 103. Se încearcă asocierea fiecărui segment cu unul din modelele de vorbitor existente, la pasul 204. Dacă asocierea există, se actualizează grupurile de segmente la pasul 205 și se actualizează modelele de vorbitor ale grupurilor respective, la pasul 206. Dacă asocierea nu există, se crează un grup nou la care segmentul este adăugat, și se generează modelul de vorbitor pentru grupul nou creat. Atât noul grup, cât și modelul său de vorbitor, sunt adăugate istoriei de diarizare.

Dacă istoria de diarizare depășește o dimensiune  $D$  prestabilită de către utilizator, se execută o procedură de curățare a acesteia:

- Se verifică dacă există grupuri în care toate segmentele au vechime mai mare de  $S$  secunde, unde  $S$  este o valoare specificată de utilizator. Aceste grupuri sunt eliminate din istoria de diarizare la pasul 207.

Dacă pasul anterior nu a dus la scăderea dimensiunii istoriei de diarizare sub dimensiunea  $D$ , se calculează o funcție de cost la pasul 208 pentru fiecare grup/model, în felul următor:

- Se calculează o valoare de cost  $CV$  a vechimii modelului, proporțională cu numărul de secunde de la ultima actualizare a modelului
- Se calculează o valoare de cost  $CD$  a dimensiunii grupului asociat modelului, invers proporțională cu numărul de segmente conținute de grupul respectiv
- Se calculează o valoare de cost total  $CT$  prin mediere ponderată a  $CV$  și  $CD$ , cu ponderi alese de utilizator

Modelele sunt ordonate în funcție de valoarea de cost total CT asociată fiecăruia, în lista 209. În mod repetitiv, modelul cu valoarea cea mai mare de cost este eliminat din istoria de diarizare, împreună cu segmentele și grupul asociate modelului, până când dimensiunea istoriei de diarizare scade sub pragul D.

Sistemul propus pentru diarizare în timp real a unui flux audio este ilustrat în Figura 3. Istoria de diarizare este conținută în memorii RAM. Gestiunea istoriei de diarizare, incluzând scrierea și citirea memoriilor 201, 202, și 203 care formează istoria de diarizare, este realizată de un automat finit de control 301 ce poate fi implementat ca un circuit sau ca un program executat pe un microcontroler. Sistemul include o memorie RAM suplimentară pentru vectori audio, 302, care conține vectori de caracteristici audio citați dintr-o memorie tampon 303, atât timp cât este necesar pentru diarizare, vectori care apoi sunt scriși în memoria tampon 304. Întregul proces este controlat de un automat finit de control ce poate fi implementat ca un circuit sau ca un program executat pe un microcontroler sau un microprocesor.

Invenția prezentată are multiple avantaje față de stadiul tehnicii:

- Metoda propusă, folosind istoria de diarizare, permite execuția în timp real a procesului de diarizare, prin menținerea unui număr relativ mic de modele de vorbitor și actualizarea acestora pe măsură ce rezultatele diarizării sunt produse. Efortul computațional pentru recunoașterea vorbitorului prezent în fiecare segment de vorbire este proporțional cu numărul de modele de vorbitor avute în vedere, prin urmare reducerea numărului de modele reduce efortul computațional.
- Metoda propusă ocupă o cantitate mai mică de resurse de memorie, prin mecanismul de gestiune care elimină modelele de vorbitor, împreună cu segmentele asociate, folosind o funcție de cost ce ține cont de vechimea modelului și ponderea sa în istoria de diarizare. Datorită faptului că necesarul de memorie pentru metoda de diarizare propusă este fix, indiferent de lungimea fluxului audio diarizat și numărul de vorbitori din acest flux, iar dimensiunea efectivă poate fi setată arbitrar de mic, metoda propusă se pretează în special sistemelor încorporate și implementărilor folosind resurse limitate.

## Revendicări

1. Metodă de diarizare în timp real a semnalelor vocale prin identificarea și marcarea de segmente din fluxul audio vocal ce aparțin aceluiași vorbitor sau aceleiași clase audio, metoda fiind **caracterizată prin aceea că** diarizarea este realizată folosind, atât modele statice pre-antrenate de vorbitor, cât și o istorie de segmente și modele de vorbitor create dinamic, istoria fiind gestionată periodic atunci când istoria de diarizare depășește o dimensiune prestabilită, menținându-se astfel un necesar fix de memorie indiferent de lungimea fluxului audio și de numărul de vorbitori.
2. Metodă de diarizare conform revendicării 1, în care se realizează gestionarea istoriei de diarizare, ce conține modele statistice caracteristice pentru vorbitorii identificați, modele ce pot fi create sau actualizate în procesul de diarizare, metoda fiind **caracterizată prin aceea că**, atunci când istoria de diarizare depășește o dimensiune prestabilită  $D$ , sunt eliminate toate modelele statistice pentru vorbitorii care nu au mai apărut în fluxul audio de cel puțin  $S$  secunde.
3. Metodă de diarizare conform revendicării 1, **caracterizată prin aceea că**, atunci când istoria de diarizare depășește o dimensiune prestabilită  $D$ , se calculează o funcție de cost  $CT$  pentru modelele statistice dinamice de vorbitor și se elimină iterativ câte un model statistic din istoria de diarizare, alegându-se de fiecare dată modelul care are costul cel mai mare, până când dimensiunea istoriei de diarizare scade sub pragul  $D$ .
4. Metodă de diarizare conform revendicării 3, folosită pentru eliminarea selectivă a modelelor statistice de vorbitor din istoria de diarizare, metoda fiind **caracterizată prin aceea că** funcția de cost  $CT$  reprezintă o medie ponderată între costul  $CV$  dat de vechimea modelului în istoria de diarizare, unde costul  $CV$  este direct proporțional cu vechimea modelului, și costul  $CD$  dat de ponderea segmentelor asociate modelului în istoria de diarizare, unde costul  $CD$  este invers proporțional cu această pondere.



5. Sistem de diarizare în timp real ce implementează metoda conform revendicării 1, **caracterizat prin aceea că** sistemul constă din memorii și automate finite de control, ce pot fi implementate ca programe software executabile pe un calculator sau circuite integrate, care mențin istoria de diarizare și comunică cu servicii externe de extragere a vectorilor de caracteristici audio, recunoaștere automată a vorbirii, actualizare a modelelor GMM și diarizare preliminară.

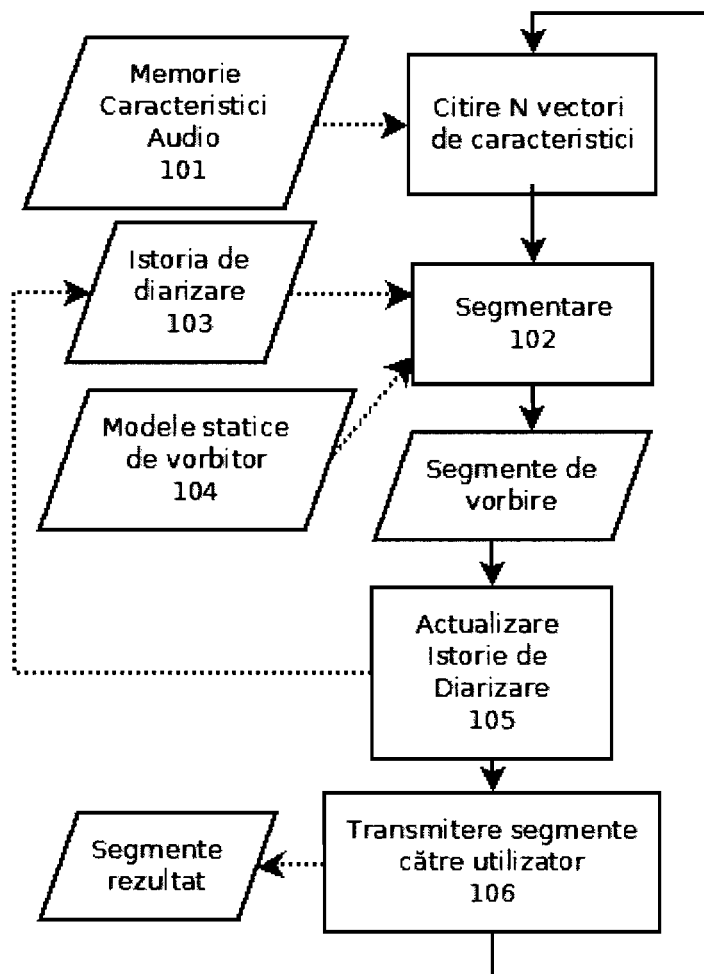


Figura 1: Metodă de diarizare

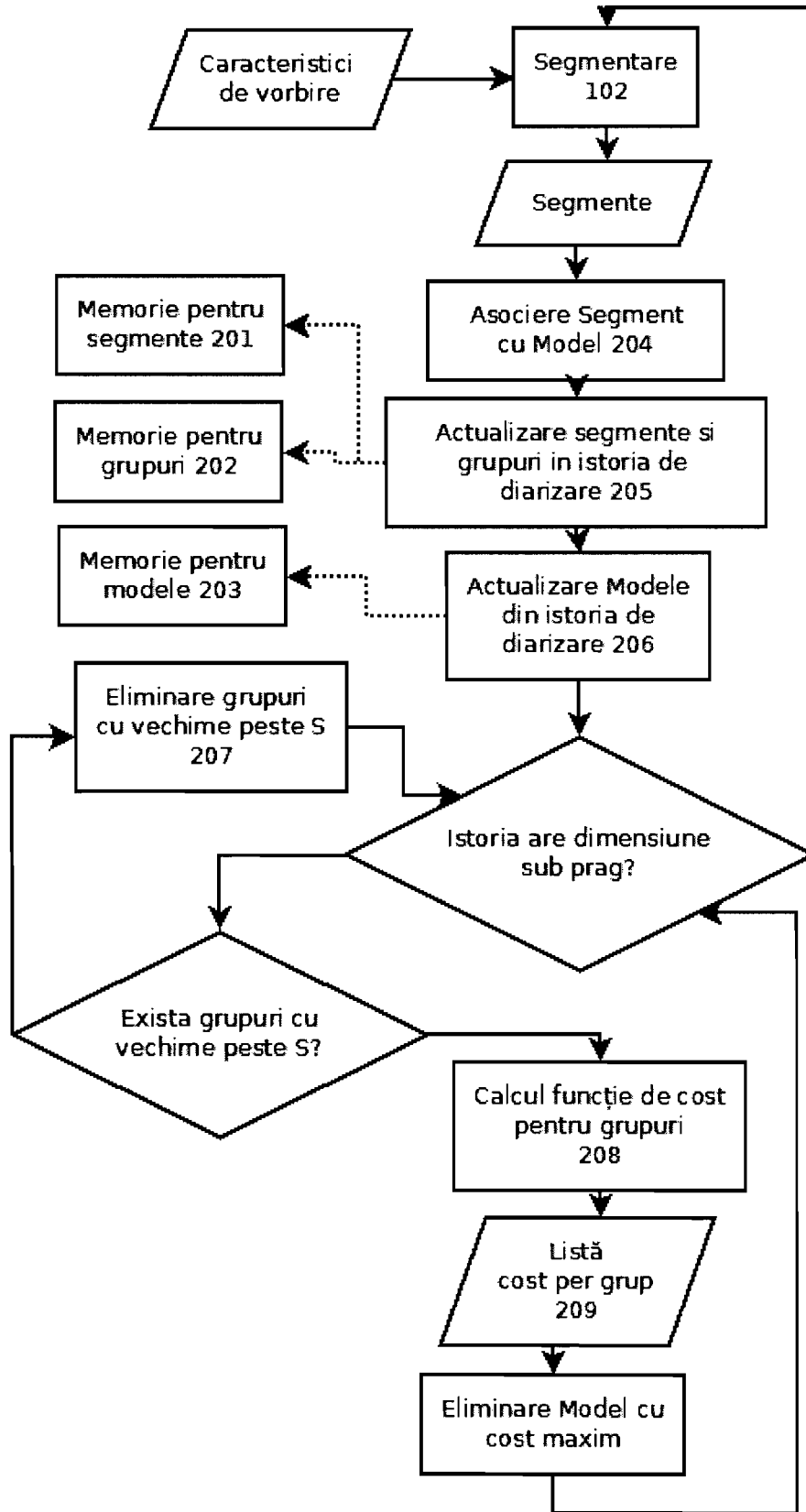


Figura 2: Gestiunea istoriei de diarizare

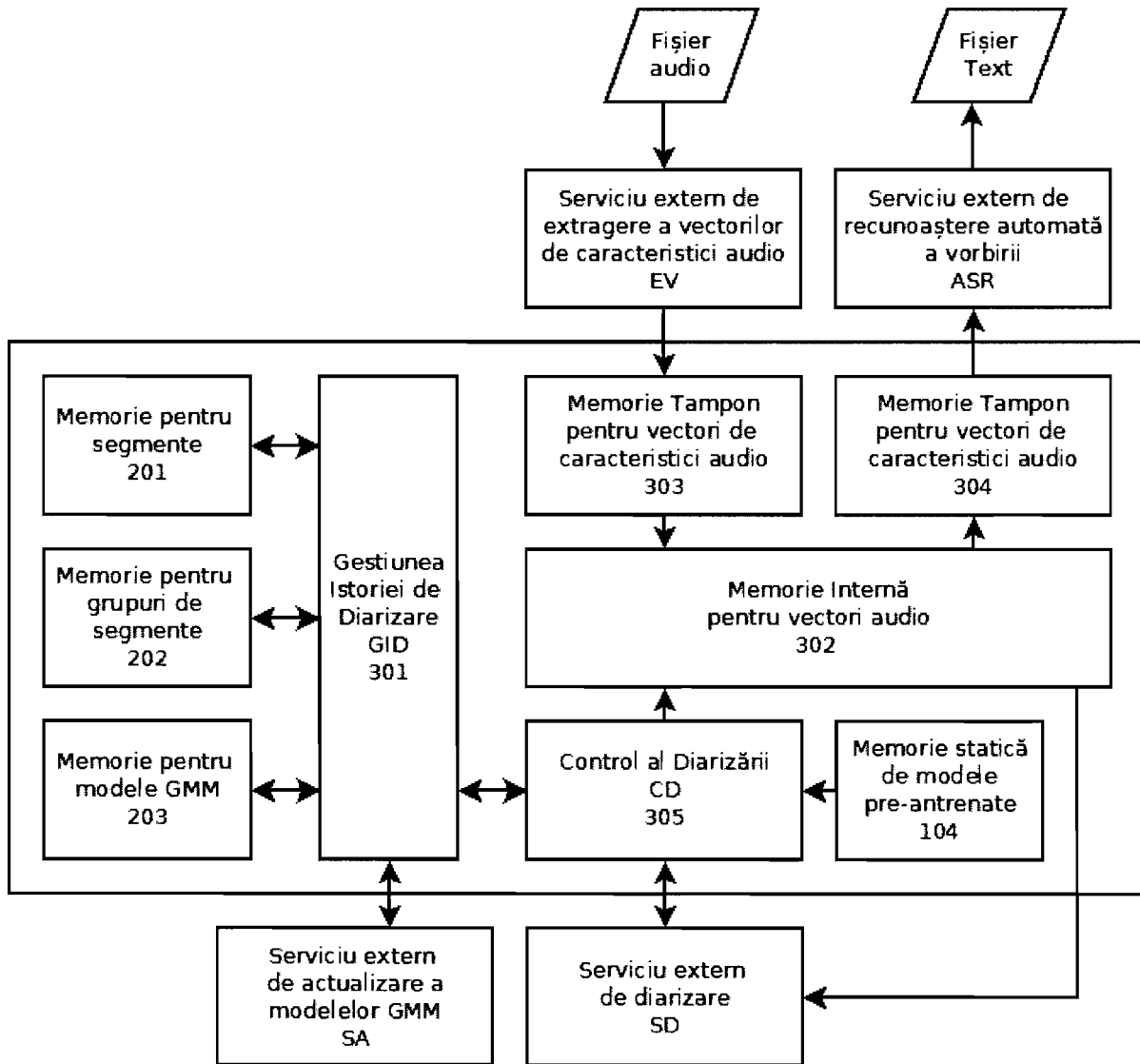


Figura 3: Sistem de diarizare în timp real