

Programul: IDEI

Tipul proiectului: Proiecte de cercetare exploratorie

Cod proiect: ID_2286

Sinteza lucrării

CERCETĂRI ÎN DIRECȚIA OPTIMIZĂRII ADAPTIVE A SISTEMELOR INFORMATICE FOLOSIND TEHNICI DE ÎNVĂȚARE AUTOMATĂ ȘI SISTEME MULTIAGENT

Etapa unică an 2011

A. Sinteza cercetării. Rezultate obținute

În cele ce urmează vom prezenta o sinteză a rezultatelor originale obținute în urma cercetărilor efectuate în cadrul proiectului în scopul îndeplinirii obiectivelor științifice propuse în planul de realizare a proiectului pe anul 2011. Vom indica, pentru fiecare obiectiv prevăzut în planul de realizare pe anul 2011, modul în care au fost îndeplinite activitățile aferente. De asemenea, vom menționa rezultatele originale obținute în vederea realizării obiectivelor.

Menționez faptul că obiectivele planificate pe anul 2011, cât și activitățile aferente acestora au fost realizate în totalitate, și desfășurate conform cu planul de realizare al proiectului. De asemenea, criteriul minim de performanță prevăzut (1 articol ISI și 3 articole BDI) a fost îndeplinit, așa cum se poate vedea în secțiunea legată de modul de diseminare al rezultatelor.

1. Elaborarea unor noi metode pentru problema introducerii șabloanelor de proiectare în sisteme informatice existente folosind tehnici de învățare automată (clustering, rețele cu autoorganizare -SOM, rețele neuronale)

Fie $S = \{C_1, C_2, \dots, C_l\}$ un sistem informatic orientat obiect, unde $C_i, 1 \leq i \leq l$ este o clasă a sistemului.

Fără a restrânge generalitatea abordării, vom considera că vrem să identificăm în sistemul S locațiile unde ar trebui introdus un anumit șablon de proiectare p . Pentru început, ne-am orientat asupra șabloanelor de proiectare structurale.

Un șablon de proiectare p din sistemul informatic S poate fi privit ca fiind o pereche, $p = (C_p, R_p)$, unde

- $C_p = \{C_1^p, C_2^p, \dots, C_{n_{C_p}}^p\}$ este o submulțime a mulțimii S , $C_p \subset S$, și reprezintă mulțimea claselor care fac parte din șablonul de proiectare p . n_{C_p} reprezintă numărul de clase din șablon.
- $R_p = \{r_1^p, r_2^p, \dots, r_{nr_p}^p\}$ este mulțimea restricțiilor (relațiilor) existente între clasele din C_p , restricții care caracterizează șablonul de proiectare p . Ca urmare, fiecare restricție $r_i^p, \forall 1 \leq i \leq nr_p$ din mulțimea R_p este o relație definită pe o submulțime de clase din C_p , iar nr_p reprezintă numărul de restricții care caracterizează șablonul de proiectare p .

Menționăm faptul că toate restricțiile din mulțimea R_p pot fi exprimate sub forma unor restricții binare (sunt două clase participante la restricție).

Problema identificării locațiilor unde ar trebui introdus șablonul p se poate reduce la problema identificării unor submulțimi de clase din S al căror comportament este similar cu cel indicat de șablonul p . Pentru a ilustra acest grad de similaritate, am introdus o măsură Sim care indică similaritatea între o mulțime de clase $\{C_{l_1}, \dots, C_{l_k}\} \subset S$ și șablonul $p = (C_p, R_p)$.

Ca urmare, problema studiată am redus-o la a determina submulțimi ale lui S a căror similaritate Sim în raport cu p este mai mare decât un prag $Prag$ dat. Pentru rezolvarea acestei probleme, care poate fi văzută de fapt ca și o problemă de satisfacere a restricțiilor, am introdus un algoritm de căutare sistematică de tip BackJumping [12].

În continuare, vom descrie pe scurt algoritmul BJ. BJ este similar cu algoritmul backtracking cronologic (BT), exceptând doar că este mai eficient atunci când nu se găsește nici o instanțiere consistentă pentru variabila curentă (dead-end). C și la BT, variabilele x_1, x_2, \dots, x_N care trebuie instanțiate pentru a determina soluția, sunt instanțiate secvențial, în ordinea x_1, x_2, \dots , astfel încât, la orice moment instanțierea parțială să fie consistentă. Când apare un dead-end, BJ sare la variabila de indice maxim care a fost conflict cu variabila curentă.

Verificările consistenței între instanțierea variabilei curente și instanțierea variabilei trecute sunt îndeplinite pe măsură ce ordinea originală a instanțierilor este (x_1, x_2, \dots) . Verificările se opresc în momentul în care o verificare a consistenței eșuează. Variabila $max_check[i]$ reține indicele maxim al variabilei care a fost verificată față de instanțierea curentă x_i . Această valoare este folosită pentru determinarea punctului de backtracking la

finalul funcției main. Problema principală cu BJ este că se fac salturi doar din nodurile terminale ale arborelui (dead-ends). Toate celelalte căutări fac parte din backtracking-ul cronologic, deci acolo este sigur o mulțime de trashing-uri. Pe de altă parte, costurile sunt mai mici la BJ.

Implementarea, evaluarea și validarea algoritmului BJ

După cum am arătat anterior, a fost implementat în JADE un prototip al sistemului *DSSEM*, în care a fost adăugată funcționalitatea de introducere a șabloanelor de proiectare în sisteme informatice existente.

Pentru evaluare experimentală, am considerat un studiu de caz simplu, un sistem informatic având 5 clase în care, în scopul îmbunătățirii design-ului ar trebui introdusă o instanță a șablonului de proiectare *Proxy*. *Proxy* este un șablon de proiectare structural care surprinde o tehnică des folosită în rețele: un obiect local comunică cu un proces aflat la distanță, dar se dorește ascunderea detaliilor despre procesul la distanță sau protocolul de comunicare. Obiectul *proxy* permite accesul la serviciile la distanță cu aceeași interfață a proceselor locale. De fapt, atunci când o anumită *Operație* este solicitată obiectului proxy, acesta delegă execuția operației solicitate obiectului *SubiectReal*. Deoarece atât *Proxy* cât și *SubiectReal* sunt subclase ale lui *Subiect*, acestea exportă aceeași interfață pentru *Operație*. Pentru a putea apela metodele obiectului *SubiectReal*, *Proxy* necesită o asociere cu acesta.

În urma aplicării algoritmului BJ pe care l-am descris anterior, a fost identificată o sumbuțime formată din 3 clase care ar trebui să reprezinte o instanță a șablonului de proiectare *Proxy*. În continuare, ne propunem să îmbunătățim și să extindem abordarea propusă în acest capitol în direcția aplicării abordării propuse pe sisteme soft reale.

2. Dezvoltarea de noi algoritmi de învățare automată (clustering, rețele cu autoorganizare -SOM, rețele neuronale) în scopul autoadaptării din punct de vedere comportamental al sistemelor informatice

2.1 Dezvoltarea de algoritmi bazați pe tehnici de învățare automată pentru autoadaptarea structurilor de date folosite la fluxul de date din sistem

Este bine știut faptul că tipurile abstracte de date reprezintă nucleul logic al aplicației soft, iar o utilizare eficientă a acestora este o cerință esențială pentru dezvoltarea unor sisteme robuste și eficiente. Structurile de date sunt esențiale în obținerea unor algoritmi eficienți, având o importanță majoră în procesul de dezvoltare de soft. Selectarea și crearea celei mai potrivite structuri de date pentru implementarea unui tip abstract de date poate avea un impact mare asupra performanței și eficienței sistemului soft. Nu este o problemă trivială pentru un dezvoltator soft, deoarece sunt greu de anticipat toate scenariile de utilizare ale aplicației analizate, iar o selectare statică înaintea execuției sistemului nu este, în general, foarte precisă.

În articolul [4] am studiat problema selectării dinamice a structurilor de date eficiente pentru implementarea tipurilor abstracte de date, folosind o abordare supervizată. În scopul selectării dinamice a celei mai potrivite reprezentări pentru un container de date în funcție de contextul curent al execuției sistemului, vom folosi o rețea neuronală. Am evaluat tehnica propusă pe un studiu de caz, evidențiind avantajele modelului propus în comparație cu abordări similare existente în literatură.

Considerații teoretice

Fie $S = \{s_1, s_2, \dots, s_n\}$ un sistem informatic orientat obiect, unde $s_i, 1 \leq i \leq n$ poate fi o clasă, o metodă dintr-o clasă sau un atribut al unei clase. Ne vom referi la un element al sistemului informatic S ca la o *entitate*. Să considerăm următoarele:

- $Class(S) = \{C_1, C_2, \dots, C_l\}$, $Class(S) \subset S$, este mulțimea claselor din sistemul informatic S .
- Fiecare clasă C_i ($1 \leq i \leq l$) este o mulțime de metode și atribute, altfel spus, $C_i = \{m_{i1}, m_{i2}, \dots, m_{ip_i}, a_{i1}, a_{i2}, \dots, a_{ir_i}\}$, $p_i \leq n$, $1 \leq r_i \leq n$, unde m_{ij} ($\forall j, 1 \leq j \leq p_i$) sunt metodele și a_{ik} ($\forall k, 1 \leq k \leq r_i$) sunt atributele din clasa C_i .
- $Meth(S) = \bigcup_{i=1}^l \bigcup_{j=1}^{p_i} m_{ij}$, $Meth(S) \subset S$, este mulțimea metodelor din toate clasele sistemului informatic S .
- $Attr(S) = \bigcup_{i=1}^l \bigcup_{j=1}^{r_i} a_{ij}$, $Attr(S) \subset S$, este mulțimea atributelor din toate clasele sistemului informatic S .

Pe baza notațiilor anterioare, sistemul informatic S poate fi definit ca în Ecuația (1):

$$S = Class(S) \cup Meth(S) \cup Attr(S). \quad (1)$$

Fie \mathcal{T} un tip abstract de date având în interfața sa o mulțime de operații \mathcal{O} care pot fi efectuate asupra unei instanțe a tipului \mathcal{T} . De asemenea, considerăm o mulțime $\mathcal{DT} = \{\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_n\}$ de structuri de date care pot fi utilizate în sistemul S pentru implementarea lui \mathcal{T} . Dacă o anumită instanță c a clasei $C \in Class(S)$

folosește \mathcal{T} , selectarea celei mai potrivite structuri de date din \mathcal{D} care să fie folosită pentru implementarea (și instanțierea) lui \mathcal{T} depinde de *contextul de execuție* curent (starea curentă a lui c). În Definiția 1 vom introduce noțiunea de *context de execuție* a lui \mathcal{T} .

Definiția 1 *Context de execuție a lui \mathcal{T} .*

Fie c o instanță a clasei $C_i \in \text{Class}(S)$ care folosește \mathcal{T} . Contextul de execuție a lui \mathcal{T} reprezintă starea clasei C_i la execuție, atunci când \mathcal{T} este inițializat. Ca urmare, este definit ca mulțimea $\mathcal{EC}_{\mathcal{T}} = \{v_{i1}, v_{i2}, \dots, v_{ir_i}\}$, unde r_i reprezintă numărul de atribute din clasa C_i și $\forall 1 \leq j \leq r_i$, v_{ij} reprezintă valoarea atributului a_{ij} în pasul curent de execuție.

În Definiția 2 vom defini noțiunea de *implementare potrivită* a lui \mathcal{T} pe baza contextului de execuție.

Definiția 2 *Cea mai potrivită implementare a lui \mathcal{T} pe baza contextului execuției $\mathcal{EC}_{\mathcal{T}}$.*

Fie \mathcal{T} un tip abstract de date și $\mathcal{EC}_{\mathcal{T}}$ contextul execuției lui \mathcal{T} într-un anumit scenariu de utilizare. Cea mai potrivită implementare a lui \mathcal{T} pe baza contextului $\mathcal{EC}_{\mathcal{T}}$ este structura de date $\mathcal{D}_i \in \mathcal{D}$ care furnizează o implementare eficientă a lui \mathcal{T} în scenariul de execuție, adică timpul necesare pentru efectuarea operațiilor lui \mathcal{T} în contextul $\mathcal{EC}_{\mathcal{T}}$ este minimizat.

De fapt, cea mai potrivită implementare a lui \mathcal{T} într-un anumit scenariu de execuție depinde de tipul și numărul de operații din \mathcal{O} care sunt efectuate asupra lui \mathcal{T} în timpul execuției. Ca urmare, putem concluda că problema selectării dinamice a celei mai potrivite implementări a lui \mathcal{T} este, de fapt, problema predicției, pe baza contextului curent al execuției, numărului și tipului operațiilor efectuate asupra lui \mathcal{T} într-un anumit scenariu de execuție. Pentru predicția supervizată a celei mai potrivite implementări a tipului abstract de date \mathcal{T} am folosit în [4] o rețea neuronală.

Evaluare experimentală

În scopul evaluării acurateței modelului ANN propus anterior, am simulat un experiment pentru simularea celei mai potrivite structuri de date pentru implementarea tipului abstract de date *List*. Setul de date considerat pentru analiză constă în rezultatele unei analize chimice a vinurilor care cresc într-o anumită regiune a Italiei, dar obținută de lădiferiți cultivatori. Analiza efectuată a determinat 13 constituente găsiți în fiecare tip de vin [13]. Pentru testare, s-a folosit o validare încrucișată de tipul 2/3-1/3 (2/3 instanțe pentru antrenare, 1/3 instanțe pentru testare). S-au folosit 50 de episoade de testare. O acuratețe globală de 0.98133 a fost obținută pe cele 50 de episoade. Rezultatele obținute de-a lungul episoadelor sunt stabile, fiind obținută o deviație standard de 0.01231079 pe setul de precizii obținut. Valoarea scăzută a deviației standard indică precizia abordării noastre.

2.2 Dezvoltarea unei metode bazate pe învățare supervizată și agenți pentru predicția comportamentului utilizatorului

În trecut, utilizatorii sistemelor informatice erau cei care le și concepeau, astfel încât nu existau probleme de utilizabilitate. În prezent, în orice domeniu al vieții economice sunt utilizate calculatoarele, iar nivelul de expertiză al utilizatorilor scade, o dată cu creșterea complexității sarcinilor. În acest context, este necesară identificarea de soluții care să ajute utilizatorii novici în îndeplinirea sarcinilor. Soluția propusă în acest sens constă în dezvoltarea unui agent inteligent, denumit **LIA** (**L**earning **I**nterface **A**gent) pentru predicția comportamentului utilizatorului. Rolul acestui agent inteligent este acela de a învăța, în baza experienței anterioare și a contextului actual al interacțiunii, să prezică următoarea acțiune pe care ar trebui să o efectueze utilizatorul pentru a realiza sarcina.

În cele ce urmează, considerăm că agentul **LIA** monitorizează interacțiunile utilizatorilor cu o aplicație interactivă \mathcal{SA} în timp ce realizează o sarcină \mathcal{T} . Vom nota prin \mathcal{A} mulțimea $\{a_1, a_2, \dots, a_n\}$ tuturor acțiunilor posibile care pot apărea în timpul interacțiunii cu \mathcal{SA} . O acțiune poate fi: apăsarea unui buton, selectarea unei opțiuni dintr-un meniu, completarea unui câmp text. În timpul interacțiunii cu \mathcal{SA} pentru efectuarea sarcinii \mathcal{T} se generează *urme utilizator* ale interacțiunii. O *urmă utilizator* este o secvență de acțiuni utilizator. Vom spune că o urmă utilizator este *cu succes* dacă sarcina este îndeplinită. În celelalte cazuri avem o urmă utilizator *fără succes*. În abordarea curentă vom lua în considerare doar urmele cu succes.

Scopul urmărit este acela de a înzestra agentul **LIA** cu capacitatea de a prezice la un moment dat următoarea acțiune pe care utilizatorul ar trebui să o efectueze pentru a realiza sarcina \mathcal{T} . Pentru a atinge acest scop, propunem o tehnică de învățare supervizată care constă în doi pași: *antrenarea*, respectiv *predicția*. Pe parcursul antrenării, agentul **LIA** monitorizează interacțiunea unui set de utilizatori reali în timp ce realizează sarcina \mathcal{T} , folosind aplicația \mathcal{SA} și își construiește o bază de cunoștințe \mathcal{KB} . Interacțiunea este monitorizată folosind AOP. O abordare mai generală poate fi obținută prin construirea a două baze de cunoștințe în timpul antrenării: una cu urme utilizator cu succes, cealaltă cu urme utilizator fără succes.

După ce antrenarea a avut loc, în pasul de predicție se urmărește predicția comportamentului pentru un nou utilizator U pe baza datelor culese în timpul antrenării, folosind un model probabilistic. După fiecare acțiune *act* efectuată de un utilizator U , exceptând prima sa acțiune, agentul **LIA** va prezice următoarea acțiune, a_r ($1 \leq r \leq n$), care ar trebui efectuată cu o anumită probabilitate $P(\text{act}, a_r)$, folosindu-se de \mathcal{KB} .

Probabilitatea $P(act, a_r)$ este dată de Egalitatea (2).

$$P(act, a_r) = \max\{P(act, a_i), 1 \leq i \leq n\} \quad (2)$$

Pentru a calcula aceste probabilități vom introduce conceptul de *scor* între două acțiuni. Scorul dintre acțiunile a_i și a_j , notat prin $score(a_i, a_j)$ indică gradul în care a_j trebuie să-i urmeze lui a_i într-o execuție cu succes a sarcinii \mathcal{T} . Acest fapt sugerează că valoarea $score(a_i, a_j)$ este maximă atunci când a_j trebuie să urmeze imediat după a_i într-o execuție cu succes a sarcinii. Scorul dintre o acțiune dată act și o acțiune utilizator a_q , $1 \leq q \leq n$, $score(act, a_q)$, se calculează ca în Egalitatea (3).

$$score(act, a_q) = \max\left\{\frac{1}{dist(t_i, act, a_q)}, 1 \leq i \leq m\right\} \quad (3)$$

unde $dist(t_i, act, a_q)$ reprezintă *distanța* dintre două acțiuni act și a_q într-o urmă t_i , calculată folosind \mathcal{KB} .

$$dist(t_i, act, a_q) = \begin{cases} length(t_i, act, a_q) - 1 & \text{dacă } \exists sub_{t_i}(act, a_q) \\ \infty & \text{altfel} \end{cases} \quad (4)$$

$length(t_i, act, a_q)$ definește distanța minimă dintre act și a_q în urma t_i .

$$length(t_i, act, a_q) = \min\{\|s\| \mid s \in SUB_{t_i}(act, a_q)\} \quad (5)$$

În abordarea noastră, $length(t_i, act, a_q)$ reprezintă numărul minim de acțiuni efectuate de utilizatorul U în urma t_i , pentru a ajunge de la acțiunea act la acțiunea a_q , adică, lungimea minimă dintre toate suburmele posibile $sub_{t_i}(act, a_q)$.

Din Egalitatea (3), rezultă că $score(act, a_q) \in [0, 1]$ și că valoarea lui $score(act, a_q)$ crește odată cu descrescerea *distanței* dintre act și a_q în urmele din \mathcal{KB} .

Pe baza *scorurilor*, $P(act, a_i)$, $\forall 1 \leq i \leq n$, se calculează după cum urmează:

$$P(act, a_i) = \frac{score(act, a_i)}{\max\{score(a_i, a_j) \mid 1 \leq j \leq n\}} \quad (6)$$

Pe baza Ecuației (6), probabilități mai mari sunt asociate acțiunilor cele mai potrivite pentru a fi executate.

Rezultatul predicției agentului este acțiunea a_r care satisface Egalitatea (2). Mențion-ăm că într-un caz nedeterministic (când mai multe acțiuni au aceeași probabilitate maximă P) se poate folosi o metodă de selecție suplimentară. Pentru a evalua acuratețea predicției agentului **LIA**, comparăm secvența de acțiuni efectuate de utilizatorul U care duc la îndeplinirea sarcinii cu secvența de acțiuni precise de agent. Considerăm că predicția unei acțiuni este precisă dacă probabilitatea predicției este mai mare decât un prag fixat. În acest scop, am definit o măsură de calitate, $ACC(\mathbf{LIA}, U)$, denumită *ACCuracy* (precizie). Evaluarea acurateții predicției agentului **LIA** s-a realizat pe un studiu de caz, folosind măsurile de calitate anterior definite. Am ales spre evaluare o aplicație interactivă de dimensiuni medii.

3. Alte rezultate obținute

În lucrarea [3] am introdus o abordare originală bazată pe rețele cu autoorganizare și algoritmul *SOMSYR* (*Self-Organizing Map for SYstems Restructuring*), al cărui scop este determinarea refactorizărilor necesare pentru îmbunătățirea structurii de clase a sistemelor soft. Pentru construcția rețelei cu autoorganizare, am folosit topologia de *tor*. Evaluarea algoritmului s-a făcut pe sistemul open source JHotDraw si un sistem informatic real.

Optimizarea combinatorială are ca obiect căutarea de soluții în spațiul discret și bine-definit al unei probleme. Metodele de optimizare au o deosebită importanță practică, în special în procesele de proiectare, în experimentele științifice, precum și în procesele decizionale în afaceri. Datorită importanței practice majore a problemelor de optimizare combinatorială, o serie de metode au fost propuse pentru abordarea lor. Una dintre acestea este generarea de soluții inteligente folosind tehnici de învățare prin întărire. Învățarea prin întărire este o abordare a instruirii automate, în care un agent poate învăța să se comporte într-un anumit fel prin alegerea unor acțiuni și primirea de pedepse, sau recompense, în funcție de acțiunile alese. Principalul obiectiv al lucrării [10] este propunerea și prezentarea unei interfețe de programare pentru rezolvarea problemelor de optimizare combinatorială, folosind tehnici bazate pe învățarea prin întărire. Această interfață permite dezvoltarea simplă de aplicații pentru rezolvarea problemelor de optimizare combinatorială, folosind învățarea prin întărire. De asemenea, lucrarea oferă și un experiment pentru rezolvarea problemei bidimensionale a plierii proteinei, proiectat cu ajutorul interfeței propuse.

În lucrarea [8] propunem un model bazat pe învățarea prin întărire pentru rezolvarea problemelor de optimizare combinatorială. Rezolvarea acestui tip de probleme în mod optim este o sarcină dificilă, motiv pentru care se încearcă îmbunătățirea soluțiilor care se pot obține. Ne concentrăm în special pe problema plierii proteinei, o bine-cunoscută problemă NP-completă, de mare importanță în multe domenii, cum ar fi bioinformatica,

biochimia, biologia moleculară și medicina. Introducem un model bazat pe învățarea prin întărire pentru determinarea structurii bidimensionale a proteinei în modelul simplificat hidrofob-polar. Modelul propus în această lucrare poate fi extins cu ușurință pentru rezolvarea altor probleme de optimizare. Oferim, de asemenea, o validare matematică a modelului propus, bazat pe învățarea prin întărire, indicând astfel potențialul acestuia.

Lucrarea [9] se axează pe cercetarea modelului bazat pe învățarea prin întărire, introdus pentru determinarea structurii bidimensionale a proteinelor, folosind modelul simplificat hidrofob-polar. Modelul propus este fondat pe o abordare de tip *Q-learning*, bazată pe agenți. Evaluarea experimentală confirmă o bună performanță a modelului bazat pe învățarea prin întărire, evidențiind potențialul propunerii noastre. Cercetările efectuate în vederea propunerii unei abordări bazate pe un model distribuit de învățare prin întărire pentru rezolvarea problemelor de optimizare combinatorială sunt prezentate în lucrarea [5]. Ne focalizăm atenția pe problema bidimensională a plierii proteinei, care se referă la determinarea structurii bidimensionale a proteinei, folosind secvența primară de aminoacizi. Modelul pe care îl propunem se bazează pe o abordare distribuită de tip *Q-learning*, care extinde către o arhitectură distribuită un model pe care l-am propus în prealabil [8]. Experimentele au fost realizate în vederea obținerii de soluții pentru problema bidimensională a plierii proteinei, iar rezultatele obținute sunt promițătoare.

Reconstrucția unei molecule de ADN, utilizând fragmente de ADN, sau, altfel spus, problema asamblării fragmentelor de ADN este o problemă importantă în numeroase domenii, cum ar fi bioinformatica, biologia computațională, biotehnologia sau medicina. Datorită faptului că ea este NP-completă, au fost propuse și dezvoltate o serie de tehnici computaționale pentru rezolvarea ei, printre acestea numărându-se și diferiți algoritmi inteligenți. Deoarece asamblarea fragmentelor de ADN este de o importanță esențială în orice tip de proiect de secvențiere a ADN-ului, o mare parte din cercetători se concentrează pe direcția dezvoltării de algoritmi performanți de asamblare. Scopul principal al lucrării [7] este propunerea unui model bazat pe învățarea prin întărire pentru rezolvarea problemei asamblării de fragmente. Ne focalizăm atenția asupra problemei asamblării fragmentelor de ADN. Noul model introdus este axat pe o abordare bazată pe agenți, de tip *Q-learning*. Rezultatele experimentelor realizate demonstrează potențialul direcției propuse, precum și faptul că aceasta merită a fi cercetată în continuare.

După cum s-a menționat mai sus, problema asamblării fragmentelor de ADN este o problemă NP-completă, prin care se încearcă reconstituirea unei molecule de ADN pornind de la un număr mare de fragmente, fiecare dintre acestea conținând câteva sute de nucleotide. Modelul distribuit propus în [2] pentru abordarea problemei asamblării fragmentelor de ADN extinde modelul de învățare prin întărire propus în [7] și este de tip *Q-learning* concurrent, în care mai mulți agenți cooperează în scopul de a obține soluții optime. Lucrarea oferă și o evaluare experimentală a sistemului introdus, ale cărei rezultate sunt încurajatoare.

De asemenea, ne-am orientat cercetările spre îmbunătățirea performanței folosind calculul paralel. Pentru aplicarea cu succes a calculului paralel se impune existența a trei componente distincte: prima este reprezentată de hardware-ul pe care se vor executa aplicațiile; a doua este mașina abstractă sau modelul de programare în care sunt scrise aplicațiile, iar a treia este reprezentată de procesul de proiectare care permite crearea aplicațiilor din specificații.

O alegere posibilă pentru un model general de programare paralelă este BSP – *Bulk Synchronous Parallelism*. Un avantaj important al BSP este dat de faptul că măsurile de cost asociate estimează în mod real costurile. Pentru că separa partea de comunicare de sincronizare este în mod particular clar și simplu. Această separare permite o dezvoltare simplă și formalizarea procesului de construcție al softului. Pentru această formalizare propunem o metodă bazată pe specificații parametrizate și dezvoltare formală [6]. Un program este considerat ca fiind un număr de procese parametrizate cu structura similară care cooperează. Metoda folosește pre- și post- condiții parametrizate și consideră distribuția datelor chiar de la începutul procesului de construcție. Posibilitatea de a estima numărul de comunicații din postcondiții ne permite să facem o evaluare a costului chiar de la primele stagii de dezvoltare și aceasta este importantă pentru luarea deciziilor de construcție eficientă.

Abordarea noastră pentru construcția unui proces parametrizat *S.q* din specificațiile funcționale este similară cu metodele de derivare folosite în programarea secvențială. Se folosesc regulile clasice de derivare pentru programarea secvențială și introducem reguli noi pentru compunere paralelă și comunicare. Se folosește metoda clasică (bazată pe calcul) de obținere a unui invariant. Programele se derivează prin calcularea condițiilor necesare menținerii invarianțului. Pentru că procesele sunt parametrizate se vor deriva și invarianti parametrizați. În modelul BSP, partiționarea datelor este crucială spre deosebire de maparea partițiilor rezultate care este irelevantă. Alegerea unei distribuții adecvate de date este un mijloc important de influențare a performanței programului. Folosim în acest scop distribuții de date simple dar și multivoce și evidențiem importanța alegerii adecvate pe mai multe exemple de dezvoltare de programe paralele.

Prin folosirea conceptului de *tip abstract de date*, o structură de date poate fi definită mai riguros și mai formal. Un pas mai departe în această direcție s-a făcut odată cu introducerea programării orientate-obiect (OOP). Folosind OOP putem defini nu doar structuri de date generice prin polimorfism și genericitate, dar putem de asemenea să separăm definițiile de implementări prin folosirea interfețelor; interfețele descriu tipuri abstracte de date specificând operațiile corespunzătoare lor. Sabloanele de proiectare au dus lucrurile mai departe permițând obținerea unui grad sporit de flexibilitate și reutilizabilitate în implementarea structurilor de date. Pentru a construi o structură de date extensibilă și reutilizabilă este necesară decuplarea comportamentului

ei primitiv și intrinsec de comportamentul specific adaptat aplicației care o manipulează.

Orice tip de structură de date de tip container este formată dintr-un număr de elemente care în general au același tip. Un container particular are proprietăți și comportament care nu este dependent de tipul elementelor constitutive.

Prin studierea diferitelor structuri de date am evidențiat următoarea clasificare a acestora [11]: structuri de date de *nivel primar* sau fundamentale; structuri de date de *nivel secundar* care sunt caracterizate de faptul că implementările lor folosesc structurile de date fundamentale. Tablourile și reprezentările înlantuite pentru liste și arbori sunt considerate structuri de date fundamentale. Pentru a implementa o mulțime sau un dicționar putem folosi un tablou sau o listă înlantuită sau un arbore; astfel mulțimile și dicționarele sunt considerate exemple de structuri de date de nivel secundar.

În proiectarea propusă structurile de date fundamentale implementează o interfață **Storage** pentru a fi evidențiată proprietatea lor de a putea fi folosite pentru stocarea elementelor structurilor de date de nivel secundar. Procesul poate fi simplificat prin folosirea șablonului *Bridge* împreună cu șabloanele *Abstract Factory* și *Singleton*. Aceste șabloane asigură independența de stocare a structurilor de date de nivel secundar [11]. Prin separarea reprezentării unei structuri de date de comportamentul ei dat de tipul ei, introducem un nou nivel de indirectare și prin urmare și de abstractizare. Avantajul este dat de posibilitatea de a implementa structuri de date bazându-ne pe diferite structuri de date fundamentale fără a fi nevoie de a construi mai multe clase. Se introduce de asemenea și un nou nivel de genericitate. Interfețele de tip - **Storage** - pot introduce o altă clasificare a structurilor de date. Orice structură de date poate fi folosită ca o posibilitate de stocare pentru altă structură de date sau ca și o componentă de stocare într-un program generic. Programarea modernă folosește medii de programare care simplifică munca programatorului foarte mult. Aceste medii oferă un nivel înalt de abstractizare și astfel se permite crearea programelor prin specificarea modului în care componentele de baza sunt compuse. În acest context un nou nivel de abstractizare a datelor manipulate de aceste programe este dorit. Putem să considerăm un scenariu în care programatorul specifică doar tipul (comportamentul) al entității de stocare: simplu, secvențial, sortat, etc. Pe baza reprezentării propuse și, de asemenea, folosind mecanisme automate cum sunt tehnicile de inteligență artificială se pot alege tipurile concrete cele mai eficiente pentru aplicația dezvoltată.

B. Activități organizatorice

Activitățile organizatorice au avut ca scop menținerea unui cadru propice activității de cercetare și îndeplinirii cerințelor de această natură stipulate în contractul de finanțare. Pagina web a proiectului, dedicată prezentării proiectului, a echipei de cercetare și a rezultatelor obținute, <http://www.cs.ubbcluj.ro/~raoss>, a fost permanent actualizată. S-au efectuat întâlniri periodice ale echipei de cercetare, în scopul realizării cu succes a activităților prevăzute în plan și diseminarea în cadrul grupului a rezultatelor obținute.

C. Diseminarea rezultatelor

Diseminarea rezultatelor științifice obținute în anul 2011 în cadrul proiectului și prezentate în secțiunile anterioare, s-a realizat prin publicarea a 12 articole de specialitate, după cum urmează. **7** publicații ISI: **4** în reviste cotate ISI - Science Citation Index Expanded (lucrările [1], [2], [3], [4]) și **3** la conferințe ISI - Conference Proceedings Citation Index (lucrările [5], [6], [7]) și **4** articole în reviste indexate BDI (lucrările [8]-[11]). De asemenea, membrii echipei au realizat periodic rapoarte tehnice conținând rezultatele obținute în vederea atingerii obiectivelor proiectului. Menționăm faptul că factorul de impact cumulat al publicațiilor în revistele cotate ISI (calculat pe 2010) este 4.985. Ca urmare, criteriul minim de performanță prevăzut pe anul 2011 (**1** articol ISI și **3** articole BDI) a fost îndeplinit.

D. Concluzii

Sintetizăm rezultatele obținute în cadrul proiectului pe anul 2011 ca fiind următoarele: elaborarea unor noi metode bazate pe satisfacerea constrângerilor pentru problema introducerii șabloanelor de proiectare în sisteme informatice existente; dezvoltarea de noi algoritmi de învățare supervizată pentru adaptarea dinamică a structurilor de date dintr-un sistem informatic; dezvoltarea unei metode bazate pe învățare supervizată și agenți pentru predicția comportamentului utilizatorului. Conform celor prezentate anterior, obiectivele planificate pe anul 2011, cât și activitățile aferente acestora au fost realizate în totalitate, și desfășurate conform cu planul de realizare al proiectului. De asemenea, criteriul minim de performanță prevăzut a fost îndeplinit. De asemenea, obiectivele proiectului pe cei trei ani de derulare a acestuia (2009-2011) au fost realizate, criteriile minime de performanță fiind îndeplinite.

Menționăm faptul că factorul de impact cumulat (pe perioada de derulare a proiectului 2009-2011) al publicațiilor în revistele cotate ISI este 6.512.

Bibliografie

- [1] Czibula, G., Cojocar, G.S., Czibula, I.G., *Evaluation Measures For Partitioning Based Aspect Mining Techniques*, International Journal of Computers, Communications and Control, 6(1), 2011, pp. 72-80 (ISI - Science Citation Index Expanded) - IF = 0.650 (pe 2010)
- [2] Bocicor, M. I., Czibula, G., Czibula, I.G., *A Distributed Q-Learning Approach to Fragment Assembly*, SIC Journal, Studies in Informatics and Control, Vol. 20, Issue. 3, 2011, to be published (ISI - Science Citation Index Expanded) (IF=0.671 pe 2010)
- [3] Czibula, G., Czibula, I.G., *Unsupervised Restructuring of Object-Oriented Software Systems using Self-Organizing Feature Maps*, International Journal of Innovative Computing Information and Control, Japan, 2011, to be published (ISI - Science Citation Index Expanded) (IF=1.664 pe 2010)
- [4] G. Czibula, I.G. Czibula, and R.D. Gaceanu. *Intelligent Data Structures Selection Using Neural Networks*. Knowledge and Information Systems, 2011, to be published (ISI - Science Citation Index Expanded) (IF=2.00 pe 2010).
- [5] Czibula, G., Bocicor, M.I., Czibula, I.G., *A Distributed Reinforcement Learning Approach for Solving Optimization Problems*, in Recent Researches in Communications and IT, Proceedings of the 5th International Conference on Communications and Information Technology (CIT '11), Greece, 2011, pp. 25-30 (ISI - Conference Proceedings Citation Index)
- [6] Niculescu, V., *Formal Refinement of BSP Programs with Early Cost Evaluation*, Proceedings of the 10th International Symposium on Parallel and Distributed Computing, ISPDC 2011, to be published (ISI - Conference Proceedings Citation Index)
- [7] Bocicor, M. I., Czibula, G., Czibula, I.G., *A Reinforcement Learning Approach for Solving the Fragment Assembly Problem*, Proceedings of the 13th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing, SYNASC 2011, to be published (ISI - Conference Proceedings Citation Index)
- [8] Czibula G., Bocicor, M.I., Czibula, I.G., *A Reinforcement Learning Model for Solving the Folding Problem*, IJCTA - International Journal of Computer Technology and Applications, Vol. 2, Issue 1, 2011, pp. 171-182 (indexed Index Copernicus)
- [9] Czibula, G., Bocicor, M.I., Czibula, I.G., *An Experiment on Protein Structure Prediction using Reinforcement Learning*, Studia Babes-Bolyai Informatica, LVI (1), 2011, pp. 25-34 (indexed Mathematical Reviews)
- [10] Czibula, I. G., Czibula, G., Bocicor, M.I., *A Software Framework for Solving Combinatorial Optimization Tasks*, Studia Universitatis "Babes-Bolyai", Informatica, Special Issue, LVI(3), pp. 3-8, 2011 (indexed Mathematical Reviews)
- [11] Niculescu, V., *Storage Independence in Data Structures Implementation*, Studia Universitatis "Babes-Bolyai", Informatica, Special Issue, LVI(3), pp. 21-26, 2011 (indexed Mathematical Reviews)
- [12] Gabriela Czibula. *Sisteme inteligente. Instruire automata*, Editura RisoPrint, Cluj-Napoca, 2008
- [13] Wine data set. <http://archive.ics.uci.edu/ml/datasets/Wine>

Cluj-Napoca, 15 noiembrie 2011

Director de proiect,
Prof. dr. CZIBULA GABRIELA