

UNIVERSITATEA POLITEHNICA TIMIȘOARA
FACULTATEA DE AUTOMATICĂ ȘI CALCULATOARE

REZUMAT TEZĂ DE DOCTORAT

cu titlul:

*Improving text accessibility and understanding
of domain-specific information*

*(Îmbunătățirea accesibilității textului
și a înțelegerii informației de specialitate)*

Autor: Vasile TOPAC

Conducător științific: prof. dr. ing. Vasile STOICU-TIVADAR

Teza acestei lucrări (ideea de bază):

Adaptarea textelor de specialitate prin identificarea și explicarea terminologiei (întâlnită atât în formă canonică cât și în formă derivată), împreună cu punerea la dispoziție a unor mijloace de a schimba aspectul textului, pot duce la o accesibilitate și înțelegere mai bună a mesajului pentru utilizatorii de rând.

Rezumat

Această teză explorează căi prin care se poate îmbunătăți accesibilitatea și înțelegerea informației textuale pentru utilizatorii de rând, luându-se în considerare diferite niveluri de limitare ale accesului (limbă, limbajul specializat și modul de prezentare <aspectul>).

Cea mai mare atenție este acordată dificultății de înțelegere a textului datorate limbajului de specialitate, alegându-se ca și domeniu de aplicare limbajul medical. Metode și instrumente pentru recunoașterea și explicarea terminologiei medicale, chiar și atunci când aceasta apare în formă derivată sunt prezentate. Recunoașterea termenilor în limbajul natural e bazată pe tehnici de tipul *fuzzy matching* combinate cu structuri de date specializate, dezvoltate de către autor. Tehnici pentru îmbunătățirea preciziei recunoașterii aproximative (*fuzzy matching*) bazate atât pe metode semi-supervizate sau autonome cât și pe metode ce au la bază inteligența umană (de exemplu platformele de *crowdsourcing*) sunt prezentate. Aceste metode de adaptare a textelor de specialitate sunt implementate în câteva scenarii de utilizare, rezultând o serie de instrumente. Impactul procesului de adaptare a limbajului medical a fost evaluat prin studii cu utilizatori, rezultatele confirmând utilitatea acestor metode și a serviciilor asociate.

În plus, mijloace de adaptare a modului de prezentare a textului (aspect) sunt explorate. Atât instrumente pentru adaptarea aspectului textului în general, cât și instrumente dedicate unor anumite tipuri de utilizatorilor (cum ar fi utilizatori cu disexie sau vedere slabă) sunt prezentate.

Toate instrumentele dezvoltate în această lucrare au fost proiectate acordându-se atenție deosebită unor factori precum accesibilitate, utilizabilitate și disponibilitate, ele fiind concepute să ruleze oriunde în mediul online, cu efort minimal.

Obiective

Această teză urmărește următoarele obiective:

- Creșterea gradului de **înțelegere a mesajului din textele de specialitate** pentru utilizatorii de rând.
Obiective derivate din acesta:
 - Recunoașterea cât mai precisă a termenilor de specialitate în limbajul natural
 - Îmbunătățirea traducerii textelor de specialitate prin adnotarea terminologiei
- **Îmbunătățirea accesibilității și a lizibilității textelor** prin dezvoltarea de metode și instrumente de adaptare a aspectului textului
 - Dezvoltarea de instrumente specializate pentru ușurarea procesului de citire la persoanele cu dislexie
- Proiectarea serviciilor de adaptare a textului amintite mai sus astfel încât acestea să prezinte un grad ridicat de accesibilitate, utilizabilitate și disponibilitate. De asemenea definirea unui model de accesibilitate universală a textului, și încadrarea serviciilor dezvoltate în acest model.

Structura tezei (capitole)

1. Introducerea

În această secțiune autorul prezintă motivația și actualitatea temei abordate: *accessibilitatea informației textuale*. Diferitele nivele de limitare a accesului la informația textuală sunt prezentate, acordându-se o atenție specială limbajului specializat și nivelului de prezentare. Sunt prezentate apoi obiectivele acestei teze, urmate de o trecere în revistă a principalelor metode și instrumente dezvoltate.

2. Stadiul curent

Stadiul curent al cercetărilor atât pe parte de limbaj specializat cât și la nivel de prezentare e expus. Ceretearea existentă e încadrată într-un model de accesibilitate universală a textului propus de autor. Sunt prezentate și servicii funcționale (din industrie) ce au ca și țintă subiecte comune cu această teză. Pe lângă studiile și proiectele de cercetare prezentate în acest capitol, fiecare din capitolele ce urmează prezintă mai în detaliu aspecte ale stadiului actual pentru respectivul subiect.

3. Limbaj specializat – studiu de limbaj și recunoaștere aproximativă (fuzzy matching)

Autorul prezintă un studiu al limbajului medical pentru limba română și engleză cu scopul identificării ratei de apariție a terminologiei de specialitate în forma canonică comparat cu forma derivată. Studiul relevă o incidență mare a terminologiei în forma derivată (mai mare în limba română decât în engleză), ceea ce confirmă nevoia de a folosi tehnici de recunoaștere aproximativă a termenilor. Apoi este prezentată structura de date proiectată de către autor, *FuzzyHashMap*, ca o soluție pentru identificarea rapidă a termenilor în limbajul natural. Aceasta este o extensie a structurii *HashMap* din limbajul Java. O serie de teste care ilustrează performanța structurii sunt apoi prezentate, împreună cu prezentarea altor scenarii în care structura a fost deja folosită.

4. Limbaj specializat – îmbunătățirea preciziei recunoașterii aproximative

Autorul vorbește în acest capitol despre faptul că, pe lângă beneficiile ce le aduce recunoașterea aproximativă, există și efecte secundare datorate naturii tolerante la erorare a acestor tehnici, și anume recunoașteri greșite, sau *fals positive*. Autorul prezintă apoi o serie de

tehnici ce au ca și scop creșterea preciziei recunoașterii aproximative. Sunt prezentate atât metode ce au la bază tehnici semi-supervizate sau autonome cât și tehnici ce folosesc inteligența umană colectivă. Tehnicile semi-supervizate includ a) antrearea sistemului pentru dezvoltarea unor dicționare de mapări incorecte ce sunt ulterior folosite pentru filtrarea mapărilor și b) generarea unui model al mapărilor aproximative și derivarea unui tipar de hashing din acesta, pentru a fi integrat în structura *FuzzyHashMap*.

Tehnicile ce au la bază inteligența umană folosesc feedback agregat de la utilizatorii sistemului propus sau integrează folosirea platformelor de *crowdsourcing*. Teste cu utilizatori ai sistemului (făcute prin chestionare) și teste cu persoane din spatele platformei de *crowdsourcing* Amazon Mechanical Turk sunt apoi prezentate. Îmbunătățirea preciziei este apoi măsurată prin folosirea metricilor specifice, iar rezultatele arată o îmbunătățire semnificativă a acesteia.

5. Accessibilitatea limbajului specializat – scenarii de utilizare și aplicații

Acest capitol prezintă o serie de scenarii de utilizare a serviciilor de adnotare a terminologiei din limbaje specializate. Următoarele instrumentele dezvoltate de către autor sunt prezentate:

- *text4all terminology interpreter*: instrumentul de bază al acestei secțiuni, folosit pentru a explica terminologia medicală din pagini web existente cu scopul îmbunătățirii înțelegerii mesajului și a facilitării procesului de *patient empowerment*,
- *text4all ITS Tagger*: instrument de adnotare a terminologiei bazat pe standartul ITS 2.0 (ce are ca și scop, printre altele, îmbunătățirea traducerii textelor de specialitate).
- integrarea serviciului de adaptare a limbajului în proiectul de tele-asistență *TELEASIS*
- *text4all term analysis*: instrument de analiză a limbajului specializat din mediul online.

Impactului folosirii instrumentului *text4all terminology interpreter* pentru îmbunătățirea înțelegerii textelor medicale a fost evaluat prin teste cu utilizatori reali. Testele au confirmat utilitatea acestui instrument. Alte aspecte și implicații corelate cu adaptarea textelor cu specific medical sunt discutate.

6. Accessibilitatea textului la nivel de prezentare

În acest capitol autorul explorează limitările accesării (citirii) textului cauzate de modul de prezentare al acestuia, mai exact de aspectul textului. Elemente precum dimensiunea textului, culori, fonturi, spațiere și altele sunt luate în considerare. Instrumente ce permit adaptarea aspectului textului sunt prezentate, și modul în care acestea pot fii utile diferitelor tipuri de utilizatori sunt explorate. Un instrument dedicat utilizatorilor cu dislexie, dezvoltat de către autor în colaborare cu cercetători specializați pe acest tip de probleme, pe baza unor studii ce au folosit tehnici de eye tracking făcute cu cititori cu dislexie, este prezentat. Alte instrumente și aspecte ale adaptării textului la nivel de aspect sunt discutate și analizate.

7. Aspecte asupra modului de proiectare a serviciilor propuse

Anumite aspecte ale modului de proiectare și a modurilor de folosire și interacțiune ale serviciilor dezvoltate sunt discutate și analizate. O particularitate cheie a serviciilor e disponibilitatea acestora, ele fiind proiectate ca și mediatori online (denumite și proiecte de *transcoding*), funcționează direct în browser, independent de platformă și nu necesită instalare sau drepturi speciale. Diferite moduri de interacționare cu instrumentele dezvoltate, precum interacțiunea prin interfața grafică sau interacțiunea direct prin adresa URL (aceasta fiind proiectată și optimizată pentru o utilizabilitate cât mai sporită) sunt prezentate. Instrumentele dezvoltate sunt amplasate într-un model reprezentând o viziune de ansamblu a accesibilității informației textuale. Alte aspecte legate de acest model sunt discutate.

8. Concluzii

E expusă o scurta trecere în revistă a celor enunțate în teză. Principalele contribuții revendicate sunt listate, urmate de o prezentare a serviciilor rezultate, accentul punându-se pe contribuțiile aduse de aceste soluții. Câteva direcții de viitor sunt de asemenea prezentate.

Contribuții

Contribuții majore:

- Dezvoltarea unor metode și instrumente asociate de recunoaștere și explicare a termenilor de specialitate, cu scopul ușurării înțelegerii mesajului. (*Metodele sunt prezentate în capitolele 3 și 4 din teza, iar instrumentele asociate metodelor sunt prezentate în capitolul 5*).
- Proiectarea structurii de date *FuzzyHashMap* (extensie a structurii *HashMap* din limbajul Java) ce permite căutări aproximative de mare viteză. O confirmare în plus a utilității acestei structuri este dată de folosirea ei de către alți cercetători în proiecte de bioinformatică sau information retrieval. (*prezentată în capitolul 3, secțiunea 3.3*)
- Definirea unui model de accesibilitate universală a textului și ancorarea tuturor metodelor și serviciilor dezvoltate în acest model. (*în capitolul 7, secțiunea 7.1*)

Alte contribuții:

- Studii de limbaj pentru aflarea incidenței terminologiei derivate. (*în capitolul 3, secțiunea 3.2*)
- Studiu asupra eficienței folosirii platformelor tip *crowdsourcing* pentru validarea recunoașterilor aproximative. (*în capitolul 4, secțiunea 4.5*)
- Proiectarea și dezvoltarea serviciilor:
 - *text4all Terminology Interpreter* pentru adaptarea textelor de specialitate. (*în capitolul 5, secțiunea 5.2*)
 - *text4all ITS Tagger* pentru adnotarea terminologiei în formatul ITS 2.0 cu scopul îmbunătățirii traducerii. (*în capitolul 5, secțiunea 5.4*)
 - *text4all DysWebxia* pentru adaptarea textului atât la nivel de prezentare cât și la nivel de limbaj (sinonime pentru cuvinte dificile) pentru persoane cu dislexie. (*în capitolul 6, secțiunea 6.3*)
 - *text4all Customizer* pentru particularizarea aspectului textului din pagini web existente. (*în capitolul 6, secțiunea 6.2*)
 - *text4all analyzer* pentru analiza limbajelor de specialitate și generarea de statistici referitoare la incidența terminologiei de specialitate în diferite forme. (*în capitolul 5, secțiunea 5.5*)