

Analiza semantică stocastică în dialogul om-calculator

Proiect de diplomă

Dan Bohuş

Iunie 2000

Departamentul de Calculatoare
Facultatea de Automatică și Calculatoare
Universitatea "Politehnica" din Timișoara
România

Conducător proiect:
ș.l. ing. Marian Boldea

**In order to understand something,
you mustn't try to understand everything.**

Weinberg – Lump Law

Cuvinte cheie: sistem de dialog, analiză semantică stocastică, modele Markov, algoritmul Viterbi, estimare de probabilitate maximă, reestimare Katz, antrenare incrementală

Rezumat

Această lucrare prezintă proiectarea, implementarea, dezvoltarea și evaluarea unui analizor semantic pentru sisteme automate de dialog vocal om-calculator, parte a unui efort mai larg de cercetare în acest domeniu.

Sistemele de dialog constituie un pas important pe drumul spre realizarea unei interfețe ideale om-calculator, iar în cadrul lor analiza semantică a limbajului natural joacă un rol fundamental: acela de a extrage și clarifica informația transmisă de utilizator și de a genera o reprezentare bine formalizată a acesteia.

În urma unui studiu al principalelor alternative de proiectare existente, a fost aleasă varianta unui analizor semantic stocatic bazat pe formalismul case-grammar. Pornind de la arhitectura analizorului semantic proiectat, sunt prezentate succesiv componentele acestuia, atât în privința suportului teoretic, cât și a soluțiilor de implementare.

Sunt ilustrate apoi modul în care implementarea rezultată poate fi folosită și performanțele posibile, prin intermediul unor experimente vizând dezvoltarea și evaluarea unui analizor semantic ce va fi integrat într-un sistem de dialog pentru informații despre orarul Departamentului de Calculatoare.

În final se prezintă concluzii desprinse din proiectarea, implementarea, dezvoltarea și evaluarea analizorului semantic, precum și posibile direcții de continuare a cercetărilor.

Cuprins

Lista tabelelor	5
Lista figurilor	7
1 Introducere	9
1.1 Sisteme de dialog	9
1.2 Structura tipică a unui sistem de dialog	11
1.3 Rolul analizorului semantic	12
1.4 Organizarea lucrării	13
2 Analiza semantică în sistemele de dialog	15
2.1 Analiza semantică	15
2.1.1 Definirea noțiunii de analiză semantică	15
2.1.2 Nivele de analiză semantică	16
2.1.3 Legături între analiza semantică și cea sintactică	17
2.2 Alternative în proiectarea unui analizor semantic	17
2.2.1 Formalisme de reprezentare a cunoștințelor	18
2.2.2 Metode de parsing	22
2.3 Stadiul actual	24
3 Un analizor semantic stocastic	27
3.1 Structura analizorului semantic	27
3.2 Modulul de preprocesare	29
3.3 Modelul Markov	31
3.3.1 Definiție. Probleme fundamentale	32
3.3.2 Modelele Markov în analiza semantică	34
3.3.3 Antrenarea modelului	34
3.3.4 Decodarea semantică	39
3.4 Constructorul de cadre	42
3.5 Detalii de implementare	44
3.5.1 Modelarea datelor	45
3.5.2 Implementarea algoritmilor	47
3.5.3 Programe utilitare	49

4	Experimente și rezultate	51
4.1	Domeniul și capabilitățile sistemului de dialog	51
4.2	Colectarea și pregătirea datelor	52
4.3	Dezvoltarea analizorului semantic	55
4.3.1	Metoda incrementală (bootstrap)	55
4.3.2	Erori și evaluarea performanțelor	56
4.3.3	Etape și evaluări intermediare	56
4.4	Evaluarea finală a analizorului	64
5	Concluzii și continuări	67
5.1	Concluzii	67
5.2	Direcții de continuare	68
	Bibliografie	71
	Anexe	73
A	Exemple de dialoguri	73
B	Structura unui corpus	77
C	Structura bazei de date orar	81
D	Fișierele de configurare	83
D.1	Controlul preprocesării	83
D.2	Controlul constructorului de cadre	88
E	Extrase de cod C++	91
E.1	Antete de clase	91
E.2	Algoritmi	111
F	Invocarea programelor utilitare	119

Lista tabelelor

1.1	Câteva exemple de sisteme de dialog	10
3.1	Pașii de preprocesare	30
4.1	Scenarii utilizate pentru colectarea datelor	53
4.2	Statistici referitoare la datele de antrenament	54
4.3	Evaluări performanțe – versiunea I	59
4.4	Evaluări performanțe – versiunea II	62
4.5	Performanțele versiunilor finale (III) pe corpusul de antrenament	64
4.6	Statistici referitoare la datele de test	65
4.7	Rezultatele evaluării finale a performanțelor	65
C.1	Tabela de săli	81
C.2	Tabela de materii	81
C.3	Tabela de profesori	81
C.4	Tabela orar	82

Lista figurilor

1.1	Exemplu de dialog cu CMU Communicator	11
1.2	Structura unui sistem de dialog pentru obținerea de informații . . .	12
2.1	Exemplu de cadru	20
2.2	Un exemplu de sistem de cadre	21
2.3	Structura unui analizor semantic bazat pe reguli	22
2.4	Structura unui analizor semantic stocastic	23
3.1	Arhitectura analizorului semantic	28
3.2	Exemplu de model Markov	33
3.3	Utilizarea modelelor Markov în decodarea semantică	34
3.4	Estimarea de probabilitate maximă	36
3.5	Principiul reestimării Katz	39
3.6	Algoritmul Viterbi	41
3.7	Pseudocod Viterbi	41
3.8	Correspondența succesiune de etichete semantice – cadru	42
3.9	Formulări fără concept	43
3.10	Clase pentru modelarea datelor	46
3.11	Clase pentru implementarea algoritmilor	47
4.1	Reducerea dicționarului prin preprocesare – versiunea I	59
4.2	Reducerea dicționarului prin preprocesare – versiunea II	61
4.3	Reducerea dicționarului prin preprocesare – versiunea III	63

Capitolul 1

Introducere

De-a lungul scurtei istorii a calculatoarelor, interfețele dintre acestea și utilizatori au cunoscut o evoluție continuă, marcată de diferite paradigme, dar guvernată permanent de necesitatea de a optimiza lucrul cu calculatorul prin eficientizarea unei componente esențiale – comunicarea cu acesta.

După panourile de comandă și comutatoarele caracteristice primelor calculatoare, paradigma care s-a impus în domeniul interfețelor cu utilizatorul a fost cea a liniei de comandă și a interfețelor-consolă. Au urmat apoi interfețele grafice 2D și paradigma point-and-click care și în acest moment sunt dominante. În prezent, cercetările în domeniul interfețelor bazate pe vorbire au condus deja la realizări impresionante, și migrarea acestor interfețe de la stadiul experimental la cel comercial a început deja. Numeroase cercetări se fac și în domeniul interfețelor 3D, bazate pe imersiune în realitate virtuală, și chiar a celor controlate prin gândire.

Cu toate acestea suntem încă departe de realizarea unei interfețe ideale. O astfel de interfață ar trebui însă să fie prin excelență una multimodală, adaptată perfect la stilul de comunicare uman. Ea ar trebui să accepte intrări direct în limbaj natural, eventual să sesizeze chiar și gesturile și mimica utilizatorului, și să prezinte ieșirile în forma cea mai adecvată și mai eficientă pentru perceperea lor de către om: grafică, scrisă, vorbită, animație, secvențe video, sunet etc.

Sistemele automate de dialog vocal om-calculator, denumite în cele ce urmează *sisteme de dialog*, constituie un pas important în această direcție. Multe din problemele, tehnicile și soluțiile utilizate în dezvoltarea sistemelor de dialog, precum recunoașterea vorbirii, analiza semantică, modelarea și controlul dialogului, sinteza vorbirii, vor fi utilizate în proiectarea unor astfel de interfețe multimodale. În plus, valoarea comercială intrinsecă a sistemelor de dialog justifică și stimulează cercetările în acest domeniu.

1.1 Sisteme de dialog

Întrucât nu există o definiție unanim acceptată a noțiunii de *sistem de dialog*, vom porni de la următoarea caracterizare [?]:

Sistemele interactive bazate pe dialog vor fi caracterizate ca sisteme de calcul ce permit oamenilor să îndeplinească cel puțin o parte din sarcini printr-o anumită formă de comunicare verbală.

Această descriere este destul de largă, incluzând o multitudine de sisteme bazate mai mult sau mai puțin pe comunicare verbală.

O primă categorie o reprezintă sistemele de furnizare de informații. Aceste sisteme operează de regulă în domenii bine definite, acționând practic ca interfețe către diverse baze de date. Ele controlează destul de rigid dialogul, libertatea de formulare a utilizatorului fiind deci redusă. Această abordare duce însă la obținerea unor sisteme robuste, cu o rată relativ redusă a erorilor, capabile să satisfacă majoritatea cererilor.

Altă categorie foarte importantă este reprezentată de sistemele ce oferă diverse categorii de servicii: rezervări de bilete pentru avioane, trenuri și alte mijloace publice de transport, servicii telefonice, bancare etc. Numeroase eforturi s-au îndreptat în această direcție, întrucât astfel de sisteme au și o pronunțată valoare comercială.

Paleta este completată de diverse sisteme cu comandă vocală și de sistemele de traducere automată, deși în aceste cazuri interactivitatea este aproape nulă.

În tabelul ?? sunt prezentate ilustrativ câteva din realizările importante în domeniul sistemelor de dialog.

Pentru a ilustra și mai bine performanțele de care sunt capabile unele sisteme de dialog la ora actuală, am inclus în figura ?? un exemplu de dialog purtat cu CMU Communicator, un sistem pentru planificarea călătoriilor și rezervarea de bilete pe liniile aeriene dezvoltat la Carnegie Mellon University.

Sistem	Tipul sistemului și domeniul
WAXHOLM	Furnizare de informații în legătură cu traficul navelor în Stockholm (KTH)
TRAINS	Planificarea călătoriilor cu trenul (University of Rochester)
JANUS/ESST	Traducerea automată a dialogurilor spontane om-om pentru planificarea întâlnirilor (Karlsruhe & CMU)
CMU Communicator	Servicii de planificare și rezervare de bilete pentru călătorii cu avionul (CMU)
TARSAN	Furnizare de informații turistice (rute, hoteluri, muzee, adrese, numere telefonice etc) (Canon)
TOSBURG	Preluare de comenzi într-un fast-food (Toshiba)

Tabelul 1.1: Câteva exemple de sisteme de dialog

```
CMC: Buna ziua, sunt CMU Communicator. Spuneti-mi numele dumneavoastra.
Alex: Sunt Alex Rudnicky.
CMC: Buna ziua doctor Rudnicky, ma bucur ca ati apelat din nou. Unde ati
dori sa calatoriti ?
Alex: As vrea sa zbor la San Francisco maine seara.
CMC: Zburati la San Francisco maine. Un zbor US Airlines pleaca din
Pittsburgh la 6:10 PM, si ajunge in San Fran...
Alex: OK, Il iau pe asta.
CMC: Va si intoarcati la Pittsburgh ?
Alex: Da, as vrea sa ma intorc vineri dupa-amiaza.
CMC: Zburati vineri, 1 octombrie. Exista un zbor ce pleaca la 1:10 PM si
soseste in Pittsburgh la 8:53 PM. Vreti sa il luati pe acesta ?
Alex: OK, Il iau pe asta.
CMC: Aveti nevoie de hotel in San Francisco ?
Alex: Fara hotel.
CMC: Aveti nevoie de o masina...
Alex: Da, te rog.
CMC: Am facut o rezervare pentru o masina in San Francisco. Doriti un sumar
al calatoriei ?
Alex: Da.
CMC: Primul zbor este US Airlines 79 cu plecare din Pittsburgh la 6:10 PM
si sosire in San Francisco la 8:27 PM. Apoi urmatorul este cu plecare
din San Francisco la 1:10 PM vineri, 1 octombrie si sosire in
Pittsburgh la 8:53 PM. Nu am facut rezervare de hotel. Am rezervat o
masina in San Francisco. Doriti sa platiti cu cartea de credit ?
Alex: Da, te rog.
CMC: Va multumesc. Va voi trimite si un sumar prin e-mail. Asta s-ar putea
sa dureze cateva minute. Va multumesc ca ati utilizat Carnegie Mellon
Communicator.
```

Figura 1.1: Exemplu de dialog cu CMU Communicator

1.2 Structura tipică a unui sistem de dialog

Realizarea unor sisteme de dialog robuste și cu performanțe ridicate a devenit posibilă relativ recent, datorită creșterii capacității de procesare și a dezvoltărilor din domeniul precum modelarea lingvistică și recunoașterea și sinteza vorbirii. Această “tinerețe” a domeniului se traduce însă într-o lipsă de metodologii de dezvoltare consacrate, lucru care face proiectarea sistemelor de dialog o activitate anevoioasă, situată deseori la limita dintre știință și artă [?].

Cu toate acestea, indiferent de tipul și domeniul sistemului de dialog, problemele ce trebuie rezolvate sunt în principiu aceleași: recunoașterea vorbirii, analiza semantică, controlul dialogului, generarea răspunsurilor și sinteza vorbirii. În consecință, majoritatea sistemelor sunt proiectate modular, putând fi identificată o arhitectură generică, dictată de aceste probleme.

O structură tipică pentru un sistem de dialog interactiv pentru furnizarea de informații este ilustrată în figura ???. În continuare se va descrie pe scurt rolul fiecăruia din modulele identificate:

- **Modulul de recunoaștere a vorbirii** preia semnalul sonor de la utilizator

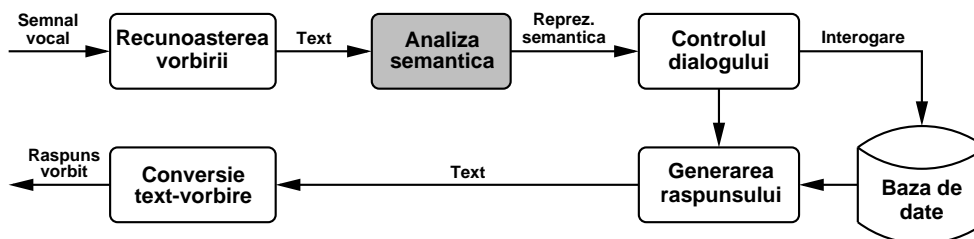


Figura 1.2: Structura unui sistem de dialog pentru obținerea de informații

(de obicei prin microfon sau linie telefonică) și generează secvența de cuvinte cea mai probabilă pornind de la acest semnal, pe baza unor modele acustice și lingvistice.

- **Modulul de analiză semantică** are ca intrare textul produs de modulul de recunoaștere și generează o reprezentare formalizată a înțelesului (sensului) acestui text. Analiza semantică are deci rolul de a extrage și clarifica conținutul de informație utilă al formulărilor utilizatorului.
- **Modulul pentru controlul dialogului** este practic nucleul sistemului de dialog, cu rolul de a realiza gestionarea dialogului și a interacțiunii cu utilizatorul. În plus, acestui modul îi revine și rolul de a menține istoricul dialogului și de a realiza analiza semantică contextuală.
- **Modulul de generare a răspunsurilor** este comandat de cel de control al dialogului în funcție de starea curentă a acestuia, și formulează răspunsul potrivit la un moment dat, folosind eventual și rezultatele unui acces la o bază de date pentru completarea informațiilor necesare.
- **Modulul de conversie text-vorbire** funcționează în sensul invers celui de recunoaștere. El preia ieșirea în formă textuală a generatorului de răspunsuri, și sintetizează forma sonoră a acesteia.

1.3 Rolul analizorului semantic

Analizorul semantic este o componentă esențială în orice sistem de dialog. Plasat între modulul de recunoaștere a vorbirii și cel de control al dialogului (vezi figura ??), analizorul semantic are rolul de a produce o reprezentare formală a sensului pronunțiilor utilizatorului.

Datorită ambiguității și varietății caracteristice limbajului natural, controlul dialogului nu se poate realiza direct pe baza textului obținut de modulul de recunoaștere a vorbirii. Analiza semantică are deci rolul de a identifica, extrage și clarifica informația utilă existentă în text, transcriind-o într-o formă riguros definită, adecvată procesării prin algoritmi specifici de control al dialogului.

Există mai multe abordări pentru realizarea analizoarelor semantice, dar majoritatea lor se încadrează într-una din două mari categorii: analizoare semantice *bazate pe reguli* și analizoare semantice *bazate pe principii stocastice*. Lucrarea de față se va axa pe a doua categorie menționată, prezentând teoretic avantajele acesteia și ilustrându-le practic prin construirea unui analizor semantic stocastic și evaluarea performanțelor lui într-un domeniu bine definit.

1.4 Organizarea lucrării

În continuare, pentru fundamentarea teoretică a lucrării, capitolul ?? trece în revistă unele aspecte ale analizei semantice în sistemele de dialog, împreună cu principalele alternative de proiectare existente (formalisme de reprezentare a cunoștințelor și metode de parsing), și face o scurtă prezentare a soluțiilor utilizate pentru analiza semantică în câteva sisteme de dialog descrise în literatura de specialitate.

Pornind de la concluziile astfel obținute, capitolul ?? prezintă realizarea unui analizor semantic stocastic care utilizează un formalism case-grammar pentru reprezentarea cunoștințelor, detaliind modulele componente și suportul teoretic pe care se bazează funcționarea lor.

Experimentele efectuate în vederea integrării acestui analizor într-un sistem de dialog pentru informații despre orarul Departamentului de Calculatoare, precum și rezultatele obținute, sunt descrise în capitolul 4.

În încheiere sunt prezentate concluzii degajate din dezvoltarea analizorului și experimentele efectuate, precum și posibile direcții de continuare a cercetărilor.

Capitolul 2

Analiza semantică în sistemele de dialog

Acest capitol este dedicat trecerii în revistă a unor aspecte teoretice ale analizei semantice, având ca finalitate alegerea între diferite opțiuni posibile în proiectarea unui analizor semantic utilizabil ca parte componentă a unor sisteme de dialog.

După definirea noțiunii de analiză semantică și stabilirea rolului analizorului semantic într-un sistem de dialog, se trece la o prezentare a metodelor existente de proiectare a analizoarelor semantice. Sunt prezentate diverse alternative, împreună cu un studiu teoretic al avantajelor și dezavantajelor fiecăreia.

Capitolul se încheie printr-o scurtă trecere în revistă a unor analizoare semantice din sisteme de dialog prezentate în literatura de specialitate, analizate prin prisma metodelor utilizate și a performanțelor lor.

2.1 Analiza semantică

2.1.1 Definirea noțiunii de analiză semantică

Dicționarul explicativ al limbii române furnizează următoarea definiție pentru noțiunea de *semantică*:

semantic, -ă, *semantici*, -ce s.f., adj. I. S.f. Ramură a lingvisticii care se ocupă cu studierea sensurilor cuvintelor și a evoluției acestor sensuri. [...]

Stabilirea sensului cuvintelor într-un anumit context nu este de loc o operație trivială, dificultatea fiind generată de câteva caracteristici esențiale ale limbajului natural. Pe de o parte, el permite exprimarea aceluiași înțeles, transmiterea aceluiași conținut informativ, prin formulări diferite. Pe de altă parte, aceeași formulare poate avea semnificații diferite în contexte diferite. În plus, mai trebuie remarcat faptul că într-o formulare nu toate cuvintele sunt purtătoare de informație utilă.

Aceste caracteristici de ambiguitate ale formei scrise sau vorbite a limbajului natural fac necesară operația de analiză semantică. Prin ea se urmărește practic extragerea informației utile într-o formă adecvată procesării ulterioare prin diverși algoritmi. Analiza semantică trebuie deci să îndeplinească următoarele sarcini:

- identificarea cuvintelor purtătoare de informație utilă;
- stabilirea sensurilor acestor cuvinte, și clarificarea lor în cazuri ambigue.

2.1.2 Nivele de analiză semantică

Sunt necesare câteva clarificări suplimentare și în ceea ce privește **nivelul** la care se încearcă identificarea sensului cuvintelor. Încă de la început trebuie precizat că stabilirea sensurilor efective ale cuvintelor în cadrul unor expresii stilistice precum metaforele, metonimiile, exprimările ironice, și altele de acest gen, depășesc complet teoriile și puterea computațională curente, și nu fac deci obiectul analizei semantice așa cum este ea înțeleasă în lucrarea de față. Acest gen de exprimări sunt suficient de rare, în special în dialoguri purtate cu sisteme de calcul, astfel încât neluarea lor în considerare nu introduce practic erori semnificative.

Lasând la o parte aceste exprimări, se pot identifica două nivele pe care putem determina sensul cuvintelor. Să ne imaginăm următoarea formulare ipotetică a unui utilizator într-o conversație cu un sistem de dialog:

– **Dar grupa întâi, când are ore tot cu el ?**

Cuvântul **ore**, dacă este considerat izolat, devine ambiguu: nu știm dacă se referă la noțiunea de timp, sau la cea de oră de curs, laborator etc. La nivelul propoziției însă, ambiguitatea dispare, evident sensul fiind cel de oră de curs.

Pe de altă parte, sensul cuvântului **el** nu poate fi determinat doar din propoziția considerată, pentru aceasta fiind necesare mai multe informații. Dacă am cunoaște și propoziția anterioară, care ar putea fi de exemplu:

– **Când are grupa doi curs cu profesorul Vancea?**

atunci ar deveni evident că **el** face de fapt referire la profesorul Vancea.

În consecință, putem defini practic două nivele de analiză semantică: la un prim nivel se poate extrage un *înțeles independent de context al propoziției* [?], caz în care cuvântul **el** este încă ambiguu; la un nivel superior, se poate obține un *sens complet al propoziției*, dependent de contextul în care a fost formulată.

Întrucât modulul de analiză semantică al unui sistem de dialog operează pe propoziții independente (vezi figura ??), rolul lui se limitează la primul nivel: obținerea sensului propoziției independent de context și exprimarea lui într-o formă riguroasă, bine definită. Analiza în context, cea care determină sensul complet al fiecărei propoziții, se realizează la nivelul modulului de control al dialogului. Acest modul menține un *istoric al dialogului*, și poate determina sensul complet pe baza acestui istoric, făcând uz de algoritmi specifici.

În concluzie, într-un sistem de dialog rolul analizorului semantic este de a extrage și clarifica sensul independent de context al formulărilor utilizatorului.

2.1.3 Legături între analiza semantică și cea sintactică

Între analiza semantică și cea sintactică poate exista o strânsă legătură: dacă semantica se ocupă cu stabilirea sensului cuvintelor și a informației transmise de ele, *sintaxa* definește regulile de îmbinare a cuvintelor în propoziții și a propozițiilor în fraze [?]. Ambele folosesc deci cuvântul ca unitate fundamentală de procesare. În plus, majoritatea limbajelor (atât naturale, cât și formale) sunt *compoziționale*, adică sensul unei propoziții derivă din sensurile părților ei componente, ceea ce reprezintă un motiv suplimentar pentru realizarea unei legături între ele [?].

Aceasta a făcut ca majoritatea tehnicilor de analiză semantică să își aibe originea în cele de analiză sintactică. Principalele instrumente utilizate sunt tot gramaticile, însă cu diverse modificări și extinderi, necesare pentru a capta aspectele caracteristice analizei semantice.

Atât analiza sintactică cât și cea semantică sunt foarte utilizate în domeniul limbajelor formale și al translatoarelor, existând o multitudine de formalisme și algoritmi specializați în acest sens. În acest domeniu, analiza sintactică precede analiza semantică, și generează informații care sunt apoi utilizate pentru analiza semantică. Deci analiza semantică în domeniul limbajelor formale se bazează și este în general dirijată de cea sintactică.

Extrapolarea acestor tehnici la nivelul limbajului natural se dovedește însă dificilă datorită varietății și ambiguităților caracteristice lui. În primul rând, construirea unei gramatici acoperitoare pentru sintaxa unei limbi naturale presupune un anumit grad de expertiză lingvistică. Apoi, chiar dacă o astfel de gramatică ar exista, în limbajul natural pot apare formulări incorecte din punctul ei de vedere, dar care să fie totuși inteligibile, ba chiar perfect acceptabile în uzul curent, datorită caracterului evolutiv al limbii, incompletitudinii gramaticii, dar mai ales datorită diverselor fenomene ce caracterizează vorbirea spontană: repetiții, starturi false, disfluențe etc. Într-un sistem de dialog, la acestea se pot adăuga eventuale erori de recunoaștere a vorbirii, care pot deasemeni compromite corectitudinea sintactică a propoziției, lăsând-o însă neschimbată din punct de vedere semantic.

În concluzie, deși majoritatea tehnicilor de analiză semantică actuale sunt extinderi ale unora de analiză sintactică, realizarea prealabilă a unei analize sintactice a limbajului natural nu este nici necesară, nici recomandabilă, putând fi chiar dăunătoare în construirea analizoarelor semantice pentru sisteme de dialog.

2.2 Alternative în proiectarea unui analizor semantic

În cele ce urmează vor fi prezentate principalele alternative de proiectare ce apar în dezvoltarea unui analizor semantic pentru un sistem de dialog.

Primul pas îl constituie alegerea unui formalism adecvat pentru reprezentarea cunoștințelor, a informației extrase. Aici există mai multe alternative, care vor fi

prezentate împreună cu o analiză a avantajelor și dezavantajelor fiecăreia.

După stabilirea formalismului de reprezentare a cunoștințelor, a doua problemă ce trebuie atacată este alegerea metodei prin care se va realiza propriu-zis extragerea informației în forma stabilită, deci a algoritmului efectiv de analiză semantică.

2.2.1 Formalisme de reprezentare a cunoștințelor

De-a lungul timpului au fost definite o multitudine de formalisme gramaticale, și diverse tehnici de parsing pentru aceste formalisme. Ele au apărut ca urmare a necesității de a realiza o analiză sintactică performantă pentru întreaga gamă de limbaje: de la cele puternic formalizate, până la limbajul natural. Datorită legăturilor care se pot stabili între analiza sintactică și cea semantică, precum și a necesității unui formalism de reprezentare a informației extrase din text, a apărut natural ideea de a extinde gramaticile utilizate în analiza sintactică pentru a capta aspectele specifice analizei semantice. Astfel, pe baza gramaticilor sintactice au apărut noi *formalisme semantice* adaptate pentru analiza semantică.

În continuare se va face o scurtă trecere în revistă a celor mai importante tipuri de gramatici sintactice, urmată de prezentarea principalelor formalisme utilizate în analiza semantică.

Gramatici sintactice

- *Gramatici Chomsky*

Noțiunea de gramatică (în sensul utilizat în teoria limbajelor formale și a translatoarelor) a fost introdusă de Noam Chomsky în anii '50. O gramatică este constituită practic dintr-un set de reguli (producții) care realizează o descriere formală a structurilor sintactice posibile în limbajul pe care îl descrie. Chomsky a realizat și o ierarhizare a acestor gramatici în funcție de tipul producțiilor lor și, implicit, de capacitatea lor de exprimare. Clasa cea mai largă este cea a gramaticilor de tipul 0, care nu au restricții în ceea ce privește scrierea producțiilor. Urmează clasa gramaticilor de tip 1, sau dependente de context. Și mai restrictive sunt gramaticile de tipul 2, sau independente de context, iar clasa cea mai restrictivă (și în același timp cea mai puțin expresivă) este cea a gramaticilor de tipul 3 sau regulate.

Gramaticile independente de context (tipul 2) au fost utilizate cu mult succes în descrierea și analiza sintactică a limbajelor formale. Ele sunt suficient de expresive pentru această sarcină și, în plus, există numeroși algoritmi eficienți (ca timp și spațiu de memorie) pentru realizarea parsingului peste aceste gramatici.

Deși ele sunt suficient de puternice și pentru descrierea majorității construcțiilor limbajului natural, realizarea unui analizor sintactic bazat pe gramatici independente de context în acest domeniu este o sarcină dificilă. O astfel de gramatică ar conține un număr mare de reguli, iar scrierea

acestora ar necesita un anumit grad de expertiză lingvistică. Dacă în plus dorim captarea robustă a tuturor aspectelor caracteristice limbajului vorbit (disfluente, repetiții, starturi false etc.), așa cum este necesar într-un sistem de dialog, numărul regulilor ar crește foarte mult, făcând impractic efortul necesar dezvoltării unei astfel de gramatici [?].

- ***Rețele de tranziții (Transition Networks)***

Rețelele de tranziții constituie un alt formalism sintactic, dar care nu folosește reguli, ci grafuri orientate pentru a descrie structurile sintactice acceptate într-un anumit limbaj. O variantă mai puternică a acestora, rețelele de tranziții recursive (Recursive Transition Networks - RTN) sunt echivalente ca și capacitate de exprimare cu gramaticile independente de context.

- ***Gramatici extinse (Augmented grammars)***

Gramaticile extinse sunt obținute plecând de la cele independente de context, prin atașarea unor atribute constituenților producțiilor lor. Astfel, se pot modela mai ușor anumite aspecte ale limbajului natural precum acordurile între cuvinte, subcategoriile etc. [?]. Din acest punct de vedere, gramaticile extinse beneficiază de o anumită flexibilitate care lipsește gramaticilor Chomsky, lucru care le face un candidat mai bun pentru analiza sintactică a limbajului natural. Trebuie remarcat că aceeași extensie – atașarea de atribute nodurilor – poate fi aplicată și rețelelor de tranziții, rezultând un nou formalism: rețelele de tranziții extinse.

- ***Gramatici arborescente (Tree adjoining grammars)***

Gramaticile arborescente reprezintă un alt formalism adaptat la particularitățile analizei sintactice a limbajului natural. Ele descriu structurile admisibile în limbaj prin intermediul unui set de arbori elementari (corespunzători celor mai simple propoziții ale limbajului), și a unei operații de “alăturare” a arborilor pe baza căreia se pot sintetiza structuri complexe [?]. Capacitatea lor de exprimare o depășește cu puțin pe cea a gramaticilor independente de context, fiind o soluție mai potrivită pentru analiza sintactică a limbajului natural.

Toate aceste gramatici sunt potrivite pentru analiza sintactică, care are însă un scop diferit de cel al analizei semantice. Pornind de la ele, au fost create noi formalisme, mai potrivite pentru sarcinile și necesitățile analizei semantice.

Formalismele pentru analiza semantică

- ***Gramatici semantice***

Noțiunea de gramatică semantică a fost introdusă de Burton în 1976, extinzând gramaticile Chomsky în sensul integrării informației semantice și sintactice în același model [?]. Dezavantajele utilizării acestor gramatici pentru analiza semantică a limbajului natural sunt, pe de o parte, lipsa

de portabilitate între domenii, pe de alta, rigiditatea lor (moștenită de la gramaticile Chomsky) față de fenomenele specifice vorbirii spontane.

- **Formalismul case-grammar (Case-grammar Formalism)**

Acest formalism este poate cel mai îndepărtat față de noțiunea tipică de gramatică (așa cum a fost introdusă de Chomsky), având totuși o mulțime de caracteristici care îl fac optim pentru utilizarea în cadrul analizei semantice a limbajului natural.

Formalismul a fost introdus de Fillmore în 1968, și extins de Bruce în 1975. La baza lui se află conceptul de *cadru* (frame), introdus de Minsky în 1975 ca o modalitate de reprezentare a cunoștințelor. Ideea este foarte simplă și în același timp puternică, găsindu-și o largă aplicabilitate în numeroase subdomenii ale inteligenței artificiale. Componentele esențiale ale unui cadru sunt un așa-numit *concept*, fixat, și o mulțime de *sloturi*, completate cu informații ce reprezintă “cunoștințele” sistemului. Activitatea de completare a sloturilor unui cadru o vom numi *instanțierea cadrului*. Pentru exemplificare, revenim la formularea ipotetică a unui utilizator în discuția sa cu un sistem de dialog furnizând informații despre orar:

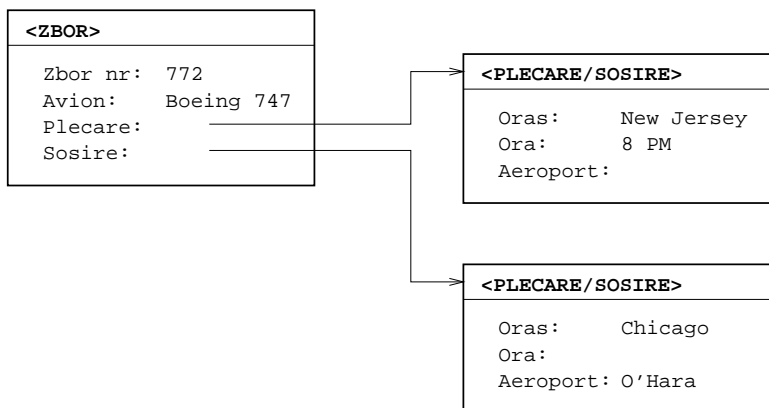
– Când are grupa doi curs cu profesorul Vancea?

Figura ?? ilustrează o instanță a unui cadru reprezentând conținutul semantic al acestei formulări. Conceptul <când>, specificat în partea superioară a cadrului, ne spune că este vorba despre o întrebare în legătură cu timpul. Restul conținutului semantic este încapsulat în valorile atribuite sloturilor **grupa**, **tip-materie** și **profesor**. Nu este obligatoriu ca toate sloturile să fie completate – informațiile sunt extrase în măsura în care ele există în text.

<când>	
an:	-
grupa:	2
subgrupa:	-
tip-materie:	curs
profesor:	Vancea

Figura 2.1: Exemplu de cadru

O a doua noțiune importantă este aceea de *marcaj* (case-marker). Un marcaj este practic un cuvânt din text care realizează o anumită constrângere în ceea ce privește completarea sloturilor. Acest gen de constrângeri lexicalizate în instanțierea cadrelor modelează practic sintaxa limbajului, limitând construcțiile interpretabile prin formalismul case-grammar. În exemplul anterior, cuvântul **grupa** este un marcaj indicând faptul că urmează probabil valoarea slotului (*case-value*) **grupa**; similar, cuvântul **profesorul** este un



Avionul Boeing 747 al cursei US Airlines 772 pleaca din New Jersey la ora 8 pm si soseste in Chicago pe aeroportul O'Hara

Figura 2.2: Un exemplu de sistem de cadre

marcaj pentru slotul profesor. Cuvântul curs este ceva mai special, el fiind în același timp marcaj și valoare pentru slotul tip-materie. Algoritmul de extragere a informației (parsing utilizând formalismul case-grammar) se bazează în mare măsură pe marcaje, acestea indicând unde se află informația și în același timp clarificând sensul ei.

Pornind de la cadre individuale, se poate realiza un sistem de cadre (case-system) care să exprime prin legături interdependențele semantice dintre cadrele ce îl compun. Legăturile se materializează prin faptul că un slot al unui cadru, în loc să fie completat cu o valoare explicită, este completat cu o legătură spre un alt cadru din sistem. Se formează astfel o rețea de cadre a cărei capacitate de reprezentare a cunoștințelor este mult mărită față de cea a unui cadru individual. Un exemplu este dat în figura ??.

Prin formalism case-grammar se înțelege deci un astfel de sistem de cadre, cu aspecte sintactice incluse prin constrângerile lexicalizate, dictate de marcaje, în instanțierea sloturilor [?].

Dintre toate variantele prezentate, formalismul case-grammar reprezintă alternativa cea mai potrivită pentru realizarea analizei semantice într-un sistem de dialog, principalul lui avantaj constând în flexibilitatea mult superioară gramaticilor sintactice și diverselor variante extinse ale acestora. Spre deosebire de gramaticile clasice, care modelează limbajul și eventual conținutul informativ al acestuia prin intermediul unor reguli bine definite, rigide, formalismul case-grammar indică doar structurile semantice acceptabile, prin intermediul cadrelor și a relațiilor dintre acestea. Aspectele sintactice sunt incluse, și dirijează procesul de extragere a informației, însă prin intermediul constrângerilor dictate de marcaje – un meca-

nism mult mai relaxat decât cel al gramaticilor clasice. Aceasta face ca formalismul case-grammar să poată fi utilizat cu succes și în cazurile în care gramaticile tipice ar eşua – cele ale formulărilor incorecte sintactic, dar în același timp valide semantic, formulări destul de des întâlnite în conversațiile dintre oameni, și care deci pot fi destul de frecvente și în conversațiile purtate cu sisteme de dialog.

Reprezentarea bazată pe cadre este foarte potrivită și din punctul de vedere al modulului de control al dialogului, deoarece simplitatea reprezentării face ca menținerea unui istoric semantic al dialogului și realizarea analizei de context să se poată face relativ ușor.

În concluzie, alegerea formalismului case-grammar este cea mai potrivită pentru reprezentarea cunoștințelor în analizorul semantic proiectat.

2.2.2 Metode de parsing

Odată definită modalitatea de reprezentare a cunoștințelor, următorul pas în proiectarea unui analizor semantic este stabilirea metodei pe care o vom utiliza pentru extragerea efectivă a informației din text. Pentru majoritatea gramaticilor clasice există numeroși algoritmi de parsing care pot genera reprezentările corespunzătoare. În ceea ce privește parsingul bazat pe un formalism case-grammar, soluțiile se grupează în două categorii: *metode bazate pe reguli* și *metode stocastice*.

Metodele bazate pe reguli, așa cum se spune și numele, utilizează un set de reguli care controlează modul în care se face identificarea conceptelor și completarea sloturilor cu valorile corespunzătoare din text. Figura ?? ilustrează structura unui astfel de analizor.

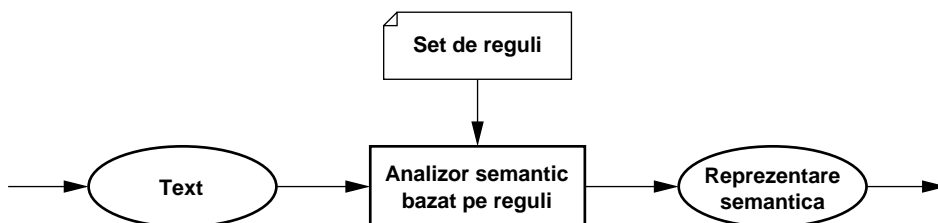


Figura 2.3: Structura unui analizor semantic bazat pe reguli

În general, regulile sunt lexicalizate: se definesc familiile de cuvinte care identifică conceptele/cadrelle, și familiile de cuvinte care joacă rolul de marcaje. În plus, trebuie definite reguli care descriu legăturile dintre marcaje și valori. Pe baza acestora se pot determina cuvintele care au conținut semantic și dau valori sloturilor cadrului identificat. Evident, corectitudinea – în sensul cel mai larg – regulilor de parsing va dicta performanța analizorului, lucru care face ca partea cea mai mare din proiectarea unui astfel de analizor să stea în definirea regulilor.

Alternativa o reprezintă metodele de parsing stocastic. Ideea fundamentală în aceste metode este de a utiliza un model probabilistic pentru decodarea semantică a formulărilor utilizatorului (figura ??).

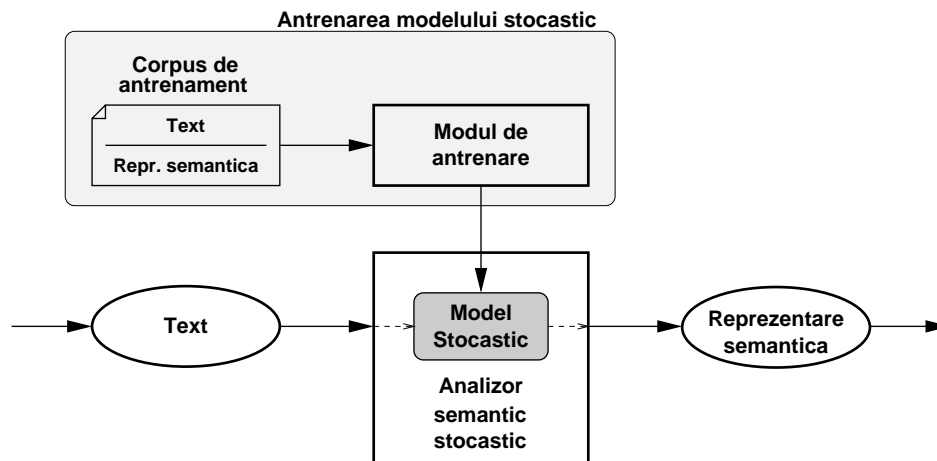


Figura 2.4: Structura unui analizor semantic stocastic

Modelul este creat într-o primă etapă, de antrenare, folosind algoritmi caracteristici tipului de model utilizat (modele Markov, rețele neuronale etc.) Resursa esențială pentru antrenare este un corpus reprezentativ de formulări utilizator care au fost în prealabil analizate semantic, astfel încât fiecare text are asociată reprezentarea semantică. Prin antrenare, modelul va “învăța”, din perechile (text, reprezentare semantică), cum se generează informația semantică din text. Se elimină astfel necesitatea de a defini explicit regulile pe baza cărora se face extracția, ele fiind practic “captate” automat în faza de antrenare, sub formă de valori ai diversilor parametri ai modelului [?].

Odată etapa de antrenare încheiată, analiza semantică decurge simplu: analizorul primește la intrare formulările utilizatorului sub formă de text, iar modelul va furniza reprezentările semantice cele mai probabile pentru intrările respective.

Acest model stocastic este practic un mecanism de învățare automată care operează în două moduri: în *modul de antrenare*, el primește la intrare perechi (text, reprezentare semantică) și își ajustează parametrii interni pe baza unor algoritmi specifici astfel încât, statistic, să capteze corespondența dintre text și reprezentarea semantică. În al doilea mod, de *decodare*, modelul primește la intrare doar textul, și furnizează reprezentarea lui semantică cea mai probabilă.

Comparație între parsingul stocastic și cel bazat pe reguli

Metoda de parsing stocastic prezintă numeroase avantaje față de metodele bazate pe reguli, lucru demonstrat în continuare printr-o scurtă analiză comparativă.

După cum s-a subliniat și anterior, la realizarea unui analizor semantic bazat pe reguli partea cea mai dificilă este proiectarea regulilor. Dezvoltarea unui astfel de set de reguli este costisitoare, și nu întotdeauna la îndemâna proiectanților unui sistem de dialog, întrucât presupune participarea unor experți lingviști.

Bazarea parsingului pe un set de reguli are și dezavantajul că reintroduce o anumită rigiditate în analiza semantică, problemă pe care am încercat să o atenuăm prin utilizarea formalismului case-grammar, iar captarea prin reguli a tuturor fenomenelor caracteristice vorbirii spontane este foarte dificilă, măbind foarte mult numărul de reguli (implicit costul dezvoltării sistemului) și reducând performanțele.

În plus, metodele bazate pe reguli suferă de o acută lipsă de portabilitate. Fie că este vorba de portarea analizorului semantic de la un domeniu de aplicație la altul, fie că este vorba de lărgirea domeniului, de cele mai multe ori este necesară cel puțin o extindere a setului de reguli (pentru o analiză detaliată vezi [?]).

În ceea ce privește analiza semantică stocastică, singura resursă necesară este un corpus de antrenament analizat semantic (cel mai probabil manual), care să includă diversele tipuri de formulări ce apar în dialogurile în care sistemul va fi implicat. Obținerea unui astfel de corpus se poate face inițial prin metoda Vrăjitorului din Oz [?], [?], sau pe parcursul funcționării unui sistem existent.

De remarcat că problema analizării semantice manuale a acestui corpus nu este foarte mare, întrucât dezvoltarea modelului se poate face *incremental (bootstrapping)*: se analizează semantic manual o porțiune a corpusului, se antrenează o primă variantă a modelului cu această porțiune, modelul este utilizat pentru a decoda semantic o nouă porțiune din corpus, decodare care se corectează manual, după care pe baza ei și a primei porțiuni se face o nouă antrenare șamd.

În plus, metodele stocastice elimină practic toate celelalte dezavantaje ale sistemelor bazate pe reguli. În cazul lor, nu trebuie dezvoltat un set de reguli, acestea fiind practic învățate automat în faza de antrenare a modelului. Dispare deci în mare măsură necesitatea expertizei lingvistice. Analiza semantică stocastică are un grad sporit de flexibilitate, modelul fiind “reglat” prin antrenare pe tipurile de formulări care apar în dialoguri. În fine, în cazul portării analizorului la alt domeniu (sau al extinderii domeniului), singurul lucru necesar este o reantrenare a modelului cu un corpus de formulări adecvat.

În concluzie, pentru sistemele automate de dialog vocal om-calculator o abordare stocastică a parsingului bazat pe formalismul case-grammar este mult mai potrivită decât cele bazate pe reguli, iar lucrarea de față se va axa în continuare pe dezvoltarea unui astfel de analizor semantic.

2.3 Stadiul actual

Vom încheia prezentarea teoretică a alternativelor de proiectare a unui analizor semantic pentru sisteme de dialog printr-o scurtă trecere în revistă a soluțiilor utilizate de câteva sisteme existente. După cum vom vedea, o analiză comparativă a metodelor folosite și a performanțelor lor reprezintă încă o justificare a alegerii soluției unui analizor semantic stocastic bazat pe formalismul case-grammar.

Sisteme cu analizoare semantice bazate pe reguli

- **CMU-PHOENIX**

Acest sistem de analiză semantică, dezvoltat la Carnegie Mellon University [?], folosește un formalism case-grammar pentru reprezentarea cunoștințelor, iar parsingul se face la nivel de expresii (nu propoziții), pe bază de reguli implementate sub forma unor rețele de tranziții recursive (vezi ??). PHOENIX este utilizat în diverse sisteme de dialog, obținând în domeniul ATIS (Air Travel Information Services) o rată a erorilor de 3.8% (cea mai bună performanță a unui analizor semantic în acest domeniu la evaluările DARPA din 1994).

- **JANUS**

JANUS [?] este un sistem de traducere automată dezvoltat în colaborare de ATR Interpreting Telephony Research Laboratories (Japonia), CMU, Universitatea din Karlsruhe și Siemens. Sistemul efectuează traduceri automate în domeniul programării de întâlniri, și utilizează în paralel două metode de analiză semantică: sistemul CMU-PHOENIX, și un parser LR generalizat. Rata erorilor sistemului este de cca 24.5% pentru traduceri din germană în engleză, și 18.6% pentru traduceri din spaniolă în engleză.

- **LIMSI-L'ATIS**

LIMSI-L'ATIS este versiunea franceză de sistem de dialog în domeniul ATIS, dezvoltat la LIMSI-CNRS [?], utilizând și el reprezentarea prin cadre și parsingul pe bază de reguli. Pentru descrierea sistemului de cadre și a setului de reguli de parsing a fost creat un limbaj declarativ dedicat. Întrucât s-a observat că regulile eșuează în anumite situații caracteristice vorbirii spontane, s-a introdus o fază de preprocesare a intrării analizorului semantic, în care se realizează o normalizare a textului. Rata de eroare măsurată a acestui sistem a fost de cca. 8.7%

- **LIMSI-MASK**

LIMSI-MASK este un sistem similar cu LIMSI-L'ATIS, dezvoltat tot la LIMSI-CNRS [?], dar care operează în domeniul rezervării de bilete și obținerii de informații referitoare la mersul trenurilor. Abordarea analizorului semantic este identică cu cea din LIMSI-L'ATIS.

Sisteme cu analizoare semantice stochastice

- **AT&T CHRONUS**

AT&T CHRONUS [?] este un analizor semantic bazat pe parsing stocastic, dezvoltat la AT&T și testat pe domeniul ATIS. Analizorul utilizează un model Markov, analiza semantică reducându-se la decodarea stărilor acestuia pentru o secvență de simboluri observate care corespunde intrării. Sistemul are performanțe foarte bune, fiind comparabil cu CMU-PHOENIX pe același domeniu – rata erorii de 3.8%

- **BBN-HUM (Hidden Understanding Model)**

Acest sistem, dezvoltat la firma BBN, utilizează tehnici de învățare automată statistică atât pentru analiza semantică la nivelul propoziției, cât și pentru analiza semantică în context [?]. În teste neoficiale pe domeniul ATIS a fost raportată o rată a erorilor de cca. 9.5%.

- **PHILIPS Train-Timetable Inquiry System**

Așa cum îi spune și numele, acesta este un sistem de dialog dezvoltat de Philips Research Laboratories [?], care operează în domeniul furnizării de informații despre mersul trenurilor. Sistemul utilizează tehnici stocastice de parsing peste o gramatică semantică (vezi ??). În această abordare, deși producțiile gramaticii sunt elaborate manual de un grup de experți, probabilitățile acestor producții sunt “învățate” automat dintr-un corpus de antrenament. Performanțele acestui sistem nu se ridică însă la nivelul celorlalte: rata raportată a erorilor este de cca. 25%.

Capitolul 3

Un analizor semantic stocastic

În capitolul anterior au fost prezentate principalele alternative de proiectare, și s-a decis utilizarea formalismului case-grammar și a parsingului stocastic pentru realizarea unui analizor semantic pentru sisteme de dialog.

În acest capitol, pornind de la schema bloc a analizorului semantic, vor fi detaliate componentele lui, prezentându-se atât suportul teoretic cât și soluțiile de implementare alese pentru fiecare modul.

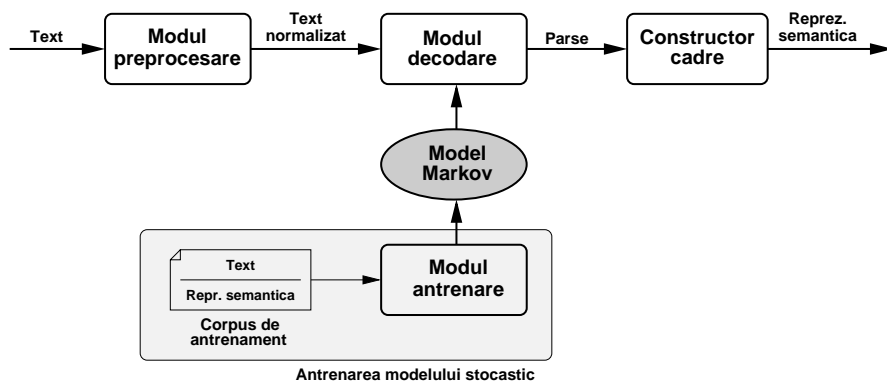
Deși analizorul este independent de domeniul sistemului de dialog în care va fi utilizat, pe parcursul acestui capitol vor fi utilizate ca exemple formulări specifice unui sistem dedicat furnizării de informații despre orarul Departamentului de Calculatoare, anticipând astfel experimentele și rezultatele din capitolul 4.

3.1 Structura analizorului semantic

Figura ?? prezintă arhitectura analizorului semantic împreună cu diversele reprezentări prin care trece textul obținut din pronunția utilizatorului pe parcursul extragerii informației și generării reprezentării semantice. Elementul central al analizorului proiectat îl constituie un model Markov (vezi ??) pe baza căruia se realizează analiza semantică. Pe lângă modulele de decodare și antrenare a modelului Markov, în componența sistemului mai intră două module: un modul pentru preprocesarea intrării și un modul pentru construirea efectivă a cadrelor. În continuare va fi prezentat pe scurt rolul fiecărei componente în cadrul arhitecturii ilustrate.

- **Modulul de preprocesare**

Modulul de preprocesare este situat la granița dintre modulul de recunoaștere automată a vorbirii și decodorul semantic propriu-zis. Rolul principal al acestuia este de a realiza o normalizare a textului primit de la modulul de recunoaștere. Preprocesarea constă în eliminarea evenimentelor non-lexicale, normalizarea numerelor, reducerea flexiunilor, unificarea expresiilor, înlocuiri de aliasuri, unificarea categoriilor și eliminarea cuvintelor din afara domeniului sistemului de dialog.



```

Text:          (aa) ce profesor tine laboratorul de s e la grupa (aa) doi
Text normalizat: ce-profesor laborator [MATERIA:"Sisteme expert"] grupa [NR:"2"]
Parse:         <ce-profesor> (m:laborator) (v:laborator) (m:grupa) (v:grupa)
Repr. semantica: <ce-profesor>
                  <identificare>
                  grupa = 2
                  <specif-materie>
                  laborator = "Sisteme expert"
  
```

Figura 3.1: Arhitectura analizorului semantic

- **Modelul Markov**

Modelul stocastic utilizat pentru realizarea parsingului în analizorul semantic proiectat este un *model Markov cu stări ascunse*, care acționează practic ca un mecanism de învățare (vezi ??). Crearea modelului se realizează printr-o etapă premergătoare de antrenare în care acesta “captează” regulile de decodare semantică sub forma unor parametri interni. Odată creat, modelul poate fi folosit într-un decodor care va genera cea mai probabilă secvență de etichete semantice (parse) corespunzătoare secvenței de intrare.

- **Modulul de decodare**

Pornind de la textul normalizat prin preprocesare, considerat ca o succesiune de simboluri observate ale modelului Markov, modulul de decodare determină cea mai probabilă succesiune de etichete semantice corespunzătoare prin aplicarea algoritmului Viterbi (vezi ??).

- **Modulul de antrenare**

Acest modul este utilizat doar în etapa premergătoare de creare/antrenare a modelului stocastic, el neavând un rol direct în analiza semantică efectivă. Modulul realizează antrenarea modelului printr-o serie de algoritmi specifici (vezi ??), pe baza unui corpus de formulări utilizator care pot fi considerate tipice pentru domeniul de aplicație al sistemului de dialog, și care în prealabil au fost etichetate semantic.

- **Constructorul de cadre**

Constructorul de cadre este ultimul modul al analizorului semantic. El primește de la modulul de decodare succesiunea de etichete semantice (parse), pornind de la care construiește reprezentarea pe bază de cadre a interpretării semantice a formulării utilizatorului.

3.2 Modulul de preprocesare

Întrucât forma în care modulul de recunoaștere a vorbirii furnizează textul nu este întotdeauna cea mai potrivită pentru realizarea analizei semantice, apare necesitatea introducerii unui modul de preprocesare între decodorul semantic și modulul de recunoaștere a vorbirii.

Rolul principal al preprocesorului este de a realiza o normalizare a textului de intrare. În plus, el realizează și o reducere a vocabularului de intrare al decodorului semantic propriu-zis, fără a afecta însă conținutul semantic al formulărilor. Această reducere a vocabularului este necesară întrucât dimensiunile modelului stocastic sunt direct dependente de dimensiunile vocabularului de intrare, iar prin reducerea dimensiunii sale se măresc șansele de a obține un model mai performant pornind de la un același corpus de antrenament.

Operațiile care compun etapa de preprocesare sunt sintetizate în tabelul ??, împreună cu exemplificări ale efectelor pe care le au asupra textului de intrare. În continuare vom detalia fiecare dintre acești pași:

1. **Eliminarea evenimentelor acustice non-lexicale**

Unele sisteme de recunoaștere automată a vorbirii pot transcrie și evenimente acustice non-lexicale care apar în semnalul sonor. Fie că este vorba de evenimente acustice produse de utilizator (tuse, ezitări etc.), fie că este vorba de diverse alte zgomote din mediu, aceste evenimente nu transportă nici un fel de informație utilă, deci pot fi eliminate fără riscuri. În această lucrare am presupus că ele sunt furnizate în paranteze de modulul de recunoaștere a vorbirii; în exemplul din tabelul ??, au fost eliminate în această fază două ezitări (aa).

2. **Normalizarea numerelor**

În acest al doilea pas al preprocesării se realizează transformarea numerelor din forma scrisă furnizată de modulul de recunoaștere în reprezentarea cu cifre arabe corespunzătoare. În exemplul dat, “doi” se transformă în 2.

3. **Reducerea flexiunilor**

Spre deosebire de alte limbi (e.g. engleza), limba română este una puternic flexionară, astfel încât poate fi utilă introducerea unui mecanism de reducere a flexiunilor. Astfel se realizează o normalizare a intrării și o reducere a vocabularului domeniului, fără însă a afecta conținutul semantic al formulării. În exemplificarea dată, “laboratorul” se transformă în forma nearticulată **laborator**, iar “grupa” în **grupă**.

Operație	Exemplificare rezultat
Recunoaștere vorbire	(aa) ce profesor ține laboratorul de s e la grupa (aa) doi
Eliminare even. non-lexicale	ce profesor ține laboratorul de s e la grupa doi
Normalizare numere	ce profesor ține laboratorul de s e la grupa 2
Reducere flexiuni	ce profesor ține laborator de s e la grupă 2
Unificare expresii	ce-profesor ține laborator de s e la grupă 2
Înlocuiri aliasuri	ce-profesor ține laborator de “Sisteme expert” la grupă 2
Unificare categorii	ce-profesor ține laborator de [MATERIA: “Sisteme expert”] la grupă [NR: “2”]
Elim. cuv. din afara domeniului	ce-profesor laborator [MATERIA: “Sisteme expert”] grupă [NR: “2”]

Tabelul 3.1: Pașii de preprocesare

4. Unificarea expresiilor

În această etapă se realizează unificarea sub formă de expresii unitare a unor grupuri de cuvinte cu un corespondent semantic unic (e.g. formularea “ce profesor” se transformă în **ce-profesor**.)

5. Înlocuiri de aliasuri

Această operație este dependentă de domeniul sistemului de dialog, și prin ea se realizează înlocuirea aliasurilor (prescurtări, sinonime etc.) unor cuvinte sau expresii (e.g. “s e” devine **“Sisteme expert”**.)

6. Unificarea categoriilor

Operația de unificare a categoriilor are ca scop gruparea cuvintelor în categorii semantice pentru a facilita analiza ulterioară și a reduce dimensiunile vocabularului decodului. În continuare decodul va lucra cu categoriile respective ca simboluri de intrare, revenindu-se la valorile concrete doar în faza de construire a reprezentării pe bază de cadre.

Putem vorbi despre existența a două tipuri de categorii. Pe de o parte, se pot defini categorii generice în care să se unifice cuvintele de același tip sau diverse sinonime. De exemplu, toate numerele care apar sunt grupate în categoria [NR] (în exemplul nostru, “2” devine [NR:**“2”**]).

Pe de altă parte, se pot defini și categorii specifice domeniului sistemului de dialog proiectat. În acest sens, se poate defini câte o categorie pentru fiecare

câmp din baza de date interogată de sistem ale cărui valori pot să apară în formulările utilizatorului. În exemplul ilustrat, “Sisteme expert”, care este o valoare a câmpului MATERIA din baza de date, este transformată prin unificarea categoriilor în [MATERIA:“Sisteme expert”]. Se observă însă că o categorie similară nu a putut fi definită și pentru câmpul GRUPA întrucât valorile acestui câmp (1, 2, 3 ...) sunt ambigue – când apare un număr în text, preprocesorul nu poate decide dacă el desemnează o grupă, o subgrupă, sau o oră. Din această cauză valorile numerice au fost unificate în categoria generică [NR], urmând ca decodorul semantic să stabilească efectiv sensul acestor numere.

În concluzie, modul în care se face definirea categoriilor este în mare măsură dictat de particularitățile domeniului sistemului de dialog.

7. Eliminarea cuvintelor din afara domeniului

Această ultimă operație implică eliminarea tuturor cuvintelor care nu sunt relevante în contextul domeniului stabilit al sistemului de dialog. În exemplul ilustrat, au fost eliminate verbul “ține” și prepozițiile “de” și “la”. Aceste eliminări se realizează pe baza unui dicționar al domeniului, dicționar care trebuie să conțină toate cuvintele relevante din punct de vedere semantic pentru domeniul respectiv. Construirea acestui dicționar se realizează în momentul analizei corpusului de antrenament al modelului stocastic (vezi capitolul 4.)

Trebuie remarcat că necesitatea utilizării tuturor acestor etape de preprocesare depinde de la un domeniu la altul, și de la un sistem de dialog la altul, fiind dictată în mare măsură și de particularitățile modulului de recunoaștere automată a vorbirii utilizat. De exemplu, unele sisteme de recunoaștere nu transcriu evenimente acustice non-lexicale, caz în care etapa 1 nu mai este necesară. Similar, anumite sisteme pot furniza numerele direct în forma cu cifre arabe, ceea ce ar face și etapa 2 inutilă. Deasemeni, funcție de domeniu, etapa de înlocuire de aliasuri poate apare ca oportună sau nu.

Pentru generalitate însă, în analizorul semantic proiectat în aceasta lucrare toate etapele de preprocesare prezentate sunt considerate necesare.

3.3 Modelul Markov

Modelele Markov pot fi considerate actualmente o metodă consacrată în recunoașterea automată a vorbirii și prelucrarea limbajului natural. Există cel puțin două aspecte care au condus la impunerea lor [?]: un prim avantaj vine din faptul că sunt bazate pe o teorie matematică riguroasă, cercetările fiind deci susținute de numeroase rezultate matematice acumulate anterior. Un al doilea avantaj este faptul că aceste modele pot fi antrenate pe baza unor date culese în condiții “reale”, lucru care face ca sistemele construite pe baza lor să prezinte o robustețe sporită.

În continuare se va face o scurtă prezentare teoretică a modelelor Markov, după care se va trece la descrierea algoritmilor și tehnicilor specifice utilizate în contextul antrenării modelului și al realizării decodării semantice pe baza lui.

3.3.1 Definiție. Probleme fundamentale

Definiția 1 Prin model Markov cu n stări ascunse și m simboluri observate vom înțelege un cvintuplu $\mathcal{M} = \langle \Omega, S, A, B, \Pi \rangle$, unde:

- $\Omega = \{o_1, o_2, \dots, o_m\}$ este o mulțime de simboluri observate sau generate;
- $S = \{s_1, s_2, \dots, s_n\}$ este mulțimea stărilor modelului;
- $A = \{a_{ij}\}$, $a_{ij} = P(s_j^{t+1} | s_i^t)$, $i, j = \overline{1, n}$, reprezintă matricea probabilităților tranzițiilor între stări, în care fiecare linie este o distribuție de probabilitate corespunzătoare unei stări;
- $B = \{b_{ik}\}$, $b_{ik} = P(o_k^t | s_i^t)$, $i = \overline{1, n}$, $k = \overline{1, m}$, reprezintă distribuțiile de probabilitate ale simbolurilor în stările modelului;
- $\Pi = \{\pi_i\}$, $\pi_i = P(s_i^1)$, $i = \overline{1, n}$, reprezintă distribuția de probabilitate a stării inițiale.

În esență un astfel de model este format dintr-o mulțime de stări S între care tranzițiile se fac cu anumite probabilități, date prin matricea A : $a_{ij} = P(s_j^{t+1} | s_i^t)$ definește probabilitatea ca modelul să ajungă în starea s_j la momentul $t + 1$ dacă se afla în starea s_i la momentul t , sau cu alte cuvinte, probabilitatea ca modelul să efectueze o tranziție din starea s_i în starea s_j . În plus, $\pi_i = P(s_i^1)$ reprezintă probabilitatea ca modelul să se afle în starea s_i la momentul inițial, deci Π modelează probabilistic starea inițială.

În fiecare din aceste stări, modelul consumă sau generează, funcție de modul în care este utilizat, un simbol din mulțimea Ω , cu o probabilitate dată de distribuția de probabilitate asociată stării respective – $b_{ik} = P(o_k^t | s_i^t)$ reprezintă deci probabilitatea în starea s_i a simbolului o_k .

Figura ?? ilustrează un model Markov cu 4 stări și 3 simboluri. După cum se observă, unele probabilități de tranziție sau ale simbolurilor pot fi nule. Mecanismul de funcționare al unui astfel de model este următorul: la momentul inițial, modelul intră într-o stare inițială conform distribuției de probabilitate a stării inițiale Π . În acea stare are loc emiterea unui simbol din mulțimea Ω conform distribuției de probabilitate a simbolurilor asociate stării respective. În continuare, se efectuează o tranziție spre o nouă stare în conformitate cu distribuția de probabilitate a tranzițiilor, urmată de o nouă emiteră a unui simbol, și așa mai departe. Specificarea unui model constă deci în definirea stărilor, a simbolurilor (în esență a numărului acestora), și a matricilor A , B , și Π .

Se poate spune deci că un model Markov cu stări ascunse (în continuare vom utiliza prescurtarea engleză HMM – Hidden Markov Model) este un model dublu

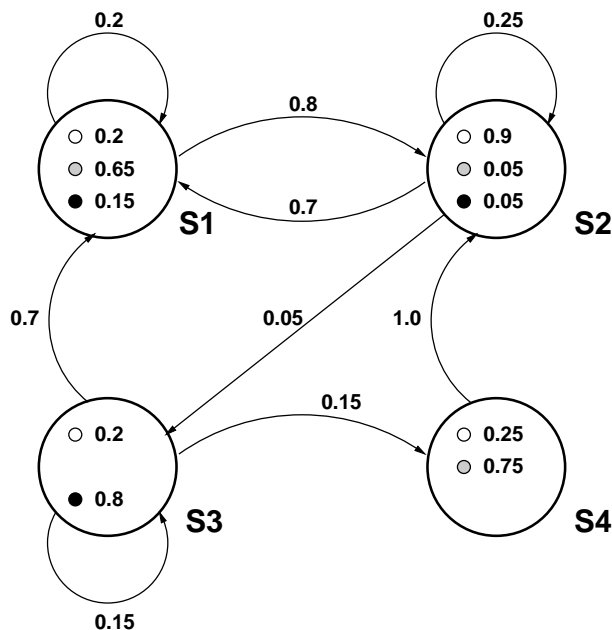


Figura 3.2: Exemplu de model Markov

stocastic [?]: el este format dintr-un proces stocastic ascuns (tranzițiile între stări), care nu poate fi observat decât prin intermediul unui alt set de procese stocastice, care generează simbolurile observate.

Utilizarea modelelor HMM în modelarea efectivă a unor procese presupune rezolvarea a trei probleme fundamentale [?]:

1. **Evaluarea HMM** - Determinarea probabilității $P(\mathcal{O}|\mathcal{M})$ ca un anumit model \mathcal{M} să fi generat o secvență de simboluri $\mathcal{O} = (q_1, q_2, \dots, q_T)$.
2. **Decodarea stărilor unui HMM** - Dată fiind o succesiune de simboluri $\mathcal{O} = (q_1, q_2, \dots, q_T)$ și un model \mathcal{M} , se cere determinarea secvenței de stări în care simbolurile ar fi putut fi generate, \mathcal{S} , optimă într-un anumit sens¹.
3. **Antrenarea HMM** - Ajustarea parametrilor unui model \mathcal{M} astfel încât să se maximizeze probabilitatea ca o secvență de simboluri $\mathcal{O} = (q_1, q_2, \dots, q_T)$ să fi fost generată de \mathcal{M} .

Pentru toate aceste probleme există soluții eficiente [?], care au impus HMM ca tehnică de modelare stocastică.

În continuare vom expune modul în care un astfel de model poate fi utilizat pentru decodarea semantică.

¹Se pot defini mai multe criterii de optim [?].

3.3.2 Modelele Markov în analiza semantică

Problema realizării analizei semantice este în esență o problemă de decodare. Scopul analizorului semantic este de a obține corespondentul semantic (într-un anumit formalism) al textului de intrare. Practic, în cazul sistemului proiectat, se cere ca, plecând de la textul de intrare normalizat, să determinăm șirul de etichete semantice corespunzătoare (parse).

În acest sens, analiza semantică privită ca o decodare a textului de intrare poate fi modelată prin intermediul unui HMM. Cuvintele din textul normalizat vor juca rolul de simboluri observate ale modelului, iar etichetele semantice vor corespunde stărilor. Astfel, dat fiind un model \mathcal{M} și o secvență de intrare \mathcal{O} (textul normalizat) se poate determina (problema 2) cea mai probabilă succesiune de stări interne \mathcal{S} care generează intrarea \mathcal{O} , cu alte cuvinte, cea mai probabilă succesiune de etichete semantice corespondente (vezi figura ??).

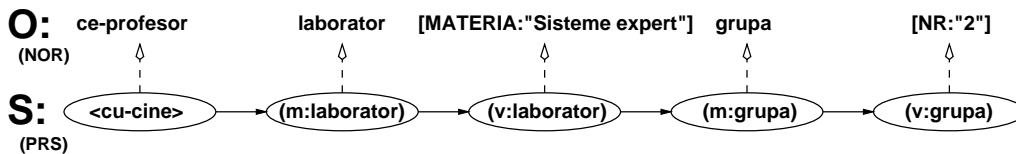


Figura 3.3: Utilizarea modelelor Markov în decodarea semantică

Rămâne problema creării modelului \mathcal{M} astfel încât decodarea să se desfășoare corect, adică respectând semantica domeniului și a limbii în discuție. Așa cum s-a menționat anterior, avantajul principal al metodei stocastice de analiză semantică este chiar faptul că modelul se crează prin antrenare pe baza unui corpus de formulări caracteristice domeniului de aplicație al sistemului. În mod evident, dimensiunile modelului (numărul de stări și numărul de simboluri observate) sunt date de numărul de etichete semantice și respectiv dimensiunea vocabularului normalizat. Odată dimensiunile stabilite, parametrii A , B și Π se pot determina pe baza corpusului de antrenament. Matricea A , reprezentând probabilitățile de tranziție între stări, va capta practic modurile posibile în care etichetele semantice se pot succeda. În mod similar, matricea B , reprezentând probabilitățile simbolurilor observate, va capta corespondența dintre etichetele semantice și lexemele textului, iar împreună, cele două matrice vor modela corespondența dintre textul normalizat și reprezentarea sa semantică.

3.3.3 Antrenarea modelului

Antrenarea modelului Markov are ca scop determinarea parametrilor lui astfel încât să modeleze cât mai corect (în sens statistic) semantica formulărilor, și se realizează pe baza unui corpus de antrenament compus din formulări normalizate și secvențele de etichete semantice corespunzătoare. Cuvintele ce apar în formulările

normalizate vor constitui simbolurile observate ale modelului, iar etichetele semantice vor corespunde stărilor.

Dimensiunile modelului sunt determinate anterior începerii antrenării: numărul etichetelor semantice distincte ce apar în corpus va fi numărul de stări (n), iar dimensiunea vocabularului formulărilor normalizate da numărul de simboluri (m). În aceste condiții, antrenarea se reduce la determinarea parametrilor a_{ij} , b_{ik} și π_i .

Antrenarea prin estimare de probabilitate maximă

Estimarea parametrilor A , B și Π se poate face pe baza numărării aparițiilor evenimentelor corespunzătoare (tranzițiile între stări și corespondențele stări-simboluri) în corpusul de antrenament. Introducem următoarele notații:

- ω_{ij} - reprezintă o tranziție între stările s_i și s_j , altfel spus o succesiune de două etichete semantice sau un *bigram* (s_i, s_j);
- $c(\omega_{ij})$ = numărul de apariții în corpus a unei tranziții între stările s_i și s_j (de câte ori apare în corpus bigram-ul (s_i, s_j));
- C_i = numărul de apariții în corpus a stării (etichetei semantice) s_i ca primă stare a unui bigram:

$$C_i = \sum_{j=1}^n c(\omega_{ij}) \quad (3.1)$$

- o_{ik} = numărul de generări ale simbolului observat o_k din starea s_i (de câte ori apare în corpus o corespondență între eticheta semantică s_i și cuvântul normalizat o_k);
- O_i = numărul total de simboluri generate în starea s_i :

$$O_i = \sum_{k=1}^m o_{ik} \quad (3.2)$$

- p_i = numărul de apariții în corpus a stării (etichetei) s_i pe prima poziție a unei formulări;
- P = numărul de formulări din corpus:

$$P = \sum_{i=1}^n p_i \quad (3.3)$$

Cu aceste notații, parametrii A , B și Π se calculează astfel:

$$a_{ij} = P(s_j^{t+1}|s_i^t) = \frac{c(\omega_{ij})}{C_i}, \quad i, j = \overline{1, n} \quad (3.4)$$

$$b_{ik} = P(o_k^t|s_i^t) = \frac{o_{ik}}{O_i}, \quad i = \overline{1, n}, \quad k = \overline{1, m} \quad (3.5)$$

$$\pi_i = P(s_i^1) = \frac{p_i}{P}, \quad i = \overline{1, n} \quad (3.6)$$

formulele ??, ?? și ?? asigurând faptul că a_{ij} , b_{ik} și π_i satisfac condiția de normalizare a distribuțiilor de probabilitate : $\sum_{j=1}^n a_{ij} = 1$, $\sum_{k=1}^m b_{ik} = 1$ și $\sum_{i=1}^n \pi_i = 1$.

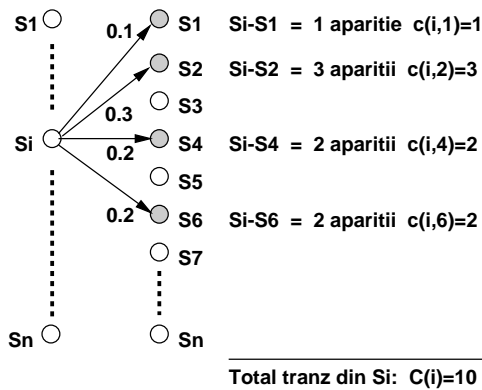


Figura 3.4: Estimarea de probabilitate maximă

de apariții în corpus a tranziției între etichetele semantice corespunzătoare. Similar, distribuțiile de probabilitate ale simbolurilor pentru fiecare stare reflectă proporțiile corespondențelor etichetă-cuvânt din corpusul de antrenament.

Observație: Dacă o anumită succesiune de etichete semantice nu apare în corpusul de antrenament, parametrul $c(\omega_{ij})$ respectiv va fi zero, și în consecință și probabilitatea calculată a_{ij} a tranziției între stările respective va fi nulă (fig. ??).

Ameliorarea modelului prin metoda Katz

Observația anterioară introduce și o problemă în ceea ce privește utilizarea cu succes a unui model Markov în decodarea semantică. În condițiile unei dimensiuni relativ reduse a corpusului de antrenament, este foarte probabil ca acesta să nu surprindă toate succesiunile valide de etichete semantice. În consecință, modelul va conține o mulțime de parametri nuli, care grevează faza de decodare (vezi

Coefficienții $c(\omega_{ij})$, o_{ik} și p_i pot fi ușor determinați prin numărare dintr-o singură parcurgere a corpusului de antrenament, antrenarea având deci complexitatea $O(n)$ în raport cu dimensiunea acestuia. Acest mod de estimare a parametrilor modelului poartă numele de *estimare de probabilitate maximă* (Maximum Likelihood Estimation), datorită faptului că maximizează probabilitatea ca simbolurile observate să fi fost generate de modelul rezultat. Așa cum ilustrează și figura ??, probabilitatea de tranziție între două stări este proporțională cu numărul

??), în sensul excluderii totale a succesiunilor de etichete semantice care nu au fost observate. Evenimentele slab reprezentate, sau care nu sunt reprezentate în corpus, nu vor fi deci modelate corect dacă utilizăm metoda de antrenare prin estimare de probabilitate maximă.

Pentru ameliorarea situațiilor de acest tip se pot utiliza o serie de *metode de netezire*, prin care se realizează o ameliorare a modelului în condiții de raritate a datelor, astfel încât să fie admisibile și evenimente neobservate. În analizorul proiectat am utilizat în acest scop metoda de reestimare Katz [?], metodă al cărei suport teoretic pentru modelele la nivel de bigrame va fi prezentat în continuare (pentru varianta generalizată la m -grame vezi [?]).

Ideea fundamentală a reestimării Katz este de a micșora într-o anumită măsură probabilitatea bigramelor observate, și de a redistribui masa de probabilitate astfel eliberată între bigramele neobservate în corpusul de antrenament.

În cele ce urmează vom nota cu n_k numărul bigramelor distincte care apar în corpusul de antrenament de k ori și cu N numărul total de bigrame din corpus. Atunci, este satisfăcută relația

$$N = \sum_k k \cdot n_k \quad (3.7)$$

Reestimarea Katz înlocuiește estimările de probabilitate maximă ale probabilităților bigramelor calculate la nivelul întregului corpus:

$$P(\omega_{ij}) = \frac{c(\omega_{ij})}{N} \quad (3.8)$$

cu estimările Turing corespunzătoare, definite prin relația:

$$P_T(\omega_{ij}) = \frac{c^*(\omega_{ij})}{N} \quad (3.9)$$

unde:

$$c^*(\omega_{ij}) = (c(\omega_{ij}) + 1) \cdot \frac{n_{c(\omega_{ij})+1}}{n_{c(\omega_{ij})}} \quad (3.10)$$

Raportul $\frac{c^*(\omega_{ij})}{c(\omega_{ij})}$ poartă denumirea de *coeficient de reducere Turing* (Turing discount coefficient), și se notează cu $d_{c(\omega_{ij})}$. Cu aceste notații este evident că:

$$P_T(\omega_{ij}) = d_{c(\omega_{ij})} \cdot P(\omega_{ij}) \quad (3.11)$$

Atunci, pe baza relațiilor de mai sus și ținând cont că:

$$\sum_{\omega_{ij}:c(\omega_{ij})>0} P(\omega_{ij}) = 1, \quad (3.12)$$

se poate arăta că masa de probabilitate F eliberată prin utilizarea estimărilor Turing este:

$$\begin{aligned}
F &= \sum_{\omega_{ij}:c(\omega_{ij})>0} \left(P(\omega_{ij}) - P_T(\omega_{ij}) \right) \\
&= \sum_{r>0} n_r \cdot (1 - d_r) \cdot \frac{r}{N} = \frac{n_1}{N}
\end{aligned} \tag{3.13}$$

În același timp:

$$\sum_{\omega_{ij}:c(\omega_{ij})=0} P_T(\omega_{ij}) = 1 - \sum_{\omega_{ij}:c(\omega_{ij})>0} P_T(\omega_{ij}) = \frac{n_1}{N} = F \tag{3.14}$$

Masa de probabilitate eliberată F poate fi deci redistribuită bigramelor care nu au apărut în corpusul de antrenament. Redistribuirea se poate face proporțional cu probabilitățile Turing de apariție în corpus a stărilor spre care se efectuează tranziția. Atunci, dacă a_{ij} sunt probabilitățile tranzițiilor calculate conform (??), ele se modifică astfel:

$$c(\omega_{ij}) > 0 : \quad a'_{ij} = d_{c(\omega_{ij})} \cdot a_{ij} \tag{3.15}$$

$$c(\omega_{ij}) = 0 : \quad a'_{ij} = \frac{1 - \sum_{\omega_{ik}:c(\omega_{ik})>0} d_{c(\omega_{ik})} \cdot a_{ij}}{\sum_{\omega_{ik}:c(\omega_{ik})=0} P_T(\omega_k)} \cdot P_T(\omega_j) \tag{3.16}$$

În cele prezentate anterior, reducerea probabilității prin înlocuirea cu estimările Turing corespunzătoare s-a făcut pentru toate bigramele reprezentate în corpus, pentru care $c(\omega_{ij}) > 0$ – vezi ecuația (??).

Metoda de reestimare Katz propune însă lăsarea intactă a probabilităților bigramelor care apar de un număr relativ mare de ori în corpus, întrucât pentru ele se poate presupune că estimările sunt de încredere (reliable) [?]. Reducerea probabilităților prin înlocuirea lor cu estimările Turing se va face doar pentru bigramele slab reprezentate în corpus, de un număr de ori cel mult egal cu un parametru fixat K (vezi figura ??). În acest caz, pentru a păstra masa de probabilitate eliberată și distribuită bigramelor inexistente în corpus la aceeași valoare dată de formula (??), este necesară o recalculare a coeficienților de reducere:

$$c(\omega_{ij}) > K : \quad d'_{c(\omega_{ij})} = 1, \tag{3.17}$$

$$1 \leq c(\omega_{ij}) \leq K : \quad d'_{c(\omega_{ij})} = \frac{d_{c(\omega_{ij})} - \frac{(K+1) \cdot n_{K+1}}{n_1}}{1 - \frac{(K+1) \cdot n_{K+1}}{n_1}} \tag{3.18}$$

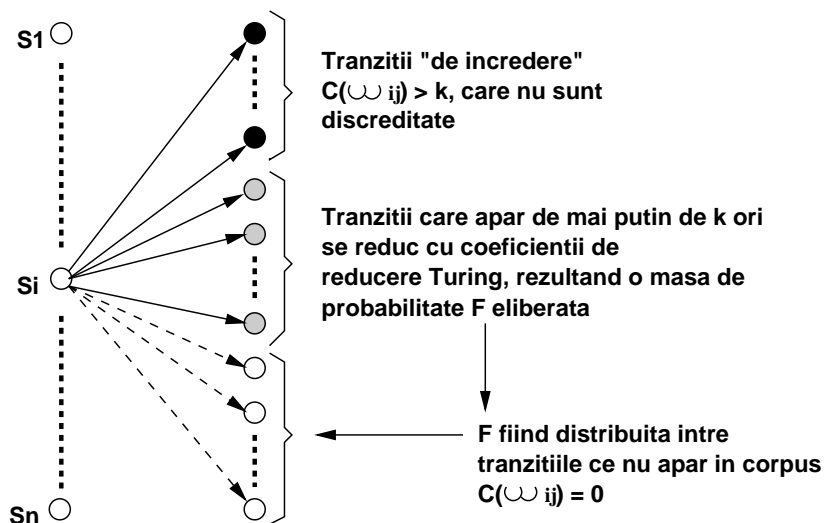


Figura 3.5: Principiul reestimării Katz

reestimarea făcându-se tot după formulele (??) și (??), dar cu coeficienții de reducere d' .

Practic, algoritmul de reestimare Katz pornește de la estimările de probabilitate maximă a_{ij} , determină coeficienții de reducere Turing d și d' , după care pe baza formulelor (??) și (??) reestimează matricea A .

În mod similar, algoritmul de reestimare Katz se poate aplica și distribuției de probabilitate a stării inițiale Π prin ajustarea estimărilor de probabilitate maximă π_i cu un set corespunzător de coeficienți de reducere.

Se elimină astfel valorile nule din A și Π , măbind robustețea și fiabilitatea modelului în fața unor evenimente (structuri semantice) care nu au apărut în corpusul de antrenament.

3.3.4 Decodarea semantică

Problema decodării semantice bazată pe modelul Markov revine la a determina succesiunea de stări \mathcal{S} cea mai probabilă să fi generat o secvență de simboluri \mathcal{O} .

Fie $\mathcal{O} = (q_1, q_2, \dots, q_T)$ o secvență de simboluri observate și $\mathcal{S} = (t_1, t_2, \dots, t_T)$ o succesiune de stări de lungime T . Atunci probabilitatea ca un model \mathcal{M} să fi generat secvența \mathcal{O} trecând prin stările \mathcal{S} este:

$$\begin{aligned}
 P(\mathcal{O}|\mathcal{S}, \mathcal{M}) &= \pi_{t_1} b_{t_1 q_1} a_{t_1 t_2} b_{t_2 q_2} a_{t_2 t_3} b_{t_3 q_3} \dots a_{t_{T-1} t_T} b_{t_T q_T} \\
 &= \pi_{t_1} b_{t_1 q_1} \prod_{k=1}^{T-1} a_{t_k t_{k+1}} b_{t_{k+1} q_{k+1}}
 \end{aligned} \tag{3.19}$$

Problema decodării se reduce deci la găsirea succesiunii de stări \mathcal{S} care maximizează această probabilitate:

$$\mathcal{S} = \arg \max_{\mathcal{S}} \{P(\mathcal{O}|\mathcal{S}, \mathcal{M})\}$$

Din ecuațiile de mai sus se poate observa cu ușurință că o căutare exhaustivă în spațiul succesiunilor de stări \mathcal{S} nu este practică, având o complexitate $O(Tn^T)$, unde n este numărul de stări ale modelului și T este lungimea secvenței observate. Un algoritm mai eficient pentru determinarea secvenței \mathcal{S} este algoritmul Viterbi [?], prezentat în continuare.

Algoritmul Viterbi este un algoritm de programare dinamică care calculează iterativ, pentru subșirurile prefix ale secvenței observate de lungimi $1, 2, \dots, T$, valorile maxime ale probabilității P , precum și succesiunile de stări corespunzătoare. Astfel, succesiunea de stări \mathcal{S} este construită incremental, pornind de la prefixele $(q_1), (q_1, q_2), (q_1, q_2, q_3)$, până la secvența completă $\mathcal{O} = (q_1, \dots, q_T)$.

Se definește variabila $\delta_k(j)$ care exprimă probabilitatea maximă ca modelul \mathcal{M} să fi generat prefixul de lungime $k - (q_1, q_2, \dots, q_k)$. Atunci δ poate fi exprimat recursiv astfel:

$$\delta_1(i) = \pi_i \cdot b_{iq_1} \quad i = \overline{1, n} \quad (3.20)$$

$$\delta_k(j) = \max_{1 \leq i \leq n} \{\delta_{k-1}(i) \cdot a_{ij}\} \cdot b_{jq_k} \quad j = \overline{1, n} \quad (3.21)$$

caz în care probabilitatea maximă căutată este:

$$\max_{\mathcal{S}} \{P(\mathcal{O}|\mathcal{S}, \mathcal{M})\} = \max_{1 \leq i \leq n} \{\delta_T(i)\} \quad (3.22)$$

Scopul algoritmului de decodare Viterbi fiind determinarea succesiunii de stări corespunzătoare acestei valori maxime, este necesară și reținerea succesiunii stărilor care maximizează δ . În acest scop, se definește:

$$\psi_k(j) = \arg \max_{1 \leq i \leq n} \{\delta_{k-1}(i) \cdot a_{ij}\} \quad j = \overline{1, n} \quad (3.23)$$

Algoritmul este prezentat în pseudocod în figura ?? și ilustrat în figura ?. El începe cu o etapă de inițializare în care sunt determinate valorile δ_{1i} . Apoi, în etapa iterativă se ia în considerare pe rând câte un nou simbol din secvența \mathcal{O} . Pentru fiecare stare, se determină probabilitatea cu care modelul poate ajunge în starea respectivă și genera secvența observată până la simbolul curent (se calculează valoarea δ_k pe baza valorilor deja calculate δ_{k-1}). În același timp se reține pentru fiecare stare care este starea anterioară care a realizat maximul probabilității (valorile ψ).

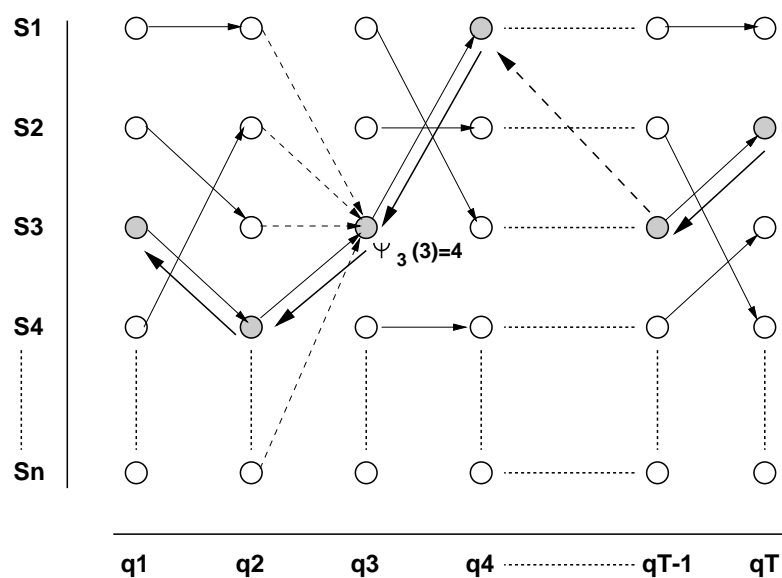


Figura 3.6: Algoritmul Viterbi

```

/* -- inițializare ----- */
pentru i de-la 1 la n
   $\delta_1(i) = \pi_i \cdot b_{iq_1}$ 
   $\psi_1(i) = 0$ 
/* -- iterații ----- */
pentru k de-la 2 la T
  pentru j de-la 1 la n
     $\delta_k(j) = \max_{1 \leq i \leq n} \{\delta_{k-1}(i) \cdot a_{ij}\} \cdot b_{jq_k}$ 
     $\psi_k(j) = \arg \max_{1 \leq i \leq n} \{\delta_{k-1}(i) \cdot a_{ij}\}$ 
/* -- terminare ----- */
 $P_{MAX} = \max_{1 \leq i \leq n} \{\delta_T(i)\}$ 
 $t_T = \arg \max_{1 \leq i \leq n} \{\delta_T(i)\}$ 
/* -- parcurgere inversă --- */
pentru k de-la T-1 la 1
   $t_k = \psi_{k+1}(t_{k+1})$ 

```

Figura 3.7: Pseudocod Viterbi

După ce toate simbolurile secvenței observate au fost prelucrate, valoarea probabilității maxime este dată, conform (??), de maximul valorilor δ_{T_i} . Starea care realizează acest maxim este practic ultima stare t_T din \mathcal{S} . Pornind de la ea, structura descrisă în figură este parcursă în sens invers pe baza indicilor ψ , identificându-se stările ce compun secvența decodată \mathcal{S} .

O analiză simplă a complexității în timp a algoritmului Viterbi arată că aceasta este de ordinul $O(Tn^2)$, deci polinomială și mult redusă față de căutarea exhaustivă care are complexitatea $O(Tn^T)$. Algoritmul Viterbi oferă deci o soluție eficientă pentru decodarea stărilor unui model Markov pornind de la o secvență observată dată. Împreună cu algoritmi de antrenare a modelului, el constituie baza pentru realizarea analizei semantice în sistemul proiectat.

3.4 Constructorul de cadre

Constructorul de cadre este ultimul modul din analizorul semantic. Întrucât decodorul furnizează doar secvența de etichete semantice corespunzătoare textului de intrare, este necesar acest modul suplimentar pentru a realiza translatarea în reprezentarea prin cadre (figura ??).

Text normalizat	cu-cine grupă [NR:“2”] [ZI-SAPT:“luni”] [PERIOADA:“dimineață”] laborator [MATERIA:“Inteligență artificială”]
Etichete semantice	<cu-cine> (m:grupa) (v:grupa) (v:zi-sapt) (v:perioada-zi) (m:laborator) (v:laborator)
Cadru	<cu-cine> <identificare> grupa = 2 <specif-timp> zi-sapt = luni perioada-zi = dimineață <specif-materie> laborator = “Inteligență artificială”
Cadru normalizat	<cu-cine> <identificare> grupa = 2 <specif-timp> zi-sapt = 1 perioada-zi = DIM <specif-materie> laborator = C024

Figura 3.8: Corespondența succesiune de etichete semantice – cadru

Formalismul case-grammar implică trei tipuri de etichete semantice: concepte – reprezentate aici sub forma <concept>, marcaje – reprezentate sub forma (m:marcaj), și valori – reprezentate sub forma (v:valoare). Pentru a realiza instanțierea cadrelor, constructorul se bazează însă numai pe concepte și valori. Pe baza conceptelor sunt determinate cadrele care trebuie instanțiate, iar valorile sunt utilizate pentru completarea sloturilor. Marcajele nu au un rol direct în construirea cadrelor, ele modelând limbajul doar la nivelul decodării semantice.

Algoritmul utilizat de constructorul de cadre este relativ simplu, însă are câteva particularități importante. Într-un prim pas se identifică printr-o parcurgere a succesiunii de etichete semantice toate conceptele care apar. Pe baza lor, cunoscând întregul sistem de cadre utilizat, constructorul instanțiază cadrele corespunzătoare, împreună cu eventualele subcadre ale acestora. Trebuie remarcat că, datorită posibilității apariției unor erori de recunoaștere și modalității

stocastice de realizare a decodării semantice, nu există nici o garanție în privința faptului că succesiunea de etichete semantice va conține un concept: într-o astfel de succesiune pot exista nici una, una, sau mai multe etichete-concept. Primul caz prezintă anumite particularități, și va fi tratat separat.

După identificarea unuia sau mai multor concepte și instanțierea cadrelor respective, constructorul realizează o nouă parcurgere a succesiunii de etichete semantice, de această dată urmărind numai valorile. Fiecare valoare este atașată la toate cadrele instanțiate care au un slot corespunzător; cu alte cuvinte, dacă există mai multe concepte în formulare, se încearcă completarea fiecăruia din cadrele instanțiate cu toate valorile ce apar în formulare.

Un ultim pas îl constituie normalizarea cadrelor. Această operație este similară normalizării realizate prin preprocesare. Rolul ei este de a adapta cadrul pentru utilizarea lui în continuare de către modulul de control al dialogului spre care va fi transmis. În principiu această operație unifică diverse sinonime care pot apărea ca valori ale sloturilor, și eventual poate realiza conversii ale acestor valori pentru a ușura utilizarea lor ulterioară, de exemplu valorile sloturilor “laborator” sunt transformate din valori textuale în coduri de acces la baza de date utilizată.

Formulări fără concept

Un caz aparte îl constituie formulările utilizator care după decodare nu conțin nici o etichetă semantică concept. Evident, această situație necesită o tratare specială prin care să se determine ce cadru trebuie instanțiat.

Pentru exemplificare, vom porni de la următoarea cerere de detalii ipotetică pe care un sistem de dialog o adresează unui utilizator:

- Dacă sunteți student, precizați anul, grupa, subgrupa.

Două posibile răspunsuri sunt ilustrate în figura ?? . Deși ambele formulări sunt valide semantic, și decodorul realizează cu succes obținerea succesiunii de etichete semantice corecte în ambele cazuri, al doilea răspuns nu poate fi transformat direct în reprezentarea prin cadre datorită lipsei unui concept.

Sistem	Dacă sunteți student, precizați anul, grupa, subgrupa.
Utilizator 1	sunt student în anul întâi grupa doi subgrupa doi
Parse 1	<identificare> (m:anul) (v:anul) (m:grupa) (v:grupa) (m:subgrupa) (v:subgrupa)
Utilizator 2	anul întâi grupa doi subgrupa doi
Parse 2	(m:anul) (v:anul) (m:grupa) (v:grupa) (m:subgrupa) (v:subgrupa)

Figura 3.9: Formulări fără concept

Această lipsă a conceptului se manifestă practic datorită lipsei unui cuvânt în formularea utilizatorului care să declanșeze un concept (cum este cuvântul “sunt”

din prima formulare). Sensul propoziției este însă clar dacă se ia în considerare și întrebarea sistemului. Soluția utilizată în aceste situații este ca modulul de control al dialogului să transmită analizorului semantic care sunt conceptele așteptate de la utilizator în momentul respectiv.

Rolul modulului de control al dialogului este de a realiza gestionarea dialogului și a interacțiunii cu utilizatorul, și, în acest sens, el este capabil să precizeze (în cel mai rău caz cu o anumită probabilitate) care este conceptul așteptat de la utilizator. Trebuie remarcat că majoritatea formulărilor fără concept apar chiar ca urmare a unor întrebări de clarificare din partea sistemului, deoarece în aceste cazuri întrebarea sistemului creează contextul în care apare răspunsul: în toate aceste situații modulul de control al dialogului a lansat întrebarea, deci el “știe” în principiu care este conceptul așteptat de la utilizator (de exemplu, după întrebarea de mai sus, este clar că sistemul așteaptă conceptul <identificare>). Ceva mai delicate sunt situațiile în care modulul de control se află într-un punct în care este rândul utilizatorului să preia inițiativa dialogului. În acest caz, el poate furniza doar probabilități asociate conceptelor așteptate.

Pentru a unifica situațiile am presupus că, la analiza semantică a fiecărei formulări utilizator, modulul de control al dialogului furnizează o listă de concepte așteptate, însoțită de probabilitățile asociate fiecăruia. Această listă este utilizată în cazul în care decodarea semantică nu reușește doar pe baza formulării utilizatorului (nu este detectat un concept). În această situație se reia decodarea introducând pe rând câte un concept din listă în fața succesiunii observate (astfel de secvențe extinse sunt utilizate și în antrenarea modelului – vezi anexa ??). De această dată conceptul va fi identificat, iar în final este aleasă varianta care maximizează produsul dintre probabilitatea conceptului și probabilitatea de observare a succesiunii de simboluri extinsă rezultată din decodare.

Dacă $\mathcal{L} = (c_1, c_2, \dots, c_t)$ este lista de concepte primite, cu probabilitățile corespundente (p_1, p_2, \dots, p_t) , iar $\mathcal{O} = (q_1, q_2, \dots, q_T)$ este succesiunea de simboluri observate care prin decodare nu duc la apariția unui concept, atunci practic se încearcă decodarea tuturor secvențelor $\mathcal{O}_i = (c_i, q_1, q_2, \dots, q_T)$. Dintre acestea este aleasă \mathcal{O}_k , unde $k = \arg \max_{1 \leq i \leq t} \{p_i \cdot P(\mathcal{O}_i | \mathcal{M})\}$.

Observație: Se poate spune că, într-o anumită măsură, prin acest mecanism de transmitere a conceptelor așteptate de către modulul de control al dialogului, analizorul semantic realizează și o analiză contextuală parțială.

3.5 Detalii de implementare

În cele ce urmează vom prezenta diverse aspecte privind implementarea analizorului semantic proiectat în secțiunile anterioare.

Implementarea efectivă a fost realizată în limbajul C++, ales în principal datorită eficienței, a capabilităților de programare orientată pe obiecte și a portabilității codului scris în acest limbaj. Dezvoltarea propriu-zisă a fost făcută pentru

sistemul de operare Linux, utilizând mediul integrat de dezvoltare KDevelop [?].

Forma finală în care analizorul semantic este distribuit este cea a unei biblioteci statice de clase, pentru a fi cu ușurință inclus într-un sistem de dialog proiectat într-un cadru mai larg. Pe baza acestei biblioteci pot fi realizate implementări ale analizorului semantic sub diverse forme: de la programe filtru care lucrează în pipe, până la servere la nivel sockets, RPC (Remote Procedure Call) sau RMI (Remote Method Invocation). Deși în figurile care urmează nu sunt ilustrați constructorii și destructorii claselor, menționăm că toate sunt scrise în forma ortodox-canonică (pentru o descriere completă a membrilor și interfețelor vezi anexa ??).

În continuare prezentăm principalele clase, relațiile dintre ele, precum și algoritmi utilizați în realizarea funcționalității lor. Se poate face o grupare a lor în două categorii: clase care modelează datele specifice problemei analizei semantice pe baze stocastice, și clase care implementează algoritmi utilizați.

3.5.1 Modelarea datelor

Prima categorie, a claselor care modelează datele, cuprinde clasele `CSequence`, `CStringSequence`, `CCaseFrame`, `CCaseFrameSystem`, `CUtterance`, `CUtteranceCorpus` și `CDictionary`, ilustrate în diagrama UML (Unified Modeling Language) din figura ?. Figura prezintă relațiile existente între aceste clase, precum și membrii și metodele lor².

- `CStringSequence` - modelează o succesiune de șiruri de caractere de lungime variabilă, fiind utilizată în principal pentru reprezentarea formei normalizate a unei formulări utilizator și a succesiunii de etichete semantice corespunzătoare. Implementarea se bazează pe tablouri alocate dinamic care se redimensionează, când este necesar, prin creșterea cu o cantitate fixată DELTA.
- `CSequence` - similară cu `CStringSequence`, modelează o secvență de numere. Este utilizată de algoritmi de antrenare și decodare, care lucrează cu indecșii șirurilor de caractere într-un dicționar în locul șirurilor respective.
- `CCaseFrame` - modelează cadrele specifice formalismului case-grammar. În esență, fiecare obiect cadru are un concept (membrul `Concept`) și o listă de atribute și valori (membrii `Attributes` și `Values`). Este de remarcat că valorile pot fi șiruri de caractere sau alte cadre, situație modelată prin unionul `CCaseValue`. Clasa oferă metode pentru citirea definiției cadrului dintr-un fișier, citirea și scrierea instanței cadrului din/într-un buffer, setarea și obținerea valorilor atributelor.
- `CCaseFrameSystem` - modelează un sistem de cadre conform formalismului case-grammar. Un astfel de sistem constituie practic o mulțime de cadre (`CaseFrames`), între care pot exista relații de tip părinte-fiu. Clasa permite

²Sunt reprezentate doar membrii și metodele considerate semnificative.

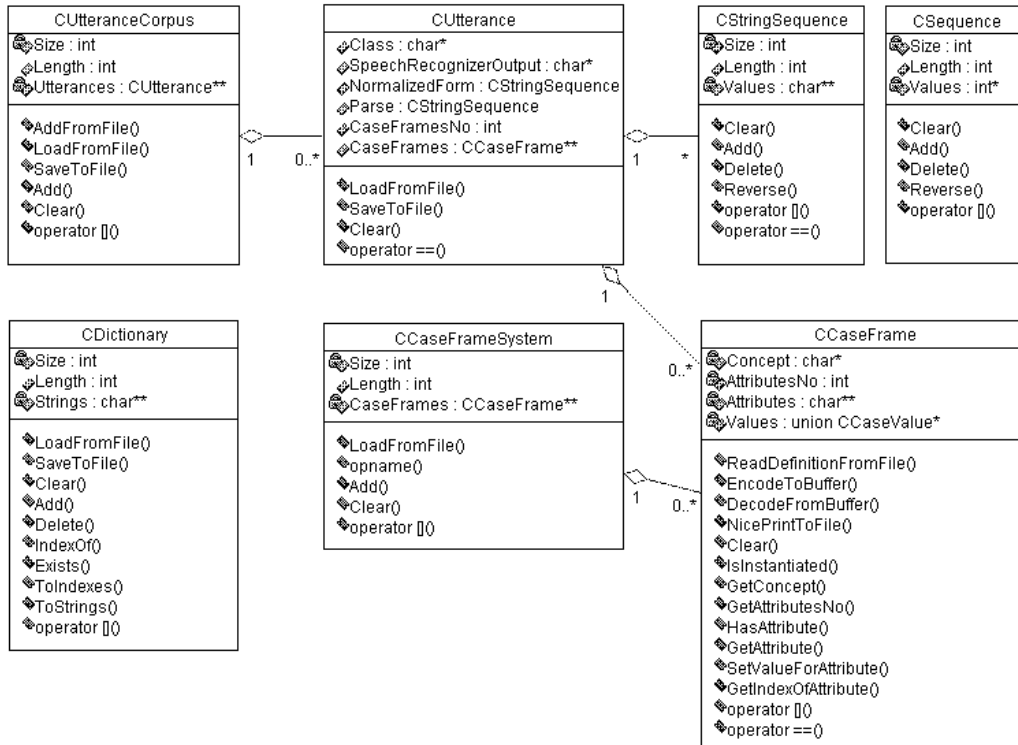


Figura 3.10: Clase pentru modelarea datelor

citirea specificației unui sistem de cadre dintr-un fișier de control, de tipul celui ilustrat în anexa ??.

- **CUtterance** - modelează o formulare utilizator, cu toate reprezentările prin care trece pe parcursul realizării analizei semantice: ieșire a sistemului de recunoaștere automată a vorbirii (*SpeechRecognizerOutput*), text normalizat (*NormalizedForm*), succesiune de etichete semantice (*Parse*) și reprezentare prin cadre (*CaseFrames*). Sunt implementate metode pentru citirea și scrierea formulării din/într-un fișier, utilizând formatul de corpus din anexa ??.
- **CUtteranceCorpus** - modelează un corpus de formulări utilizator, oferind metode pentru citirea și scrierea lui din/într-un fișier. Implementarea se bazează pe tablouri alocate dinamic de obiecte de tipul *CUtterance*.
- **CDictionary** - modelează un dicționar, implementând metodele caracteristice acestei structuri de date. Implementarea este bazată pe un tablou alocat dinamic, iar căutarea este liniară³.

³Dimensiunile mici ale dicționarelor implicate fac suficientă o astfel de implementare.

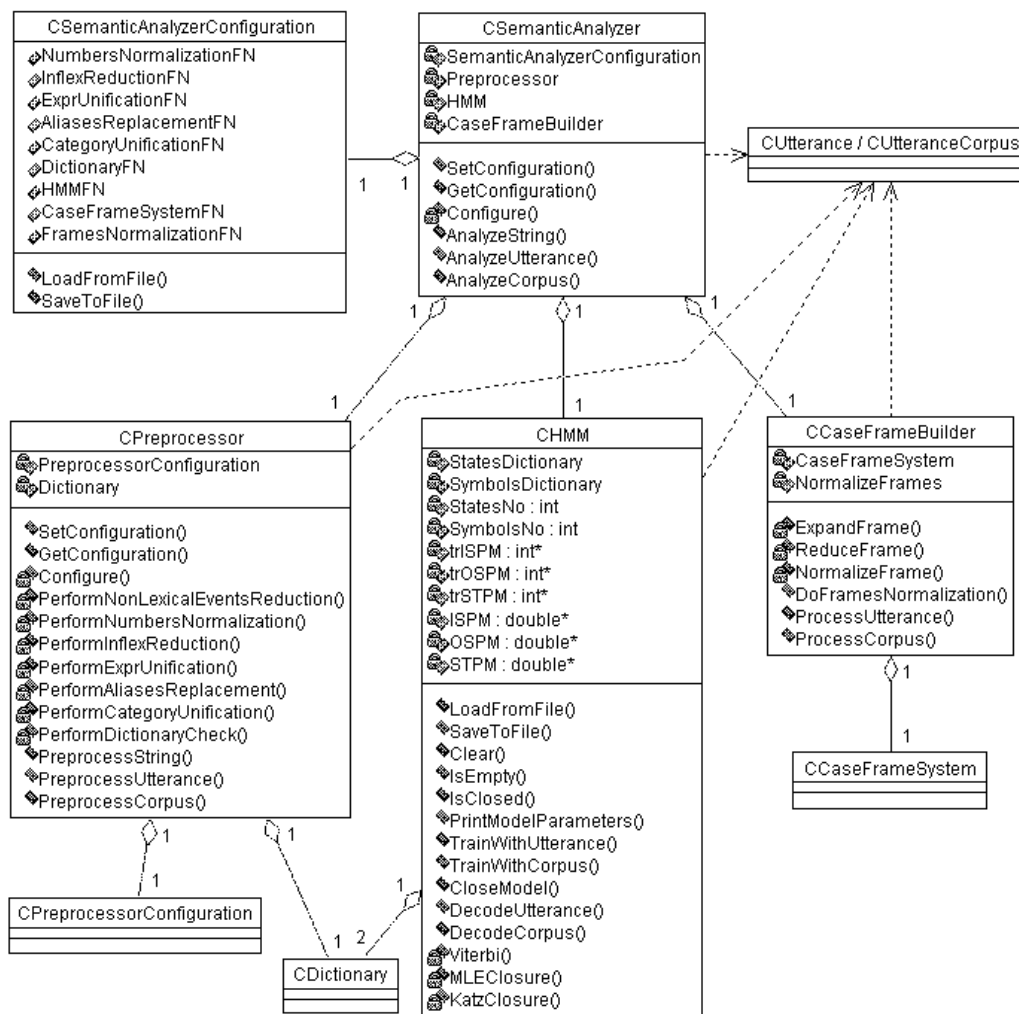


Figura 3.11: Clase pentru implementarea algoritmilor

3.5.2 Implementarea algoritmilor

Clasele care implementează algoritmi implicați în rezolvarea problemelor analizei semantice pe baze stocastice sunt prezentate în diagrama UML din figura ???. Pe lângă acestea, mai există și o serie de clase auxiliare, introduse pentru o mai bună modularizare și pentru creșterea lizibilității codului.

- **CPreprocessor** - implementează funcționalitatea modului de preprocesare. Controlul preprocesorului se realizează pe baza unui obiect membru de tip **CPreprocessorConfiguration** prin care se specifică pașii de preprocesare ce trebuie efectuați și fișierele care îi controlează. Clasa oferă metode pentru preprocesarea unui șir de caractere – **PreprocessString**, a unei formulări uti-

lizator – `PreprocessUtterance`, sau a unui întreg corpus – `PreprocessCorpus`. Realizarea efectivă a preprocesării se face prin metodele private `Perform`. Eliminarea evenimentelor non-lexicale se realizează prin simpla parcurgere a textului de intrare și eliminarea oricărui text în paranteze. Etapele de normalizare a numerelor, reducere a flexiunilor, unificare a expresiilor, înlocuire a aliasurilor, și unificare a categoriilor, se realizează prin înlocuiri de șiruri de caractere, specificate prin intermediul unor fișiere de control (vezi anexa ??). A fost construit un analizor lexical și sintactic pentru aceste fișiere, și prin analizarea lor se generează tabelele de înlocuiri care vor fi utilizate pentru efectuarea fiecăruia din pașii de preprocesare enumerați.

- **CHMM** - implementează modelul Markov, incluzând parametrii acestuia și algoritmi de antrenare și decodare utilizați. În componența modelului intră dicționarele de simboluri și etichete semantice (`SymbolsDictionary` și `StatesDictionary`), precum și distribuțiile de probabilitate A (`STPM` – State Transition Probability Matrix), B (`OSPM` – Observation Symbol Probability Matrix) și Π (`ISPM` – Initial State Probability Matrix). Variantele `trISPM`, `trOSPM`, `trSTPM` sunt utilizate în timpul antrenării modelului pentru a număra evenimentele respective. Clasa `CHMM` este serializabilă, existând metode pentru salvarea și citirea modelelor în/din fișiere. Clasa oferă metode pentru antrenarea modelului folosind o formulare sau un întreg corpus (`TrainWithUtterance` și `TrainWithCorpus`). Metoda `CloseModel` realizează închiderea modelului și calculează matricile de probabilitate `ISPM`, `OSPM` și `STPM` prin estimare de probabilitate maximă sau reestimare Katz (pe baza metodelor private `MLEClosure` și `KatzClosure`). Deasemeni clasa oferă metode pentru realizarea decodării semantice (`DecodeUtterance` și `DecodeCorpus`) pe baza algoritmului Viterbi.
- **CCaseFrameBuilder** - această clasă corespunde modulului de construcție a cadrelor. Ea utilizează un sistem de cadre (`CaseFrameSystem`) care specifică toate cadrele posibile, și pune la dispoziție o metodă pentru construcția reprezentării prin cadre a unei formulări utilizator (`ProcessUtterance`).
- **CSemanticAnalyzer** - modelează întreaga linie de analiză semantică, formată din modulele de preprocesare, decodare semantică și construcție a cadrelor. Configurarea analizorului se realizează prin intermediul unui obiect serializabil de tipul `CSemanticAnalyzerConfiguration`, care specifică parametrii necesari pentru fiecare din modulele implicate. Clasa oferă metode pentru analiza semantică a unei formulări utilizator dată ca șir de caractere sau ca obiect `CUtterance`. Aceste metode (`AnalyzeString` și `AnalyzeUtterance`) au ca parametru opțional, pe lângă textul sau formularea utilizatorului, o listă de concepte cu probabilitățile corespunzătoare, utilizată conform algoritmului descris în secțiunea ?? pentru tratarea formulărilor fără concept. Deasemeni, clasa `CSemanticAnalyzer` oferă o metodă (`AnalyzeCorpus`) pentru realizarea analizei semantice a unui întreg corpus de formulări.

3.5.3 Programe utilitare

Pe baza bibliotecii menționate au fost dezvoltate o serie de programe ce pot fi utilizate pe parcursul dezvoltării unui sistem de dialog pentru procesarea corpusurilor și antrenarea și evaluarea modelelor Markov implicate. În continuare se va face o scurtă descriere a fiecăruia (anexa ?? prezintă sumar instrucțiunile de utilizare, afișate la apelul fără parametri în linia de comandă).

- **hmm_build**
Acest program constituie practic modulul de antrenare al modelului Markov. El operează în trei etape: crearea, antrenarea și închiderea modelului. La crearea este necesară specificarea vocabularului textului normalizat și a mulțimii etichetelor semantice. Etapa de antrenare se realizează pe baza unui corpus de antrenament (pentru structura corpusului, vezi anexa ??). În această etapă se realizează practic doar numărarea evenimentelor. În etapa de închidere a modelului sunt calculați efectiv parametrii A , B și Π prin estimarea de probabilitate maximă, eventual cu aplicarea reestimării Katz (caz în care trebuie specificați și parametrii corespunzători).
- **hmm_print**
Afișează parametrii A , B și Π ai unui model Markov într-o formă lizibilă pentru proiectanți. Dacă modelul nu este închis, sunt afișate numărul tranzițiilor și al observațiilor, împreună cu distribuțiile bigramelor și stărilor inițiale (utile pentru studiul reestimării Katz). Pentru modelele închise sunt listate distribuțiile de probabilitate din A , B și Π .
- **hmm_decode**
Realizează decodarea semantică a unui corpus de formulări utilizator preprocesate (incluzând forma NOR - vezi anexa ??). Trebuie specificat fișierul conținând modelul Markov pe baza căruia se realizează decodarea.
- **preproc**
Acesta este modulul de preprocesare a intrării. Programul operează asupra unui corpus de formulări utilizator în forma furnizată de modulul de recunoaștere a vorbirii (SRO), și realizează preprocesarea lor în conformitate cu instrucțiunile transmise prin linia de comandă. Se pot specifica etapele de preprocesare care să fie efectuate, iar controlul lor se realizează prin intermediul unor fișiere de configurare (vezi anexa ??).
- **frame_build**
Operează asupra unui corpus de formulări decodate (în forma PRS - vezi anexa ??), și generează reprezentarea prin cadre pentru fiecare formulare. Specificațiile sistemului de cadre utilizat sunt date prin intermediul unui fișier de control a cărui structură este ilustrată în anexa ?? . Se poate specifica și un fișier pentru controlul normalizării cadrelor.

- **corpus_stats**
Calculează diverse statistici ale corpusurilor de formulări utilizator: numărul de formulări și distribuțiile lor pe clase sau funcție de formele (SRO, NOR, PRS și FRM) disponibile ale fiecăreia. Deasemeni calculează numărul de bigrame din corpus și distribuția etichetelor semantice.
- **corpus_filter**
Acesta este un utilitar cu ajutorul căruia pot fi efectuate diverse operații asupra corpusurilor de formulări utilizator. Se pot specifica prin expresii regulate ce clase de formulări să se păstreze sau să se șteargă din corpus. Deasemeni, pot fi șterse anumite reprezentări (NOR, PRS, FRM) ale formulărilor utilizator dintr-un corpus.
- **corpus_diff**
Compară două corpusuri, identifică și eventual marchează diferențele dintre ele. Este utilizat în etapele de dezvoltare și evaluare a modelului stocastic, când se realizează comparații între decodările manuale și automate ale corpusurilor.
- **dict_build**
Este utilizat pentru construirea dicționarelor de simboluri observate și stări ale unui model Markov, plecând de la corpusul de antrenament. O altă facilitate este posibilitatea de a înregistra evoluția dimensiunii dicționarului pe măsura parcurgerii corpusului, lucru util pentru a observa eventuala ei stabilizare.

Capitolul 4

Experimente și rezultate

Acest capitol va ilustra modul în care implementarea software descrisă în capitolul anterior poate fi utilizată efectiv pentru dezvoltarea unui analizor semantic dedicat pentru un sistem de dialog cu un anumit domeniu de aplicație, precum și performanțele care se pot obține.

După descrierea domeniului și capabilităților sistemului ales ca exemplu, și a colectării și pregătirii datelor folosite pentru dezvoltarea și evaluarea analizorului semantic, sunt prezentate etapele succesive de dezvoltare a acestuia, urmate de evaluarea performanțelor sale.

4.1 Domeniul și capabilitățile sistemului de dialog

Prima etapă în construcția unui sistem de dialog îl constituie analiza domeniului în care acesta va opera și delimitarea capabilităților sale. Unul din rezultatele acestei etape va fi o definiție clară a noțiunilor cu care va lucra analizorul semantic.

Lucrarea de față a urmărit integrarea analizorului într-un sistem care furnizează informații despre orarul Departamentului de Calculatoare. Aceasta alegere este justificată de faptul că domeniul este suficient de complex pentru a permite evaluări realiste ale tehnicilor utilizate în diversele module ale sistemului (recunoașterea vorbirii, analiza semantică, controlul dialogului), și în același timp abordabil în condițiile resurselor limitate disponibile.

Analiza domeniului de aplicație a pornit de la baza de date din care se va face extragerea informațiilor despre orar, care trebuie să conțină practic toate informațiile disponibile în legătură cu orarul (o posibilă structură a unei astfel de baze de date este dată în anexa ??). Au fost astfel identificate principalele tipuri de informații cu care va opera sistemul de dialog (materii, profesori, săli de clasă, ani de studii etc.), și care stabilesc o primă delimitare a domeniului de aplicație.

Pentru a defini încă de la început capabilitățile sistemului, au fost stabilite următoarele caracteristici ale acestuia:

- sistemul furnizează informații despre orar;

- baza de date conține orarul de pe semestrul curent;
- sistemul nu tratează interogări care implică prelucrări suplimentare ale informațiilor din baza de date (exemplu: câte ore am în total luni?);
- sistemul lucrează cu noțiuni precum: materii, tipuri de ore (cursuri, seminarii, proiecte, laboratoare), profesori (nume și grad), săli de clasă, ore, ani de studii, specializări, grupe, subgrupe;
- sistemul nu cunoaște numele studenților, identificarea acestora făcându-se doar la nivel de an, specializare, grupă, subgrupă;
- sistemul acceptă majoritatea specificatorilor de timp: zile ale săptămânii, perioade ale zilei (dimineața, după-amiza, seara), ore exacte, intervale orare, exprimări relative ale zilei (azi, mâine, ieri);
- sistemul lansează subdialoguri de clarificare/obținere de detalii suplimentare în momentele în care acest lucru este necesar.

Câteva dialoguri considerate tipice pentru acest sistem, înregistrate folosind mediul specializat pentru dezvoltarea sistemelor de dialog MDOWZ [?], [?], sunt prezentate în anexa ??.

După definirea domeniului și a capacităților sistemului de dialog în care va fi inclus, a doua etapă în dezvoltarea unui analizor semantic stocastic o constituie colectarea datelor ce vor fi utilizate în acest scop.

4.2 Colectarea și pregătirea datelor

Resursa principală utilizată în dezvoltarea unui analizor semantic stocastic este un corpus de formulări utilizator, tipice sistemului de dialog în care va fi integrat, pe baza căruia se elaborează fișierele de control și se antrenează modelul Markov.

Datele folosite pentru experimentele descrise în continuare au fost colectate folosind mediul specializat de dezvoltare a sistemelor de dialog MDOWZ [?], [?]. Au fost obținute astfel patru corpusuri totalizând 1375 formulări utilizator. Primele trei, notate woz-i, woz-ii și woz-iii, au fost utilizate pentru dezvoltarea analizorului, iar ultimul (woz-test) a fost păstrat pentru evaluarea acestuia.

Cele patru corpusuri au fost colectate în patru etape succesive, în perioada martie–iunie 2000. Pentru a asigura atât o acoperire minimă a domeniului de aplicație, cât și o varietate cât mai mare a formulărilor utilizatorilor, au fost înregistrate dialoguri bazate pe scenarii (tabelul ??) și dialoguri libere. În primul caz, fiecare utilizator a interpretat 2-3 scenarii. În dialogurile libere, utilizatorii au putut formula orice fel de interogări, fiind informați doar asupra domeniului sistemului. De menționat că dacă în cazul scenariilor utilizatorii au fost constrânși să rezolve doar problema dată, în dialogurile libere aceștia au putut formula un număr mai mare de interogări diverse în cadrul aceluiași dialog.

Scenariul 1	Încercați să aflați când aveți un anumit curs/seminar/proiect/laborator
Scenariul 2	Încercați să aflați unde aveți un anumit curs/proiect/seminar/laborator
Scenariul 3	Încercați să aflați cu cine aveți un anumit curs/seminar/proiect/laborator
Scenariul 4	Încercați să aflați ce ore aveți într-o anumită zi
Scenariul 5	Încercați să aflați ce ore aveți într-o anumită perioadă a unei zile
Scenariul 6	Încercați să aflați unde are ore o anumită grupă dintr-un an
Scenariul 7	Încercați să aflați ce ore are o anumită grupă dintr-un an într-un anumit interval
Scenariul 8	Încercați să aflați ce ore se țin într-o anumită sală într-un anumit interval

Tabelul 4.1: Scenarii utilizate pentru colectarea datelor

Au fost astfel înregistrate dialoguri purtate de 42 de subiecți, majoritatea studenți și cadre didactice de ambele sexe din cadrul Facultății de Automatică și Calculatoare. Vârsta acestora a avut o medie de 23 ani, iar majoritatea lor nu utilizaseră anterior un sistem de dialog pentru obținerea de informații, putând fi deci considerați un grup potrivit pentru obținerea unor date pe baza cărora să se poată realiza dezvoltarea și evaluarea analizorului.

Dialogurile înregistrate au fost analizate folosind diverse criterii. Pe de o parte, dialogurile au fost grupate în dialoguri finalizate și întrerupte (prin aceasta înțelegând terminarea prematură a dialogului, cu diverse cauze: renunțarea utilizatorului, instabilitatea MDWOZ etc.) Pe de altă parte, dialogurile au fost apreciate ca reușite sau eșuate – un dialog eșuat fiind unul în care sistemul nu a reușit să furnizeze răspunsurile corecte, deși întrebările utilizatorului respectau domeniul și capacitățile sistemului, iar în baza de date existau informațiile cerute. Rezultatele acestor clasificări și alte statistici referitoare la dialogurile din care au fost extrase datele folosite pentru dezvoltarea analizorului sunt cuprinse în tabelul ??.

Pentru ca dezvoltarea analizorului semantic (inclusiv antrenarea modelului Markov) să nu fie afectată de erorile care pot apare într-un sistem de dialog la ieșirea modului de recunoaștere automată a vorbirii, mediul MDWOZ permite înregistrarea formulărilor utilizatorilor sub forma unor fișiere audio care pot fi ulterior transcrise manual, evitându-se astfel erorile de recunoaștere menționate.

Următorul pas a constat deci în transcrierea manuală a formulărilor utilizatorilor înregistrate în fișiere audio într-un format (SRO – prescurtare de la Speech Recognizer Output) care ar fi putut fi generat de un modul ideal de recunoaștere automată a vorbirii. Au rezultat astfel 3 fișiere corpus de tip text (mai multe

detalii despre structura lor sunt date în anexa ??).

Corpus	woz-i	woz-ii	woz-iii	TOTAL
Dimensiuni corpus				
Număr de dialoguri	59	62	61	182
Număr de formulări utilizator	371	377	340	1088
Distribuția dialogurilor după sexul utilizatorilor				
Bărbați	47	52	55	154 (85%)
Femei	12	10	6	28 (15%)
Distribuția dialogurilor după ocupația utilizatorilor				
Studenti	40	54	52	146 (80%)
Profesori	16	8	5	29 (16%)
Alții	3	0	4	7 (4%)
Distribuția dialogurilor pe tipuri (scenarii/libere)				
Scenariul 1	8	6	2	16
Scenariul 2	3	8	5	16
Scenariul 3	5	13	0	18
Scenariul 4	6	3	7	16
Scenariul 5	11	1	4	16
Scenariul 6	6	2	8	16
Scenariul 7	2	8	6	16
Scenariul 8	0	7	9	16
Total scenarii	41	48	41	130 (71%)
Dialoguri libere	18	14	20	52 (29%)
Dialoguri finalizate/întrerupte				
Finalizate	56	60	58	174 (96%)
Întrerupte	3	2	3	8 (4%)
Dialoguri reușite/eșuate				
Reușite	44	46	55	145 (80%)
Eșuate	15	16	6	37 (20%)
Distribuția dialogurilor reușite/eșuate după tip				
Scenarii reușite	32	34	36	102 (78%)
Scenarii eșuate	9	14	5	28 (22%)
Libere reușite	12	12	19	43 (83%)
Libere eșuate	6	2	1	9 (17%)

Tabelul 4.2: Statistici referitoare la datele de antrenament

4.3 Dezvoltarea analizorului semantic

Dezvoltarea analizorului semantic pentru un sistem de dialog având un anumit domeniu de aplicație urmărește, pe de o parte, construirea fișierelor de configurare care controlează preprocesarea și construcția cadrelor (secțiunile ?? și ?? și anexa ??), pe de alta, definirea și antrenarea modelului Markov (secțiunea ??).

Pentru a putea însa antrena modelul pe baza unui corpus, este necesar ca acesta să fie deja etichetat semantic (secțiunea ??), ceea ce impune realizarea manuală cel puțin a decodării semantice (presupunând că preprocesarea poate fi făcută automat pe baza unor fișiere de configurare generate manual). În continuare vom prezenta metoda utilizată pentru a trata această problemă.

4.3.1 Metoda incrementală (bootstrap)

Etichetarea semantică manuală a celor 1088 formulări utilizator din corpusul de antrenament ar fi fost o activitate de lungă durată, monotonă și predispusă la erori (și costisitoare, dacă ar fi fost plătită), astfel că pentru a atenua aceste probleme am folosit o abordare *incrementală* (așa-numita metodă *bootstrap*).

Ideea principală a metodei este de a construi versiuni succesive ale analizorului (fișiere de configurare și model Markov) pe baza unor porțiuni din ce în ce mai mari din corpusul de antrenament, și de a le utiliza apoi pentru a eticheta automat noi porțiuni din acesta. Astfel, în etapa de *inițializare*, se generează prima versiune a fișierelor de configurare și se etichetează semantic manual prima porțiune din corpusul de antrenament, pe baza căreia se crează un prim model. Folosind fișierele de configurare și modelul astfel obținute, se analizează semantic o nouă porțiune din corpus (de această dată automat), după care se realizează corecția manuală a etichetării respectivei porțiuni, precum și eventuale actualizări ale fișierelor de configurare. Cu porțiunea deja etichetată din corpus (prima parte manual, a doua automat și corectată manual), se construiește o nouă versiune a modelului, utilizată în continuare la etichetarea automată a unei noi porțiuni din corpusul de antrenament șamd. Procesul continuă astfel în *iterații* până la etichetarea întregului corpus, pe baza căruia se construiește versiunea finală a analizorului.

Avantajele metodei incrementale sunt multiple: în primul rând, ea elimină în mare măsură munca de etichetare manuală a unui număr foarte mare de formulări și reduce problema la etichetarea manuală a unei prime porțiuni din corpus și doar corectarea manuală a etichetărilor restului formulărilor, ceea ce reduce timpul necesar dezvoltării modelului Markov. În plus, această metodă permite evaluări intermediare ale performanțelor modelului și ale întregului analizor, prin care se pot urmări îmbunătățirile obținute pe baza corpusului de antrenament. Metoda permite deci și detectarea momentului în care performanțele se stabilizează, moment din care este puțin probabil ca utilizarea unor noi formulări să le îmbunătățească spectaculos.

4.3.2 Erori și evaluarea performanțelor

Înainte de a trece la prezentarea în detaliu a etapelor dezvoltării analizorului, sunt necesare câteva precizări cu privire la tipurile de erori de care acesta poate fi afectat și evaluarea performanțelor sale.

În primul rând, trebuie menționat că în corpusurile de antrenament pot exista și formulări invalide, în sensul că acestea fie sunt incorect formulate semantic, fie depășesc domeniul prestabilit al sistemului de dialog propus (pentru câteva exemple, vezi anexa ??). Acest gen de formulări nu sunt utilizate nici în fazele de antrenare (ar influența negativ modelul creat) și nici în realizarea evaluărilor.

Lăsând la o parte aceste formulări invalide, în cadrul analizei semantice pot să apară erori la nivelul fiecăruia dintre modulele componente (preprocesor, decodor, constructor de cadre). De cele mai multe ori, aceste erori se vor propaga prin lanțul de analiză (e.g. o eroare de preprocesare va genera în continuare o eroare de decodare și apoi de construcție a cadrelor), însă pot exista și situații în care o eroare apărută la nivelul unui modul este “corectată” de unul din modulele următoare. Cel mai des acest comportament apare la nivelul constructorului de cadre: întrucât corespondența dintre succesiunile de etichete semantice și cadrele corespunzătoare nu este biunivocă, aceeași reprezentare poate fi obținută pornind de la succesiuni de etichete semantice diferite, de unde posibilitatea “corecției” unor erori de etichetare semantică.

În evaluările efectuate în continuare vor fi prezentate atât erorile care apar la nivelul fiecărui modul (evaluare glass-box), ele fiind utile pentru a aprecia robustețea și performanța fiecărei componente, cât și cele globale (evaluare black-box), ele determinând în final performanțele analizorului dezvoltat.

4.3.3 Etape și evaluări intermediare

Dezvoltarea analizorului s-a realizat în trei etape succesive, utilizând corpusurile woz-i, woz-ii și woz-iii transcrise în format SRO. Rezultatul fiecărei etape este o versiune a analizorului semantic definită prin fișierele de configurare (vezi anexa ??) și modelul Markov corespunzătoare.

Cele trei etape au fost:

- **Inițializarea**

- Crearea versiunii I pe baza corpusului woz-i
- Evaluarea pe corpusul de antrenament (woz-i)
- Evaluarea pe un nou corpus (woz-ii)

- **Iterația 1**

- Crearea versiunii II pe baza corpusurilor woz-i și woz-ii
- Evaluarea pe corpusul de antrenament (woz-i și woz-ii)
- Evaluarea pe un nou corpus (woz-iii)

- **Iterația 2**

- Crearea versiunii III (finală) pe baza woz-i, woz-ii și woz-iii
- Evaluarea pe corpusul de antrenament (woz-i, woz-ii și woz-iii)

În continuare vom detalia fiecare din aceste etape și vom prezenta rezultatele evaluărilor intermediare.

Inițializarea

- **Crearea versiunii I pe baza corpusului woz-i**

- S-a realizat preprocesarea manuală a corpusului woz-i. Simultan, ținând cont și de analiza preliminară a domeniului, au fost elaborate fișierele de control al preprocesării – fișierele pentru normalizarea numerelor, reducerea flexiunilor, unificarea expresiilor, înlocuirea aliasurilor și unificarea categoriilor (vezi anexa ??). În ceea ce privește unificarea categoriilor, în această fază au fost identificate 8 categorii, din care 3 generice (ZI-RELATIV, PERIOADA-ZI, NR) și 5 corespunzătoare câmpurilor din baza de date (MATERIA, ZI-SAPT, GRUPA, PROFESOR, SPECIALITATEA). Pe baza formei preprocesate a corpusului woz-i, s-a construit și dicționarul formei normalizate (dicționarul de simboluri observate ale modelului Markov).
- S-a realizat etichetarea semantică manuală a corpusului woz-i. Pe parcursul acestei etichetări, au fost marcate formulările fără concept (NC). Deasemeni, au fost identificate 19 (5.1%) formulări invalide (NEG), care au fost eliminate din corpusul de antrenament, dar au fost totuși păstrate separat pentru o eventuală analiză a extinderii domeniului și capabilităților sistemului (vezi anexa ??). Tot în baza etichetării semantice manuale a fost construită mulțimea de etichete semantice utilizate.
- Pe baza algoritmului de estimare de probabilitate maximă, s-a antrenat un prim model Markov (h-i) folosind corpusul etichetat semantic manual woz-i (352 formulări). Modelul rezultat are 40 de stări și 46 de simboluri observate (dimensiuni dictate practic de numărul de etichete semantice și dimensiunile dicționarului normalizat).
- Pe baza analizei preliminare a domeniului și a formulărilor apărute în primul corpus de antrenament, a fost conceput un sistem de cadre pentru reprezentarea semantică a formulărilor utilizator, format în primă instanță din 13 cadre: <da>, <nu>, <identificare>, <specif-timp>, <specif-materie>, <ce-materie>, <ce-curs>, <ce-laborator>, <ce-seminar>, <ce-proiect>, <cand>, <unde> și <cu-cine> (pentru forma finală a sistemului de cadre utilizat, vezi anexa ??).

La sfârșitul acestei etape de inițializare s-a obținut o primă versiune (I) de analizor semantic, definită de fișierele de control al preprocesării, modelul Markov h-i, și fișierul de descriere a sistemului de cadre.

- **Evaluarea versiunii I pe corpusul de antrenament (woz-i)**

În continuare, a fost realizată o evaluare a performanțelor analizorului pe setul de formulări de antrenament (woz-i). În mod evident, acest gen de evaluare nu va identifica nici erori de preprocesare și nici de construcție a cadrelor, deoarece fișierele de control respective au fost elaborate pornind chiar de la acest corpus. Pe de altă parte, datorită naturii stocastice a modelului Markov, pot apărea erori de decodare semantică, deci putem realiza o primă evaluare a modelului Markov și pe baza setului de antrenament.

Calitatea preprocesării poate fi deasemeni evaluată prin măsura în care se realizează o stabilizare a dicționarului normalizat. Așa cum am precizat în secțiunea ??, unul din rolurile preprocesorului este de a realiza o reducere a dicționarului domeniului, fără a pierde informații importante. În ultima etapă a preprocesării toate cuvintele care nu fac parte din dicționarul normalizat sunt eliminate din formulare. În concluzie, dacă această etapă nu este atent controlată prin fișierul de configurare aferent, cuvinte semnificative pot fi eliminate din formulare în mod eronat. Printr-un studiu al evoluției dicționarului normalizat pe măsura parcurgerii corpusului de antrenament, se poate estima calitatea preprocesării efectuate.

În consecință:

- S-a realizat o statistică a reducerii dicționarului domeniului prin cele 7 etape succesive de preprocesare, ilustrată prin graficul din figura ??. După cum se observă pe acest grafic, cele mai pronunțate reduceri ale dicționarului sunt realizate în ultimele două etape, de unificare a categoriilor și de eliminare a cuvintelor necuprinse în dicționar. Deși se observă o tendință de stabilizare a dimensiunilor dicționarului normalizat (curba cea mai de jos), nu putem vorbi încă despre o stabilizare completă a acestuia.
- S-a realizat etichetarea semantică automată a corpusului woz-i folosind versiunea I a analizorului. Verificarea etichetării semantice a evidențiat 14 erori de decodare, dar după construcția cadrelor s-a constatat o reducere a numărului acestora la numai 6 (vezi tabelul ??).

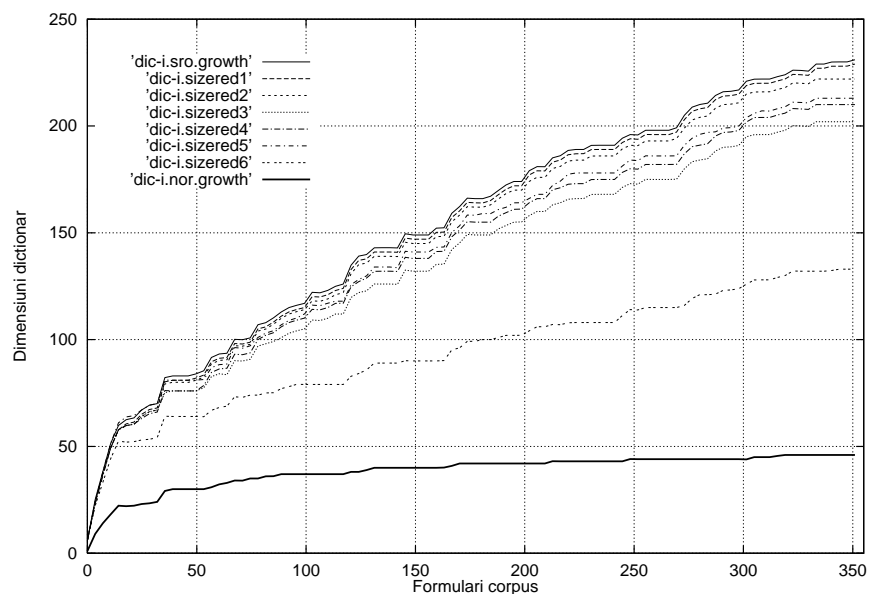


Figura 4.1: Reducerea dicționarului prin preprocesare – versiunea I

EVALUARE PE CORPUS DE ANTRENAMENT (woz-i)	
Număr formulări de test	352
Erori de decodare	14 4.0 %
Erori globale (black-box)	6 1.7 %
EVALUARE PE CORPUS NOU (woz-ii)	
Număr total de formulări	377
Formulări invalide (NEG)	8 2.1 %
Număr formulări de test	369
Erori de preprocesare	35 9.5 %
Erori de decodare	44 11.9 %
Erori de construcție a cadrelor	1 0.3 %
Erori globale (black-box)	80 21.7 %

Tabelul 4.3: Evaluări performanțe – versiunea I

- **Evaluarea versiunii I pe un corpus nou (woz-ii)**

- S-a realizat analiza semantică (preprocesarea, decodarea și construcția cadrelor) a corpusului woz-ii folosind versiunea I a analizorului.
- Rezultatele au fost corectate și totodată s-a realizat evaluarea performanțelor modelului. În acest scop, prin inspecția manuală a rezultatelor analizei semantice au fost identificate și marcate corespunzător erorile de preprocesare, erorile de decodare, erorile de construcție a cadrelor și formulările invalide (NEG). Evident, în majoritatea cazurilor are loc o propagare a erorilor, adică erorile de preprocesare vor genera în continuare erori de decodare și de construcție a cadrelor. Per total, rata erorii obținută în această evaluare a fost de 21.7 %, statisticile concludente în acest sens fiind cuprinse în partea a II-a a tabelului ??.

Iterația 1

- **Crearea versiunii II pe baza corpusurilor woz-i și woz-ii**

- S-a realizat analiza erorilor de preprocesare apărute la analiza semantică a corpusului woz-ii (vezi subsecțiunea anterioară). Pe baza acestei analize au fost îmbunătățite fișierele de control al preprocesării pentru a elimina respectivele erori. În esență au fost adăugate noi expresii, aliasuri, și chiar o nouă categorie – SALA – corespunzătoare câmpului respectiv din baza de date. Întrucât s-a observat că numeroase erori au apărut datorită incompletitudinii fișierelor de configurare, a fost realizat un studiu detaliat al acestora, și au fost aduse și alte modificări pentru completitudine și omogenitate. În continuare, la dicționarul normalizat au fost adăugate câteva noi cuvinte apărute pe parcursul preprocesării celui de-al doilea corpus.
- S-a realizat etichetarea semantică a variantei deja preprocesate a corpusului woz-ii utilizând versiunea I a analizorului. Erorile de decodare semantică au fost corectate și, pe baza celor 721 formulări etichetate semantic din primele două corpusuri (woz-i și woz-ii) a fost antrenată prin algoritmul de estimare de probabilitate maximă o a doua variantă a modelului Markov – h-ii (46 stări și 56 simboluri observate).
- Eroarea de construcție a cadrelor apărută la evaluarea versiunii I a condus la o reevaluare a sistemului de cadre utilizat, modificat pentru a trata și întrebări despre existența unor anumite cursuri/laboratoare etc. (e.g.: E laborator de SDA astăzi?).

Noile variante ale fișierelor de configurare a preprocesării, modelului Markov și descrierii sistemului de cadre utilizat definesc o nouă versiune de analizor

semantic – versiunea II, care a fost evaluată în continuare atât pe corpusul său de antrenament, cât și pe un nou corpus de test.

• **Evaluarea versiunii II pe corpusul de antrenament (woz-i și woz-ii)**

- Calitatea preprocesării a fost din nou analizată prin curba de evoluție a dimensiunilor dicționarului normalizat, ilustrată în figura ???. În urma modificărilor aduse fișierelor de control al preprocesării și a analizei corpusului woz-ii, dimensiunea vocabularului normalizat a crescut de la 46 la 56 de cuvinte, deci încă nu putem vorbi despre o stabilizare a acestuia.

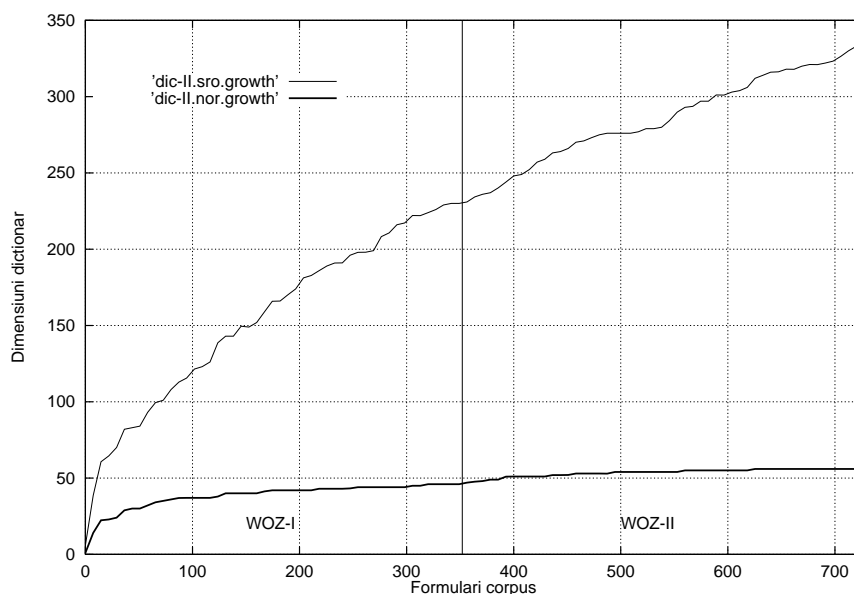


Figura 4.2: Reducerea dicționarului prin preprocesare – versiunea II

- S-a realizat analiza semantică automată a corpusului woz-ii folosind versiunea II a analizorului. Verificarea manuală a acesteia a pus în evidență 22 de erori de etichetare semantică, rezultând în 14 erori de reprezentare pe bază de cadre. O statistică mai detaliată este prezentată în prima parte a tabelului ???.

• **Evaluarea versiunii II pe un corpus nou (woz-iii)**

- S-a realizat analiza semantică (preprocesarea, decodarea și construcția cadrelor) a corpusului woz-iii folosind versiunea II a analizorului.
- După corectarea rezultatelor s-a realizat evaluarea performanțelor. Este de remarcat că de această dată nu au mai apărut erori de preprocesare

EVALUARE PE CORPUS DE ANTRENAMENT (woz-i+ii)	
Număr formulări de test	721
Erori de decodare	22 3.1 %
Erori globale (black-box)	14 1.9 %
EVALUARE PE CORPUS NOU (woz-iii)	
Număr total de formulări	340
Formulări invalide (NEG)	10 2.9 %
Număr formulări de test	330
Erori de preprocesare	0 0 %
Erori de decodare	46 13.9 %
Erori de construcție a cadrelor	0 0 %
Erori globale (black-box)	45 13.6 %

Tabelul 4.4: Evaluări performanțe – versiunea II

și nici de construcție a cadrelor, lucru care constituie un prim indiciu în direcția stabilizării fișierelor de control și a sistemului de cadre utilizate. Au fost identificate 46 erori de decodare, din care 45 s-au propagat și după construcția cadrelor. Per total, rata erorii versiunii II a fost de 13.6%, semnificativ mai redusă decât a versiunii I (rezultatele sunt prezentate în detaliu în partea a doua a tabelului ??).

Iterația 2

- **Crearea versiunii III pe baza corpusurilor woz-i, woz-ii și woz-iii**
 - S-a preprocesat corpusul woz-iii folosind versiunea II a analizorului. Întrucât nu au apărut erori (vezi subsecțiunea anterioară), nu a fost necesară intervenția manuală și nici efectuarea de modificări în fișierele de control respective.
 - S-a realizat etichetarea semantică a variantei preprocesate a corpusului woz-iii utilizând versiunea II. Erorile de decodare semantică au fost corectate și, pe baza celor 1051 formulări etichetate semantic din primele trei corpusuri (woz-i, woz-ii și woz-iii), au fost antrenate două variante (finale) de model Markov. Prima (h-iii.mle) a fost antrenată utilizând ca și până acum estimarea de probabilitate maximă. La a doua variantă

(h-iii.katz) a fost aplicată reestimarea Katz atât pentru probabilitățile tranzițiilor între stări ($k = 4$ – au fost discreditate bigramele care au apărut de maximum 4 ori), cât și pentru probabilitățile stărilor inițiale ($k = 2$).

- Întrucât nu au existat erori de construcție a cadrelor la analiza semantică a corpusului woz-iii (vezi secțiunea anterioară), fișierul de descriere al sistemului de cadre a rămas nemodificat.

În urma acestei etape au fost construite versiunile finale ale analizorului semantic stocastic proiectat: versiunea III-MLE (cu modelul Markov construit prin estimare de probabilitate maximă) și versiunea III-Katz (cu modelul Markov construit aplicând și reestimarea Katz).

• **Evaluarea versiunii III pe corpusul de antrenament (woz-i+ii+iii)**

- Calitatea preprocesării a fost din nou analizată prin intermediul curbei de evoluție a dimensiunilor dicționarului normalizat, ilustrată în figura ???. De această dată, la parcurgerea ultimului corpus de antrenament (woz-iii) nu au apărut erori de preprocesare și nici cuvinte noi în dicționarul normalizat, acesta rămânând deci la dimensiunea de 56 cuvinte. În concluzie, putem spera că din acest moment avem o stabilizare a dicționarului simbolurilor observate.

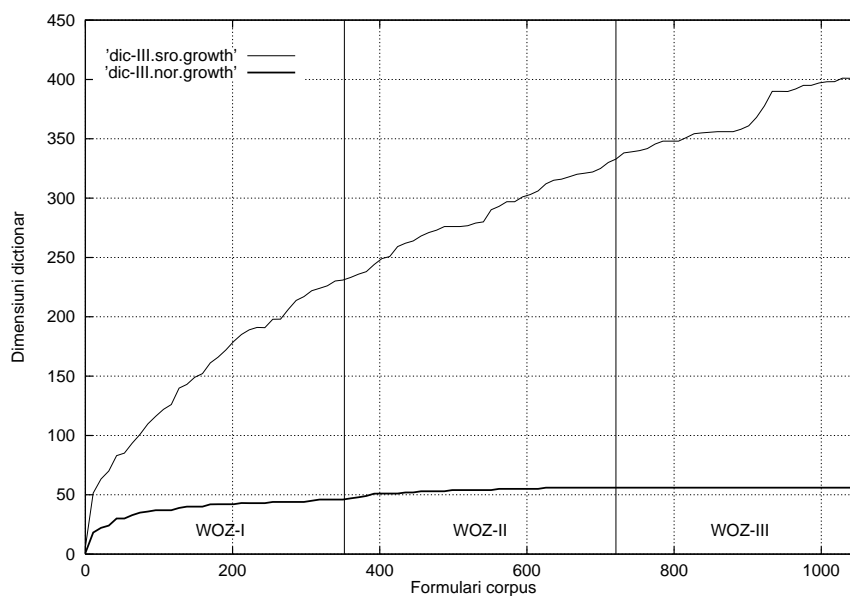


Figura 4.3: Reducerea dicționarului prin preprocesare – versiunea III

- S-a realizat analiza semantică automată a primelor 3 corpusuri folosind versiunile III-MLE și III-Katz ale analizorului. Din statistica din tabelul

EVALUARE PE CORPUS DE ANTRENAMENT (woz-i+ii+iii)		
Număr formulări de test		1051
MLE	Erori de decodare	31 2.9 %
	Erori globale (black-box)	20 1.9 %
Katz	Erori de decodare	38 3.6 %
	Erori globale (black-box)	26 2.5 %

Tabelul 4.5: Performanțele versiunilor finale (III) pe corpusul de antrenament

??, rezultă că modelul reestimat prin metoda Katz a produs sensibil mai multe erori de etichetare semantică decât modelul MLE. Acest lucru se explică prin faptul că testul a fost făcut pe corpusul de antrenament. Întrucât scopul reestimării Katz este de a modela mai bine evenimentele care nu apar în corpusul de antrenament, este clar că reestimarea nu va aduce îmbunătățiri ale performanței în cazul unui astfel de test. Mai mult, datorită faptului că metoda Katz discreditează anumite evenimente ce au apărut în corpusul de antrenament, performanța modelului reestimat poate fi chiar mai slabă decât a celui inițial la o testare pe formulările din corpusul de antrenament.

Această etapă încheie dezvoltarea incrementală a analizorului semantic, urmând ca în continuare să prezentăm evaluările performanțelor versiunilor finale.

4.4 Evaluarea finală a analizorului

În scopul evaluării comparative a performanțelor celor două versiuni de analizor dezvoltate (III-MLE și III-Katz), a fost utilizat cel de-al patrulea corpus de formulări utilizator (woz-test). Câteva statistici privind acest corpus de test sunt prezentate în tabelul ???. El totalizează 287 de formulări, însă înainte de a trece la realizarea analizei lui semantice au fost identificate și eliminate 4 formulări invalide (NEG), testarea urmând a se efectua pe 283 de formulări.

După efectuarea analizei semantice a corpusului woz-test folosind cele două versiuni de analizor, a fost realizată inspecția manuală a rezultatelor și marcarea diverselor tipuri de erori apărute (o statistică a lor este prezentată în tabelul ??).

În testele efectuate nu a mai apărut nici o eroare la nivelul modulelor de preprocesare și de construcție a cadrelor. La nivelul decodării semantice, modelul Markov obținut după aplicarea reestimării Katz s-a comportat semnificativ mai bine decât cel antrenat prin estimare de probabilitate maximă, generând numai 19 erori, față de cele 26 generate de al doilea. Construcția cadrelor a eliminat însă

Corpus	woz-test
Dialoguri	45
Dialoguri după scenarii	23 (51%)
Dialoguri libere	22 (49%)
Număr total de formulări	287
Formulări invalide (NEG)	4 (1.4%)
Număr formulări de test	283

Tabelul 4.6: Statistici referitoare la datele de test

Versiune	III-MLE	III-Katz
Număr formulări de test	283	283
Erori de preprocesare	0	0
Erori de decodare (procentaj)	26 9.18%	19 6.71%
Erori de decodare corectate la construcția cadrelor	3	1
Erori de construcție a cadrelor	0	0
Total erori (black-box) (procentaj)	23 8.12%	18 6.36%

Tabelul 4.7: Rezultatele evaluării finale a performanțelor

câteva erori de decodare, astfel încât per total rata erorii versiunii III-Katz a fost de 6.36%, față de 8.12% pentru versiunea III-MLE.

În concluzie, aplicând reestimarea Katz, numărul de erori de decodare a fost redus de la 26 la 19, ceea ce reprezintă o îmbunătățire relativă cu cca. 27% a ratei erorilor de decodare.

Rata globală a erorilor de 6.36% clasifică analizorul semantic dezvoltat pe domeniul propus ca unul de performanță ridicată (reamintim că performanțele cele mai bune menționate în literatură sunt de 3.8%) și confirmă oportunitatea abordării stocastice în dezvoltarea lui.

Observație: Întrucât nu s-a observat încă o stabilizare a performanțelor (rata erorilor a scăzut de la 21.7% la 13.6%, 8.12%, și în cele din urmă la 6.36%), se poate încerca o îmbunătățire în continuare a acestora prin utilizarea unor noi seturi de formulări. În acest sens se poate apela într-o primă instanță chiar la ultimul corpus utilizat în evaluarea analizorului (woz-test).

Capitolul 5

Concluzii și continuări

În capitolele anterioare au fost descrise în detaliu proiectarea, implementarea software, dezvoltarea și evaluarea performanțelor unui analizor semantic stocastic pentru sisteme automate de dialog vocal om-calculator.

Acest ultim capitol este dedicat prezentării concluziilor desprinse din această experiență și identificării câtorva posibile direcții de continuare a cercetărilor.

5.1 Concluzii

Această lucrare a prezentat dezvoltarea unui modul de analiză semantică pentru sisteme de dialog. Etapele parcurse în acest sens au fost:

- realizarea unui studiu al diverselor alternative de proiectare, pe baza căruia a fost luată decizia utilizării formalismului case-grammar pentru reprezentarea cunoștințelor și a unei abordări stocastice bazate pe modele Markov pentru realizarea parsingului;
- proiectarea arhitecturii analizorului semantic și a funcționalității fiecăruia dintre modulele componente;
- implementarea software efectivă a analizorului proiectat;
- dezvoltarea resurselor (fișiere de configurare și modele Markov) necesare pentru integrarea analizorului într-un sistem de dialog pentru furnizarea de informații despre orarul Departamentului de Calculatoare;
- evaluarea performanțelor analizorului semantic construit.

Pe baza etapelor sintetizate mai sus se pot desprinde următoarele concluzii:

- Deși nu se pot face comparații directe (datorită diferențelor existente între domeniile de dialog și datele utilizate), rata erorilor de 6.36%, obținută la evaluarea finală a analizorului, este totuși apropiată de cele mai bune (3.8%) prezentate anterior în literatură.

- Abordarea stocastică aleasă a permis dezvoltarea rapidă și eficientă a unui analizor semantic pentru sisteme automate de dialog vocal om-calculator.
- Utilizarea unor algoritmi specifici de îmbunătățire a modelului stocastic a permis construirea unui analizor cu o robustețe sporită în fața diverselor fenomene specifice vorbirii spontane, concretizată într-o reducere relativă a ratei erorilor de circa 27%.
- Analizorul semantic proiectat în această lucrare este unul general, independent atât de domeniul cât și de limba sistemului de dialog în care urmează a fi integrat. Adaptarea analizorului la un anumit sistem de dialog (ilustrată în capitolul 4) presupune un efort mult mai redus decât elaborarea unui eventual set de reguli pentru modelarea analizei semantice. La schimbarea sau extinderea domeniului de dialog trebuie doar înregistrat un nou corpus de antrenament, caracteristic noului sistem, și repetate etapele de dezvoltare prin care se construiesc modelul Markov și setul de fișiere de configurare a preprocesării și construcției cadrelor. Analizorul proiectat este deci caracterizat de un grad ridicat de portabilitate, oferind o soluție generală pentru problema analizei semantice în sistemele de dialog.

Toate acestea arată că abordarea stocastică, deja dominantă în recunoașterea automată a vorbirii, este pe deplin justificată și în analiza semantică a limbajului natural pentru sisteme automate de dialog vocal om-calculator, pentru care oferă o soluție performantă, robustă, și cu un grad ridicat de portabilitate.

5.2 Direcții de continuare

Continuările cercetărilor prezentate în această lucrare se pot canaliza pe două direcții principale: pe de o parte, îmbunătățiri ale implementării și performanțelor analizorului, pe de alta, integrarea acestuia într-un sistem de dialog complet.

În privința îmbunătățirii implementării și performanțelor:

- La nivelul implementării pot fi îmbunătățite unele dintre soluțiile alese, în special în ceea ce privește modulul de preprocesare. Este vorba în primul rând despre normalizarea numerelor, care în faza curentă se realizează pe baza unor simple înlocuiri de șiruri de caractere, soluție acceptabilă în sistemul de dialog proiectat, dar insuficientă într-un caz general. Aceasta presupune scrierea unei gramatici a numerelor și realizarea normalizării lor pe baza ei. Etapa de reducere a flexiunilor poate fi deasemeni realizată într-un cadru general pe baza unui dicționar morfo-sintactic conținând atât formele flexionare ale cuvintelor, cât și lemele respectivelor forme [?]. Deasemeni, se pot realiza optimizări ale structurilor de date și algoritmilor din biblioteca de clase (e.g. implementare cu hash-table a dicționarelor etc.)

- O altă extensie utilă ar fi dezvoltarea unui mecanism pentru detectarea și semnalarea erorilor apărute pe parcursul analizei semantice. În momentul de față analizorul nu dispune de un astfel de mecanism, cadrele construite fiind transmise către modulul de control al dialogului chiar dacă reprezentarea obținută nu este cea corectă (din punct de vedere al semanticii limbajului natural). Realizarea unui astfel de mecanism este însă problematică și presupune o analiză în contextul modulelor adiacente de recunoaștere a vorbirii și control al dialogului, constituind în sine un domeniu interesant de studiu pe viitor.
- În cadrul domeniului de dialog propus, se poate încerca îmbunătățirea performanței analizorului prin continuarea antrenării modelului Markov pe baza unor noi seturi de formulări. Întrucât pe parcursul etapelor prezentate nu s-a observat încă o stabilizare a performanței, este posibil să obținem astfel o nouă creștere a ei.

O a doua direcție în care se pot focaliza eforturile de cercetare viitoare o constituie integrarea analizorului semantic dezvoltat într-un sistem de dialog complet pentru furnizarea de informații despre orar. În acest sens, menționăm ca lucrarea de față este parte integrantă a unui efort pe scară mai largă de dezvoltare a unor sisteme automate de dialog vocal om-calculator [?], [?], [?].

Bibliografie

- [1] J. Allen. *Natural Language Understanding*. Benjamin/Cummings Publishing Company, 1995.
- [2] H. Aust și M. Oerder. The Philips automatic train timetable information system. *Speech Communication*, 17:249–62, 1995.
- [3] S.K. Bennacef, H. Bonneau-Maynard, J.L. Gauvain, L.F. Lamel, și W. Minker. A Spoken Language System for Information Retrieval. În *Proceedings IC-SLP'94*, pag. 1271–74, 1994.
- [4] D. Bohuş și M. Boldea. A Web-based Text Corpora Development System. În *Proceedings Second International Conference on Language Resources and Evaluation*, Atena, Grecia, 2000.
- [5] M. Boldea. SADVOC: Sisteme Automate de Dialog Vocal Om-Calculator. Propunere de proiect major de cercetare înaintată CNCSIS, noiembrie 1999.
- [6] I. Coteanu, L. Seche, și M. Seche. *Dicționarul explicativ al limbii române*. Editura Academiei, 1975.
- [7] J.L. Gauvain, S. Bennacef, L. Devillers, L.F. Lamel, și S. Rosset. The Spoken Language Component of the Mask Kiosk. În K. Varghese și S. Pflieger, *Human Comfort and Security of Information Systems*, pag. 93–103. Springer, 1997.
- [8] S.M. Katz. Estimation of Probabilities from Sparse Data for the Language Model Component of a Speech Recognizer. *IEEE Transactions on Acoustics, Speech and Signal Processing*, 1987.
- [9] KDevelop. www.kdevelop.org.
- [10] W. Minker, A. Waibel, și J. Mariani. *Stochastically-Based Semantic Analysis*. Kluwer Academic Publishers, 1999.
- [11] R. De Mori, editor. *Spoken Dialogues with Computers*. Academic Press, 1998.
- [12] C. Munteanu. Mediu pentru dezvoltarea sistemelor de dialog prin metoda Vrăjitorului din Oz. Lucrare de diplomă, Universitatea "Politehnica" din Timișoara, iunie 1999.

-
- [13] C. Munteanu. Modelarea și controlul dialogului om-calculator. Dizertație de studii aprofundate, Departamentul de Calculatoare, Universitatea "Politehnica" din Timișoara, iunie 2000.
- [14] C. Munteanu și M. Boldea. MDWOZ: A Wizard of Oz Environment for Dialog Systems Development. În *Proceedings Second International Conference on Language Resources and Evaluation*, Atena, Grecia, 2000.
- [15] L.R. Rabiner și B.H. Juang. An Introduction to Hidden Markov Models. *IEEE Acoustics, Speech and Signal Processing Magazine*, 1986.
- [16] S. Russell și P. Norvig. *Artificial Intelligence: A Modern Approach*. Prentice Hall, 1995.
- [17] N.O. Bernsen și H. Dybkjaer și L. Dybkjaer. *Designing Interactive Speech Systems: From First Ideas to User Testing*. Springer, 1997.
- [18] E. Tzoukermann, R. Pieraccini, și Z. Gorelov. Natural Language Processing in the CHRONUS System. În *Proceedings ICSLP'92*, 1992.
- [19] W. Ward. Understanding Spontaneous Speech: the PHOENIX System. În *Proceedings ICASSP'91*, pag. 365–68, 1991.
- [20] M. Woszczyna et al. Recent Advances in JANUS: A Speech Translation System. În *Proceedings EUROSPEECH'93*, pag. 1295–98, Berlin, 1993.

Anexa A

Exemple de dialoguri

În continuare sunt prezentate câteva dialoguri considerate tipice pentru sistemul de dialog la care se referă experimentele din capitolul 4. Prefixul “S:” indică formulările sistemului, iar “U:” – formulările utilizatorilor.

Dialog 1 (scenariul 1)

S: Buna ziua, ati mai utilizat acest sistem ?

U: Nu

S: Acest sistem furnizeaza informatii despre orar. Va rog formulati intrebarile si raspunsurile dumneavoastra cat mai clar si nu interveniti in timp ce sistemul va raspunde.

Daca sunteti student precizati anul, grupa, subgrupa. Daca sunteti cadru didactic precizati numele. In caz contrar precizati ca nu doriti sa va identificati.

U: Sunt student in anul doi grupa doi subgrupa doi unu.

S: Va rog formulati intrebarea dumneavoastra.

U: As vrea sa aflu cand am laboratorul de SDA.

S: La anul 2 grupa 2 subgrupa 1, laboratorul de ‘‘Structuri de date si algoritmi’’ se tine marti de la 16 la 18 in sala A305. Mai doriti alte informatii ?

U: Nu mersi ok, asta e tot.

S: Va multumesc ca ati folosit acest sistem. La revedere.

Dialog 2 (scenariul 5)

S: Buna ziua, ati mai utilizat acest sistem ?

U: Da

S: Daca sunteti student precizati anul, grupa, subgrupa. Daca sunteti cadru didactic precizati numele. In caz contrar precizati ca nu

doriti sa va identificati.
U: Student anul intai grupa trei subgrupa unu, Cosmin.
S: Va rog formulati intrebarea dumneavoastra.
U: As dori sa aflu ce ore am intre ... vineri intre ora 8 si 12
S: La anul 1, cursul de ‘‘Matematici speciale’’ se tine vineri de la 8
la 10 in sala CA2. Mai doriti alte informatii ?
U: Nu multumesc.
S: Va multumesc ca ati folosit acest sistem. La revedere.

Dialog 3 (liber)

S: Buna ziua, ati mai utilizat acest sistem ?
U: Da.
S: Daca sunteti student precizati anul, grupa, subgrupa. Daca sunteti
cadru didactic precizati numele. In caz contrar precizati ca nu
doriti sa va identificati.
U: Anul intai.
S: Va rog formulati intrebarea dumneavoastra.
U: Cand are grupa doi semigrupa unu seminar de Analiza Matematica 2 ?
S: La anul 1 grupa 2 subgrupa 1, seminarul de ‘‘Analiza Matematica 2’’
se tine miercuri de la 16 la 18 in sala A314. Mai doriti alte
informatii ?
U: Da.
S: Va rog formulati intrebarea dumneavoastra.
U: Cand am timp liber marti, dimineata de la opt pana seara la opt ?
S: Imi pare rau, nu am inteles. Va rog reformulati.
U: Ce ore am marti ?
S: Exista mai multe raspunsuri, sistemul le va preciza pe toate daca nu
specificati criteriile suplimentare pentru clarificarea intrebării. Doriti
sa specificati astfel de criterii ?
U: Da.
S: Precizati grupa si subgrupa.
U: Grupa 2 semigrupa 1.
S: De la 8 la 10 laboratorul de ‘‘Analiza si sinteza disp. numerice’’.
De la 10 la 12 seminarul de ‘‘Matematici speciale’’.
De la 14 la 16 seminarul de ‘‘Sport’’.
Mai doriti alte informatii ?
U: Da.
S: Va rog formulati intrebarea dumneavoastra.
U: Cu cine si unde se tine laboratorul de la 8 la 10 ?
S: La anul 1 grupa 2 subgrupa 1, laboratorul de ‘‘Analiza si sinteza disp.
numerice’’ se tine marti de la 8 la 10 in sala B624 cu sl. Mihaescu.
Mai doriti alte informatii ?
U: Nu.
S: Va multumesc ca ati folosit acest sistem. La revedere.

Dialog 4 (liber)

- S: Buna ziua, ati mai utilizat acest sistem ?
U: Da.
S: Daca sunteti student precizati anul, grupa, subgrupa. Daca sunteti cadru didactic precizati numele. In caz contrar precizati ca nu doriti sa va identificati.
U: Anul 4 grupa 4 subgrupa 1.
S: Va rog formulati intrebarea dumneavoastra.
U: Ce curs este de la ora 11 in sala A109 ?
S: La anul 4, cursul de ‘‘Sisteme de intrare-iesire’’ se tine luni de la 11 la 14 in sala A109.
La anul 4, cursul de ‘‘Rețele locale de calculatoare’’ se tine joi de la 11 la 14 in sala A109.
Mai doriti alte informatii ?
U: Da.
S: Va rog formulati intrebarea dumneavoastra.
U: Cine tine cursurile respective ?
S: La anul 4, cursul de ‘‘Sisteme de intrare-iesire’’ se tine cu prof. Strugaru.
La anul 4, cursul de ‘‘Rețele locale de calculatoare’’ se tine cu prof. Strugaru. Mai doriti alte informatii ?
U: Da.
S: Va rog formulati intrebarea dumneavoastra.
U: Ce laboratoare are in cursul zilei de marti grupa patru ?
S: Exista mai multe raspunsuri, sistemul le va preciza pe toate daca nu specificati criteriile suplimentare pentru clarificarea intrebării. Doriti sa specificati astfel de criterii ?
U: Da.
S: Precizati la ce subgrupa.
U: Patru unu.
S: La anul 4 grupa 4 subgrupa 1, laboratorul de ‘‘Sisteme de intrare-iesire’’ se tine marti de la 8 la 10 in sala B413.
La anul 4 grupa 4 subgrupa 1, laboratorul de ‘‘Ingineria sistemelor industriale’’ se tine marti de la 12 la 14 in sala SPM.
Mai doriti alte informatii ?
U: Da.
S: Va rog formulati intrebarea dumneavoastra.
U: As vrea sa stiu daca Radu Marinescu are laborator de IP marti.
S: Exista mai multe raspunsuri, sistemul le va preciza pe toate daca nu specificati criteriile suplimentare pentru clarificarea intrebării. Doriti sa specificati astfel de criterii ?
U: Da.
S: Precizati grupa.
U: Grupa trei.
S: La anul 4 grupa 3 subgrupa 2, laboratorul de ‘‘Ingineria programării’’ se tine marti de la 16 la 18 in sala B528 cu asistent Radu Marinescu.
La anul 4 grupa 3 subgrupa 1, laboratorul de ‘‘Ingineria programării’’

se tine marti de la 18 la 20 in sala B528 cu asistent Radu Marinescu.

Mai doriti alte informatii ?

U: As vrea sa stiu daca exista totusi laborator de IP de la 12 marti.

S: La anul 4 grupa 4 subgrupa 2, laboratorul de ‘‘Ingineria programarii’’
se tine marti de la 12 la 14 in sala B528. Mai doriti alte informatii ?

U: Nu.

S: Va multumesc ca ati folosit acest sistem. La revedere.

Anexa B

Structura unui corpus

Dezvoltarea fișierelor de configurare și a modelului Markov utilizate în analiza semantică se fac pe baza unui corpus de formulări utilizator tipice pentru sistemul de dialog proiectat. În continuare vom prezenta formatul în care formulările sunt stocate într-un fișier conținând un asemenea corpus.

Fiecare formulare utilizator trece prin mai multe reprezentări de-a lungul lanțului de analiză semantică: text de la modulul de recunoaștere, text normalizat, succesiune de etichete semantice și reprezentare pe bază de cadre. Evident, toate aceste forme trebuie reprezentate în corpusul de antrenament, drept pentru care a fost stabilită următoarea structură pentru o formulare utilizator:

```
%<clasă formulare>
SRO:<forma obținută de la modulul de recunoaștere>
NOR:<forma normalizată, obținută după preprocesare>
PRS:<succesiunea de etichete semantice, obținută după decodare>
FRM:<reprezentare pe bază de cadre>
$
```

Clasa formulării (e.g. NEG – formulare invalidă, NC – formulare fără concept etc.) este folosită pentru operații selective asupra corpusului și poate fi vidă. Formele SRO (Speech Recognizer Output), NOR (NORmalized), PRS (PaRSe) și FRM (FRaMe) deasemeni pot să existe sau nu pentru o anumită formulare.

Pentru exemplificare, prezentăm în continuare o porțiune dintr-un corpus.

```
%
SRO:da am mai utilizat acest sistem
NOR:da
PRS:<da>
FRM:(<da>)
$
```

```

%
SRO:(aa) sunt student în anul doi grupa trei
NOR:sunt student an [NR:"2"] grupă [NR:"3"]
PRS:<identificare> (null) (m:anul) (v:anul) (m:grupa) (v:grupa)
FRM:(<identificare> (anul "2") (grupa "3"))
$

%
SRO:(aa) aş vrea să ştiu ce ore am marţi
NOR:ce-ore [ZI-SAPT:"marţi"]
PRS:<ce-materie> (v:zi-sapt)
FRM:(<ce-materie> (<specif-timp> (zi-sapt "marţi")))
$

%
SRO:în ce zile şi la ce ore îl găsesc pe domnul profesor Jian
NOR:ce-zi şi la ce-ore profesor [PROFESOR:"Ionel Jian"]
PRS:<când> (null) (null) (null) (m:profesor) (v:profesor)
FRM:(<când> (<identificare> (profesor "Ionel Jian")))
$

%NC
SRO:anul cinci grupa management
NOR:<identificare> an [NR:"5"] grupă management
PRS:<identificare> (m:anul) (v:anul) (m:grupa) (v:grupa)
FRM:(<identificare> (anul "5") (grupa "management"))
$

%
SRO:(aa) sunt în anul întâi grupa a doua semigrupa întâi şi aş dori să-mi
    spui când am laborator de electronică doi
NOR:sunt an [NR:"1"] grupă [NR:"2"] subgrupă [NR:"1"] şi când
    laborator [MATERIA:"Electrotehnică II"]
PRS:<identificare> (m:anul) (v:anul) (m:grupa) (v:grupa) (m:subgrupa)
    (v:subgrupa) (null) <cand> (m:laborator) (v:laborator)
FRM:(<identificare> (anul "1") (grupa "2") (subgrupa "1")) (<când>
    (<identificare> (anul "1") (grupa "2") (subgrupa "1")))
    (<specif-materie> (laborator "Electrotehnică II")))
$

%
SRO:în ce zile sunt laboratoare de circuite integrate la anul trei
NOR:ce-zi sunt laborator [MATERIA:"Circuite integrate"] la an [NR:"3"]
PRS:<când> (null) (m:laborator) (v:laborator) (null) (m:anul) (v:anul)
FRM:(<când> (<identificare> (anul "3")) (<specif-materie> (laborator
    "Circuite integrate")))
$

%

```

```
SRO:unde se ține cursul de proiectarea translatoarelor
NOR:unde curs [MATERIA:"Proiectarea translatoarelor"]
PRS:<unde> (m:curs) (v:curs)
FRM:(<unde> (<specif-materie> (curs "Proiectarea translatoarelor")))
$

%
SRO:ce predă profesorul Popa în anul trei
NOR:ce profesor [PROFESOR:"Mircea Popa"] an [NR:"3"]
PRS:<ce-materie> (m:profesor) (v:profesor) (m:anul) (v:anul)
FRM:(<ce-materie> (<identificare> (anul "3") (profesor "Mircea Popa")))
$

%
SRO:în ce sală ține grupa a treia din anul patru calculatoare laboratorul
de inteligență artificială
NOR:ce-sală grupă [NR:"3"] an [NR:"4"] [SPECIALITATEA:"calculatoare"]
laborator [MATERIA:"Inteligență artificială"]
PRS:<unde> (m:grupa) (v:grupa) (m:anul) (v:anul) (v:specialitatea)
(m:laborator) (v:laborator)
FRM:(<unde> (<identificare> (anul "4") (specialitatea "calculatoare")
(grupa "3"))) (<specif-materie> (laborator "Inteligență artificială")))
$

%
SRO:am cursuri în ziua de marți între ora opt și zece
NOR:curs zi [ZI-SAPT:"mari"] între oră [NR:"8"] și [NR:"10"]
PRS:<ce-curs> (m:zi-sapt) (v:zi-sapt) (m:ora-start) (null) (v:ora-start)
(m:ora-stop) (v:ora-stop)
FRM:(<ce-curs> (<specif-timp> (zi-sapt "marți") (ora-start "8")
(ora-stop "10")))
$

%
SRO:aș vrea să știu cu cine se ține laboratorul de programare în limbaj
de asamblare
NOR:cu-cine laborator [MATERIA:"Programarea în limbaj de asamblare"]
PRS:<cu-cine> (m:laborator) (v:laborator)
FRM:(<cu-cine> (<specif-materie> (laborator "Programarea în limbaj de
asamblare")))
$

%
SRO:ce ore am duminică noaptea
NOR:ce-ore [ZI-SAPT:"duminică"] [PERIOADA-ZI:"noapte"]
PRS:<ce-materie> (v:zi-sapt) (v:perioada-zi)
FRM:(<ce-materie> (<specif-timp> (zi-sapt "duminică") (perioada-zi
"noapte")))
$
```

```
%
SRO:sunt student în anul șase (aa) la grupa de soft
NOR:sunt student an [NR:"6"] la grupă
PRS:<identificare> (null) (m:anul) (v:anul) (null) (m:grupa)
FRM:(<identificare> (anul "6"))
$
```

```
%
SRO:mă numesc Pop Traian
NOR:<identificare> [PROFESOR:"Traian Pop"]
PRS:<identificare> (v:profesor)
FRM:(<identificare> (profesor "Traian Pop"))
$
```

Prezentăm și câteva formulări invalide (din clasa NEG), eliminate atât din antrenamentul cât și din testrea analizorului. Ele fie depășesc domeniul sistemului de dialog, fie sunt greșit formulate, nereprezentând intrări valide pentru analizor.

```
%NEG
SRO:ce înseamnă a se de ne
NOR:ce [MATERIA:"Analiza și sinteza dispozitivelor numerice"]
$
```

```
%NEG
SRO:în ce zi am curs dimineața și laborator după masă cu același profesor
    în săli diferite
NOR:ce-zi curs [PERIOADA-ZI:"dimineață"] și laborator [PERIOADA-ZI:
    "după-amiază"] profesor sală
$
```

```
%NEG
SRO:de ce se face pedagogie
NOR:ce [MATERIA:"Pedagogie"]
$
```

```
%NEG
SRO:care este primul laborator din săptămână
NOR:care laborator săptămână
$
```

```
%NEG
SRO:care este următorul laborator
NOR:care laborator
$
```

Anexa C

Structura bazei de date orar

Această anexă ilustrează o posibilă structură a unei baze de date conținând informații despre orarul facultății noastre, stabilită în cursul definirii domeniului sistemului de dialog.

SALI		
I	ID.SALA	NUMERIC
	NUME	STRING
	LOCATIE	STRING

Tabelul C.1: Tabela de săli

MATERII		
I	ID.MATERIE	NUMERIC
	NUME	STRING
	TIP	CHAR
	DESCRIERE	STRING

Tabelul C.2: Tabela de materii

PROFESORI		
I	ID.PROFESOR	NUMERIC
	NUME	STRING
	GRAD	STRING

Tabelul C.3: Tabela de profesori

ORAR		
I	ID_MATERIE	NUMERIC
I	ID_PROFESOR	NUMERIC
I	ID_SALA	NUMERIC
	ORA_START	NUMERIC
	ORA_STOP	NUMERIC
	ZI	NUMERIC
	PARITATE_SAPT	NUMERIC
	AN	NUMERIC
	SPECIALIZARE	CHAR
	GRUPA	NUMERIC
	SUBGRUPA	NUMERIC

Tabelul C.4: Tabela orar

Anexa D

Fișierele de configurare

Comportamentul fiecăruia din cele 3 module ale analizorului semantic proiectat este controlat prin intermediul unor fișiere: modulul de preprocesare se bazează pe o serie de fișiere care controlează modalitatea în care se realizează etapele de preprocesare, decodorul semantic se bazează pe modelul Markov, iar constructorul de cadre – pe un fișier cu descrierea sistemului de cadre utilizat pentru generarea reprezentării semantice, și eventual unul pentru controlul normalizării cadrelor. Structura acestor fișiere de control este descrisă în continuare.

D.1 Controlul preprocesării

Cinci din cele șapte etape de preprocesare (normalizarea numerelor, reducerea flexiunilor, unificarea expresiilor, înlocuirea aliasurilor și unificarea categoriilor) se realizează prin înlocuiri de șiruri de caractere specificate în fișiere de configurare. Acestea sunt simple fișiere text, compuse din intrări de forma:

```
<șir_înlocuitor>: <șir_1>,<șir_2>, ..., <șir_N>;
```

unde `șir_1, ..., șir_N` vor fi înlocuite cu `șir_înlocuitor`.

Dacă șirurile de caractere conțin spații, ele vor fi delimitate prin ghilimele. În fișiere sunt permise și comentarii, acestea fiind introduse prin caracterul `#` și terminându-se la sfârșitul liniei curente.

În continuare sunt prezentate porțiuni din cele 5 fișiere utilizate pentru controlul etapelor de preprocesare menționate în varianta finală a analizorului.

numbers.pp - normalizarea numerelor

```
# Fișier control preprocesare - normalizarea numerelor
```

```
1: unu, întâi, "a întâia";
```

```
2: doi, două, "al doilea", "a doua";
3: trei, "al treilea", "a treia";
4: patru, "al patrulea", "a patra";
5: cinci, "al cincilea", "a cincea";
6: şase, "al şaselea", "a şasea";
7: şapte, "al şaptelea", "a şaptea";

# ... şi așa mai departe ... (până la 24 + numerele de săli)
```

influx - reducerea inflexiunilor

```
# Fisier control preprocesare - reducerea inflexiunilor
```

```
an: anul, anii;
grupă: grupa, grupe, grupele;
subgrupă: subgrupa, subgrupe, subgrupele, semigrupă,
    semigrupa, semigrupe, semigrupele;
serie: seria, seriile, serii;
secție: secția, secțiile, secții;

laborator: laboratorul, laboratoarele, laboratoare;
curs: cursul, cursuri, cursurile;
seminar: seminarul, seminarele, seminare, seminarii, seminariile;
proiect: proiectul, proiectele, proiecte;

sală: sala, sălile, săli, sale;

student: studentul, studentă, studenta, studenți, studenții;
profesor: profesorul, profesoară, profesoara, profesori, profesorii,
    "cadru didactic", "cadrul didactic";
asistent: asistentul, asistentă, asistenta, asistenți, asistenții;

zi: zile, ziua, zilele, zilei;
dimineață: dimineața, diminețile;
seară: seara, serile, deseară;
noapte: noaptea, nopțile;
semestru: semestrul, semestre, semestrele;

astăzi: azi;

orar: orarul, orare, orarele;
ore: orele;
oră: ora;

duminică: dumineca, duminica;
luni: luna;
marți: marțea;
```



```
miercuri: miercurea;  
joi: joia;  
vineri: vinerea;  
sâmbătă: sâmbăta;  
  
impar: impară, impare;
```

expr.pp - unificarea expresiilor

```
# Fisier control preprocesare - unificarea expresiilor
```

```
ce-ore: "ce ore", "ce orar";  
ce-oră: "ce oră";  
ce-curs: "ce curs";  
ce-laborator: "ce laborator";  
ce-seminar: "ce seminar";  
ce-proiect: "ce proiect";  
ce-profesor: "ce profesor", "ce asistent";  
  
ce-sală: "ce sală";  
cu-cine: "cu cine", "cu ce profesor", "cu ce asistent";  
  
de-la: "de la oră", "de la";  
la: "la oră";  
  
după-amiază: "după masă", "după amiază", "după masa", "după amiaza";  
ce-zi: "ce zi";  
ce-am: "ce am";
```

alias.pp - înlocuirea aliasurilor

```
# Fisier control preprocesare - înlocuiri alias-uri
```

```
# Aliasuri denumiri materii -----
```

```
"Dispozitive și circuite electronice":  
  "d c e",  
  "de ce e",  
  "dispozitive și circuite electronice",  
  "circuite electronice";  
  
"Matematici speciale":  
  "matematici speciale",  
  "me se",
```

```
"m s",
"matematici";

"Analiza și sinteza dispozitivelor numerice":
  "analiza și sinteza dispozitivelor numerice",
  "a s d n",
  "a se de ne";

# ... și așa mai departe

# Aliasuri sali de clasa -----
A101: "a 101", "a 1 0 1", "101";
A106: "a 106", "a 1 0 6", "106";
A109: "a 109", "a 1 0 9", "109";
A110: "a 110", "a 1 1 0", "110";

# ... și așa mai departe

# Aliasuri specialitati-----
H1: "hard 1", "h 1";
H2: "hard 2", "h 2";
S1: "soft 1", "s 1";
S2: "soft 2", "s 2";

# Aliasuri profesori -- probabil nu foarte acurate-----
"Sorin Babii":
  "Babii",
  "Babii Sorin";

"Dan Cosma":
  "Cosma",
  "Cosma Dan";

"Vladimir Crețu":
  "Crețu",
  "Crețu Vladimir";

# ... și așa mai departe
```

categ.pp - unificarea categoriilor

```
# Fişierul control preprocesare - unificarea categoriilor
```

```
MATERIA:
```

```
"Dispozitive şi circuite electronice",  
"Matematici speciale",  
"Sport",  
"Engleza",  
"Analiza şi sinteza dispozitivelor numerice",  
"Programarea calculatoarelor",  
"Electrotehnică II",
```

```
# ... şi aşa mai departe ...
```

```
ZI-SAPT:
```

```
luni,  
marţi,  
miercuri,  
joi,  
vineri,  
sâmbătă,  
duminică;
```

```
ZI-RELATIV:
```

```
astăzi,  
măine,  
poimăine,  
răspoimăine,  
ieri,  
alaltăieri;
```

```
GRUPA:
```

```
H1, H2, S1, S2;
```

```
NR:
```

```
1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24;
```

```
PROFESOR:
```

```
"Sorin Babii",  
"Nicolae Bogoevici",  
"Radu Boraci",  
"Carmen de Sabata",  
"Horia Ciocârlie",  
"Dan Cosma",
```

```
# ... şi aşa mai departe
```

```
SPECIALITATEA:
```

```

calculatoare,
automatică,
colegiu;

```

PERIOADA-ZI:

```

dimineață,
după-amiază,
seară,
prânz,
noapte;

```

SALA:

```

A101, A106, A109, A110, A117, A204, A212,
A305, A314, A315, ASPC, B011, B020, B028,
B303, B305, B413, B414, B419, B425, B513,
B514, B520, B528, B528B, B529, B623, B624,
C312, CA2, SPM;

```

D.2 Controlul constructorului de cadre

Modulul de construcție a cadrelor citește definiția sistemului de cadre utilizat pentru realizarea analizei semantice dintr-un fișier de configurare. Structura acestui fișier este asemănătoare cu cea a fișierelor de control al preprocesării, fiind de tip text și având intrări de forma:

```
<concept>: <slot_1>,<slot_2>, ...,<slot_N>;
```

unde fiecare concept definește un cadru. Sloturile pot fi șiruri de caractere, caz în care valorile lor sunt simple, sau pot fi la rândul lor concepte (încadrate între caracterele “<” și “>”), caz în care valoarea slotului respectiv indică o legătură spre un subcadru dat de respectivul concept.

Pentru exemplificare, prezentăm în continuare fișierul de descriere a sistemului de cadre utilizat în varianta finală a analizorului semantic proiectat.

frames.fb - construcția cadrelor

```
<da>; # cadru fără sloturi
```

```
<nu>; # cadru fără sloturi
```

```
<identificare>:
```

```
  anul,
```

```
specialitatea,  
seria,  
grupa,  
subgrupa,  
profesor;
```

```
<specif-timp>:  
  zi-relativ,  
  zi-sapt,  
  perioada-zi,  
  ora-start,  
  ora-stop;
```

```
<specif-materie>:  
  tip-materie,  
  materie,  
  laborator,  
  seminar,  
  proiect,  
  curs;
```

```
<ce-materie>:  
  <identificare>, # subcadru  
  <specif-materie>, # subcadru  
  <specif-timp>, # subcadru  
  sala; # cadru simplu
```

```
<ce-curs>:  
  <identificare>,  
  <specif-materie>,  
  <specif-timp>,  
  sala;
```

```
<ce-laborator>:  
  <identificare>,  
  <specif-materie>,  
  <specif-timp>,  
  sala;
```

```
<ce-seminar>:  
  <identificare>,  
  <specif-materie>,  
  <specif-timp>,  
  sala;
```

```
<ce-proiect>:  
  <identificare>,  
  <specif-materie>,  
  <specif-timp>,  
  sala;
```

```
<cand>:  
  <identificare>,  
  <specif-materie>,  
  <specif-timp>,  
  sala;
```

```
<unde>:  
  <identificare>,  
  <specif-materie>,  
  <specif-timp>,  
  sala;
```

```
<cu-cine>:  
  <identificare>,  
  <specif-materie>,  
  <specif-timp>,  
  sala;
```

Anexa E

Extrase de cod C++

E.1 Antete de clase

CStringSequence

```
/* *****
                                     stringsequence.h
*****
Models a sequence of strings
*****
begin                               : Wed Apr 5 2000
copyright                             : (C) 2000 by Dan Bohus
email                                  : danbo@ear.utt.ro
***** */

#ifndef STRINGSEQUENCE_H
#define STRINGSEQUENCE_H

/**
 * @author Dan Bohus
 */

#define STRINGSEQUENCE_DEFAULT_SIZE 10
#define STRINGSEQUENCE_DELTA 10

// CStringSequence class - models a sequence of strings data
// structure. The strings are represented as a dynamically
// allocated array that grows with STRINGSEQUENCE_DELTA
// when necessary.
class CStringSequence {
private:
    int Size;           // allocated array size
    char **Values;     // strings array
};
```

```
public :
    int Length;      // number of strings
public :
    // Orthodox-canonical class construction
    CStringSequence();
    CStringSequence(const CStringSequence& aStringSequence);
    CStringSequence& operator=(const CStringSequence&
                                aStringSequence);
    bool operator==(const CStringSequence& aStringSequence);
    ~CStringSequence();

    // string sequence manipulation & access routines
    void Clear();
    int Add(const char* aValue);
    int Insert(int aIndex, const char* aValue);
    void Delete(int aIndex);
    void Reverse();
    char*& operator [] (const int aIndex);
};

#endif
```


CSequence

```

/*****
                                sequence.h
*****
Models a sequence of integers
*****
begin                          : Tue Apr 4 2000
copyright                        : (C) 2000 by Dan Bohus
email                            : danbo@ear.utt.ro
*****/

#ifndef SEQUENCE_H
#define SEQUENCE_H

/**
 * @author Dan Bohus
 */

#define SEQUENCE_DEFAULT_SIZE 20
#define SEQUENCE_DELTA 10

// CSequence class - models a sequence of integers data
// structure. The integers are stored as a dynamically
// allocated array that grows with SEQUENCE_DELTA when
// necessary.
class CSequence {
private:
    int Size;           // allocated array size
    int* Values;       // integers array
public:
    int Length;        // number of integers
public:
    // Orthodox-canonical class construction
    CSequence();
    CSequence(const CSequence& aSequence);
    CSequence& operator=(const CSequence& aSequence);
    ~CSequence();

    // sequence manipulation & access routines
    void Clear();
    int Add(const int aValue);
    void Delete(int aIndex);
    void Reverse();
    int& operator[](const int aIndex);
};

#endif

```

CCaseFrame

```

/*****
                                caseframe.h
*****/
Models the case-frame data structure.
*****/
begin                                : Wed May 3 2000
copyright                             : (C) 2000 by Dan Bohus
email                                  : danbo@ear.utt.ro
*****/

#ifndef CASEFRAME_H
#define CASEFRAME_H

/**
 * @author Dan Bohus
 */

#include <stdio.h>

// forward CCaseFrame class declaration
class CCaseFrame;

// union used to encompass String and SubFrame values for
// case-values
typedef union TCaseValue {
    char* String;
    CCaseFrame* SubFrame;
};

// result type for case-frame definition reading routine
typedef enum CReadFrameDefinitionResult {rfdROK,
    rfdRError, rfdREndOfFile};

// CCaseFrame class - models a case frame.
class CCaseFrame {
private:
    char* Concept;                // frame concept
    int AttributesNo;             // number of attributes
    char** Attributes;           // attributes array
    TCaseValue* Values;          // values array
public:
    // Orthodox-canonical Class Construction
    CCaseFrame();
    CCaseFrame(CCCaseFrame& aCaseFrame);
    CCaseFrame& operator=(CCCaseFrame& aCaseFrame);
    bool operator==(CCCaseFrame& aCaseFrame);
    ~CCaseFrame();

```

```
// case-frame clearing routines
void Clear(); // clears the entire frame
void ClearValues(); // clears the case-values

// seralization routines
CReadFrameDefinitionResult ReadDefinitionFromFile(FILE*
                                                fileP);
bool DecodeFromBuffer(const char* buffer, int& pointer);
void EncodeToBuffer(char* buffer);
void NicePrintToFile(FILE* fileP, bool printVoidValues,
                    int tab=0);

// frame manipulation & access routines
bool IsInstantiated();
char* GetConcept();
int GetAttributesNo();
bool HasAttribute(char* attribute);
char* GetAttribute(int aIndex);
int GetIndexOfAttribute(char* attribute);
TCaseValue* operator [] (char* attribute);
TCaseValue* operator [] (int aIndex);

bool SetValueForAttribute(char* attribute, char* value);
};

#endif
```

CCaseFrameSystem

```

/*****
                                caseframesystem.h
*****/
Models a case-frame system
*****/
begin                                : Wed May 3 2000
copyright                            : (C) 2000 by Dan Bohus
email                                : danbo@ear.utt.ro
*****/

#ifndef CASEFRAMESYSTEM_H
#define CASEFRAMESYSTEM_H

/**
 * @author Dan Bohus
 */

#include "caseframe.h"

#define CASEFRAMESYSTEM_DEFAULT_SIZE 10
#define CASEFRAMESYSTEM_DELTA 3

// CCaseFrameSystem class - models a set of case frames
// forming a case-frame system. The set is implemented as a
// dynamically allocated array of CCaseFrame objects
class CCaseFrameSystem {
private:
    int Size;                        // allocated array size
    CCaseFrame** CaseFrames;        // case-frames
public:
    int Length;                      // number of case-frames in the system
public:
    // Orthodox-canonical class construction
    CCaseFrameSystem();
    CCaseFrameSystem(CCaseFrameSystem& aCaseFrameSystem);
    CCaseFrameSystem& operator=(CCaseFrameSystem&
                                aCaseFrameSystem);
    ~CCaseFrameSystem();

    // read definition from file
    bool LoadFromFile(const char* fileName);
    // case frame system acces & manipulation routines
    int Add(CCaseFrame& aCaseFrame);
    CCaseFrame& operator [] (const int aIndex);
    void Clear();
};

#endif

```

CUtterance

```

/*****
                                     utterance.h
*****
Models a user utterance, with all possible forms through
decoding: speech recognizer output, normalized form,
parse & case-frame representation
*****
begin                               : Wed Apr 5 2000
copyright                             : (C) 2000 by Dan Bohus
email                                  : danbo@ear.utt.ro
*****/

#ifndef UTTERANCE_H
#define UTTERANCE_H

/**
 * @author Dan Bohus
 */

#include <stdio.h>

#include "stringsequence.h"
#include "caseframe.h"

// CUtterance class - models an utterance
class CUtterance {
public:
    char* Class; // utterance class
    char* SpeechRecognizerOutput; // speech recognizer output
    CStringSequence NormalizedForm; // normalized form
    CStringSequence Parse; // parse
    int CaseFramesNo; // number of case-frames
    CCaseFrame* CaseFrames; // case-frames
public:
    // Orthodox-canonical class construction
    CUtterance();
    CUtterance(char* aClass,
               char* aSpeechRecognizerOutput,
               CStringSequence* aNormalizedForm = NULL,
               CStringSequence* aParse = NULL,
               int aCaseFrameNo = 0,
               CCaseFrame* someCaseFrames = NULL);
    CUtterance(FILE* fileP);
    CUtterance(CUtterance& aUtterance);
    CUtterance& operator=(CUtterance& aUtterance);
    bool operator==(CUtterance& aUtterance);
    ~CUtterance();

```

```
// serialization routines
bool LoadFromFile(FILE* fileP);
bool SaveToFile(FILE* fileP);

// utterance access & manipulation routines
bool HasClass();
bool HasSpeechRecognizerOutput();
bool HasNormalizedForm();
bool HasParse();
bool HasCaseFrames();

void Clear();
void ClearCaseFrames();
void DeleteCaseFrame(const int aIndex);
void SetClass(char* aClass);
void SetSpeechRecognizerOutput(char*
                                aSpeechRecognizerOutput);
};

#endif
```

CUtteranceCorpus

```

/*****
                                     utterancecorpus.h
*****/
Models a corpus of utterances
*****/
begin                               : Thu Apr 6 2000
copyright                             : (C) 2000 by Dan Bohus
email                                  : danbo@ear.utt.ro
*****/

#ifndef UTTERANCECORPUS_H
#define UTTERANCECORPUS_H

/**
 * @author Dan Bohus
 */

#include "utterance.h"

#define UTTERANCECORPUS_DEFAULT_SIZE 500
#define UTTERANCECORPUS_DELTA 50

// CUtteranceCorpus class - models a corpus of utterances.
// The utterances are stored as a dynamically allocated
// array that grows with UTTERANCECORPUS_DELTA whenever
// necessary
class CUtteranceCorpus {
private:
    int Size;                          // allocated array size
    Utterance** Utterances;           // utterances array
public:
    int Length;                        // number of utterances
public:
    // Orthodox-canonical class construction
    CUtteranceCorpus();
    CUtteranceCorpus(const char* fileName);
    CUtteranceCorpus(CUtteranceCorpus& aUtteranceCorpus);
    CUtteranceCorpus& operator=(CUtteranceCorpus&
                                aUtteranceCorpus);
    ~CUtteranceCorpus();

    // serialization routines
    bool AddFromFile(const char* fileName);
    bool LoadFromFile(const char* fileName);
    bool SaveToFile(const char* fileName);

    // utterance corpus access & management routines
    int Add(CUtterance& aUtterance);

```

```
    void Clear();  
    CUtterance& operator [] (const int aIndex);  
};  
  
#endif
```


CDictionary

```

/*****
                                dictionary.h
*****/
Modells a dictionary data structure
*****/
begin                               : Tue Apr 4 2000
copyright                             : (C) 2000 by Dan Bohus
email                                  : danbo@ear.utt.ro
*****/

#ifndef DICTIONARY_H
#define DICTIONARY_H

/**
 * @author Dan Bohus
 */

#include "sequence.h"
#include "stringsequence.h"

#define DICTIONARY_DEFAULT_SIZE 500
#define DICTIONARY_DELTA 50

// CDictionary class - models a typical dictionary data
// structure. The words are represented as a dynamically
// allocated array that grows with DICTIONARY_DELTA when
// necessary
class CDictionary {
private:
    int Size;           // allocated size
    char** Strings;    // strings array
public:
    int Length;        // number of strings in dictionary
public:
    // Orthodox-canonical class construction
    CDictionary();
    ~CDictionary();
    CDictionary(const CDictionary& aDictionary);
    CDictionary& operator=(const CDictionary& aDictionary);

    // serialization routines
    bool LoadFromFile(const char* fileName);
    bool SaveToFile(const char* fileName);

    // dictionary access routines
    char* operator [] (const int aIndex);
    int IndexOf(const char* aString);
    bool Exists(const char* aString);

```

```
// dictionary manipulation routines
int Add(const char* aString);
void Delete(const char* aString);
void Clear();

// dictionary-based conversion routines
bool ToIndexes(CStringSequence& ss, CSequence& s);
bool ToStrings(CSequence& s, CStringSequence& ss);
};

#endif
```

CPreprocessor

```

/*****
                                preprocessor.h
*****
Implements the preprocessing module
*****
begin                          : Sat Apr 8 2000
copyright                       : (C) 2000 by Dan Bohus
email                           : danbo@ear.utt.ro
*****/

#ifndef PREPROCESSOR_H
#define PREPROCESSOR_H

/**
 * @author Dan Bohus
 */

#include "preprocessorconfiguration.h"
#include "utterancecorpus.h"
#include "sequence.h"
#include "stringsequence.h"
#include "dictionary.h"
#include "repltable.h"

// preprocessor stage type definition
typedef enum TReplType {rtNumbersNormalization,
                       rtInflexReduction,
                       rtExprUnification,
                       rtAliasesReplacement,
                       rtCategoryUnification};

// CPreprocessor class – models semantic preprocessing
class CPreprocessor {
private:
    CPreprocessorConfiguration PreprocessorConfiguration;
    TReplTable ReplTables[5];           // replacement tables
    CStringSequence Replacements;      // replacements
    CDictionary Dictionary;            // normalized dictionary
private:
    // internal configuration routines
    bool Configure();
    bool ConfigureType(TReplType type, const char* fileName);
    void StringToSequence(char* from, CStringSequence& sequence);
public:
    // Orthodox-canonical class construction
    CPreprocessor();
    CPreprocessor(CPreprocessor& aPreprocessor);
    CPreprocessor& operator=(CPreprocessor& aPreprocessor);

```

```
~CPreprocessor();

// set / get preprocessor configuration routines
bool SetConfiguration(CPreprocessorConfiguration&
                    aPreprocessorConfiguration);
CPreprocessorConfiguration& GetConfiguration();

// preprocessing routines
CStringSequence PreprocessString(char* string);
void PreprocessUtterance(CUtterance& aUtterance);
void PreprocessCorpus(CUtteranceCorpus& aUtteranceCorpus);

// preprocessing stages routines
void PerformNonLexicalEventsReduction(char* from, char* to);
void PerformNumbersNormalization(char* from, char* to);
void PerformInflexReduction(char* from, char* to);
void PerformExprUnification(char* from, char* to);
void PerformAliasesReplacement(char* from, char* to);
void PerformCategoryUnification(char* from, char* to);
void PerformDictionaryCheck(char* from,
                          CStringSequence& sequence);
};

#endif
```

CHMM

```

/*****
                                hmm.h
*****
Models the Hidden Markov Model. The A, B and Pi probability
distributions are represented as matrices. Training is
accomplished by counting events in the training corpus.
Routines for Katz reestimation and Viterbi decoding
are implemented.
*****
begin                               : Thu Apr 20 2000
copyright                             : (C) 2000 by Dan Bohus
email                                  : danbo@ear.utt.ro
*****/

#ifndef HMM_H
#define HMM_H

/**
 * @author Dan Bohus
 */

#include "dictionary.h"
#include "utterancecorpus.h"

// type of HMM closing used
typedef enum TClosingType {ctNone,
                           ctMaximumLikelihood,
                           ctKatzReestimation};

// CHMM class
class CHMM {
private:
    int StatesNo;    // number of internal (hidden) states
    int SymbolsNo;  // number of symbols

    int* trISPM;    // training initial state prob. matrix
    int* trOSPM;    // training observation symbol prob. matrix
    int* trSTPM;    // training state transition prob. matrix

    double* ISPM;   // initial state probabilities matrix
    double* OSPM;   // observation symbol probabilities matrix
    double* STPM;   // state transition probabilities matrix

    CDictionary SymbolsDictionary; // obs. symbols dictioanry
    CDictionary StatesDictionary;  // states dictionary

public:
    // Orthodox-canonical class construction

```

```

CHMM();
CHMM(const char* aSymbolsDictionaryFileName,
      const char* aStatesDictionaryFileName);
CHMM(char* fileName);
CHMM(CHMM& aHMM);
CHMM& operator=(CHMM& aHMM);
~CHMM();

// serialization routines
bool LoadFromFile(const char* fileName);
bool SaveToFile(const char* fileName);

// HMM model management routines
void Clear();
bool IsEmpty();
bool IsClosed();

// public training routines
int TrainWithUtterance(CUtterance& aUtterance, FILE* logFP);
void TrainWithCorpus(CUtteranceCorpus& aCorpus, FILE* logFP);

// public closure routines
void CloseModel(TClosingType ct, int k=0);

// public state decoding routines
double DecodeUtterance(CUtterance& aUtterance);
void DecodeCorpus(CUtteranceCorpus& aCorpus);

// printing routine
void PrintModelParameters(FILE* fp);

private:
// closure routines:
void MLEClosure(); // maximum likelihood estimation
void KatzClosure(int k); // Katz re-estimation

// decoding (Viterbi) algorithm
double Viterbi(CSequence& obs, CSequence& states);
};

#endif

```



```
~CCaseFrameBuilder();

// case-frame builder settings control
void SetCaseFrameSystem(CCaseFrameSystem& aCaseFrameSystem);
CCaseFrameSystem& GetCaseFrameSystem();
void DoFramesNormalization(const char* fileName);

// public frame construction routines
TFrameBuilderResult ProcessUtterance(CUtterance& aUtterance);
void ProcessCorpus(CUtteranceCorpus& aUtteranceCorpus);
};

#endif
```


CSemanticAnalyzer

```

/*****
                                semanticanalyzer.h
*****/
*****/
Models the semantic analyzer by aggregating separate
preprocessing, decoding and frame-building modules
*****/
begin                               : Fri May 19 2000
copyright                             : (C) 2000 by Dan Bohus
email                                  : danbo@ear.utt.ro
*****/

#ifndef SEMANTICANALYZER_H
#define SEMANTICANALYZER_H

/**
 * @author Dan Bohus
 */

#include <stdio.h>

#include "semanticanalyzerconfiguration.h"
#include "utterance.h"
#include "utterancecorpus.h"
#include "conceptarray.h"
#include "preprocessor.h"
#include "hmm.h"
#include "caseframebuilder.h"

// CSemanticAnalyzer class
class CSemanticAnalyzer {
private:
    CSemanticAnalyzerConfiguration SemanticAnalyzerConfiguration;
    CPreprocessor Preprocessor;           // preprocessor module
    CHMM HMM;                           // HMM decoding module
    CCaseFrameBuilder CaseFrameBuilder; // frame build. module

private:
    // internal configuration routine
    bool Configure();

public:
    // Orthodox-canonical class construction
    CSemanticAnalyzer(const char* fileName = NULL);
    CSemanticAnalyzer(CSemanticAnalyzer& aSemanticAnalyzer);
    CSemanticAnalyzer& operator=(CSemanticAnalyzer&
                                aSemanticAnalyzer);
    ~CSemanticAnalyzer();

```

```
// get / set analyzer configuration
bool SetConfiguration(CSemanticAnalyzerConfiguration&
                      aSemanticAnalyzerConfiguration);
CSemanticAnalyzerConfiguration& GetConfiguration();

// semantic analysis routines
char* AnalyzeString(char* string,
                    CConceptArray* aConceptArray = NULL);
void AnalyzeUtterance(CUtterance& aUtterance,
                     CConceptArray* aConceptArray = NULL);
void AnalyzeCorpus(CUtteranceCorpus& aUtteranceCorpus);
};

#endif
```

E.2 Algoritmi

Reestimarea Katz

```

// D: perform Katz Closure
// Katz discounts all events occurring less than k times,
// considering them unreliable. The freed probability mass
// is redistributed among unseen events.
// LittleK - used for reestimation of initial state prob.
// BigK - used for reestimation of state transition prob.

void CHMM::KatzClosure(int LittleK, int BigK){
    if (!IsClosed()) {
        // allocate floating point matrices
        ISPM = new double [StatesNo];
        OSPM = new double [StatesNo*SymbolsNo];
        STPM = new double [StatesNo*StatesNo];

        // normalize trOSPM - just like MLE (Katz only works on
        // state transition and initial state probabilities)
        for (int i=0; i<StatesNo; i++)
            for (int j=0; j<SymbolsNo; j++)
                if (trOSPM[i*(SymbolsNo+1)+j]==0)
                    OSPM[i*SymbolsNo+j] = 0;
                else
                    OSPM[i*SymbolsNo+j] = trOSPM[i*(SymbolsNo+1)+j]/
                        (double)trOSPM[i*(SymbolsNo+1)
                        +SymbolsNo];

        // if LittleK = 0, perform simple MLE
        if (LittleK == 0) {
            // normalize ISPM
            for (int i=0; i<StatesNo; i++)
                if (trISPM[i]==0) ISPM[i]=0;
                else ISPM[i] = trISPM[i]/(double)trISPM[StatesNo];
        } else {
            // else perform Katz reestimation for initial state
            // probability distribution

            // first compute initial state density distribution
            // find out most occurent initial state
            int maxValue = -1;
            int maxI = 0;
            for (int i=0; i<StatesNo; i++)
                if (trISPM[i]>maxValue) {
                    maxValue = trISPM[i];
                }
        }
    }
}

```

```

    maxI = i;
}

// allocate space and initialize density distribution
int *nis = new int [maxValue+1];
for (int i=0; i<=maxValue; nis[i++] = 0);

// compute initial state density distribution
for (int i=0; i<StatesNo; i++)
    nis[trISPM[i]]++;

// compute normalizing factor for discount coefficients
double fis = 1 - ((LittleK+1)*nis[LittleK+1])/
    (double)(nis[1]);

// allocate initial state discount coefficients
double* dcis = new double [LittleK];

// precompute initial state discount parameters
for (int i=1; i<=LittleK; i++)
    // compute discount parameter for i, store it at
    // dcis[i-1]
    dcis[i-1] = (((i+1)*nis[i+1])/((double)(i*nis[i]))+
        fis-1)/fis;

// now for the initial states
// compute A = freed probability mass
// compute B = normalization factor for freed probability
// mass redistribution
double A = 1; // start from full probability
double B = 0; // start from zero
bool noDiscount = 1;
for (int i=0; i<StatesNo; i++) {
    if (trISPM[i]>0) {
        // positive initial state occurrence
        if (trISPM[i]<=LittleK) { // discounted
            A -= dcis[trISPM[i]-1]*trISPM[i]/
                (double)trISPM[StatesNo];
            noDiscount = 0;
        } else // not discounted
            A -= trISPM[i]/(double)trISPM[StatesNo];
    }
    else B += trSTPM[i*(StatesNo+1)+StatesNo];
}

// if no discount occurred, A must be at or pretty
// close to 0, so manually normalize it to avoid small
// non-zero values
if (noDiscount) A = 0;

```

```

// now, finally we can compute the reestimated Katz
// probabilities for the initial state distribution
for (int i=0; i<StatesNo; i++) {
    if (trISPM[i]>0) {
        if (trISPM[i]<=LittleK)
            ISPM[i] = dcis[trISPM[i]-1]*trISPM[i]/
                (double)trISPM[StatesNo];
        else ISPM[i] = trISPM[i]/(double)trISPM[StatesNo];
    }
    else ISPM[i] = (A/B)*trSTPM[i*(StatesNo+1)+StatesNo];
}

// deallocate dcis and nis
delete [] nis;
delete [] dcis;
}

// now go to trSTPM - state transition probabilities
// if no BigK, then perform simple MLE
if (BigK == 0) {
    // normalize trSTPM
    for (int i=0; i<StatesNo; i++) {
        int total = 0;
        for (int j=0; j<StatesNo; j++)
            total += trSTPM[i*(StatesNo+1)+j];
        for (int j=0; j<StatesNo; j++)
            if (trSTPM[i*(StatesNo+1)+j]==0)
                STPM[i*StatesNo+j] = 0;
            else STPM[i*StatesNo+j] = trSTPM[i*(StatesNo+1)+j]/
                (double)total;
    }
} else {
    // else perform Katz reestimation on state transition
    // probabilities

    // first determine bigram density distribution
    // first find out most occurent bigram
    int maxValue = -1;
    int maxI = 0, maxJ = 0;
    for (int i=0; i<StatesNo; i++)
        for (int j=0; j<StatesNo; j++)
            if (trSTPM[i*(StatesNo+1)+j]>maxValue) {
                maxValue = trSTPM[i*(StatesNo+1)+j];
                maxI = i;
                maxJ = j;
            }

    // allocate space and initialize bigram density distr.

```

```

int *n = new int [maxValue+1];
for (int i=0; i<=maxValue; n[i++] = 0);

// compute bigram density distribution
for (int i=0; i<StatesNo; i++)
    for (int j=0; j<StatesNo; j++)
        n[trSTPM[i*(StatesNo+1)+j]]++;

// compute normalizing factor for discount coefficients
double f = 1 - ((BigK+1)*n[BigK+1])/(double)(n[1]);

// allocate discount coefficients
double* dc = new double [BigK];

// precompute discount parameters
for (int i=1; i<=BigK; i++)
    // compute discount parameter for i, store it
    // at dc[i-1]
    dc[i-1] = (((i+1)*n[i+1])/(double)(i*n[i]))+f-1)/f;

// now for each state
for (int i=0; i<StatesNo; i++) {
    // ci - total occurrences of state i as the first state
    // in a transition
    int ci = 0;
    for (int j=0; j<StatesNo; j++)
        ci += trSTPM[i*(StatesNo+1)+j];

    // A - freed probability mass
    // B - normalizing constant
    double A = 1; // start from full probability
    double B; // start from zero

    bool noDiscount = 1;

    // for each transition compute freed probability mass
    for (int j=0; j<StatesNo; j++) {
        int cij = trSTPM[i*(StatesNo+1)+j];
        if (cij>0) {
            if (cij<=BigK) {
                A -= dc[cij-1]*cij/(double)ci;
                noDiscount = 0;
            } else A -= cij/(double)ci;
        }
        else B += trSTPM[j*(StatesNo+1)+StatesNo];
    }

    // if no discount occurred, A must be at or pretty
    // close to 0, so manually normalize it to avoid small

```

```

        // non-zero values
        if (noDiscount) A = 0;

        // now, finally we can compute the reestimated Katz
        // probabilities
        for (int j=0; j<StatesNo; j++) {
            int cij = trSTPM[i*(StatesNo+1)+j];
            int cj = 0;
            for (int l=0; l<StatesNo; l++)
                cj += trSTPM[j*(StatesNo+1)+l];
            if (cij>0) {
                if (cij<=BigK)
                    STPM[i*StatesNo+j] = dc[cij-1]*cij/(double)ci;
                else
                    STPM[i*StatesNo+j] = cij/(double)ci;
            } else
                STPM[i*StatesNo+j] = (A/B)*
                    trSTPM[j*(StatesNo+1)+StatesNo];
        }
    }

    // don t forget to deallocate n and dc
    delete [] n;
    delete [] dc;
}

// finally deallocate training matrices
delete [] trISPM;
delete [] trOSPM;
delete [] trSTPM;
trISPM = trOSPM = trSTPM = NULL;
}
}

```

Algoritmul Viterbi

```

// D: performs Viterbi decoding on the specified
// observation sequence
// returns the maximum path probability obtained

double CHMM::Viterbi(CSequence& obs, CSequence& states){
    // get observation sequence length
    int m = obs.Length;
    int n = StatesNo;

    // allocate delta(symbol, state), phi(state, nextstate)
    double* delta = new double[n*m];
    int* phi = new int[n*m];

    // INITIALIZATION
    for(int state=0; state<n; state++) {
        delta[/(n*0)+*/state] = ISPM[state]*
                                OSPM[state*SymbolsNo+obs[0]];
        phi[/(n*0)+*/state] = 0;
    }

    // RECURSION
    for(int symbol=1; symbol<m; symbol++) {
        for(int state=0; state<n; state++) {
            double maxvalue = -1;
            int maxarg = -1;
            for(int prevstate=0; prevstate<n; prevstate++) {
                double newvalue = delta[(n*(symbol-1))+prevstate]*
                                    STPM[prevstate*StatesNo+state];
                if(newvalue>maxvalue) {
                    maxvalue = newvalue;
                    maxarg = prevstate;
                }
            }
            delta[(n*symbol)+state] = maxvalue *
                                    OSPM[state*SymbolsNo+obs[symbol]];
            phi[(n*symbol)+state] = maxarg;
        }
    }

    // TERMINATION
    double maxvalue = -1;
    int maxarg = -1;
    for(int state=0; state<n; state++) {
        double newvalue = delta[n*(m-1)+state];
        if(newvalue>maxvalue) {
            maxvalue = newvalue;
            maxarg = state;
        }
    }
}

```



```
    }  
  }  
  
  // PATH BACKTRACKING  
  states.Clear();  
  for (int i=0; i<m; i++) states.Add(i);  
  
  states[m-1] = maxarg;  
  for (int i=m-2; i>=0; i--)  
    states[i] = phi[(n*(i+1))+states[i+1]];  
  
  return maxvalue;  
}
```


Anexa F

Invocarea programelor utilitare

hmm_build

```
hmm_build      [-create|-train|-close filename] [-l|-la [filename]]  
               <specific_options>
```

Options:

- create filename - creates a new HMM model to filename
- train filename - trains an existing HMM model
- close filename - closes the training of an existing HMM model
- l filename - create log file -> filename
- la filename - append to log file -> filename

Specific CREATE options:

- sym filename - symbols dictionary filename
- st filename - states dictionary filename

Specific TRAIN options:

- corpus filename - training corpus filename

Specific CLOSE options:

- mle - close for maximum likelihood estimations
- katz k1 - close by Katz re-estimation of state transition probabilities with parameter k1
- katzI k2 - close by Katz re-estimation of initial state probabilities with parameter k2

hmm_print

```
hmm_print      filename [-o filename]
```

Options:

- filename - HMM model to print
- o filename - file to direct output to

hmm_decode

```
hmm_decode      hmm_filename -i filename [-o filename]
Options:
hmm_filename    - HMM filename
-i filename     - utterance corpus filename to decode
-o filename     - utterance corpus filename to write to
                  if absent, output is directed to input file
```

preproc

```
preproc         filename [-o filename] <tasks>
Options:
filename        - utterance corpus filename to process
-o filename     - utterance corpus filename to write to
                  if absent, output is directed to input file

Tasks:
-nler           - do non-lexical events reduction
-nn filename   - do numbers normalization as specified by filename
-ir filename   - do inflexions reduction as specified by filename
-eu filename   - do expression unification as specified by filename
-a filename    - do aliases replacement as specified by filename
-cu filename   - do category unification as specified by filename
-dc filename   - do dictionary check as specified by filename
```

frame_build

```
frame_build     frames_filename [-nor filename] -i filename [-o filename]
Options:
frames_filename - filename containing frame-system description
-nor filename  - normalize frames according to filename
-i filename    - utterance corpus filename to process
-o filename    - utterance corpus filename to write to
                  if absent, output is directed to input file
```

corpus_stats

corpus_stats filename [-o filename] [-sld dic_filename]

Options:

filename - utterance corpus filename to process
 -o filename - statistics filename to write to
 if absent, output is directed to stdout
 -sld dic_filename - compute semantic labels density (dic_filename
 is the semantic labels dictionary)

corpus_filter

corpus_filter filename [-o filename] [-keep|-discard regexp] [-nonor]
 [-noprs] [-nofrm]

Options:

filename - utterance corpus filename to filter
 -o filename - utterance corpus filename to write to
 if absent, output is directed to input file
 -keep regexp - keep only utterances with names matching regexp
 -discard regexp - discard utterances with names matching regexp
 -nonor - eliminate NOR form from utterances
 -noprs - eliminate PRS form from utterances
 -nofrm - eliminate FRM form from utterances

corpus_diff

corpus_diff -c1 filename -c2 filename -utt|-sro|-nor|-prs|-frm
 [-markdiff mark]

Options:

-c1 filename - first corpus filename
 -c2 filename - second corpus filename
 -utt - compare whole utterances
 -sro - compare SRO forms
 -nor - compare NOR forms
 -prs - compare PRS forms
 -frm - compare FRM forms
 -markdiff mark - mark utterances that are different by appending
 -<mark> to the name of the utterance

dict_build

dict_build [-sro|-nor|-prs filename] [-i|-c filename] [-s filename]

Options:

-sro filename	- extracts words from speech recognizer output
-nor filename	- extracts words from normalized form
-prs filename	- extracts words from parse
-i filename	- increment dictionary
-c filename	- create new dictionary
-s filename	- log dictionary statistics