

UNIVERSITATEA BABEȘ-BOLYAI, CLUJ-NAPOCA
FACULTATEA DE MATEMATICĂ ȘI INFORMATICĂ

**METODE FORMALE ÎN DEZVOLTAREA
CORECTĂ A PROGRAMELOR**

Rezumatul tezei de doctorat

CONDUCĂTOR ȘTIINȚIFIC:

PROF. DR. MILITON FRENȚIU

DOCTORAND:

MIHIȘ ANDREEA-DIANA

CLUJ-NAPOCA

2011

MULȚUMIRI

Mulțumesc pe această cale tuturor celor care mi-au fost alături în elaborarea acestei teze, în special Doamnei Prof. dr. Tătar Doina, Domnului Prof. dr. Militon Frențiu, Domnilor Prof. dr. Ștefan Mărușter și Prof. dr. ing. Ioan Salomie, Familiei Mihăilă, Familiei Mihiș, Colegilor de birou, catedră și departament.

Această lucrare a fost suportată parțial din: Grant UBB tip TP nr. 2/2006, Grant CNCSIS TD 400/2007, Proiect CNMP 91-037/19.09.2007 CONTO, Proiect CNMP 42-117/2008 I-GLOB și Proiect IDEI ID_2412/2009.

Cuprinsul rezumatului:

I. Introducere	1
II. Metode formale în dezvoltarea programelor	4
2.1. Ciclul de viață a unui produs soft	4
2.2. Specificarea cerințelor	4
2.2.1. Specificații informale	4
2.2.2. Tehnici semiformale de specificare	4
2.2.3. Tehnici formale de specificare	4
2.3. Rafinarea în pași succesivi	4
2.4. Metode de procesare a limbajului natural utilizate în analiza cerințelor	4
2.5. Ontologiile	5
III. Contribuții în domeniul prelucrării limbajului natural	6
3.1 Contribuții în dezambiguarea automată a unui text	6
3.1.1 Identificarea sensului cuvintelor	6
3.1.2. Identificarea părților de vorbire a cuvintelor	7
3.1.3. Relația de consecință logică dintre două texte	8
3.1.4. O simplă ontologie bazată pe relația de consecință logică dintre două texte	11
3.2. Noi abordări ale segmentării unui text	11
3.2.1. Segmentare liniară	11
3.2.2. Alte variante de segmentare Logical TextTiling	13
3.2.3. Segmentare genetică cu număr predefinit de segmente bazată pe implicații	14
3.2.4. O nouă metodă de segmentare top-down bazată pe lanțuri lexicale	15
3.3. Noi abordări ale sumarizării unui text	16
3.3.1. Extragerea sumarului dintr-un text segmentat	16
3.3.2. Sumare de lungime variabilă	18
3.3.3. Sumarizare bazată pe lanțuri lexicale	20
IV. Utilizarea ontologiilor în prelucrarea cerințelor	22
4.1. Potrivirea unei ontologii la un text	22
Evaluarea ontologiilor	22
Metrici de evaluare a potrivirii unei ontologii la un text scris în limbaj natural	22
Metrici cantitative	22
Metrici bazate pe relația de consecință logică dintre două texte	23

4.2. Învățarea unei ontologii dintr-un text de limbaj natural pe baza arborelui de analiză sintactică a unei propoziții	23
Analiza sintactică a unei fraze	24
Identificarea tripletelor	24
4.3. Similaritatea utilizată în identificarea nevoii de a refactoriza o ontologie	25
Algoritmul de identificare a elementelor cel mai probabil plasate greșit	26
4.4. Rolul dezambiguării în evaluarea ontologiilor	27
Ontologiile și dezambiguarea	27
Algoritmul de dezambiguare	27
4.5. Prelucrarea cerințelor asistată de ontologii	28
Selectarea semi-automată a cerințelor	28
Concluzii	29
V. Aspecte ale utilizării metodelor formale pentru dezvoltarea corectă a programelor	31
5.1. Asigurarea calității softului – o activitate continuă	31
5.2. Abordarea multiformală în specificarea sistemelor soft	32
5.3. Reutilizarea specificațiilor formale	33
5.4. O aplicație care asistă rafinarea schemelor Z	34
5.5. O aplicație care asistă procesul de rafinare a codului din specificații	36
5.6. Simplificarea codului prin prelucrarea automată a expresiilor condiționale	36
VI. Concluzii	39
VII. Bibliografie	40

Cuprinsul tezei de doctorat:

Lista publicațiilor proprii	1
I. Introducere	5
II. Metode formale în dezvoltarea programelor	18
2.1. Ciclul de viață a unui produs soft	18
2.2. Specificarea cerințelor	20
2.2.1. Documentul de specificare a cerințelor	20
2.2.2. Specificații informale	23
2.2.3. Tehnici semiformale de specificare	23
Tehnici grafice orientate pe acțiuni	23
Tehnici grafice orientate pe date	24
Alte tehnici semiformale de specificare a cerințelor	24
2.2.4. Tehnici formale de specificare	25
Mașini cu stări finite	25
Rețele Petri	26

Limbajul Z	26
Alte tehnici formale	28
2.3. Rafinarea în pași succesivi	29
2.4. Metode de procesare a limbajului natural utilizate în analiza cerințelor	32
Câteva tehnici de colectare a cerințelor	33
Mecanisme de procesare a limbajului natural care pot asista procesul de analiză	34
2.5. Ontologiile	38
III. Contribuții în domeniul prelucrării limbajului natural	41
3.1 Contribuții în dezambiguarea automată a unui text	41
3.1.1 Identificarea sensului cuvintelor	42
3.1.2. Identificarea părților de vorbire a cuvintelor	47
3.1.3. Relația de consecință logică dintre două texte	51
3.1.4. O simplă ontologie bazată pe relația de consecință logică dintre două texte	61
3.2. Noi abordări ale segmentării unui text	64
3.2.1. Segmentare liniară	64
3.2.2. Alte variante de segmentare Logical TextTiling	70
3.2.3. Segmentare genetică cu număr predefinit de segmente bazată pe implicații	73
3.2.4. O nouă metodă de segmentare top-down bazată pe lanțuri lexicale	80
3.3. Noi abordări ale sumarizării unui text	81
3.3.1. Extragerea sumarului dintr-un text segmentat	82
3.3.2. Sumare de lungime variabilă	86
3.3.3. Sumarizare bazată pe lanțuri lexicale	92
IV. Utilizarea ontologiilor în prelucrarea cerințelor	96
4.1. Potrivirea unei ontologii la un text	96
Evaluarea ontologiilor	97
Metrici de evaluare a potrivirii unei ontologii la un text scris în limbaj natural	98
Metrici cantitative	98
Metrici bazate pe relația de consecință logică dintre două texte	99
Precizia și recall-ul	100
Evaluarea	100
4.2. Învățarea unei ontologii dintr-un text de limbaj natural pe baza arborelui de analiză sintactică a unei propoziții	102
Analiza sintactică a unei fraze	103
Identificarea tripletelor	105
Evaluarea	108
4.3. Similaritatea utilizată în identificarea nevoii de a refactoriza o ontologie	109
Similaritatea	110
Elemente similare dintr-o ontologie	111
Plasări greșite	111
Algoritmul de identificare a elementelor cel mai probabil plasate greșit	111
Un exemplu de funcționare a algoritmului	113
Aplicația care implementează algoritmul	113
4.4. Rolul dezambiguării în evaluarea ontologiilor	114
Ontologiile și dezambiguarea	114
Câteva exemple	115
Algoritmul de dezambiguare	118

4.5. Prelucrarea cerințelor asistată de ontologii	119
Selectarea semi-automată a cerințelor	120
Experimente	125
Concluzii	125
V. Aspecte ale utilizării metodelor formale pentru dezvoltarea corectă a programelor	127
5.1. Asigurarea calității softului – o activitate continuă	127
5.2. Abordarea multiformală în specificarea sistemelor soft	130
5.3. Reutilizarea specificațiilor formale	134
5.4. O aplicație care asistă rafinarea schemelor Z	139
5.5. O aplicație care asistă procesul de rafinare a codului din specificații	143
5.6. Simplificarea codului prin prelucrarea automată a expresiilor condiționale	146
VI. Concluzii	152
VII. Bibliografie	155
VIII. Anexe	172
Anexa 1 (a)	172
Anexa 1 (b)	174
Anexa 2	175
Anexa 3	179
Anexa 4	182
Anexa 5	184
Anexa 6	185

I. Introducere

Scopul principal al dezvoltatorilor de soft este de a îndeplini cerințele clienților, indiferent dacă aceste cerințe se referă la funcționalitate, timp sau bani. În general, se dorește ca produsele realizate să fie de calitate superioară, dar principala cerință este legată de corectitudinea softului produs. Prin corectitudinea softului se înțelege calitatea lui de a-și respecta specificațiile. În anumite domenii, existența unui soft care să se comporte exact după cum era prevăzut este esențială. Din păcate cerințele clienților sunt de cele mai multe ori scrise în limbaj natural, limbaj ambiguu și dificil de prelucrat.

O soluție este eliminarea ambiguității cerințelor, fapt posibil datorită metodelor formale. Metodele formale pot fi privite ca un mijloc formal de a descrie problema sau de a modela sistemul [Gab06]. Practic orice aplicare a matematicii în probleme de ingineria sistemelor reprezintă o metodă formală. Metodele formale pornesc de la o bază matematică precisă, din care sunt derivate metode de modelare și analiză.

Prin metode formale se înțelege utilizarea notațiilor, tehnicilor și metodelor matematice în specificarea, proiectarea și implementarea softului [Bur95, Cla96, Tho95], cu alte cuvinte, „Metodele formale sunt matematica informaticii” [Hol96].

Există mai multe tipuri de formalizări. Cea mai simplă este o utilizare a matematicii (notații, metode) pe alocuri la specificare, proiectare, implementare, testare sau întreținerea produsului.

Următorul nivel de formalizare îl reprezintă limbajele formale de specificare, cum sunt limbajul Z [Abr80], Vienna definition method (VDM) [Jon86] sau Gist [Bal85]. Specificațiile formale au o formă predefinită, utilizează notații matematice și este posibilă și verificarea consistenței lor. Valoarea limbajelor formale de specificare este pe deplin recunoscută, deoarece deși elaborarea specificațiilor necesită timp, costul crescut al fazei de specificare, fază timpurie a ciclului de viață a unui produs soft, este recuperat ulterior, întrucât costul proiectării, dezvoltării și întreținerii scad în mod semnificativ [Sch05].

Totodată se pot deriva programe corecte prin rafinare în pași succesivi [Dij75, Gri85, Mor90]. Rafinarea constă în aplicarea pas cu pas a unui set de reguli de rafinare. Procesul pornește de la un program abstract și se termină la obținerea codului. Codul astfel obținut își va respecta specificațiile, nefiind necesare testări ulterioare.

Există însă și metode formale utilizate în demonstrare automată și verificarea consistenței modelelor [Hei96, Rus95]. Acesta este nivelul superior de aplicare a metodelor formale.

Metodele formale s-au aplicat cu succes în specificare, verificare, demonstrare automată dar și în generarea automată de teste.

După cum s-a putut observa, pentru a asigura faptul că atât clientul cât și dezvoltatorul au înțeles același lucru din documentul de specificare a cerințelor, acesta ar trebui scris utilizând un limbaj formal. Din păcate nu întotdeauna acest lucru este posibil. Dacă totuși se dorește conversia cerințelor scrise în limbaj natural de către client în specificații formale, dezvoltatorii pot apela la ajutorul oferit de diverse unelte capabile să prelucreze limbajul natural.

Structura tezei. Teza este structurată în șase capitole (unul introductiv, unul care prezintă noțiuni fundamentale, trei conținând contribuții originale și ultimul concluzii), are o bibliografie care include un număr de 219 de referințe și cinci anexe.

Capitolul 1 introduce contextul, motivațiile și scopul tezei, rezumă contribuțiile din aceasta și furnizează o descriere pe scurt a ei.

Capitolul 2 prezintă succint noțiuni fundamentale de ingineria softului, metode formale, dar și legătura existentă între metodele de prelucrare a limbajului natural și metodele formale. Acest capitol conține și o scurtă prezentare a limbajului Z, dar și a metodei de rafinare în pași succesivi, singura care garantează obținerea de programe care să-și respecte specificația.

Capitolul 3 descrie contribuțiile în domeniul prelucrării limbajului natural. El este împărțit în trei subcapitole. În primul subcapitol sunt prezentate contribuțiile în dezambiguare și verificarea relației de implicație dintre două texte. Al doilea conține contribuțiile în domeniul segmentării unui text, iar ultimul cele de extragere a sumarului dintr-un text.

Capitolul 4 realizează trecerea de la limbaj natural la specificații formale, prezentând diferite aspecte ale utilizării ontologiilor în acest sens. Sunt cuprinse aici problema potrivirii unei ontologii la un text scris în limbaj natural, a extragerii unei ontologii dintr-un text scris în limbaj natural utilizându-se arborele de analiză sintactică a unei propoziții. De asemenea este analizat rolul dezambiguării elementelor dintr-o ontologie și utilizării similarității pentru a îmbunătăți calitatea unei ontologii. În final este folosită o ontologie pentru a prelucra cerințele cu scopul de a obține structuri cât mai apropiate de specificații formale.

Capitolul 5 prezintă diferite contribuții în domeniul metodelor formale. Printre acestea se numără o aplicație care asistă elaborarea specificațiilor Z, operarea cu ele și transformarea lor în programe abstracte, o aplicație care asistă procesul de rafinare a programelor abstracte la cod, respectiv o aplicație care simplifică structurile condiționale și câteva considerente cu privire la reutilizarea specificațiilor formale și multiformalism.

Capitolul 6 furnizează concluziile studiilor desfășurate pe parcursul elaborării tezei.

Cuvinte cheie: specificații, metode formale, corectitudine, prelucrarea limbajului natural, ontologii.

Contribuțiile originale introduse în această teză sunt conținute în capitolele 2, 3, 4 și 5 și sunt următoarele:

- O dezbateră a aplicabilității metodelor de prelucrare a limbajului natural în prelucrarea cerințelor [Mih08c] (subcapitolul 2.4).
- Un nou algoritm de dezambiguare [Tăt07a] (subcapitolul 3.1.1).
- Un nou algoritm de dezambiguare și identificare a părții de vorbire a cuvintelor [Mih07] (subcapitolul 3.1.2).
- Trei metode direcționale de evaluare a relației de consecință logică [Tăt07b, Tăt09a] (subcapitolul 3.1.3).
- O metodă de folosire a relației de implicație logică dintre două texte pentru obținerea unei ontologii [Mih08b] (subcapitolul 3.1.4).
- Trei metode logice de segmentare a unui text [Tăt08e, Tăt08b] (subcapitolul 3.2.1).
- Alte două variante ale algoritmului de segmentare logică [Mih08d] (subcapitolul 3.2.2).
- O metodă de segmentare genetică cu număr predefinit de segmente [Mih08a] (subcapitolul 3.2.3).
- O metodă de segmentare top-down bazată pe lanțuri lexicale [Tăt08b, Tăt08c, Tăt08d, Tăt09b] (subcapitolul 3.2.4).
- O metodă de sumarizare post segmentare [Tăt08e, Tăt08b] (subcapitolul 3.3.1).

- Trei metode de extragere a unui sumar de lungime variabilă [Tăt08d] (subcapitolul 3.3.2).
- O metodă de sumarizare pe lanțuri lexicale [Tăt08b, Tăt08c, Tăt08d, Tăt09b] (subcapitolul 3.3.3).
- Mai multe metrice pentru evaluarea gradului în care o ontologie corespunde unui text de limbaj natural [Mih10b] (subcapitolul 4.1).
- Două metode de prelucrare a arborelui de analiză sintactică a unei propoziții și de extragere automată a tripletelor ontologice [Mih10c] (subcapitolul 4.2).
- Un algoritm de identificare a elementelor plasate greșit dintr-o ontologie [Mih10d] (subcapitolul 4.3).
- O dezbateră cu privire la necesitatea utilizării dezambiguării în evaluarea ontologiilor [Mih10e] (subcapitolul 4.4).
- O metodă de extragere a specificațiilor formale dintr-un text reprezentând cerințele [Mih10f] (subcapitolul 4.5).
- Un studiu privind continuitatea activității de asigurare a calității softului [Șer05] (subcapitolul 5.1).
- O dezbateră pe marginea necesității abordării multiformale a specificării produselor soft complexe [Cio04] (subcapitolul 5.2).
- O analiză a refolosirii specificațiilor formale [Mih05] (subcapitolul 5.3).
- O unealtă care permite rafinarea asistată a schemelor Z [Mih10a] (subcapitolul 5.4).
- O unealtă care asistă procesul de rafinare a programelor abstracte în programe cod [Mih06a] (subcapitolul 5.5).
- O unealtă care asistă în mod automat simplificarea expresiilor condiționale [Mih06b, Lup08, Lup09] (subcapitolul 5.6).

II. Metode formale în dezvoltarea programelor

2.1. Ciclul de viață a unui produs soft

Procesul de dezvoltare a softului parcurge un „ciclu de viață” caracterizat de trei momente importante, numite faze ale ciclului: definiția, dezvoltarea și exploatarea.

2.2. Specificarea cerințelor

Specificarea cerințelor este ultima fază a procesului de analiză a cerințelor. Rezultatul analizei cerințelor este **Documentul de specificare a cerințelor softului** care reprezintă contractul client-dezvoltator.

2.2.1 Specificații informale

În multe proiecte de dezvoltare, documentul de specificare constă în pagină după pagină scrise în limbaj natural, care însă este ambiguu.

2.2.2 Tehnici semiformale de specificare

Există mai multe tipuri de tehnici semiformale de specificare: **tehnici grafice orientate pe acțiuni** [DeM78, Gan79, You79], **tehnici grafice orientate pe date** [Che76] și **alte tehnici semiformale de specificare a cerințelor** (PLS/PLA [Tei77], SADT [Ros85], SREM [Alf85]).

2.2.3 Tehnici formale de specificare

Dintre tehnicile formale de specificare se pot aminti: **Mașini cu stări finite** [Kam87], **Rețele Petri** [Pet62] și **limbajul Z** [Abr80]. O specificație Z este compusă din patru părți: mulțimea datelor de intrare, tipurile datelor și constantele, definiții de stări, starea inițială și operațiile. **Alte tehnici formale** sunt: Anna [Luc85], Gist [Bal85], VDM [Jon86] și CSP [Hoa85].

2.3. Rafinarea în pași succesivi

Primul care a propus o metodă care asigură corectitudinea programelor a fost Floyd [Flo67]. În același timp s-a considerat mult mai important să se scrie programe corecte prin construcție. În această direcție un rol deosebit l-au avut Dijkstra [Dij75], Gries [Gri81], Dromey [Dro89] și Morgan [Mor90]. Ideea lui Dijkstra, de a deriva în mod formal programe din specificații a fost continuată apoi de Gries [Gri81], care consideră că este mai important să se dezvolte programe corecte decât să se demonstreze ulterior corectitudinea lor [Fre06]. Pentru a se obține programe corecte, specificațiile trebuie transformate în cod folosind un set de reguli bine definite.

2.4. Metode de procesare a limbajului natural utilizate în analiza cerințelor

La conceperea unui produs nou, clienții și dezvoltatorii se întâlnesc pentru a specifica cerințele noului sistem. Pentru a nu omite anumite cerințe, se folosesc diferite **tehnici de colectare a cerințelor**. Astfel, se generează un volum mare de cerințe exprimate în limbaj natural. Prelucrarea lor poate beneficia de pe urma folosirii unor **mecanisme de prelucrare a limbajului natural**, după cum se poate observa din figura 1:

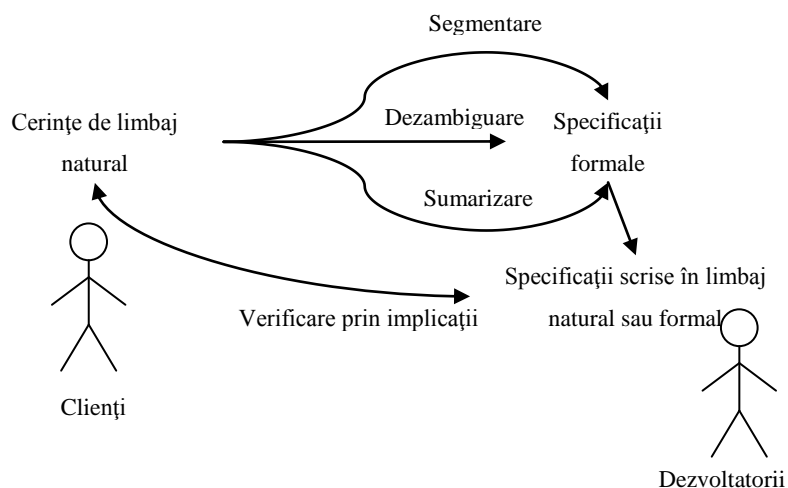


Figura 1. Schema utilizării mecanismelor de prelucrare a limbajului natural în analiza cerințelor

2.5. Ontologiile

O ontologie este „O specificare explicită a unei conceptualizări” [Gru93]. Astăzi coloana vertebrală a Internetului Semantic se consideră a fi OWL (Web Ontology Language) și RDF (Resource Description Framework) [Pol09]. Ele reprezintă cele mai bune limbaje de modelare a unei ontologii. Cel mai general mod de definire a unei ontologii este utilizând triplete [Bra85]. Un triplet e compus dintr-un *concept*, un *predicat* care reprezintă o relație direcțională spre *obiect*, care de cele mai multe ori este o caracteristică a conceptului. Pentru uniformitate, toate componentele ontologiei conceptele, relațiile și obiectele sunt reprezentate de un URI (Uniform Resource Identifiers) [Ber98].

III. Contribuții în domeniul prelucrării limbajului natural

3.1. Contribuții în dezambiguarea automată a unui text

Pentru a putea dezambigua rapid specificațiile scrise în limbaj natural de către clienți, este nevoie de o unealtă capabilă de a identifica automat sensul cuvintelor. Cum este de așteptat, când se urmărește identificarea sensului corect al cuvintelor, cea mai de încredere sursă este dicționarul.

3.1.1 Identificarea sensului cuvintelor

Una dintre cele mai cunoscute metode de dezambiguare bazată pe dicționar este cea a lui Lesk [Les86]. În [Ban03] algoritmul lui Lesk a fost dezvoltat utilizând dicționarul de limbă engleză WordNet [***WNT]. Algoritmul dezvoltat de Banerjee și Pedersen adaugă definiției din WordNet sinonime, exemple și definiții ale cuvintelor înrudite. Totodată se introduce o nouă măsură „overlap” care favorizează intersecțiile multiple de cuvinte.

Algoritmul original CHAD [Tăt07a] seamănă cu algoritmul prezentat în [Ban03], doar că dezambiguesază toate cuvintele unui text, printr-o singură parcurgere a textului și nu doar un cuvânt țintă. Algoritmul se aplică în mod repetat pe grupuri de trei cuvinte, dintre care primele două au fost deja dezambiguate, iar al treilea urmează să fie dezambiguit. Astfel putem privi algoritmul ca pe un lanț, ale cărui ochiuri sunt reprezentate de grupurile de trei cuvinte. Primele trei cuvinte sunt însă dezambiguate simultan.

Algoritmul folosește trei măsuri pentru calculul scorului tripletului de sensuri $s_{w_1}^i, s_{w_2}^j, s_{w_3}^k$ al celor trei cuvinte w_1, w_2, w_3 : Dice [Dic45], overlap și Jaccard [Jac901]. Aceste trei măsuri evaluează în mod diferit numărul de cuvinte comune existente în definițiile corespunzătoare sensurilor $s_{w_1}^i, s_{w_2}^j, s_{w_3}^k$ ale celor trei cuvinte w_1, w_2 , respectiv w_3 , definiții notate cu $D_{w_1}, D_{w_2}, D_{w_3}$:

$$\text{Dice: scor}(s_{w_1}^i, s_{w_2}^j, s_{w_3}^k) = 3 * \frac{|D_{w_1} \cap D_{w_2} \cap D_{w_3}|}{|D_{w_1}| + |D_{w_2}| + |D_{w_3}|}$$

$$\text{overlap: scor}(s_{w_1}^i, s_{w_2}^j, s_{w_3}^k) = \frac{|D_{w_1} \cap D_{w_2} \cap D_{w_3}|}{\min(|D_{w_1}|, |D_{w_2}|, |D_{w_3}|)}$$

$$\text{Jaccard: scor}(s_{w_1}^i, s_{w_2}^j, s_{w_3}^k) = \frac{|D_{w_1} \cap D_{w_2} \cap D_{w_3}|}{|D_{w_1} \cup D_{w_2} \cup D_{w_3}|}$$

Incluzând subalgoritmul de dezambiguare a unui triplet de cuvinte, se obține algoritmul CHAD prezentat în figura 2.

Deoarece deseori intersecția definițiilor celor trei cuvinte este egală cu 0, există cuvinte ale căror scoruri sunt egale cu 0 pentru toate sensurile lor. Acestor cuvinte li s-a atribuit primul sens din WordNet, dat fiind faptul că în WordNet sensurile cuvintelor sunt ordonate după frecvența lor.

Algoritmul a fost testat pe zece texte alese în mod aleator din corpusul Brown [***BrC, Nel79] (A01, A02, A11, A12, A13, A14, A15, B13, B20 și C01) care au deja identificate părțile de vorbire ale cuvintelor. Astfel rezultatele au putut fi comparate cu cele

din corpusul SemCor [***Sem], unde pe lângă partea de vorbire, este indicat și sensul cuvântului. Algoritmul CHAD a fost rulat prima dată pentru substantive, apoi doar pentru verbe și a treia oară pentru substantive, verbe, adjective și adverbe.

Algoritmul CHAD (W, n, S) este:

date W {un text format din cuvintele w_1, w_2, \dots, w_n }, n {numărul cuvintelor din text}

@Dezambiguesă_triplet ($w_1, w_2, w_3, s_{w_1}^*, s_{w_2}^*, s_{w_3}^*$)

Pentru $i \leftarrow 4, n$ **execută**

$p \leftarrow 1$

$max \leftarrow scor(s_{w_{i-2}}^*, s_{w_{i-1}}^*, s_{w_i}^1)$

{sensurile identificate deja pentru cuvintele w_{i-2} și w_{i-1} sunt $s_{w_{i-2}}^*, s_{w_{i-1}}^*$ }

Pentru fiecare sens $s_{w_i}^m$ **execută**

Dacă $max < scor(s_{w_{i-2}}^*, s_{w_{i-1}}^*, s_{w_i}^m)$

atunci $p \leftarrow m$

$max \leftarrow scor(s_{w_{i-2}}^*, s_{w_{i-1}}^*, s_{w_i}^m)$

sfd

sfp

$s_{w_i}^* \leftarrow s_{w_i}^p$

sfp

rezultate S {sensurile $s_{w_i}^*, i = 1, n$ }

sf CHAD

Figura 2. Algoritmul CHAD

Chiar dacă precizia nu depășește precizia primului sens din WordNet (diferența medie pentru substantive fiind 0,0338, pentru verbe de 0,0340, iar pentru toate cuvintele 0,0491), aceste rezultate se încadrează în rezultatele obținute până acum unde la a doua competiție Senseval [***Sen], doar două din șapte echipe (cu metode nesupervizate) au obținut rezultate mai bune față de primul sens din WordNet.

Algoritmul CHAD nu este dependent de limba în care e scris textul, el putând fi aplicat cu succes și pentru limba română. Acest algoritm poate fi aplicat pentru a îmbunătăți traducerea [Tăt07a], dar și pentru a identifica părțile de vorbire [Mih07] sau pentru a verifica relația de implicație dintre două texte.

3.1.2 Identificarea părții de vorbire a cuvintelor

În momentul în care o persoană întâlnește un cuvânt necunoscut, mai ales într-o limbă pe care nu o cunoaște suficient a cărei gramatică nu și-a însușit-o încă, va încerca să caute cuvântul respectiv în dicționar, să-i afle sensul corespunzător contextului. Dar unele cuvinte constituie mai multe părți de vorbire diferite, depinzând de context. În cazul în care cuvântul căutat are astfel de caracteristici, i se va identifica și partea de vorbire corespunzătoare contextului. Acest fapt a stat la baza unui nou algoritm de identificare simultană a părților de vorbire și a sensului cuvintelor derivat din algoritmul CHAD aflat în figura 3 [Mih07].

Experimentele s-au efectuat pe aceleași zece texte din corpusul Brown. Deoarece în SemCor, care de asemenea a fost folosit ca standard pentru evaluare, există cuvinte a căror parte de vorbire nu este identificată (Notag), acele cuvinte nu au fost considerate în calculul preciziei metodei. Valorile minime, maxime și medii se află în tabelul 1.

Metoda este însă promițătoare în ciuda rezultatelor care se află sub rezultatele obținute utilizând mai multe resurse și reguli gramaticale, deoarece ea poate fi aplicată pentru mai multe limbi și se realizează în același timp și dezambiguarea cuvintelor.

Algoritmul PV_CHAD (W, n, S, P) este
date W {un text format din cuvintele w_1, w_2, \dots, w_n }, n {numărul de cuvinte din text}
@Dezambiguesă_și_identifică_PV_triplet ($w_1, w_2, w_3, S_{1 p_1}^{*w_1}, S_{2 p_2}^{*w_2}, S_{3 p_3}^{*w_3}$)
Pentru $i \leftarrow 1, n$ **execută**
 $r \leftarrow$ prima parte de vorbire din dicționar corespunzătoare cuvântului w_i
 $k \leftarrow 1$ {primul sens al primei părți de vorbire a cuvântului w_i }
 $max \leftarrow scor(S_{i-2 p_{i-2}}^{*w_{i-2}}, S_{i-1 p_{i-1}}^{*w_{i-1}}, S_{i r}^{w_i})$
Pentru fiecare parte de vorbire q a cuvântului w_i **execută**
Pentru fiecare sens $S_{i q}^m$ al cuvântului w_i cu partea de vorbire q **execută**
Dacă $max < scor(S_{i-2 p_{i-2}}^{*w_{i-2}}, S_{i-1 p_{i-1}}^{*w_{i-1}}, S_{i q}^m)$
atunci $r \leftarrow q$
 $k \leftarrow m$
 $max \leftarrow scor(S_{i-2 p_{i-2}}^{*w_{i-2}}, S_{i-1 p_{i-1}}^{*w_{i-1}}, S_{i q}^m)$
sfd
sfp
sfp
 $p_i^* \leftarrow r$
 $S_{i p_i}^{*w_i} \leftarrow S_{i r}^k$
sfp
rezultate S {sensurile $s_i^*, i \leftarrow 1, n$ }, P {părțile de vorbire $p_i^*, i \leftarrow 1, n$ }
sf PV_CHAD

Figura 3. Algoritmul PV_CHAD

Precizia	Valoarea minimă	Valoarea maximă	Valoarea medie
Identificarea tuturor părților de vorbire	20,16%	64,13%	42,78%
Dezambiguarea tuturor cuvintelor	7,00%	41,43%	21,78%
Substantive	53,04%	82,83%	66,12%
Verbe	67,83%	92,02%	84,45%
Adjective	29,55%	75,34%	56,39%
Adverbe	46,87%	80,65%	69,41%

Tabelul 1. Valorile minimă, maximă și medie ale preciziei în funcție de părțile de vorbire.

Acest algoritm poate fi folosit cu succes pentru identificarea automată a părților de vorbire din cerințe, astfel se pot identifica mai ușor entitățile și relațiile pentru o elaborare a unui model Entitate-Relație.

3.1.3 Relația de consecință logică dintre două texte

Principalele raționamente care se efectuează relativ la un text sau un grup de texte se bazează pe demonstrarea sau verificarea relației de consecință logică sau implicația dintre două texte. Relația de implicație logică dintre două texte, T textul și H ipoteza, reprezintă un fenomen fundamental al limbajului natural. Ea se notează cu $T \rightarrow H$ și înseamnă ceea ce se înțelege din H poate fi dedus din T .

Chiar dacă problema nu e nouă, majoritatea abordărilor automate sunt rezultatul direct al competițiilor Pascal Textual Entailment Challenges [***RTE]. Cu toate acestea,

doar puțini autori au folosit caracterul direcțional al acestei relații: dacă are loc $T \rightarrow H$ este puțin probabil să aibă loc și $H \rightarrow T$. Dintre metodele propuse la competiția RTE-1, metoda cu cele mai bune rezultate a fost cea direcțională a lui Glickman [Gli05] 58,5%, iar metoda direcțională a lui Kõuylekov [Kõu06] a obținut 56%.

Prima metodă de evaluare a relației de consecință logică propusă în [Tăt07b] utilizează similaritatea textelor. În [Cor05] se definește următoarea măsură a similarității dintre două texte, similaritatea dintre textele T_i și T_j în raport cu T_i :

$$sim(T_i, T_j)_{T_i} = \frac{\sum_{pos} (\sum_{w_k \in WS_{pos}^{T_i}} (\max Sim(w_k) \times idf_{w_k}))}{\sum_{pos} \sum_{w_k \in WS_{pos}^{T_i}} idf_{w_k}} \quad (3.1.3.1)$$

Mulțimea cuvintelor din textul T_i care au o anumită parte de vorbire pos se notează cu $WS_{pos}^{T_i}$. Pentru un cuvânt w_k cu o anumită parte de vorbire pos în T_i , cea mai mare similaritate a cuvintelor cu aceeași parte de vorbire pos din celălalt text T_j , se notează cu $\max Sim(w_k)$. Deoarece ipoteza H conține mai puțină informație față de textul T , pentru o implicație între cele două texte va avea loc următoarea relație [Tăt09a]:

$$sim(T, H)_T < sim(T, H)_H \quad (3.1.3.2)$$

Această formulă este aplicată textelor dezambiguate cu algoritmul CHAD prezentat în subcapitolul anterior. În formula (3.1.3.1) s-au folosit doar substantivele și verbele ca și părți de vorbire pos , iar similaritatea dintre două cuvinte s-a considerat a fi 1 dacă cele două cuvinte sunt identice sau sunt sinonime și 0 altfel. Totodată s-au folosit câteva reguli de identificare a falselor implicații, notate cu *COND*.

Pentru a se ușura scrierea algoritmului de verificare a relației de consecință logică bazat pe similaritatea textelor, vom nota textul T cu T_1 , iar ipoteza H cu T_2 și vom nota cu: NP_1 și NP_2 – numele proprii din T_1 respectiv T_2 , I_c – cuvintele comune care nu sunt nume proprii din T_1 și T_2 , $SYN(T_1)_{T_2}$ – cuvinte din T_1 care nu sunt nume proprii, nu apar în T_2 , dar sunt substantive sau verbe sinonime cu unele cuvinte din T_2 , analog $SYN(T_2)_{T_1}$, $C_1 = |SYN(T_1)_{T_2}|$, $C_2 = |SYN(T_2)_{T_1}|$, $W_{T_1} = NP_1 \cup I_c$, $W_{T_2} = NP_2 \cup I_c$.

Din (3.1.3.1) și (3.1.3.2) se obține următoarea condiție: $C_1 \leq C_2$, adică $|SYN(T_1)_{T_2}| \leq |SYN(T_2)_{T_1}|$. Relația nu este strictă datorită definiției mulțimilor $SYN(T_1)_{T_2}$ și $SYN(T_2)_{T_1}$.

Funcția Consecință_logică (T_1, T_2) este
date T_1, T_2 {două texte}
dacă $W_{T_1} \subseteq W_{T_2}$
atunci dacă $T_2 = NP_2 \cup I_c$
atunci dacă @ se verifică toate regulile din *COND*
atunci Consecință_logică \leftarrow False {nu are loc $T_1 \rightarrow T_2$ }
altfel Consecință_logică \leftarrow True {cazul I: $T_1 \rightarrow T_2$ }
altfel dacă $C_1 \leq C_2$
atunci Consecință_logică \leftarrow True {cazul II: $T_1 \rightarrow T_2$ }
altfel Consecință_logică \leftarrow False {nu are loc $T_1 \rightarrow T_2$ }
altfel Consecință_logică \leftarrow False {nu are loc $T_1 \rightarrow T_2$ }
rezultate True {are loc $T_1 \rightarrow T_2$ }, respectiv False {nu are loc $T_1 \rightarrow T_2$ }
sf Consecință_logică

Figura 4. Funcția de verificare a relației de consecință logică pe baza similarităților

Un caz important este acela în care T_2 conține doar nume proprii care apar și în T_1 , de aceea algoritmul de verificare a relației de consecință logică pe baza similarităților va începe cu verificarea acestei condiții, vezi figura 4.

Aplicația construită pornind de la acest algoritm necesită ca intrare texte în care s-au identificat părțile de vorbire, inclusiv numele proprii. Algoritmul a fost evaluat folosind baza de 800 de perechi de texte din competiția RTE-1 [***RTE].

Precizia globală a algoritmului propus este de 52,625%, iar cea medie 56,60647%.

A doua metodă de evaluare a relației de consecință logică se bazează pe măsura cosinus [Tăt09a]. S-au definit trei măsuri cosinus pentru a evalua distanța dintre vectorii de cuvinte $T = t_1, t_2, \dots, t_m$ și $H = h_1, h_2, \dots, h_n$. Prima măsură cosinus, $\cos_T(T, H)$ măsoară distanța dintre vectorul m -dimensional $\vec{T} = (1, 1, \dots, 1)$ și vectorul m -dimensional \vec{H} care reprezintă proiecția lui H pe T , adică $\vec{H}_i = 1$, dacă t_i este cuvânt din H și $\vec{H}_i = 0$ altfel. A doua măsură cosinus, $\cos_H(T, H)$ măsoară distanța dintre vectorul n -dimensional $\vec{H} = (1, 1, \dots, 1)$ și vectorul n -dimensional \vec{T} care reprezintă proiecția lui T pe H , adică $\vec{T}_i = 1$, dacă h_i este cuvânt din T și $\vec{T}_i = 0$ altfel. Pentru ultima măsură cosinus, $\cos_{H \cup T}(T, H)$, primul vector se obține înlocuind cuvintele din $H \cup T$ conținute de T cu 1 și restul cu 0, iar al doilea prin înlocuirea cuvintelor din $H \cup T$ conținute de H cu 1 și cele din $T \setminus H$ cu 0.

Dacă notăm cu c numărul de cuvinte comune celor doi vectori, cele trei măsuri

$$\text{sunt: } \cos_T(T, H) = \sqrt{\frac{c}{m}}, \quad \cos_H(T, H) = \sqrt{\frac{c}{n}} \quad \text{și} \quad \cos_{H \cup T}(T, H) = \sqrt{\frac{4c^2}{(n+c)(m+c)}}.$$

Pentru a se îndeplini condiția: *T implică H dacă și numai dacă H conține mai puțină informație decât T* [Mon01], similaritățile dintre T și H calculate în raport cu T respectiv $H \cup T$ trebuie să fie foarte apropiate. În mod analog similaritățile dintre T și H calculate în raport cu H respectiv $H \cup T$ trebuie să fie foarte apropiate. De asemenea toate cele trei similarități (cosinusuri) trebuie să fie mai mari decât o anumită limită. Deci trebuie respectate următoarele condiții: $|\cos_{H \cup T}(T, H) - \cos_T(T, H)| \leq \tau_1$, $|\cos_H(T, H) - \cos_{H \cup T}(T, H)| \leq \tau_2$ și $\max\{\cos_T(T, H), \cos_H(T, H), \cos_{H \cup T}(T, H)\} \geq \tau_3$. Valorile identificate pentru cele trei limite sunt: $\tau_1 = 0,095$, $\tau_2 = 0,15$ și $\tau_3 = 0,7$.

În urma evaluării rezultatelor obținute s-a observat că precizia pentru perechile adevărate este 68,92230576%, iar pentru cele false este de 46,36591479%. Precizia globală este de 57,62%.

A treia metodă originală se bazează pe o variantă modificată a distanței Levenshtein [Tăt09a]. Notăm cu $LD(w_1, w_2)$ distanța Levenshtein modificată dintre cuvintele w_1 și w_2 . Această distanță se definește ca numărul minim de transformări (ștergeri, inserări și substituții) astfel încât cuvântul w_1 se transformă în cuvântul w_2 . Notăm cu T_{word} „cuvântul” obținut din propoziția T , unde spațiul se consideră a fi o literă nouă și toate cuvintele lui T sunt juxtapuse. În mod analog se obține „cuvântul” H_{word} . Se consideră că $LD(T_{word}, H_{word})$ reprezintă caracteristica informativă a lui H în raport cu T . Această distanță Levenshtein nu e o distanță reală, deoarece $LD(w_1, w_2) \neq LD(w_2, w_1)$.

Deoarece *T implică H dacă și numai dacă H conține mai puțină informație decât T*, atunci ar trebui să se verifice următoarea relație: $LD(T_{word}, H_{word}) < LD(H_{word}, T_{word})$. Aplicând acest criteriu celor 800 de perechi din baza de texte RTE-1, s-a obținut precizia globală de 53,19%.

Costurile transformării cuvântului w_1 în cuvântul w_2 sunt următoarele: transformarea literei mari în mici sau invers 1, costul de inserare 3, costul ștergerii 3, costul substituției 5, costul interschimbării 2.

Rezultatele obținute se încadrează între rezultatele obținute la ultimele competiții RTE, metoda a doua fiind cea mai bună. Pe lângă faptul că a decide dacă un text implică un alt text este de folos și în această formă dezvoltatorilor, această relație a fost folosită cu succes și pentru a segmenta și sumariza un text, care de asemenea ușurează sarcina de

analiză și prelucrare a specificațiilor scrise în limbaj natural, după cum se va putea observa din secțiunile următoare.

3.1.4 O simplă ontologie bazată pe relația de consecință logică dintre două texte

În [Mih08b] s-a prezentat un mod de a crea o ontologie informală bazată pe relația direcțională de consecință logică dintre texte. S-a folosit metoda cu cel mai bun scor de identificare a implicațiilor dintre texte, metoda cosinusurilor, prezentată în subcapitolul anterior.

Pornindu-se de la un text scris în limbaj natural, se verifică toate relațiile de consecință logică existente între propozițiile textului. Astfel se obține un graf direcțional între propozițiile textului. Graful poate reprezenta o ontologie informală simplă a textului respectiv în care entitățile sunt reprezentate de propoziții, iar relațiile dintre entități de relația de consecință logică.

Aceeași relație de consecință logică poate fi de asemenea folosită pentru a relaționa texte întregi, obținându-se un graf al relațiilor dintre texte. Dacă respectivele texte reprezintă definițiile entităților, atunci graful relațiilor dintre entități poate fi graful relațiilor de consecință logică.

3.2. Noi abordări ale segmentării unui text

În cazul implementării unui sistem informatic de mari dimensiuni volumul cerințelor va fi de asemenea mare. Astfel prelucrarea sa va necesita munca în paralel a mai multor persoane sau echipe. Bineînțeles se impune împărțirea specificațiilor după funcționalități, ceea ce corespunde cu operația de segmentare specifică limbajului natural.

3.2.1 Segmentare liniară

În [Tăt08e, Tăt09c] s-a propus o metodă de segmentare liniară a unui text folosindu-se relația de consecință logică prezentată în subcapitolul anterior. Cel mai simplu algoritm de segmentare utilizând implicațiile se numește „PureEntailment” (PE), unde un nou segment se consideră că începe în momentul în care ultimul segment nu mai implică noua propoziție.

O a doua abordare a segmentării, Logical TextTiling (LTT) [Tăt08e, Tăt09c], este asemănătoare cu algoritmul de detectare a structurii subtemei TextTiling (TT) [Hea93], dar LTT detectează structura logică a textului. Principalele diferențe sunt:

- În LTT unitățile de bază sunt propozițiile, iar în TT sunt pseudopropozițiile, toate având aceeași dimensiune predefinită;
- În TT fiecare limită dintre pseudopropoziții se scorează cu similaritatea pseudopropozițiilor de la $i-k$ la i față de pseudopropozițiile de la $i+1$ la $i+k+1$. În LTT se scorează propozițiile, fiecare propoziție s_i primind ca și scor $scor(s_i)$ numărul de propoziții din întregul text implicate de s_i . O propoziție importantă va avea un scor mare, ceea ce este în concordanță cu criteriul implicației textelor: *Un text T implică un alt text H dacă și numai dacă conține mai multă informație decât H* [Tăt07b].

În LTT o limită dintre două segmente, determinată de schimbarea subiectului discursului este cauzată de schimbarea configurației importanței logice a propozițiilor (vezi figura 5). Astfel există o bază fundamentală de descompunere a discursului într-o mulțime de segmente mai mici, fiecare cu o configurație diferită a importanței propozițiilor componente.

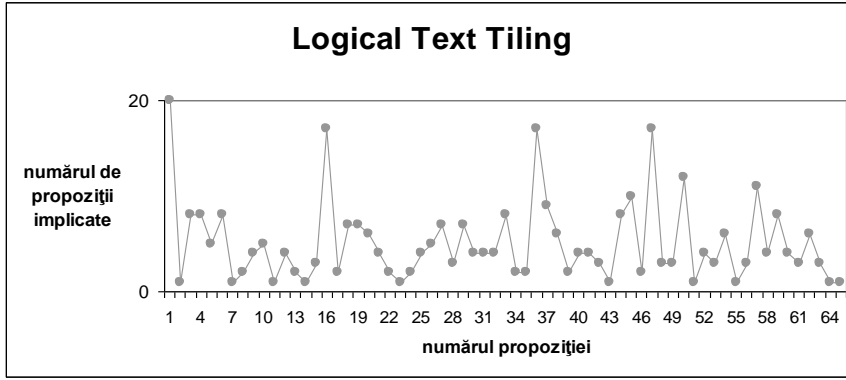


Figura 5. Structura logică a textului

Spre deosebire de alte metode de segmentare, LTT este o metodă liniară:

Algoritmul LTT(S, n, SC, SEG) este

date S {un text, s_i – propozițiile textului, $i \in \{1, n\}$, n {numărul de propoziții din text},

SC { s_i = scor(s_i) scorurile propozițiilor din text}

$j \leftarrow 1, p \leftarrow 1$, inițializează(Seg_j, s_1) {inițial segmentul Seg_j conține doar propoziția s_1 }

$dir \leftarrow$ „urcă”

Cât timp $p < n$ execută

Dacă $sc_{p+1} > sc_p$

atunci Dacă $dir =$ „urcă”

atunci adaugă(Seg_j, s_{p+1}) {adaugă segmentului Seg_j propoziția s_{p+1} }

altfel $j \leftarrow j+1$

inițializează(Seg_j, s_{p+1})

$dir \leftarrow$ „urcă”

sf

altfel adaugă(Seg_j, s_{p+1})

Dacă $scor(s_{p+1}) < scor(s_p)$

atunci $dir \leftarrow$ „coboară”

sf

sf

$p \leftarrow p + 1$

sf

rezultate $SEG = \{Seg_1, \dots, Seg_j\}$

sf LTT

Având ca punct de pornire măsurarea importanței unei propoziții în funcție de numărul de alte propoziții cu care aceasta se află în relație (sunt implicate de ea), s-au propus două noi variante de determinare a scorului propozițiilor:

$$ArcInt_k = \sum_{i=1}^k \sum_{j=k, i \neq j}^n implică_{i,j} + \sum_{i=k}^n \sum_{j=i, i \neq j}^k implică_{i,j} \quad \text{unde}$$

$$ArcReal_k = \sum_{i=1}^k \sum_{j=k, i \neq j}^n \frac{implică_{i,j}}{|i-j|} + \sum_{i=k}^n \sum_{j=i, i \neq j}^k \frac{implică_{i,j}}{|i-j|} \quad implică_{i,j} = \begin{cases} 1, & \text{dacă } s_i \rightarrow s_j \\ 0, & \text{altfel} \end{cases}$$

O altă abordare a problemei segmentării este cea a aplicării unui algoritm de programare dinamică. Această metodă realizează mai întâi o sumarizare (vezi subcapitolul 3.3.1). Limitele dintre segmente sunt acele propoziții care au un scor minimal și sunt situate între două propoziții din sumar (vezi figura 6). Astfel numărul de segmente obținute este corelat cu lungimea x a sumarului.

Pentru a evalua metodele de segmentare propuse s-a utilizat textul literar „The Koan” [Ric91], text care a fost segmentat, sumarizat și i s-au rezolvat anaforele în mod manual, pe baza criteriilor lingvistice de către un specialist [Tăt08e].

Conform cu [Hea93], compararea cu segmentarea manuală constă în a evalua:

- câte dintre limitele găsite de metodă sunt limite manuale exacte (sau sunt foarte apropiate de ele) (precizia);
- câte dintre limitele adevărate au fost identificate (recall-ul).

Totodată s-a folosit și măsura $F = \frac{2 * Precizia * Recall}{Precizia + Recall}$.

Algoritmul Segmentare_post_sumarizare($S, n, SC, P, x, SEG, nrseg$) este:

date S {un text, s_i – propozițiile textului, $i \leftarrow 1, n$ }, n {numărul de propoziții din text},
 SC { $sc_i = scor(s_i)$ scorurile propozițiilor din text}, P { p_i – pozițiile din text ale
propozițiilor sumarului, $i \leftarrow 1, x$ }, x {numărul de propoziții din sumar}

$nrseg \leftarrow 1$

$begpoz \leftarrow 1$

Pentru $i \leftarrow 1, x-1$ **execută**

$endpoz \leftarrow \arg \min_{j, j \in \{p_i, p_{i+1}-1\}} scor(s_j)$

@ Seg_{nrseg} începe cu $begpoz$ și se termină pe $endpoz$

$nrseg \leftarrow nrseg + 1$

$begpoz \leftarrow endpoz + 1$

sfp

rezultate SEG { $Seg_1, \dots, Seg_{nrseg}$ – lista segmentelor}, $nrseg$ {numărul segmentelor}

sf Segmentare_post_sumarizare

Figura 6. Algoritmul Segmentare post sumarizare

Câteva dintre rezultatele evaluării comparative a diferitelor metode de segmentare (Logică/Programare dinamică) se află în figura 7. Dintre metodele logice, cele mai bune rezultate le-a obținut LTT.

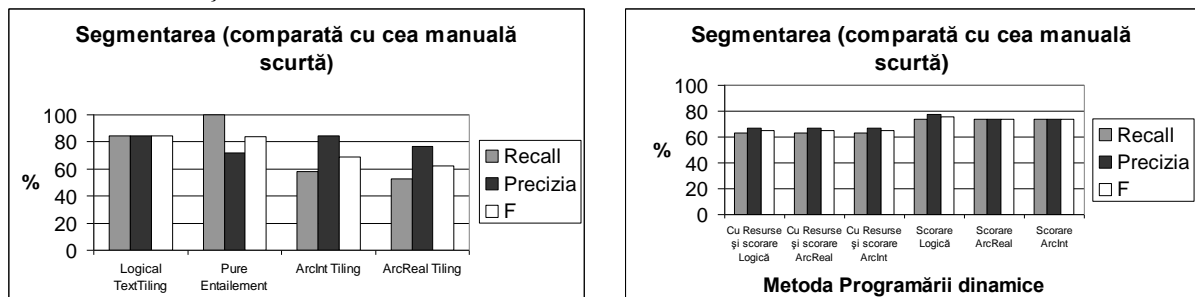


Figura 7. Evaluarea segmentării

3.2.2 Alte variante de segmentare Logical TextTiling

Metoda Logical TextTiling necesită o scorare în prealabil a fiecărei propoziții, ceea ce conduce la o complexitate de ordinul n^2 . Pentru a reduce complexitatea s-a propus un LTT pe vecinătăți, în care fiecare propoziție va fi scorată într-o vecinătate de dimensiune predefinită, nu din întregul text și respectiv un LTT ponderat care face prioritare implicațiile locale pentru a aplatiza grafurile care reprezintă structura logică a textului [Mih08d]. Structura logică a textului obținută în prima variantă generează aceleași segmente ca și LTT-ul pentru o vecinătate de minim 20% din lungimea textului inițial (pentru textul „The Koan” [Ric91]). Și structura logică a textului obținută din a doua variantă generează aceleași segmente ca și LTT-ul. În consecință cele două variante de LTT pot înlocui cu succes metoda originală LTT.

3.2.3 Segmentare genetică cu număr predefinit de segmente bazată pe implicații

Metoda următoare folosește interpretarea logică a textului, adică relația de consecință logică dintre două propoziții pentru a realiza acest scop: o coeziune internă ridicată a segmentelor și o conexiune scăzută între segmente vecine [Mih08a].

Dacă pentru metodele anterioare nu se poate impune numărul de segmente care se vor obține, metoda propusă în cele ce urmează va împărți textul într-un număr impus k de segmente. Prin urmare dacă avem textul $T = \{t_1, t_2, \dots, t_n\}$ care conține n propoziții, a-l împărți în k segmente înseamnă a împărți mulțimea $\{1, 2, \dots, n\}$ în k submulțimi, $\{1, 2, \dots, n_1\}$, $\{n_1+1, \dots, n_2\}$, ..., $\{n_{k-1}+1, n_{k-1}+2, \dots, n_k\}$. Această partiționare a mulțimii inițiale poate fi reprezentată ca un vector $b = \{n_1, \dots, n_k\}$, unde $n_k = n$, iar n_1, \dots, n_k reprezintă indicii ultimelor propoziții din cele k segmente.

Primul pas al metodei de segmentare este construirea matricii implicațiilor. Textului T i se va asocia o matrice E de dimensiune $n \times n$, care este definită după cum urmează:

$$e_{i,j} = \begin{cases} 1, & \text{daca propozitia } t_i \text{ implica propozitia } t_j, i, j = \overline{1, n} \\ 0, & \text{altfel} \end{cases} \quad (3.2.3.a)$$

o a doua definiție a aceleiași matrici este:

$$e_{i,j} = \begin{cases} \frac{1}{|i-j|+1}, & \text{daca propozitia } t_i \text{ implica propozitia } t_j, i, j = \overline{1, n} \\ 0, & \text{altfel} \end{cases} \quad (3.2.3.b)$$

Implicațiile sunt calculate utilizând măsura \cos , metodă prezentată anterior în subcapitolul 3.1.3.

Fiecare segmentare va fi evaluat separat, din punctul de vedere al coeziunii:

$$\text{scor}(\text{segment}_{i,j}) = \sum_{k=i}^j \sum_{l=i}^j e_{k,l} \quad (3.2.3.1)$$

normalizarea scorului:

$$\text{scor}(\text{segment}_{i,j}) = \frac{\sum_{k=i}^j \sum_{l=i}^j e_{k,l}}{(|i-j|+1)^2} \quad (3.2.3.1')$$

și varianta în care s-au eliminat auto-implicațiile:

$$\text{scor}(\text{segment}_{i,j}) = \frac{\sum_{k=i}^j \sum_{l=i}^j e_{k,l} - (|i-j|+1)}{(|i-j|+1)^2} \quad (3.2.3.1'')$$

Din punctul de vedere al implicațiilor, două segmente vecine sunt deconectate unul față de celălalt dacă între propozițiile lor sunt cât mai puține implicații:

$$\text{scor}(\text{segment}_{i,j}, \text{segment}_{j+1,k}) = \sum_{l=i}^j \sum_{m=j+1}^k e_{l,m} \quad (3.2.3.2)$$

Suma scorurilor locale (3.2.3.2) trebuie minimizată pentru a obține o segmentare cu o interconectivitate mică, iar suma scorurilor locale (3.2.3.1), (3.2.3.1') respectiv (3.2.3.1'') trebuie maximizată, pentru a obține o segmentare în care segmentele au o coeziune internă ridicată.

Problema determinării unui număr predefinit de segmente poate fi privită ca o problemă combinatorială, iar cei mai potriviți algoritmi pentru astfel de probleme sunt algoritmi genetici. Algoritmul propus pornește de la o populație de cromozomi generați aleator, care reprezintă posibile segmente de dimensiune k . Pentru a genera următoarea generație, se aplică următorii operatori genetici: „binary tournament selection” pentru selecția părinților și „one point crossover” și „flip mutation” [Bac00] pentru încrucișarea și mutația care generează noii cromozomi. Procesul de evaluare este similar cu schema de evaluare a unui algoritm genetic standard. În plus se va folosi o selecție a celor mai buni indivizi. Soluția problemei de segmentare este reprezentată în algoritmul genetic numit GATTS. Algoritmul GATTS a folosit următorii parametri: dimensiunea populației de 200, 50 de generații, probabilitatea de încrucișare de 0,7 și cea de mutație de 0,1.

Evaluarea s-a realizat pe două dintre textele folosite anterior, un text literar, „The Koan” [Ric91] și un articol de ziar, Hirst [Mor91] din care s-au eliminat propozițiile scurte (cu mai puțin de 5 cuvinte inclusiv) celelalte 41 de propoziții rămase fiind renumerotate. Rezultatele au fost evaluate comparativ cu cele obținute prin metoda Programării dinamice (DP) prezentată anterior, deoarece ambele metode obțin un număr predefinit de segmente.

Rezultatele medii obținute sunt comparabile cu cele mai bune rezultate ale metodei programării dinamice, iar cele maxime le depășesc, unele fiind chiar egale cu 100%, ceea ce indică faptul că s-a identificat chiar segmentarea manuală.

Se observă că pentru Koan cele mai bune rezultate se obțin pentru varianta (3.2.3.a) a matricii E , iar dintre definițiile scorurilor locale (3.2.3.1”) și (3.2.3.2) conduc la cele mai bune rezultate. În ceea ce privește matricea (3.2.3.b), combinația ei cu scorul (3.2.3.1’) pare a fi cea mai potrivită. Pentru textul Hirst, matricea (3.2.3.b) conduce la rezultatele cele mai bune și din nou se evidențiază scorurile (3.2.3.1”) și (3.2.3.2).

Ca o direcție viitoare de cercetare, se va căuta un mod de a corela cele două scoruri locale, (3.2.3.1”) și (3.2.3.2), pentru a îmbunătăți metoda de segmentare.

3.2.4 O nouă metodă de segmentare top-down bazată pe lanțuri lexicale

O metodă de segmentare care oferă rezultate mai bune decât metoda LTT este cea top-down bazată pe lanțuri lexicale. Un lanț lexical este o secvență de cuvinte astfel încât sensul fiecărui cuvânt din secvență poate fi obținut precis pornind de la sensurile celorlalte cuvinte [Mor91, Bar99, Har97, Sil02]. Toate aceste lanțuri lexicale furnizează o reprezentare a structurii lexicale de coeziune a textului. În aceste abordări lanțurile lexicale au fost obținute folosind o abordare bottom-up, luându-se fiecare cuvânt candidat al textului și identificând relațiile pe care acest cuvânt le are cu vecinii lui, pe baza unui tezaur cum este Rodget-ul [Mor91] sau WordNet-ul [Bar99, Sto04]. Dacă există astfel de relații, cuvântul respectiv este inserat cu sensul corespunzător în lanțul curent, iar sensurile celorlalte cuvinte din lanț sunt modificate dacă este necesar. Dacă cuvântul nu este legat deloc de celelalte cuvinte, atunci se inițializează un nou lanț.

Metoda prezentată în continuare [Tăt08b/c/d] funcționează în ordine inversă: mai întâi se dezambiguează întregul text și apoi se construiesc lanțurile lexicale de dimensiuni maxime. Pentru dezambiguare s-a folosit algoritmul CHAD prezentat anterior. El identifică printr-o singură parcurgere a textului cuvintele care nu sunt legate ca sens de cuvintele anterioare, iar aceste cuvinte vor primi în mod „forțat” primul sens din WordNet. Dacă scorul fiecărei propoziții va măsura numărul de cuvinte care au fost „forțate” să ia primul sens din WordNet, se va obține structura lexicală coezivă a textului. Se notează cu F numărul de cuvinte „forțate” să ia primul sens din WordNet dintr-o propoziție și atunci în

punctele de minim ale graficului construit cu valorile $1/F$ se vor afla marginile lanțurilor lexicale (vezi figura 8) [Tăt08b/c].

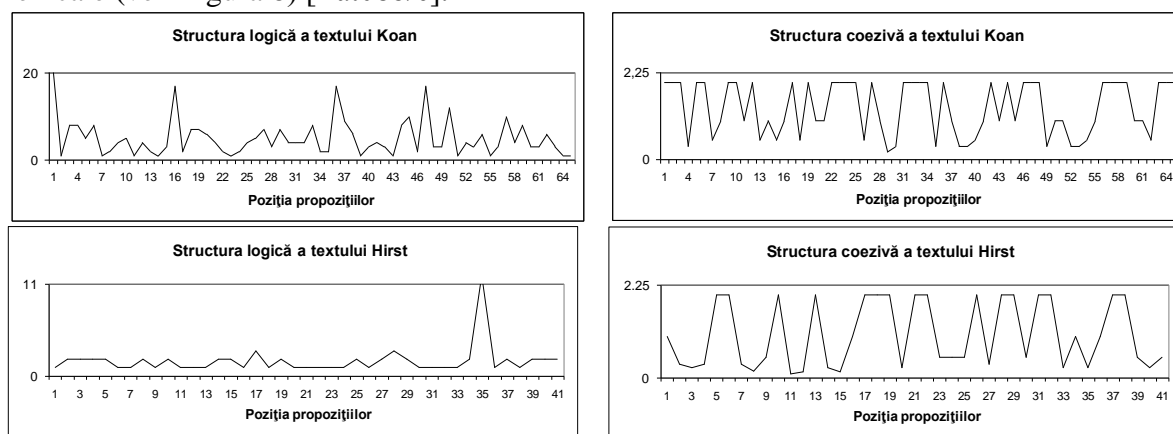


Figura 8. Structura logică și coezivă a textelor Koan respectiv Hirst

Această idee de identificare liniară a lanțurilor lexicale stă la baza unui algoritm liniar de segmentare CTT (Cohesive TextTiling), similar cu algoritmul LTT prezentat anterior cu singura diferență a formulei de scorare a fiecărei propoziții, $Scor(S_i) = \frac{1}{nuw_i}$,

unde nuw_i este numărul de cuvinte „forțate” la primul sens din WordNet din propoziția S_i . Dacă $nuw_i = 0$, atunci $Score(S_i) = 2$. CTT a fost testat pe textul Hirst [Mor91] și „The Koan” [Ric91], obținând rezultate mai bune față de LTT.

Metodele de segmentare prezentate mai sus pot fi de un real ajutor echipei de dezvoltatori, în special metodele de segmentare cu un număr predefinit de segmente – numărul de membri ai echipei. Astfel s-ar putea analiza și rafina mai ușor și mai rapid specificațiile clientului, într-un proces paralel.

3.3. Noi abordări ale sumarizării unui text

Sistemele care pot în mod automat să extragă sumarul dintr-un text au devenit din ce în ce mai studiate și utilizate. Sumarul este o variantă mai scurtă a unui text sau a mai multor texte, cu lungimea de maxim 50% din lungimea textului original, dar care păstrează ideile principale ale textului. În literatură se întâlnesc două feluri de sumare: extras și abstract [Hov03]. Majoritatea sumarizatoarelor automate realizează un extras. Pentru aceasta în literatură se cunosc mai multe abordări [Mar97]. Propozițiile sunt selecționate după criteriile: conțin cele mai frecvente cuvinte, conțin cuvintele din titluri și subtitluri, se află la începutul sau sfârșitul unor secțiuni, utilizează expresii care accentuează importanța pe care o au în text, sunt cel mai strâns legate de alte părți ale textului. Dintre aceste abordări ultima este cel mai dificil de realizat. Conectivitatea poate fi estimată utilizând numărul cuvintelor, sinonimelor sau anaforelor comune [Oră06, Rad02].

3.3.1 Extragerea sumarului dintr-un text segmentat

Există mai multe metode de obținere a sumarului unui text, dar dacă sumarul este obținut dintr-un text care a fost anterior împărțit în segmente, calitatea lui ar trebui să crească, deoarece selectarea unor propoziții care aparțin unor segmente diferite trebuie să îmbogățească sumarul cu noi informații relativ independente. În [Tăt08e, Tăt09c] pe lângă cele câteva metode de segmentare liniară prezentate anterior, se demonstrează faptul că segmentarea îmbunătățește calitatea sumarizării.

Un sumar obținut simplu ca și consecință a relației de implicație dintre două texte este sumarul „PureEntailment” (PE) și se obține adăugând în mod secvențial sumarului fiecare propoziție care nu este implicată de sumarul curent.

Pornind de la segmentările PE, LTT, ArcInt și respectiv ArcReal descrise în subcapitolul 3.2.1, se poate construi un sumar alegând cele mai importante propoziții nu din tot textul, ci din fiecare segment în parte astfel ca în final sumarul să aibă numărul dorit de propoziții. Numărul de propoziții care se aleg din fiecare segment și modul de alegere este dat de algoritmul din figura 9.

Algoritmul *SumPostSeg*(S, n, SEG, j, Sum, x) **este:**

date S {un text; s_1, s_2, \dots, s_n – propozițiile textului}, n {numărul de propoziții din text}, SEG { Seg_1, \dots, Seg_j segmentele în care este împărțit textul}, j {numărul de segmente}, x {mărimea sumarului}

@ se calculează importanța fiecărui segment și se ordonează segmentele după importanță:

$Seg_{i_1}, \dots, Seg_{i_j}$

@ se calculează numărul c_{i_s} de propoziții „esențiale” pentru fiecare segment

$Seg_{i_s}, s = i_1, i_j$

@ se selectează primele k ($< j$) segmente astfel încât $\sum_{Seg_{i_s}} c_{i_s} = x$

@ se reordonează propozițiile „esențiale” din cele k segmente selectate, după ordinea în care apar în text, Seg'_1, \dots, Seg'_k

@ propozițiile din Seg'_1, \dots, Seg'_k sunt propozițiile sumarului SUM

rezultate SUM {sumarul de lungime x }

sf *SumPostSeg*

Figura 9. Algoritmul de obținere a sumarului dintr-un text segmentat

În această primă abordare a construirii sumarului pornind de la o segmentare, s-a considerat că toate segmentele sunt la fel de importante și din fiecare segment s-a ales cea mai importantă (cu cel mai mare scor) propoziție (câte o singură propoziție „esențială”).

Algoritmul *DP_Sum*($S, n, SC, penalitate, Sum, x$) **este:**

date S {un text, s_i – propozițiile textului, $i \leftarrow 1, n$ }, n {numărul de propoziții din text},

SC { $sc_i = scor(s_i)$ scorurile propozițiilor din text}, $penalitate$ {constantă}, P { p_i – pozițiile din text a propozițiilor sumarului, $i \leftarrow 1, x$ }, x {numărul de propoziții din sumar}

pentru $i \leftarrow 1, n$ **execută**

$\delta_i^1 \leftarrow scor(s_i)$

pentru $k \leftarrow 2, x$ **execută**

$$\delta_i^k = \begin{cases} \max_{j(i < j)} (scor(s_i) + \delta_j^{k-1}) & , \text{dacă } s_i \text{ și } s_j \text{ au cuvinte comune} \\ \max_{j(i < j)} (penalitate * (scor(s_i) + \delta_j^{k-1})) & , \text{altfel} \end{cases}$$

$$h_i^k = \begin{cases} \arg \max_{j(i < j)} (scor(s_i) + \delta_j^{k-1}) & , \text{dacă } s_i \text{ și } s_j \text{ au cuvinte comune} \\ \arg \max_{j(i < j)} (penalitate * (scor(s_i) + \delta_j^{k-1})) & , \text{altfel} \end{cases}$$

sfp

sfp

$i \leftarrow \arg \max_j (\delta_j^x)$

inițializează(Sum) {inițial sumarul Sum nu conține nici o propoziție}

pentru $k \leftarrow x, 1$ **execută**

adaugă(Sum, s_i) {adaugă sumarului propoziția s_i }

$i \leftarrow h_i^k$

sfp

rezultate Sum {sumarul de lungime x }

sf *DP_Sum*

Figura 10. Algoritmul de obținere a sumarului dintr-un text segmentat

O altă abordare a problemei sumarizării este cea a aplicării unui algoritm de programare dinamică (vezi figura 10). Cu scopul de a păstra coerența sumarului, acest algoritm va selecta un lanț de propoziții astfel încât două propoziții consecutive să aibă cel puțin un cuvânt în comun, ceea ce corespunde cu principiul continuității din teoria centrării care stipulează ca două unități ale discursului să aibă cel puțin o entitate comună [Oră06].

Se presupune că fiecare propoziție are asociat un scor al reprezentativității ei. Se vor folosi cele trei scoruri logice definite în subcapitolul 3.2.1: suma implicațiilor, a arcelor întregi respectiv ponderate. Scorul unui sumar este suma scorurilor tuturor propozițiilor din sumar. Sumarul va fi construit astfel încât acest scor să fie maxim. Pentru metoda Programării dinamice s-a folosit o penalitate de 1/10 și lista cu 571 de cuvinte „stop” LYRL2004, dezvoltată pentru proiectul SMART [Lew04]. Totodată doar substantivele, verbele, adjectivele și adverbele au fost folosite drept cuvinte comune și s-a utilizat Stemmer-ul Porter [***Por] pentru a compara cuvintele.

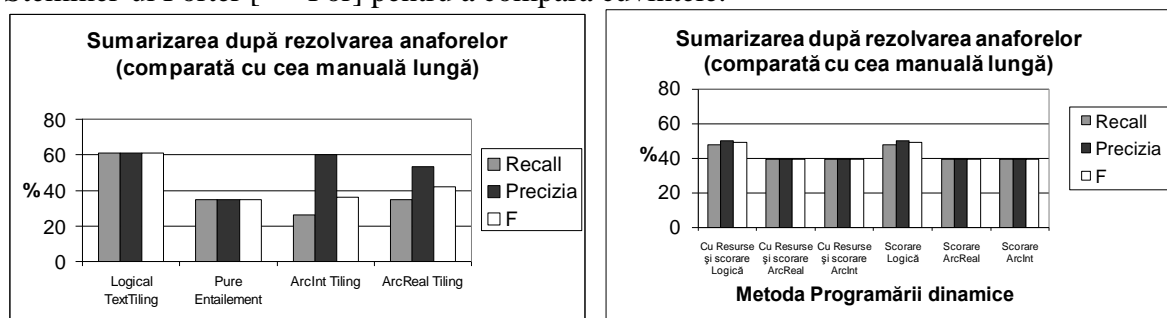


Figura 11. Evaluarea sumarizării

Pentru a evalua metodele de sumarizare propuse s-a utilizat textul literar „The Koan” [Ric91], care pe baza criteriilor lingvistice a fost sumarizat și i s-au rezolvat anaforele în mod manual [Tăt08e]. Din figura 11 se observă că segmentarea îmbunătățește sumarizarea, iar combinarea metodei programării dinamice cu scorarea logică duce la rezultate bune. Dintre metodele logice, cele mai bune rezultate se obțin pentru LTT.

3.3.2 Sumare de lungime variabilă

Sumarul propus anterior se bazează pe segmentare. Și tot pornind de la segmentarea LTT (vezi subcapitolul 4.1) se va obține și următorul sumar de lungime variabilă. Segmentarea LTT pornește de la structura logică a textului, care are o formă fixă, cu un număr fix de minime, astfel că și numărul de segmente obținute este tot fix. De aceea abordarea anterioară care selecta din fiecare segment propoziția cu scorul cel mai mare, unde în prealabil segmentele au fost ordonate după importanța lor, până când s-a obținut numărul dorit de propoziții din sumar nu poate fi aplicată tot timpul, cum ar fi în cazul în care numărul de propoziții din sumar este mai mare decât numărul de segmente. În [Tăt08a] s-au propus trei metode de sumarizare de lungime variabilă care pornesc de la o segmentare existentă, metode care vor fi prezentate în continuare. Testarea s-a realizat pe textul Hirst [Mor91] (din care s-au exclus propozițiile de cinci sau mai puține cuvinte), iar rezultatele au fost comparate cu cele obținute de autori prin lanțuri lexicale.

Un prim pas pentru sumarizare constă în scorarea segmentelor. Sumarul se va obține prin selectarea unui număr de propoziții proporțional cu scorul segmentului respectiv. Se vor folosi următoarele notații: $Scor(s_i) = \text{numărul de propoziții implicate de}$

$$s_i, \quad Scor(Seg_j) = \frac{\sum_{s_i \in Seg_j} Scor(s_i)}{|Seg_j|}, \quad Scor_{final}(s_i) = Scor(s_i) \times Scor(Seg_j), \quad \text{unde } s_i \in Seg_j,$$

$$\text{Scor}(\text{Text}) = \sum_{j=1}^n \text{Scor}(\text{Seg}_j), \quad \text{Ponderea unui segment} : c_j = \frac{\text{Scor}(\text{Seg}_j)}{\text{Scor}(\text{Text})}; \quad 0 < c_j < 1,$$

$n =$ numărul segmentelor obținute de algoritmul LTT, $x =$ lungimea prestabilită a sumarului, $NSenSeg_i =$ numărul de propoziții selectate din segmentul Seg_j .

Algoritmul de sumarizare de lungime prestabilită, notat cu AL (Arbitrary Length) se află în figura 12. S-ar putea ca numărul de propoziții selectate în sumar dintr-un segment, $NSenSeg_j > 1$. Dacă $x < n$, atunci pentru anumite segmente, $NSenSeg_j = 0$.

Algoritmul $AL(\text{Text}, \text{SEG}, n, \text{Lseg}, \text{SC}, \text{SelM}, \text{Sum}, x)$, este:

date Text {un text} SEG {segmentele textului Seg_1, \dots, Seg_n }, n {numărul de segmente},
 Lseg {vectorul cu lungimile (în număr de propoziții) a celor n segmente}, SC
{scorurile fiecăror propoziții $sc_i = \text{Scor}_{\text{final}}(s_i)$ }, SelM {metodă de selecție}, x
{mărimea sumarului}

@ se calculează ponderea fiecărui segment $c_j, j=1, n$

@ se ordonează segmentele după ponderea c_j

{se calculează numărul $NSenSeg_j, j=1, n$ }

Cât timp @ numărul propozițiilor selectate deja în sumar nu a depășit numărul
prestabilit x **execută**

Dacă $[x \times c_j] \geq 1$ {parte întreagă}

atunci $NSenSeg_j = \min(\text{lseg}_j, [x \times c_j])$

altfel $NSenSeg_j = 1$

sfd

sfc

@ se inițializează sumarul Sum

@ se selectează conform metodei de selecție SelM din fiecare segment $Seg_j -$
 $NSenSeg_j$ propoziții și se adaugă sumarului Sum

@ se reordonează propozițiile sumarului, în funcție de ordinea în care sunt în Text
rezultate Sum {sumarul de lungime x }

sf AL

Figura 12. Algoritmul de obținere a sumarului dintr-un text segmentat

Metoda de selecție a propozițiilor dintr-un segment (SelM) este decisivă pentru calitatea sumarului. S-au propus trei metode de selecție a celor mai importante propoziții din fiecare segment, fiecare încercând să excludă cât mai puține informații.

Definiția 1 Se dă un text segmentat $T = \{Seg_1, \dots, Seg_n\}$. Din fiecare segment $Seg_i, i = \overline{1, n}$, se selectează primele $NSenSeg_i$ propoziții astfel încât

$$\sum_{i=1}^n NSenSeg_i = x. \text{ Se obține sumarul } Sum_1 = \{s'_1, \dots, s'_x\}.$$

Definiția 2 Dându-se o segmentare a textului inițial $T = \{Seg_1, \dots, Seg_n\}$, sumarul se obține în felul următor: $Sum_2 = \{s'_1, \dots, s'_x\}$, unde $\forall i = \overline{1, n}, Sum_2 \cap Seg_i = \text{SelPropSeg}_i$, $|\text{SelPropSeg}_i| = NSenSeg_i$, $\sum_{i=1}^n NSenSeg_i = x$ și $\forall s_j \in \text{SelPropSeg}_i$ și $\forall s_k \in Seg_i \setminus \text{SelPropSeg}_i, \text{scor}(s_j) \geq \text{scor}(s_k)$.

Definiția 3 Dându-se o segmentare a textului inițial $T = \{Seg_1, \dots, Seg_n\}$, sumarul se obține în felul următor: $Sum_3 = \{s'_1, \dots, s'_x\}$, unde $s'_1 = s_1$, Seg_0 conține doar propoziția s_1 și $\forall i = \overline{1, n}, Sum_3 \cap Seg_i = \text{SelPropSeg}_i, |\text{SelPropSeg}_i| = NSenSeg_i$, $\sum_{i=1}^n NSenSeg_i = x$ și $\forall s_j \in$

SelPropSeg_i și $\forall s_k \in Seg_i \setminus \text{SelPropSeg}_i, \text{sim}(s_j, Seg_{i-1}) \leq \text{sim}(s_k, Seg_{i-1}), \text{sim}(s_j, Seg_{i-1})$ reprezintă similaritatea dintre s_j și ultima propoziție selectată în Seg_{i-1} . Similaritatea a două propoziții, s, s' , fiind calculată cu metrica $\text{cos}(s, s')$.

Similaritatea sumarelor de diferite lungimi obținute cu AL și cele trei definiții Sum_1 , Sum_2 , Sum_3 , respectiv metoda programării dinamice, calculată relativ la textul inițial este prezentată în figura 13.

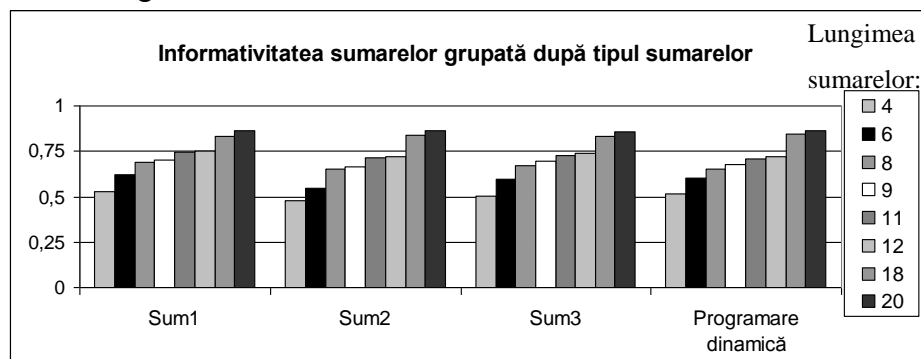


Figura 13. Caracteristica informativă a sumarelor de diferite lungimi, obținute cu AL și cele trei definiții Sum_1 , Sum_2 , Sum_3 , respectiv metoda programării dinamice

Se observă că cele mai bune rezultate au fost obținute pentru aplicarea algoritmului AL cu definiția Sum_1 . Totodată se observă influența pozitivă a segmentării asupra sumarizării. Pentru lungimile 9 și 12, în 5 din 6 cazuri metoda sumarizării precedată de segmentare are rezultate mai bune față de metoda sumarizării directe.

3.3.3 Sumarizare bazată pe lanțuri lexicale

Sumarizarea folosește același algoritm AL prezentat în subcapitolul anterior, doar că algoritmul AL inițial se consideră a fi prima variantă $Var1$, unde scorul unui segment

era calculat după formula $Scor(Seg_j) = \frac{\sum_{S_i \in Seg_j} Scor(S_i)}{|Seg_j|}$. A doua variantă a aceluiași

algoritm, $Var2$ diferă doar prin formula de calcul a scorului unui segment: $Scor(Seg_j) = \sum_{S_i \in Seg_j} Scor(S_i)$. Ambele variante au avantajele și dezavantajele lor,

variantele $Var1$ poate dezavantaja segmentele de lungimi mari în care doar câteva propoziții au un scor ridicat, pe când $Var2$ acordă o prea mare importanță segmentelor lungi.

Disponând de două segmentări obținute prin LTT (subcapitolul 3.2.1) și CTT (subcapitolul 3.2.4), de variantele $Var1$ și $Var2$ de aplicare a algoritmului de sumarizare cu lungime variabilă AL, cât și de cele trei definiții Sum_1 , Sum_2 respectiv Sum_3 , s-a încercat o combinare a sumarelor rezultate într-un potențial sumar Ideal IdS, prin alegerea propozițiilor cele mai frecvente din sumarele obținute.

Pentru evaluare s-au utilizat textele „The Koan” [Ric91], Hirst [Mor91], Tucker1 și Tucker2 [Tuc99]. S-a calculat similaritatea fiecărui sumar cu textul din care a provenit (vezi tabelul 2). Se poate ajunge la concluzia că segmentările LTT și CTT favorizează obținerea unor sumare de calitate. Dintre cele două metode, CTT conduce la rezultate mai bune, parțial datorită faptului că precizia verificării relației de implicație logică este mai mică față de cea a dezambiguării cuvintelor (algoritmului CHAD).

Textul Koan							
seg.	tipul sumarului	lungimea sumarului					
		5		6		10	
		Var1	Var2	Var1	Var2	Var1	Var2
LTT	Sum ₁	0,402941	0,357519	0,490186	0,419314	0,597022	0,614695
	Sum ₂	0,357519	0,303239	0,427179	0,376552	0,548034	0,583095
	Sum ₃	0,427179	0,548034	0,508666	0,531751	0,583095	0,587805
CTT	Sum ₁	0,449629	0,514595	0,531751	0,587805	0,654998	0,662474
	Sum ₂	0,449629	0,442326	0,463739	0,463739	0,568535	0,542697
	Sum ₃	0,483779	0,463739	0,502625	0,508666	0,627334	0,635489
	IdS	0,419314		0,47724		0,631438	
Textul Hirst							
seg.	tipul sumarului	lungimea sumarului					
		5		6		10	
		Var1	Var2	Var1	Var2	Var1	Var2
LTT	Sum ₁	0,579	0,6205	0,6229	0,644	0,7301	0,7126
	Sum ₂	0,5345	0,5313	0,5439	0,5649	0,6978	0,6939
	Sum ₃	0,5377	0,5706	0,5845	0,6253	0,6901	0,6997
CTT	Sum ₁	0,5281	0,4635	0,5439	0,5345	0,6761	0,6529
	Sum ₂	0,5182	0,501	0,5345	0,559	0,672	0,6615
	Sum ₃	0,562	0,501	0,6253	0,5248	0,7089	0,6761
	IdS	0,4975		0,5561		0,6802	
Textul Tucker1							
seg.	tipul sumarului	lungimea sumarului					
		5		6		10	
		Var1	Var2	Var1	Var2	Var1	Var2
LTT	Sum ₁	0,5102	0,517512	0,564498	0,579741	0,67433	0,67433
	Sum ₂	0,5102	0,517512	0,564498	0,579741	0,67433	0,67433
	Sum ₃	0,528176	0,548496	0,554993	0,579741	0,65143	0,665371
CTT	Sum ₁	0,570682	0,538497	0,588551	0,588551	0,678714	0,67433
	Sum ₂	0,558193	0,528176	0,582704	0,582704	0,6873	0,711701
	Sum ₃	0,561361	0,57675	0,602715	0,579741	0,703782	0,691504
	IdS	0,5102		0,594292		0,699744	
Textul Tucker2							
seg.	tipul sumarului	lungimea sumarului					
		5		6		10	
		Var1	Var2	Var1	Var2	Var1	Var2
LTT	Sum ₁	0,637022	0,637022	0,664775	0,664775	0,805387	0,806
	Sum ₂	0,637022	0,637022	0,664775	0,664775	0,805387	0,805387
	Sum ₃	0,594385	0,5547	0,671345	0,581695	0,807441	0,708069
CTT	Sum ₁	0,645	0,610558	0,680946	0,668078	0,740819	0,740819
	Sum ₂	0,644194	0,637022	0,680946	0,690257	0,755929	0,755929
	Sum ₃	0,644194	0,568481	0,730297	0,696311	0,82717	0,713782
	IdS	0,637022		0,664775		0,727607	

Tabelul 2. Similaritățile sumarelor obținute automat

IV. Utilizarea ontologiilor în prelucrarea cerințelor

Cu toate că tehnicile de prelucrare a limbajului natural sunt utile în prelucrarea cerințelor pentru obținerea specificațiilor, de preferat formale sau într-o formă apropiată, sunt necesare resurse suplimentare. Cea mai folosită resursă este ontologia, de vreme ce ontologia permite asocierea unei informații semantice unui text de limbaj natural, text care nu poate fi înțeles altfel de calculator.

4.1. Potrivirea unei ontologii la un text

În prezent a crescut cantitatea de informație semantică existentă în internet. Totuși mai există o cantitate enormă de informație scrisă în limbaj natural fără nici un fel de suport semantic, pe care „făpturile umane nu le-ar putea organiza pe toate” [Pol09]. Și cum există deja ontologii de diverse tipuri, obținute în diverse moduri, s-ar pune mai degrabă problema asocierii unei ontologii existente unui text, decât extragerea ontologiei din text.

Evaluarea ontologiilor

Considerându-se diferitele tipuri de ontologii și modul în care ele au fost obținute, se impune un mijloc de garantare a calității ontologiilor. În [Gan05b] metricile de asigurare a calității ontologiei au fost clasificate în trei tipuri: metrici structurale, măsuri funcționale și măsuri de prelucrare a ușurinței de utilizare. Potrivirea unei ontologii la un text aparține de categoria măsurilor funcționale [Gan05b], astfel că următoarele metrici recomandate vor fi precizia și recall-ul criteriilor de potrivire propuse [Mih10b].

Nevoia de a identifica ontologia corespunzătoare unei anumite sarcini a fost discutată în diferite articole. În [Eng05] ontologiile sunt căutate și evaluate pe baza unui set de cuvinte cheie. În [Tan09] e propus un cadru de selecție a celei mai adecvate ontologii pentru o anumită aplicație de exploatare a textului din domeniul medical. O altă lucrare similară este [Doa03], unde sunt utilizate anumite măsuri de similaritate și tehnici de învățare automată pe bază de nume, context și etichete.

Metriци de evaluare a potrivirii unei ontologii la un text scris în limbaj natural

Există nenumărate ontologii, multe fiind în mod continuu îmbogățite. Dar din punctul de vedere al prelucrării limbajului natural, revine următoarea întrebare: Care ontologie este cea mai potrivită pentru un anumit text și care parte din respectiva ontologie? Pentru a putea răspunde la această întrebare, s-a propus în [Mih10b] utilizarea unui set de metriци specifice.

Prima metriца propusă va evalua câte concepte se regăsesc în textul scris în limbaj natural. Deoarece toate tipurile de ontologii se pot reprezenta ca o colecție de triplete (concept – predicat - obiect), ontologia analizată s-a presupus că se află reprezentată în această formă.

Metriци cantitative

Cea mai simplă tehnică de evaluare a potrivirii unei ontologii la un text e de a căuta conceptele în text. Astfel prima metriца propusă, similară cu o metriца de evaluare a calității unei ontologii și cu o tehnică de aliniere a ontologiilor este *numărul de concepte identificate în text*. Deoarece conceptele pot să apară în text în diferite forme (singular, plural, ...) s-a considerat că s-a găsit conceptul dacă s-a identificat cel puțin un cuvânt care începe cu conceptul căutat.

Câteva metrici derivate pot fi evaluate în mod similar: *numărul de rădăcini găsite în text*, *numărul de frunze găsite în text*, iar lista poate continua. Rădăcinile sunt acele concepte care apar doar în părțile stângi ale tripletelor, iar frunzele sunt conceptele care apar doar în părțile drepte.

Metriți bazate pe relația de consecință logică dintre două texte

În ultimii ani tehnicile de evaluare a implicației dintre două texte sau a relației de consecință logică dintre două texte au evoluat. De aceea s-a propus folosirea acestei tehnici și în [Mih10b] pentru evaluarea potrivirii unei ontologii cu un text. Și aceasta mai ales datorită faptului că un triplet este similar cu o propoziție simplă de limbaj natural (subiect, predicat, complement).

Utilizarea relației de implicație dintre două texte în evaluarea potrivirii unei ontologii la un text este realizată prin verificarea cazului în care ontologia, sau părți din ea implică textul scris în limbaj natural. S-au propus următoarele metrici: *implicația* (prin care se verifică dacă ontologia implică textul), respectiv *numărul de implicații* (prin care se verifică câte triplete ale ontologiei implică textul). În același timp se poate verifica câte propoziții din text sunt implicate de un triplet al ontologiei și astfel s-a propus o a treia măsură, *numărul de propoziții implicate* (se numără fiecare implicație dintre un triplet și o propoziție a textului), respectiv *numărul propozițiilor distincte implicate* (se numără propozițiile distincte implicate de cel puțin un triplet). Pentru a verifica relația de consecință logică s-a ales metoda cosinusurilor (vezi 3.1.3).

S-a evaluat **precizia** (numărul total de potriviri posibile) și **recall-ul** (câte potriviri au fost găsite din numărul total de elemente testate). Pentru a **evalua** metricile propuse s-a folosit un set de nouă ontologii din [***WnB], ontologii care au toate ca rădăcină termenul „business”, dar unde cuvântul „business” are pentru fiecare ontologie un alt sens. Cele nouă ontologii au fost comparate cu prima parte a unui articol din Wikipedia despre „small business” [***WSM]. Din cele nouă ontologii primele patru și cea de a șasea par a fi cele mai potrivite cu textul.

Din păcate metricile care au la bază implicația nu sunt suficient de precise, în special datorită conceptului „business” care e comun celor nouă ontologii. Metricile „cantitative” propuse confirmă presupunerea făcută, conform căreia doar primele patru ontologii și cea de a șasea se potrivesc cu textul (pentru ontologia 1 s-a obținut un recall de 10,5%, pentru cea de a doua 8,1%, a patra 20%, iar a șasea 6,5%).

4.2. Învățarea unei ontologii dintr-un text de limbaj natural pe baza arborelui de analiză sintactică a unei propoziții

Rezultatul analizei gramaticale a unei propoziții este reprezentat de obicei prin intermediul unui arbore. Iar între cuvintele unei propoziții se pot identifica diferite relații de dependență [Mar08], relații care sunt surprinzător de asemănătoare tripletelor RDF (Resource Description Framework – cadrul de descriere al relațiilor), cel mai simplu mijloc de a reprezenta o ontologie. Pentru analiza gramaticală a unei propoziții s-a folosit unealta on-line dezvoltată de grupul de procesare a limbajului natural de la Stanford, Natural Language Processing Group [***STO], care prezintă o precizie mai mare de 87% [Kle03].

În procesul de învățare a ontologiilor s-a recunoscut importanța analizei sintactice a textului. Dar relațiile gramaticale au fost utilizate doar ca și șabloane pentru recunoaștere în [Cim05, May09] sau ca și constrângeri pentru identificarea relațiilor dintre conceptele ontologiei în [Kaw04].

Analiza sintactică a unei fraze

Analizorul Stanford [***STS] furnizează arborele de analiză sintactică, dependențele de tip și dependențele condensate (de exemplu: *nsubjpass*(abbreviated-6, title-3)). Pornind de la dependențe, se poate genera o structură sub formă de graf (vezi figura 14, corespunzătoare propoziției „This function title is often abbreviated to the initials 'HR'.” [***WHR]).

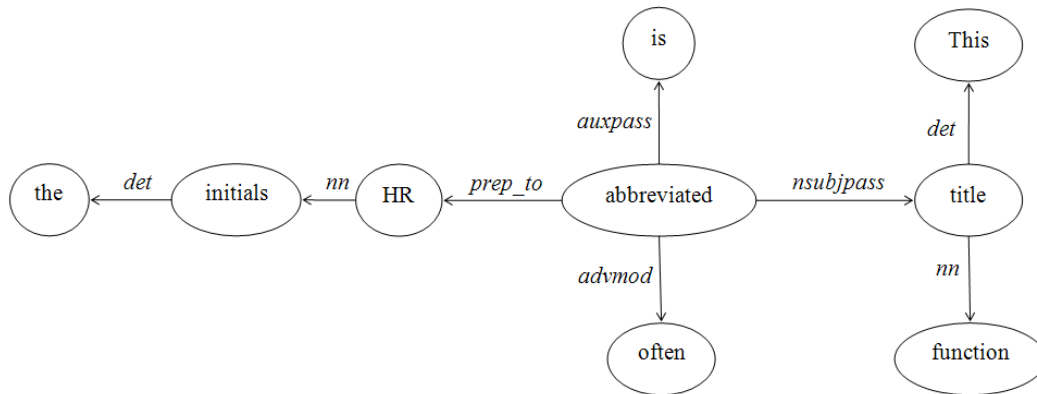


Figura 14. Graful obținut din dependențele de tip

Identificarea tripletelor

Prima metodă propusă de extragere a tripletelor este cea de **extragere a tripletelor din graful dependențelor** (vezi figura 15) [Mih10c].

@elimină toate dependențele „dep”

Pentru @fiecare dependență D₁ execută

Pentru @fiecare dependență D₂ diferită de D₁ execută

Dacă @D₁ și D₂ au cuvinte comune atunci

Dacă @cuvântul comun este verb iar cuvintele distincte sunt substantive atunci

@adaugă tripletul (cuvântul distinct din D₁, cuvântul comun, cuvântul distinct din D₂) la lista tripletelor

sfd
sfd
sfp
sfp

Figura 15. Extragerea tripletelor din graful de dependențe

A doua metodă de extragere a tripletelor este cea de **extragere a tripletelor din arborele de analiză sintactică** (vezi figura 16).

Pentru testarea metodelor propuse, s-a folosit același text și aceleași ontologii ca și în subcapitolul precedent. S-a evaluat dacă două cuvinte sunt similare sau nu (sunt incluse unul în altul), cuvintele stop [***SPW] și cuvântul „business” neluându-se în considerare.

Metoda/Ontologia	prima	a doua	a patra	a șasea
I	0/0	2+2(inversele)/1	0/0	0/0
II	2/2	2/3	1/1	2/2

Tabelul 3. Tripletele similare extrase din text/apartinând de ontologie

Rezultatele **evaluării** se pot observa în tabelul 3. Deoarece cele mai multe triplete extrase din text sunt similare cu triplete din cea de a doua ontologie pentru ambele metode, articolului testat despre o afacere mică i se potrivește cel mai bine ontologia în care „business” are sensul de afacere (ceea ce se aștepta).

Subalgoritmul IdentificăTriplet(S) este {S este o propoziție}
 @citește simbolul SB
Dacă @SB este NN atunci @conceptul din triplet este cuvântul următor sfd
Dacă @SB este NP atunci @analizează constituenții sfd
Dacă @SB este S atunci IdentificăTriplet(@noua propoziție S) @folosește ca și concept conceptul din noua S sfd
 @citește simbolul SB
Dacă @SB este VB, atunci @predicatul tripletului este următorul cuvânt sfd
Dacă @SB este VP atunci @analizează constituenții sfd
 @citește simbolul SB
Dacă @SB este NN atunci @obiectul tripletului este următorul cuvânt sfd
Dacă @SB este NP sau ADJP atunci @analizează constituenții sfd
Dacă @SB este S atunci IdentificăTriplet (@noua propoziție S) @și folosește ca și obiect conceptul din noua S sfd
sf IdentificăTriplet
 {în „analizează constituenții”, are loc un proces similar: se citește următorul simbol, se identifică tipul său până la epuizarea ramurii curente, iar în cazul unei conjuncții sau a unei disjuncții de termeni (pentru concepte și obiecte, sau disjuncție pentru predicate), atunci fiecare element al conjuncției va fi considerat separat, iar în final, în loc să se obțină un triplet, se va obține o serie de triplete}

Figura 16. Extragerea tripletelor din arborele de analiză sintactică

4.3. Similaritatea utilizată în identificarea nevoii de a refactoriza o ontologie

Ontologiile existente s-au obținut în diverse forme, de la modul manual, produsul muncii diferiților specialiști sau voluntari [Wil06], la modul semi-automat [Mor06] dar și automat [Dah06]. Mai mult, nevoia de îmbunătățire a unei ontologii existente a fost discutată în [Suk08, Che06]. Unele sunt dezvoltate într-o perioadă lungă de timp sau au un număr impresionant de elemente. Când într-o lucrare de mari dimensiuni sunt implicate diferite surse, mai mult când între momentul conceperii ontologiei și momentul final al livrării ei trece mult timp, atunci crește probabilitatea ca în produsul final (ontologia finală) să existe anumite discrepanțe, să existe anumite elemente plasate greșit.

Diferite măsuri ale similarității au fost utilizate pentru evaluarea unei ontologii [Bra05], pentru alinierea ontologiilor [Euz04] sau pentru potrivirea ontologiilor [Euz07]. Dar de obicei măsurarea similarității utilizează similaritatea cuvintelor, evaluată prin intermediul măsurilor de evaluare a șirurilor. Omonimele, deși sunt scrise în mod identic, au sensuri diferite și sunt recunoscute în mod greșit ca fiind identice. Un alt tip de similaritate utilizat în prelucrarea limbajului natural este similaritatea textelor. Pentru aceasta sunt folosite aproape toate cuvintele care apar în text (mai puțin cuvintele stop), iar similaritatea textului se calculează folosind similaritatea cuvintelor din text.

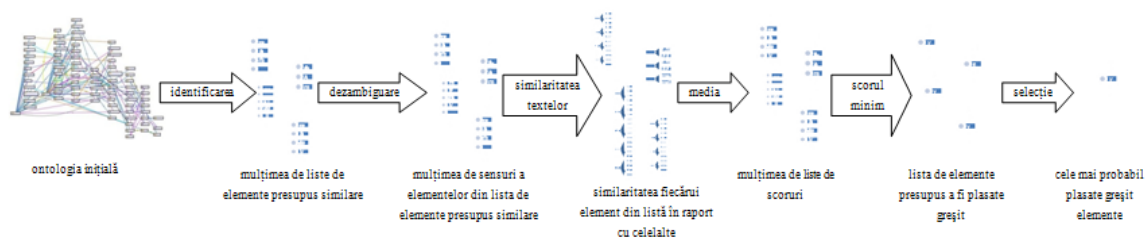


Figura 17. Procesul de selectare a elementelor cel mai probabil plasate greșit

Dacă se consideră cazul în care ontologia este reprezentată de o listă de triplete, atunci se pot extrage subliste de elemente care sunt legate prin același predicat de același concept sau predicate care leagă aceleași două concepte. Astfel de **elemente** ar trebui să fie **similare**. Dacă se evaluează similaritatea a două câte două elemente, se face o medie a

similarităților pentru fiecare element în parte, atunci elementul care obține cel mai mic scor este elementul cel mai probabil a fi **plasat greșit** (vezi figura 17) [Mih10d].

Algoritmul de identificare a elementelor cel mai probabil plasate greșit

Algoritmul Selectarea_Celor_mai_Probabil_Plasate_Greșit_Elemente (O, n) este

Date O- ontologia dată prin lista ei de triplete

n- numărul de elemente probabil plasate greșit căutate

@ extrage lista de elemente (concepte)

Pentru @fiecare element distinct din ontologie **execută**

@identifică lista de elemente distincte care sunt legate prin același predicat de elementul curent

@ identifică lista de predicate care leagă același element de elementul curent

sfp

@păstrează într-o mulțime S doar listele cu cel puțin trei elemente

@inițializează o mulțime vidă de liste S'

Pentru @fiecare element distinct e din fiecare listă care aparține de mulțimea S **execută**

@ ia toate predicatele distincte și elementele distincte relaționate cu e și se construiește lista L cu ele

@adaugă și elementul e listei L

sensul_selectat=1

sim_min=1

Pentru @fiecare sens i din dicționar, sens al numelui elementului curent e **execută**

@calculează similaritatea s a textului compus din numele elementelor din L și sensul i

Dacă @acesta este primul sens **atunci**

sim_min=s

altfel Dacă sim_min>s **atunci**

sim_min=s

sensul_selectat=i

sfd

sfd

sfp

@adaugă în S' într-o poziție analoagă poziției curentului element e sensul selectat

sfp

i=0

@inițializează mulțimea vidă S''

Pentru @fiecare listă L' din S' **execută**

mins=1

poz_min=1

i=i+1

Pentru @fiecare text t din lista L' **execută**

nr=0

sum=0

Pentru @fiecare text t' din lista L'

Dacă e ≠ e' **atunci**

@adaugă la sum similaritatea lui t relativă la t'

nr=nr+1

sfd

sfp

med=sum/nr

sfp

Dacă i=1 **atunci**

mins=med

altfel Dacă mins>med

mins=med

poz_min=i

sfd

sfd

@adaugă la S'' perechea formată de elementul de pe poziția i din lista L a mulțimii S (unde lista L este „părintele” listei de liste L') și scorul ei, med

sfp

@ordonează elementele mulțimii S'' în ordinea crescătoare a scorurilor și se păstrează primele n

Rezultate: primele n elemente cu cele mai mici scoruri
sf Selectarea_Celor_mai_Probabil_Plasate_Greșit_Elemente

Observație: în algoritm, prin elemente sau relații distincte se înțelege cele care au URI distincte.

Ca exemplu de funcționare a algoritmului, se presupune că într-o ontologie cuvintele „man”, „women”, „elderly” și „driver” reprezintă toate subclase ale clasei „adult”, iar în urma dezambiguării s-a identificat primul sens din WordNet [***WNT] pentru toate cele patru cuvinte:

man – man, adult male (an adult person who is male (as opposed to a woman)) "there were two women and six men on the bus"

women – woman, adult female (an adult female person (as opposed to a man)) "the woman kept house while the man hunted"

elderly – aged, elderly (people who are old collectively) "special arrangements were available for the aged"

driver – driver (the operator of a motor vehicle)

În urma calculelor elementul „driver” obține cel mai mic scor.

4.4. Rolul dezambiguării în evaluarea ontologiilor

Dezambiguarea a fost deja utilizată în domeniul ontologiilor pentru a îmbogăți o ontologie [Ste02] sau pentru a construi o ontologie [San07]. În cele ce urmează se va discuta rolul pe care îl are dezambiguarea pentru potrivirea ontologiilor, respectiv potrivirea unei ontologii cu un text de limbaj natural.

Ontologiile și dezambiguarea

Diversitatea mijloacelor prin care s-au obținut ontologii nu garantează faptul se folosesc aceiași URI pentru concepte identice aparținând unor ontologii diferite. Însă fiecare concept are asociat un nume, care poate fi folosit pentru a decide dacă conceptele sunt identice sau nu, chiar și în cazul în care au URI diferit. Doar că o simplă măsură a similarității cuvintelor nu e utilă în cazul omonimelor. De aceea trebuie folosite și elementele direct legate de conceptul de dezambiguit din ontologie [Mih10e]:

Algoritmul de dezambiguare

funcția DezambiguareRecursivă(*WORD*, *SET*, *ONTOLOGY*) *este*

Date: *WORD*, *SET*, *ONTOLOGY*

Precondiții: *WORD* este cuvântul de dezambiguit

SET este mulțimea vecinilor săi din ontologia *ONTOLOGY*

SENSE ← 0

max ← 0

pentru @fiecare sens i a cuvântului *WORD* din dicționar **execută**

@evaluează scorul „overlap” al mulțimii de cuvinte *SET* față de glosa i din dicționar

dacă @scorul > *max* **atunci**

SENSE ← i

sfd

sfp

dacă $max = 0$ **atunci**

dacă @există vecini în ontologia *ONTOLOGY* ai elementelor din *SET* neincluși în *SET* **atunci**

@adaugă vecinii la *SET*

SENSE ← DezambiguareRecursivă(*WORD*, *SET*, *ONTOLOGY*)

sfd

sfd

Rezultate: *SENSE*

Postcondiții: *SENSE* este indexul sensului corect al cuvântului *WORD*, 0 altfel

sf DezambiguareRecursivă

Abordările prezentate accentuează rolul pe care tehnicile de prelucrare ale limbajului natural le joacă în rezolvarea problemei asocierii informațiilor semantice unui text de limbaj natural, dar și a îmbunătățirii calității unei ontologii.

4.5. Prelucrarea cerințelor asistată de ontologii

De vreme ce ontologiile îi permit calculatorului să înțeleagă limbajul natural, ele au o mare aplicabilitate în domeniul extragerii informațiilor. Acest lucru a fost evidențiat încă din 2003 în [Mae03], unde autorii au propus o metodă pas cu pas de îmbogățire a ontologiei existente pentru a o adapta extragerii informației dintr-un nou text.

Într-o lucrare mai recentă [Yil07], este prezentată o metodă de extragere a informației supervizată de o ontologie, fără a se utiliza nici o altă resursă de cunoaștere. Totuși precizia rezultatului va depinde de calitatea ontologiei folosite. În același timp metoda identifică elementele neutilizate ale ontologiei și în acest mod calitatea ontologiei se poate îmbunătăți, precum și informația extrasă.

Selectarea semi-automată a cerințelor

Scopul lucrării [Mih10f] este de a se extrage un program abstract dintr-un text scris în limbaj natural care reprezintă cerințele programului. Cu alte cuvinte de a identifica precondiția, postcondiția și lista de variabile.

Subalgoritmul *PrePostSelection*(*g*, *pre_list*, *post_list*) *este*:

date: *g* – un graf de dependențe direcționale

@Identifică lista de verbe *VB*: *VB_list*

@Inițializează lista listelor de cuvinte, *Word_ll* și pune în fiecare listă pe prima poziție verbul *VB*

@Inițializează o listă vidă de cuvinte distincte *Var_l*

Pentru @fiecare listă *L* din *Word_ll* **execută**:

@Identifică traseele din graf care se termină cu o dependență „sub” pentru un *VB* sau un *VBD*, iar dacă nu există nici o astfel de dependență, traseele care se termină cu o dependență „obj” pentru *VB*, respectiv cele care se termină cu o dependență „amod” pentru *VCN*

@ Adaugă toate cuvintele de pe traseu la *L*

@ Cuvântul *NN* care se află în relația „sub”, „obj” sau „amod” trebuie adăugat listei de variabile *Var_l*

@ Adaugă la *L* toate cuvintele conectate direct sau indirect cu ultimul cuvânt (*NN*)

sfp

rezultate: *pre_list* – lista *L* din *Word_ll* ce începe cu *VB*

post_list – celelalte liste din *Word_ll*

Var_l – lista variabilelor

sf *PrePostSelection*

Figura 18. Algoritmul de analiză a frazei

Se pot identifica două tipuri de cerințe: cerințe care sunt compuse din mai multe fraze, dar și mai multe cerințe care compun o singură propoziție sau frază. În primul caz se pare că problema identificării precondiției și a postcondiției e mai ușor de rezolvat, de vreme ce ea poate fi redusă la problema selectării frazelor. În cel de al doilea caz dintr-o singură propoziție ar trebui extrase una sau mai multe precondiții și una sau mai multe postcondiții. În cel de al doilea caz, însăși **fraza** trebuie **analizată**. Se va utiliza din nou Stanford Parser-ul [***STS]. Precondițiilor le corespunde în limbaj natural o propoziție la trecut sau prezent, iar postcondițiilor una la viitor. Lista de variabile este reprezentată de subiectele respectivelor propoziții. Algoritmul corespunzător se află în figura 18. În cazul în care cerințele sunt exprimate prin mai multe fraze, înainte de a se aplica algoritmul ar trebui analizate frazele, extrase triplete, identificate sinonimele și cuvintele subordonate,

unindu-se astfel grafurile de dependențe (vezi subcapitolul 4.2). Dacă volumul cerințelor e mare, acestea se vor segmenta mai întâi.

Pentru a se obține într-adevăr specificații Z din cerințe scrise în limbaj natural, nu este suficientă selectarea listelor de cuvinte corespunzătoare precondiției/postcondiției. Pentru a se obține scheme Z este necesară utilizarea unei ontologii, care are noduri-bază tipurile de bază și operațiile de bază pentru aceste tipuri. De exemplu în „*Generate the first prime number larger than a given natural number n .*” [***Adr], se poate observa că este utilizat tipul de bază *număr natural* (natural number), iar *prim* (prime) și *mai mare decât* (larger than) sunt operații definite pe mulțimea numerelor naturale, doar *primul* (first) nu este un predicat care aparține de mulțimea funcțiilor matematice definite pe mulțimea numerelor naturale.

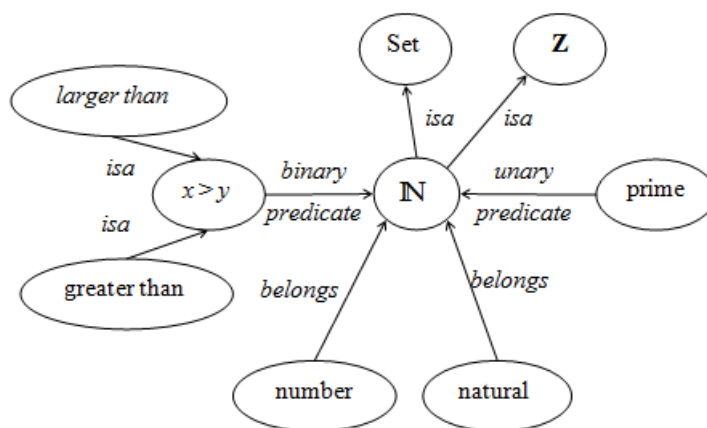


Figura 19. Fragment din ontologia numerelor naturale

În cazul utilizării unei ontologii, prima dată se vor căuta în ontologie cuvintele care reprezintă variabile (vezi Figura 19) și se vor înlocui cu tipul cel mai apropiat (în exemplu, după ce s-a identificat nodul „natural”, cel mai apropiat tip este \mathbf{N}). Apoi se caută și celelalte cuvinte din listele de cuvinte care exprimă precondiția și postcondiția în sub-ontologiile care au rădăcina tipul identificat. Dacă cuvintele respective sunt identificate în ontologie, atunci predicatul de bază conectat direct la tipul identificat va înlocui cuvântul curent.

În exemplul discutat „*Generate the first prime number larger than a given natural number n .*”, lista de cuvinte identificate care exprimă postcondiția este *Generate + larger + number* (the subject) + *the + first + prime*. *Număr* (Number) poate fi înlocuit de $x \in \mathbf{N}$ și de asemenea *mai mare* (larger) $x > n$. Calculatorul poate identifica singur aceste tipuri sau predicate de bază și în acest mod poate asista procesul de transformare a cerințelor scrise în limbaj natural în specificații formale.

Concluzii

Au fost propuse două seturi de metrici pentru evaluarea potrivirii unei ontologii cu un text scris în limbaj natural. Primele se bazează pe evaluarea statistică a numărului de concepte care apar în text, iar cele din urmă pe relația de implicație logică dintre texte.

Cu toate că din testele efectuate metricile bazate pe implicația logică dintre două texte au fost neconcludente, parțial datorită marii diferențe între dimensiunea ontologiilor analizate și a textului testat, primul set de metrici a condus spre rezultatele așteptate.

S-au prezentat două metode de conversie a rezultatului analizei gramaticale ale unui text scris în limbaj natural în triplete ale unei ontologii. Prima metodă propusă furnizează mai puține triplete dar mai concise, apropiate de calculator, iar cea de a doua

poate furniza triplete mai apropiate de utilizatorul uman. Ele accentuează rolul pe care gramatica îl joacă în procesul de construire a unei ontologii sau de potrivire a unei ontologii la un text.

Totodată s-a prezentat o metodă capabilă să identifice precis cele mai probabil plasate greșit elemente dintr-o ontologie și să ajute dezvoltatorii să îmbunătățească calitatea ontologiei.

În final s-a prezentat o metodă semi-automată de extragere a specificațiilor formale dintr-un text scris în limbaj natural. Ea utilizează Stanford Parser-ul pentru a obține graful de dependențe, urmat de un proces de unire a ontologiilor pentru a conecta mini-ontologii, iar la sfârșit pe baza unor principii semantice, sunt extrase propoziții scrise în limbaj natural care reprezintă precondiții și postcondiții.

V. Aspecte ale utilizării metodelor formale pentru dezvoltarea corectă a programelor

5.1. Asigurarea calității softului – o activitate continuă

Este bine cunoscut faptul că din efortul asociat dezvoltării și utilizării produsului soft, întreținerea lui reprezintă 60%. Pentru a reduce acest cost, calitatea produsului trebuie garantată încă de la primele faze ale ciclului de dezvoltare. În cazul produselor de mari dimensiuni, întreținerea lor este imposibilă dacă nu sunt îndeplinite anumite cerințe de calitate. Măsurarea calității softului se face cu ajutorul metricilor soft, existând metrici specifice pentru diferitele faze de dezvoltare a produsului.

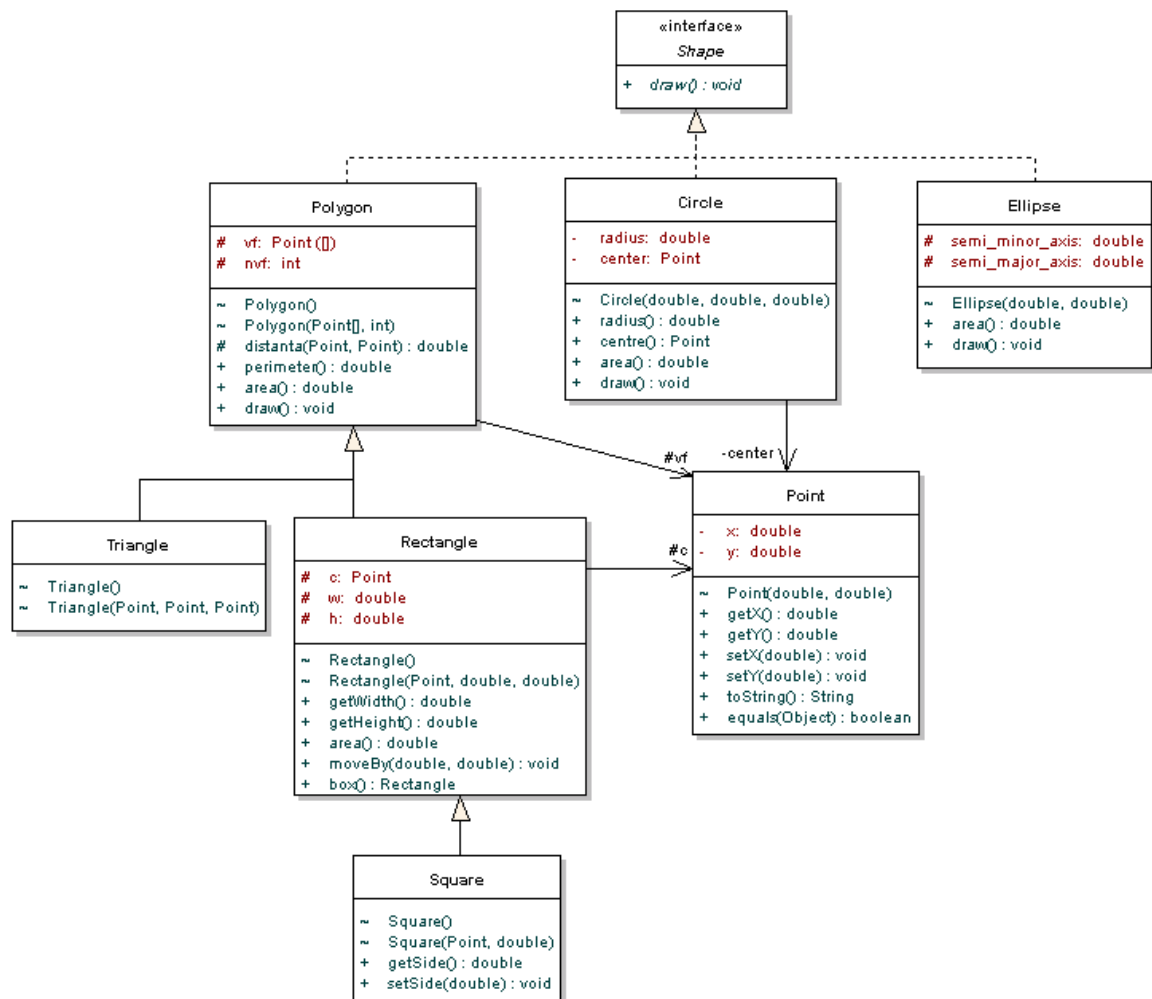


Figura 20. Ierarhia de clase figuri geometrice

O metodă de a asigura calitatea produsului soft este de a-i evalua calitatea în fiecare fază de dezvoltare și a îmbunătăți produsul. Pentru a demonstra importanța evaluării calității softului cât mai devreme în procesul de dezvoltare, s-a proiectat o ierarhie de clase [Șer05] (vezi figura 20). În tabelul 4 se află rezultatele evaluării claselor din ierarhie. Se observă că se poate îmbunătăți această ierarhie de clase, prin derivarea clasei Circle din clasa Ellipse, caz în care complexitatea va scădea prin reducerea

numărului de operații și operatori ai clasei Circle. Totodată valoarea metricii COM ar trebui să fie mai mare, ceea ce înseamnă că ar trebui introduse câteva comentarii suplimentare.

Clasa/Metrica	SIZE (numărul liniilor de cod)	COM (procentul comentariilor)	WMC (greutatea metodelor unei clase)	DIT (adâncimea arborelui de moștenire)	NOC (numărul de „copii”)	LCOM (lipsa coeziunii metodelor)	CBO (cuplare a dintre clase)
Shape	6	33%	0	1	0	-	0
Polygon	37	8%	8	1	2	0	1
Circle	17	17%	5	1	0	0,5	1
Ellipse	34	5%	3	1	0	0	0
Triangle	15	7%	2	2	0	0	1
Rectangle	46	20%	7	2	1	0,46	1
Square	22	9%	4	3	0	0	1
Point	37	8%	9	1	0	0,33	0

Tabelul 4. Rezultatele evaluării metricilor pentru clase

O altă metodă este de a dezvolta soft de calitate superioară cu ajutorul șablonelor de proiectare sau al metodelor formale. Din păcate nu există un model general valabil pentru a garanta calitatea softului. Pentru fiecare proces de dezvoltare ar trebui specificate cerințele de calitate de la început, urmărite pe tot parcursul procesului de dezvoltare și rezolvate discrepanțele în primele faze, când costul de modificare este redus.

5.2. Abordarea multiformală în specificarea sistemelor soft

Sistemele soft de mari dimensiuni prezintă mai multe aspecte. Fiecare aspect necesită aplicarea unui anumit formalism și verificarea sa formală utilizând unelte specifice. De aceea specificațiile obținute prin aplicarea diferitelor metode formale trebuie integrate, rezultând o specificație multiformă (integrată sau heterogenă).

Metodele se clasifică în: metode formale orientate pe stare (sau model), orientate pe proprietăți (axiomatice sau algebrice) și hibride. Metode formale orientate pe stare sunt: Z [Abr80], VDM [Jon86], B [Joh73], automatele cu stări finite [Gil62], rețelele Petri [Pet62], CCS [Mil80], CSP [Hoa85]. Metodele formale orientate pe proprietăți sunt: CASL [Ast02], ACT ONE [Ehr83], Anna [Luc85], Larch [Gar93], CLEAR [Bid91], OBJ [Fut85], LOTOS [Eij89]. În general integrarea metodelor formale constă în combinarea a două sau mai multe metode formale de specificare complementare.

În cazul specificării unui sistem complex, se poate aplica o abordare omogenă sau eterogenă. Abordarea omogenă descrie multiplele aspecte ale sistemului soft utilizând un singur formalism, capabil să exprime toate aspectele sistemului. Astfel de specificații omogene pot fi realizate sau specificând noi limbaje, sau extinzând unul existent. Abordările eterogene, multiformale, folosesc mai multe formalisme existente pentru a acoperi toate aspectele sistemului specificat. În funcție de modul de îmbinare sintactică a limbajelor, acest tip de abordare poate fi o integrare puternică a formalismelor, sau o compunere/cuplare a părților independente ale specificației, fiecare parte fiind scrisă folosind un alt limbaj.

În cazul combinării sintactice a limbajelor există trei abordări principale. Prima abordare utilizează reprezentări grafice ale comportamentului (rețele Petri, diagrame de stare, sisteme etichetate ale tranzițiilor) și diferite tipuri de date (specificații algebrice, B, Z, VDM). A doua abordare constă în combinarea unei algebre a procesului, cum este CCS sau CSP cu un formalism algebric. Un astfel de multiformalism este LOTOS. A treia abordare încearcă integrarea algebrei procesului (CCS, CSP) cu un formalism orientat pe stări (Z, Object-Z [Smi99], B, VDM). Un astfel de formalism integrat este ZCCS [Gal96].

Premisele unei abordări multiformale sunt: învățarea metodelor formale, utilizarea metodelor integrate și a uneltelor necesare. Unul dintre principalele dezavantaje ale metodelor formale, acela de a fi dificil de învățat datorită notațiilor, conceptelor și a metodelor matematice, poate constitui un impediment în cazul unei abordări multiformale. În plus combinând diferite notații, pot fi introduse ambiguități, incompatibilități și pot să rezulte specificații dificil de înțeles. Dar s-a argumentat că prin folosirea mai multor metode, fiecare în cele mai potrivite situații, se vor obține specificații mai scurte, clare și concise față de cazul utilizării unui singur formalism [Cio04].

5.3. Reutilizarea specificațiilor formale

În momentul în care se produce o nouă variantă a unui soft existent, o parte din cod este reutilizat, dar și specificațiile, documentațiile sau alte produse secundare ale procesului de dezvoltare soft pot fi de asemenea reutilizate. Dintre toate tipurile de specificații, specificațiile formale sunt cele mai utile, dar sunt și cel mai greu de dezvoltat. Reutilizarea lor ar fi foarte utilă. Un alt caz de reutilizare a specificațiilor formale este cel de construire a unui produs pentru mai multe platforme.

Totodată specificațiile formale pot fi folosite pentru a identifica elementele reutilizabile. De exemplu în cazul în care există o librărie de componente și se dorește producerea uneia noi, dar se urmărește reutilizarea câtor mai multe componente din vechea librărie, atunci identificarea componentelor analoage este cel mai ușor de realizat dacă cele două componente au fost specificate formal cu ajutorul aceluiași formalism.

Avantajul folosirii unor specificații formale este acela că ele sunt într-o formă matematică și astfel este ușor de identificat existența unei funcții bijective între două specificații, demonstrând similaritatea lor. K. Periyasamy și J. Chidambaram, [Per96] au definit relația de identitate sau analogie a specificațiilor Z din punctul de vedere al declarației, semnăturii și proprietăților. Forma standard a specificațiilor formale garantează faptul că aceeași specificație se va afla tot timpul în aceeași formă.

Totodată analogia este mijlocul cel mai simplu de învățare. Astfel specificații formale similare pot fi folosite pentru a ușura procesul de învățare a respectivei metode formale [Mih05].

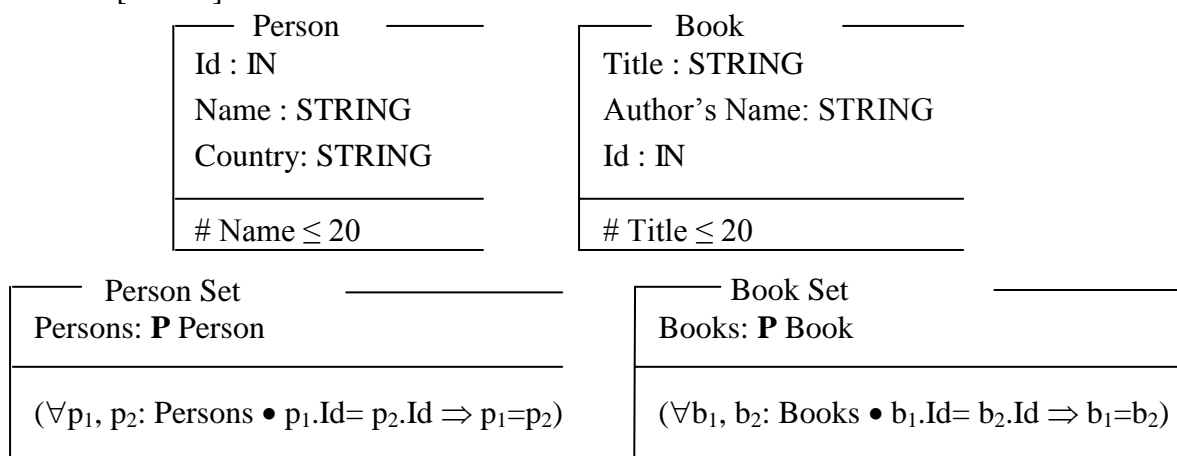


Figura 21. Scheme de date Z

Să presupunem că avem specificațiile pentru tipurile de date Persoană și pentru o Mulțime de Persoane, după cum pot fi observate din figura 21. Se dorește specificarea în Z a entității Carte (Book) și a unei mulțimi de Cărți (Book Set). În urma studierii cerințelor legate de Carte, se observă că cele două entități sunt similare. Pornind de la modelul

Persoanei, se scriu imediat specificațiile pentru Carte, după cum se poate observa din figura 21.

Un alt caz de reutilizare a specificațiilor formale este cazul în care este necesară combinarea a două componente. Pentru a evita introducerea unor erori, se recomandă ca mai întâi să se compună specificațiile lor. Dacă aceste specificații sunt formale, compunerea se va face simplu prin aplicarea unei reguli de compunere, obținându-se astfel specificația componente rezultată.

Elementele similare vor avea specificații similare. Dacă se observă că două elemente sunt similare, se va specifica doar unul, iar apoi în mod analog prin reutilizare se va obține specificația celuilalt element. Similaritatea celor două specificații poate constitui un semnal de alarmă pentru client, dacă cele două elemente nu sunt de fapt similare, caz în care cerința clientului nu a fost deplin înțeleasă de dezvoltatori.

5.4. O aplicație care asistă rafinarea schemelor Z

Având ca scop principal încurajarea folosirii metodelor formale în general și a limbajului Z în special, s-a realizat o aplicație care asistă în mod semi-automat rafinarea schemelor Z , aplicație care va fi descrisă în cele ce urmează [Mih10a].

Pentru a putea rafina o schemă, aplicația trebuie să permită în primul rând definirea unei scheme Z . În ceea ce privește editarea schemelor Z , există deja o astfel de aplicație [Gao09]. Aplicația proprie permite editarea schemelor Z în mod grafic, utilizatorul putând edita o singură schemă la un moment dat, modificarea realizându-se componentă cu componentă (nume, declarații, predicate). Simbolurile speciale pot fi inserate dintr-o listă cu categorii de simboluri [Dil99]. În același timp pentru a sprijini procesul de învățare a limbajului Z , în momentul în care un utilizator alege un caracter din lista de simboluri, va apare o scurtă descriere a simbolului împreună cu un exemplu în care respectivul simbol este utilizat.

În momentul adăugării unei noi componente, are loc o analiză sintactică a ei care are ca scop principal identificarea elementelor de bază, conform schemei de notație [Mih10a], a listei componentelor identificate deja și a simbolurilor speciale „?” și „!”. În cazul în care utilizatorul observă că analizatorul sintactic a identificat în mod greșit un element, el poate interveni pentru a corecta eroarea.

Numele unei scheme noi definite va fi adăugat listei cu nume de scheme. Din această listă schemele pot fi selectate pe rând, câte una sau câte două, în vederea rafinării lor [Woo96].

Dacă se selectează două scheme, acestora li se poate aplica o operație binară, cum ar fi conjuncția, disjuncția și compunerea (vezi figura 22). Dacă se selectează o singură schemă, acesteia i se poate aplica o operație unară, cum este negarea sau decorarea. Concomitent întrucât majoritatea componentelor unei scheme sunt elemente matematice, acestora li se pot aplica diferite teoreme sau proprietăți matematice. Toate aceste operații se pot efectua în mod asistat, aplicația oferind un exemplu de folosire a operației selectate, utilizatorul trebuind să introducă doar elementele noi. Colecția de operații de rafinare a elementelor unei scheme se îmbogățește în mod continuu, prin adăugarea de către utilizator a unor operații noi. Prima folosire a unei operații noi va fi memorată și oferită ca exemplu aplicațiilor viitoare. Operațiile de rafinare care pot fi aplicate componente selectate din schema curentă se identifică pe baza relației de similaritate dintre două componente.

Definiție: Două componente ale unor scheme Z sunt similare dacă sunt compuse din același număr de elemente de bază, aflate în aceeași ordine și având aceleași tipuri.

Cea mai importantă operație de rafinare care poate fi realizată în mod asistat este transformarea unei scheme Z într-un program abstract. Pentru a transforma o schemă Z într-un program abstract, trebuie identificate: cadrul, care conține lista tuturor variabilelor

din schemă, precondiția și postcondiția. Cadrul a fost identificat anterior, prin analiza sintactică a elementelor constituente ale unei scheme Z. Trebuie doar selectate variabilele din lista de elemente. Precondiția este conjuncția predicatelor din schemă care conțin variabile al căror nume include simbolul „!”, iar postcondiția e conjuncția predicatelor în care numele variabilei se termină cu „?”. Bineînțeles utilizatorul poate interveni pentru corectarea erorilor. Programul abstract obținut poate fi salvat într-un fișier text, în vederea prelucrării ulterioare.

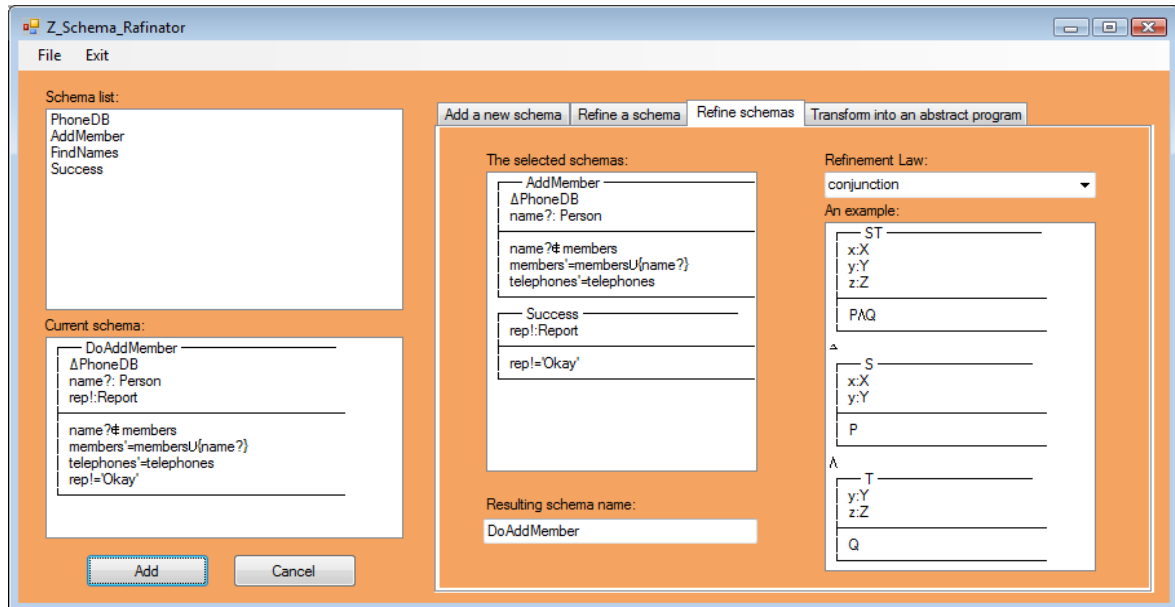


Figura 22. Conjuncția a două scheme Z

Aplicația are în principal scop didactic și a fost dezvoltată în C#. Programele abstracte obținute din schemele Z pot fi rafinate în continuare, după cum se poate observa din subcapitolul următor.

5.5. O aplicație care asistă procesul de rafinare a codului din specificații

În [Mih06a] s-a prezentat un mini instrument CASE, „Asistent rafinare”, care permite rafinarea asistată a codului din programe abstracte, instrument descris în cele ce urmează. Programele abstracte sunt precizate prin lista explicită a variabilelor care primesc valori w , predicatul de intrare sau precondiția, notat cu pre , și predicatul de ieșire sau postcondiția, notat cu pos : „ $w:[pre,pos]$ ”. Rafinarea presupunea trecerea de la o stare la alta prin aplicarea unei reguli de rafinare. Regulele care pot fi aplicate în procesul asistat de rafinare sunt: regula atribuirii, regula alternanței, regula compunerii secvențiale și regula iterației.

Aplicația oferă următoarele facilități: utilizarea automată a regulilor de mai sus pentru un subprogram abstract, după ce utilizatorul a identificat toate elementele noi necesare procesului de rafinare. Utilizatorii pot modifica însă și manual liniile de cod. Întregul proces de rafinare este memorat într-un fișier istoric „log.txt”, ceea ce permite reutilizarea procesului de rafinare pentru cazuri similare. Totodată se poate salva forma curentă a programului sau se poate încărca un program salvat anterior.

Această aplicație ușurează munca dezvoltatorilor, mai ales în cazul unor programe de dimensiuni medii sau mari, care necesită aplicarea repetată a regulilor rafinării.

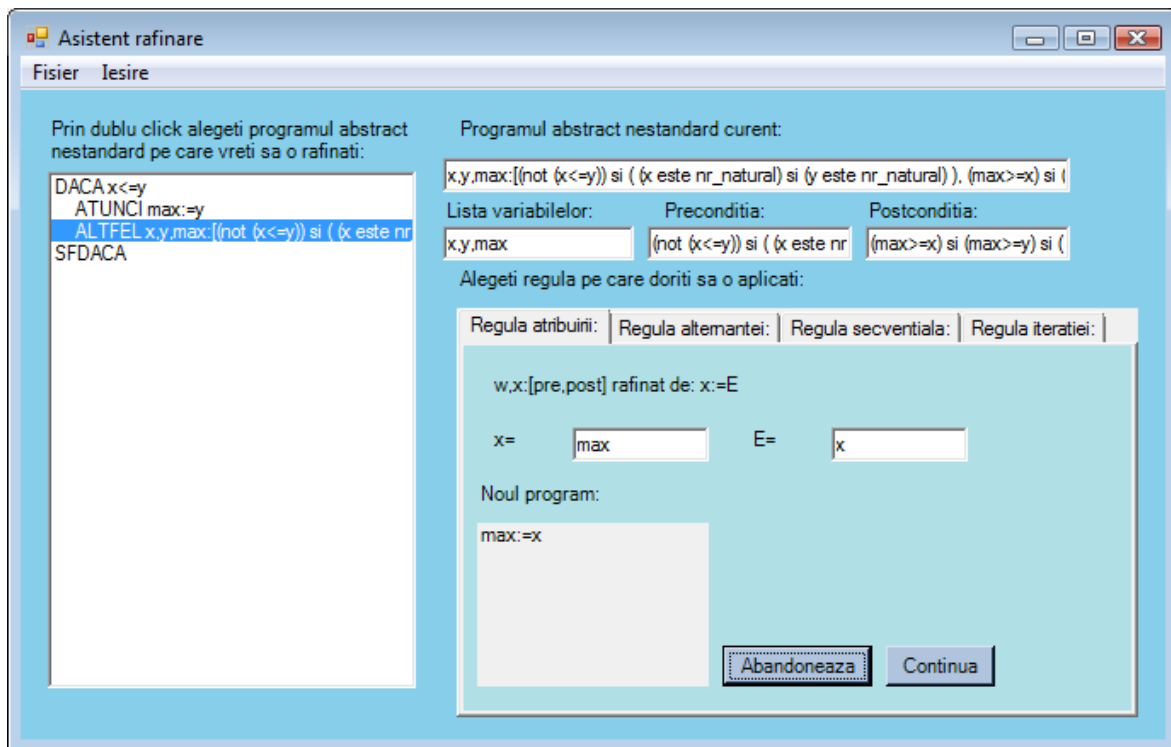


Figura 23. Fereastra de bază a Asistentului rafinării

5.6. Simplificarea codului prin prelucrarea automată a expresiilor condiționale

Un alt mini instrument CASE este propus în [Mih06b] și prezentat într-un subcapitol din [Lup08/09] favorizează îmbunătățirea calității codului prin aplicarea simplificării funcțiilor booleene asupra structurilor condiționale complicate existente în cod. Aceasta deoarece fiecare expresie condițională simplă are doar două valori de adevăr: adevărat, care poate fi notat cu 1 , respectiv fals, notat cu 0 . Dacă se notează toate expresiile condiționale simple cu variabile x_i , $i = \overline{1, n}$, unde n reprezintă numărul condițiilor, expresia condițională complexă poate fi înlocuită de funcția booleană $f(x_1, x_2, \dots, x_n)$. Așadar problema simplificării unei expresii condiționale se reduce la simplificarea unei funcții booleene, proces care a fost automatizat în unealta numită BOOFS (BOOlean Function Simplificator).

Metoda folosită pentru specificare este metoda lui Quine-McClusky combinată cu metoda lui Moasil [Tăt99, Lup08/09]. Această metodă se aplică însă doar formei canonice disjunctive a unei funcții booleene, formă la care se ajunge după aplicarea unui proces de normalizare.

Pentru a putea aprecia eficient îmbunătățirile obținute prin aplicarea metodei de simplificare, unealta BOOFS permite evaluarea automată a câtorva metrici care se bazează pe numărul operatorilor (conectivelor) și al operanzilor, metrici care sunt definite în continuare: *numărul de operanzi distincți*, *numărul total de operanzi*, *numărul de operatori distincți*, *numărul total de operatori*, *suma priorităților operatorilor distincți*, *suma priorităților tuturor operatorilor*.

Valoarea acestor metrici este un număr natural, care ar trebui să descrească în cazul simplificării. Pentru metricile 1, 3 și 5 descreșterea valorii metricii este evidentă, însă pentru metricile 2 și 4, în unele cazuri valoarea lor crește. Cea mai bună metrică pare a fi metrica 6. Prioritățile folosite în metricile 5 și 6 sunt de la 1 la 5, corespunzătoare

operatorilor: \neg , \wedge , \vee , \rightarrow , \leftrightarrow . Aceste priorități sunt folosite ca și criteriu și în codul aplicației BOOFS.

De exemplu, pentru expresia condițională $((x < y) \text{ and } (y \leq z) \text{ and } (z < 5)) \text{ or } ((x < y) \text{ and } (y \geq 5) \text{ and } (z < 5))$, care s-a simplificat la $((x < y) \text{ and } (z < 5))$, valorile celor 6 metrice au scăzut de la 3, 6, 3, 6, 6, 12, la 2, 2, 1, 1, 2, 2.

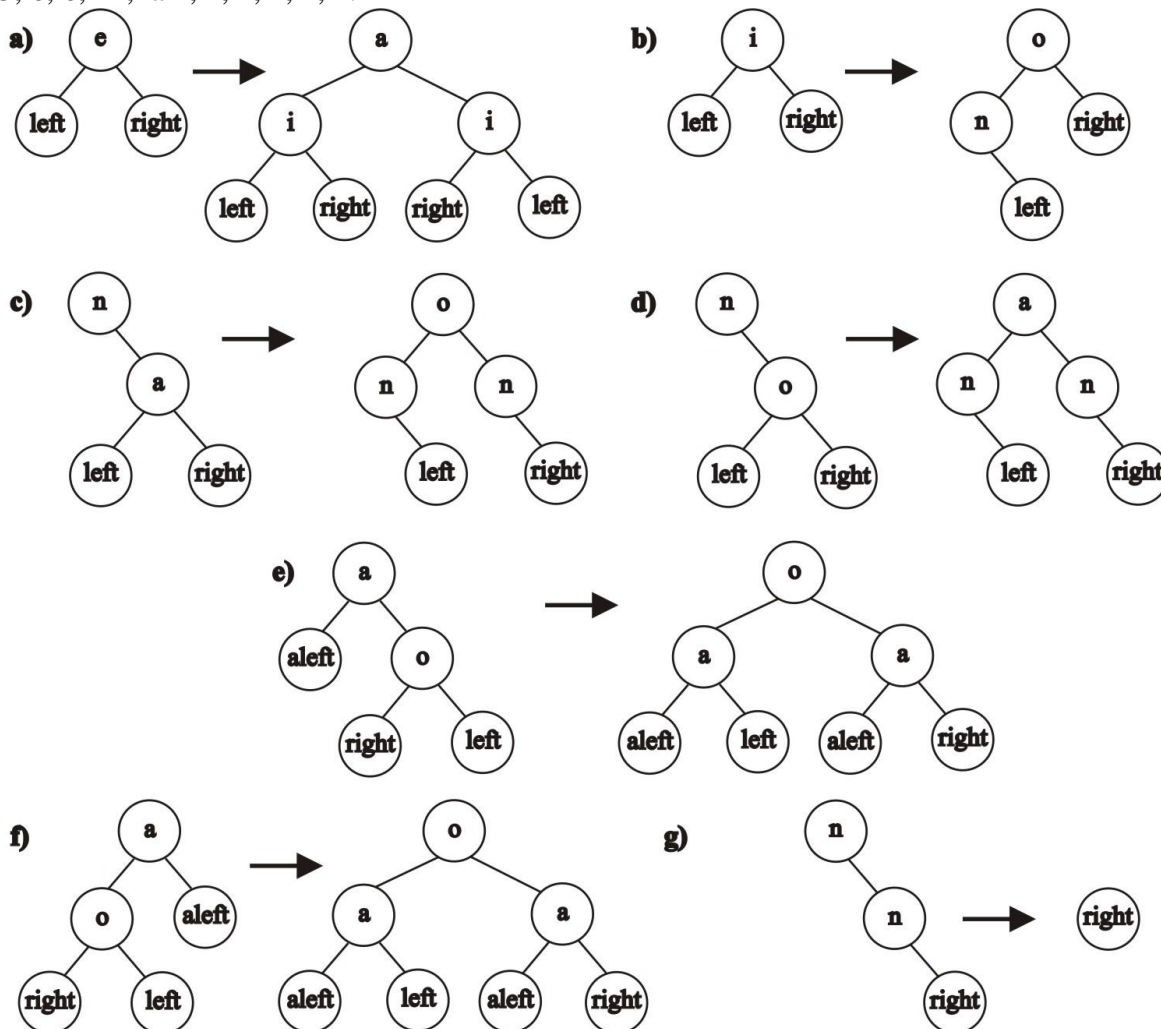


Figura 24. Operațiile de aducere a funcției booleene la forma disjunctivă pe arbore binar: a) asigură eliminarea echivalenței, b) elimină implicația, c) și d) sunt legile lui DeMorgan, e) și f) reprezintă distributivitatea (la dreapta respectiv la stânga), iar g) reprezintă legea dublei negații.

Aplicația BOOFS a fost realizată în Java. Execuția începe prin deschiderea unui fișier de cod sursă, din care se va putea selecta o expresie condițională care va fi simplificată. Deoarece aplicația nu poate identifica automat expresiile condiționale simple, ele vor trebui indicate manual. Simultan, întrucât aplicația funcționează indiferent de limbajul în care este scris codul sursă, utilizatorul va indica modul în care apar cele cinci conective logice. Conectivile se memorează prin litere: $n(\neg)$, $a(\wedge)$, $o(\vee)$, $i(\rightarrow)$, $e(\leftrightarrow)$, iar operanzii prin x_1 , x_2 , x_3 și așa mai departe. Funcția booleană obținută trebuie normalizată. Normalizarea s-a efectuat pe arbore, prin aplicarea operațiilor din figura 24, într-un algoritm back-tracking. Pentru a transforma funcția din forma infixată în arbore, se va transforma mai întâi în forma poloneză postfixată prin intermediul unei stive. Tot utilizând o stivă se va transforma forma poloneză postfixată în arbore. După aducerea arborelui la

forma sa disjunctivă, se va obține forma infixată disjunctivă a funcției booleene printr-o parcurgere în ordine a arborelui binar.

Se construiește mulțimea suport, care conține n -uple cu componente egale cu 0 sau 1. Metoda lui Quine-McClusky este folosită pentru obținerea mulțimii monoamelor maximale. În n -uplele corespunzătoare s-a folosit cifra 2 pentru a marca variabilele simplificate. Metoda lui Moisiil este folosită pentru a identifica forma simplificată. Metoda lui Moisiil se bazează pe transformarea unei forme normale conjunctive într-una disjunctivă, transformare realizată din nou pe arbore.

VI. Concluzii

În prezent nevoia de corectitudine este din ce în ce mai pregnantă. Dacă un produs soft nu își respectă specificațiile, el nu va fi acceptat de client. Existența unor specificații complete, corecte și clare este esențială. Astfel de specificații sunt specificațiile formale.

Din păcate învățarea și utilizarea unei metode formale nu este ușoară și necesită timp. Este de un real ajutor existența unei unelte capabile să asiste elaborarea unor astfel de specificații. O astfel de unealtă, care permite construcția specificațiilor Z, compunerea lor și transformarea lor în programe abstracte este prezentată în această lucrare.

Există mai multe tipuri de specificații formale, de exemplu unele sunt orientate pe procese, altele pe date. Sistemele complexe necesită o specificație formală care să surprindă toate aspectele aplicației, necesitând o abordare multiformală a specificației. Specificațiile formale pot fi și reutilizate cu succes, mai ales în cazul dezvoltării unei noi versiuni a unui program soft a cărui specificație formală există deja.

Pornindu-se de la specificații formale, se poate genera cod corect prin rafinare pe bază de reguli. Deși procesul de specificare formală și rafinare poate părea destul de greu de realizat, costurile totale ale dezvoltării sunt reduse. Existența unor unelte care să asiste și să automatizeze o parte din acest proces este binevenită. O astfel de unealtă este prezentată în această lucrare. O altă unealtă realizată simplifică prin metode logice expresiile condiționale din cod.

Pe de altă parte s-ar putea ca specificarea formală a cerințelor clienților să nu poată fi realizată de aceștia împreună cu dezvoltatorii. În acest caz dezvoltatorii pot utiliza aplicații de prelucrare a limbajului natural pentru a-și ușura munca, cum ar fi aplicații capabile să dezambigueze cuvintele, să împartă textul în segmente sau să realizeze sumarul.

Ontologiile sunt candidate excelente pentru o prelucrare mai complexă a limbajului natural, cum e cea de selectare a cerințelor pentru obținerea specificațiilor. Totuși, ontologia folosită în acest scop trebuie aleasă cu grijă.

Lucrarea de față prezintă diferite aspecte ale utilizării metodelor formale în dezvoltarea corectă a programelor, dar conține și o parte importantă de prelucrare a limbajului natural și ontologii. Însă toate metodele și aplicațiile prezentate au ca scop facilitarea utilizării metodelor formale în dezvoltarea corectă a programelor.

VII. Bibliografie

- [Abr80] Abrial J.-R., *The Specification Language Z: Syntax and Semantics*, Oxford University Computing Laboratory, Programming Research Group, Oxford, U.K., aprilie 1980.
- [Alf85] Alford M., *SREM at the Age of Eight: The Distributed Computer Design System*, IEEE Computer, Vol. 18, Nr. 4, aprilie 1985, pg. 36-46.
- [Ast02] Astesiano E., Bidot M., Krieg-Brückner B., Mosses P.D., Sannella D. și Tarlecki A., *CASL: The Common Algebraic Specification Language*, Theoretical Computer Science, 286(2), 2002, pg. 153-196.
- [Bac00] Back T., Fogel D.B. și Michalewicz Z. (editori), *Evolutionary Computation: Basic Algorithms and Operators*, Vol. 1 și *Evolutionary Computation: Advanced Algorithms and Operators* Vol. 2, Institute of Physics Publishing, Philadelphia, PA, 2000.
- [Bal85] Balzer R., *The establishment and use of sound engineering principles in order to obtain economically software that is reliable and works efficiently on real machines.*, NATO Science Committee Conference, 1969.
- [Ban03] Banerjee S. și Pedersen T., *Extended Gloss Overlaps as a Measure of Semantic Relatedness*, Proceedings of the Eighteenth International Joint Conference on Artificial Intelligence, Acapulco, Mexic, 9-15 august 2003, pg. 805-810.
- [Bar99] Barzilay R. și Elhadad M., *Using lexical chains for Text summarization*, editori J. Mani și M. Maybury, Advances in Automatic Text Summarization, MIT Press, 1999.
- [Ber98] Berners-Lee T., Fielding R. T. and Masinter, L., *Uniform resource identifiers (URI): Generic syntax*. RFC 2396, IETF, 1998.
- [Bid91] Bidoit M., Kreowski H.-J., Lescanne P., Orejas F. și Sannella D., *Algebraic System Specification and Development*, Lecture Notes in Computer Science, Vol. 501, Springer-Verlag, 1991.
- [Bra85] Brachman R. și Schmolze J., *An Overview of the KL-ONE Knowledge Representation System*, Cognitive Science, vol. 9, nr. 2, 1985, pg. 171-216, <http://nlp.shef.ac.uk/kr/papers/klone.ps>.
- [Bra05] Brank J., Grobelnik M. și Mladenić D., *A Survey of Ontology Evaluation Techniques*, Proceedings of the Conference on Data Mining and Data Warehouses (SIKDD 2005).
- [Bur95] Burgess C.J., *The Role of Formal Methods in Software Engineering Education and Industry*, Proceedings of the 4th Software Quality Conference <http://www.cs.bris.ac.uk/Tools/Reports/Abstracts/1995-burgess-3.html>.
- [Che76] Chen P., *The Entity-Relationship Model – Towards a Unified View of Data*, ACM Transactions on Database Systems, Vol. 1, Nr. 1, martie 1976, pg. 9-36.
- [Che06] Chen R., Lee Y. și Pan R., *Adding New Concepts on the Domain Ontology Based on Semantic Similarity*, International Conference on Business și Information (BAI 2006), Singapore, 12-14 iulie 2006, <http://bai2006.atissr.org/CD/Papers/2006bai6169.pdf>.
- [Cim05] Cimiano P. și Voelker J., *Text2Onto - A Framework for Ontology Learning and Data-driven Change Discovery*, Proceedings of the 10th International Conference on Applications of Natural Language to Information Systems (NLDB), Alicante, Spania, 2005.

- [Cio04] Ciobotariu-Boer, V. și **Mihș, A.D.**, *The Multiformalism Approach in Software Specifications*, Proceedings of the Symposium „Zilele Academice Clujene”, Computer Science Section, Faculty of Mathematics and Computer Science, "Babes-Bolyai" University, Cluj-Napoca, România, Litografia Universității "Babeș-Bolyai", Editor: Prof. dr. Militon Frențiu, 2004, pg. 21-26.
- [Cla96] Clarke E.M. și Wing J.M., *Formal Methods: State of the Art and Future Directions*, ACM Computing Surveys, 1996, <http://www.cs.cmu.edu/> CMU-CS-96-178.
- [Cor05] Corley C. și Mihalcea R., *Measuring the semantic similarity of texts*, Proceedings of the ACL Workshop on Empirical Modelling of Semantic Equivalence and Entailment, Ann Arbor, iunie 2005, pg. 13-18.
- [Dah06] Dahab M., Hassan H. și Rafea A., *TextOntoEx: Automatic Ontology Construction from Natural English Text*, International Conference on Artificial Intelligence and Machine Learning (AIML-06), Sharm El Sheikh, Egipt, 13-15 iunie 2006, http://www.icgst.com/con06/aiml06/Final_Articles/P1120615104.pdf.
- [DeM78] DeMarco T., *Structured Analysis and System Specification*, Yourdon Press, New York, 1978.
- [Dic45] Dice L.R., *Measures of the amount of ecologic association between species*, Ecology, 26(3), 1945, pg. 297–302.
- [Dil99] Diller A., *Z: An Introduction to Formal Methods*, a doua ediție, John Wiley & Sons, aprilie 1999.
- [Dij75] Dijkstra E.W., *Guarded commands, nondeterminacy and formal derivation of programs*, Comm. ACM, Vol. 18(1975), Nr. 8, pg. 453-457.
- [Doa03] Doan A., Madhavan J., Dhamankar R., Domingos P. și Halevy A., *Learning to match ontologies on the Semantic Web*, The VLDB Journal – The International Journal on Very Large Data Bases, vol. 12, issue 4, noiembrie 2003, pg. 303-319, <http://www.cs.washington.edu/homes/pedrod/papers/vldb04.pdf>.
- [Dro89] Dromey G., *Program Derivation. The Development of Programs from Specifications*, Addison Wesley, 1989.
- [Ehr83] Ehrig H., Fey W. și Hansen H., *ACT ONE: An algebraic specification language with two levels of semantics*, Technical Report No. 83-03, Technische Universität Berlin, 1983.
- [Eij89] van Eijk P.H.J., Vissers C.A. și Diaz M., *The formal description technique LOTOS*, Elsevier Science Publisher B.V., 1989.
- [Eng05] Engel L., Jaeger M. și Mühl G., *Search and Evaluation of Ontologies for Semantic Web Services in the Internet*, IADIS International Conference WWW/Internet 2005 Lisabona, Portugalia, 19-22 octombrie 2005, http://www.iadis.net/dl/final_uploads/200507C050.pdf.
- [Euz04] Euzenat J. și Valtchev P., *Similarity-based ontology alignment in OWL-lite*, Proc. 16th European Conference on Artificial Intelligence (ECAI 2004), 2004, pg. 333–337.
- [Euz07] Euzenat J. și Shvaiko, P., *Ontology Matching*, Springer, New York, 2007.
- [Flo67] Floyd R.W., *Assigning meanings to programs*, Proc. Symposium in Applied Mathematics, Schwartz J.T. (Ed.), Am.Math.Soc., Vol. 19, 1967, pg. 19-32.
- [Fre06] Frențiu M., Pop H.F., *Fundamentals of Programming*, Presa Universitară Clujeană, Cluj-Napoca, 2006.
- [Fut85] Futatsugi F., Goguen J.A., Jouannaud J.P. și Meseguer J., *Principles of OBJ2*, Annual Symposium on Principles of Programming Languages, Proceedings of 12th

- ACM SIGACT-SIGPLANT symposium on Principles of programming languages, New Orleans, Louisiana, Statele Unite ale Americii, 1985, pg. 52–66.
- [Gab06] Gabbar H.A. (editor), *Modern Formal Methods and Applications*, Springer, 2006.
- [Gal96] Galloway A. și Stoddart B., *Integrated Formal Methods*, Research Report, Institut de Recherche en Informatique de Nantes, 1996.
- [Gan79] Gane C. și Sarson T., *Structured Analysis: Tools and Techniques*, Prentice Hall, Englewood Cliffs, NJ, 1979.
- [Gan05b] Gangemi A., Catenacci C., Ciaramita M. și Lehmann J., *Ontology Evaluation and Validation. An integrated formal model for the quality diagnostic task*, Technical report, ISTC-CNR, Lab. for Applied Ontology, http://www.loa-cnr.it/Files/OntoEval4OntoDev_Final.pdf.
- [Gao09] Gao X., *The Design and Implementation of Z Language Editor*, Algorithms and Architectures for Parallel Processing, vol. 5574/2009, Springer, iulie 2009, pg. 684-692.
- [Gar93] Garland S.J., Guttag J.V. și Horning J.J., *An Overview of Larch*, Lecture Notes in Computer Science, Vol. 693, Functional Programming, Concurrency, Simulation and Automated Reasoning: International Lecture Series 1991-1992, McMaster University, Hamilton, Ontario, Canada, 1993, pg. 329-348.
- [Gil62] Gill A., *Introduction to the Theory of Finite-state Machines*, McGraw-Hill, 1962.
- [Gli05] Glickman, O., Dagan, I. și Koppel, M., *Web Based Probabilistic Textual Entailment*, Proceedings of the PASCAL Challenges Workshop on Recognizing Textual Entailment, 2005, pg. 33-36.
- [Gri81] Gries D., *The Science of Programming*, Texts and Monographs in Computer Science, 1st ed. 1981. 5th printing, Springer Verlag, Berlin, 1981.
- [Gru93] Gruber T., *A Translation Approach to Portable Ontology Specifications*, Knowledge Acquisition, vol. 5, nr. 2, 1993, pg. 199–220.
- [Har97] Harabagiu S. și Moldovan D., *TextNet – a textbased intelligent system*, Natural Language Engineering, 3(2), 1997, pg. 171-190.
- [Hea93] Hearst M., *TextTiling: A Quantitative Approach to Discourse Segmentation*, Technical Report 93/24, University of California, Berkeley, 1993.
- [Hei96] Heitjmeier C.L., Jeffords R.D. și Labaw B.G., *Automated Consistency Checking of Requirements Specifications*, ACM Trans. on Software Engineering and Methodology, 5(1996), nr. 3, pg. 231-261.
- [Hoa85] Hoare C.A.R., *Communicating Sequential Processes*, Prentice Hall International, Englewood Cliffs, NJ, 1985.
- [Hol96] Holloway C.M., *Why Engineers Should Consider Formal Methods*, NASA Langley Research Center, <http://shmesh.larc.nasa.gov/cmh.html>, 1996.
- [Hov03] Hovy E., *Text summarization*, The Oxford Handbook of Computational Linguistics, editor Mitkov R., Oxford University Press, Chapter 32, 2003.
- [Jac901] Jaccard P., *Distribution de la flore alpine dans le bassin des Dranses et dans quelques régions voisines*, Bulletin de la Société Vaudoise des Sciences Naturelles 37, 1901, pg. 241-272.
- [Joh73] Johnson S.C. & Kernighan B.W., *The Programming Language B*, Technical Report CS TR 8, Bell Labs, ianuarie 1973, Prentice Hall, Englewood Cliffs, NJ, 1986.
- [Jon86] Jones C.B., *Systematic Software Development Using VDM*, Prentice Hall, Englewood Cliffs, NJ, 1986.

- [Kam87] Kampen G. R., *An Eclectic Approach to Specification*, Proceedings of the Fourth International Workshop on Software Specification and Design, Monterey, CA, aprilie 1987, pg. 178-182.
- [Kaw04] Kawtrakul A., Suktarachan M. și Insombut A., *Automatic Thai Ontology Construction and Maintenance System*, Workshop on Papillon 2004, Grenoble, Franța, <http://www.moac.go.th/knowledgebase/uploadfile/42808973.pdf>
- [Kle03] Klein D. și Manning C. D., *Fast Exact Inference with a Factored Model for Natural Language Parsing*, Advances in Neural Information Processing Systems 15 (NIPS 2002), Cambridge, MA: MIT Press, 2003, pg. 3-10.
- [Kõu06] Kõulekov M. și Magnini B., *Tree Edit Distance for Recognizing Textual Entailment: Estimating the Cost of Insertion*, Proceedings of the Second PASCAL Challenges Workshop on Recognizing Textual Entailment, Veneția, Italia, 2006.
- [Les86] Lesk M., *Automatic sense disambiguation using machine readable dictionaries: how to tell a pine cone from an ice cream cone*, Proceedings of the 1986 SIGDOC Conference, Association for Computing Machinery, New York, 1986, pg. 24–26.
- [Lew04] Lewis D.D., Yang Y., Rose T. și Li F., *RCV1: A New Benchmark Collection for Text Categorization Research*, Journal of Machine Learning Research, Vol. 5, 2004, pg. 361-397.
- [Luc85] Luckham D.C. și von Henke F.W., *An Overview of Anna, a Specification Language for Ada*, IEEE Software, Vol. 2, Nr. 2, martie 1985, pg. 9-22.
- [Lup08] Lupea M., **Mihîș A.D.**, *Logici clasice și Circuite logice. Teorie și exemple*, S.C. Albastră Casa de Editura S.R.L., Cluj-Napoca, 2008, 223 pg.
- [Lup09] Lupea M. și **Mihîș A.D.**, *Logici clasice și Circuite logice. Teorie și exemple*, S.C. Albastră Casa de Editura S.R.L., Cluj-Napoca, 2009, 223 pg.
- [Mae03] Maedche, A., Neumann, G. și Staab, S., *Bootstrapping an Ontology-based information extraction system*, in Intelligent exploration of the web, Editori: Piotr S. Szczepaniak, Javier Segovia, Janusz Kacprzyk, Lotfi A. Zadeh, Physica-Verlag GmbH Heidelberg, Germania, 2003, pg. 345 - 359.
- [Mar97] Marcu D., *From discourse structure to text summaries*, Proceedings of the ACL/EACL '97 Workshop on Intelligent Scalable Text Summarization, Madrid, Spania, pg. 82-88.
- [Mar08] de Marneffe M.-C., și Manning C. D., *The Stanford typed dependencies representation*, COLING Workshop on Cross-framework and Cross-domain Parser Evaluation, Manchester, United Kingdom, 2008, <http://nlp.stanford.edu/pubs/dependencies-coling08.pdf>.
- [May09] Maynard D., Funk A. and Peters W., *Using Lexico-Syntactic Ontology Design Patterns for ontology creation and population*, WOP 2009 – ISWC Workshop on Ontology Patterns, Washington, 2009.
- [Mih08a] Mihăilă A.A., **Mihîș A.D.** și Mihăilă C.F., *A Genetic Algorithm for Logical Topic Text Segmentation*, International Conference on Digital Information Management, IEEE Computer Society Press, 978-1-4244-2917-2/08, IEEE Xplore, 2008, pg. 500-505.
- [Mih05] **Mihîș A.D.**, *Formal Specification Reuse*, Proceedings of Symposium „Colocviul Academic Clujean de Informatică”, 2005, pg. 202-206.
- [Mih06a] **Mihîș A.D.**, *An Application that Assist StepWise Refinement*, Proceedings of Symposium „Zilele Academice Clujene”, 2006, pg. 143-147.

- [Mih06b] **Mihiş A.D.**, Chisăliţă-Creţu C., Mihăilă C.A., Şerban C.A., *BOOFS-A Tool That Supports Simplifying Conditional Expressions using Boolean Functions Simplification Methods*, Studii şi Cercetări Ştiinţifice, Seria Matematică, 16 (2006), Supplement, Proceedings of ICMI 45, editori: Mocanu Marcelina & Nimineţ Valer, Bacău, 18-20 Septembrie, 2006, pg. 493-502.
- [Mih07] **Mihiş A.D.**, *Chain Algorithm used for Part of Speech Recognition*, Mathematical Reviews, www.ams.org/mathscinet/, Studia Universitatis "Babes-Bolyai", Informatica: KEPT 2007, pg. 89-95.
- [Mih08b] **Mihiş A.D.**, *A Simple Ontology based on Text Entailment Directional Relationship*, Proceedings of KM-03 Knowledge Management: Projects, Systems, and Technologies, 23-25 octombrie 2008, Bucureşti, România, Editura ASE Bucureşti, ISBN: 978-606-505-124-9, pg. 531-534.
- [Mih08c] **Mihiş A.D.**, *Natural Language Processing Methods Used in Requirement Analysis*, „Zilele Informaticii Economice Clujene“, Mediamira Science Publisher, Cluj-Napoca, editori: prof. dr. Niţchi Ştefan şi alţii, 2008, pg. 251-258.
- [Mih08d] **Mihiş, A.D.**, *Various Types of Logical Text Tiling*, Zilele Academice Clujene, Presa Universitară Clujeană, editor: Prof. dr. Frenţiu Militon, 2008, pg. 129-133.
- [Mih10a] **Mihiş A.D.**, *A Tool for Refinement of Z Schemas*, Proceedings of Symposium „Zilele Academice Clujene”, 2010, pg. 64-67.
- [Mih10b] **Mihiş A.D.**, *The Evaluation of Ontology Matching versus Text*, Informatica Economică/Economy Informatics, Categ. CNCSIS B+, vol. 14, nr. 4, 2010, pg. 147 – 155.
- [Mih10c] **Mihiş A.D.**, *Ontology Learning from Text Based on the Syntactic Analysis Tree of a Sentence*, The 5th International Conference on Virtual Learning (ICVL - 2010), Târgu-Mureş, Editura Universităţii Bucureşti, Editor: Vlada Marian, Albeanu Grigore, Popovici Dorin Mircea, ISSN: 1844-8933, 1842-4708, <http://c3.icvl.eu/2010>, 2010, pg. 128-134, articol ce a obţinut premiul special „Intel@Education”.
- [Mih10d] **Mihiş A.D.**, *Similarity Used In The Identification Of The Need To Refactoryze An Ontology*, Knowledge Management: Projects, Systems and Technologies (KM – conference), "Carol I" National Defence University Publishing House, Editor: Toma Pleşanu, Constanţa Bodea, Luiza Kraft, 2010, pg. 23-28.
- [Mih10e] **Mihiş A.D.**, *The Role of Disambiguation in Ontology Evaluation*, The Seventh International Conference on Applied Mathematics (ICAM7), Minisymposium 6- Software Engineering - Principles and Practices, Baia Mare, 1-4 septembrie 2010.
- [Mih10f] **Mihiş A.D.**, *Ontology Assisted Formal Specification Extraction from Text*, Studia Universitatis Babeş-Bolyai, Informatica, Vol. 55, No. 4, 2010, Cluj-Napoca, pg. 103-113.
- [Mil80] Milner R., *A Calculus of Communicating Systems*, Springer-Verlag, 1980.
- [Mon01] Monz C. şi de Rijke M., *Light-Weight Entailment Checking for Computational Semantics*, Proceedings of the third workshop on inference in computational semantics (IcoS-3), editori Blackburn P. şi Kohlhase M., 2001.
- [Mor06] Morita T., Fukuta N., Izumi N. şi Yamaguchi T., *DODDLE-OWL: A Domain Ontology Construction Tool with OWL*, Proceedings of the 1st Asian Semantic Web Conference, Lecture Notes in Computer Science, vol. 4185, Beijing, China, 2006, pg. 537-551,
<http://iws.seu.edu.cn/resource/Proceedings/ASWC/2006/papers/4185/41850537.pdf>.

- [Mor91] Morris J., Hirst G., *Lexical Cohesion Computed by Thesaural Relations as an Indicator of the Structure of Text*, Computational Linguistics, Vol. 17, Nr. 1, 1991, pg. 21-48.
- [Mor90] Morgan C., *Programming from Specifications*, Programming Research Group, University of Oxford, Prentice Hall International (UK), 1990.
- [Nel79] Nelson W.F. și Kucera H., *Brown corpus manual*, Dept. of Linguistics, Brown University 1979, <http://icame.uib.no/brown/bcm.html>.
- [Oră06] Orășan C., *Comparative evolution of modular automatic summarization systems using CAST*, Teză de doctorat, University of Wolverhampton, UK, 2006.
- [Per96] Periyasamy K., Chidambaram J., *Software Reuse Using Formal Specification of Requirements*, Proceedings of the 1996 conference of the Centre for Advanced Studies on Collaborative Research, IBM Press, Toronto, Canada, 1996 <http://www.cs.ubc.ca/local/reading/proceedings/cascon96/pdf/periyaasc.pdf>
- [Pet62] Petri C.A., *Kommunikation mit Automaten*, Ph. D. Dissertation, University of Bonn, Germany, 1962.
- [Pol09] Pollock J. T., *Semantic Web for Dummies*, Wiley Publishing, Indianapolis, 2009.
- [Rad02] Radev D., Hovy E. și McKeown K., *Introduction to the Special Issues on Summarization*, Computational Linguistics, Vol. 28, Nr. 8, 2002, pg. 399-408.
- [Ric91] Richie D., *The Koan*, Zen Inklings. Some Stories, Fables, Parables and Sermons, New York and Tokyo: Weatherhill, 1991, pg. 25-27.
- [Ros85] Ross T. D., *Applications and Extensions of SADT*, IEEE Computer, Vol. 18, No. 4, aprilie 1985, pg. 25-34.
- [Rus95] Rushby J., *Model Checking and Other Ways of Automating Formal Methods*, <http://www.csl.sri.com/Reports/html/SQW95.html>.
- [San07] Sanchez D. și Moreno A., *Semantic disambiguation of taxonomies*, Proceeding of the 2007 Conference on Artificial Intelligence Research and Development, Frontiers in Artificial Intelligence and Applications, vol. 163, 2007, pg. 245-254.
- [Sil02] Silber H., McCoy K., *Efficiently computed lexical chains, as an intermediate representation for automatic text summarization*, Computational Linguistics, 28(4), 2002, pg. 487-496.
- [Sch05] Schach S.R., Vanderbilt University, *Object-Oriented and Classical Software Engineering*, sixth edition, McGraw-Hill, New York, 2005.
- [Smi99] Smith G., *The Object-Z Specification Language*, Kluwer Academic Publishers, 1999.
- [Ste02] Stevenson M., *Combining Disambiguation Techniques to Enrich an Ontology*, Proceedings of the ECAI 2002 Workshop on Machine Learning and Natural Language Processing for Ontology Engineering, Lyon, Franța, 2002.
- [Sto04] Stokes N., *Applications of Cohesion Analysis in the Topic Detection and Tracking Domain*, Teză de doctorat, Faculty of Science, National University of Ireland, Dublin, 2004.
- [Suk08] Suktarachan M., Thamvijit D., Rajbhandari S., Noikongka D., Mahasarakram P., Yongyuth P., Kawtrakul A. și Sini M., *Workbench with Authoring Tools for Collaborative Multilingual Ontological Knowledge Construction and Maintenance*, Proceedings of the Sixth International Language Resources and Evaluation (LREC'08), Marrakech, Maroc, 2008, pg. 2501-2508, http://www.lrec-conf.org/proceedings/lrec2008/pdf/624_paper.pdf.

- [Şer05] Şerban, A. C., **Mihiş, A.D.**, *Software Quality Assurance*, Proceedings of the Symposium „Zilele Academice Clujene”, 2005, pg. 207-212.
- [Tan09] Tan H. și Lambrix P., *Selecting an Ontology for Biomedical Text*, Mining Human Language Technology Conference, Proceedings of the Workshop on BioNLP, Association for Computational Linguistics, Boulder, Colorado, 4-5 iunie 2009, pg. 55-62, <http://www.aclweb.org/anthology/W/W09/W09-1307.pdf>.
- [Tăt99] Tătar D., *Bazele matematice ale calculatoarelor*, Litografia UBB, 1999.
- [Tăt07a] Tătar D., Şerban G., **Mihiş A.D.**, Lupea M., Lupşa D. și Militon F., *A Chain Dictionary Method for Word Sense Disambiguation and Applications*, Mathematical Reviews, www.ams.org/mathscinet/, Studia Universitatis "Babeş-Bolyai", Informatica: KEPT 2007, 6-8 iunie 2007, pg. 41-49.
- [Tăt07b] Tătar D., Şerban G., **Mihiş A.D.**, Mihalcea R., *Textual Entailment as a Directional Relation*, CALP 2007, INCOMA Ltd., editor: Orăşan Constantin și alții, 2007, pg. 53-58.
- [Tăt08a] Tătar D., **Mihiş A.D.** și Lupşa D., *Text Entailment for Logical Segmentation and Summarization*, BDI, Proceedings of the 13th international conference on Natural Language and Information Systems: Applications of Natural Language to Information Systems, Document Processing and Text Mining, <http://www.informatik.uni-trier.de/~ley/db/conf/nldb/nldb2008.html>, Lecture Notes in Computer Science, Vol. 5039, 2008, pg. 233-244.
- [Tăt08b] Tătar D., **Mihiş A.D.**, Şerban G., *Lexical Chains Segmentation in Summarization*, SYNASC, IeAT Technical Report, 08-11, BDI, 2008, Timișoara, pg. 95-101.
- [Tăt08c] Tătar D., **Mihiş A.D.**, Şerban G., *Top-down Cohesion Segmentation in Summarization*, BDI, <http://www.aclweb.org/anthology/W/W08/W08-2232.bib>, Research in Computational Semantics, vol. 1(2008), pg. 389-397.
- [Tăt08d] Tătar D., **Mihiş A.D.**, Şerban G., *Top-down Cohesion Segmentation in Summarization*, Step 2008, Veneția, Italia, 22-24 sept 2008, College Publications, editori: Bos J. și Delmonte R., 978-1-904987-93-2, BDI, 2008, pg. 145-151.
- [Tăt08e] Tătar D., Tămăianu-Morita E.S., **Mihiş A.D.**, și Lupşa D., *Summarization by Logic Segmentation and Text Entailment*, CICLing 2008, DBLP, <http://www.informatik.uni-trier.de/~ley/db/journals/index-r.html>, Research in Computing Science, Vol. 33(2008), pg. 15-26.
- [Tăt09a] Tătar D., **Mihiş A.D.**, Şerban G., Mihalcea R., *Textual Entailment as a Directional Relation*, Journal of Research and Practice in Information Technology, 41(2009), 1, pg. 53 – 64.
- [Tăt09b] Tătar, D., **Mihiş, A. D.**, Şerban, G., *Lexical Chains Segmentation in Summarization*, SYNASC 2008, IEEE Computer Society, Editor: Viorel Negru, Tudor Jebelean, Dana Petcu, Daniela Zaharie, 978-0-7695-3523-4, <http://www.computer.org/cps>, 2009, pg. 95-101.
- [Tăt09c] Tătar D., **Mihiş A.D.**, Lupşa D., Tămăianu-Morita E.S., *Entailment-Based Linear Segmentation in Summarization*, International Journal of Software Engineering and Knowledge Engineering, 19(2009), 8, pg. 1023 – 1038.
- [Tho95] Muffy T., *Formal methods and their role in developing safe systems*, www.iee.org.uk/PAB/SCS/wrkshop.htm, Workshop report, 20 martie 1995.
- [Tei77] Teichroew D. și Hershey E.A. III, *PSL/PSA: A Computer-Aided Technique for Structured Documentation and Analysis of Information Processing Systems*, IEEE Transactions on Software Engineering, SE-3, ianuarie 1977, pg. 41-48.

- [Tuc99] Tucker R. I., *Automatic summarising and the CLASP system*, Teză de doctorat, University of Cambridge, Computer Laboratory, 1999.
- [Wil06] Wilkinson M. și Good B., *Construction and Evaluation of OWL-DL Ontologies*, Microsoft Research Faculty Summit, Redmond, USA, 17-18 iulie 2006, http://research.microsoft.com/en-us/um/redmond/events/fs2006/agenda_mon.aspx.
- [Woo96] Woodcock J. și Davies J., *Using Z, Specification, Refinement and Proof*, Prentice Hall, 1996.
- [Yil07] Yildiz, B. and Miksch, S., *ontoX - a method for Ontology-driven information extraction*, Proceedings of the 2007 international conference on Computational science and its applications - Volume Part III, Workshop on CAD/CAM and web based collaboration (CADCAM 07), Springer-Verlag Berlin, Heidelberg 2007, pg. 660-673.
- [You79] Yourdon E. și Constantine L.L., *Structured Design: Fundamentals of a Discipline of Computer Program and Software Design*, Prentice Hall, Englewood Cliffs, NJ, 1979.
- [***Adr] <http://www.cs.ubbcluj.ro/adriana/Teaching.html>, decembrie 2010
- [***BrC] <http://www.archive.org/details/BrownCorpus>
- [***Por] http://ir.dcs.gla.ac.uk/resources/linguistic_utils/
- [***RTE] <http://www.pascal-network.org/Challenges/RTE/>
- [***Sem] http://www.gabormelli.com/RKB/SemCor_Corpus
- [***Sen] <http://www.senseval.org/>
- [***SPW] Probably the most widely used stopword list, <http://www.lextek.com/manuals/onix/stopwords1.html>
- [***STO] The Stanford Natural Language Processing Group, Stanford Parser, <http://nlp.stanford.edu:8080/parser/>
- [***STS] The Stanford Natural Language Processing Group, The Stanford Parser: A statistical parser, <http://nlp.stanford.edu/software/lex-parser.shtml>
- [***WHR] A Wikipedia article about Human Resources, http://en.wikipedia.org/wiki/Human_resources
- [***WnB] The definition of the noun business <http://www.wordnet-online.com/business.shtml>.
- [***WNt] <http://wordnetweb.princeton.edu/perl/webwn>
- [***WSM] http://en.wikipedia.org/wiki/Small_business