

# Abordări în dezvoltarea sistemelor ”întrebare – răspuns”

Liviu Carcea

Technical University of Moldova

carcea@mail.utm.md

**Rezumat** – acest articol este o încercare de analiză a abordărilor ale sistemelor ”întrebare-răspuns”(SIR). În ultimii ani se observă un interes deosebit pentru sistemele ”întrebare-răspuns”, care a condus la o schimbare calitativă a funcţionării acestor sisteme. Direcţiile principale de cercetare ţin de formularea întrebărilor în limbaj natural, utilizarea motoarelor de căutare, evaluarea calităţii răspunsurilor şi au provocat apariţia abordărilor deosebite de cele tradiţionale.

**Cuvinte-cheie:** abordare, sisteme ”întrebare-răspuns”, regăsirea informaţiei, interogare, motor de căutare.

## I. INTRODUCERE

Sistemele de tip ”întrebare-răspuns” sunt complexe şi utilizează în general mai multe moduri de analiză şi tipuri de resurse. Crearea suportului informatic pentru acest gen de sisteme necesită implicare a specialiştilor din domeniile atât informaticii cât şi a lingvisticii.

Obiectivele ştiinţifice ale cercetărilor din acest domeniu sunt atât cercetarea cât şi analiza abordărilor problemei regăsirii informaţiei relevante la interogarea expusă şi crearea unui sistem automat de informare a populaţiei în limbaj natural, inclusiv şi în limba română. Sistemul de informare va prelucra întrebările formulate de către persoane conform gramaticii limbii române (propoziţii interogative) şi va extrage răspunsul din textele relevante.

Evaluarea calităţii răspunsurilor şi respectiv asigurarea găsirii lor sunt însă probleme de cercetare, pentru care abordările tradiţionale au devenit insuficiente [ 7 ]. Astfel, implementările motoarelor de căutare moderne recurg din ce în ce mai des la tehnicile PLN (Procesarea Limbajului Natural), acestea fiind utilizate în toate etapele fluxului de prelucrare, începând de la nivelul specificării întrebării şi până la extragerea fragmentului de text relevant. Concomitent au apărut şi campaniile de evaluare în domeniul regăsirii inteligente a a informaţiei. Acestea constituie astăzi priorităţi ale cercetării de avangardă dedicată spaţiului digital al cunoaşterii. Ele au fost organizate iniţial în SUA (*MUC-Message Understanding Conference*, *TREC-Text Retrieval Conference*, *DUC-Document Understanding Conference*). În Europa, manifestarea similară este *CLEF(Cross Language Evaluation Forum)*, care are ca subiect în primul rând limbile Uniunii Europene, inclusiv şi cele cu resurse electronice limitate. În acest context proiectul naţional din România SIR-RESDEC lansat în 2007 a răspuns unei priorităţi europene, propunându-şi realizarea unui sistem de ”întrebare-răspuns” în limbaj natural.

Conform planului de acţiuni „eEurope – O societate informaţională pentru toţi”, în Uniunea Europeană a fost propusă lista de 20 de servicii publice în formă electronică prestate cetăţenilor. În Moldova la fel se acordă o atenţie sporită dezvoltării serviciilor electronice. Sistemele automatizate de informare prezintă o parte-componentă a setului de servicii electronice

## II. METODE FOLOSITE ÎN DEZVOLTAREA SISTEMELOR ”ÎNTREBARE-RĂSPUNS”

Cele mai eficiente metode de descoperire şi de achiziţie a informaţiei o reprezintă, în prezent, **motoarele de căutare**. Scopul acestora este de a servi utilizatorului un set de articole în care acesta să poată găsi informaţia care îi este necesară. De multe ori, articolele servite de motoarele de căutare nu îndeplinesc dezideratul de a răspunde cerinţei utilizatorului, de a conţine un răspuns care să poată fi considerat satisfactor. Deasemenea, ele nu oferă răspunsul concret la problema utilizatorului, ci doar un set de articole din care utilizatorul trebuie să extragă singur informaţia căutată.

Pasul următor în domeniul achiziţiei informaţiei e constituit de dezvoltarea sistemelor capabile să răspundă la întrebări formulate de utilizator în limbaj natural. Dezideratul principal al unui astfel de sistem este să asigure un răspuns la întrebarea utilizatorului care să îndeplinească următoarele trei condiţii:

- să fie corect,
- să fie formulat tot în limbaj natural uman şi
- să fie suficient de succint.

Un sistem de răspuns la întrebări necesită o procesare mai complexă a limbajului natural decât sistemele de achiziţie de documente (cum sunt unele motoare de căutare a informaţiei).

În teorie, procesarea limbajului natural e un subiect foarte atractiv, datorită aplicabilităţii sale în domenii ca cel al interacţiunii om-maşină. În practică se constată, însă, o serie de limitări majore, datorate, mai ales, modului diferit în care o afirmaţie în limbaj natural poate fi interpretată şi a multitudinii de sensuri pe care cuvintele constituente le pot lua.

Pentru construirea unui sistem de răspuns la întrebări avem două variante:

- Abordare de tip **shallow**, bazată pe cuvinte cheie.

În această metodă se folosesc cuvinte cheie pentru a găsi pasaje şi propoziţii în text care ar putea candida drept răspunsuri valide la întrebări. Aceste potenţiale răspunsuri urmează să fie analizate mai apoi în profunzime pentru a se stabili dacă sunt răspunsuri reale sau nu. Această metodă poate fi folosită cu succes în cazul întrebărilor scurte, factuale, când se caută nume, date, locaţii, cantităţi.

- Abordarea de tip **deep**, ce implică o analiză mai sofisticată, o procesare sintactică, semantică și contextuală. Există o serie de metode ce pot fi încadrate în această categorie: *abduction, named-entity recognition, relation detection, etc.*

Alegerea unuia dintre cele două modele depinde de complexitatea întrebărilor ce vor fi formulate și de gradul de performanță dorit de la sistem. Este clar că sistemele din cea de-a doua categorie sunt superioare primelor.

### III. MODELE DE ABORDĂRI

Una din cele mai importante etape în sistemele întrebare/răspuns este analiza întrebărilor. De fapt întrebările sunt exprimate în limbaj natural spre deosebire de căutările clasice în care interogările sunt prezentate prin cuvinte cheie [3]. Rezultă că una din etapele de bază constă în extragerea cuvintelor cheie din întrebare pentru a formula una sau mai multe interogări. Altă funcție a analizei întrebării este determinarea tipului răspunsului așteptat. De obicei tipul este formalizat avînd forma unei entități. De exemplu: *PERSOANA, NUMARUL, LOCUL* etc.

Astfel întrebarea “*Cîte facultăți are UTM ?*” va conduce la cuvintele cheie “*facultăți/UTM*” și se va aștepta un răspuns de tip “*numarul (de facultăți)*”.

În general analiza întrebării se face pe trei nivele. Pentru început sunt detectate entitățile întrebării. Apoi se va determina tipul întrebării și deci și tipul răspunsului. În final se aplică o analiză sintactică pentru a determina legăturile dintre elementele întrebării.

În [6] este prezentat un sistem care răspunde la întrebări formulate în limbaj natural în baza unui corpus de date cu posibilitățile următoare:

- Dezambiguizarea întrebărilor.
- Se va ține cont de contextul general al întrebării.
- Extragera răspunsurilor.
- Asigurarea interacțiunii cu utilizatorul ce va ajuta la justificarea răspunsului.
- Reformularea răspunsului în caz de fuziune a documentelor.
- Răspunsul va fi obținut în timp real.

Pentru a obține aceste obiective s-a făcut apel la tehnicile următoare:

- Căutarea informației și procesarea limbajului natural
- Raționament în bază de caz (Case-Based Reasoning)
- Tehnicile inteligenței artificiale (multiagent)

Menționăm faptul că sistemul nu ține cont de profilul utilizatorului și care reduce posibilitățile sale destul de avansate.

Aplicația PIQUANT prezentată la TREC12 [2] a păstrat toate performanțele versiunilor precedente la care a adăugat noi posibilități. Prelucrarea începe cu analiza întrebării, care la rîndul său implică mai multe acțiuni cum ar fi analiza sintactică, identificarea entităților și altele.

Analiza întrebării produce un QFRAME, care conține tipul răspunsului dorit, tipul întrebării, cuvintele cheie din întrebare și o formă sintactică simplă. La etapa următoare sistemul transmite QFRAME generatorului QPLAN, care de fapt generează lista agenților disponibili. Această listă este transmisă modulului “*Rezoluție și*

*răspuns*”, care alege răspunsul adecuat în funcție de tipul întrebării.

PIQUANT utilizează o tehnică numită *QA-par Dossier* pentru a justifica răspunsul. Totuși sistemul este limitat datorită lipsei de interacțiune cu utilizatorul ceea ce nu-i permite să ia o decizie în caz de ambiguitate.

Sistemul descris în [1] are ca principiu de funcționare utilizarea multiplelor surse de date, de fapt doi agenți independenți de căutare: unul pentru Internet și altul pentru corpusul TREC. Acest lucru îi permite sistemului să obțină scoruri înalte de certitudine pentru răspunsuri. Sistemul conține modulele următoare: analiza întrebării, selecția documentului, identificarea entităților și extragerea răspunsului. Dat fiind numărul foarte mare de documente pe Web și redundanța lor ridicată sistemul are șanse mari de a găsi răspunsul exact. Totuși lipsa interacțiunii cu utilizatorul și aplicarea unui algoritm unic pentru a răspunde la întrebări de tip diferit prezintă un obstacol pentru rezolvarea problemei evaluării răspunsului. În [4] se ține cont de rolul cuvîntului în întrebare. Autorii menționează următoarele trei roluri:

- **Focus** – este entitatea asupra căreia punctează întrebarea și pentru care o caracteristică sau o definiție este căutată. Prin definiție acest element trebuie să fie prezent și în textul răspunsului. În întrebarea “*Ce sport practică Zinedine Zidane?*” focusul este “*Zinedine Zidane*”;
- **Domeniul** – care nu este întotdeauna prezent în întrebare, dar de obicei precizează tipul răspunsului. În întrebarea precedentă domeniul este “*sport*”;
- **Verbul principal** – este vorba de verbul prezent în întrebare și care are un rol important în răspuns atunci când introduce un fapt sau o acțiune. În întrebarea precedentă verbul principal este “*practică*”.

Prelucrarea întrebărilor în [7] conține următorii pași. După primirea întrebării este apelat serviciul web care se ocupă de clasificarea întrebărilor pentru a obține tipul de răspuns care trebuie căutat. La următorul pas folosind informația adnotată după procesare, întrebarea este transformată în interogări într-un limbaj formal înțeles de motorul de căutare. Autorii folosesc doi algoritmi pentru a genera două interogări diferite, ambele conținând ca termen de căutare clasa întrebării. Se menționează faptul că în faza de indexare au fost indexate odată cu paragrafele și clasele corespunzătoare acestora, pentru ca reducerea spațiului de căutare să se facă direct din faza de regăsire documentară.

Etapa de regăsire documentară este similară cu cea descrisă în [5], adică cuvintele documentelor au fost filtrate în funcție de descrierea lor morfo-sintactică și ulterior normalizate la forma lor lemă. Autorii au folosit 2 algoritmi de generare a interogărilor formale către motorul de căutare pentru ca sistemul să poată întoarce pentru unele întrebări mesajul că niciun răspuns nu a fost găsit.

Analiza abordărilor de dezvoltare a SIR demonstrează utilizarea complexă și profundă a tehnologiilor moderne de prelucrare a textelor în limbaj natural.

#### IV. CONCLUZII

În societatea modernă, care este o societate informațională, informației îi revine un rol important. Obținerea operativă a informației a devenit o necesitate zilnică pentru mulți cetățeni ai Moldovei.

În cadrul proiectelor instituționale de cercetare științifică și în scopul creării unei societăți informaționale în Republica Moldova este lansat proiectul "Cercetarea în domeniul de Regăsire a Informației în scopul creării sistemului electronic de informare publică", realizat de către cercetătorii Laboratorului de Inginerie a Limbajului Uman (LILU) al catedrei Informatica Aplicată facultatea Calculatoare, Informatica și Microelectronica a Universității Tehnice a Moldovei.

Obiectivele științifice ale proiectului propus sunt de a cerceta problema regăsirii informației relevante la interogarea expusă și de a crea un sistem automat de informare a populației în limba română. Sistemul de informare va prelucra întrebările formulate de către persoane conform gramaticii limbii române (propoziții interogative) și va extrage răspunsul din textele relevante.

Crearea sistemului necesită o implicare a specialiștilor în lingvistică și informatică, care activează în echipa noastră. Pentru sistemul dat vor fi utilizate metodele lingvistice de procesare a propozițiilor, dar și metode statistice de regăsire a informației în volume mari de texte.

Deși există deja sisteme/medii care oferă posibilitatea de a pune întrebări în limbaj natural și de a obține răspunsuri extrase dintr-un volum de texte, nu există deocamdată un mediu complet care să satisfacă toate cerințele în aplicațiile practice. Mai mult, nu cunoaștem nici o aplicație care să folosească limba română. Astfel, sistemul propus prezintă un grad de noutate și complexitate semnificativ și are ca scop dezvoltarea de aplicații reale, și anume, un sistem de informare publică. Potențialii beneficiari sunt practic toți cetățenii Republicii Moldova din țară sau din afara granițelor, dar și cetățeni străini care vor avea nevoie de informație din țara noastră.

În baza abordărilor de elaborare și dezvoltare a SIR, inclusiv și cele netradiționale, dar și a cerințelor formulate de către comandatarii proiectului un astfel de sistem trebuie să conțină:

- Unitatea de analiză a întrebărilor în limba română și de formare a interogării pentru căutarea răspunsului;
- Unitatea de căutare a textelor relevante pentru întrebarea dată;
- Unitatea de detectare a fragmentelor textelor ce conțin răspunsul posibil;
- Unitatea de formare și vizualizare a răspunsului;
- Unitatea de analiză a întrebărilor, răspunsurilor și de evaluare a lor de către utilizatori.

Sistemul necesită implicarea specialiștilor din diferite domenii. Unitatea de analiză a întrebărilor și de formare a răspunsului, sunt realizate aplicând metode de analiză lingvistică, pe când unitatea de căutare a textelor relevante și unitatea de detectare a fragmentelor textelor ce conțin răspunsul se realizează prin metode statistice.

Unitatea regăsirii informației relevante la interogarea expusă este cea principală în sistemul dat și afectează imens performanța lui. Astfel, în proiectul propus, se va:

1) cerceta metodele de regăsire a informației (Information Retrieval - IR),

2) elabora un sistem de informare care va căuta răspunsuri în colecția de documente, utilizând tehnologiile avansate de regăsire a informației,

3) examina impactul mai multor factori care afectează performanța sistemului prin exploatarea lui experimentală on-line.

Sistemele automate de informare a populației pot fi utilizate pe larg de orice instituție guvernamentală și non-guvernamentală, în departamentele de relații cu publicul în scopul maximizării volumului informației obținute de populație în paralel cu minimizarea cheltuielilor de timp. În acest mod, persoanele interesate se vor putea documenta în domeniul solicitat, vor putea primi răspuns la întrebările apărute în procesul întocmirii documentelor necesare pentru, de exemplu, înregistrarea întreprinderilor mici, privatizarea locuințelor, tranzacțiilor comerciale legate de imobil etc.

#### BIBLIOGRAFIA

- [1] Chalendar, G., Dalmas, T. LIMSI-CNRS, The Question. Answering System QALC at LIMSI Experiments in using Web and Wordnet, In Proceedings of the 11 th Text REtrieval Conference (TREC-10), 2002.  
[http://trec.nist.gov/pubs/trec11/t11\\_proceedings.html](http://trec.nist.gov/pubs/trec11/t11_proceedings.html)
- [2] Chu-Carrol, J., Prager, J., Welty, C., Czuba, K., Ferrucci, D. IBM T.J. Watson Research Center IBM's PIQUANT in TREC2003 (TREC-12), 2003  
[http://trec.nist.gov/pubs/trec12/t12\\_proceedings.html](http://trec.nist.gov/pubs/trec12/t12_proceedings.html)
- [3] Crestan, E., Lemaire, E., Claude de Loupy, Ressources pour un sytème de Question/Reponse. Actes de TALN 2004.  
<http://edutice.archives-ouvertes.fr>
- [4] Grappy, A., Ligozat, A., Grau, B. Evaluation de la reponse d'un systeme de question -reponse et de sa justification. COnference en Recherche d'Information et Application, CORIA 2008, Traitement automatique de langues et résumé automatique p. 273-288.
- [5] Ion, R., Ștefănescu, D., Ceaușu, A. and Tufiș D. (2009) RACAIș QA System at the Romanian-Romanian QA@CLEF2008 Main Task, Lecture Notes in Computer Science, vol. 5706, september, p. 393-400.
- [6] Merdaoui, B., Frasson, C. QUERI: Un systeme de question-reponse collaboratif et interactif.  
<http://edutice.archives-ouvertes.fr>
- [7] Ștefănescu, D., Ion, R., Ceaușu, A., Tufiș, D. Sistem întrebare-răspuns antrenabil pentri limba română. În lucrările conferinței "Resurse lingvistice și instrumente pentru prelucrarea limbii române" București, 6-7 mai 2010 p.153-164.