

# Premise pentru o tehnologie de recunoaștere automată a vorbirii în limba română aplicată domeniului juridic

Melania Duma<sup>1</sup>, Corina Giurcea<sup>2</sup>, Mihaela Ordean<sup>2</sup>, Paula Zălhan<sup>2</sup>

<sup>1</sup>Universitatea Babeș-Bolyai din Cluj-Napoca – UBB  
Str. M. Kogălniceanu nr. 1, Cluj-Napoca, 060042  
E-mail: [duma\\_melania@yahoo.com](mailto:duma_melania@yahoo.com)

<sup>2</sup>iQuest Technologies Cluj-Napoca – iQuest  
Str. Motilor nr.4-6, Cluj-Napoca, 40001  
E-mail: [corina.giurcea@iquestgroup.com](mailto:corina.giurcea@iquestgroup.com); [mihaela.ordean@iquestgroup.com](mailto:mihaela.ordean@iquestgroup.com);  
[paula.zalhan@iquestgroup.com](mailto:paula.zalhan@iquestgroup.com);

**Rezumat.** Gândit pentru a veni în întâmpinarea unui obiectiv din strategia IT a Ministerului Român de Justiție, proiectul JustASR, în curs de derulare, urmărește concepția și proiectarea elementelor tehnice necesare implementării unui sistem de interacțiune om-calculator bazat pe detectarea automată a limbajului pentru limba română având un scop restricționat la un domeniu particular. Proiectul își propune să dezvolte nucleul unei tehnologii de recunoaștere automată a vorbirii cu aplicabilitate directă în sălile de judecată ale României. Articolul de față se dorește a fi o prezentare succintă a premiselor și a unor elemente lingvistice de natură fonetică care țin de concepția proiectului. După o introducere în problematica tehnologiei de recunoaștere automată a vorbirii, articolul conține o trecere în revistă a eforturilor naționale din domeniul recunoașterii vorbirii. Sunt prezentate în continuare provocările specifice limbii române care constituie punctul de pornire al efortului științific și tehnic al tehnologiei care face obiectul proiectului JustASR. Cea mai consistentă secțiune a articolului, conținând partea aplicată a principiilor teoretice, este cea dedicată aspectelor fonetice care țin de realizarea dicționarului care va fi utilizat în proiect. Sunt prezentate criteriile generale utilizate pentru selecția cuvintelor care vor fi incluse în dicționar, precum și categoriile de cuvinte din care s-a făcut extragerea.

**Cuvinte cheie:** recunoașterea vorbirii, corpus, dicționar, procese fonetice.

## 1.Introducere

Vorbirea, cea mai naturală formă de comunicare între oameni, a reprezentat încă de acum 80 de ani o provocare în proiectarea și construirea unei mașini

care să mimeze comportamentul uman sub aspectul capacității de a vorbi natural și de a răspunde corect limbii vorbite conform Rabiner & Huang (2005).

Un interes major pentru studiul interacțiunii om-calculator prin intermediul limbajului a manifestat și manifestă nu doar comunitatea științifică, ci și mediul de afaceri. În condițiile fuziunii dintre sistemele de calcul și cele de comunicații, tehnologiile de procesare a limbajului natural capătă o importanță deosebită.

Mai mult, astăzi, aceste tehnologii sunt chemate să soluționeze problemele societății contemporane. Astfel, o serie de evenimente de tipul intervențiilor armate sau al calamităților naturale din ultimii ani au reliefat greutatea misiunilor umanitare sau a celor de salvare datorită necunoașterii de către participanții la aceste misiuni a limbilor vorbite de populația afectată.

Pe de altă parte, se estimează că, până la sfârșitul acestui secol, peste o jumătate din limbile existente vor dispărea (Besacier et al, 2014), existând un risc de diminuare a diversității lingvistico-culturale a omenirii.

Tehnologiile de procesare a limbajului uman pot și trebuie să aducă un aport major în soluționarea problemelor creionate mai sus prin: îmbunătățirea utilităților de procesare a textului, realizarea de dicționare electronice, continua dezvoltare de sisteme avansate de procesare de tip TTS (eng. *Text to Speech*) sau SR (eng. *Speech Recognition*). Spre exemplificare, sistemele de traducere automată pot constitui o alternativă la translatorii umani, insuficienți la nivelul impus de evenimente de tipul celor amintite mai sus în timp ce sistemele de recunoaștere automată a vorbirii ASR (eng. *Automatic Speech Recognition*) pot contribui la salvarea unor limbi amenințate cu dispariția (cele mai multe vorbite și nu scrise) prin transcrierea acestora.

În mare, tehnologia vorbirii în cadrul interacțiunii om-calculator presupune una sau ambele din următoarele: (i) o tehnologie care sintetizează automat vorbirea având ca intrare un text și (ii) o tehnologie care recunoaște automat limbajul. Recunoașterea automată a vorbirii are ca scop determinarea, cu o cât mai mare exactitate a șirului de cuvinte pronunțat de un vorbitor și este utilă în aplicații de tip comandă-și-control, dictare automată.

Cele două tehnologii prezentate mai sus stau la baza dezvoltării de sisteme avansate de procesare de tip TTS respectiv ASR, integrate într-o măsură tot mai mare în produse comerciale (în special în produse software)

destinate unor domenii precum telecomunicațiile, centrele de informare (call-center-e) automatizate, învățarea electronică (e-learning), sau medicină și justiție.

Tehnologia ASR este matură pentru limbile răspândite în lume: engleza, chineza, spaniola sau franceza. Marea majoritate a limbilor din țări în curs de dezvoltare au fost, până nu de mult, sub un con de umbră. Mai mult, studiile din cadrul unor proiecte ample în care au fost implicate rețele de excelență și/sau asociații profesionale (Besacier et al, 2014) au reliefat faptul că o serie de limbi cu un puternic potențial economic, multe dintre ele aflate în topul celor mai vorbite 20 de limbi ale lumii, intră în categoria așa numitelor limbi cu deficit de resurse (Cucu et al, 2014).

Termenul de **limbă cu deficit de resurse** (eng. *under-resourced language*) a fost introdus de Krauwer în anul 2003 și de Berment în 2004 (Besacier et al, 2014) pentru a desemna o limbă caracterizată prin unele, dacă nu prin toate aspectele care urmează: absența unui sistem de scriere unic sau cu ortografie aplicată consistent, prezența limitată pe web, lipsa de expertiză lingvistică, lipsa de resurse pentru procesarea limbajului natural (corpusuri monolingve, dicționare electronice bilingve, dicționare pentru pronunție, corpusuri adnotate fonetic/sintactic).

Cercetările efectuate în acest domeniu trebuie să răspundă la două mari provocări. Prima dintre acestea vizează rezolvarea unor probleme specifice limbii cu resurse limitate, probleme care nu au apărut în cazul limbilor cu resurse lingvistice dezvoltate. De exemplu, sistemele fonetice asociate unor limbi cu resurse limitate și dificultatea de a crea dicționare pentru pronunție constituie una dintre problemele care trebuie rezolvată în vederea dezvoltării de sisteme ASR sau TTS. A doua mare provocare o constituie optimizarea costurilor atașate dezvoltării unor sisteme de procesare a limbajului, costuri, de multe ori, mari.

În ultimii 7 ani, echipele de cercetare multidisciplinare s-au focalizat pe crearea de resurse lingvistice scrise și acustice, pe portarea tehnologiilor, pe dezvoltarea de algoritmi și metode pentru adaptarea modelelor acustice și lingvistice pentru limbi cu deficit de resurse. Începând cu 2008 se organizează bianual un workshop internațional, *Spoken Language Technologies for Under-resourced languages*, iar revista *Speech Communication* a avut în 2014 un număr special (56/2014) dedicat procesării limbilor cu resurse limitate (Besacier et al, 2014 b).

Cu toate că este limbă oficială a Uniunii Europene, limba română este considerată ca făcând parte din categoria limbilor cu resurse limitate. Încadrarea în această categorie s-a făcut pe baza rezultatelor studiilor efectuate în cadrul rețelei de excelență METANET. În cadrul acestei rețele reunind 60 de centre de cercetare din peste 34 de țări, s-au întocmit, pe baza unor criterii unitare, o serie de studii care analizează resursele și tehnologiile lingvistice disponibile pentru 30 de limbi europene și prezintă situația sprijinului acordat tehnologiilor limbajului. Rezultatelor studiilor privind limba română le-a fost dedicată o monografie ce se constituie într-un volum distinct al seriei (Trandabăț et al, 2012).

După știința noastră, până la această oră nu există pe piața românească o soluție comercială de recunoaștere automată a vorbirii dezvoltată în România. Trebuie menționat aici ca există pe piață implementări care utilizează recunoașterea vorbirii fără ca acestea să menționeze însă aplicațiile comerciale sursă utilizate.

Când vine vorba de limba română, majoritatea companiilor care dezvoltă produse comerciale care integrează o soluție de recunoaștere automată a vorbirii cer o bază de date cu înregistrări pe baza cărora implementările sunt personalizate. Un exemplu de companie care a lansat pe piață o aplicație comercială de recunoaștere a vorbirii în limba română este Nuance Communications, Inc. Conform unui comunicat de presă postat pe situl Nuance (<http://www.nuance.com>), compania a lansat în anul 2012 aplicația comercială Dragon Dictation pentru iPhone, iPad și iPod. Aplicația utilizează recunoașterea vorbirii în vederea dictării rapide a mesajelor de tip SMS și poșta electronică, a informației care se dorește a fi actualizată pe rețelele sociale sau obținută prin căutare pe web.

În anul 2010, în vederea îndeplinirii condițiilor legale și a recomandărilor Comisiei Europene privind transparența actului public de justiție, Ministerul de Justiție din România a organizat o licitație publică urmărind achiziționarea unei soluții de recunoaștere a vorbirii pentru sălile de judecată, oferind peste 2.5 Mil EUR. Pe piața din România, nici un furnizor de software nu a fost capabil să livreze un astfel de produs. În urma acestui „eșec” în strategia IT a ministerului pentru perioada 2013-2017, publicată în aprilie 2013, Ministerul de Justiție reiterează interesul pentru un produs care să realizeze transcrierea automată a discuțiilor din sălile de judecată.

În acest context, în dorința de a răspunde provocării și totodată oportunității comerciale descrise mai sus, a fost gândit proiectul JustASR. Proiectul, aflat în curs de derulare în cadrul Programului National de

Cercetare – Parteneriate în Domenii Prioritare, are ca obiective concepția și proiectarea elementelor tehnice necesare implementării unui sistem de interacțiune om-calculator bazat pe detectarea automată a limbajului pentru limba română, având un scop restricționat la un domeniu particular: realizarea unei transcrieri automate, de calitate, pentru audierile publice din sălile de judecată din România.

Articolul de față nu își propune să intre în aspectele tehnice legate de dezvoltarea nucleului tehnologiei de recunoaștere a vorbirii în limba română aplicată domeniului juridic, care face obiectul proiectului JustASR. Aceste aspecte vor constitui subiectul unei alte publicații propuse de Sofroni & Stan (2015).

Prezentul articol investighează o serie de aspecte practice, de natură fonetică, dar necesare în realizarea unui dicționar, respectiv a unui corpus de antrenament pentru un sistem de recunoaștere automată a vorbirii pentru limba română.

## **2. Sinoptic privind cercetările efectuate la nivel național în vederea dezvoltării unui sistem ASR pentru limba română**

Deși, la această oră, nu este încă disponibil pe piața românească un produs dezvoltat în România care să integreze tehnologia de recunoaștere a vorbirii pentru limba română, demersuri în acest sens au fost întreprinse de mai multe grupuri de cercetare din România, cu rezultate demne de luat în calcul.

În teza lui de doctorat, Boldea (2003) a colectat, pentru prima dată, o bază de date fonetice, acoperind direct sau prin combinații, fonemele limbii literare române. El a dezvoltat un algoritm pentru identificarea limitelor unităților de semnal și a efectuat adnotarea majorității semnalelor înregistrate. Se cuvine a preciza aici că fonemul reprezintă unitatea minimală a nivelului fonetic ce individualizează un sunet articulat și care e capabilă să diferențieze sensul.

La Universitatea Tehnică din Cluj-Napoca, grupul de cercetare coordonat de către Prof. G. Todorean a dezvoltat un sistem bazat pe Modele Markov Ascunse Discrete (eng. *Hidden Markov Models*) pentru recunoașterea cuvintelor izolate descris în Todorean & Buza (2008). În teza sa de doctorat, elaborată sub îndrumarea Prof. G. Todorean, Margit Antal (2006) a

investigat recunoșterea vorbirii și a vorbitorului pe baza fonemelor și a propus modele pentru recunoșterea și clasificarea fonemelor.

Grupul Prof. H. N. Teodorescu de la Universitatea Tehnică Iași, a dezvoltat „Sunetele Limbii Române” (SroL) conținând 1500 de înregistrări distincte cu diferite formate de precizie și codificări descrise în (Teodorescu et al, 2010). În privința SroL, lucrarea Teodorescu (2010) furnizează exemple despre modul în care se pot extrage elemente relevante din semnale destinate recunoașterii limbajului natural.

O largă contribuție în procesarea limbajului natural este oferită la Universitatea Politehnică București sub îndrumarea Prof. Corneliu Burileanu, printr-o muncă cumulată de-a lungul anilor. Lucrarea (Popescu et al, 2008) prezintă o metodă computațională simplă, pentru urmărirea vorbitorului în dialogurile multi-laterale, bazată pe traversarea unui arbore decizional construit în funcție de o regresie liniară neghidată de tip Maximum Likelihood. Rezultatele arată o performanță de 80% în procesul de urmărire a vorbitorului.

Menționăm aici ca pentru evaluarea performanței unui sistem ASR trebuie calculată rata de eroare (word-error-rate-WER), obținută prin compararea textului vorbit cu cel recunoscut de către software. Această rată de eroare este calculată ca număr total de substituiri, inserări și ștergeri de cuvinte necesare pentru a face cele două texte identice, raportat la numărul total de cuvinte. Acuratețea procesului ASR este definită ca opusul ratei de eroare, adică Acuratețea = 1-WER.

Grupul condus de Prof. Burileanu a început să avanseze către problematica recunoașterii limbajului în vorbirea spontană în limba română în cadrul proiectului PNII IDEI nr. 114/2007. Numeroase articole au fost publicate în acest sens (Petrea et al, 2009), (Petrea et al, 2010). Utilizând Hidden Markov Modelling Toolkit (HTK) (HTK toolkit <http://htk.eng.cam.ac.uk/>), experimentele preliminare au arătat o performanță de 70% acuratețe (Burileanu et al, 2010). Prin redefinirea metodei și modelarea incertitudinii, prin detectarea de Modele Gaussiene Mixte (eng. *Gaussian Mixture Models*) și introducerea log-probabilităților, ei au reușit să înregistreze o performanță de 18.5% WER, utilizând un corpus de 10000 de cuvinte distincte, pronunțate de numai 5 vorbitori.

Într-un articol recent (Cucu et al, 2014) se raportează construirea, de către echipa din cadrul laboratorului de cercetare Speed din Universitatea Politehnică din București (UPB), a celui mai mare vocabular al limbii romane (peste 60000 cuvinte) destinat dezvoltării unui sistem ASR.

Articolul citat este centrat pe prezentarea metodologiei propuse de autori pentru portarea/adaptarea unui corpus asociat unui domeniu specific dintr-o limbă cu resurse dezvoltate (în cazul de față, un corpus asociat jurnalismului de știri în limba franceză) într-o limbă cu resurse limitate (în cazul de față limba română), utilizând metode de traducere automată bazate pe statistica oferită de Mașini Statistice de Traducere sau SMT (eng. *Statistical Machine Translation*). În același articol, modelele pentru restaurarea diacriticelor limbii române, respectiv cel pentru conversia grafem-fonem, propuse, dezvoltate și testate de autori, sunt evaluate și comparate cu alte modele elaborate pentru limba română.

Studiul lui Schiopu (2010) propune o comparație între sistemele ASR dezvoltate la Cluj și la Iași (prezentate mai sus), dar nu ia în considerare dezvoltările Prof. C. Burileanu.

Grupul condus de către Acad. Dan Tufiș de la Academia Română a desfășurat o muncă de cercetare în direcția dezvoltării resurselor computaționale lingvistice pentru limba română. Merită menționate în acest sens: Proiectul European CLARIN (proiectul stabilește o infrastructură de cercetare formată din resurse și tehnologii în scopul diminuării fragmentării curente și oferirea unei infrastructuri stabile, persistente, accesibile și extensibile), proiectul eDTLR realizat în cooperare cu grupul Prof. Dan Cristea de la Universitatea Al. Ioan Cuza din Iași – acest proiect a condus procesul de digitizare a Dicționarului Tezaur al limbii române eDTLR (2013) sau proiectul SEE – ERA.NET (construirea resurselor lingvistice și a unor modele de traducere automatizată cu ajutorul calculatorului în limbile balcanice și de origine sud-slavonă).

Toate aceste resurse merită luate în considerare, pentru că ele au reușit să clarifice aspecte controversate și variate ale modelării limbii române, reprezentând o bază reală pentru construirea unei tehnologii comerciale pentru recunoașterea automată a limbajului natural. Ca o remarcă, se observă că aceste grupuri de cercetare au muncit izolate unele față de celelalte, fiecare utilizând propria bază de date cu semnale de voce înregistrate.

Până foarte recent, nici unul dintre grupurile de cercetare amintite mai sus, nu a reușit să înregistreze și să proceseze semnale vocale provenind de la un număr mare de vorbitori care utilizează cuvinte extras dintr-un vocabular de scară medie. Proiectul RASC - Romanian Anonymus Speech Corpus, derulat în cadrul Institutului de Cercetări pentru Inteligență

Artificială al Academiei Române, își propune ca obiectiv principal tocmai crearea unei colecții mari de resurse care să poată fi utilizate în aplicații care implică prelucrarea limbajului vorbit. Colectarea corpusului de voce se face într-un mod inedit, prin „crowd-sourcing” conform (Dumitrescu et al, 2014).

Portalul RASC (<http://rasc.racai.ro/>) a fost deschis tuturor celor care doresc să se înregistreze citind propoziții dintr-un set prestabilit, asigurându-se astfel atât variabilitatea condițiilor de înregistrare reflectate în calitatea semnalului obținut (medii cu sau fără zgomot de fundal, înregistrări cu sau fără microfon) cât și variabilitatea în vorbire (înregistrări provenind de la persoane de sex diferit, de vârste diferite, provenind din diversele regiuni ale țării). Se cunoaște faptul că fiecare persoană are un tract vocal diferit și că un vorbitor al unei limbi având caracteristici generale clare poate rosti una și aceeași propoziție cu mici variații. Toate aceste variații se vor transpune și pot fi observate în reprezentările digitale ale semnalelor de vorbire. La ora redactării acestui articol, conform informației disponibile pe site-ul proiectului, corpusul conținea 3056 de propoziții înregistrate, avea o mărime 1060 MB și era pus la dispoziția comunității de cercetare (ca open-source) în vederea dezvoltării unor sisteme de recunoaștere a vorbirii continue în limba română.

### **3. Aspecte lingvistice ale unui proiect de recunoaștere a vorbirii în limba română**

Particularitățile gramaticale distinctive ale limbii române în cadrul limbilor romanice și deficitul de resurse în ceea ce privește utilizarea românei în contexte informatice determină o dublă provocare pentru orice demers ce presupune digitalizare și aplicare de algoritmi.

Dezvoltarea unei tehnologii de recunoașterea automată a vorbirii pentru limba română necesită, în primul rând, o înțelegere a particularităților fonetice și morfologice ale limbii, iar, în al doilea rând o modelizare a regulilor sintactice derivate din acestea în funcții.

#### **3.1. Specificul limbii române**

La nivel sintactic limba română se individualizează în cadrul limbilor romanice prin: nelexicalizarea subiectului prenominal, dublarea de către un



pronume personal a unui grup nominal lexical, dublarea cliticelor, concordanța negativă și negația dublă.

La nivel morfologic, limba română se evidențiază printr-un sistem flexionar bogat, în care substantivele, adjectivele și pronumele au cinci cazuri și două numere. Verbele au trei persoane atât la singular, cât și la plural pentru cele cinci moduri personale, la care se adaugă modurile nepersonale: infinitivul, participiul și gerunziul, în total un verb putând număra treizeci și trei de forme. Adjectivele și substantivele pot avea forme articulate și nearticulate, un substantiv putând avea în total cinci forme, iar un adjectiv șase.

La nivel fonetic, flexiunile cuvintelor limbii române pot să aibă alternanțe în interiorul rădăcinii, la care se adaugă sufixe morfologice și desinențe (Trandabăț et al, 2012).

O caracteristică a limbii române cu impact fonetic pentru ASR este folosirea diacriticelor: *ă*, *â*, *î*, *ș* și *ț*, care apar în aproape 30-40% dintre cuvintele dintr-un text (Cucu et al, 2014).

În cazul în care diacriticele lipsesc, mesajul transmis s-ar putea să fie neclar, aceasta determinând atât lizibilitate scăzută, cât și ambiguitate. Deși pe Internet sunt disponibile pentru descărcare diferite corpusuri românești de știri fără diacritice, acestea nu pot fi utilizate în dezvoltarea unui sistem ASR pentru limba română (Cucu et al, 2014). O adnotare ASR trebuie să conțină diacritice și să elimine ambiguitatea cuvintelor care pot avea forme distincte cu și fără diacritice.

Pentru transcrierea fonetică automată a unui vocabular, există în momentul de față două abordări posibile. Primul tip de abordări se bazează pe reguli de conversie grafem-fonem. Pentru aceste abordări contează numărul de reguli necesare în dezvoltarea sistemului ASR. Dacă pentru limba engleză, numărul acestora este foarte mare (în jur de 1500 de reguli), pentru limba română (care este o limba fonetică) numărul regulilor poate fi considerabil mai mic. Al doilea tip de abordări sunt cele bazate pe mașini de învățare. Sistemele care folosesc mașini de învățare se bazează, în principal, pe cunoștințe dobândite dintr-un set de învățare, care sunt mai apoi aplicate datelor din aplicația ASR. Aceste cunoștințe sunt de obicei exprimate prin arbori de decizie sau prin rețele neuronale.

### 3.2. Principii lingvistice în recunoașterea automată a vorbirii

Recunoașterea automată a vorbirii presupune interpretarea unui flux de intrare auditiv și transpunerea sa în text. Realizarea unei interpretări cât mai fidele este condiționată de câteva etape premergătoare:

- Realizarea unui corpus de antrenament și a unui corpus de testare;
- Realizarea unei parsări, adică a unei analize discrete a elementelor ce individualizează un anumit nivel de limbă.

Analiza lingvistică implicată în recunoașterea automată a vorbirii se concretizează în parsări de tip fonetic, morfologic și sintactic. Articolul de față prezintă o parte din criteriile fonetice de care trebuie să se țină cont în evidențierea interacțiunii dintre nivelele limbii pentru un sistem de recunoaștere automată a vorbirii pentru limba română.

Dicționarul nu este doar o înșirare de cuvinte, ci un set de cuvinte ordonat după criterii fonetice, morfologice sau sintactice, așa cum evidențiază și Cristea (2005).

Criteriul fonetic, conform cu DOOM (2005) are la bază interacțiunea dintre unitățile sale minimale, fonemele. Un dicționar realizat pe astfel de principii va include cuvintele și transcrierea lor fonetică. Selecția cuvintelor va avea în vedere acoperirea inventarului de foneme ale limbii române, mai precis a instanțelor acestor foneme ce poartă numele de alofone. Alofonele sunt rezultatul interacțiunii dintre sunet și contextul său diagnostic.

Criteriul morfologic, în conformitate cu GALR (2008), are la bază interacțiunea dintre unitățile sale minimale, morfemele. Morfemele sunt de două tipuri: lexicale (partea stabilă a vocabularului) și gramaticale (partea flectivă a vocabularului). Dicționarul va include partea lexicală a cuvântului și variațiile sale în funcție de context: poziționare sau regulă sintactică.

Criteriul sintactic, corespunzător cu GALR (2008), are la bază interacțiunea dintre grupuri sintactice, iar acest fapt nu generează un dicționar în sine, ci reglementează gradul de filtrare și de acces al intrărilor din dicționarul ce conține deja elemente de fonetică și de morfologie.

În practică dicționarul va conține cuvinte și un cumul de criterii, fonetice și morfologice. Regulile sintactice se adaugă în calitate de constrângeri ce acționează la nivel macro-sintagmatic.

### 3.3. Realizarea dicționarului și a corpusului

Pentru un proiect de recunoaștere a vorbirii, dicționarul include intrările

lexicale atent selectate după o serie de criterii. Cum proiectul de față se adresează serviciilor publice din zona juridică, selectarea contextelor reprezentative pentru acest domeniu va ține cont atât de jargonul de specialitate, cât și de cuvintele uzuale.

#### Criteriul frecvenței

Cuvintele trebuie selectate pe baza unui corpus de texte juridice (comunicate, sentințe, transcrieri din cadrul proceselor), pe baza corpusurilor realizate ad-hoc și a celor existente prezentate în Dascălu (2002), Dascălu & Pop (2003).

#### Criteriul tematic

Corpusul va acoperi texte ce țin de domeniile: apărare, condamnare, ordin, rugăminte, solicitare informații și expunere argumente. Realizarea unui dicționar pentru domeniul juridic trebuie să conțină două categorii de cuvinte:

a. cuvinte ce țin de limbajul comun, deci de o zonă invariabilă a vocabularului:

- *prepoziții*: cu Acc (cu, pe, peste, în, din, la etc.), cu G-D (împotriva, contra, asupra, conform, potrivit, contrar etc);
- *conjunții*: coordonatoare (dar, iar, însă, ci, și, sau, ori) și subordonatoare (pentru că, din moment ce, ca să, dacă etc.);
- *pronume și adverbe relative*: care, cine, ce, unde, când, cum, cât;
- *pronume*: personale (eu, tu, el), demonstrative (acesta, acela, același), nehotărâte (oricare, fiecare, altul), negative (nimeni, nimic);
- *nume*: zile, luni, ora, instituții/locații, nume persoane, grade de rudenie;
- *expresii de mod* (așa, repede, încet, corect, incorect, conform procedurii, contrar procedurii);
- *expresii de timp* (acum, atunci, cândva, acum <x> ore/zile/luni/ani, ieri, ora/z ziua/săptămâna/luna trecută, anul trecut, când ....., de <x> ore/zile/săptămâni/luni/ani, din <data/anul> până în <data/anul>);
- *expresii de loc* (aici, acolo, în+<loc>, la+<loc>).

b. cuvinte specializate:

- *verbe/expresii de stare*: a fi, a exista, a sta, a se mira, a se învecina, a se înrudi;

- *verbe/expresii de opinie*: a crede, a considera, a pune sub semnul întrebării, a contesta, a admite, a fi de acord, a nu fi de acord;
- *verbe/expresii de acțiune*: a lua, a da, a crea, a realiza, a evidenția, a convoca, a elibera, a pune (sub acuzare), a trimite, a primi, a fura, a comite o infracțiune, a defăima, a omorî, a vătăma, a răni, a trece, a supraveghea, a învăța, a repara, a vinde, a cumpăra, a se căsători, a divorța, a strânge, a elimina, a introduce, a facilita, a media, crește, a descrește, a îmbătrâni, a îngriji, a acorda atenție, a ezita, a minți, a spune adevărul, a aduce dovezi, a utiliza, a contesta;
- *verbe/expresii performative*: mă scuzați, mă/ne angajăm să, te/vă acuz, te/vă autorizez, te/vă anunț, te/vă avertizez, te/vă declar, te/vă felicit, te/vă insult, te/vă întreb, te/vă previn, te rog/vă rog, te/vă sfătuiesc, îți/vă interzic, îți/vă jur, îți/vă mulțumesc, îți/vă ordon, îți/vă promit, îți/vă recomand, îți/vă sugerez;
- *verbe/expresii de senzație*: i se urăște, i se face frică, i se face teamă, a se teme, i se face rău, îi place, îi pare bine, îi pasă, îl doare, îl uimește, îl îngrozește, a iubi, a urî;
- *verbe/expresii de aspect/modale*: a începe, a continua substantiv/de+Participiu, a termina substantiv/de+Participiu, a putea să, trebuie să, e nevoie de/să, a avea nevoie de/să, e posibil să, e obligatoriu să, e necesar să, e interzis să, e permis să, e recomandat să;
- *substantive de tip actanți*: inculpat, reclamant, pârât, petent, avocatul apărării, avocatul acuzării, magistrat, judecător, grefier, procuror, polițist, organele de poliție, martor;
- *substantive de tip acuzații*: corupție, înșelăciune, fals, uz de fals, crimă, defăimare, evaziune;
- *substantive de tip sancțiuni*: pedeapsa cu închisoare între <x> zile/luni/ani și <x> zile/luni/ani, plata retroactivă amendă, prestarea de muncă în folosul comunității, privare de libertate, arest - arestare – în stare de arest - arest la domiciliu;
- *substantive de tip nume locuri*: închisoare, penitenciar, tribunal;
- *substantive de tip elemente ale procesului juridic*: acuzare, dovedirea nevinovăției, chemarea în instanță a martorilor/pârâților, amânarea procesului/a sentinței.

Corpusul reprezintă un set de enunțuri înregistrate care au fost în prealabil:

- Selectate sau construite pentru a acoperi inventarul relevant fonetic și morfologic, din dicționar;
- Parsate în conformitate cu regulile fonetico-morfologice convenite de specialiști.

Cu cât enunțurile sunt rostite de mai multe persoane în diverse contexte de claritate, comportamentul elementelor gramaticale și interpretarea acestora se poate evidenția mai bine. Acuratețea și relevanța corpusului sunt dublate de gradul de individualizare a regulilor gramaticale pentru limba respectivă.

### **3.4. Nivelul fonetic. Generalități**

Pentru o analiză corectă și coerentă a nivelului fonetic, trebuie să se înțeleagă de la bun început o serie de aspecte fundamentale ale acestui nivel:

1. Sunetele ce alcătuiesc un cuvânt sunt individualizate prin intermediul fonemelor, adică prin intermediul reprezentării lingvistice asociate unui sunet.
2. Fonemul este o reprezentare ideală a sunetului. Ceea ce interesează este ocurența sunetului într-un context diagnostic, adică alofonul. Alofonul unui sunet este rezultatul influenței exercitate asupra acestuia de contextul precedent și/sau succedent.
3. Decelarea inventarului de alofone specifice unei limbi se realizează prin identificarea proceselor fonetice specifice acelei limbi.
4. Gradul variabil de sonoritate a fluxului vocal determină existența vocalelor și a semivocalelor.
5. Adiacența vocală-vocală se realizează doar în silabe diferite și poartă numele de hiat.
6. Adiacența vocală-semivocală (în diverse poziții) se realizează în aceeași silabă și poartă numele de diftong (o vocală+o semivocală) sau de triftong (o vocală+2 semivocale).

Din această prezentare schematică, se poate ușor observa necesitatea înțelegerii și a interpretării valorii sunetului: (i) în contextul adiacent sau diagnostic din care face parte; (ii) în contextul apartenenței sale la un tip de proces fonetic care îi modifică atributele; (iii) în contextul silabei din care face parte.

În continuarea acestui articol, vom face o scurtă prezentare a patru aspecte esențiale pentru o bună realizare a analizei fonetice pentru limba română și pentru sistemele ce o utilizează în prelucrare: problema diftongilor și a triftongilor, problema hiatului, probleme de omofonie și omografie, problema stabilirii unui inventar minimal de procese fonetice.

### 3.5. Diftongii și triftongii vs. hiatul în limba română

Elementele *a, ă, î, ă* sunt mereu vocale în limba română. Rezultă de aici că sunetele *e, i, o, u* vor putea fi fie vocale, fie semivocale în funcție de contextul din care fac parte.

Diftongul este asocierea unei vocale și a unei semivocale în aceeași silabă. În funcție de poziționarea semivocalei în raport cu vocala, rezultă:

- 14 diftongi descendenți:
  - Semivocala *i*: ai (e.g. hai), ai (e.g. a destăinui), îi/âi (e.g. pâine), oi (e.g. doi), ui (e.g. a uita), ei (e.g. idei), ii (e.g. fii);
  - Semivocala *u*: au (e.g. lucrau), ău (e.g. rău), îu/âu (e.g. grâu), ou (e.g. lingou), iu (e.g. viu), eu (e.g. meu), uu (e.g. perpetuu).
- 10 diftongi ascendenți:
  - Semivocala *i*: ia (e.g. iarbă), ie (e.g. iertare), io (e.g. niciodată), iu (e.g. iubire);
  - Semivocala *e*: ea (e.g. începea), eo (e.g. vreo);
  - Semivocala *u*: ua (e.g. ploua), uă (e.g. nouă), uî/uâ (e.g. plouând);
  - Semivocala *o*: oa (e.g. moară).

Triftongul este asocierea unei vocale și a două semivocale în aceeași silabă. În funcție de poziționarea semivocalelor în raport cu vocala, rezultă:

- 11 triftongi descendenți: eau (e.g. aveau), eai (e.g. primeai), iai (e.g. suiai), iau (e.g. miau), oai (e.g. leoaică), iei (e.g. miei), ieu (e.g. maieu), iou (e.g. maiou), uăi (e.g. rouăi), uai (e.g. înșeuai), uau (e.g. înșeuau);
- 3 diftongi ascendenți: eoa (e.g. learcă), ioa (e.g. lăcrămioară), uea (e.g. înșeuează).

În momentul în care ocurențele de mai sus, adică grupările de sunete de mai sus apar în silabe diferite acestea nu mai sunt interpretate ca diftongi sau triftongi, ci ca hiat. Hiatul este deci prezența a două vocale în silabe diferite. Hiatul implică și fenomenul fonetic de stop glotal cu apendice palatal prezentat mai jos în cadrul proceselor fonetice.

(1) Ați mai furat doi cai cu un an înainte.

Cuvântul *cai* conține diftongul descendent [ai], în timp ce *i-na-in-te* conține hiatul [a-<sup>l</sup>i] unde sunetul *i* este concretizat ca alofon cu stop glotal și apendice palatal. Modul în care gruparea *ai* se pronunță diferă, dar în transcrierea vorbirii, gruparea trebuie redată identic ca grafie.

Manifestarea grupurilor de vocale și semivocale ca hiat vs. diftong/triftong implică trăsături articulatorii diferite, cu implicații pentru scriere, motiv pentru care, pentru fiecare grupare, vor trebui identificate cazurile de omofonie și omografie. De multe ori, aceste situații sunt motivate sintactic sau semantic.

### 3.6. Probleme de omofonie și omografie

Fenomenul în care o singură grafie sau scriere are drept corespondent două sau mai multe pronunții diferite poartă numele de omografie. Invers, în situațiile în care o anumită grupare se pronunță identic, dar scrierea e diferită, vorbim de omofonie.

(2) Mai știi când m-ai lovit?

Cuvântul *mai* și expresia *m-ai* se pronunță identic. Cu toate acestea, grafia este diferită din motive ce țin de morfologie și sintaxă: *mai* este adverb, pe când *m-ai* este o combinație de pronume personal în Acc și auxiliar verbal. Se poate observa astfel motivul pentru care dicționarul trebuie să conțină elemente din toate palierele limbii.

Oferim în cele ce urmează câteva exemple în care omofonia, omografie și schimbarea de sens conlucrează, fapt ce are drept consecință o complexitate sporită a regulilor lingvistice la nivel fonetic.

(3) Grup identic - Pronunție identică – Scriere diferită:

(3.1.) sau vs. s-au

Sau inculpații au adus noi dovezi, *sau* avocații lor s-au schimbat.

(3.2.) mai vs. m-ai

În luna *mai* m-ai contactat să mai povestim despre afaceri.

(3.3.) cai vs. c-ai

C-ai văzut, că n-ai văzut niciunul din cei doi *cai* nu mă privește.

(3.4.) nai vs. n-ai

N-ai cumva un *nai*?

(3.5.) iar vs. i-ar

I-ar oferi o nouă șansă *iar*.

- (3.6.) Iași vs. i-aș  
I-aș fi dat documentele la *Iași* dacă aș fi avut suficient timp.
- (3.7) cei vs. ce-i  
Ce-i poți oferi când *cei* mai bogați nu au soluții?
- (3.8) Mia-i vs. mia-i vs. mi-ai  
 Tot asta mi-ai spus: „*Mia-i* de încredere, iar mia-i la mine.”
- (4) Grup identic - Pronunție diferită – Scriere diferită:  
 ie vs. i-e  
 Cu această *ie* i-e rece în celulă.
- (5) Grup identic – Pronunție diferită – Scriere identică:
- (5.1.) Diftongul ascendent ii vs. Diftongul descendent iu vs. Hiatal iu  
*iulie* vs. vișiniu vs. el știu  
*iunie* vs. eu știu vs. el știu
- (5.2.) Diftongul descendent ai vs. Hiatal ai  
*haină* (substantiv) vs. haină (adjectiv)
- (5.3.) Diftongul descendent ai vs. Hiatal ai  
 eu mă *bâlbâi* (acum) vs. el se bâlbâi (ieri)
- (5.4.) Diftongul descendent ai (context anterior: consoană) vs. diftongul descendent ai (context anterior: vocală)  
*cai* vs. știi
- (5.5.) Pronunția hiatului ca *diftong* în funcție de subiectivitatea vorbitorului (rapiditate în vorbire)  
 a-gre-si-u-ne vs. a-gre-*siu*-ne  
ști-ut vs. *știut*  
 glo-ri-e vs. glo-*rie*  
 a-bre-vi-e-re vs. a-bre-*vie*-re

### 3.7. Procese fonetice

Procesele fonetice se referă la transformările pe care un fonem le suportă în contextul său diagnostic. Rezultatul proceselor fonetice vor fi alofonele, adică instanțele clasei fonem. Procesele fonetice ale limbii române au fost studiate și documentate de Rosetti (1963), Stati (1972) și Neamțu (2002).

#### Palatalizarea și labializarea

Palatalizarea și labializarea acționează doar în cadrul unei silabe, iar fenomenele presupun influența vocalelor și/sau a semivocalelor *e* și *i* (cazul palatalizării), respectiv a vocalelor și/sau a semivocalelor *o* și *u* (cazul



labializării) asupra consoanei care le precedă și asupra vocalei care le succedă.

(6) *via-ță* – alofonele palatalizate sunt [v] și [a] sub influența lui *i*.

(7) *no-uă* – alofonele labializate sunt [n] sub influența lui *o* și [ă] sub influența lui *u*.

Palatalo-labializarea reprezintă influența semivocalelor *eo*, *io* și a diftongilor ascendenți *io*, *eo*, *iu* asupra consoanei care îi precedă și/sau asupra vocalei care le succedă.

(8) *pleoa-pă* – alofonele palatalo-labializate sunt [l] și [a] sub influența lui *eo*

#### Nazalizarea

Nazalizarea este un proces ce influențează vocalele care sunt urmate de consoanele nazale *m* sau *n*, atunci când *m* sau *n* sunt urmate de o consoană. Consoana nazală și consoana care îi urmează trebuie să facă parte din aceeași silabă doar în cazul lui *m*.

(9) *an-te-ri-or* – alofonul nazalizat e [a] sub influența lui *n+cons*.

#### Fricativizarea

Fricativizarea reprezintă efectul consoanelor *f* sau *v* asupra consoanelor *m* sau *n* care le precedă.

(10) *in-for-ma-ți-e* – alofonul fricativizat e [n] sub influența lui *f*.

Fricativizarea este produsă însă și când consoanele *s*, *z*, *ș* sau *j* afectează consoana *n* care le precedă.

(11) *în-jo-si* – alofonul fricativizat e [n] sub influența lui *j*.

#### Velarizarea

Velarizarea este procedeul fonetic care afectează consoana *n*, atunci când *n* este urmată de sunetele *k*, *g* sau *h*.

(12) *en-gle-ză* – alofonul velarizat e [n] sub influența lui *g*.

#### Bilabializarea

Bilabializarea este procesul fonetic în urma căruia consoana *n* se pronunță *m*. Acest fenomen apare când *n* este urmat de *p*, *b* sau *m*.

(13) *în magazie* – alofonul bilabializat este [n] sub influența lui *m*.

#### Desonorizarea

Desonorizarea constă în asurzirea sau amuțirea sunetelor *m*, *n*, *l* când sunt precedate de o consoană în poziție finală de cuvânt.

(14) *calm* – alofonul desonorizat e [m] sub influența poziției sale finale și a lui *l*.

#### Aspirarea

Aspirarea este procesul fonetic prin care sunetele *p, b, t, d, k, g* sunt însoțite de un ușor sunet *h*.

(15) cap – alofonul aspirat e [p] sub influența poziției sale finale.

#### Neutralizarea

Neutralizarea este procesul fonetic în urma căruia sunetele *s, z* se confundă când sunt urmate de *m, n, l*.

(16) snob – alofonul neutralizat e [s] sub influența lui *n*.

#### Anularea stopului glotal

Anularea stopului glotal apare atunci când două vocale se află în hiat. În fața celei de-a doua vocale se aude un sunet mai slab (asemănător unei semivocale) asemănător lui *i* (notat <sup>l</sup>) sau lui *u* (notat <sup>w</sup>).

(17) am-bi-ți-<sup>l</sup>e (stop glotal cu apendice palatal)

sca-<sup>w</sup>un (stop glotal cu apendice labial)

#### Scurtarea vocalică postconsonantică

Scurtarea vocalică postconsonantică este procesul fonetic în urma căruia *i*, aflat la final de cuvânt, are o pronunțare scurtă, lipsită de intensitate sonoră

(18) Tu vezi pom<sub>i</sub>.

Ceea ce procesele fonetice scot în evidență este faptul că vocala *a*, de exemplu, va avea asociat diferite spectre acustice și diferite reprezentări ale trăsăturilor sale fonetice în funcție de procesul fonetic care acționează asupra sa contextual. Aceste ocurențe trebuie avute în vedere în realizarea corelațiilor dintre reprezentarea grafică a lui *a*, în acest caz, și multitudinea modurilor de pronunție ale sunetului corespondent în vorbirea articulată.

## 4. Concluzii

Perspectivile din punctul de vedere al aplicabilității unui sistem automat de recunoaștere a vorbirii sunt nelimitate. Obținerea unor rezultate caracterizate de acuratețe și eficiență în transcrierea fluxului oral pentru limba română depinde de modul discret de organizare a nivelurilor ierarhice ale limbii. Analiza fonetică este primul pas într-un astfel de demers, dar limitările sale, în absența introducerii unor reguli de morfologie și de sintaxă, sunt evidente. Fenomenele de omografie și omofonie, spectrul proceselor fonetice și prezența grupărilor diftong/triftong/hiat determină un model lingvistic complex care necesită cercetări colaborative de lungă durată.

## Mulțumiri

Studiul de față a fost posibil ca urmare a finanțării de către UEFISCDI a proiectului JustASR în cadrul programului Parteneriate, prin contractul 14/01.11.2013, cod înregistrare proiect PN-II-PT-PCCA-2013-4-1644.

## Referințe

- Antal, M., *Contributions to Speech and Speaker Recognition*. Cluj-Napoca, Technical University of Cluj-Napoca, 2006.
- Besacier, L., Barnard, E., Karpov, A., Schultz, T. 2014. Automatic speech recognition for under-resourced languages: A survey. *Speech Communication*, 85-100.
- Besacier, L., Barnard, E., Karpov, A., Schultz, T. 2014 b. Introduction to the special issue on processing under-resourced languages. *Speech Communication*, 83-84.
- Boldea, M. *Contributions to Automatic Continuous Speech Recognition for the Romanian Language*. Timișoara: Politehnica University, 2003.
- Burileanu, C., Popescu, V., Buzo, A., Petrea, C., Ghelmez-Hanes, D. 2010. Spontaneous Speech Recognition for Romanian in Spoken Dialogue Systems. *Proceedings of the Romanian Academy series A*, 83-91.
- Cinque, G., Giusti, G. *Advances in Romanian Linguistics*. John Benjamins Publishing, 1995.
- Cristea, D. 2005. Resurse lingvistice si tehnologii ale limbajului natural. Cazul limbii române, *Prelegeri Academice*, Academia Română - Filiala Iași, vol. III, nr. 3, Iași.
- Cucu, H., Buzo, A., Besacier, L., & Burileanu, C. (2014). SMT-based ASR domain adaptation methods for under-resourced languages: Application to Romanian. *Speech Communication*, 195-212.
- Dascălu Jinga, L., Pop, L. *Dialogul în limba română vorbită*. Oscarprint, 2003.
- Dascălu Jinga, L. *Corpus de română vorbită (CORV)*. Eșantioane. Oscarprint, 2002.
- \*\*\* *DOOM (Dicționarul ortografic și ortoepic al limbii române)*. Editura Universul Enciclopedic, 2005.
- eDTLR – The Thesaurus Dictionary of the Romanian Language, 2013, <http://metashare.info.uaic.ro/repository/browse/edtlr-the-thesaurus-dictionary-of-the-romanian-language/c05a74c063d611e28e8252540060617d5d02cd135a8b4759904cb7a956d41aa7/>
- Dumitrescu, S.D., Boros, T., Ion, R., 2014, Crowd - sourced, automatic speech - corpora collection - building the Romanian Anonymous Speech Corpus, presentation at LREC Workshop, [http://www.ilc.cnr.it/ccurl2014/LREC2014-Workshops\\_CCURL2014-Oral-7\\_Presentation.pdf](http://www.ilc.cnr.it/ccurl2014/LREC2014-Workshops_CCURL2014-Oral-7_Presentation.pdf)

- \*\*\* *GALR (Gramatica Academiei Limbii Române)*. Editura Academiei Române, 2008.
- \*\*\* *HTK Speech Recognition Toolkit*, <http://htk.eng.cam.ac.uk/>, 2009.
- Neamțu, G.G. *LRC. Fonetica limbii române*. curs susținut la Facultatea de Litere, Universitatea Babeș-Bolyai, Cluj-Napoca, 2002.
- \*\*\**NuanceCommunications,Inc.* Nuance Dragon Dictation and Dragon Search App Now Available in the Czech Republic and Romania. Press release from May 11, 2012 <http://www.nuance.com/company/news-room/press-releases/czechromaniaapps.doc>
- Petrea, C. S., Buzo, A., Cucu, H., Pasca, M., & Burileanu, C., 2010. Speech Recognition Experimental Results for Romanian Language. *Proceedings of the 6th European Conference on Intelligent Systems and Technologies*.
- Petrea, S., Ghelmez-Hanes, D., Buzo, A., Burileanu, C. 2009. Statistical Results in the Context of Romanian Spontaneous Speech Recognition . *Proceedings of the 13rd International Conference Speech and Computers SPECOM 2009*. St. Petersburg.
- Popescu, V., Burileanu, C., Caelen, J. 2008. Unsupervised Speaker Tracking in a Speech Recognition Module for Multi-Party Human-Computer Dialogue . *16th European Signal Processing Conference EUSIPCO 2008*.
- Rabiner, L., Huang, Y. (2005). Automatic Speech Recognition – A Brief History of the Technology, în K. Brown, *Encyclopedia of Language and Linguistics Second Edition*, Elsevier, 523-541.
- Rosetti, Al. *Introducere în fonetică*. Editura Științifică, 1963.
- \*\*\* *RASC (Romanian Anonymous Speech Corpus)*, <http://rasc.racai.ro/>, 2015.
- Schiopu, D. 2010. A Comparative Study of Three Speech Recognition Systems for Romanian Language, *The 5th International Conference on Virtual Learning*.
- Sofroni, V., Stan, A., (2015), O Descriere generală a Arhitecturii Sistemelor Automate de Recunoaștere Vocală. *Revista Romana de Interactiune Om-Calculator 8(1)*, 45-64.
- Stati, S. (1972). Unitățile limbii, în A. Graur, S. Stati, L. Wald, *Tratat de lingvistică generală*. Editura Academiei. 221-233
- Theodorescu, H. 2010. AI Tools for Speech Analysis Applied to the Romanian Language. *Proceedings of the 4th Conference on European Computing*. WSEAS Press.
- Theodorescu, H., Pistol, L., Feraru, M., Zbancioc, M., & Trandabat, D. (2010). *Sounds of the Romanian Language Corpus (SRoL)*. Retrieved from [http://www.etc.tuiasi.ro/sibm/romanian\\_spoken\\_language/index.html](http://www.etc.tuiasi.ro/sibm/romanian_spoken_language/index.html)
- Todorean, G., Buza, O. (2008). *Research report on project Interactive System based on Voice for Blind People*. Cluj Napoca: Technical University of Cluj Napoca.
- Trandabăț, D., Irimia, E., Barbu Mititelu, V., Cristea, D., Tufis, D. *The Romanian Language in the Digital Age-Limba Romana in Era Digitala*. Springer, 2012.