

# Sistemul Întrebare-Răspuns Online

Constantin BREAHA  
Universitatea Technica a Moldovei  
constantin.breahna@yahoo.com

Victoria BOBICEV  
vika@rol.md

**Abstract** — lucrarea dată prezintă un sistem întrebare-răspuns online. Acest sistem răspunde utilizatorilor la întrebări din domeniul medicinei, ce țin de boli, simptome și metode de tratament. Întrebările sunt formulate de către utilizatori în limba română, prin interfața web. Apoi, acestea sunt analizate de către sistem și răspunsul este prezentat utilizatorului. Problema regăsirii automate și formulării răspunsului corect este una din cele mai dificile în domeniul de cercetare numit procesarea limbajului natural. Deseori, analiza sintactică a întrebării nu este suficientă, deoarece este necesară înțelegerea mai detaliată a cerințelor utilizatorului. Din cauza dificultăților de acest gen, sistemele existente răspund doar la un număr limitat de tipuri de întrebări.

**Index Terms** — procesarea limbajului natural, lingvistica computațională, regăsirea informației, sisteme întrebare-răspuns, extragerea informației.

## I. INTRODUCERE

În prezent oamenii au acces la o cantitate extrem de mare de informație și necesită instrumente care ar permite selecția datelor potrivite din acest volum uriaș de cunoștințe. Lingvistica computațională încearcă să ofere sisteme informatice menite să ajute utilizatorul să se descurce în multitudinea documentelor disponibile.

Întrebare – Răspuns (ÎR) este un domeniu de cercetare care combină diferite domenii, și anume: Căutarea Informației, Extragerea Informației și Procesarea Limbajului Natural [2].

## II. STRUCTURA GENERALĂ A SISTEMELOR ÎR

Un sistem ÎR conține, de regulă, 3 module: modulul de procesare a întrebării, modulul de căutare a informației și modulul de extragere a răspunsului [1].

Modulul de procesare a întrebării o analizează sintactic și o clasifică în baza unei taxonomii deja existente. Tipul întrebării la rândul său indică tipul răspunsului așteptat. Prin urmare, modulul dat transformă întrebarea în interogare pentru un motor de căutare. Tipurile de întrebări posibile sunt exemplificate în Tabelul 1.

Clasificarea întrebării și cunoașterea tipului ei nu sunt suficiente pentru găsirea răspunsului la toate întrebările. Întrebarea "ce?" este ambiguă din punctul de vedere al informației prezentate în întrebare. Pentru a aborda această ambiguitate, este necesară o componentă care identifică centrul. Centrul întrebării este un cuvânt sau set de cuvinte care indică ce informație este cerută în întrebarea dată. De exemplu, întrebarea "Care este cel mai lung râu din SUA?" are ca centru "cel mai lung râu".

Odată ce au fost identificate "centrul" și "tipul întrebării", modulul formează o listă de cuvinte-cheie care vor fi trimise componentului de căutare a informației.

Pot fi utilizate și metode de extindere a setului de cuvinte-cheie a întrebării, de exemplu, adresarea la sursele lexicale online, cum ar fi ontologia WordNet. Seturile de sinonime din WordNet pot fi folosite pentru a extinde seturile de cuvinte-cheie prin cuvintele legate semantic. Acestea pot apărea în documente conținând răspunsul corect la întrebare.

TABELUL 1. EXEMPLE TIPURI DE ÎNTREBĂRI

Model	Întrebare	Tipul răspunsului
cine X?	Cine este președintele României?	Persoană
unde X?	Unde se afla Mariana Stanciu pe 17 aprilie?	Loc
cînd X?	Cînd a fost adoptată constituția?	Data
cum X?	Cum au fost construite piramidele din Egipt?	Manieră
dece X?	De ce ninge?	Motiv

Modulul de căutare a informației primește la intrare cuvintele cheie formate de modulul precedent și returnează un set de documente relevante.

Sistemele de căutare a informației sunt, de regulă, evaluate în baza următoarelor două valori – precizie (precision) și retragere (recall). Precizia se referă la raportul dintre documentele relevante returnate din numărul total și numărul total de documente obținute. Retragera este numărul de documente relevante returnate din numărul total de documente disponibile în colecția de documente căutate. În general, scopul sistemelor de extragere a informației este de a optimiza atât precizia cât și retragera.

$$\text{Precizia} = \frac{\text{numărul de răspunsuri corecte}}{\text{numărul de întrebări răspunse}}$$

$$\text{Retragera} = \frac{\text{numărul de răspunsuri corecte}}{\text{numărul de întrebări care urmează a fi răspunse}}$$

Pentru sistemele întrebare-răspuns, optimizarea diferă considerabil. Un sistem ÎR efectuează post-procesarea documentelor returnate, importanța retragerii sistemului este mult mai semnificativă decât precizia lui. În alte cuvinte, modulul dat trebuie să returneze cât mai multe documente în scopul de a nu pierde vre-un document cu răspunsul corect la această etapă.

Prin urmare, numărul de documente returnate de sistemul de căutare a informației poate fi foarte mare. Documentele returnate sunt filtrate și ordonate de următorul modul – cel de extragere a răspunsului. Obiectivul principal al modulului dat este de a crea un set de paragrafe candidate, care conțin răspunsul sau răspunsurile posibile. Filtrarea paragrafelor este folosită

pentru a reduce numărul de documente candidate, și de a reduce cantitatea textului din fiecare document. Paragrafele sunt ordonate după următorul principiu: cu cât mai aproape se află în cadrul lui cuvintele cheie, cu atât el este plasat mai aproape de începutul listei. Dacă cuvintele cheie dintr-o întrebare sunt dispersate prin tot volumul textului din paragraf, acesta va fi plasat mai departe în lista ordonată [3].

În fază finală, modulul de procesare a răspunsului realizează identificarea, extragerea și validarea răspunsului din setul de paragrafe ordonate transmise din modulul de procesare a documentelor.

Există diverse reguli de extragere a răspunsului corect din răspunsurile candidate. Extragerea poate fi bazată pe măsurarea distanței dintre cuvintele cheie, numărul de cuvinte-cheie potrivite sau alte reguli similare.

Astfel, sistemul ÎR efectuează următoarele operațiuni:

1. Utilizatorul postează întrebarea în sistemul ÎR.
2. Analizatorul de întrebare determină centrul întrebării pentru a spori eficacitatea sistemului ÎR.
3. Clasificarea întrebării joacă un rol important în sistemul ÎR prin identificarea tipului întrebării și tipului răspunsului așteptat.
4. În cadrul reformulării întrebării, aceasta este parafrazată și extinsă.
5. Componenta căutării informației este folosită pentru extragerea documentelor relevante cuvintelor-cheie din întrebare.
6. Documentele obținute sunt filtrate și împărțite în paragrafe ce pot conține răspunsul.
7. Paragrafele filtrate sunt ordonate și trecute prin modulul de procesare a răspunsului.
8. Bazându-se pe tipul întrebării și alte tehnici de recunoaștere, răspunsurile candidate sunt identificate.
9. Un set de reguli definite este utilizat pentru a extrage doar frazele relevante care răspund la întrebare.
10. Răspunsul extras este validat și prezentat utilizatorului.

### III. SISTEME DE ÎNTREBARE – RĂSPUNS CUNOSCUTE

Sisteme BASEBALL [7] și LUNAR [6] au fost printre primele de acest fel. BASEBALL putea răspunde la întrebările legate de jocul baseball din SUA; sistemul LUNAR răspundea la întrebările cu privire la analiza geologică a rocilor returnate de misiunile lunare Apollo. Ambele sisteme au fost destul de eficiente în domeniul lor. LUNAR a fost demonstrat la conferința științifică lunară în 1971 și a fost capabil să răspundă la 90% din întrebările puse. În următorii ani au fost dezvoltate mai multe sisteme ÎR circumscrise. Trăsătura comună a acestor sisteme era folosirea bazelor de date sau sistemelor de cunoștințe scrise de mână de către experți în domeniul ales. Abilitățile lingvistice ale sistemelor BASEBALL și LUNAR erau asemănătoare celor folosite de ELIZA, primul program – robot pentru chat [8].

START este primul sistem întrebare-răspuns bazat pe internet. Acest sistem a apărut, pentru prima dată, în 1993 și acum poate fi accesat pe adresa <http://start.csail.mit.edu>. Sistemul Start a fost dezvoltat în Laboratorul de inteligență artificială MIT de către Boris Katz. În căutarea răspunsului

la întrebare, sistemul folosește atât baza de date locală, cât și un șir de resurse din internet. Spre deosebire de motoarele de căutare, sistemul START are scopul de a oferi utilizatorului doar informația corectă, în loc de a furniza o simplă listă cu rezultate găsite. În prezent, sistemul poate răspunde la milioane de întrebări în limba engleză bazate pe localități (orașe, țări, coordonate, etc.), filme (titluri, actori, regizori, etc.), definiții și multe altele din următoarele categorii:

- Întrebări explicative (What is fractal?)
- Întrebări legate de fapte (Who invented telegraph?)
- Întrebări legate de relații (What country is bigger, Russia or USA?)
- Întrebări cu răspunsuri multiple (Show me some poems by Alexander Puskin?)

Sistemul START folosește "adnotarea limbajului natural" pentru a conecta căutătorii de informație la sursele de date. Această tehnică utilizează propoziții și fraze în limbajul natural – adnotări – ca descrieri ale conținutului, care sunt asociate cu segmente de informație. Un segment de informație este preluat atunci când adnotarea lui corespunde întrebării. Această metodă permite sistemului START să manipuleze toate tipurile de date, cum ar fi: texte, diagrame, imagini, clipuri video sau audio, seturi de date, pagini web și altele.

Componenta de procesare a limbajului natural este formată din două module care folosesc aceeași gramatică. Modulul de înțelegere analizează textul și codifică informația găsită în baza cunoștințelor existente. Având segmentul potrivit al bazei de cunoștințe, modulul de generare creează propoziții.

Sistemul RACAI a fost dezvoltat de către Institutul de Cercetări pentru Inteligența Artificială a Academiei Române. Pentru prima dată sistemul RACAI a fost prezentat în competiția de întrebare-răspuns CLEF în anul 2006.

Sistemul RACAI [4] poate răspunde la întrebările formulate în limba română căutând răspunsul în documentele Wikipedia. Sarcina sistemului este de a oferi șirul cel mai scurt și bine format sintactic, care va răspunde la întrebarea utilizatorului. Caracteristicile specifice ale sistemului sunt:

1. Pentru varianta cross-linguală, traducerea interogării este realizată pe baza ontologiei lexicale RO-WordNet aliniată la nivel de serie sinonimică cu ontologia pentru limba engleză WordNet.
2. Pentru indexare și regăsire este folosită platforma open-source LUCENE. Precizia ei este superioară unui motor general de căutare ce se reflectă substanțial în performanța sistemului ÎR.
3. Clasificarea întrebărilor se realizează de un sistem bazat pe modelul maximizării entropiei, antrenabil în raport cu o tipologie dată. Sistemul de clasificare este extrem de precis (la CLEF2007, din cele 200 de întrebări, 199 au fost corect clasificate în una din cele 8 categorii definite de organizatori – temporal, time interval, definition, measure, list, location, names și explanations).
4. Pentru prelucrările lingvistice (segmentarea lexicală, adnotarea morfo-sintactică, lematizarea, analiza sintactică parțială) este folosită platforma de servicii web

dezvoltată de către Institutul de Cercetări pentru Inteligența Artificială a Academiei Române.

5. Extragerea celui mai probabil răspuns la o întrebare se realizează cu ajutorul unui algoritm recursiv de potrivire structurală între arborele de dependențe al întrebării și arborele de dependență al candidatului de răspuns (metoda este originală și în mare măsură responsabilă pentru performanța foarte bună a sistemului).

#### IV. DESCRIEREA SISTEMULUI CREAT

Sistemul creat în cadrul lucrării date răspunde la întrebările din domeniul medicinei. Acestea țin de boli, simptome și metode de tratament.

Sistemul este constituit din modulele clasice descrise anterior, și anume: modulul de analiză a întrebării, modulul de căutare a informației și modulul de extragere a răspunsului.

Întrebarea este inițial procesată, iar apoi analizată din două aspecte: determinarea tipului întrebării și a centrului. Cele două atribute sunt folosite pentru a extrage răspunsul. Centrul întrebării va fi folosit pentru formarea interogării pentru motorul de căutare. Din documentele obținute de către motorul de căutare vor fi extrase paragrafe mici de text ca răspunsuri conform tipului întrebării, și prezentate utilizatorului.

Primul pas este analiza morfologică. Sistemul creat folosește algoritmul elaborat de către Universitatea Alexandru Ioan Cuza din Iași (UAIC) - UAIC Romanian Part of Speech Tagger. Acest Part Of Speech tagger românesc combină un model statistic cu unul bazat pe reguli. Dicționarul morfologic a fost extras în mare parte din DexOnline și conține 1,25 milioane de cuvinte distincte. POS tagger-ul are o precizie de 96.6% testat pe o variantă corectată a corpusului "1984". Un exemplu de întrebare analizată morfologic este prezentat în figura 1.

	<b>Cum</b>	<b>se</b>	<b>tratează</b>	<b>afazia</b>	<b>?</b>
	↓	↓	↓	↓	
<i>leme</i>	cum	sine	trata	afazia	
<i>părți de vorbire</i>	adverb	pron. reflexiv	verb predicativ	nume comun	quest

FIGURA 1. ANALIZA MORFOLOGICĂ A CUVINTELOR ÎN ÎNTREBARE.

Următorul pas este clasificarea întrebării, sau determinarea tipului întrebării. Acest pas joacă un rol important, fiindcă depistarea corectă a tipului întrebării va determina tipul răspunsului oferit de program.

Pentru rezolvarea acestei sarcini este folosit un tabel de cuvinte-cheie, ce pot descrie clasa sau tipul întrebării. Cuvintele-cheie din întrebare se compară cu cuvintele-cheie a claselor din tabel și astfel se determină tipul întrebării.

Modul funcționării acestei metode este demonstrat în baza întrebării: *Cum se tratează afazia?* Se analizează fiecare cuvânt din propoziția dată:

- Cuvântul "cum" este găsit în lista cuvintelor cheie a clasei *manner*, respectiv, clasei *manner* i se

mărește scorul.

- Cuvântul "se" nu este găsit printre cuvintele cheie și este sărit.
- Cuvântul "tratează" este găsit în cuvintele cheie a clasei *manner*, respectiv, clasei *manner* i se mărește scorul.

În final clasa *manner* obține scorul cel mai mare și este considerată clasa întrebării date. Tabelul 2 conține câteva exemple de clase și cuvinte-cheie asociate claselor date.

TABEL 2. CLASELE ÎNTREBĂRILOR ȘI CUVINTELE CHEIE ASOCIATE CU CLASELE DATE

Clasa întrebării	Cuvintele cheie
<i>Definiton</i>	ce, definiția, explică
<i>Cause</i>	care, cauzele, cauza, pricina
<i>Manner</i>	cum, tratează, tratament, tratarea

Următoarea etapă în căutarea răspunsului este definirea centrului întrebării. Acesta este definit de setul de cuvinte extrase din întrebare. Pentru determinarea centrului întrebării și formarea interogării sunt parcurși următorii pași:

- excluderea semnelor de punctuație;
- excluderea cuvintelor interogative;
- excluderea cuvintelor funcționale;
- termenii rămași formează centrul;
- crearea interogării.

După determinarea centrului întrebării poate fi formulată interogarea.

Inițial sistemul caută răspunsul în baza de date a răspunsurilor stocate. Aceasta conține următoarele câmpuri: titlul documentului; descrierea documentului; cuvintele-cheie ale documentului; sursa documentului; data adăugării/modificării; tipul întrebării; cuvintele-cheie care definesc tipul întrebării; întrebarea trimisă sistemului ÎR de către utilizator; numărul de adresări cu întrebarea dată. La momentul dat baza de date se completează manual, însă se preconizează crearea unui modul care va completa baza de date automat cu întrebările utilizatorilor și răspunsurile corecte în procesul utilizării sistemului online.

Dacă în baza de date nu a fost găsit un răspuns potrivit, sistemul se adresează online la Wikipedia, o sursă vastă de informații din diverse domenii. Wikipedia oferă un set de API în diferite limbaje de programare și formate de date.

Sistemul este creat cu ajutorul limbajului PHP. Datele din Wikipedia sunt extrase în formatul SOAP care se bazează pe limbajul XML. Documentul este analizat paragraf după paragraf, căutând cuvintele-cheie și indicatorii tipului de întrebare. Paragraful cu cea mai înaltă rată de corespondență este selectat și returnat utilizatorului.

#### V. CONCLUZIE

Lucrarea dată prezintă un sistem de întrebare-răspuns online. Sistemul dat răspunde la întrebările utilizatorilor din domeniul medicinei, și anume despre bolile, simptomele și metodele de tratare. Întrebările se formulează în limba română de către utilizatori prin interfața web. Sistemul analizează întrebările puse, caută răspunsul și oferă utilizatorului.

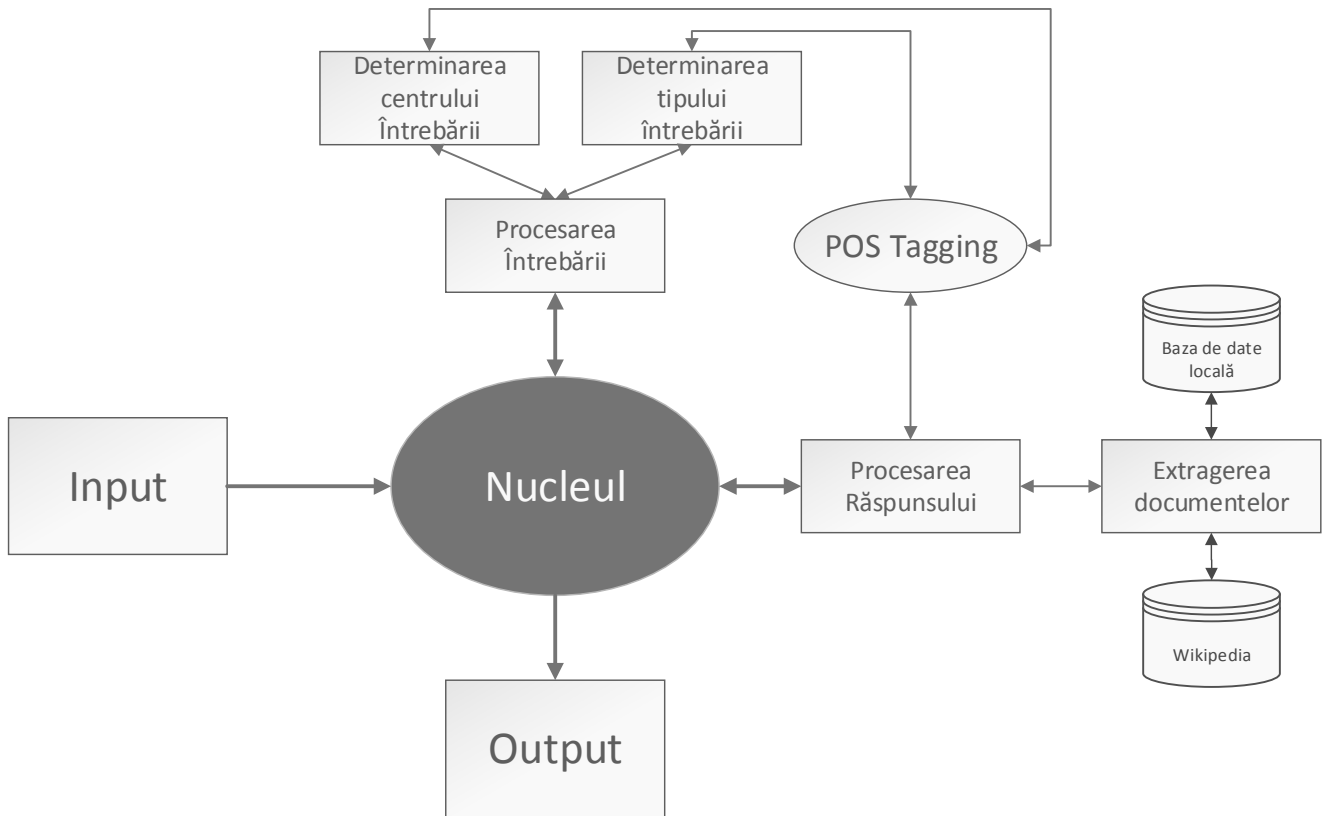


FIGURA 2. SCHEMA GENERALĂ A SISTEMULUI CREAT

Problema regăsirii automate și formulării răspunsului corect este una din cele mai dificile în domeniul de cercetare numit procesarea limbajului natural. În multe cazuri, analiza sintactică a întrebării nu este suficientă. Astfel, este necesară înțelegerea tipului de răspuns solicitat de utilizator. Din cauza dificultăților de acest gen, sistemele existente răspund doar la un număr limitat de tipuri de întrebări

La momentul dat sistemul creat răspunde la întrebările de trei tipuri: definiția, simptomele și tratamentul posibil a bolilor descrise în Wikipedia în limba română.

#### BIBLIOGRAFIA

- [1] Constantin Orăsan, Doina Tatar, Gabriela Șerban, Dana Lupsa, Andrian Oneț. How to build a QA system in your back-garden: application for Romanian
- [2] Dan Tufiș. Sisteme de întrebare – răspuns în limbaj natural pentru spații de căutare deschise - Seminarul internațional "Instrumente pentru asistarea traducerii".
- [3] Radu Ion, Alexandru Ceaușu, Dan Ștefănescu, Dan Tufiș, Elena Irimia, Verginica Barbu Mititelu. Monolingual and Multilingual Question Answering on European Legislation – Research Institute for Artificial Intelligence, Romanian Academy, no.13, Romania.
- [4] Radu Ion, Dan Ștefănescu, Alexandru, Ceaușu, Dan Tufiș. RACAI'S QA System at the Romanian-Romanian QA @ CLEF2008 Main Task - Research Institute for Artificial Intelligence, Romanian Academy, Septembrie, no.13, Romania.
- [5] Boris Katz, Gary Borchardt and Sue Felshin. Natural Language Annotations for Question Answering. Proceedings of the 19th International FLAIRS Conference (FLAIRS 2006), 2006.
- [6] Woods, W. A., Kaplan, R. M., and Nash-Webber, B. The lunar sciences natural language system: final report. Rep. No. 2378, Bolt Beranek and Newman Inc., Cambridge, Mass., 1972.
- [7] Bert F.Green, Jr., Alice K.Wolf, Carol Chomsky, and Kenneth Laughery. BASEBALL: an Automatic Question-Answerer. Western joint IRE-AIIEE-ACM computer conference, 1961.
- [8] Weizenbaum, Joseph (January 1966), "ELIZA - A Computer Program For the Study of Natural Language Communication Between Man And Machine", Communications of the ACM 9 (1), 1966.