

Raport științific și tehnic in extenso

Program: Parteneriate în domenii prioritare

Proiect

Titlul: Sistem de Asistenta pentru Clădiri Inteligente, controlat prin Voce și Vorbire Naturală

Acronim: ANVSIB

Data: 25.09.2017

Etapa: 3/2016

Activitatea / activitățile:

- Activitatea 4.1. Integrarea sistemului de sinteza a vorbirii cu controlerul de voce
- Activitatea 4.2. Integrarea sistemului de recunoaștere a vorbirii cu controlerul de voce
- Activitatea 4.3. Scenarii de utilizare si configurare
- Activitatea 4.4. Integrarea hardware a infrastructurii KNX
- Activitatea 4.5. Evaluarea prototipului
- Activitatea 4.6. Diseminarea rezultatelor proiectului (partea II)
- Activitatea 4.7. Dezvoltarea unor modele acustice multilingve si pentru limba engleză (partea II)
- Activitatea 4.8. Integrarea algoritmului de căutare bazat pe DTW
- Activitatea 4.9. Implementarea infrastructurii si algoritmilor de detecție multilingva a termenilor
- Activitatea 4.10. Compararea si evaluarea performanțelor tehniciilor implementate
- Activitatea 4.11. Crearea unui prototip de sistem de recunoaștere multilingvă a vorbirii
- Activitatea 4.12. Crearea unui sistem de sinteza a vorbirii pentru limba romana
- Activitatea 4.13. Crearea unui sistem de sinteza a vorbirii pentru limba engleză
- Activitatea 4.14. Cercetarea si implementarea unui nou modul de prozodie
- Activitatea 4.15. Integrarea cu HTS
- Activitatea 4.16. Optimizarea si configurarea sistemului de sinteza a vorbirii

Număr contract: 32 / 2014

Acord de colaborare: 10677/24.06.2014 UPB, 1744/03.06.2014 IWAVE, 164/19.06.2014 RACAI

Autoritatea contractantă: Unitatea Executivă pentru Finanțarea Învățământului Superior, a Cercetării, Dezvoltării și Inovării

Conducător de proiect: UPB

Director de Proiect: Horia Cucu

Cuprins

1	Obiectivele activităților desfășurate	3
2	Perioada de desfășurare și personalul implicat.....	3
3	Rezumat activități.....	3
3.1	Activitatea 4.1. Integrarea sistemului de sinteza a vorbirii cu controlerul de voce	3
3.2	Activitatea 4.2. Integrarea sistemului de recunoaștere a vorbirii cu controlerul de voce	5
3.3	Activitatea 4.3. Scenarii de utilizare si configurare.....	5
3.4	Activitatea 4.4. Integrarea hardware a infrastructurii KNX	6
3.5	Activitatea 4.5. Evaluarea prototipului.....	7
3.6	Activitatea 4.6. Diseminarea rezultatelor proiectului (partea II).....	8
3.7	Activitatea 4.7. Dezvoltarea unor modele acustice multilingve si pentru limba engleza (partea II) .	9
3.8	Activitatea 4.8. Integrarea algoritmului de căutare bazat pe DTW	11
3.9	Activitatea 4.9. Implementarea infrastructurii și algoritmilor de detecție multilingva a termenilor vorbiți	12
3.10	Activitatea 4.10. Compararea si evaluarea performantelor tehniciilor implementate	15
3.11	Activitatea 4.11. Crearea unui prototip de sistem de recunoaștere multilingvă a vorbirii.....	16
3.12	Activitatea 4.12. Crearea unui sistem de sinteza a vorbirii pentru limba romana	18
3.13	Activitatea 4.13. Crearea unui sistem de sinteza a vorbirii pentru limba engleză.....	19
3.14	Activitatea 4.14. Cercetarea si implementarea unui nou modul de prozodie	19
3.15	Activitatea 4.15. Integrarea cu HTS	19
3.16	Activitatea 4.16. Optimizarea si configurarea sistemului de sinteza a vorbirii	20
3.17	Sumar al realizărilor / rezultatelor.....	20

1 Obiectivele activităților desfășurate

În ultima etapă a proiectului ANVSIB, obiectivele UPB au fost a) crearea prototipului de sistem de recunoaștere multilingvă a vorbirii și b) integrarea algoritmului DTW de căutare a cuvintelor cheie într-un sistem complet de detectie a termenilor vorbiti.

RACAI a avut ca obiectiv crearea, configurarea și optimizarea a două modele funcționale de sinteză a vorbirii pentru limbile română, respectiv engleză și integrarea lor într-un prototip multilingv de sinteză a vorbirii.

Obiectivele IWAVE au fost a) integrarea infrastructurii hardware de tip KNX cu controlerul de voce și sistemele de recunoaștere, respectiv sinteză a vorbirii și b) evaluarea prototipului final.

2 Perioada de desfășurare și personalul implicat

Nr.	Denumire activitate	Nume și prenume	Perioada
1.	Activitatea 4.1. Integrarea sistemului de sinteză a vorbirii cu controlerul de voce	Ştefan Daniel Dumitrescu; Tiberiu Boroş Marian Bădulescu	noi 2016 – sept 2017
2.	Activitatea 4.2. Integrarea sistemului de recunoaștere a vorbirii cu controlerul de voce	Marian Bădulescu; Horia Cucu; Lucian Georgescu; Alexandru Caranica; Andi Buzo; Corneliu Burileanu; Florentina Mincă	noi 2016 – sept 2017
3.	Activitatea 4.3. Scenarii de utilizare și configurare	Marian Bădulescu; Horia Cucu; Lucian Georgescu; Alexandru Caranica; Andi Buzo; Corneliu Burileanu; Florentina Mincă; Ştefan Daniel Dumitrescu; Tiberiu Boroş	noi 2016 – sept 2017
4.	Activitatea 4.4. Integrarea hardware a infrastructurii KNX	Marian Bădulescu;	noi 2016 – sept 2017
5.	Activitatea 4.5. Evaluarea prototipului	Marian Bădulescu;	noi 2016 – sept 2017
6.	Activitatea 4.6. Diseminarea rezultatelor proiectului (partea II)	Horia Cucu; Lucian Georgescu; Alexandru Caranica; Andi Buzo; Ştefan Daniel Dumitrescu; Tiberiu Boroş	noi 2016 – sept 2017
7.	Activitatea 4.7. Dezvoltarea unor modele acustice multilingve și pentru limba engleză (partea II)	Horia Cucu; Lucian Georgescu; Alexandru Caranica; Andi Buzo; Corneliu Burileanu; Florentina Mincă	noi 2016 – sept 2017
8.	Activitatea 4.8. Integrarea algoritmului de căutare bazat pe DTW	Horia Cucu; Lucian Georgescu; Alexandru Caranica; Andi Buzo; Corneliu Burileanu; Florentina Mincă	noi 2016 – sept 2017
9.	Activitatea 4.9. Implementarea infrastructurii și algoritmilor de detectie multilingvă a termenilor	Horia Cucu; Lucian Georgescu; Alexandru Caranica; Andi Buzo; Corneliu Burileanu; Florentina Mincă	noi 2016 – sept 2017
10.	Activitatea 4.10. Compararea și evaluarea performanțelor tehniciilor implementate	Horia Cucu; Lucian Georgescu; Alexandru Caranica; Andi Buzo; Corneliu Burileanu; Florentina Mincă	noi 2016 – sept 2017
11.	Activitatea 4.11. Crearea unui prototip de sistem de recunoaștere multilingvă a vorbirii	Horia Cucu; Lucian Georgescu; Alexandru Caranica; Andi Buzo; Corneliu Burileanu; Florentina Mincă	noi 2016 – sept 2017
12.	Activitatea 4.12. Crearea unui sistem de sinteză a vorbirii pentru limba română	Ştefan Daniel Dumitrescu; Tiberiu Boroş	noi 2016 – sept 2017
13.	Activitatea 4.13. Crearea unui sistem de sinteză a vorbirii pentru limba engleză	Ştefan Daniel Dumitrescu; Tiberiu Boroş	noi 2016 – sept 2017
14.	Activitatea 4.14. Cercetarea și implementarea unui nou modul de prozodie	Ştefan Daniel Dumitrescu; Tiberiu Boroş	noi 2016 – sept 2017
15.	Activitatea 4.15. Integrarea cu HTS	Ştefan Daniel Dumitrescu; Tiberiu Boroş	noi 2016 – sept 2017
16.	Activitatea 4.16. Optimizarea și configurarea sistemului de sinteză a vorbirii	Ştefan Daniel Dumitrescu; Tiberiu Boroş	noi 2016 – sept 2017

3 Rezumat activități

3.1 Activitatea 4.1. Integrarea sistemului de sinteză a vorbirii cu controlerul de voce

Integrarea sistemului de sinteză a vorbirii cu controllerul de voce este prezentată în contextul întregului sistem, deoarece este parte integrantă.



Arhitectural, sistemul poate fi descris din mai multe puncte de vedere, însă în acest capitol vom descrie modul de interconectare al diferitelor componente și module, cu accent mai puțin pe detalierea fiecărui.

Figura de mai sus prezintă principalele componente ale sistemului. În partea stângă apar toate dispozitivele inteligente care vor fi controlate (lumini, aer condiționat, uși, etc.). Serverul este punctul central, creierul casei. Totul va fi interconectat aici; tot aici are loc procesarea audio. Serverul poate să fie conectat la Internet, a.î. utilizatorul poate controla din exteriorul casei, spre exemplu, aerul condiționat. Interfețele de control (partea dreaptă) reprezintă puncte de interacțiune ale utilizatorului cu spațiul intelligent. În cadrul acestui proiect sunt disponibile două metode de control: dispozitive mobile (telefon/tabletă/etc.) și așa numitele "dispozitive fixe". Aceste dispozitive sunt cutii mici (de dimensiunea unei doze electrice) care sunt conectate la server și prin cablu Ethernet, și au un microfon care ascultă permanent evenualele comenzi ale utilizatorului. De asemenea, sunt conectate la un set de difuzoare în aceeași cameră prin care comanda dată va fi confirmată.

Prezentăm în continuare fiecare componentă separat, cu accent pe modul în care modulele discută între ele (protocolul de comunicație folosit).

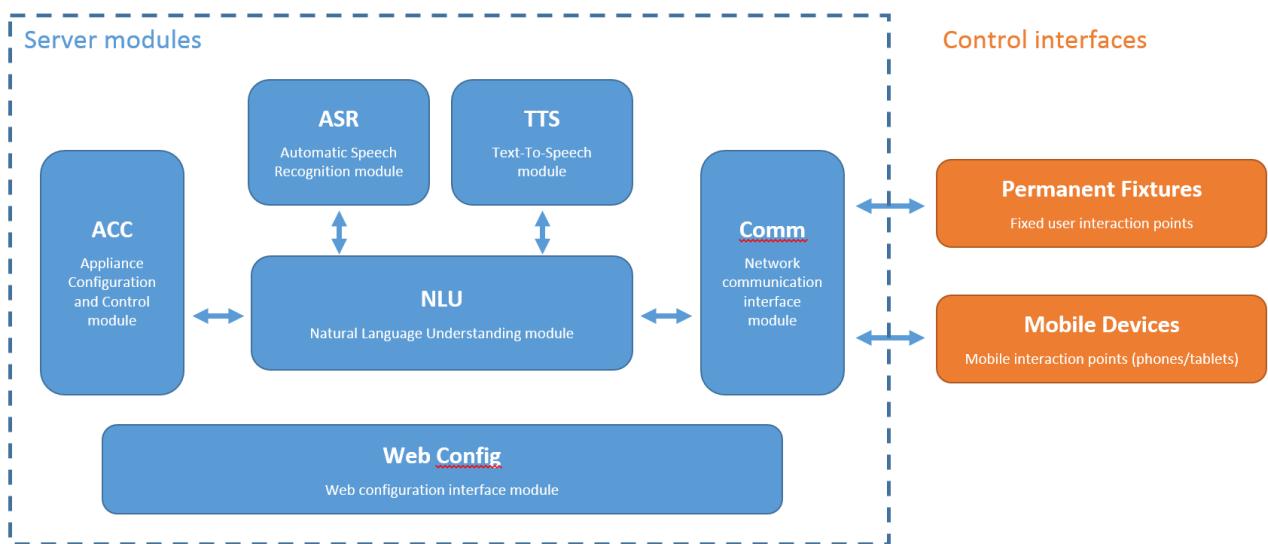


Diagrama de mai sus prezintă fiecare modul al sistemului. Serverul, din punct de vedere logic, are următoarele module:

- **Modulul NLU – Natural Language Understanding**

Acest modul este creierul spațiului intelligent. Primește și trimite mesaje de la/către toate celelalte module. Sarcina principală este de a înțelege mesajele primite (mod text) pentru ca să le poată „traduce” mai departe către modulul ACC (Appliance Configuration and Control).

- **Modulul TTS – Text-to-Speech**

Aici are loc sinteza mesajelor text pentru a obține voce. Modulul primește un mesaj text, limba în care se dorește sinteza (engleză sau română) și întoarce un stream wave audio direcționat către punctul de la care s-a primit mesajul text.

- **Modulul ASR – Automatic Speech Recognition**

În acest modul are loc procesul de recunoaștere de voce, primind un stream wave audio și returnând un mesaj text. Este practic inversul procesului de sinteză din modulul precedent.

- **Modulul ACC – Appliance Configuration and Control module**

În acest modul are loc configurarea inițială a dispozitivelor existente în spațiul intelligent. Protocolul de comunicație folosit este KNX, interconectarea fiind efectuată prin TCP/IP (Ethernet sau Wireless acolo unde dispozitivul permite configurare).

- **Modulul Web**

Acest modul este de fapt un server standard Apache, având instalat un site care permite utilizatorului să configureze toate device-urile din casă, alături de metodele de apel către aceste device-uri (adică utilizatorul își poate defini comanda proprie de control asupra fiecărui dispozitiv).

Interfețele de control sunt fie dispozitive mobile fie fixe:

- **Dispozitive mobile**

Utilizatorul poate controla casa inteligentă folosind telefonul sau tableta proprie, sau orice alt dispozitiv intelligent care poate rula aplicații Android și are acces la server prin Wireless local sau prin conexiune 3G/4G la server.

- **Dispozitive fixe**

Un punct fix reprezintă un punct de interacțiune plasat în cameră. Acest dispozitiv ascuns, de exemplu, în perete într-o doză electrică sau direct aplicat pe perete, are un microfon care stă permanent pornit, ascultând eventualele comenzi. Acest dispozitiv fix ascultă un cuvânt cheie prin care se începe comanda către casă, și, în funcție de configurație, fie procesează semnalul audio local (adică efectuează procesul de ASR local trimițând către server doar textul comenzi) fie trimită fluxul audio către server, serverul ocupându-se mai departe de procesare. Primește înapoi de la server confirmarea comenzi. Dispozitivul în sine este un sistem x86 sau ARM tip System-On-A-Chip (ex: Raspberry Pi) de dimensiuni reduse.

3.2 Activitatea 4.2. Integrarea sistemului de recunoaștere a vorbirii cu controlerul de voce

Modulul ASR (Automatic Speech Recognition), alături de modulul TTS prezentat în activitatea anterioară, a fost prezentat în contextul întregului sistem (ca mod de legătură), nemai fiind necesară o descriere în acest punct. Însă, pentru a înțelege modul în care ASR-ul și TTS-ul sunt integrate cu sistemul, considerăm oportună prezentarea unui exemplu de comunicație între aceste module, alături de toate celelalte module prezente pe server.

Exemplu de comunicație între module.

Să presupunem că avem un scenariu simplu în care utilizatorul dorește *închiderea luminilor din bucătărie*. Vom exemplifica în paralel pentru ambele tipuri de interfețe de control.

- **Pasul 1.** Utilizatorul dă comanda „Casandra închide lumina în bucătărie” pe telefonul său mobil. Mobilul înregistrează comanda și o trimită mai departe către NLU sub formă de audio raw (fișier wav). Dacă comanda a fost trimisă prin intermediul unui dispozitiv fix și dacă acel dispozitiv fix rulează procesul de ASR local, audio-ul va fi transformat direct pe dispozitivul fix în mesaj text, către server trimițându-se doar mesajul text, sărind peste pasul 2.
- **Pasul 2.** Audio-ul este trimis de către NLU către modulul de ASR. De acolo primește înapoi un mesaj text.
- **Pasul 3.** Modulul NLU analizează mesajul și îl transformă într-o comandă efectivă către un dispozitiv inteligent din casă. În exemplul curent, mesajul transformat arată astfel: {target_appliance="kitchen_lights", action="off"}.
- **Pasul 4.** Modulul ACC efectuează comanda, și anume închide luminile din bucătărie. În acest pas, este posibil ca dispozitivul către care se trimită comanda să ofere un feedback către utilizator de tipul: lumina a fost închisă, sau Ușa a fost închisă.
- **Pasul 5.** Dacă modulul NLU așteaptă un feedback, în funcție de statusul primit va trimite către dispozitivul de interacțiune un mesaj care poate fi fie mesajul standard “comandă confirmată” sau un mesaj corespunzător comenzi utilizatorului (ex: “volumul a fost setat la 60%”).
- **Pasul 6.** Dispozitivul de interacțiune (fie mobil fie fix) va primi mesajul text și îl va sintetiza prin trimiterea directă către modulul TTS a mesajului. Modulul TTS trimite direct înapoi audio-ul raw care este redat fie de către telefonul mobil, fie prin difuzeoarele la care este legat dispozitivul fix.

3.3 Activitatea 4.3. Scenarii de utilizare și configurare

Sistemul NLU este bazat pe scenarii. În orice moment, utilizatorul poate modifica configurația sistemului (Cassandra) prin adăugarea, modificarea sau ștergerea de scenarii de utilizare.

Un scenariu este definit de un nume și de trei componente:

1. O lista de exemple de propoziții în care fiecare parametru este precedat de un semn special (\$). De exemplu: “Setează temperatura la * grade” unde * este un parametru variabil (utilizatorul poate spune “zece”, “două zeci șișapte”, etc.).
2. O listă de acțiuni pe care o va efectua sistemul în urma identificării comenzi definite anterior. Acțiunile sunt definite în format JSON astfel: a. id-ul dispozitivului inteligent (ex: o anume unitate A/C, un set de becuri, o anumită ușă, etc.); b. Parametrii comenzi – o structură fără format fix JSON care îi spune sistemului ce parametrii să transmită mai departe dispozitivului inteligent. (ex: aprinde lumina on=1, stinge lumina off=0, setează A/C la X grade, etc.) Valorile acestor parametri pot fi fie constante, fie variabile predefinite, sau valori efective extrase din text (adică string-uri, numere, etc.).
3. Feedback. După procesare, JSON-ul definit aici este pasat dispozitivului de control care a inițiat sesiunea de comenzi, fiind folosit pentru a transmite informație utilizatorului. Acest JSON are de asemenea două atribute: a. Friendly_response – un răspuns text care, dacă există, va fi sintetizat și redat utilizatorului; răspunsul poate conține variabile predefinite sau extrase din comanda inițială, sau generate de dispozitiv. (exemplu: temperatura efectivă din cameră ca răspuns oferit de aparatul de aer condiționat la întrebarea utilizatorului asupra temperaturii ambientale). b. Launch_intent – o structură care informează dispozitivul de interacțiune (telefonul mobil în acest caz) ce aplicație trebuie să lanseze (tip Android Intent) cu un anumit set de parametrii. Cassandra poate genera acțiuni pentru previzualizarea de imagini, audio și video pe telefoane mobile.

Metodologia pentru identificarea scenariilor si extragerea de parametrii este efectuata într-o serie de pași secvențiali.

- a. Scenariul este identificat utilizând un arbore de decizie tip ID3. Există și posibilitatea de a utiliza o rețea convecțională special creată, aceasta însă necesitând o serie de date lingvistice pe baza cărora se generează o codare a cuvintelor (folosind word2vec). În producție, se folosește arborele de decizie pentru a identifica scenariul, deoarece nu este necesară nicio altă resursă lingvistică care să scadă din universalitatea metodei de identificare.

Odată definite toate scenariile, se extrag propozițiile pentru fiecare scenariu. Generarea arborelui de decizie se efectuează extrăgând într-un format bag-of-words fiecare cuvânt din fiecare comandă, având asociat scenariul efectiv. Așadar, pentru fiecare scenariu avem cel puțin o propoziție (de obicei minim 3-4). Algoritmul de training analizează fiecare propoziție (comandă) în parte, generând un arbore. Arborele este redus ca dimensiune pentru a garanta un anumit nivel de universalitate.

- b. Extragerea parametrilor din comandă este efectuată tot cu un arbore de decizie tip ID3. În primul rând trebuie identificată poziția parametrilor, acest lucru fiind efectuat în doi pași. În primul pas este identificată poziția de start (indicele cuvântului care începe parametrul – un parametru poate avea mai mult de un cuvânt, de exemplu “temperatura la două zeci și șapte de grade” are un singur parametru de 4 cuvinte) folosind un arbore antrenat pentru a prezice, pe fiecare cuvânt, fie clasa Start (ex: cuvântul “două”) fie Null (toate celelalte cuvinte). Odată identificat un start de parametru, se folosește un alt arbore de decizie, antrenat pentru a prezice Stop (ex: cuvântul “șapte”) sau Null. Parametrul este cel între Start-Stop, inclusiv. Odată marcat un cuvânt cu Stop, se pornește căutarea cu arborele care prezice Start, deoarece este posibil ca o comandă să aibă mai mulți parametri. Contextul folosit pentru arbori este de 4 cuvinte înaintea cuvântului analizat pentru arborele Start, respectiv 4 cuvinte după pentru arborele Stop.

Odată identificată parametrii, trebuie găsit tipul fiecărui. Acest lucru se face recurgând la un alt arbore de decizie, care stabilește dacă parametrul este valoare numerică, string, predefinit, etc.

În continuare prezentăm câteva scenarii practice de utilizare.

Scenariul principal: utilizatorul controlează aparatele electrocasnice folosind vocea. Există mai multe sarcini predefinite cum ar fi controlul multimedia, iluminare, climatizare sau securitate. De exemplu: “Casandra, pornește sistemul de climatizare” (en. Casandra, turn on the air-conditioning system), “Casandra, aprinde toate luminile” (en. Casandra, turn on all the lights), “Casandra, activează stropitorile” (en. Casandra, start the irrigation system), “Casandra, mărește temperatura cu trei grade în sufragerie” (en. Casandra, raise the living room temperature by three degrees). În acest scenariu se încadrează toate interacțiunile cu dispozitivele inteligente din casă.

Scenariul 2: reprezintă un sistem standard de interogare (întrebare-răspuns sau QA), în care utilizatorul îi poate adresa Casandrei diverse întrebări (de ex. “cum este vremea”) la care sistemul va răspunde folosind o bază de cunoștințe (KB), similară cu ce ar răspunde Google Now, Stiri sau Amazon Echo. Baza de cunoștințe (KB) a fost construită folosind Wikipedia pentru limba română. Acest modul este de tip proof-of-concept, nefiind prevăzut în cadrul proiectului dar fiind implementat pentru a extinde funcționalitatea sistemului în afara controlului direct al dispozitivelor inteligente.

Scenariul 3: este un set de întrebări scurte / răspunsuri. Acest set de Q/A conține un joc care se numește “psihologul”. Aceasta permite sistemului să efectueze asociații de cuvinte într-un mod similar în care un psiholog ar cere unui pacient să o facă. Acest joc s-a dovedit a fi foarte atrăgător pentru utilizatorii din grupul nostru de testare, în primul rând pentru că a făcut sistemul să pară mai “asemănător unui om”. Ideea de bază este că, folosind Wikipedia și word2vec, am extras seturi de cuvinte similare. Rostind un cuvânt, “psihologul” va să rostească un alt cuvânt cu aproape același sens (dpd. Semantic), jocul continuând în runde.

3.4 Activitatea 4.4. Integrarea hardware a infrastructurii KNX

Asociația KNX este creațoarea și proprietara tehnologiei KNX - STANDARD mondial pentru toate aplicațiile de control pentru case și clădiri, variind de la controlul iluminatului și al jaluzelor până la variate sisteme de securitate, încălzire, ventilație, aer condiționat, monitorizare alarme, controlul instalațiilor sanitare, managementul energiei, contorizare, precum și al dispozitivelor de uz casnic, audio/video și multe altele. Tehnologia KNX poate fi utilizată atât în clădiri noi cât și în cele deja existente. KNX este standardul global pentru controlul clădirilor cu o singură platformă software, independentă de un producător KNX, pentru design și punere în funcțiune (ETS), cu o gamă completă de medii de transmisie disponibile (TP, PL, RF și IP), precum și cu o gamă completă de moduri de configurare (Modul System și Modul Easy). KNX este aprobat ca Standard European (CENELEC EN 50090 și CEN EN 13321-1) și ca Standard Internațional (ISO/IEC 14543-3). Pentru a asigura transferul datelor între toate componente de management al clădirii, este necesar un sistem care să rezolve problema comunicării cu dispozitivele izolate, asigurându-se că toate componentele comunică într-un limbaj comun unic: pe scurt, un sistem precum Magistrala KNX, independent de producător și de domeniul de aplicație. Acest standard se bazează pe mai mult de 23 de ani de experiență pe piață, printre altele, cu sistemele predecesoare KNX: EIB, EHS și BatiBUS. În mediul KNX toate dispozitivele sunt conectate la același BUS (pereche torsadată, RF, linii electrice sau IP/Ethernet) cu posibilitatea de a comunica între ele. Dispozitivele KNX pot fi: senzori sau elemente de acționare, necesare controlului diferitelor sisteme ale clădirii,

precum: iluminat, jaluzele/storuri, sisteme de securitate, managementul energiei, încălzire, ventilație și sisteme de aer condiționat, sisteme de semnalizare și monitorizare, interfețe cu sistemele de service și control al clădirii, control de la distanță, contorizare, control audio/video, electrocasnice, etc. Toate aceste funcții sunt controlate, monitorizate și semnalizate prin intermediul unui sistem unitar fără a fi necesare centre de control suplimentare. (descriere preluată de pe www.knx.ro)

Deoarece în standardul KNX nu există un controller (fiecare echipament este independent și ascultă pe bus), practic efortul de integrare a fost majoritar dozat în configurația software a echipamentelor. Din punct de vedere hardware tot ce trebuie făcut este conectarea lor în cadrul unei rețele. Există mai multe metode suportate de KNX pentru interconectare (în funcție de interfața echipamentului):

1. KNX TP: Twisted Pair, adică o pereche de fire de cupru subțiri, este mediu vechi de comunicare, preluat din standardul EIB, viteza 9600bps, oferă contact punct-la-punct.
2. KNX RF: Radio Frequency, adică conexiune prin radio. Pachetele se transmit pe distanță mică, frecvență de 868MHz, puterea maximă permisă fiind de 25mW.
3. KNX PL: Power Line, conexiune directă pe linia de curent, viteza de 1200bps.
4. KNX IP: mesajele sunt transmise pe rețea de date TCP/IP. Marea parte a dispozitivelor inteligente oferă această interfață.

Echipamentele conectate/controlate de I-WAVE în cadrul proiectului au fost: serie de lumini, dispozitiv audio multimedia intelligent (conectat la internet), lumini colorate, sistem climatizare automată, sistem de securitate cu mai multe partiții activabile individual. Toate aceste echipamente au fost conectate, configurate și controlate de către Cassandra. Unele dintre ele se leagă în KNX direct (dacă au interfață KNX), în timp ce marea majoritate cu Comfortclick prin HTTP API; unele comunică prin infraroșu IR (ex aerul conditionat individual).

Comfortclick integrează mai multe protocoale (HUE, KNX, MODBUS, HTTP API, IR) și le pune la dispozitie pe portul 84 prin JSON .

Efortul de configurație hardware a constat în adăugarea fiecărui dispozitiv în Comfortclick, identificarea printr-un ID unic și preluarea aceluia ID și introducerea în interfață de configurație a Cassandrei. În acea interfață au trebuit definite acțiunile pentru fiecare dispozitiv. De exemplu, pentru lumini sunt necesari doi parametri: starea luminii (aprins/stins – variabilă booleană) și poziția luminii (adică locația ei, de exemplu bucătărie este lumina numărul 1, sufragerie numărul 2, etc. – variabilă număr întreg). Pentru lumini colorate mai sunt necesari încă doi parametri, și anume culoarea (variabilă tip RGB) și intensitatea (ex: 100%, 60%, etc. – variabilă număr întreg). În final, trebuie setate un număr de propoziții parametrizabile pe care utilizatorul le va folosi pentru a efectua o acțiune. De exemplu, pentru a aprinde lumina, utilizatorul va defini propoziția „aprindă lumina în LOC” unde LOC este oricare din locațiile anterior definite, de tipul bucătărie = 1, sufragerie = 2, etc, iar acțiunea este definită prin a păsa o valoare de tip bool “value=true” dispozitivului bec cu id-ul LOC.

Așadar, această activitate a presupus studierea fiecărui dispozitiv în parte; conectarea fizică a fiecărui dispozitiv în funcție de interfață de comunicație; asignarea în Comfortclick a unui ID unic; configurația în interfață de management a Cassandrei a fiecărui dispozitiv și definirea parametrilor în JSON-ul propriu alături de acțiunea propriu-zisă.

3.5 Activitatea 4.5. Evaluarea prototipului

Serverul este “creierul” casei. În sine este un calculator de dimensiuni reduse. Serverul se ocupă cu toată partea de procesare efectivă, așadar are nevoie de resurse de calcul. Testele în cadrul proiectului s-au desfășurat pe un calculator cu procesor Intel i3 cu 4 GB RAM și un SSD pentru stocare. Dimensiunea SSD-ului nu este relevantă deoarece nu stocăm practic nimic în afară de sistemul de operare reprezentat de un Linux, alături de modulele software care în sine au o dimensiune nesemnificativă comparativ cu spațiul oferit de orice SSD modern, adică cel puțin 100GB). Serverul funcționează separat, nefiind nevoie de acces la Internet. Dacă totuși este conectat, serverul poate primi comenzi de la dispozitive mobile din afara razei Wireless sau a rețelei locale.

Controllerul de voce, reprezentat prin modulul NLU (Natural Language Understanding) este punctul central al serverului, deoarece leagă toate celelalte module între ele. Misiunea acestui modul este de a traduce mesajele generate de ASR în mesaje către modulul de control al dispozitivelor inteligente din casă.

Evaluarea prototipului constă în implementarea unor funcționalități practice, bazare pe dispozitive atât reale cât și virtuale. Astfel, au fost implementate următoarele scenarii (bazate pe tipul 1 de scenariu prezentat în secțiunea precedentă – * reprezintă valoarea dorită, LOC reprezintă o locație predefinită de tipul “bucătărie”, “dormitor”, “baie”, etc.):

HVAC – Read temperature casandra care este temperatura din LOC casandra câte grade sunt în LOC casandra citește temperatura din LOC	Lights – Set decrease intensity casandra micșorează luminozitatea din LOC la * la sută casandra micșorează luminozitatea din LOC până la * la sută casandra scade luminozitatea din LOC la * la sută casandra micșorează intensitatea luminii din LOC până la * la sută
HVAC – Set absolute temperature casandra setează temperatura din LOC la * de grade casandra schimbă temperatura din LOC la * de	

grade	casandra dă aerul condiționat pe * grade casandra pornește aerul condiționat pe * grade casandra schimbă aerul condiționat pe * grade casandra pune aerul condiționat pe * grade casandra setează aerul condiționat pe * grade	Security - Call police casandra sună la poliție casandra cheamă poliția casandra sună poliția
HVAC - Set increase temperature	casandra crește temperatura din LOC cu * grade casandra setează aerul condiționat mai cald casandra dă aerul condiționat mai cald	Security - Call ambulance casandra sună la salvare casandra sună la ambulanță casandra cheamă salvarea casandra cheamă ambulanță casandra sună salvarea casandra sună ambulanță
HVAC - Set decrease temperature	casandra scade temperatura din LOC cu * grade casandra setează aerul condiționat mai rece casandra dă aerul condiționat mai rece	Security - Call firefighters casandra sună la pompieri casandra cheamă pompierii casandra sună pompierii
HVAC - Set windows off	casandra închide ferestrele din LOC casandra închide geamurile din LOC	Security - Set efractation state on casandra armează centrala de alarmă casandra armează partitia cu numărul * casandra armează partitia * casandra activează centrala de alarmă casandra activează partitia cu numărul * casandra activează partitia *
HVAC - Set windows on	casandra deschide ferestrele din LOC casandra deschide geamurile din LOC	Security - Set efractation state off casandra dezarmează centrala de alarmă casandra dezarmează partitia cu numărul * casandra dezarmează partitia * casandra dezactivează centrala de alarmă casandra dezactivează partitia cu numărul * casandra dezactivează partitia *
HVAC - Set window blinds off	casandra închide obloanele din LOC casandra închide jaluzelele din LOC	Security - Get efractation state casandra care este statusul partiei * casandra ce status are partitia * casandra cum este partitia * casandra status partie * casandra statusul partiei *
HVAC - Set window blinds on	casandra deschide obloanele din LOC casandra deschide jaluzelele din LOC	Multimedia - Set absolute sound intensity casandra setează volumul televizorului la * la sută casandra schimbă volumul televizorului la * la sută casandra setează volumul muzicii la * la sută casandra schimbă volumul muzicii la * la sută
Lights - Set on	casandra aprinde lumina din * casandra pornește lumina din * casandra deschide lumina din *	Multimedia - Set increase sound intensity casandra mărește volumul televizorului casandra crește volumul televizorului
Lights - Set off	casandra stinge lumina din * casandra oprește lumina din * casandra inchide lumina din *	Multimedia - Set decrease sound intensity casandra micșorează volumul televizorului casandra scade volumul televizorului
Lights - Set on all	casandra aprinde toate luminile casandra pornește toate luminile casandra deschide toate luminile	Multimedia - Set on casandra pornește muzica casandra pornește radio casandra pornește radioul casandra pornește televizorul
Lights - Set off all	casandra oprește toate luminile casandra închide toate luminile casandra stinge toate luminile	Multimedia - Set off casandra oprește muzica casandra oprește radio casandra oprește radioul casandra oprește televizorul
Lights - Set color	casandra schimbă culoarea luminii din * în * casandra aprinde lumina * în * casandra pornește lumina * în * casandra deschide lumina * în *	Multimedia - Set radio station casandra vreau să ascult radio *
Lights - Set absolute intensity	casandra setează luminozitatea din LOC la * la sută casandra schimbă luminozitatea din LOC la * la sută casandra setează intensitatea luminii din LOC la * la sută casandra schimbă intensitatea luminii din LOC la * la sută	
Lights - Set increase intensity	casandra mărește luminozitatea din LOC la * la sută casandra mărește luminozitatea din LOC până la * la sută casandra crește luminozitatea din LOC la * la sută casandra crește luminozitatea din LOC până la * la sută	

La fiecare scenariu Casandra va efectua ac iunea dorită și va confirma respectiva ac iune printr-un mesaj sintetizat. Prototipul a fost testat pentru fiecare din scenariile de mai sus, confirmându-se funcționarea corectă.

3.6 Activitatea 4.6. Diseminarea rezultatelor proiectului (partea II)

Pentru această etapă finală, activitatea de diseminație a RACAI a rezultatelor proiectului s-a efectuat prin participarea la 2 conferințe unde au fost prezentate 3 articole:

Conferința SPED17: The 9th Conference on Speech Technology and Human-Computer Dialogue, Iulie 6-9, 2017 – București, Romania

- Articolul 1: A “Small-Data” - Driven Approach to Dialogue Systems for Natural Language Human Computer Interaction , Tiberiu Boroș, Ștefan Daniel Dumitrescu. Acest articol prezinta modul de configurare al Cassandra din punct de vedere al scenariilor (interacțiune om-casa inteligentă).
- Articolul 2: Cassandra Smart-Home System Description, Ștefan Daniel Dumitrescu. Acest articol prezinta o descriere a Cassandra din punct de vedere al arhitecturii logice și al implementării modulelor.

Conferința ACL (Association for Computational Linguistics), workshop CONLL (The SIGNLL Conference on Computational Natural Language Learning), Vancouver, Canada, 3-4 August 2017

- Articol 1: RACAI’s Natural Language Processing pipeline for Universal Dependencies, Stefan Daniel Dumitrescu, Tiberiu Boroș and Dan Tuși. Acest articol prezinta preprocesatorul de text prezent în Casandra și modul în care a fost extins la un număr arbitrar de limbi (unul din întele acestui proiect), utilizând resurse lingvistice puse la dispoziție în cadrul competiției CONLL 2017.

Pentru această etapă finală, activitatea de diseminare a UPB a rezultatelor proiectului s-a efectuat prin participarea la o conferință ISI și transmiterea spre publicare a unui articol de jurnal ISI:

- Alexandru Caranica, Horia Cucu, Corneliu Burileanu, François Portet, Michel Vacher, “Speech Recognition Results for Voice-controlled Assistive Applications,” in the Proceedings of the 9th Conference on Speech Technology and Human-Computer Dialogue (SpeD), Bucharest, 2017, 8p, ISBN 978-1-5090-6496-0;
- Alexandru Caranica, Horia Cucu, Andi Buzo, Corneliu Burileanu, “Survey on Multilingual Spoken Term Detection,” in Romanian Journal of Information Science and Technology (trimisă spre publicare).

3.7 Activitatea 4.7. Dezvoltarea unor modele acustice multilingve și pentru limba engleză (partea II)

În cadrul acestei activități, au fost realizate experimente ce au avut ca scop crearea unor modele acustice, atât pentru limba engleză, cât și pentru limba română. Corpusurile audio au fost împărțite în seturi de antrenare și de evaluare, fiind utilizate modele de limbă probabilistice de tip n-gram. Pentru antrenarea modelelor s-a folosit atât utilitarul Kaldi cât și Sphinx. Sphinx permite parametrizarea semnalului vocal și extragerea coeficienților de tip MFCC, PLP sau LPC, în plus grupul nostru integrând și parametrii robusti la zgomot, de tip PNCC, pentru care am prezentat rezultatele în raportul anterior. Modelarea acustică folosind acest toolkit presupune crearea unor sisteme statistice de tip HMM-GMM.

În plus, față de Sphinx, Kaldi oferă antrenarea pe baza de rețele neuronale, ceea ce a oferit un spor de performanță pentru anumite modele, după cum am prezentat anterior. Avantajul major este reprezentat de existența algoritmilor de antrenare ce folosesc rețele neuronale adânci (DNN), cu posibilitatea utilizării unui cluster de stații de lucru, echipate cu GPU-uri Tesla. În urma testelor, s-a observat faptul că această platformă de tip DNN oferă o acuratețe mai bună în comparație cu sistemele HMM-GMM din Sphinx.

În continuare, în cadrul grupului, am continuat achiziția de resurse lingvistice, în special pentru limba engleză. O nouă bază de date audio, intitulată homeAuto4, a fost achiziționată, în cadrul grupului SpeeD, pentru evaluarea performanțelor sistemului multilingv, în limba engleză. Un set de 154 de transcrieri, în limba engleză, au fost realizate, ce reprezintă comenzi ce pot fi înțelese și interpretate de sistemul de recunoaștere a vorbirii. Un număr de 11 vorbitori de limba engleză, non-nativi, au realizat înregistrări audio, folosind noua versiune a aplicației “Speech Recorder”, disponibilă la adresa <https://speech-recorder.speed.pub.ro> (Figura 1).

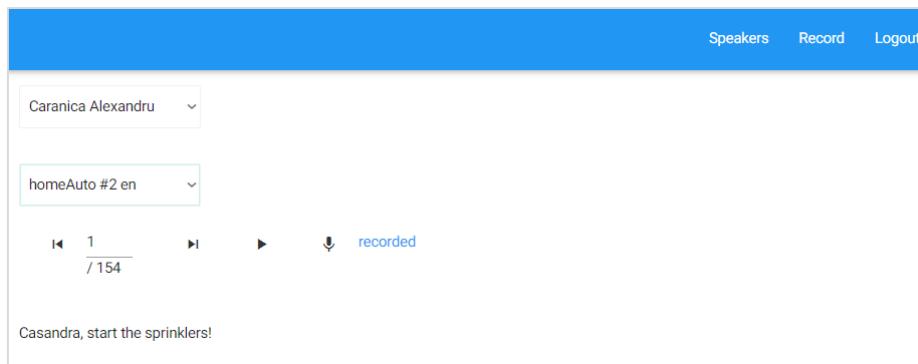


Figura 1 Noua versiune a aplicației Speech Recorder, a grupului SpeeD

Noul corpus de evaluare totalizează un număr de aproximativ 2 ore și 10 minute de înregistrări, în limba engleză, înregistrările fiind realizate de către studenți și personal al grupului de cercetare, cu vârste cuprinse între 24 și

35 de ani, procentul de vorbitori feminini fiind de aproximativ 27%. Aceste informații, summarizate, se regăsesc în Tabelul 1.

	Evaluare
Număr conversații	154
Durata totală	2 h 10 m
Vorbitori masculini	~95m
Vorbitori feminini	~35m
Durata medie	11 m 8 s
Număr de vorbitori unici	11
Vârstă	24-35 ani

Tabelul 1 Corpusul de evaluare homeAuto 4 achiziționat

Pentru realizarea modelului acustic, în limba engleză, s-a folosit corpusul de vorbire TED-LIUM, versiunea 1 [Rousseau, 2012], ce conține înregistrări audio și transcrierile aferente de la conferințele Ted (<https://www.ted.com>). Corpusul a fost împărțit în 3 seturi distincte ale căror caracteristici se regăsesc în Tabelul 2.

	Antrenare	Dezvoltare	Testare
Număr de conversații	774	19	11
Durata totală	118 h 4 m 48 s	4 h 12 m 55 s	3 h 4 m 5 s
• Vorbitori masculini	81 h 53 m 7s	3 h 13 m 57 s	2 h 21 m 35 s
• Vorbitori feminini	36 h 11 m 41 s	58 m 58 s	42 m 29 s
Durata medie	9 m 9 s	13 m 18 s	16 m 43 s
Număr de vorbitori unici	666	19	11

Tabelul 2 Corpusul TED-LIUM

Construirea unui model multilingual RO-EN s-a realizat în CMU-Sphinx, prin antrenarea unui sistem folosind corpusurile de antrenare în limba engleză și cele din limba română.

Pentru realizarea și evaluarea modelului acustic în limba română s-a folosit un corpus ce conține vorbire citită în condiții de liniște, dar și înregistrări ale unor emisiuni de televiziune și radio, unele dintre acestea fiind afectate de diverse zgomote de fundal (muzică, răsete, etc.). Corpusul a fost împărțit conform Tabelul 3. Fonemele au fost antrenate în comun, pentru limba engleză și română la un loc.

Set	Subset	Durata	Conținut
Antrenare	Omega train	94 h 46 m 14 s	Vorbire citită, fără zgomot
	TV train1 (TVDatabase)	27 h 27 m 21 s	
	TV train2 (TVautomaticAquisition)	103 h 17 m 56 s	
Testare	Omega test	5 h 29 m 6 s	Vorbire citită, fără zgomot
	TV test (News all)	3 h 29 m 3 s	

Tabelul 3 Corpusul audio în limba română

Evaluarea modelului acustic multilingual RO-EN, în Sphinx, s-a realizat pe bazele de date homeAuto2, homeAuto3 și homeAuto4.

Suplimentar, gramatica cu stări finite folosită în decodarea sistemului prototip conține și un model de tip "garbage", pentru rejecția cuvintelor din afara vocabularului, scopul fiind de a identifica comenzi greșite, sau propozițiile ce nu conțin comenzi. Astfel, în momentul în care o comandă nu este recunoscută, sistemul poate fi programat să sintetizeze un mesaj audio, prin care utilizatorului îi se cere rostirea unei comenzi corecte. Figura 2 evidențiază o porțiune din gramatica în limba română, ce conține această regulă de tip "garbage".

```

53 <numar_ac> = zece | unsprezece | doisprezece | douasprezece | treisprezece | paisprezece | șaisprezece
| șaptesprezece | optspreezece | nouăspreezece | douăzeci | douăzeci_ și_unu | douăzeci_ și_doi | douăzeci_ și_trei |
douăzeci_ și_patru | douăzeci_ și_cinci | douăzeci_ și_șase | douăzeci_ și_șapte | douăzeci_ și_opt | douăzeci_ și_nouă |
treizeci | treizeci_ și_unu | treizeci_ și_doi | treizeci_ și_trei | treizeci_ și_patru | treizeci_ și_cinci |
treizeci_ și_șase | treizeci_ și_șapte | treizeci_ și_opt | treizeci_ și_nouă | patruzeci;
54
55
56 //-----
57
58 <din_in_camere> = (din | în) (bucătărie | baie | dormitor | sufragerie | cameră);
59
60 //-----
61
62 public <comandaCasa> = (/0.0001/garbage*) | (/0.999/cassandra (<securitate> | <multimedia> | <hvac> | <electric>));

```

Figura 2 Gramatica cu stări finite, în limba română, ce conține "garbage" model.

Bibliografie

[Rousseau, 2012] A. Rousseau, P. Deléglise, and Y. Estève, "TED-LIUM: an automatic speech recognition dedicated corpus", in Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12), May 2012.

3.8 Activitatea 4.8. Integrarea algoritmului de căutare bazat pe DTW

Sistemul propus de căutare directă, în cadrul fișierelor audio, folosește un model acustic multilingv cu un algoritm scalabil de tip Dynamic Time Warping (DTW). Pentru a încerca rezolvarea problemei căutării direct în cadrul fișierelor audio, fără a cunoaște limba, am folosit un sistem de indexare cu o recunoaștere de fonem multilingvă. De asemenea, am experimentat cu mai multe tipuri de caracteristici, pentru a vedea care se comportă mai bine și în ce condiții, pentru această sarcină.

Implementarea finală a sistemului se bazează pe arhitectura propusă de NIST. Am decis separarea indexării și a modulelor de căutare, în detrimentul căutării directe în corpus a fiecărui termen al interogării, pentru a face căutarea mai rapidă, deoarece metoda de indexare simplifică căutarea [Buzo, 2013]. Astfel abordarea constă în două etape:

- Indexarea, adică recunoașterea de fonem pe baza de date audio, de conținut
- Căutarea, adică găsirea unui sir similar de foneme din conținutul indexat, care se potrivește cu cel al interogării, prin utilizarea unui algoritm de căutare DTW.

Recunoașterea fonemelor este folosită pentru indexare, astfel tot conținutul vorbit este transformat în siruri de foneme. Pentru limbile slab-dotate, unde resursele sunt limitate, recunoașterea de foneme este o alegere bună, deoarece nu avem suficiente date pentru a construi modele complexe de limbă de tip n-gram [Caranica, 2015].

Pentru experimente, am început cu două tipuri de caracteristici, pentru a reprezenta parametric vorbirea: MFCC și PNCC, pentru o robustețea îmbunătățită, teoretic, la zgomot. În ceea ce privește modelul acustic (AM), în prima abordare, folosind decodarea de foneme bazate pe HMM, ne-am dorit să comparăm efectul utilizării de resurse multilingve față de cele monolingve și pentru a realiza acest lucru, am construit șase modele acustice descrise în Tabelul 4. Am început cu modelele acustice pentru fiecare limbă, apoi am folosit clasificarea IPA pentru a combina fonemele comune în modelele de limbă (LM) antrenate cu toate datele. Combinarea fonemelor comune este motivată de numărul mare de foneme obținute în AM4, în cazul în care fonemele din diferite limbi sunt antrenate separat și sunt privite ca entități diferite. Acest lucru ne-a permis să reducem numărul de foneme la 98, pentru modelul în mai multe limbi AM5. Am folosit un număr moderat de date de antrenare pentru limbile individuale, pentru a avea o antrenare echilibrată între limbi. Pentru comparație, am antrenat și un model acustic suplimentar pentru limba română (AM6), cu o cantitate mare de date (64h) și un set relativ mic de foneme (34).

ID	Language	LM no. phonemes	Training data [h]
AM1	Romanian	34	8.7
AM2	Albanian	36	4.1
AM3	English	75	3.9
AM4	Multilingual separate phones	145	16.7
AM5	Multilingual common phones (IPA)	98	16.7
AM6	RomanianBig	34	64

Tabelul 4 Seturile de date de antrenare folosite pentru abordarea HMM, pentru algoritmul de căutare

Suplimentar, am încercat introducerea unor noi limbi în fază de antrenare, pentru a vedea modul în care acestea se comportă, împreună cu o rețea neuronală bazată pe recunoaștere de foneme, de la BUT [BUT, 2015]. Acest decodor de foneme utilizează o divizare dependentă de context (STC) pentru extracția de caracteristici, cu clasificatorii bazati pe rețele neuronale pentru a obține posteriogramă, în timp ce algoritmul Viterbi este folosi pentru decodare. Putem folosi ieșirea acestui instrument în algoritmul de căutare DTW ca și caracteristici de intrare, pentru a face căutarea. Cu scopul de a utiliza limbi suplimentare pentru a construi parametrii caracteristici pentru recunoaștere fonemelor, am folosit sistemele pre-antrenate disponibile la [BUT, 2015], descrise în Tabelul 5.

ID	Limba	# foneme in LM	WER[%]
AM7	Czech	45	24.24
AM8	Hungarian	61	33.32
AM9	Russian	52	39.27

Tabelul 5 Sistemele antrenate folosind STC

Limbile utilizate pentru antrenarea sistemelor descrise în Tabelul 5 fac parte din baza de date SpeechDat-E a limbilor vorbite în Europa de Est [Speechdat-E, 2015]. Un alt motiv pentru a folosi aceste sisteme o reprezintă existența unor simboluri non-vorbire antrenate, care s-ar dovedi utile în medii acustice dificile, cum ar fi:

- „int” pentru zgomot intermitent
- „spk” pentru zgomotul vorbitorului (tuse etc)

- „pau” pentru zonele de liniște

Metoda de căutare propusă folosește un algoritm de Dynamic Time Warping (DTW) pentru a alinia un sir (o interogare) într-un anumit conținut. Căutarea nu se face pe întregul conținut, ci numai pe o parte a acestuia prin intermediul unei ferestre glisante proporțională cu lungimea interogării, unde atât interogarea cât și conținutul în care se caută sunt un sir de foneme. Termenul este considerat detectat dacă scorul DTW trece de un prag. În plus, față de algoritmul clasic DTW, am inclus, în formula distanței, efectul lungimii interogării și a împrăștierii detecției DTW. Efectul lor este ponderat, cu scopul de a găsi o configurație optimă. Precizia sistemului RVC utilizat pentru indexare joacă un rol important, deoarece algoritmul de căutare trebuie să compenseze pentru rata destul de mare de eroare a fonemelor (*PhER*).

Deoarece lungimea conținutului este de obicei mai mare decât lungimea interogării, comparația se face cu ajutorul unei ferestre glisante a cărei lungime este proporțională cu lungimea interogării. Pentru fiecare fereastră, alinierea este dată de scorul *s*:

$$s = (1 - PhER) \quad (1)$$

unde *s* este un scor de similitudine. Detectia se bazează pe un prag, care se determină empiric.

Metoda de detecție este îmbunătățită prin introducerea unei penalizări pentru interogările scurte și împărtierea detecției DTW. Penalizările sunt motivate de presupunerea că, pentru două interogări de diferite lungimi care se potrivesc cu conținutul lor la aceeași rată de eroare a fonemelor (*PhER*), potrivirea interogării mai lungi este mai probabil să fie cea corectă. Formula pentru scorul *s* este acum dată de ecuația [Buzo, 2013]:

$$s = (1 - PhER)(1 + \alpha \frac{L_Q - L_{Qm}}{L_{QM} - L_{Qm}})(1 + \beta \frac{L_w - L_s}{L_Q}) \quad (2)$$

unde L_Q este lungimea interogării, L_{QM} și L_{Qm} sunt valorile maxime și minime pentru L_Q , L_w este mărimea ferestrei, L_s este împărtierea detecției DTW, iar α și β sunt coeficienții care controlează nivelul penalizărilor introduse. Penalizările introduse în ecuația de mai sus sunt motivate de presupunerea că, pentru două interogări de diferite lungimi care se potrivesc cu conținutul lor la aceeași rată de eroare a fonemelor (*PhER*), potrivirea interogării mai lungi este mai probabil să fie cea corectă. În mod similar, detectiile mai compacte DTW se presupun a fi mai probabile decât cele mai lungi.

Bibliografie

[BUT, 2015] BUT Phoneme recognizer based on long temporal context, Brno University of Technology, <http://speech.fit.vutbr.cz/software/phoneme-recognizer-based-long-temporal-context>, accessed august 2015.

[Buzo, 2013] Buzo A., Cucu H., Safta M., Burileanu C., “Multilingual Query by Example Spoken Term Detection for Under-Resourced Languages”, in Speech Technology and Human - Computer Dialogue (SpeD), (2013).

[Caranica, 2015] Alexandru Caranica, “Optimizations in spoken language recognition”, PhD Thesis, University “Politehnica” of Bucharest, Nov 2015 (scientific coordinator: prof. Cornelius Burileanu).

3.9 Activitatea 4.9. Implementarea infrastructurii și algoritmilor de detecție multilingvă a termenilor vorbiți

Sistemele de detecție a recurenței termenilor vorbiți (Spoken Term Discovery) identifică fragmente de vorbire repetitive din discursul ce conține vorbire, fără nici o informație despre limba acestuia [Park, 2008], pentru a construi clase de fragmente de vorbire similară. Abordările actuale pentru detecția recurenței termenilor vorbiți se bazează pe variante de tip Dynamic Time Warping (DTW), pentru a efectua în mod eficient o căutare într-un corpus de vorbire, cu scopul de a descoperi repetiții (clase similară) în vorbire (numite în continuare "termeni" sau "motive").

Metodologia propusă folosește un algoritm open-source de detecție a recurențelor audio, denumit MODIS [Catanese, 2013], bazat pe sistemul propus în [Muscarello, 2012]. Funcționarea sistemului urmează așa numitul principiu descoperire a „fragmentelor”, adică de a căuta un segment audio mai scurt într-un segment mai mare, căutarea fiind realizată prin intermediul unei „variante de avansare” prin segmente de tip „dynamic time warping” (DTW). În acest sistem segmentul mai scurt este numit „seed”, în timp ce fragmentul mai mare e denumit „buffer”.

Algoritmul inspectează secvențele acustice prezente în buffer pentru a evalua dacă acesta conține orice repetiție a fragmentelor, apoi o decizie de asemănare este luată prin comparația scorului DTW cu un prag similitudine DTW a unei librării de motive. Dacă nu s-au găsit asemănări între fragmente, procesul „sare” la următorul fragment din buffer, pentru a căuta asemănări cu librăriile de motive. Dacă o asemănare este găsită, atunci libraria de motive este actualizată sau „învățată”, în final obținându-se „clustere” de vorbire similară. Algoritmul continuă până la finalizarea căutării tuturor fragmentelor în librăriile disponibile, prin repetiția proceselor de mai sus. După cum am menționat anterior, modelul de cluster corespunzător unei librării de motive este actualizat continuu. După ce întregul corpus a fost analizat, motivele găsite sunt comparate între ele, în termeni de suprapunere, iar elemente suprapuse sunt unite într-o singură

librărie. Acest proces este ilustrat în Figura 3. MODIS a fost proiectat pentru a lucra cu mai multe tipuri de caracteristici și reprezentări parametrice a vorbirii. Astfel sunt implementate două tipuri de distanțe: distanța euclidiană, în general utilizată împreună cu reprezentările spectrale, iar distanța definită în ecuația de mai jos, pentru reprezentările cu probabilități posterioare:

$$d(a,b) = -\log \left(\sum_{i=0}^N (a_i \cdot b_i) \right) \quad (3)$$

Termenii implicați în ecuația de mai sus reprezintă cei doi vectori caracteristici pentru care se calculează distanța (a și b) și lungimea lor (N).

Pentru experimentele de căutare a recurențelor termenilor vorbiți, am variat pragul de similitudine, păstrând în același timp restul parametrilor constanți. Durata fragmentului a fost stabilită la 0.25 secunde, iar dimensiunea minimă a termenilor returnați e de 0.5 secunde, pentru a putea găsi cuvinte întregi. Lungimea bufferului e de 600 secunde, deoarece, în general, fișierele audio din bazele noastre de date pot avea mai puține repetiții ale aceluiași cuvânt. Modelul ales pentru a reprezenta clusterele termenilor a fost modelul median, în timp ce o matrice de verificare similitudini [Muscarello, 2011] a fost folosită pentru a potrivi modelul din cluster, de fragmentele analizate. În ceea ce privește distanțele utilizate, caracteristicile posteriorgram folosesc distanța prezentată în ecuația anterioară (3), în timp ce MFCC și vectorii de foneme utilizează distanța euclidiană.

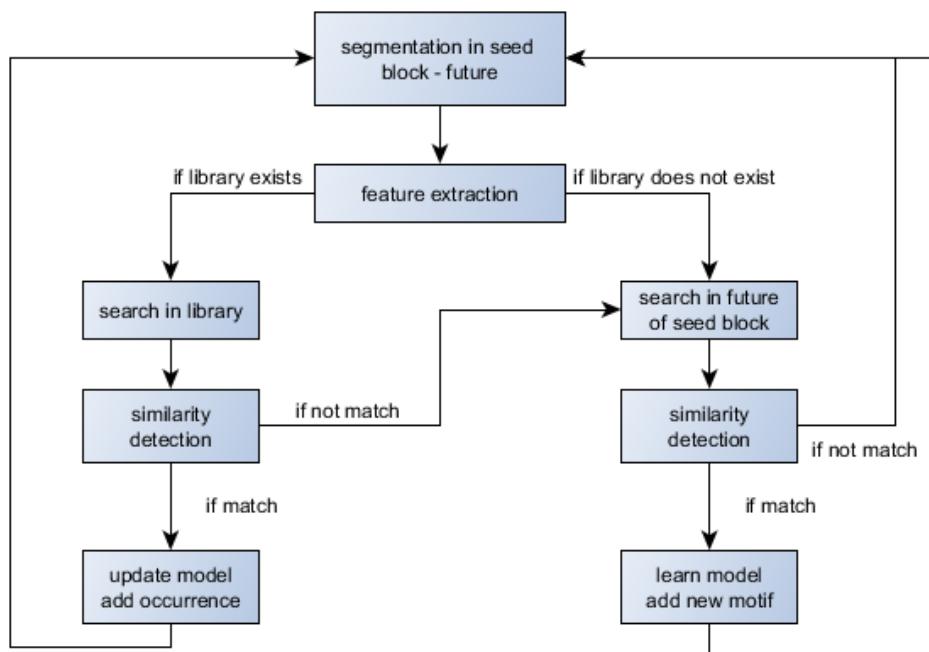


Figura 3 Arhitectura algoritmului MODIS, propus de [Catanese, 2013]

Pentru efectuarea sarcinii de detecție a cuvintelor cheie într-o secvență de vorbire, a fost creat un serviciu web ce poate fi accesat online la adresa <http://keywords.speed.pub.ro>. Acesta este parte componentă a unui ansamblu de servicii, a cărui arhitectură poate fi observată în Figura 4 Arhitectură SpeeD Web Services.

Din punct de vedere funcțional, această aplicație este compusă din următoarele componente:

- I. Interfață grafică (frontend) - permite accesul utilizatorului la mai multe facilități:
 - crearea și accesarea unui cont de utilizator, precum și editarea datelor personale;
 - încărcare a fișierelor audio sau text în contul de utilizator;
 - crearea unor cereri de transcriere a fișierelor audio existente în contul de utilizator;
 - crearea unor cereri de căutare a cuvintelor cheie ce se pot găsi în fișierele audio existente în contul de utilizator;
 - crearea unor cereri de restaurare a diacriticelor în fișierele text existente în contul de utilizator.
- II. Componente de management
 - Baza de date: stochează informații despre utilizatori, cât și date necesare efectuării proceselor descrise sau date rezultate în urma acestora;

- Manager de cozi: stochează cererile trimise de utilizatori prin API Gateway, precum și informații despre starea acestor cereri (efectuate, în curs de procesare).

III. Componente de procesare a datelor

- Diarizator: etichetează datele audio, oferind informații temporale asupra segmentelor rostite de către un vorbitor. De asemenea, oferă informații despre natura segmentelor audio (vorbire sau non-vorbire);
- Speech Transcriber: realizează transcrierea segmentelor audio, oferind textul aferent acestora. Este bazat pe sistemul de recunoaștere automată a vorbirii antrenat cu utilitarul Kaldi.
- Transcription Postprocessor: realizează procesarea textelor obținute în urma transcrierii. Acest proces include formatarea paragrafelor, inserarea semnelor de punctuație, inserarea majusculelor;
- Diacritics Restorer: realizează restaurarea diacriticelor în fișiere text;
- Keyword Spotter: realizează detecția unor cuvinte cheie introduse de utilizator în textul aferent transcrierii fișierelor audio.

Pentru a efectua o operație de detecție a cuvintelor cheie, utilizatorul va iniția o cerere în care va specifica fișierul audio în care se va face căutarea, precum și lista de cuvinte cheie. Cererea va fi introdusă în sistemul de cozi, iar de aici va fi preluată de către Diarizator, Speech Transcriber și Keyword Spotter.

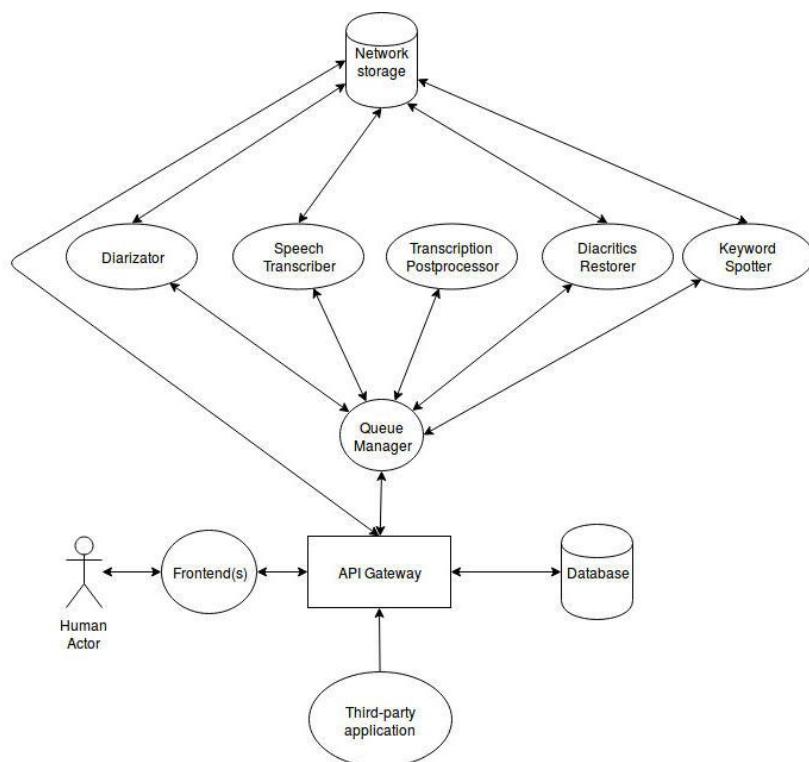


Figura 4 Arhitectură SpeeD Web Services

Pentru a efectua o operație de detecție a cuvintelor cheie, utilizatorul trebuie mai întâi să încarce un fișier audio în aplicație. Astfel, va fi efectuată o cerere de tipul HTTP POST către API Gateway pe terminalul /audio, care va conține tokenul obținut în urma autentificării, precum și octetii audio în format Base64. API Gateway va stoca fișierul audio în Network Storage, iar informații despre acesta vor fi înregistrate în baza de date, fiind asociat fișierului și un ID unic.

Mai departe, utilizatorul va iniția o nouă cerere de tip HTTP POST către API Gateway pe terminalul /kws, care va conține ID-ul fișierului audio în care se va face căutarea cuvintelor cheie, domeniul căruia aparține vorbirea din fișier (general, economic, politic, sportiv, etc.), precum și lista cuvintelor cheie ce se doresc a fi detectate. ID-ul fișierului, domeniul de vorbire, tipul sarcinii (detecția cuvintelor cheie), precum și statusul sarcinii, vor fi inserate în baza de date, în tabelul AudioJobs. În tabelul Keywords vor fi inserate cuvintele cheie.

API Gateway va pune o cerere de diarizare ce conține ID-ul fișierului audio în coada corespunzătoare componentei Diarizator. Aceasta va prelua fișierul audio din Network Storage, va realiza operația de diarizare și va pune informațiile rezultate în coada corespunzătoare componentei API Gateway. Componenta va prelua informațiile din coadă și va actualiza tabelul AudioSegments din baza de date, completând cu informații despre fiecare segment (etichete temporale, genul vorbitorului, etc.). Apoi, vor fi puse în coada corespunzătoare componentei Speech Transcriber, N cereri de transcriere, asociate celor N segmente. Cерерile conțin etichetele temporale ale fiecărui segment, precum și domeniul de vorbire, folosit ulterior pentru selecția modelului acustic potrivit. Componenta Speech

Transcriber va prelua fișierul audio din Network Storage, va transcrie fiecare segment, iar la final va pune transcrierea rezultată în coada corespunzătoare componentei API Gateway. Acesta va prelua informațiile din coadă și va actualiza tabelul Transcriptions din baza de date. Totodata, va pune în coada corespunzătoare componentei Keyword Spotter, N cereri de detecție a cuvintelor cheie. Cererile conțin transcrierile segmentelor și ID-urile cuvintelor cheie. Keyword Spotter va prelua aceste informații din coadă și va căuta fiecare cuvânt cheie în transcrierile furnizate. La final, va pune etichetele temporale între care se găsesc cuvintele cheie detectate, în coada corespunzătoare componentei API Gateway. Aceasta va citi informațiile din coadă, va actualiza tabelul KeywordPositions, va actualiza tabelul AudioJobs, modificând statusul sarcinii de detecție a cuvintelor cheie, acesta fiind acum terminat. De specificat însă faptul că la început, se va verifica dacă au fost efectuate în trecut operațiile de diarizare sau transcriere pentru fișierul audio în cauză. În această situație, se va sări peste pașii respectivi.

Pentru a obține rezultatul final, utilizatorul va iniția o cerere de tip HTTP GET către API Gateway pe terminalul /kws/{job_id}, primind în răspuns o listă a cuvintelor cheie detectate, împreună cu momentele de timp la care acestea apar în fișierul audio.

Bibliografie

[Catanese, 2013] L. Catanese, N. Souviraà-Labastie, B. Qu, S. Campion, G. Gravier, E. Vincent, and F. Bimbot, “MODIS: an audio motif discovery software,” in Proc. of INTERSPEECH, 2013.

[Park, 2008] A. Park and R. Glass, “Unsupervised pattern discovery in speech,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 16, no. 1, pp. 186–197, 2008.

[Muscariello, 2012] A. Muscariello, G. Gravier, and F. Bimbot, “Unsupervised motif acquisition in speech via seeded discovery and template matching combination,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 20, no. 7, pp. 2031–2044, 2012.

3.10 Activitatea 4.10. Compararea și evaluarea performanțelor tehnicilor implementate

Această activitate prezintă rezultatele decodării bazei de date, achiziționate în limba Engleză (homeAuto4), în cadrul grupului de lucru, folosind modelele acustice antrenate pe corpusul TED-LIUM, în cele două utilitare integrate de către grupul SpeeD: CMU Sphinx și Kaldi. Se compară modelele folosind metrica standard WER (Word Error Rate).

Sistem	Configuratie	WER [%]
Kaldi	Tri 1	55
Kaldi	Tri 2 (LDA+MLLT)	54
Kaldi	Tri 3 (LDA+MLLT+SAT)	18
Kaldi	Tri3 (MMI)	17
Kaldi	NNET2	16
Sphinx	GMM-LM	49
Sphinx	GMM-FSG	49

Tabelul 6 Evaluarea performanțelor sistemului în limba Engleză

Se poate observa, cu ușurință, că utilitarul Kaldi oferă performanțe mult mai bune, cel puțin în limba engleză, datorită următorilor factori:

- Antrenarea pe baza de rețele neuronale (Kaldi), ceea ce a oferit un spor de performanță pentru baza de date achiziționată în limba engleză, după cum se observă. Avantajul major este reprezentat de existența algoritmilor de antrenare ce folosesc rețele neuronale adânci (DNN), ce pot capta mai bine variațiile frazelor rostite de vorbitorii non-nativi de limba engleză.
- Kaldi realizează și o adaptare la vorbitor („speaker adaptation”), pentru a capta variațiile din cadrul vorbirii, cum ar fi accentele non-native, stilul discursului, etc.
- CMU Sphinx nu poate astfel recupera diferențele majore în vorbire, dintre baza de date de antrenare și cea de decodare (pronunție, stil, accentul limbii, etc).

Pentru această evaluare, s-a folosit o gramatică cu stări finite (FSG – Finite State Grammar), pentru ambele toolkit-uri, în cadrul sistemului prototip. Figura 5 prezintă o porțiune din acest model FSG, ce conține, ca și în limba română, un model de tip “garbage” pentru rejecția propozițiilor ce nu conțin comenzi.

```

grammar functii;

//-----
<multimedia> = <multimedia_on_off> | <multimedia_on_off_1> | <volume_settings> | <volume_change_1> | <radio_selection> | <source>;
<multimedia_on_off> = (start | stop | turn off | play) ([the] music |[the] radio |[the] tv);
<multimedia_on_off_1> = (play | run) ([the] music list | [the] music playlist | [the] playlist);

<volume_settings> = (set | change | make) (volume | sound) ([the] tv | music) to <decimal_number> percent;
<volume_change_1> = (decrease|increase|lower|make [the]) (volume | tv [quieter] | tv [volume]);

<radio_selection> = (i want to listen [to] | play) (radio europa fm | radio zu | radio pro fm);
<source> = (change|switch) (source | to hdmi | the tv | plug-in) [(source | source to hdmi | [source] to vga | source to plug-in)];

//-----
<exterior> = <sprinklers> | <front_gate> | <garage_door> | <front_door> | <modes>;
<sprinklers> = (start|stop|activate|deactivate|turn) ( [the] sprinklers | [the] irrigation | [the] irrigation system [on]);
<front_gate> = (lock|unlock|close|open) (the entrance gate | entrance);
<garage_door> = (lock|unlock|close|open) (the garage door | garage door | the garage);
<front_door> = (lock|unlock|close|open) ([the] front door | [the] entrance);

<modes> = (start|activate|stop|deactivate) (reading|read|dinner|movie) mode;

```

Figura 5 Gramatica cu stări finite, în limba Engleză

Pentru CMU Sphinx s-a realizat și o etapă de “tuning” a parametrilor OOGP (Out-of-grammar Probability), RBW (Relative Beam Width) și WIP (Word Insertion Probability), pentru gramatica de tip FSG.

3.11 Activitatea 4.11. Crearea unui prototip de sistem de recunoaștere multilingvă a vorbirii

Prototipul sistemului de recunoaștere multilingvă a vorbirii constă într-un kit de dezvoltare Raspberry PI 3, unde cu ajutorul utilitarului Sphinx, a fost implementat un sistem mixt, ce realizează atât detectia cuvintelor cheie (DCC), cât și recunoașterea automată a vorbirii (RAV).

Odată cu alimentarea kitului, este pornit un proces de DCC ce preia date audio de la microfon și încearcă să stabilească dacă printre acestea apare cuvântul cheie „Casandra”. Atunci când cuvântul cheie a fost detectat, va începe un proces de RAV ce utilizează o gramatică cu stări finite, care va stabili tipul comenzii rostită de utilizator. Comenzile sunt cele definite în etapele anterioare ale proiectului.

Integrarea modulului de recunoaștere a vorbirii cu controlerul de voce se realizează prin trimiterea către acesta din urmă a unei cereri de tip HTTP POST, ce conține transcrierea comenzii. Controlerul de voce înțelege comanda primită și o efectuează fizic (aprindere lumina, modifică temperatură, etc.), iar apoi întoarce în răspunsul HTTP un text ce informează asupra efectuării cu succes a comenzii.

Integrarea modulului de recunoaștere a vorbirii cu sistemul de sinteză a vorbirii se realizează prin trimiterea către acesta din urmă a unei cereri de tip HTTP GET, ce conține textul întors de către controlerul de voce. Acest text va fi sintetizat, obținându-se un fișier audio care va fi întors în răspunsul HTTP.

În final, sistemul implementat pe Raspberry PI redă la difuzor fișierul audio sintetizat, informând utilizatorul despre efectuarea comenzii. Apoi, face comutarea de la modul de funcționare RAV la modul de funcționare DCC. Astfel, sistemul așteaptă ca utilizatorul să rostească din nou cuvântul cheie „Casandra” urmat de o nouă comandă vocală; tot procesul descris mai sus se reia. Acest proces este ilustrat în Figura 6 Organograma sistemului DCC și RAV.

Pentru implementarea procesării multilingve, au fost propuse două variante. Varianta 1 presupune:

- crearea și utilizarea mai multor modele (acustice, fonetice și lingvistice) distincte pentru fiecare limbă în parte
- existența unui mecanism de selecție statică a limbii de interacțiune cu sistemul
- comutarea manuală pe limba dorită și apoi interacțiunea cu sistemul în acea limbă

Concret, această variantă implică conectarea la kitul Raspberry PI a unui buton cu două stări, ce are rolul de a selecta limba dorită (română sau engleză). În funcție de poziția butonului, modulul de recunoaștere automată a vorbirii va porni având încărcate modelele corespunzătoare uneia dintre limbi. Ulterior, dacă se dorește interacțiunea cu sistemul în celalătă limbă, este necesară comutarea manuală a butonului de selecție.

Varianta 2 presupune crearea și utilizarea antrenarea unei singur set de modele (acustice, fonetice și lingvistice) folosind o bază de date de vorbire multilingvă. Astfel, fonemele sunt antrenate la comun, a cum a fost descris în Activitatea 4.7. Dezvoltarea unor modele acustice multilingve și pentru limba engleză (partea II) Modelul de limbă de tip gramatică cu stări finite se obține prin concatenarea regulilor corespunzătoare gramaticilor fiecărei limbi.

Varianta 1 are avantajul de a fi mai simplu de gestionat la momentul dezvoltării. Practic, pentru fiecare limbă în parte se crează modele acustice, fonetice și lingvistice distincte, utilizând baze de date de vorbire monolingve și cunoștințe lingvistice specifice limbii în cauză. La limită, în cazul în care sunt multe limbi, dezvoltarea se poate realiza de echipe separate (fiecare fiind responsabilă de o singură limbă). Evaluarea acestor sisteme se face de asemenea mult

mai simplu (cu baze de date de vorbire monolingve). Un alt avantaj major al acestei variante este acela că acuratețea sistemului nu depinde de numărul de limbi suportate și este fără îndoială mai mare decât în cazul variantei 2. Dezavantajul principal al acestei variante este reprezentat de mecanismul manual de selecție a limbii, mecanism ce ar putea părea rudimentar într-un sistem intelligent precum cel prezentat mai sus. Cu toate acestea, în contextul curent (casă inteligentă cu care se poate interacționa prin voce) este puțin probabil ca schimbarea limbii de interacționare să fie realizată prea frecvent. Cazul real de funcționare este acela în care limba este configurată în momentul instalării sistemului (dintr-o listă de limbi posibile), iar utilizatorul interacționează întotdeauna cu sistemul folosind acea limbă.

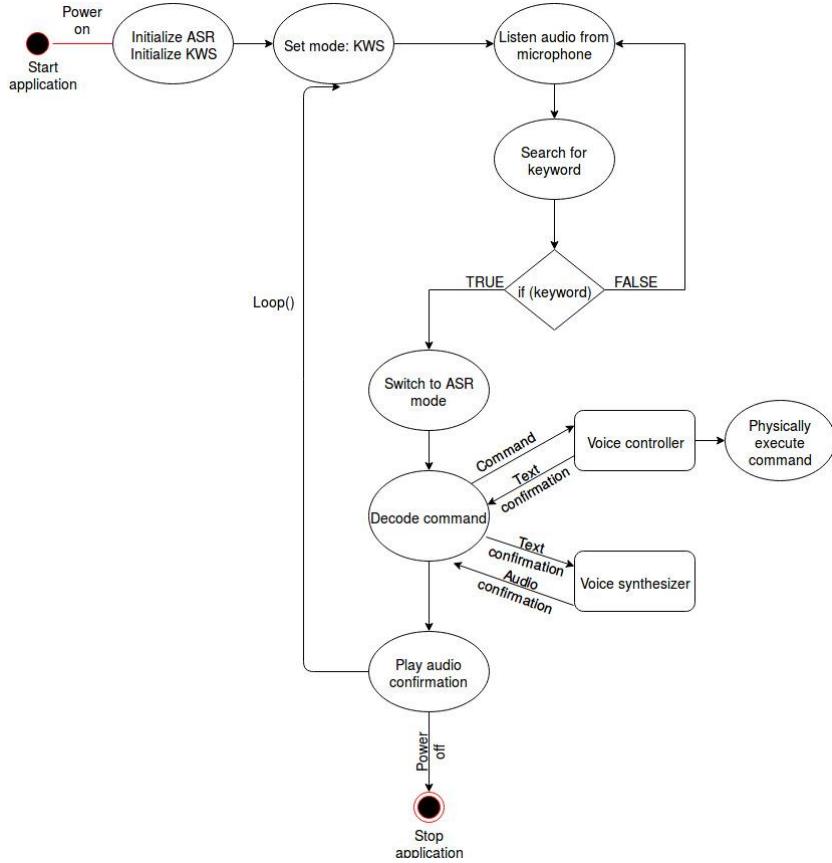


Figura 6 Organigrama sistemului DCC și RAV

Varianta 2 are avantajul de a fi mai scalabilă și mai simplu de gestionat la momentul utilizării pentru că întotdeauna va exista un singur model acustic, un singur model fonetic și un singur model lingvistic (indiferent cătă limbi va permite sistemul). Dezavantajul major al acestei variante este acela că acuratețea sistemului (acuratețea de recunoaștere a comenziilor) este mai mare. Motivul este simplu de înțeles: cu cât sunt mai multe comenzi posibile, în mai multe limbi posibile, cu atât probabilitatea de confuzie a unei comenzi cu o alta crește.

În prototipul realizat a fost implementată varianta 1 de procesare multilingvă a vorbirii.

Întregul sistem de recunoaștere multilingvă a vorbirii a fost implementat și într-o variantă demonstrativă, fizic funcțională ("demo bench"), ce a constat în configurația și integrarea următoarelor dispozitive și periferice:

- Kit Philips Hue, ce conține 3 becuri LED de 10W A19 și E27, controlabil prin intermediul unui API de tip REST (GET/POST)
- Sursă KNX + Router IP KNX, pentru controlul unor sisteme de automatizare clădiri, cum ar fi un termostat
- Termostat KNX cu afișaj digital, pentru reglarea automată a temperaturii
- Intel NUC N3700, cu 8GB ram și SSD de 120GB, pentru partea centrală (server) a aplicației demonstrative
- Router WiFi Mikrotik RB951Ui-2HnD, cu suport USB, ce permite conectarea la modem-uri 3G / 4G, pentru a oferi conectivitate IP sistemului, independent de locație
- Două Raspberry Pi 3 Model B, ce rulează aplicația de recunoaștere a cuvintelor cheie "Casandra" și a comenziilor, ce le trimit prin intermediul unor variabile POST către controller
- Boxe pentru Raspberry Pi, cu scopul de a reda comenzi sintetizate
- Două microfoane USB, pentru captarea sunetului

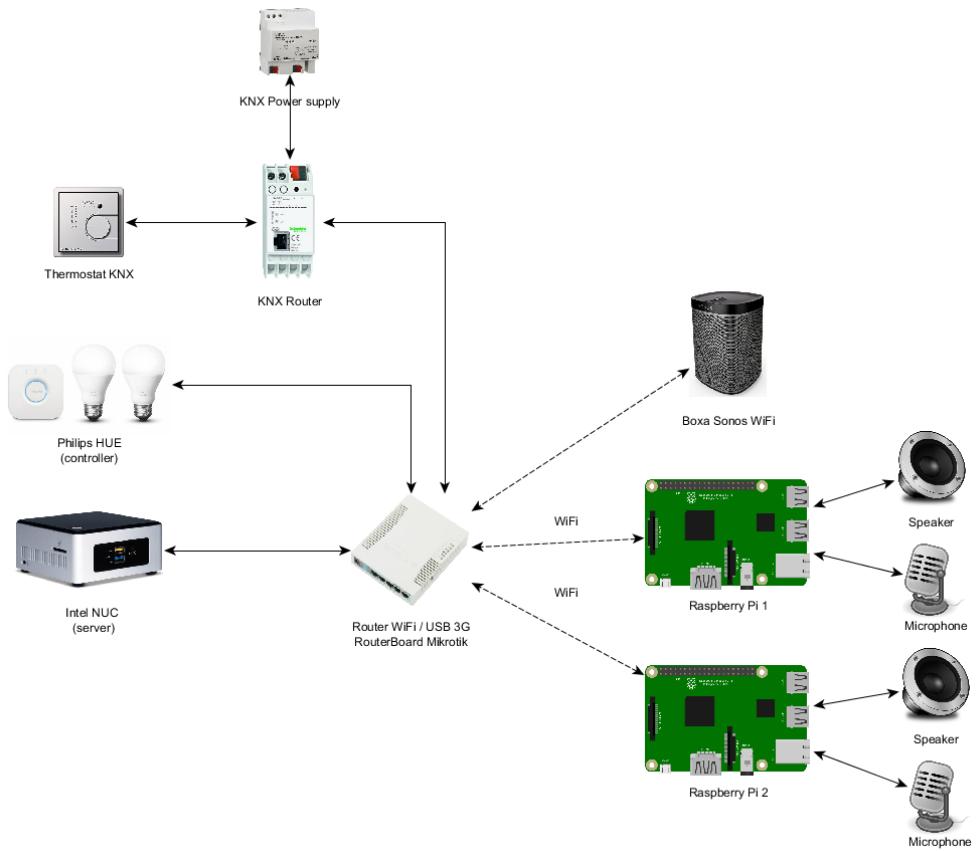


Figura 7 Arhitectura prototipului hardware, demonstrativ

Scopul final al acestui prototip hardware este de a demonstra, într-un spațiu cât mai compact și cu un consum de energie cât mai mic, un sistem intelligent de control al unei încăperi, prin intermediul vorbirii, în două limbi (română + engleză). Figura 7 prezintă arhitectura prototipului hardware și modul de conectare al dispozitivelor și perifericelor.

Partea centrală (serverul) o reprezintă mini PC-ul Intel NUC, pe care rulează proxy-ul de voce (ComfortClick), sub sistem de operare Windows. Sintetizatorul de vorbire (modulul TTS) și controllerul de voce (modulul NLP) sunt integrate într-o mașină virtuală Linux (Ubuntu Server 16.04 LTS), virtualizată cu ajutorul proiectului open-source Oracle VirtualBox. Tot în această mașină virtuală se află și modelele acustice, de limbă, pentru sintetizare, etc. Routerul WiFi Mikrotik, pe lângă asigurarea conectivității IP la internet, oferă și parte de switch și routare pentru toate dispozitivele TCP/IP din Figura 7. Dispozitivele Raspberry Pi realizează detectia de cuvinte cheie și recunoașterea comenzielor rostită de utilizator și comunică, prin intermediul REST, cu controllerul de voce și sintetizatorul de vorbire. Din punct de vedere hardware la dispozitivele Raspberry Pi sunt conectate microfoane pentru preluarea comenzielor și difuzoare pentru redarea răspunsurilor furnizate de casa inteligentă. Proxy-ul de voce (ComfortClick) asigură controlul, tot prin comenzi de tip REST, a becurilor inteligente Philips HUE, iar router-ul IP KNX controlează temperatura ambientală, prin intermediul termostatului, tot KNX.

Astfel, acest sistem prototip simulează o întreagă arhitectură "inteligentă", controlată prin voce și vorbire naturală.s

3.12 Activitatea 4.12. Crearea unui sistem de sinteza a vorbirii pentru limba română

Sistemul de sinteză a vorbirii dezvoltat în cadrul proiectului ANVSIB este unul bazat pe corpus, având la bază (a) un sistem de analiză a textului, responsabil pentru extragerea de informații lingvistice relevante pentru sinteză vorbirii (explicații în cele ce urmează) și (b) un sistem de sinteză a vorbirii pornind de la informațiile lingvistice obținute de sistemul de analiză a textului. Extragerea de informații lingvistice relevante pentru sinteza vorbirii este un proces complex, prin care se urmărește conversia unei propoziții/paragraf într-o secvență de atrbute discrete ce servește procesului de selectare/generare de unități acustice necesare în procesul de generare de voce. Pornind de la sistemul de analiză a textului deja existent, în cadrul acestei activități, RACAI a revizuit întreg lanțul de procesare și a introdus un set nou de unelte de prelucrare a limbajului natural având ca scop îmbunătățirea modelelor deja existente, micșorarea dimensiunii modelelor și reducerea complexității de calcul. Astfel, noul set de instrumente folosește un perceptron cu decodator Viterbi pentru etichetare morfosintactică și o cascadă de clasificatori bazați pe arbori de decizie pentru

stabilirea atributelor morfolexicale. Întreg lanțul de prelucrare a fost testat pe mai mult de 50 de limbi în cadrul competiției CONLL Shared Task on Universal Dependencies Parsing, rezultatele fiind disponibile online¹.

De asemenea, RACAI a lucrat la îmbunătățirea și extinderea sintetizatorului vocal, mai multe detalii fiind oferite în cadrul secțiunilor (3.13, 3.14 și 3.15).

3.13 Activitatea 4.13. Crearea unui sistem de sinteza a vorbirii pentru limba engleză

Modelul de analiză a textului pentru limba engleză are la bază modelele antrenate în cadrul competiției CONLL, având în plus un sistem de transcriere fonetică pentru limba engleză ce a fost antrenat folosind lexiconul CMUDict. Modelul acustic folosit în experimentele noastre a fost pus la dispoziție în cadrul Blizzard Challenge 2016, modelele acustice putând fi folosite în condițiile impuse de organizatori.

3.14 Activitatea 4.14. Cercetarea și implementarea unui nou modul de prozodie

Generarea elementelor prozodice este în continuare un proces extrem de complicat în sinteza vorbirii pornind de la text. Principalul impediment constă în faptul că, pe lângă relația de suprafață dintre text și elementele prozodice, vorbitorii umani folosesc prozodia atât pentru a transmite propriile stări emoționale către interlocutori cât și pentru a facilita înțelegerea mesajelor complexe prin scoaterea în evidență a unor cuvinte cu încărcătură cognitivă sau a celor cuvinte care ar putea fi înțelese greșit de către interlocutori din motive subiective sau obiective (de exemplu cuvinte care se pronunță asemănător sau cuvinte rare ce nu sunt utilizate în vorbirea de zi cu zi).

Astfel, principala strategie în sinteza prozodică, în lipsa unor reguli clare și bine studiate de către foneticieni, este folosirea unui set exhaustiv de atribute extrase din text, ce pot contribui la construirea unei linii melodice corelată cu o cadență variabilă în rostirea cuvintelor ce ține în special de stilul particular al vorbitorului folosit în crearea corpusului. Conform rezultatelor raportate în cercetare precum și a propriei experiențe în realizarea de sisteme de sinteză a vorbirii, setul de trăsături folosit de noi este:

- (a) Contextul fonetic: fonemul curent, însăși de două foneme aflate la dreapta și două foneme aflate la stânga (însoțite de atribute specifice: consoană/vocală, mod de articulare etc.), precum și poziția fonemului în silabă (început, final, mijloc)
- (b) Silaba curentă: însăși silaba, în cazul în care aceasta este frecventă (mai mult de 5² aparții în corpus), accent lexical (da/nu), a câtei silabe din cuvânt, a câtei silabă din propoziție, numărul de silabe din cuvânt, numărul de silabe din propoziție, numărul de silabe de la ultimul semn de punctuație până la poziția curentă, numărul de silabe până la următorul semn de punctuație
- (c) Următorul semn de punctuație/tipul propoziției (declarativă, exclamativă, interrogativă)
- (d) Informații morfosintactice: poziția cuvântului curent în grupuri sintactice (Ion, 2007) (început, mijloc, final), partea de vorbire a cuvântului curent.

Stabilirea modului în care aceste atribute influențează prozodia se face în mod automat folosind arbori de decizie construți pe baza principiului Minimum Description Length (MDL). Astfel, sistemul construiește 3 arbori de decizie, având ca scop modelarea spectrală, a frecvenței fundamentale a vocii precum și a duratei de rostire. Analiza se face astfel:

- (a) În prealabil se folosește un aliniator la nivel de fonem pentru a obține corespondența între text și vocea înregistrată (în cazul nostru am utilizat HTK)
- (b) Fiecare fonem este ulterior modelat folosind un Modele Markov cu 5 stări (doar trei emit)
- (c) Parametrii spectrali (MLSA sau STRAIGHT), durata și frecvența fundamentală sunt extrase din fișierele acustice
- (d) Pe baza alinierilor se construiesc arborii de decizie, având ca noduri întrebări referitoare la atributele lingvistice (de exemplu „Fonemul curent este ,A’: DA/NU) și în frunze media și deviația standard calculate pe toate eșantioanele acustice care ajung în frunză în faza de antrenare
- (e) În timpul rulării sistemul generează automat o secvență de durate/frecvență fundamentală și spectru, prin parcurgerea arborilor folosind atribute lingvistice extrase din propoziția ce trebuie sintetizată

3.15 Activitatea 4.15. Integrarea cu HTS

În cadrul acestei activități RACAI a dezvoltat un modul de sinteză parametrică statistică ce suportă atât analiză de tip MLSA cât și STRAIGHT. Modelarea spectrală, a duratei de rostire și a frecvenței fundamentale a vocii este realizată în prezent folosind arborii de decizie descriși anterior, menționând însă faptul că orice altă metodă poate și integrată cu sistemul actual (în experimentele noastre am folosit o rețea bidirecțională de tip LSTM pentru generarea frecvenței fundamentale, model care însă nu a fost folosit în sistemul final, datorită complexității crescute de calcul).

¹ <http://universaldependencies.org/conll17/results.html>

² Valoare stabilită experimental

Sistemul este integral scris în limbajul de programare JAVA, fără dependințe externe de funcționare, cu mențiunea că, pentru generarea arborilor de decizie folosim sistemul HTS.

3.16 Activitatea 4.16. Optimizarea și configurarea sistemului de sinteza a vorbirii

În această etapă am experimentat cu un număr mare de atribute lingvistice și praguri de includere a acestora în modelele prozodice, având ca scop obținerea unor rezultate performante atât din punct de vedere al calității vocii (naturalitate și inteligență) cât și din punct de vedere al cerințelor computaționale:

- (a) modelele noastre sunt reduse ca dimensiunea datorită utilizării modelelor parametrice statistice
- (b) oferă rezultate stabile, fără degradări în calitatea vorbirii și fără a introduce artefakte sonore neplăcute
- (c) complexitate redusă: cele mai mari cerințe de calcul vin din partea etichetării morfosintactice ($O(k^2n)$) și din procesul de generare a vocii folosind STRAIGHT (motiv pentru care sistemul poate rula și pe baza filtrului MLSA)
- (d) Pachetele utilizate în momentul rulării nu necesită dependențe externe, putând fi folosite în orice mediu de producție (inclusiv sisteme ARM)
- (e) Cerințele de memorie sunt extrem de mici, sistemul putând rula sinteză MLSA în doar 16MB RAM

Dimensiunea totală a modelelor este redată în Tabelul 7.

Tabelul 7 Dimensiunea modelelor de sinteză a vorbirii

Limbă	Dimensiune model PLN	Dimensiune model acustic
Engleză	257MB	7.5MB
Română	2.9MB	3.8MB

3.17 Sumar al realizărilor / rezultatelor

Nr.	Denumirea rezultatului	Denumire activitate	Procent de realizare
1.	D4.1 Modele acustice multilingve și pentru limba engleză	Activitatea 4.7. Dezvoltarea unor modele acustice multilingve și pentru limba engleză (partea II)	100%
2.	D4.2 Model funcțional de detectie multilingva a termenilor vorbiti	Activitatea 4.9. Implementarea infrastructurii și algoritmilor de detectie multilingva a termenilor	100%
3.	D4.3 - Prototip de sistem de recunoaștere multilingva a vorbirii	Activitatea 4.11. Crearea unui prototip de sistem de recunoaștere multilingva a vorbirii	100%
4.	D4.4 - Model funcțional de sinteza a vorbirii pentru limba romana	Activitatea 4.12. Crearea unui sistem de sinteza a vorbirii pentru limba romana	100%
5.	D4.5 - Model funcțional de sinteza a vorbirii pentru limba engleză	Activitatea 4.12. Crearea unui sistem de sinteza a vorbirii pentru limba engleză	100%
6.	D4.6 - Prototip multilingv de sinteza a vorbirii	Activitatea 4.16. Optimizarea și configurarea sistemului de sinteza a vorbirii	100%
7.	D4.7 - Prototip de camera inteligenta controlabil prin voce	Activitatea 4.5. Evaluarea prototipului	100%