

Modelarea, analiza și clasificarea conversațiilor în medii colaborative

Alexandru Bardaş¹, Mihai Dascălu¹, Ștefan Trăușan-Matu^{1,2}

¹Universitatea Politehnică din București, Facultatea de Automatică și Calculatoare
Bd. Splaiul Independenței, Nr. 313, 060042, București, România

²Institutul de Cercetări în Inteligență Artificială
Calea 13 Septembrie, Nr. 13, 050711, București, România

E-mail: alex.bardas@gmail.com, mihai.dascalu@cs.pub.ro, stefan.trausan@cs.pub.ro

Rezumat. Clasificarea conversațiilor în medii colaborative este necesară pentru o mai bună înțelegere a subiectelor discutate. Ontologiile reprezintă o modalitate eficientă și reprezentativă de conceptualizare a unui domeniu. Plecând de la un set predefinit de ontologii specifice mai multor domenii, se dorește clasificarea cât mai precisă a unei conversații purtate între mai mulți participanți. Pentru a facilita acest proces, este nevoie ca limba în care mesajele sunt scrise să fie detectată și intervențiile utilizatorilor să nu conțină greșeli gramaticale. Astfel, termenii utilizați de către participanții la conversații sunt regăsiți în conceptele ce definesc ontologiile, iar conceptele predominante și centrale determină domeniile pe baza cărora conversațiile sunt clasificate. Pentru a realiza această analiză, am construit o platformă web colaborativă ce permite schimbul de mesaje între utilizatori în timp real, vizualizarea termenilor regăsiți în ontologii cu ajutorul unui graf radial și exportul arhivei de discuții într-un format ușor accesibil care să faciliteze clasificarea lor ulterioară. Aplicația este construită pentru a permite accesul de pe orice tip de dispozitiv, afișând informațiile minime necesare pentru un mobil și redimensionându-se automat pentru un ecran mai mare. În articol sunt descrise implementarea aplicației, avantajele folosirii acesteia și ale clasificării automate, precum și oportunitatea integrării ei în diverse contexte educaționale.

Cuvinte cheie: mesagerie instantanee, conversație, colaborare, clasificare, corecție ortografică, ontologie.

1. Introducere

Web-ul social este centrat pe interactivitatea dintre utilizatori, precum și pe activități de publicare, căutare, regăsire și interpretare a informațiilor. Implementarea unei infrastructuri care să permită colaborarea cât mai eficientă și mai rapidă sau care să facă atât posibilă extragerea de date cu privire la profilele și schemele de comunicare ale utilizatorilor, cât și

prelucrarea acestora, a devenit o necesitate în cazul marilor comunități online. Scopul acestor facilități este de a-și ajuta utilizatorii să transmită informații relevante, să cunoască alte persoane cu aceleași interese sau să găsească soluții la problemele pe care le întâmpină.

Deși world wide web-ul a fost proiectat inițial ca un mediu pur informațional, acesta a evoluat spre un mediu al aplicațiilor, marcând tranziția către un web semantic (Berners-Lee et al., 2001). În prezent, aceste aplicații reprezintă sisteme software complexe care oferă servicii interactive și personalizabile.

Companii cu renume investesc sume considerabile în atragerea cu orice preț a unui număr cât mai mare de utilizatori și oferirea posibilității de colaborare între aceștia. În noiembrie 2011 Microsoft a achiziționat Skype pentru 8.5 miliarde de dolari (Bright, 2011), iar în februarie 2014, Facebook a achiziționat aplicația de mesagerie instantanee WhatsApp pentru o tranzacție estimată la peste 19 miliarde dolari (Stone, 2014). Aceste achiziții ilustrează cerința foarte mare care există pentru comunicații în timp real și retenția utilizatorilor. Astfel, relaționând datele deja existente sau înțelegând contextul conversațiilor purtate între utilizatori, precum și graful rețelei de prieteni, domeniile de interes ale acestora pot fi identificate mult mai bine.

Informațiile sunt atât de numeroase încât prelucrarea sau clasificarea lor nu mai poate fi făcută manual, de către oameni. Conform McAfee & Brynjolfsson (2012), în 2012 aproximativ 2.5 exabytes de date erau creați în fiecare zi, iar volumul era estimat că se dublează la fiecare 40 luni. Suplimentar, conform Cisco (2013), în 2017 peste 3.6 miliarde de utilizatori (mai bine de 48% din populația preconizată la acel moment) vor fi conectați la Internet. Vor fi aproximativ 1.4 dispozitive mobile pe cap de locuitor, fapt ce reliefează o dependență din ce în ce mai mare vizavi de tehnologie.

O adopție mai mare va determina și o creștere în ceea ce privește cantitatea de informație disponibilă pe Internet. Este însă interesant de observat ce se va întâmpla cu toate datele noi care apar, deoarece și interesul pentru gruparea și înțelegerea semantică a documentelor a crescut foarte mult. Mai mult, majoritatea motoarelor de căutare indexează în continuare documentele în funcție de conținutul lexical, și doar parțial în funcție de semantica aferentă.

Astfel, trebuie găsită o alternativă ca aceste date să fie prezentate atât într-un format accesibil utilizatorilor, cât și sistemelor de calcul ce urmează să le prelucreze și să le clasifice. Din acest motiv, definirea lor semantică

devine foarte importantă. Aceasta va facilita managementul cunoștințelor pe web și în particular, conectarea și reutilizarea cunoașterii (sindicalizarea conținutului), dar și localizarea cunoștințelor noi și relevante, cu ajutorul sistemelor de recomandare.

Lucrarea de față realizează o trecere în revistă a aplicațiilor colaborative existente, detaliază arhitectura soluției propuse și facilitățile integrate de prelucrare a limbajului natural, precum și aplicabilitatea aplicației dezvoltate în diverse scenarii educaționale.

2. Platforme web colaborative

În general, colaborarea facilitează munca în echipa și oferă posibilitatea creșterii productivității, alocării mai bune a resurselor și dobândirii mai ușoare de informații noi (Stahl, 2006). O aplicație colaborativă poate fi văzută ca un program accesat din Internet sau dintr-o rețea, care oferă o serie de funcționalități menite să ajute la dezvoltarea și îndeplinirea rapidă a sarcinilor de lucru.

Principalele categorii de software colaborativ sunt destinate:

- coordonării: calendare, foi de calcul online, sisteme de management al proiectelor;
- comunicării: wiki, email, blog;
- consfătuirii: mesageria instantanee, forum.

Pentru a argumenta utilizarea unei anumite tehnologii în cadrul experimentelor educaționale derulate în cadrul facultății, s-a considerat necesară analiza minuțioasă a celor mai bune alternative web disponibile. Multe platforme online au rolul de a pune în legătură utilizatorii prin oferirea unor soluții eficiente, cât mai ușor de utilizat, furnizând totodată seturi diferite de funcționalități. Modalitățile de colaborare sunt diferite, de la soluții de tipul întrebi și primești răspuns (Stackoverflow, Quora), la aplicații de management de proiect (Asana, Jira, Trello), până la soluții ce oferă mesagerie instantanee (Campfire, HipChat, Slack - Figura 1 și Appear.in - Figura 2).

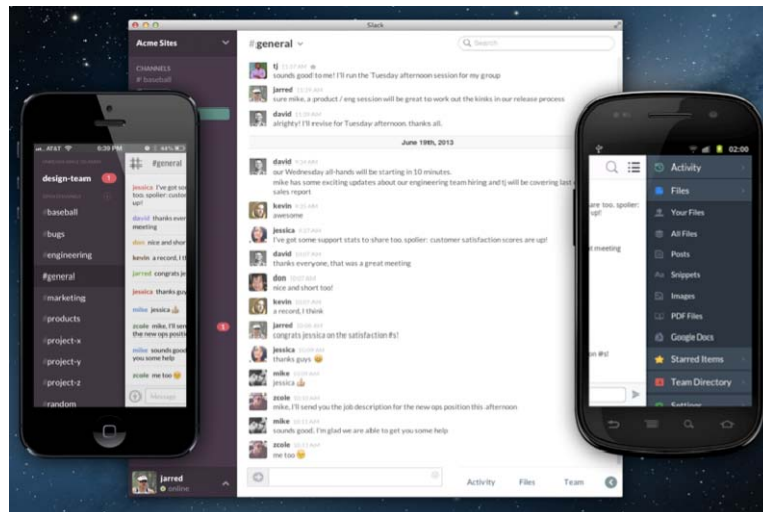


Figura 1. Interfața aplicației Slack reliefând portabilitatea platformei

O soluție populară este reprezentată de chatul prin canale IRC care reprezintă o alternativă matură, testată în timp, simplă și eficientă. Totodată, pe măsura trecerii timpului și o dată cu apariția alternativelor, limitările au apărut și s-au accentuat tot mai clar. Multe canale nu sunt arhivate, discuțiile purtate fiind practic pierdute. De asemenea, căutarea unor topice sau termeni poate fi uneori dificilă, alții imposibilă. De cele mai multe ori, domeniul de discuție al camerei este greu de identificat. Comunicarea prin Internet Relay Chat (IRC – RFC 1459) este o soluție depășită pentru cei care vor să analizeze subiectele discutate sau să colaboreze cât mai eficient.

Slack este o soluție relativ nouă de mesagerie instantanee, cu o adopție foarte bună încă de la început, de la creatorul platformei Flickr. A apărut din nevoia de a oferi ceva nou față de aplicațiile disponibile pe piață (ex: Campfire, Hipchat). Este rapid și foarte ușor de utilizat, oferind o interfață simplistă și intuitivă. Este destinat exclusiv colaborării într-o organizație, pentru un grup ridicat de persoane. Oferă căutare de termeni atât în cadrul discuțiilor, dar și al documentelor încărcate. Căutarea unui topic se efectuează ușor, putând, de exemplu, căuta mai întâi dacă au mai fost semnalate probleme asemănătoare înainte de a le mai adresa încă o dată. Este disponibil atât în versiune web, dar și ca aplicație nativă pe toate sistemele de operare. Succesul care îl are dovedește că multe companii sunt

dispușe să plătească pentru o aplicație destinată atât comunicării, cât și organizării discuțiilor unui grup.

Appear.in este o platformă de ultimă generație care folosește cele mai noi tehnologii și API-uri web (WebRTC) pentru a oferi pe lângă funcționalitatea de chat, și facilități video și audio. O cameră poate găzdui o conferință video cu până la 8 participanți. Un mare avantaj pe care îl oferă este că înregistrarea pe sit este opțională, oricine putând accesa o cameră de discuții doar știind numele acesteia. Soluția este compatibilă doar cu cele mai noi versiuni ale browserelor web și demonstrează faptul că web-ul este un mediu matur și suficient de puternic pentru a oferi funcționalitatea și performanța unui mediu nativ.

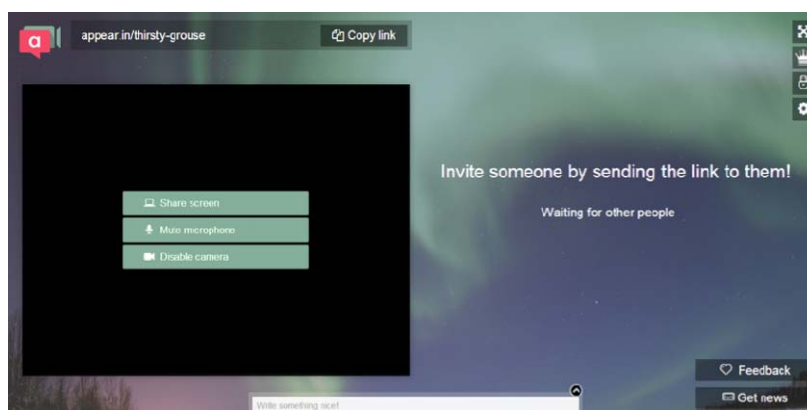


Figura 2. Interfața aplicației Appear.in

3. Facilități integrate de prelucrare a limbajului natural

În vederea clasificării discuțiilor, detecția limbii este un prim pas important. Cu ajutorul acesteia, se pot oferi sugestii adecvate de corecție și se pot căuta diverși termeni în ontologiile specifice anumitor domenii. Astfel, conversațiile purtate într-un mediu colaborativ pot fi clasificate în funcție de domeniul de discuție. Există mai mulți algoritmi folosiți pentru această clasificare, cei mai mulți fiind bazați pe metode statistice. Detalierea acestora se regăsește în secțiunea 3.1.

Totodată, întrucât în cadrul unei conversații pot exista vorbitori non-nativi ai limbii utilizate, am optat pentru integrarea unui sistem de traducere, din limba detectată într-o alta limbă selectată de către utilizator. S-a

constatat că cele mai bune astfel de motoare sunt disponibile ca servicii web (Google Translate, Bing Translator, Babel Fish). Implementarea în aplicație a unui algoritm eficient nu este un aspect trivial, iar rezultatele ar fi fost sub cele ale serviciilor oferite de cele web menționate. În final, am decis integrarea serviciului *Bing Translator*, ce oferă suport pentru 45 limbi. Spre deosebire de alternative, *Bing Translator* este disponibil gratuit și pune la dispoziție o documentație amplă, aspecte ce au favorizat integrarea acestuia.

Ulterior, procesul de corecție ortografică este aplicat înainte de a realiza o clasificare automată a discuțiilor purtate. Ambele facilități sunt descrise în detaliu în sub-secțiunile următoare.

3.1 Detectia limbii bazată pe modelul n-gram

Una dintre cele mai eficiente metode statistice de detectare a limbii presupune crearea unui model n-gram (secvențe de n caractere). Metoda descrisă de Cavnar & Trenkle (1994) se bazează pe faptul că pentru fiecare limbă există un set de cuvinte dominante din prisma frecvenței ridicate de apariție a lor în cadrul textelor. Astfel, textele vor avea distribuții ale frecvenței n-gram asemănătoare dacă sunt din același domeniu și limbă.

Modelul inițial este construit dintr-un text de învățare specific fiecărui limbaj. Pentru un text nou se creează un alt model, iar acesta este comparat cu fiecare model stocat anterior (aferețt fiecării limbi). Limba cea mai probabilă este cea care are modelul cel mai asemănător cu modelul creat pentru noul text. Cu cât textul este mai lung, cu atât detectia este mai precisă. Problemele apar mai ales în cazul textelor scurte sau când sunt folosite mai multe limbi în același text, caz în care algoritmul poate fi ineficient.

Algoritmul funcționează rapid și în cele mai multe cazuri (mesaje corecte, mai lungi de 100 caractere) oferă rezultate relevante. Astfel, acuratețea metodei este foarte bună, detectia având o rată de succes între 96 și 100% pentru diverse texte și limbi.

3.2 Corecție ortografică

Pentru detectarea greșelilor de ortografie și pentru oferirea de sugestii este folosit un algoritm statistic asemănător cu cel descris de Norvig (2014). Pentru a antrena acest algoritm este nevoie de un set destul de mare de texte și de dicționarul specific unei limbi. Detectia limbii unui mesaj efectuată în

prealabil este necesară pentru a selecta setul de date corespunzător corecției ortografice. Astfel, pornind de la un corpus de texte, sunt extrase toate cuvintele, transformate în varianta cu minuscule și este antrenat un model de probabilitate prin numărarea frecvenței de apariție a fiecărui cuvânt. Distanța de editare a 2 cuvinte sau distanța Levenshein (1966) este definită ca numărul de editări ce trebuie făcut într-un cuvânt pentru a-l transforma în cel de-al doilea cuvânt. Pentru un cuvânt w , sunt enumerate toate corecțiile posibile c care poate reprezenta un set foarte mare de cuvinte, iar apoi sunt păstrate doar corecțiile valide, care sunt regăsite în dicționarul limbii. Se consideră că 80-90% dintre greșelile de ortografie diferă de cuvântul corect doar printr-o literă.

Astfel, fiecare cuvânt care nu este găsit direct în dicționar va fi supus acestui algoritm pentru o distanță de editare egală cu 1. Dacă nu este găsită nicio corecție, va fi aplicat algoritmul și pentru o distanță de editare egală cu 2, iar corecțiile întoarse sunt oferite ca sugestii de editare. Totodată, specific conversațiilor chat sunt abrevierile (spre ex., "brb", "lol") care sunt înlocuite automat pe baza unor liste predefinite de mapări.

3.3 Clasificarea discuțiilor

Fiecare conversație purtată într-o cameră virtuală de discuții este analizată ulterior pentru a fi clasificată într-un domeniu de interes. Algoritmul de clasificare este bazat pe frecvența aparițiilor cuvintelor cheie din arhiva inițială de discuții pe baza unui set de ontologii predefinite. Cu cât fiecare set este mai mare, cu atât crește probabilitatea regăsirii a cât mai multor cuvinte în ontologii. Analiza se bazează pe faptul că atât mesajele, cât și ontologiile sunt scrise și definite în aceeași limbă.

Algoritmul funcționează astfel: pentru o arhivă A , se cunosc toate mesajele $m[1:N]$, iar pentru un mesaj m_i se știe autorul m_{ia} , limba mesajului m_{il} și data trimiterii lui m_{id} . Se găsește limba predominantă L din lista de mesaje, iar replicile scrise în altă limbă sunt traduse în L (folosind serviciul Bing Translator). Pentru o listă inițială de ontologii $O[1:M]$ definite în L , se creează 2 dicționare: (1) H_s , unde cheile sunt conceptele prezente în ontologii, iar valorile sunt liste cu ontologiile în care apar și nivelul ierarhic pe care apar (2 ontologii diferite pot conține același concept) și (2) H_o , unde cheile sunt domeniile ontologiilor, iar valorile vor fi setate ulterior.

Dacă un cuvânt w din m_i nu este găsit în dicționarul limbii L , atunci se aplică algoritmul de corecție cu o valoare a distanței de editare egală cu 1 și

rezultatele sunt căutate în H_s . Dacă sunt găsite, pentru fiecare valoare v de la cheia din H_s găsită, $H_o[v]$ va fi incrementat cu o valoare $v_2 \cdot i_v$. Dacă w este găsit în L , este căutat în H_s și dacă este găsit, pentru fiecare valoare de la cheia din H_s găsită, $H_o[v]$ va fi incrementat cu o valoare $v_1 \cdot i_v$, unde i_v este determinat în funcție de nivelul ierarhic în care termenul v apare în ontologie. Astfel, dacă se află pe nivelul n , i_v va fi $1/n$. Cu alte cuvinte, cu cât conceptul este mai apropiat de rădăcina ontologiei, cu atât acesta este mai general și relevanța vizavi de un domeniul specific este diminuată (Resnik, 1995).

Algoritmul a fost testat pe diverse seturi de date constând în arhive de conversații cu sute/mii de mesaje și ontologii definite în limba engleză. Pentru stabilirea valorilor optime v_1 , v_2 s-a rulat algoritmul cu valori incrementale ale variabilelor. Astfel, s-au obținut rezultatele descrise în Tabelul 2, din care se observă că termenii regăsiți direct în ontologii sunt mai relevanți decât cei găsiți prin aplicarea corecției ortografice. Ontologiile adăugate în sistem sunt preluate dintr-o sursă publică, ce poate fi regăsită la <http://semanticweb.org/wiki/Ontology> și descriu numeroase concepte și domenii de activitate (ex: muzică, sporturi, date biografice, filme, etc). Totodată, în cadrul aplicației a fost integrată și ontologia de Human-Computer Interaction dezvoltată în cadrul cursului de Interacțiune Om-Calculator din cadrul Universității Politehnica din București, Departamentul Calculatoare.

Tabelul 2. Rezultatele experimentale pentru clasificarea mesajelor

Mesaje în arhivă	Concepte în ontologii	v_1	v_2	Clasificare adecvată a arhivei într-un domeniu
< 200	< 3000	0.25	0.5	Nu
< 200	< 3000	0.5	0.75	Nu
< 200	> 3000	1	0.5	Da
< 200	> 3000	1	0.75	Da
> 200	> 3000	0.25	0.75	Nu
> 200	> 3000	0.75	0.25	Da
> 200	> 3000	0.75	0.5	Da
> 200	> 3000	1	0.5	Da

În final, se aleg valorile maxime din H_o , iar cheile corespunzătoare acestor valori sunt domeniile care au apărut cel mai frecvent în A . În cazul în care valorile sunt sub o anumită limită sau egale toate cu 0, se consideră că arhiva nu a putut fi clasificată.

4. Descrierea aplicației dezvoltate

Aplicația web dezvoltată facilitează comunicarea între doi sau mai mulți participanți și este bazată pe camere de discuție, cu posibilitatea creării unui număr nelimitat de astfel de camere, dar și a invitării de prieteni. Utilizatorii se pot înregistra în aplicație pentru a lua parte la discuțiile la care au acces sau pot iniția sesiuni în camere noi. Astfel se pot crea 2 tipuri de canale de discuție: (1) *chat*, având scopul de a facilita comunicarea între 2 utilizatori (1:1), respectiv (2) *chatroom* pentru comunicarea între mai mulți utilizatori (1:N). Fiecare cameră are propriul istoric al discuțiilor, iar acestea se pot exporta oricând în format XML sau JSON, pentru a putea fi prelucrate și analizate independent de aplicația curentă.

Spre deosebire de marea majoritate a aplicațiilor care oferă astfel de soluții, o funcționalitate aparte o reprezintă posibilitatea adăugării de legături explicite, facilitate regăsită și în ConcertChat. În acest mod, orice mesaj dintr-o cameră poate fi referențiat în momentul adăugării unei noi replici, iar toate referințele unui mesaj pot fi vizualizate în interfață, fiind marcate corespunzător. Un utilizator își poate edita un mesaj imediat după postare, putând aplica astfel corecții ortografice sau putând aduce o îmbunătățire semantică.

Interfața și serverul au fost dezvoltate folosind JavaScript. Întrucât experiența oferită utilizatorului, precum și viteza de încărcare și navigare sunt criteriile care au fost considerate necesare, s-a optat pentru dezvoltarea unei aplicații care încarcă o singură pagină la început, apoi îi modifică dinamic conținutul în funcție de interacțiunea utilizatorului cu aceasta (Single Page Application) prin intermediul websockets. Serverul pune la dispoziție un API (interfață destinată programatorilor) securizat, bazat pe permisiunile utilizatorilor și acționează ca un proxy între interfață și baza de date. Beneficiind de o astfel de arhitectură, aplicația poate fi transferată și ca aplicație nativă pe sisteme desktop, prin Chromium Embedded Framework sau pe sisteme mobile, prin Apache Cordova

Aplicația este construită din mai multe module: autentificare, panou de control, camere de discuție, traducere mesaje, vizualizare termeni regăsiți în ontologie. Prin modulul de autentificare se pot crea noi conturi sau iniția sesiuni persistente cu serverul. După autentificare, se realizează redirectarea către Dashboard (Panoul de control). De aici, un utilizator poate vizualiza camerele de discuție la care are acces, prietenii și poate adăuga noi camere de discuție.

Datorită creșterii tot mai mari a accesului la Internet prin intermediul dispozitivelor mobile, s-a considerat necesară și posibilitatea adaptării automate a interfeței la rezoluțiile specifice acestor dispozitive. Astfel, s-a decis că dezvoltarea trebuie să se axeze în primul rând pe dispozitive mobile. În acest mod, sunt accesibile doar funcționalitățile minime necesare prin care un utilizator poate folosi aplicația la capacitate maximă (Figura 3.a). Când aplicația este accesată de pe un ecran mai mare (Figura 3.b), datorită spațiului adițional, sunt prezente mai multe elemente, fiind sporită astfel interacțiunea cu utilizatorul.

Pentru interfață au fost folosite următoarele biblioteci:

- bootstrap (dezvoltat de Twitter). S-a folosit sistemul grid (grilă) pentru a permite adaptarea automată a aplicației pe diferitele rezoluții ale dispozitivelor pe care va fi accesată.
- angular (dezvoltat de Google). S-a folosit pentru a permite dezvoltarea întregii arhitecturi de structurare, navigare și încărcare dinamică a datelor.

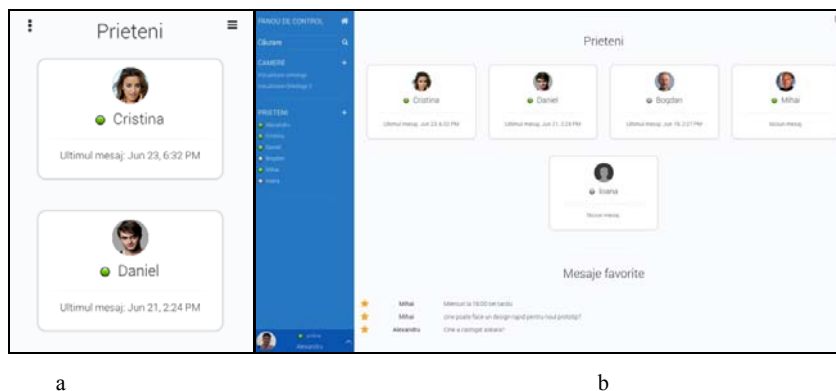


Figura 3. Diferența de prezentare a interfeței utilizator în funcție de rezoluție: a) Rezoluție 320 x 480, meniul adițional apare la apăsarea butonului de meniu, b) Rezoluție 1366 x 768, mai multă informație vizibilă

Serverul a fost dezvoltat cu ajutorul următoarelor tehnologii: nodejs, redis, mongodb. Redis este folosit ca mediu distribuit de stocare al sesiunilor utilizatorilor, iar mongodb pentru persistența datelor. Aplicația permite comunicarea între client și server prin websockets și trimiterea de notificări în timp real către clienți ca urmare a actualizărilor făcute în modelele de date. Astfel, în momentul în care un utilizator adaugă o nouă resursă (de exemplu, scrie un nou mesaj), acesta este mai întâi salvat în baza de date. Apoi, este transmisă o notificare către toate modelele ce sunt abonate la actualizările resursei tocmai modificate, în cazul de față al mesajului tocmai salvat. Aceste actualizări sunt trimise către celelalte aplicații client ce sunt conectate în acel timp la server. Întrucât aplicația este dinamică, actualizările sunt reflectate automat în regiunile din interfață ce depind de modelele de date modificate.

4.1 Detecția limbii și traducerea mesajelor

Înțelegerea limbii folosite în discuție este necesară pentru detecția și corecția automată a greșelilor de ortografie comise de utilizatori. Astfel, participanții la dialog pot transmite un mesaj mai clar și corect, iar contextul discuțiilor poate fi ulterior mai bine înțeles.

Posibilitatea traducerii mesajelor în alte limbi elimină barierele lingvistice care pot împiedica participarea anumitor persoane la dialog. Oferirea unei modalități cât mai simple de traducere a unui mesaj în limba dorită este o funcționalitate ce va ajuta un număr mai mare de participanți să contribuie activ, dar și alte persoane să înțeleagă subiectul discutat.

Serverul realizează o detecție a limbii fiecărui mesaj trimis, folosind algoritmul n-gram. Limba detectată este folosită ulterior dacă se dorește traducerea mesajului sau dacă acesta are greșeli ortografice. În stadiul actual sunt suportate 52 de limbi prin intermediul serviciului web pus la dispoziție de Bing.

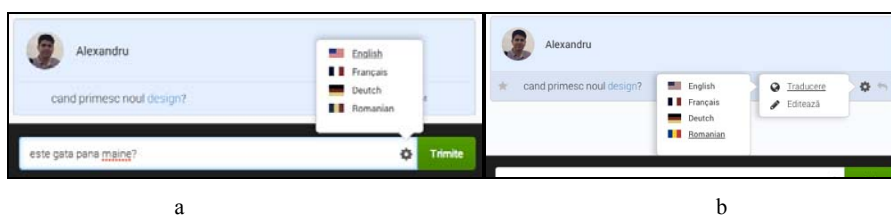


Figura 4. Modalitățile de traducere ale unui mesaj oferite de aplicație

Mesajul se poate traduce fie înainte de a fi trimis, fie dacă este prezent în conversație. Mesajul poate fi tradus automat înainte de a fi postat (figura 4.a), sau poate fi selectat orice mesaj din camera pentru traducere ulterioară (figura 4.b).

4.2 Vizualizarea ontologiilor

Într-o cameră de discuție, în momentul în care este introdus un nou mesaj (Figura 5), acesta este analizat și, dacă acesta se regăsește în ontologiile predefinite, acesta poate fi vizualizat, folosind biblioteca JavaScript theJit, sub forma unui graf radial generat pe baza relațiilor dintre conceptele ontologiei selectate (Figurile 6 și 7).



Figura 5. Detecția în ontologie a unui cuvânt din mesaj

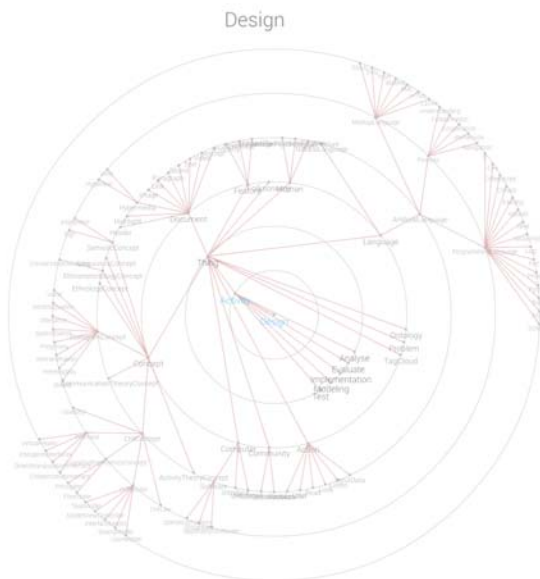


Figura 6. Vizualizarea grafului radial în care a fost regăsit cuvântul *Design*

De asemenea, identificarea limbii folosite într-o conversație ajută la o mai bună înțelegere a contextului. Cu cât detecția este mai bună, se pot oferi sugestii mai precise pentru corecția ortografică a cuvintelor. Algoritmii statistici prezentați, atât cel de corecție cât și cel pentru detecția limbii au o precizie mai bună cu cât volumul inițial de date de antrenare este mai mare. Pentru corecție s-au folosit doar date de antrenare în limba engleză, dar principiile sunt aplicabile la nivelul celorlalte limbi suportate.

Combinând acești algoritmi și având un set cât mai mare de ontologii disjuncte la nivelulul cărora se efectuează regăsirea termenilor conversațiilor purtate, se poate determina subiectul discuțiilor. Existența mai multor ontologii nu duce obligatoriu la o clasificare mai bună întrucât elementele comune în mai multe ontologii introduc ambiguitate și interpretabilitate. Astfel, acestea pot fi clasificate cât mai precis, în funcție de domeniul din care fac parte ontologiile predominante.

Ca direcții de dezvoltare ulterioară, se consideră necesară traducerea ontologiilor pentru o mai bună mapare a acestora pe termenii fiecărei limbi, precum și extinderea setului curent de ontologii folosite. Desigur, utilizarea platformei în scenarii educaționale integrate în cadrul unei platforme de e-learning este considerată o extensie naturală a eforturilor actuale (Iftene & Rotaru, 2010). Adicional, avem în vedere două funcționalități suplimentare. Pe de o parte, pentru camerele de discuții care nu au subiectul clasificat, se va cere clasificarea manuală de către utilizator, facilitând astfel crearea ulterioară a unui algoritm de învățare automată pentru îmbunătățirea predicției. Pe de altă parte, pornind de la exportul conversației, ne dorim integrarea cu platforma *ReaderBench* (Dascălu et al., 2013; Dascălu, 2014) în vederea evaluării automate a gradului de participare și a colaborării participanților. Astfel, aplicația de derulare a conversațiilor integrată cu platforma *ReaderBench* va permite derularea de multiple experimente centrate pe evaluarea conversațiilor chat derulate în diverse scenarii educaționale (spre ex., rezolvarea colaborativă de problemelor, dezbateri).

Mulțumiri

Rezultatele prezentate în acest articol au fost obținute cu sprijinul Ministerului Fondurilor Europene prin Programul Operațional Sectorial Dezvoltarea Resurselor Umane 2007-2013, Contract nr. POSDRU/159/1.5/S/134398.

Referințe

- Angular, disponibil online la <https://angularjs.org/>, accesat la data de 18.06.2014
- Apache Cordova, disponibil online la <http://cordova.apache.org/>, accesat la data de 18.06.2014
- Appear.i, disponibil online la <https://appear.in/>, accesat la data de 18.06.2014
- Asana, disponibil online la <https://asana.com/>, accesat la data de 18.06.2014
- Babel Fish, disponibil online la <http://www.babelfish.com/>, accesat la data de 18.06.2014
- Berners-Lee, T., Hendler, J., Lassila, O. (2001), Scientific American, *The Semantic Web*
- Bing Translator, disponibil online la <http://www.bing.com/translator/>, accesat la data de 18.06.2014
- Bootstrap, disponibil online la <http://getbootstrap.com/>, accesat la data de 18.06.2014
- Bright, P. (2011), *Microsoft Buys Skype for \$8.5 Billion. Why, Exactly?*, Wired, publicat online la <http://www.wired.com/2011/05/microsoft-buys-skype-2/>, accesat la data de 18.06.2014
- Campfire, disponibil online la <https://campfirenow.com/>, accesat la data de 18.06.2014
- Cavnar, W.B., Trenkle, J.M. (1994), *N-Gram-Based Text Categorization*, Environmental Research Institute of Michigan
- Chromium Embedded Framework, disponibil online la <https://code.google.com/p/chromiumembedded/>, accesat la data de 18.06.2014
- Cisco, *Cisco's Visual Networking Index Forecast Projects Nearly Half the World's Population Will Be Connected to the Internet by 2017*, The Network Cisco's Technology News Site, disponibil online la <http://newsroom.cisco.com/release/1197391/>, 2013, accesat la data de 18.06.2014
- ConcertChat, disponibil online la <http://sourceforge.net/projects/concertchat/>, accesat la data de 18.06.2014
- Dascălu, M. (2014), *Analyzing discourse and text complexity for learning and collaborating*, Studies in Computational Intelligence. Springer, Switzerland
- Dascălu, M., Trăușan -Matu, S., Dessus, P. (2013), *Cohesion-based analysis of CSCL conversations: Holistic and individual perspectives*, In 10th Int. Conf. on Computer-Supported Collaborative Learning (CSCL 2013), N. Rummel, M. Kapur, M. Nathan and S. Puntambekar Eds. ISLS, Madison, USA, pp. 145–152
- Deshpande, Y., Murugesan, S., Ginige, A., Hansen, S., Schwabe, D., Gaedke, M., White, B. (2002), *Web Engineering*, Journal of Web Engineering, 1 (1), pp. 3-17
- Google Translate, disponibil online la <https://translate.google.com/>, accesat la data de 18.06.2014
- HipChat, disponibil online la <https://www.hipchat.com/>, accesat la data de 18.06.2014
- Iftene, A., Rotaru, A. 2010. *Îmbunătățirea unui sistem de eLearning folosind elemente specifice unui sistem de tip întrebare-răspuns*. In Revista Română de Interacțiune Om-Calculator (RRIOC), Vol. 3 (2010) Numar special: Interacțiune Om-Calculator 2010,

- ISSN 1843-4460, Pp. 71-74, Bucuresti, 2-3 Septembrie, 2010.
- Internet Working Group, *Internet Relay Chat Protocol*, disponibil online la <https://tools.ietf.org/html/rfc1459>, accesat la data de 18.06.2014
- Jira, disponibil online la <https://www.atlassian.com/software/jira>, accesat la data de 18.06.2014
- Levenshtein, V.I. (1966), "*Binary codes capable of correcting deletions, insertions, and reversals*". Soviet Physics Doklady, 10 (8), pp. 707–710
- McAfee, A., Brynjolfsson, E. (2012), *Big Data: The Management Revolution*, Harvard Business Review, publicat online la <http://hbr.org/2012/10/big-data-the-management-revolutionar>, accesat la data de 18.06.2014
- Mongodb, disponibil online la www.mongodb.org/, accesat la data de 18.06.2014
- Mozilla Developer Network, *WebRTC*, disponibil online la <https://developer.mozilla.org/en-US/docs/Web/Guide/API/WebRTC>, accesat la data de 18.06.2014
- Nodejs, disponibil online la <http://nodejs.org/>, accesat la data de 18.06.2014
- Norvig, *How to Write a Spelling Corrector*, disponibil online la <http://norvig.com/spell-correct.html>, accesat la data de 18.06.2014
- Powell, T., Jones, D., Cutts, D. (1998), *Web Site Engineering: Beyond Web Page Design*, Prentice Hall
- Quora, disponibil online la <http://www.quora.com/>, accesat la data de 18.06.2014
- Reddit, disponibil online la <http://redis.io/>, accesat la data de 18.06.2014
- Resnik, P., 1995. Using information content to evaluate semantic similarity in a taxonomy. In 14th Int. Joint Conf. on Artificial Intelligence (IJCAI'95) Morgan Kaufmann, Montreal, Canada, 448–453.
- Single Page Application, disponibil online la <http://www.johnpapa.net/pageinspa/>, accesat la data de 18.06.2014
- Slack, disponibil online la <https://slack.com/>, accesat la data de 18.06.2014
- Stackoverflow, disponibil online la <http://stackoverflow.com/>, accesat la data de 18.06.2014
- Stahl, G., Group cognition. Computer support for building collaborative knowledge. MIT Press, Cambridge, MA, 2006.
- Stone, B. (2014), *Facebook Buys WhatsApp for \$19 Billion*, Bloomberg Businessweek, publicat online la <http://www.businessweek.com/articles/2014-02-19/facebook-acquires-whatsapp-for-19-billion>, accesat la data de 18.06.2014
- theJit, disponibil online la <http://philogb.github.io/jit/>, accesat la data de 18.06.2014
- Trăușan -Matu, S., Rebedea, T. (2010), *A polyphonic model and system for inter-animation analysis in chat conversations with multiple participants*, In 11th Int. Conf. Computational Linguistics and Intelligent Text Processing (CICLing 2010), A.F. Gelbukh Ed. Springer, New York, pp. 354–363
- Trello, disponibil online la <https://trello.com/>, accesat la data de 18.06.2014