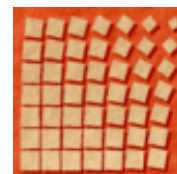


Universitatea "Petru Poni" din Timișoara
Facultatea de Automatică și Calculatoare
Departamentul Calculatoare

2, Vasile Pârvan Bv., 300223 – Timișoara, Romania
Tel: +40 256 403261, Fax: +40 256 403214
Web: <http://www.cs.upt.ro>



CLONEDETECTOR

PROGRAM PENTRU DEPISTAREA

PROIECTELOR JAVA COPIATE

Proiect de diplomă

Adrian-Vlad LEP

Conducător:
Conf.dr.ing. Ioana ȘORA



Timișoara,
2011

Cuprins

I. Introducere	5
1. Domeniul temei	5
2. Problema	5
3. Obiective	6
II. Fundamente teoretice	7
1. Categorii de clone	7
2. Categorii de algoritmi	8
III. JTransformer	12
IV. Descrierea aplicatiei	
1. Arhitectura aplicației	
2. Principiul algoritmul de comparare	
3. Justificarea metricilor alese	
4. Compararea metricilor	
5. Stabilirea probabilitatii de copiere	
V. Implementarea programului	
1. Module	
2. Impartirea pe fisiere	
3. Diferentierea claselor proiectelor de clasele librariilor	
4. Implementarea metricilor	
5. Grafuri de apel	
6. Luarea deciziilor	

7. Compararea a doua proiecte

VI. Manualul utilizatorului

VII. Results

VIII. Concluzii

I. Introducere

1. Domeniul Temei

Înainte de apariția calculatoarelor, copierea proiectelor presupunea efort din partea înfăptuitorului, acesta fiind nevoit să caute în cărți și să scrie manual un proiect chiar și copiat, fără adaos de originalitate. Depistarea proiectelor copiate se realiza exclusiv de către o persoană fizică, prin analizarea tuturor candidaților.

Odată cu apariția calculatorului și transformarea informațiilor în format electronic copierea a fost mult ușurată, determinând creșterea numărului de plagiat-uri. Proiectele software și cele școlare au suferit și ele această problemă, un factor suplimentar fiind și apariția proiectelor open source.

Depistarea clonelor devine astfel tot mai grea cu ochiul liber, fiind nevoie de tool-uri special create pentru aceasta..

2. Problema

Problema încălcării drepturilor de autor a fost serios dezbătută și luată în considerare. În direcția proiectelor software au fost dezvoltate numeroase unelte și încercate diverse abordări pentru depistarea plagiat-urilor. Rezultatele obținute de aceste programe sunt încurajatoare dar mai sunt mulți pași de făcut în această direcție până a putea spune că avem un tool ce poate depista cu o rată de succes acceptabilă proiect software copiate.

Problemele întâmpinate în găsirea unei soluții generice sunt :

- *Cantitatea mare de informații.* Proiectele pot fi copiate între studenții din aceeași an de studiu, de la studenții din anul anterior, de acum doi ani s.a.m.d De asemenea ele pot fi găsite pe internet, diverse site-uri găzduind cod sursă pentru proiecte : sourceforge.net, code.google.com, etc.
- *Clonele sunt greu de depistat.* Compararea codului și depistarea duplicării pare o abordare a problemei, care în contextul unei cantități mari de informații devine o provocare de implementat. Totuși această abordare este insuficientă, depistarea fiind imposibilă în cazul unor modificări (redenumiri de exemplu), modificări ușor de realizat cu uneltele actuale folosite în dezvoltare.
- *Greu de decis dacă sunt clone.* Depistarea unui grup de două sau mai multe proiecte software similare este doar prima etapă. Confirmarea presupunerilor și stabilirea cu certitudine că acestea au fost copiate este o sarcină care se rezolvă prin inspecție umană. Între un proiect copiat și o metodă asemănătoare de rezolvare de multe ori există o graniță foarte greu de stabilit. Dacă cerința determină o abordare consacrată, apariția asemănarilor între proiecte va fi evidentă și deci depistarea a ceea ce este copiat sau nu devine foarte greu. De exemplu pentru implementarea căutării binare metoda este binecunoscută iar rezolvarea problemei nu va diferi foarte mult de la un proiect la altul

3. *Objective*

Deoarece metoda tradițională de compare manuală a proiectelor nu mai este fezabilă datorită cantității de informații ce trebuie analizate, pentru depistarea clonelor au fost elaborate, de diverse universități în special, programe pentru analizarea automată a proiectelor. Acestea sunt de un real ajutor pentru profesori și nu numai. Aceste instrumente se diferențiază între ele prin :

- *tipul de date analizat* : proiecte literare, proiecte software : Java, C, C# etc.
- *aria de căutare a duplicatelor* : proiectele dintr-o locație /dintr-un director, proiectele de pe un calculator sau căutarea pe internet a posibilei perechi.
- *tipul de clone depistate* : există diferite tipuri de clone pe care fiecare program le poate depista
- *dimensiunea maximă a datelor ce o pot analiza* : proiecte școlare, proiecte (software sau nu) de dimensiuni mai mari, internetul.
- *algoritmul de depistare* : categorii de algoritmi vor fi descriși în capitolul următor.

Platforma CloneDetector analizează static codul sursă și depistează proiectele Java copiate pe baza unui set de reguli Prolog. CloneDetector folosește o abordare originală, combinând metodele de depistare bazate pe metrici și cele bazate pe grafuri de apeluri (call graph) create din arbori sintactici abstracți(ĂST). Prin acest instrument se dorește depistarea clonelor de tipul 3, adică a proiectelor copiate în care pe lângă redenumire s-au operat schimbări de natură sintactică, s-au adăugat sau șters porțiuni de cod. Această platformă s-a dovedit un succes pe proiectele școlare testate, având potențialul de a deveni un instrument larg folosit în depistarea copiilor.

Lucrarea de față descrie în detaliu algoritmul de detecție folosit pentru depistarea clonelor în secțiunea IV. În secțiunea următoare vom introduce noțiunile teoretice folosite în elaborarea lucrării actuale. Secțiunea III vă prezentată modul de obținere a ĂST-ului din codul sursă Java. Secțiunea V prezintă detalii de implementare a algoritmului și a metricilor folosite. Manualul de utilizare se găsește în secțiunea VI iar rezultatele experimentale sunt descrise în secțiune VII.

II. Fundamente teoretice

1. Categorii de clone

În această secțiune ne propunem să definim și să clasificăm clonele sau copiile software pe baza cercetărilor efectuate în acest domeniu.

Conform [11] plagiatul este reproducerea muncii altei persoane fără a recunoste proveniență informației. Cele mai întâlnite cazuri apar în mediul academic unde studenții copiază materiale din cărți, de pe internet sau de la colegi fără să citeze sursa lor de proveniență. Plagiatul este o practică și în alte domenii, un exemplu celebru fiind acela dintre Steven Spielberg și Chase-Riboud care îl acuză pe primul că personaje, scene și alte aspecte din cartea sa “Echo of Lion” au fost copiate în filmul Amistad.

Un alt exemplu, este procesul pornit între IBCOS Computers Ltd și Barclays Mercantile Highland Finance Ltd. pentru copierea ilegală de software. Cazul era legat de un programator care s-a mutat de la o firmă la alta și a refolosit codul dezvoltat la prima în produsele noii sale companii. Instanța a decis că drepturile de copyright au fost încălcate.

Copierea proiectelor software presupune reutilizarea codului scris de altă persoană, în întregime sau parțial, încălcând astfel drepturile de autor. Conform [1] clonele software se împart în trei categorii în funcție de natura lor :

Tipul 1 sunt clone identice, fără modificări (doar eventuale comentarii sau spații modificate).

Tipul 2 sunt copii identice din punct de vedere sintactic. Sunt posibile redenumiri de variabile, tipuri sau identificatori de funcții.

Tipul 3 sunt clone care au suferit mai multe schimbări; prin modificarea, ștergerea sau adăugarea unor declarații, instrucțiuni.

După cum se observă tipul 1 și 2 de clone sunt clar și precis definite. Tipul 3 este mai vag definit reprezentând o categorie mai mare de posibile clone. Din acest motiv unele unelte consideră unitar clonele de tipul 1 și 2, în timp ce clonele de tipul 3 sunt tratate separat, datorită diferențelor de abordare.

Clonele de tipul 1 și 2 pot fi declarate automat că fiind copiate, dar pentru clonele 3 nu se pot formula metrici care să stabilească cu certitudine dacă sunt sau nu copiate. Astfel clonele de tipul 3 declarate suspecte trebuie analizate de către o persoană înainte de a fi pronunțate adevărate.

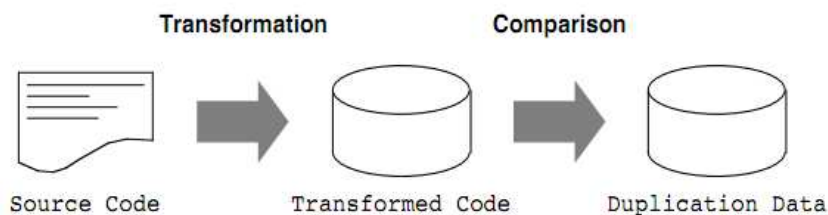
Există și alte clasificări ale clonelor oferite de Balazinska în [12] sau de Kasper în [13], [14] , acestea fiind considerate cele mai detaliate clasificări actuale. Totuși în această lucrare am adoptat clasificarea prezentată mai sus, fiind simplă, ușor de înțeles și suficient de sugestivă.

În urmă analizei proiectelor, instrumentul creează o listă de posibili suspecti. Aceasta poate fi completă, în cazul în care toate clonele adevărate apar în listă sau incompletă dacă unele lipsesc. Din această raport unele clone vor fi confirmate în timp ce altele se vor dovedi presupuneri false.

Un program eficient, se dorește să raporteze toate sau cât mai multe clone ce sunt adevărate (să nu scape nimic) și în același timp să raporteze cât mai puține clone false. Un instrument ce raportează prea multe clone false este nefolositor, sporind cantitatea de informație ce trebuie analizată.

2. Categoriile de Algoritmi

Pentru depistarea proiectelor software copiate au fost abordate mai multe tehnici. Aceste se diferențiază prin tehnică de comparare(algortmul) și nivelul la care este analizată informația. Figura 2.1 este sugestivă, prezentând câteva tehnici și caracteristicile lor.



Author	Level	Transformed Code	Comparison Technique
Johnson 94	Lexical	Substrings	String-Matching
Ducasse 99	Lexical	Normalized Strings	String-Matching
Baker 95	Syntactical	Parameterized Strings	String-Matching
Mayrand 96	Syntactical	Metrics Tuples	Discrete comparison
Kontogiannis 97	Syntactical	Metrics Tuples	Euclidean distance
Baxter 98	Syntactical	AST	Tree-Matching

Fig.2.1 [15] Tipuri de algoritmi

Vom prezenta succint câteva dintre ele oferind exemple de programe ce le implementează.

▪ Comparare de șiruri

Presupune interpretarea liniilor de cod ca niște șiruri de caractere și compararea lor pe baza diversilor algoritmilor. Ea se împarte la rândul ei în două categorii :

Comparare textuală : Compară linii întregi de cod. Pentru a crește performanța sunt calculate funcții hash ale liniilor. Această abordare a fost introdusă de Ducasse în [2]

Comparare simbolică : Compararea este la nivel de simboluri, o linie de cod având de obicei mai multe simboluri. În locul unei analize a șirurilor de caractere se crează un arbore de sufixuri și se compară pe baza de functori. Propusă de Baker în '95 [3].

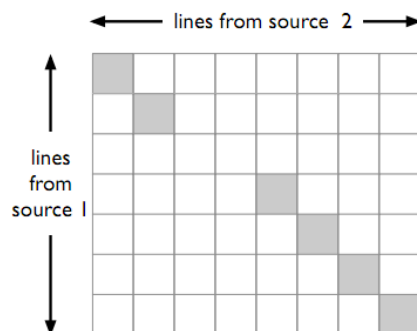


Fig. 2.2 [15] Comparare textuala

Practic în programele ce adoptă acesta variantă vor compara linii de cod dintr-un fișier al unui proiect cu cele din alt fișier dintr-un proiect diferit. Dacă găsesc o portine exact la fel ea este memorată ca fiind copiată, în figura de mai sus, Fig.2.2, asta realizându-se prin colorarea gri a pătratului corespunzător. Dacă proiectele ar fi identice în totalitate ar trebui să obținem diagonală principală a matricei în totalitate gri, iar restul alb. Cum acest caz este mai rar întâlnit matricea poate arată ca în figura 2.3. Se observa că există porțiuni din primul fișier care nu se regăsesc în al doilea, primul spațiu lipsă pe diagonală principală. Se poate datora modificării codului. De asemenea există porțiuni care se vor găsi decalat în al doilea fișier, cum este cazul celui de-al doilea spațiu de pe diagonală principală. Acesta decalare poate apărea datorită adăugării unei noi porțiuni de cod.

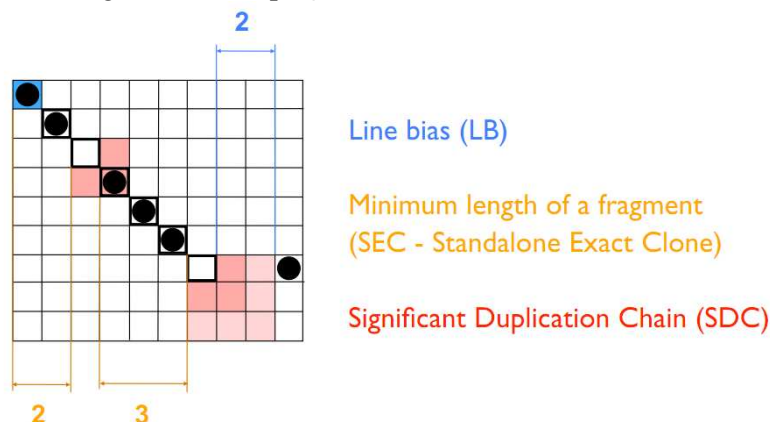


Fig. 2.3 [15] Interpretarea rezultatelor

Pentru a declara o porțiune de cod suspectă de copiat, întregul fișier sau proiect, trebuie să existe o porțiune de cod copiat de lungime minimă și cu un decalaj maxim între bucăți.

Pe baza acestor metode au fost dezvoltate numeroase platforme pentru depistarea copiilor software dar și pentru depistarea celor literare sau a duplicării de cod. Câteva exemple ar fi :

- *MOSS (Measure Of Software Similarity)* Este unul din cele mai cunoscute unelte în domeniu, dezvoltat de Universitatea Stanford în 1994. Este folosit și astăzi mai ales pentru depistarea temelor software copiate. Programul poate analiza codul sursă pentru următoarele limbaje de programare : C, C++, Java, C#, Python, Visual Basic, Javascript, FORTRAN, ML, Haskell, Lisp, Scheme, Pascal, Modula2, Ada, Perl, TCL, Matlab, VHDL, Verilog, Spice, MIPS assembly, a8086 assembly, a8086 assembly, MIPS assembly, HCL2. Este disponibil ca serviciu web pe <http://theory.stanford.edu/~aiken/moss/>.

- *JPlag* Este o platformă dezvoltată de Universitatea din Karlsruhe. Se bazează pe algoritmul de compărare de șiruri cunoscut sub numele :Karp-Rabin Greedy String Tiling. Oferă suport pentru limbajele Java, C#, C și C++. Este disponibil la adresa dată ca serviciu web <https://www.ipd.uni-karlsruhe.de/jplag/>

- *YAP* : Dezvoltat de universitatea din Sydney, are trei versiuni schimbând în fiecare dintre ele algoritmul de comparare. Acum folosesc același algoritm ca și JPlag, trecând de la vechiul algoritm Heckel folosit în versiunea 2. Este folosită ca serviciu web <http://luggage.bcs.uwa.edu.au/~michaelw/YAP.html>

După cum observăm principalii contribuitori în acest domeniu au fost universitățile, ele încercând să găsească soluții la problema copiatului. Majoritatea au adoptat varianta serviciilor web prin care orice profesor se poate înscrie și trimite un set de teme, urmând să primească rezultatul. Această metodă este sigură, studenții fiind împiedicați să încerce să păcălească

unealta. Totuși pentru noi a fost dificil să obținem acces pentru a compara metoda noastră cu aceste abordări și am ales ca tool de comparare din aceasta categorie DUDE, dezvoltat de Radu Marinescu [10].

▪ Compararea pe baza de metrici

Presupune obținerea unor metrici ce caracterizează o bucată de cod. Se vor compara vectorii de metrici obținuți pentru diverse porțiuni și vor fi considerate egale cele ce au sub o anumită distanță admisă. Diferite metrici au fost propuse de către : [3],[4] și [5],[17],[18].

- Halstead :
 - Numărul de operatori distincti
 - Numărul de operanzi distincti
 - Numărul total de referiri(folosiri) ale operatorilor
 - Numărul total de referiri ale operanzilor.
- McCabe
 - Complexitatea ciclomatică
- Dunsmore
 - Adâncimea ierarhiilor de clase dintr-un modul
- Kolmogorov complexity

Exemple de sisteme ce implementează acest tip de algoritmi :

- SID descris în detaliu în [19] folosește metrica Kolmogorov.
- Platforma de depistat programe copiate dezvoltată de Ottenstein [33] de la Universitatea Purdue, bazată pe metricile lui Halstead

▪ Compararea arborelui sintactic abstract (AST)

Această abordare se bazează pe obținerea arborelui sintactic abstract al programului, în prima fază, iar pe baza acestuia analizându-se similaritatea codului. Programul e partitionat în bucăți de AST care sunt comparați prin diverși algoritmi.

Analiza AST-ului obținut poate fi realizată în mai multe feluri. O variantă este *compararea grafului de apeluri a programului*. Dependențele unei funcții față de date sau alte funcții pot fi reprezentate printr-un graf. Clonele ar trebui să aibă subarbori identici. Diferiți algoritmi au fost propuși și implementați în [6], [7] și [8].

Code Digger este un exemplu de program dezvoltat bazându-se pe această abordare de către Peter Bulychev și Marius Minea. Compararea arborilor de bazează pe algoritmul anti-unification.

Această categorie de algoritmi sunt mult mai exacti ei fiind capabili să depisteze și bucăți de cod copiate care au suferit modificări mai avansate. Pe de altă parte, acești algoritmi fiind mai complexi au un timp de rulare mai mare, fiind greu de folosit pe proiecte software mai mari.

▪ Alți algoritmi

Pe lângă aceste tipuri clasice de abordări au apărut programe ce folosesc algoritmi de căutare pe internet, de explorarea datelor și alții.

- Siff compara fișiere mari similare într-un sistem. Se bazează pe aproximarea fișierului printr-un fingerprint, un sir de caractere scurt specific fiecărui fișier. Se calculează checksumul fiecărui fingerprint, iar acestea vor fi apoi rapid comparate între ele.
- SCSSB (System Call Short Sequence Birthmark) și IDSCSB (Input Dependant System

Call Subsequence Birthmark) descris în [16] propun o metoda de detectare a duplicarii de cod în proiecte mari software bazata pe apelerile de sistem. Aceasta abordare nu se preteaza sistemelor software mici, dar la sisteme mari succesiunea de apeluri catre kernel-ul sistemului de operare sunt o metrica ce caracterizeaza sistemul sau portiuni ale lui în întregime.

<pre> 1. S₀; 2. If (i==1) { 3. S₁; 4. for (j=0;j<3;j++) 5. { 6. S₂; 7. S₃; 8. S₄; 9. } 10. } else { 11. S₆; 12. }</pre>	<pre> S₀ S₁ S₂ S₃ S₄ S₂ S₃ S₄ S₂ S₃ S₄ S₅</pre>
(a)	(b)
	<pre> S = { S₀S₁S₂S₃, S₁S₂S₃S₄, S₂S₃S₄S₂, S₃S₄S₂S₃, S₄S₂S₃S₄, S₂S₃S₄S₅}</pre>
	(c)

Fig.2.4 a) Exemplu de cod b) Apelurile catre sistem generate pentru intrarea i=1 c) Lista de apeluri generate pentru i=4

În figura de mai sus avem un scurt exemplu pentru a ilustra această tehnică. Se observă astfel lista de apeluri ce se va obține executând codul dat pentru diferitele valori de intrarea ale lui i. Prin compararea acestor liste se pot depista porțiuni de cod clonate. Această comparare este una dinamică, deoarece pentru a obține lista apelurilor către sistem proiectele analizate vor fi executate.

Întrucât plagiatul proiectelor software și duplicarea de cod sunt similare ca metode de detecție, diferență fiind făcută de baza de informații ce e comaprata, algoritmi descriși pot fi folosiți pentru a atinca oricare din cele două scopuri.

Observăm existența numeroaselor abordări pentru depistarea clonelor software. Eficientă lor fiind analizată în funcție de timpul de execuție și de acuratețea cu care depistează anumite tipuri de clone. În urma analizării unui set de platforme cunoscute s-a realizat că fiecare are avantaje pentru unele categorii de clone. Din punct de vedere al rapidității, uneltele ce implemeteaza compararei bazate pe șiruri de caractere sau simboluri sunt mult mai rapide decât cele ce analizează arbori sintactici dar acuratețea lor fiind mai redusă. O descriere mai detaliată a avantajelor și dezavantajelor fiecărei platforme se poate găsi în publicațiile [22] și [23] care revizuiesc o mare parte a uneltelor prezentate.

III. JTransformer

Pentru analizarea proiectelor suspecte de copiat suntem nevoiți să transformăm codul sursă, scris în Java, într-un arbore sintactic, reprezentabil în Prolog. Această transformare a fost realizată folosind JTransformer 3.0.1.

JTransformer este un motor de interogare și de transformare a codului sursă Java, disponibil ca plugin pentru Eclipse. Cu acesta am creat Arborele Sintactic Abstract (AST) pentru un proiect Java, luând în considerare codul sursă complet, tipurile de fișiere și membrii clasei. AST-ul Java este reprezentat intern ca o bază de date Prolog, facilitând analize puternice și consturirea unor reprezentări abstracte ale proiectului într-un mod ușor, în doar câteva linii de cod.

JTransformer are un nivel detaliat de interpretare a codului Java, luând în considerare atribute, metode, instrucțiuni, chiar și comentarii pentru clasele din proiect. Pentru clasele din librării nu sunt luate în considerare detaliile de interpretare. Se reprezintă în baza de date prolog doar interfața claselor.

Pentru fiecare nod din AST se creează o axiomă(faptă). Aceste axiome sunt legate între ele prin ID-uri unice. De exemplu, pentru o clasa se creează o axiomă classT. Lista complete de axiome posibile se găsește la [9].

În continuare vom lua în considerare un exemplu pentru a vedea modul în care este transformat codul Java într-o bază de cunoștințe Prolog.

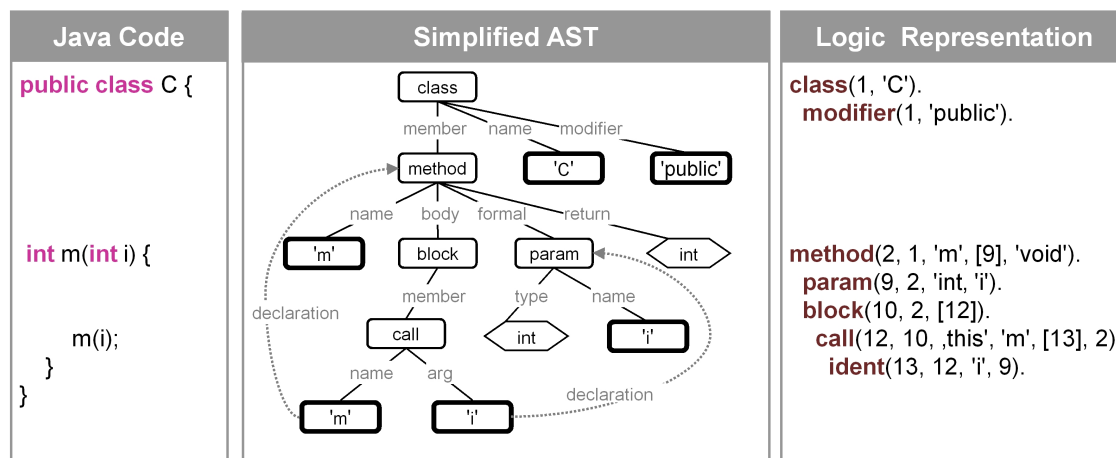


Fig. 3.1 [9] Exemplu de transformare

În imaginea de mai sus avem în partea stângă codul sursă Java, iar în partea dreaptă reprezentarea Prolog creată. Se observă că pentru definirea clasei se creează două intrări în baza de date, creindu-se o intrare separată pentru modificatorul claselor. De asemenea metoda Java e reprezentată în prolog prin axiome de tipul method, param – care conțin parametrii metodei și block – corpul metodei. Block va fi părinte pentru celelalte instrucțiuni din interiorul metodei. Ierahia stabilită între axiome, prin identificatori, este prezentată în schema din mijlocul imaginii. Transformarea prezentată în Figura 3.1 este simplificată comparativ cu cea generată de JTransformer. Inea nu sunt incluse identificatori care fac legătura din spre părinți sper copii.

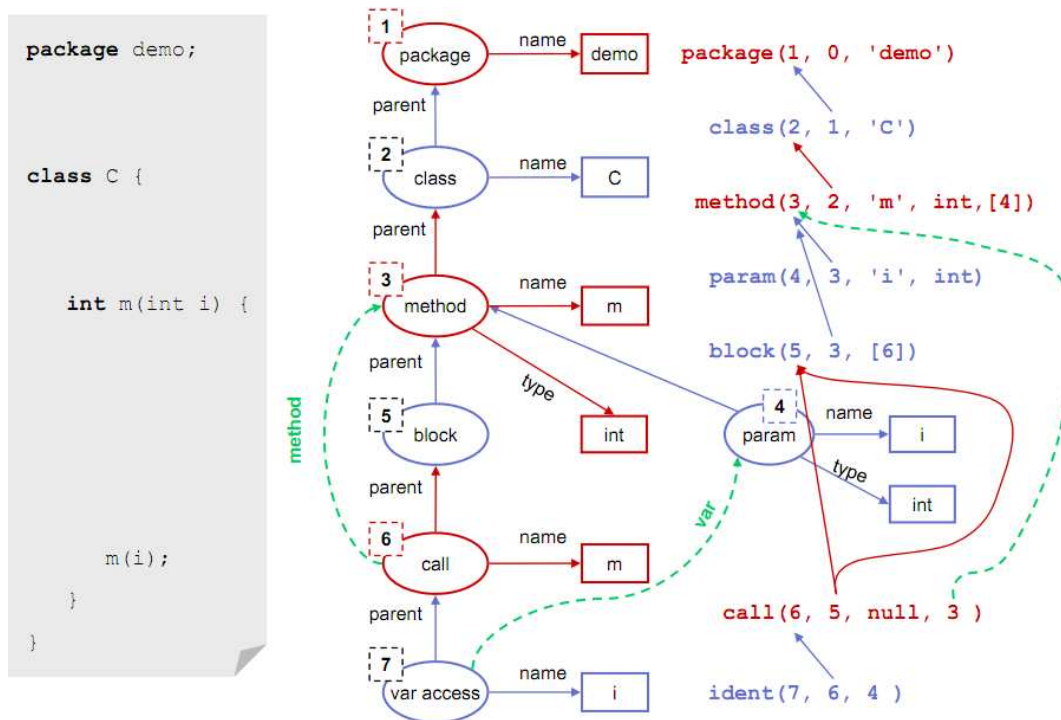


Fig. 3.2.[9] Relatiile între clauze

Figura 3.2. prezintă relațiile din interiorul bazei de date Prolog, modul în care un identificator leagă două clauze între ele. O clauză poate fi legată de mai multe clauze în mod direct sau prin redirectări.

În reprezentarea reală Prolog generată de plugin vor fi mai multe detalii. În continuare sunt exemplificate câteva clauze :

```

classT(#id, #parent, 'name', [#def_1,...])
methodT(#id, #classT, 'name', [#param_1,...], TYPE, [#exception_1,...], #body)
paramT(#id, #parent, TYPE, 'name')
blockT(#id, #parent, #enclMethod, [#statement_1,...])
callT(#id, #parent, #encl, #expr, 'name', [#arg_1,...], #method)
identT(#id, #parent, #encl, 'name', #symbol)

```

Întotdeauna primul număr este identificatorul unic. Al doilea parametru este identificatorul părintelui, în cazul metodei va fi id-ul clasei de exemplu. Aceștia pot fi urmați de identificatori ce relaționează clauzele sau de parametrii caracteristici axiomei curente. Parametrii caracteristici fiecărei axiome sunt numele, de exemplu. După cum putem observa un termen al unei clauze poate fi o listă, cum avem la metoda lista de parametrii de exemplu. Termenul #encl trimite către metoda din care face, deci clauzele pot referii pe lângă părinții direcți și părinții părinților ei, mergând până la nivel de clase. Mai multe detalii despre modul de reprezentare al informații vor fi prezentate și în capitole ulterioare în momentul în care vom fi nevoiți.

JTransformer a fost alegerea potrivită pentru noi, datorită capacității sale de a genera o reprezentare Prolog detaliată din codul sursă Java. De asemenea, platforma este ușor de folosit datorită integrării cu eclipse, cel mai popular mediu de dezvoltare în Java. Viteza de transformare este impresionantă, fiind capabilă să transforme proiecte de dimensiuni mare, chiar de 1,000,000 LOC în câteva secunde..

IV. Descrierea aplicatiei

În secțiunea următoare vom descrie platforma propusă de noi, discutând arhitectura aplicației, algoritmul de detecție folosit și interfața cu utilizatorul.

CloneDetector a fost dezvoltat pentru a depista proiecte școlare Java. El a fost implementat în Prolog, un limbaj folosit cu precădere în domeniu inteligenței artificiale. Paradigma de programare din care face parte și Prolog face ușor de reprezent probleme din viața reală prin fapte și reguli de legătură. Sunt diferențe mari de implementare între Prolog și limbaje procedurale sau orientate pe obiecte, datorate în mare parte mașinii de inferență. Aceasta implementează algoritmul backtracking încercând să genereze toate soluțiile posibile pentru o anumită regulă. Pentru noi această abordare a fost un avantaj putând reprezenta ușor și intuitiv regulile algoritmului și beneficiind apoi de mașină de inferență care caută soluțiile pentru regulile descrise.

Compilerul folosit a fost SWI-Prolog acesta fiind probabil cel mai complet de pe piață, oferind un mediu de depanare grafic foarte folositor și o gama largă de librării. SWI-Prolog suportă execuția multi-threading o facilitare importantă care a îmbunătățit timpul de execuție al aplicației. Este o platformă bine menținută ce suportă toate marile tipuri de sisteme de operare (Windows, Linux, Mac OSX) pe arhitecturi de 32 de biți cât și pe 64.

1. Arhitectura aplicației

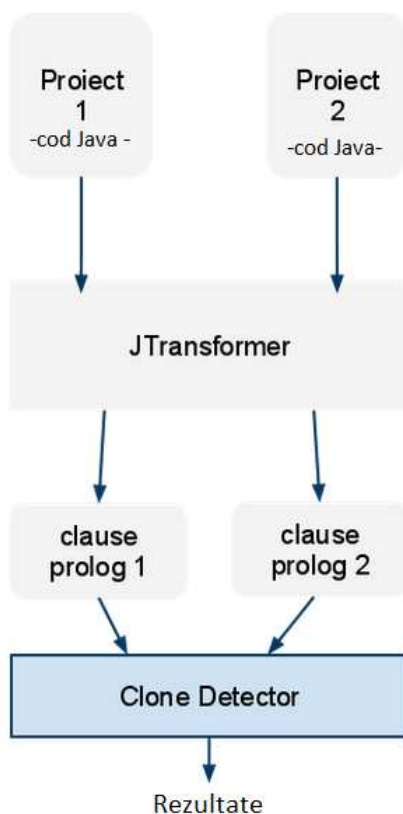


Fig. 4.1 Arhitectura sistemului

Ținând cont că sursele proiectelor analizate sunt scrise în Java, iar noi dorim o reprezentare Prolog a acestora am folosit plugin-ul eclipse JTransformer pentru a transforma acest cod într-o reprezentare Prolog. Codul va fi transformat în fapte care prin relațiile dintre ele alcătuiesc un arbore sintactic abstract. Ca ieșire, JTransformer generează un fișier pentru fiecare proiect cu reprezentare proiectului. Pentru a îmbunătății performanțele am ales ca acesta să genereze un fișier gata compilat, cu extensia qlf (quick load file). Încărcarea acestor fișiere fiind de aproximativ 50 de ori mai rapidă.

În Fig.4.1 avem o descriere generală a sistemului. Pentru compararea a două proiecte, din codul Java se va genera cu ajutorul plugin-ului Eclipse JTransformer arborele sintactic abstract (AST) al fiecărui proiect. AST-ul va fi format din clauze Prolog și în fișiere compilate Prolog.

CloneDetector va încărca în baza de cunoștințe aceste fișiere și pe baza algoritmului descris în continuare va analiza dacă proiectele au fost copiate.

Întrucât compararea manuală a proiectelor două câte două ar fi fost inefficientă, consumând mult timp din partea profesorului, s-a încercat automatizarea acestei sarcini. Programul suportă posibilitatea de a primi ca parametru un director, comparând toate proiectele, două câte două, găsite în acesta.

2. Principiul algoritmul de comparare

Algoritmul de comparare propus de acesta lucrare este unul original combinând abordarea bazată pe metrici cu cea bazată pe grafuri de apeluri. Metricile descriu porțiuni de cod, iar combinarea lor pot diferenția unic aceste porțiuni, ca un fel de amprentă a codului. Grafurile de apeluri sunt obținute din AST, luând în calcul apelurile realizate de o metodă către alte metode. Obiectivul nostru era de a obține un tool capabil să depisteze clone de tipul 1, 2 dar mai ales cele de tipul 3, imune la redenumiri, cu un timp de procesare acceptabil. Cele două abordări au avantaje și dezavantaje iar noi am încercat să luăm avantajele fiecăruia și să minimizăm dezavantajele. Metricile sunt mai rapide de calculat, dar sunt inexacte, fiind greu să descrii o porțiune de cod doar prin metrici, în timp ce arborii de apeluri sunt exacti dar costisitori din punct de vedere computațional.

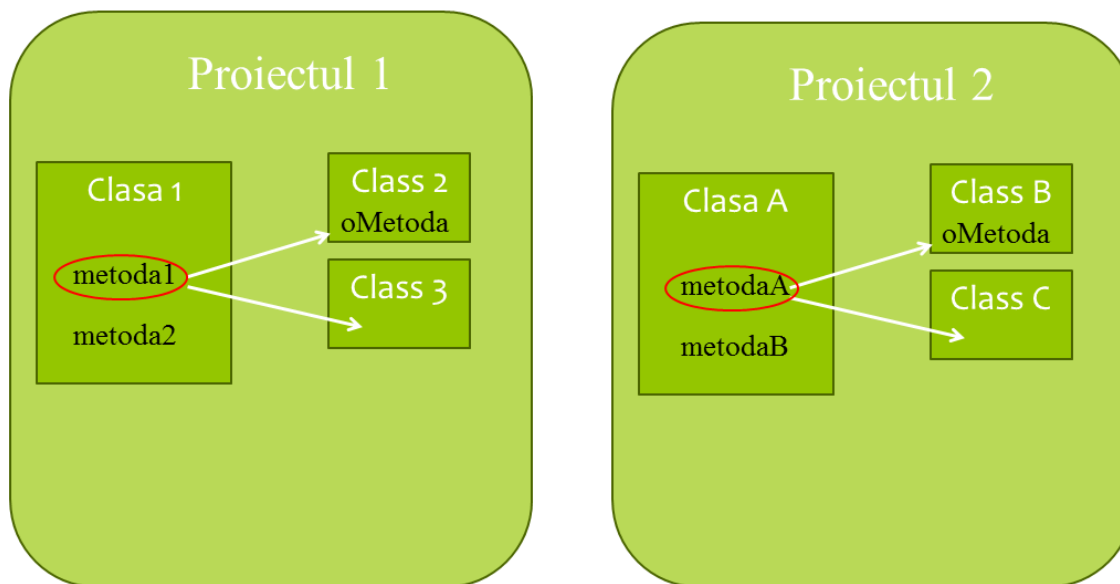


Fig. 4.2 Descrierea algoritmului

În continuare vom explica pașii algoritmului folosindu-ne de figura 4.2 pentru exemplificare. Pentru compararea a două proiecte instrumentul începe să genereze toate combinațiile posibile dintre clasele primului proiect și cel de al doilea. Analiza dorește să stabilească la final care clasă este identică cu care, memorând perechile de clase copiate și probabilitatea lor de a fi copiate.

Două clase vor fi considerate identice dacă vor avea similarități la nivel de clasă, dacă metodele lor vor fi similare și dacă clasele care vor fi apelate din metodele celor două clase vor fi asemănătoare. Deci clasele 1 respectiv A din figura vor fi considerate identice dacă în prima etapă metricile la nivel de clasă sunt similare. Urmează compararea la nivel de metodă unde încercăm să împerechem metodele similare din cele 2 clase. Metoda 1 va fi similară cu metoda 2 dacă metricile la nivelul celor două metode sunt apropiate, după care se analizează la nivel de apeluri. Cele două metode vor fi similare dacă metoda oMetoda din Class2 va fi similară cu metoda oMetoda din ClassB și dacă Class2 va fi similară cu ClassB. La al doilea nivel, compararea claselor și metodelor apelate(Class2, ClassB, oMetoda), se analizează doar metricile clasei și a metodei.

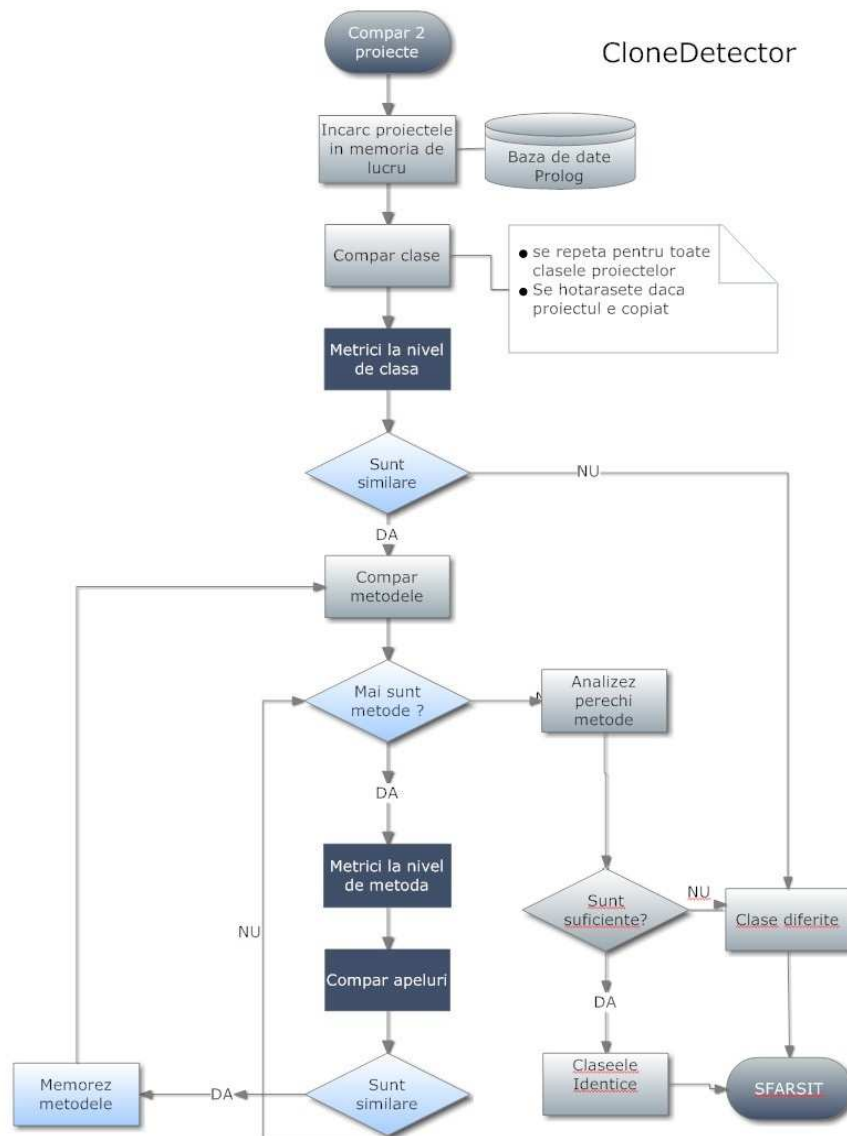


Fig. 4.3 Reprezentare schematica algoritm

O abordare recursivă la nivel de apeluri (luare în calcul și a apelurilor pentru metodele apelaate) ar conduce la o aprofundare a analizei și deci creșterii duratei de timp. S-ar genera de asemenea numeroase cazuri excepționale, apeluri ciclice, problema moștenirii etc, pentru tratarea cărora ar trebui complicat algoritmul. O astfel de abordare ar atrage atât avantajele dar mai ales dezavantajul algoritmilor bazați pe grafuri de apeluri, adică timpul lung de procesare. Analiza doar a primului nivel a grafului a fost suficientă. Combinând-o cu metricile ea obținând rezultate bune.

După analiza primului apel a metodei 1 din cele două clase se vor analiza și restul apelurilor către clasele 3 respectiv C, după care se va trece la următoarea metodă.

Practic, analiza claselor trece printr-un set de filtre. Trecerea de un filtru presupune o probabilitate mai mare ca cele două clase să fie copiate. Trecerea filtrelor la nivel de clasă determină asignarea acestora ca fiind suspecte puțin probabile de copiat. Dacă filtrele la nivel de metodă sunt satisfăcute, clasele devin suspecte cu probabilitate medie de a fi copiate, iar trecerea de filtrele comportamentale vor determina clasele să fie considerate suspecte foarte probabile.

O reprezentare schematică a ceea ce am exemplificat mai sus este realizată în figura 4.3. Metricile considerate reprezentative și alese pentru compararea la nivel de clasă și metodă sunt enumerate în continuare, implementarea lor fiind arată în capitolul următor.

A.Metrici la nivel de clasă :

1. numărul de atribute
2. numărul de metode
3. numărul de interfete implementate
4. moștenește sau nu o superclasa.
5. numărul de metode similare

B.Metrici la nivel de metodă :

1. numărul de cicluri while, do while, for
2. numărul de condiții if, if else
3. numărul de operatori -, *, /
4. modificatorul static
5. semnătura metodei

3. Justificarea metricilor alese

Înainte de a ne hotărâ asupra acestor metrici am testat pe un set de proiecte dacă acestea sunt suficiente pentru a descrie unic clase sau metode. Doream ca prin această cercetare să vedem dacă proiectele diferite au metrici diferite, dacă cele similare au metrici asemănătoare și care este diferență între metricile proiectelor copiate sau nu. Imaginea 4.4 prezintă o bucată din studiu, acesta fiind disponibil și online pe docs.google.

1	Nume proiect	Nume clasă	nr atribut	nr metode	nr interf	nr superclass	Nume metoda	nr op (+, -, *, /)	cicli while/for	if	nr param	ret val
2	passc - xml	XMLReader	6	7	1	1	startDocument	(0,0,0,0)	0		0	void
3							endDocument	(0,0,0,0)	0		0	void
4							startElement	(5,0,0,0)	0	1	4	void
5							endElement	(1,0,0,0)	0	0	3	void
6							characters	(1,0,0,0)	0	1	3	void
7							storeList	(0,0,0,0)	0	0	1	void
8							loadList	(3,0,0,0)	1	1	0	List
9		Student	3	6	0	0	Student/constructor	(0,0,0,0)	0	0	3	
10							getNrMaticol	(0,0,0,0)	0	0	0	int
11							getMedie	(0,0,0,0)	0	0	0	double
12							toString	(3,0,0,0)	0	0	0	String
13							crestMedie	(2,0,0,0)	1	1	0	bool
14							valZBuna	(1,0,0,0)	0	1	0	bool
15		Main	0	1	0	0	main	(6,0,0,0)	0	0	1	void
16		Horoscop	2	7	0	0	constructor	(0,0,0,0)	0	0	1	
17							changeSource	(0,0,0,0)	0	0	1	void
18							loadList	(0,0,0,0)	0	0	0	void
19							printList	(0,0,0,0)	0	0	0	void
20							findByNr_matr	(0,0,0,0)	1	0	1	student
21							findByMedie(double)	(0,0,0,0)	1	0	1	List<Student>
22							catVoriAveallMediaScaz	(0,0,0,0)	1	0	0	int
23	<<interfata>>	DataAccessObject	0	2	0	0	loadList					List
24							storeList					void
25												
26	webservice	ServerGUI	4	2	0	1	cosnstructor					
27		BrowseListener	6	8	0	0						
28												
29	Passc 1	Boala	3	5	0	0						
30		FisaPacient	4	6	0	0						
31		Main	0	1	0	0						
32		Medic	3	5	0	0						
33		Pacient	1	2	0	0						
34		Simptom	1	3	0	0						
35		TraseuPacient	2	4	0	1						

Fig. 4.4 Studiu realizat pe un set de proiecte

Studiu a relevat că încă de la nivel de clasă există puține coliziuni, adică puține clase care au aceleași metrici, acestea fiind de obicei cele de pe același nivel al unei ierarhii de moștenire. Totuși comparând la nivel de metoda clasele vor fi diferite total pentru proiecte diferite. Și atunci vine întrebarea de ce nu am adoptat varianta doar cu metrici fără graf de apeluri ? Proiectele ar fi copiate dacă metricile ar fi identice. Problema ce ne-a determinat să combinăm abordările este dată de comportamentul studenților care copiază. Aceștia au tendința să redenumescă, lucru care nu ne afecta, dar câteodată mai schimbă sau adaugă cod creând clone de tipul 3. Acestea nu ar putea fi depistate dacă metricile ar fi comparate identic, de aceea noi vom compara asemănarea metricilor nu egalitatea lor, acceptând unele diferențe.

Dacă ne uităm în tabelul din figura 4.4 observăm că acceptând aceste diferențe la nivel de clasă și metodă vom obține coliziuni, astfel analiza apelurilor fiind justificată.

Pe lângă aceste metrici, enumerate mai sus, care sunt calculate direct din reprezentarea AST a proiectelor, sunt setate filtre care iau în considerare, la nivel de clasă, numărul metodelor ce au fost grupate și nivelul lor de asemenare. Pe baza acestor filtre asociem nivelul de asemenarea între doua clase.

4. Compararea metricilor

Compararea valorilor obținute pentru metrice nu se face exact, după cum am menționat anterior, adică numărul de if-uri dintr-o metodă nu va trebui să fie exact același și în cealaltă. Se acceptă aceste diferențe tocmai pentru a permite prinderea plagiaturilor schimbate prin adăugarea de noi condiții sau alte modificări. Am adoptat două variante de a defini aceste toleranțe, procentual sau prin diferență absolută.

Procentual presupune calcularea procentului dintre două metrice, de exemplu 80% între numărul metodelor din două clase reprezintă că prima clasă are 80% din numărul metodei celei de a doua, fără a ști dacă sunt sau nu egale. În acest caz am putea stabili că acest procent, 80% să fie minimul acceptat pentru a continua analiza a două clase. Astfel dacă o clasă are doar 50% ca număr de metode față de alta, nu se va încerca să se grupeze cele două. Această abordare procentuală este avantajoasă pentru proiecte, clase sau metode mai mari. Pentru clase mici procentele au variații mari. Între o clasă cu 1 metodă și alta cu 2 metode avem 50% procentul, la fel ca între una cu 10 metode și alta cu 20. În primul caz sunt șanse ca acele clase să fie copiate, în vreme ce în cel de-al doilea putem spune aproape cu certitudine ca cele două clase sunt total diferite.

Diferența absolută între metrice. În acest caz calculăm diferența absolută între metrice și impunem ca aceasta să fie mai mică ca o valoare maximă tolerată. Astfel dacă avem 2 metode respectiv 7 în alta clasă diferența (delta) între ele va fi 5. O astfel de diferență o considerăm prea mare pentru a continua analiza claselor, dacă valoarea maximă admisă ar fi 3. Această metodă de a exprima toleranțele este bună pentru proiecte mici, în care valorile metricilor variază puțin, dar pentru proiecte mari ar putea fi problematică. Dacă avem 20 și 24 de metode în două clase ar trebui să le comparăm, cele 4 metode putând fi obținute prin spargerea unor metode mari.

CloneDetector folosește cu precădere cea de-a doua metodă de a defini toleranțele între două metrice, datorită proiectelor mici analizate, dar în anumite cazuri am considerat utilă și prima metodă, cea procentual, înlănțuind condiții procentuale sau diferențiale prin și/sau logic.

Alegerea acestor diferențe maxime/procente minime a fost externalizată în fișierul `profiles.pl`. Am creat aici două profile unul `loose` și altul `tight`, cu toleranțe diferite. Se poate alege între cele două în funcție de grupul de proiecte analizate. Profilul `loose` este mai permisiv, permițând variații mai mari între metrice. Este ideal pentru proiecte a căror rezolvări pot diferi mai mult, nefiind sugerate anumite căi de la început. Profilul `tight` este folositor pentru teme a căror rezolvare este clasică, cunoscută sau a fost sugerată. Dacă am rula pe același set de teme profilul `loose` ar genera mai mulți suspecti, în timp ce profilul `tight` ar genera mai puțini dar cu șanse mai mari de a fi copiate.

5. Stabilirea probabilitatii de copiere

Pe baza tolerantelor setate în fișierul de configurare, se vor lua în continuare decizii legate de nivelul de similitudine între anumite componente (metode, clase) acestea putând fi : diferite sau cu anumite grade de asemănare.

Metodele pot avea două grade de similitudine : *High* sau *Low*.

Metodele grupate cu grad de asemănare *Low* sunt metode care sunt considerate suspecte de copiat, dar cu o probabilitate mai mică. Ele sunt declarate astfel dacă trec de filtrele la nivel de metoda dar nu trec de compararea grafurilor de apeluri.

Metodele grupate cu grad de asemănare *High* sunt metode care sunt considerate suspecte de copiat cu o probabilitate foarte mare. Ele au trecut atât de filtrele la nivel de metodă cât și de cele la nivel de grafuri de apeluri.

Clasele au trei grade de asemănare : *Low*, *Medium*, *High*

Clasele suspecte Low sunt clasele care trec de filtrele la nivel de clasă, dar nu și de cele la nivel de metodă. Sunt puțin probabile să fie copiate, numărul de clase asemănătoare la acest nivel putând fi mare. Nu au o pondere importantă în luarea deciziei dacă două proiecte trebuie afișate sau nu ca suspecte de copiat.

Clasele suspecte Medium sunt clase care trece de filtrele la nivel de clasă și care au metode asemănătoare. Totuși numărul de metode ce sunt foarte probabil copiate (grupate *High*) nu este suficient de mare pentru ca aceste să fie declarate foarte suspecte. Acest număr se setează în fișierul de profile. Putem considera că sunt clase ce trec de filtru la nivel de clasă și cel la nivel de metodă, dar nu îndeplinesc condițiile la nivel de graf de apeluri.

Clasele suspecte High sunt clase foarte probabil copiate. Ele au trecut de filtrele la nivel de clasă, iar metodelor lor sunt asemănătoare atât la nivel de metrici cât și la nivel de graf de apeluri. Existența unor astfel de grupuri de clase între două proiecte este luat în considerare, acele proiecte putând fi copiate.

O facilitate oferită de CloneDetector este aceea de a exclude unele clase, pachete sau foldere din analiza. Aceasta se dovedește foarte utilă în cazul în care toate proiectele au porțiuni de cod date de la început. De exemplu, pentru unele teme de laborator se oferă un set de clase utilitare de la care să se pornească. Acestea ar fi cu siguranță în toate proiectele și tool-ul le-ar grupa. Dacă numărul lor ar fi semnificativ ele ar influența verdictul, multe proiecte devenind astfel suspecte datorită acestor clase. Nu dorim acest lucru, deoarece s-ar pierde mult timp pentru analiza suspjecțiilor neadevărați. Pentru excluderea claselor am creat un fișier `commonClasses.pl` în care putem scrie clasele ce sunt comune între toate proiectele și pe care dorim să le excludem din analiză.

V. Implementarea programului

Primul lucru pe care trebuie să îl facem este să încărcăm cele două proiecte în baza de cunoștințe Prolog. Ele vor fi salvate pe disc în format qlf. Încărcarea lor în Prolog se realizează cu comanda `consult` care compiliează fișierele .pl iar pe cele qlf le încarcă direct.

Totuși simplă încărcare în baza de date generează o problemă, coliziuni între identificatori. După cum am menționat anterior fiecare faptă are ca prim termen un identificator unic în cadrul unui proiect. Aceștia încep de la 10001 și cresc în funcție de dimensiunile proiectului. Totuși acești identificatori nu vor fi unici dacă unim două baze de cunoștințe în același spațiu. Se vor suprapune identificatorii și ne va fi imposibil să deducem care din ce proiect face parte. De aceea încărcarea simplă în baza de cunoștințe, cu `consult` nu este o soluție.

Soluția găsită presupune folosirea modulelor din Prolog.

1. Module

Un modul Prolog este o colecție de predicate care definesc o interfață publică prin intermediul unui set de predicate furnizate și a operatorilor. Module Prolog sunt definite de un standard ISO. InSWI-Prolog sistemul e modul este derivat din sistemul de Quintus modulul Prolog.

În sistemele clasice Prolog, toate predicatele sunt organizate într-un singur namespace și orice predicat poate apela orice predicat. Deoarece fiecare predicat într-un fișier poate fi apelat de oriunde în program, devine foarte greu să găsim dependențele și să modificăm un predicat, fără să riscăm să stricăm programul general. Acest lucru este valabil pentru orice limbaj, dar chiar și mai rău pentru Prolog din cauza nevoii sale frecvente de predicate de "ajutor".

Un modul Prolog încapsulează un set de predicate și definește interfață. Modulul poate importa alte module, ceea ce creează dependențe explicite. Având în vedere dependențele explicite și o interfață bine definită, devine mult mai ușor să schimbăm organizarea internă a unui modul fără a strica aplicația generală.

Prolog conține două module speciale. Primul este modulul `system`. Acesta conține toate predicatele din librăria standar Prolog, cea care se încarcă la lansare programului. Acesta nu importă alte module. Al doilea modul special se numește `user`, acesta fiind spațiul de lucru implicit la pornirea Prolog-ului. Inițial este gol, dar importă modulul `system`. Astfel toate predicatele predefinite devin disponibile.

Modulele au două proprietăți declarate, numele modului și predicatele pe care acel modul le exportă. Astfel noi am creat două module și le-am folosit pentru a încărca în fiecare din ele un proiect. Predicatele din proiectele încărcate nu vor fi exportate, declarate vizibile. Vom declara vizibile niște predicate de acces la date pentru fiecare modul. Am făcut convenția ca predicatele exportate din primul modul să se termine în 1 iar cele din cel de-al doilea modul să se termine cu 2.

În proiectul nostru avem create două module `mod1` și `mod2`. Inacestea se vor încărca pe rând proiectele comparate. Modulele se creează în mod static, declarând modulul cu `module/2`, într-un fișier pe prima linie ca în Fig.5.1. După cum se observă declarația va conține numele, urmată de predicatele exportate. Nu se vor exporta predicatele generate de `JTransformet`, cum ar fi `classT`, `methodT` s.a pentru a evita problema descrisă. Pentru încărcare celor două module în modulul de start, `user`, folosim apelul `use_module/2`.

După cum se observă din figura în acest fișier nu am plasat logica de interogare a bazei de cunoștințe ci o apelăm prin redirectare, implemmentarea fiind realizată în alt fișier.

Am ales această abordare deoarece cele două proiecte încărcate în module diferite, vor fi interogate la fel și ar fi fost greșit să implementăm aceleași lucru în două locuri. Ar fi mult mai greu de modificat și întreținut astfel. Logica de interogare a bazei de cunoștințe va fi implementată în `utilities.pl`, fișier importat de ambele module.

Modulul doi va fi identic cu modulul 1 schimbându-se doar numărul de la finalul fiecărui predicat din 1 în 2.

```
:-module(mod1,[myClass1/2,load1/1,calcNrAtrib1/2,clearDatabase1/0,calcNrMet1/2,
               calcNrInterf1/2,calcNrExtends1/2,methodsOfClass1/3,nrOfCycles1/2,
               nrOfIf1/2,nrOperators1/3,callT1/3,isInterface1/1,isStatic1/1,
               paramMet1/2,param1/3]).

:-include('./utilities.pl').
:-include('./commonClasses.pl').

load1(Prj):- load(Prj).

myClass1(ClassId,Name):- myClass(ClassId,Name).

calcNrAtrib1(IdClasa,Nr):- calcNrAtrib(IdClasa,Nr) .

clearDatabase1:-clearDatabase.

calcNrMet1(IdClasa,Nr):-calcNrMet(IdClasa,Nr) .

calcNrInterf1(IdClasa,Nr):-calcNrInterf(IdClasa,Nr) .

calcNrExtends1(IdClasa,Nr):-calcNrExtends(IdClasa,Nr) .

methodsOfClass1(IdClasa,IdMethod,Name):-methodsOfClass(IdClasa,IdMethod,Name) .

nrOfCycles1(MethodId,Nr):-nrOfCycles(MethodId,Nr) .

nrOfIf1(MethodId,Nr):-nrOfIf(MethodId,Nr) .

nrOperators1(IdMethod,Operator,Nr):-nrOperators(IdMethod,Operator,Nr) .

callT1(MethodId1,CalledClassId1,CalledMetId1):-
    callDep(MethodId1,CalledClassId1,CalledMetId1) .

isInterface1(ClassId):-isInterface(ClassId) .
```

Fig. 5.1 Interfata catre modulul 1

În concluzie, ar trebui să privim modulele ca două bucăți din toată baza de cunoștințe Prolog. Toată baza de cunoștințe este încărcată în modulul `user`, Fig. 5.2, predicatele predefinite, predicatele implementate de noi ce conțin logică programului și cele două module. Din modulul `user` nu se pot accesa decât predicatele vizibile ale modulului 1 respectiv 2, celelalte pot fi apelate doar prin intermediul celor vizibile, adică din cadrul modulului.

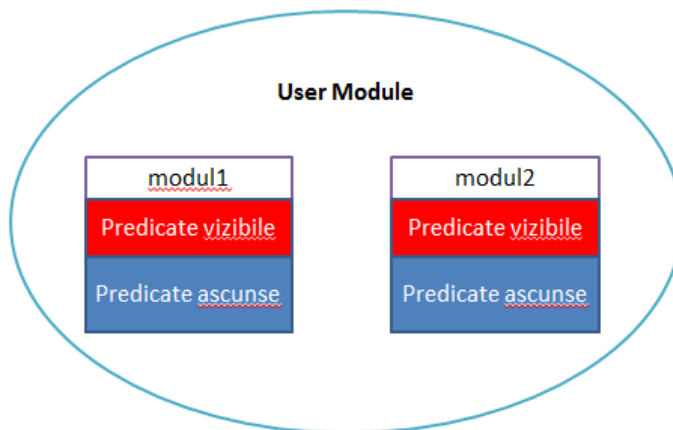


Fig. 5.2 Reprezentare module

2. Împărțirea pe fișiere

În continuare vom descrie împărțirea pe fișiere în cadrul programului și o scurtă descriere a fiecăruia.

Proiectul CloneDetector este format din 4 fișiere de bază și două de configurare :

- cloneDetector.pl
- utilities.pl
- IProject1.pl
- IProject2.pl
- commonClasses.pl
- profiles.pl

Fișierul cloneDetector.pl este fișierul de lansare a programului și cel mai complex. Conține logica de nivel înalt a unelei, luând deciziile pe bază informațiilor obținute despre cele două proiecte. Se încarcă automat la execuție în modulul user și va încărca cele două module : mod1, mod2.

Utilities.pl este fișierul încărcat în ambele module. El conține predicate de interogare a bazei de cunoștințe, obținând informații despre proiecte, calculând metrici. Este al doilea ca și complexitate. Nici un predicat din el nu va fi vizibil în modulul user, predicatele fiind apelate prin redirectare : apelăm predicatul unui modul care va apela predicatele din acest fișier.

IProject1.pl respectiv IProject2.pl sunt fișierele de modul. Întrucât în baza de cunoștințe exportată de JTransformet nu se pun automat directive de modul, am creat aceste două fișiere care creează modulul și interfața către cele două proiecte încărcate la un moment dat în baza de cunoștințe. Figura 5.1 este o poză a fișierului IProject1.pl. După cum spuneam ele definesc modulul la început. Tot în aceste fișiere se încarcă în module fișierul comun utilities.pl prin directivă *include*.

Fișierul commonClasses.pl este un fișier de configurare în care ne putem defini clasele ce dorim să le excludem din analiză. Se folosește când avem proiecte care au cod dat de la început și nu vrem să ne afecteze compararea, generând prea mulți suspecti. O descriere a modului de completare a acestor fișiere se găsește în capitolul VI. Manualul utilizatorului.

În profiles.pl sunt setate profilele de execuție, aspect asupra căruia vom reveni. Și acesta este un fișier de configurare.

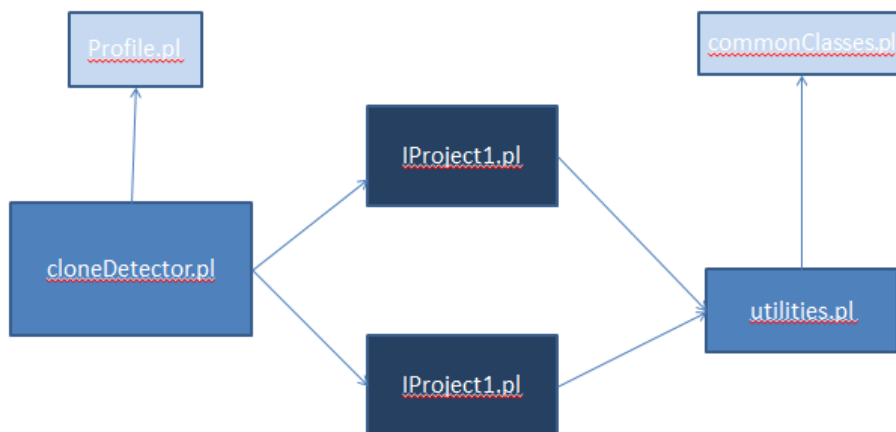


Fig. 5.3 Relații între fișier

În Fig. 5.3 prezentăm relațiile dintre fișiere. Săgeata reprezintă ce fișier interoghează pe altul, vârful săgeții fiind îndreptat spre cel folosit. Nuanțele de albastru sugerează asemănările dintre ele, fișierele deschise fiind de configurare, cele închise definind modulele iar celelalte două cuprinzând logica programului de nivel înalt și de interogare a bazei de date.

3. Diferențierea claselor proiectelor de clasele librăriilor

JTransformer generează fapte atât pentru clasele scrise de noi cât și pentru clasele din librăriile incluse, chiar și din librarial standard Java. Nu dorim să grupăm clase care nu sunt scrise de noi sau să le luăm în considerare în calcularea unor metrici. De aceea am creat predicatul `myClass` care exclude din analiză clasele librăriilor și clasele pe care noi le declarăm comune pentru mai multe proiecte.

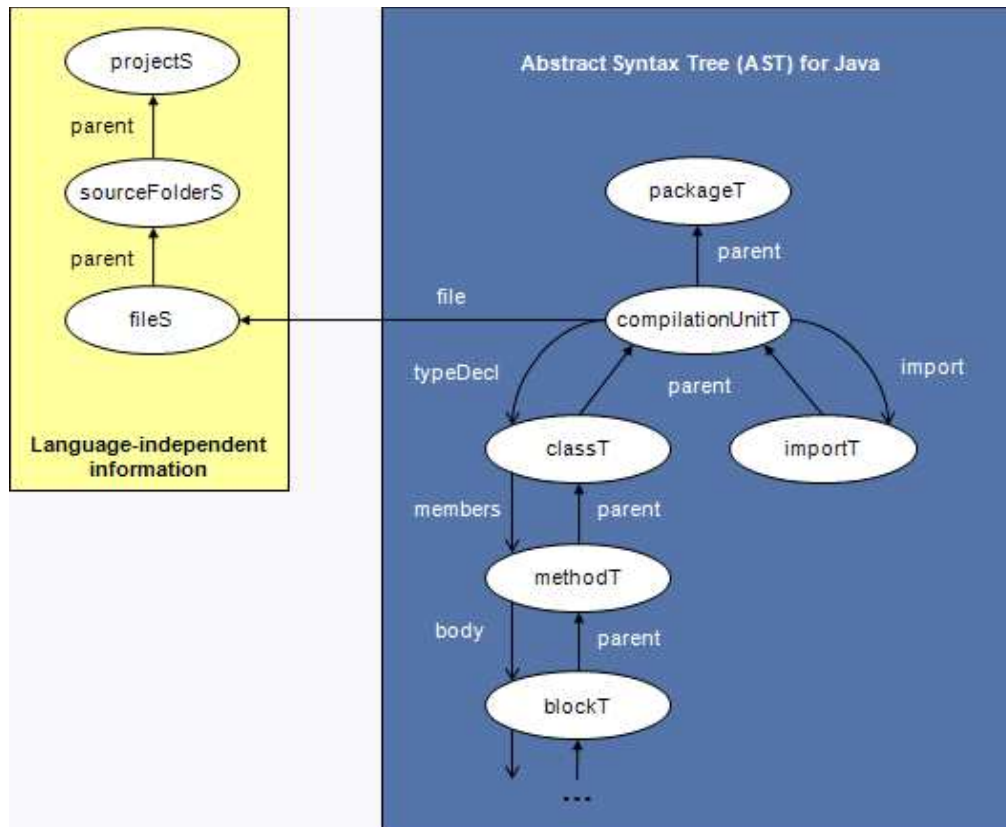


Fig. 6.4 Legături între fapte specifice Java și independente de limbaj

JTransformer generează două tipuri de fapte, independente de limbaj și dependente, după cum se poate vedea și în figura 6.4. Relațiile dintre ele sunt prezentate prin săgeți, arcul săgeții fiind îndreptat spre părinte. Între majoritatea clauzelor există arce în ambele direcții. Arcele se stabilesc în baza de cunoștințe prin identificatori.

Noi am folosit aceste relații între fapte pentru a determina care sunt clasele noastre. Pentru clasele din librării nu se generează clauze independente de limbaj. Predicatul `myClass` generează toate clasele unui proiect pornind de la axioma `classT` și verificând ca acestea să aibă axiomele independente de limbaj generate.

```
myClass(AllCls,Name) :-
  classT(AllCls,Cu,Name,_),
  compilationUnitT(Cu,_Fid,_),
  fileS(Fid,Src,Path),
  sourceFolderS(Src,Pid,_),
  projectS(Pid,_,_,_),
  findall(PartPath,common(Path,PartPath),Result),
  count(Result,Nr), Nr=0.
```


Pentru clasele ce fac parte din proiectul nostru dar dorim noi să le excludem, de obicei ele fiind date la începutul proiectului și deci le regăsim în toate proiectele, se generează aceste axiome independente de limbaj. Pentru a le distingem am creat un predicat special *common* ce verifică dacă o clasă a fost definită ca fiind comună.

```
common(Path,PartPath):-commonClass(PartOfPath),
    ( atom_concat(PartPath,PartOfPath,Path);
      atom_concat(PartOfPath,PartPath,Path) ).
```

4. Implementarea metricilor

În continuare vom prezenta modul de implementare a metricilor propuse cu exemplificari din codul sursa. Întrucât pentru implementare lor sunt apelate predicate din mai multe fișiere vom specifica, unde este cazul din ce fișier fac parte. Este important să reținem că predicatele care se termină în 1 sunt din fișierul *IProject1*, ce descrie modulul 1, încărcând în el primul proiect iar cele ce se termină în 2 sunt din fișierul *IProject2* care îndeplinește aceeași funcționalitate pentru proiectul 2.

A. Metrici la nivel de casa

Fiecare metrică este un filtru de care clasele suspecte pot trece sau nu. Trecerea lor presupune existența asemănărilor și continuarea către următorul filtru. Picarea unui filtru întrerupe compararea acestor clase. Ordinea în care sunt puse filtrele are importanță, putând îmbunătăți viteza de execuție a programului. Dacă punem filtre ușor de calculat și care caracterizează mai bine clasele, vom opri compararea claselor diferite mai repede.

În cod filtrele la nivel de clasă sunt apelate de predicatul *compareClassLevel/2* care primește ca termeni identificatorii unici ai celor două clase comparate.

```
compareClassLevel(ClassID1,ClassID2):-
    areInterfaces(ClassID1,ClassID2),
    compareNrAtrib(ClassID1,ClassID2),
    compareNrMet(ClassID1,ClassID2),
    compareNrInterf(ClassID1,ClassID2),
    compareNrSuperClass(ClassID1,ClassID2).
```

Ordinea filtrelor aleasă de noi este prezentată în porțiunea de cod de mai sus.

Primul filtru *areInterfaces/2* verifică ca id-urile primite să fie ambele a unor clase sau a unor interfețe. *JTransformer* va genera atât pentru clase cât și pentru interfețe o faptă de tipul *classT*, iar pentru interfețe va mai genera una în plus de tipul *interfaceT*.

```
%returns true if class is an interface.
```

```
isInterface(ClassId):-
    interfaceT(ClassId).
```

```
areInterfaces(ClassID1,ClassID2):-
    ( isInterface1(ClassID1),
      isInterface2(ClassID2));

    ( not(isInterface1(ClassID1)),
      not(isInterface2(ClassID2))).
```


Filtrul al doilea numara numarul de attribute al unei clase si verifica ca doua clase au diferenta dintre numarul de attribute mai mic ca MaxDifference si procentul dintre cele doua mai mare ca MinProcent. Acest tolerante se initializeaza în functie de profilul setat.

```
compareNrAtrib(Class1,Class2):-
    profile(Profile),
    classDelta(Profile,[MaxDifference,MinProcent,_,_,_,_]),
    calcNrAtrib1(Class1,Nr1),
    calcNrAtrib2(Class2,Nr2),
    (
        (delta(Nr1,Nr2,Difference),   Difference =< MaxDifference),
        (procent(Nr1,Nr2,Procent),   MinProcent =< Procent)
    ).
```

Numărul de attribute din fiecare clasă se obține apelând interfețele pentru fiecare modul în care sunt încărcate cele două proiecte comparate. Apelurile vor fi redirectate către predicatul din fișierul utilities.pl `calcNrAtrib/2`. Pentru înțelegerea predicatului, vom explica mai întâi cum se utilizează *findall/3*. Findall este un predicat oferit de Prolog pentru a genera toate soluțiile unei reguli. Generează o listă de soluții, memorată în ultimul termen. Primul termen este cel ce va fi adăugat în listă pentru fiecare soluție găsită, iar al doilea este regula care va genera soluțiile. Astfel pe exemplu de mai jos, căutăm toate faptele `fieldT` care au ca părinte clasa primită ca parametrul. Identificatorul , notat cu X, va fi adăugat pentru fiecare soluție în lista `Result`. Numărul de elemente din lista rezultată va fi numărat prin clauza `count` și returnat prin termenul `Nr`. Această abordare se adopta la multe predicate ce calculează metrici, iar înțelegerea modului de funcționare a predicatului predefinit `findall` este esențială. Dacă nu am fi pus `findall`, ci am fi apelat `fieldT` pur și simplu am fi obținut tot aceste soluții dar nu deodată, s-ar fi așteptat intereacțiunea utilizatorului pentru a genera următoarea soluție.

```
calcNrAtrib(IdClasa,Nr) :-
    forall(X,fieldT(X,IdClasa,_,_,_),Result),
    count(Result,Nr).
```

Trebuie să justificăm de ce am ales doar să numărăm attributele unei clase, în loc să verificăm și tipul lor. Verificarea tipului unui atribut și existența unuia de același tip în clasa comparată este ușor pentru tipurile elementare (`String`, `int`, `float...`), dar dacă avem ca attribute clase proprii lucrurile se complică deoarece nu ne putem baza pe numele lor în comparare. Ar trebui să verificăm asemănarea lor nu după nume ci după metrici, eventual după attributele lor, metodele s.a.m.d. Astfel am merge recursiv dintr-o clasă în alta, iar în unele cazuri am obține chiar cicluri dacă o clasă conține ca atribut o alta care o va conține la rândul ei pe ea, dependentă dublă între ele. Astfel ne trebuiau niște metrici independente de tip, care să nu genereze cicluri infinite, pentru compararea claselor. Deși mai slabe, prin necompararea tipurilor, metricile alese la nivel de clasă nu generează cicluri și sunt temelia analizei. Verificare tipurilor se va face ulterior la nivel de metodă și ne vom baza pe aceste metrici pentru asta.

Predicatul `compareNrMet(ClassID1,ClassID2)` este asemănător în implementare cu cel ce calculează numărul de attribute, calculând numărul de metode. De această dată toleranța metodelor este exprimată doar prin diferența absolută între numărul lor, în funcție de profilul ales.

```
compareNrMet(ClassID1,ClassID2):-
    profile(Profile),
    classDelta(Profile,[_,_,MaxDifference,_,_,_]),
    calcNrMet1(ClassID1,Nr1),
    calcNrMet2(ClassID2,Nr2),
    delta(Nr1,Nr2,Difference),
    Difference =< MaxDifference.
```

Următoarele două predicate se leagă de relațiile de moștenire ale claselor, o caracteristică importantă în programarea orientată pe obiecte. Verificăm că numărul de interfețe pe care o clasă îl implementează să fie aproximativ același iar dacă o clasă moștenește pe alta atunci și posibilă copie va trebui să extindă o clasa.

```
compareNrSuperClass(ClassID1,ClassID2):-
    calcNrExtends1(ClassID1,Nr1),
    calcNrExtends2(ClassID2,Nr2),
    Nr1 = Nr2.
```

```
compareNrInterf(ClassID1,ClassID2):-
    profile(Profile),
    classDelta(Profile,[_,_,MaxDifference,_,_,_]),
    calcNrInterf1(ClassID1,Nr1),
    calcNrInterf2(ClassID2,Nr2),
    delta(Nr1,Nr2,Difference),
    Difference =< MaxDifference.
```

Nu am mai prezentat cum este interogată baza de cunoștințe pentru ultimele filtre pentru că este asemănător cu calcularea numărului de atribute. Se folosește *findall* dar pe fapte diferite *methodT*, *extendsT*, *implementsT*.

B. Metrici la nivel de metoda

În compararea a două clase, după filtrele la nivel de clasă, CloneDetector urmează să compare fiecare metodă a claselor încercând să le grupeze. Primul nivel de analiză a metodelor este la nivel de metrici. Predicatul din care pornește analiza la nivel de metodă este *compareMethodLevel/4* care apelează predicatul *methodMetrics/2* pentru filtrele la nivel de metodă.

Metricile la nivel de metoda sunt inlantuite astfel. Am pus primul filtru cel ce ia în considerare ciclurile din programe deoarece aceste sunt rareori modificate și reprezintă o semnatura puternică a metodei.

```
methodMetrics(MethodId1,MethodId2):-
    whileFilter(MethodId1,MethodId2),
    operatorsFilter(MethodId1,MethodId2),
    ifFilter(MethodId1,MethodId2),
    areStatic(MethodId1,MethodId2),
    compareSigniture(MethodId1,MethodId2).
```

whileFilter verifică numărul de cicluri dintr-o metodă, numărând câte while, do while și for există. Este foarte greu să schimbi într-o clonă o astfel de metrică, înlocuirea ciclurilor repetitive putând fi făcută doar prin recursivitate. De aceea în acest caz nu admitem toleranțe, numărul ciclurilor trebuie să fie exact la fel în două metode copiate.

```
whileFilter(MethodId1,MethodId2):-
    nrOfCycles1(MethodId1,Nr1),
    nrOfCycles2(MethodId2,Nr2),
    Nr1=Nr2.
nrOfCycles(MethodId,Nr):-findall(MethodId,whileT(_,_ ,MethodId,_ ,_),While),
    count(While,NrWhile),
    findall(MethodId,doWhileT(_,_ ,MethodId,_ ,_),DoWhile),
    count(DoWhile,NrDoWhile),
    findall(MethodId,forT(_,_ ,MethodId,_ ,_),For),
    count(For,NrFor),
    Nr is (NrWhile + NrDoWhile + NrFor).
```

Filtru operatorsFilter implementează metrică propusă de Halstead, care spune că o bucată de cod poate fi definită de numărul de operatori și de tipul lor. Operatori propuși sunt de adunare, scădere, înmulțire și împărțire, operațiile matematice elementare.

În urma analizei acestei metrici am observat că ea ajută la distingerea claselor, dar în cazul Java operatorul “+” nu poate fi analizat el fiind folosit de numeroase ori pentru concatenarea șisurilor. Astfel prin simplă schimbare a mesajelor afișate, și ștergerea operatorului de concatenare s-ar putea împiedica detecția clonei. De aceea noi am analizat doar ultimii trei operatori.

```
operatorsFilter(MethodId1,MethodId2):-
    nrOperators1(MethodId1,-,NrSub1),
    nrOperators2(MethodId2,-,NrSub2),
    NrSub1=NrSub2,

    nrOperators1(MethodId1,/,NrDiv1),
    nrOperators2(MethodId2,/,NrDiv2),
    NrDiv1=NrDiv2,

    nrOperators1(MethodId1,*,NrMul1),
    nrOperators2(MethodId2,*,NrMul2),
    NrMul1=NrMul2.
```

Se observă că se cere o egalitate exactă a operatorilor între două metode copiate. Toate tipurile de operatori sunt memorați în baza de cunoștințe prin fapte de tipul *operationT*.

```
ifFilter(MethodId1,MethodId2):-
    profile(Profile),
    methodDelta(Profile,[MaxDifference]),
    nrOfIf1(MethodId1,Nr1),
    nrOfIf2(MethodId2,Nr2),
    delta(Nr1,Nr2,Delta),
    Delta =<= MaxDifference.
```

Filtrul la nivel de condiții compară câte if-uri sunt într-o metodă. Am admis o tolerantă în acest caz, pentru posibile adăugări sau ștergeri de conditii, prin tratarea unor cazuri speciale sau renunțarea la ele.

Filtrul *areStatic* este asemănător ca implementare cu filtrul de la nivel de clasă *areInterface*. Acesta va verifica dacă o metodă este sau nu statică și va impune ca și perechea ei să fie.

Un filtru complex care diferențiază destul de mult metodele este ultimul de la acest nivel, *compareSignature/2* care compară semnătura a două metode. Deși este un filtru puternic a fost lăsat la sfârșit pentru că el consumă mai mult timp de procesare, comparațiile făcute de el fiind similare cu cele de la grafurile de apeluri.

Filtrul compară ca numărul de parametri a două metode să fie egal dar și ca pentru fiecare parametru de un anumit tip să existe unul de același tip în metoda copiată. Adică dacă am avea un parametru String în prima metodă să existe unul de același tip și în cea suspectă.

```
compareSignature(MethodId1,MethodId2):-
    paramMet1(MethodId1,ListParam1),
    paramMet2(MethodId2,ListParam2),
    count(ListParam1,NrParam1),
    count(ListParam2,NrParam2),
    NrParam1 = NrParam2,
    findall(_,paramCompare(MethodId1,MethodId2),_),
    findall(Id,signMatch(Id,_),ResultList),
    count(ResultList,NrMatch),
    retractall(signMatch(_,_)),
    NrMatch = NrParam1 .
```

Se obține la început listă parametrilor celor două metode și se verifică dacă numărul lor este egal. În caz contrar se oprește compararea. Dacă au același număr de parametri se vor compara pentru toți tipuri și se va încerca să se grupeze doi câte doi. Memorarea perechilor se va face dinamic prin fapte *signMatch* care au ca termeni id-urile celor doi parametri. Se folosește predicatul predefinit *assert*.

```
paramCompare(MethodId1,MethodId2):-
    param1(MethodId1,ClsId1,IdParam1),
    param2(MethodId2,ClsId2,IdParam2),
    not(signMatch(IdParam1,_)),
    not(signMatch(_,IdParam2)),
    compareClassLevel(ClsId1,ClsId2),
    assert(signMatch(IdParam1,IdParam2)).
```

Dacă un parametru a fost deja împerecheat cu unul, el nu va mai fi grupat cu altul, lucru verificat la începutul predicatului *paramCompare*. Observăm că tipul parametrilor se verifică comparând cele două clase cu metricile definite la nivel de clasă. Această abordare este corectă atât pentru parametri ce au tipuri elementare sau definite de noi, JTransformer generând interfețele claselor și pentru clasele din librării.

După ce s-a încercat să se grupeze toți parametrii, se verifică ca numărul celor împerecheați să coincidă cu numărul total de parametri. Pentru ca baza de cunoștințe să rămână intactă pentru următoarea verificare se șterg toate clauzele alocate dinamic prin predicatul predefinit *retractall*.

5. Grafuri de apel

Dacă o pereche de două clase suspectate de copiat trece de primele două mari categorii de filtre, metrici la nivel de clasă și la nivel de metodă, urmează să se analizeze la nivel de metodă toate apelurile pe care o metodă le face către o alta. După cum am explicat în capitolul de descriere a algoritmului, metodele dar și clasele apelate trebuie să fie asemănătoare. Această asemănare este stabilită dacă clasele și metodele analizate trec de metricile la nivel de clasă respectiv de metodă.

În baza de cunoștințe apelurile sunt reprezentate sub forma de fapte *callT* iar id-ul metodei și clasei apelate se obține cu predicatul *callDep*. Acesta analizează doar apelurile între clasele proiectului, nu și către librării.

```
callDep(MethodId,CalledClassId,CalledMetId):-
    callT(____,MethodId,____,CalledMetId),
    methodT(CalledMetId,CalledClassId,____,____),
    myClass(CalledClassId,_____).
```

Predicatul *callDependencies* este cel din care se pornește analiza și cel care stabilește condițiile pentru a trece de acest filtru. Ținând cond că ar fi complicat să se schimbe apelurile către metode în proiecte copiate și este rareori întâlnit, nu am permis toleranțe cerând o egalitate exactă. La început verificăm, înainte de a încerca să grupăm clasele, dacă au același număr de apeluri. După aceea generăm grupurile de apeluri similare, le numărăm și verificăm să fie egale cu numărul total de apeluri.

```
callDependencies(MethodId1,MethodId2):-
    findall(CalledMet1,callT1(MethodId1,____,CalledMet1),ListMet1),
    uniqueList(ListMet1,UniqueList1),
    count(UniqueList1,NrCalledMet1),
    findall(CalledMet2,callT2(MethodId2,____,CalledMet2),ListMet2),
    uniqueList(ListMet2,UniqueList2),
    count(UniqueList2,NrCalledMet2),
    NrCalledMet1 = NrCalledMet2,
    findall(MethodId1,compareAllCallT(MethodId1,MethodId2),MatchingCalls),
    count(MatchingCalls,Nr),
    retractall(partialMatch(____)),
    retractall(partialMetMatch(____)),
    Nr = NrCalledMet1.
```

Predicatul care generează și testează grupările de apeluri posibile este *compareAllCallT*. Se verifică înainte de a încerca să grupăm două apeluri, dacă unul dintre ele nu este deja grupat. Se observă că grupurile de metode sunt alocate dinamic prin fapte de tipul *partialMetMatch* iar clasele prin fapte de tipul *partialMatch*. Aceste fapte memorează identifiatorii celor două perechi.

```
compareAllCallT(MethodId1,MethodId2):-
    callT1(MethodId1,CalledClassId1,CalledMetId1),
    callT2(MethodId2,ClassId2,CalledMetId2),
    not(partialMetMatch(CalledMetId1,____)),
    not(partialMetMatch(____,CalledMetId2)),

    methodMetrics(CalledMetId1,CalledMetId2),
```

```
(
  (
    not(partialMatch(CalledClassId1,_) ,
    not(partialMatch(_,ClassId2)),
    compareClassLevel(CalledClassId1,ClassId2),
    assert(partialMetMatch(CalledMetId1,CalledMetId2)),
    assert(partialMatch(CalledClassId1,ClassId2))
  );
  (
    partialMatch(CalledClassId1,ClassId2),
    not(partialMetMatch(CalledMetId1,CalledMetId2)),
    assert(partialMetMatch(CalledMetId1,CalledMetId2))
  )
).
```

Testând programul de-al lungul dezvoltării sale, am realizat că această filtru are o importantă deosebită ducând la o creștere semnificativă a corectitudinii grupării claselor și a verdictelor date. Această filtru se poate îmbunătăți, surprinzător, odată cu dezvoltarea metricilor la nivel de clasă sau metode, deoarece ea folosește aceste metrice la compararea entităților apelate.

6. Luarea deciziilor

Până acum am prezentat în mare modul de implementare al algoritmului explicând partea de metrice la nivel de clasă și metodă și modul de verificare a apelurilor metodelor. Totuși intercalat sau după aceste filtre sunt numeroase puncte de decizie în program care au un rol crucial. Toate rezultatele generate de filtre fiind interpretate în aceste puncte urmând să se enunțe unele verdicte. Putem să ne imaginăm această unealtă ca pe un proces juridic din viața reală, unde probele sunt datele generate de filtre, iar verdictele intermediare sau finale date de judecător sunt deciziile date de aceste puncte de control.

Vom prezenta modul de luare al deciziilor pentru clase. După cum am mai amintit există trei nivele de asemănare între clase *LOW*, *MEDIUM*, *HIGH*. Compararea claselor a două proiect începe în predicatul *generateAllMatchingClasses*. Se generează perechi de clase care nu au fost încă memorate în baza de date cu probabilitate ridicată de comparare. Acestea sunt comparate în predicatul *compare2Classes*. Dacă clasele trec de toate filtrele sunt declarate foarte probabil copiate și memorate în baza de cunoștințe.

```
generateAllMatchingClasses:-
  myClass1(Id1,Name1),
  myClass2(Id2,Name2),
  not(match(Id1,_,_,high)),
  not(match(_,_,Id2,high)),
  compare2Classes(Id1,Name1,Id2,Name2),
  retractall(match(Id1,_,_,_)),
  retractall(match(_,_,Id2,_)),
  assert(match(Id1,Name1,Id2,Name2,high)).
```

```
compare2Classes(Id1,Name1,Id2,Name2):-
  compareClassLevel(Id1,Id2),
  assertClassMatchLow(Id1,Name1,Id2,Name2),
  compareMethodLevel(Id1,Name1,Id2,Name2).
```

Dacă în predicatul *compare2Classes* se trece de compararea la nivel de clasă vom memora dinamic în baza de cunostințe perechea de clase ca având probabilitate scăzută de asemanare. Urmează să fie analizate la nivel de metodă, probabilitate putând fi schimbată.

```
assertClassMatchLow(Id1,Name1,Id2,Name2):-
(
  not(match(Id1,_,_,medium)),
  not(match(_,Id2,_,medium)),
  not(match(Id1,_,_,low)),
  not(match(_,Id2,_,low)),
  assert(match(Id1,Name1,Id2,Name2,low)),!
);1=1.
```

Continuăm prin a prezenta cum sunt luate deciziile în continuare, în predicatul *compareMethodLevel*. Aici practic se analizează toate metodele dintr-o clasă. Clauza începe prin a genera toate grupurile de metode cu cele două tipuri de probabilități scăzute sau ridicate. Odată ce avem aceste perechi verificăm ca diferențele între numărul lor și numărul total de metode ale celor două clase să fie în limitele stabilite de profil, atât din punct de vedere al diferențelor absolute cât și procentual. Predicatul e gândit că dacă nu sunt suficiente metode grupate să nu se schimbe probabilitatea de copiere a celor două clase și să returneze fals. Dacă sunt suficiente perechi dar nu sunt cu probabilități mari de a fi copiate, atunci clasele devin medium și predicatul returnează fals. Dacă se ajunge la finalul predicatului înseamnă că nu suntem în cele două cazuri, deci clasele sunt foarte probabil copiate și returnăm true. Returnând true, clasele devin foarte probabil copiate în predicatul de un este apelat *compareMethodLevel*.

```
compareMethodLevel(ClassId1,Name1,ClassId2,Name2):-
  findall(_generateAllMatchingMethods(ClassId1,ClassId2),_),
  findall(Id,methodMatch(Id,_,_,high),ResultHigh),
  count(ResultHigh,NrOfMatchesHigh),
  findall(Id,methodMatch(Id,_,_,low),ResultLow),
  count(ResultLow,NrOfMatchesLow),
  NrOfMatches is NrOfMatchesHigh + NrOfMatchesLow,
  calcNrMet1(ClassId1,NrMet1),
  calcNrMet2(ClassId2,NrMet2),
  retractall(methodMatch(_,_,_,_)),
  profile(Profile),
  classDelta(Profile,[_,_,_,_,MaxNrOfUnmatchingMethods,MinProcent]),
  delta(NrOfMatches,NrMet1,Delta1),
  Delta1=<MaxNrOfUnmatchingMethods,
  delta(NrOfMatches,NrMet2,Delta2),
  Delta2<MaxNrOfUnmatchingMethods,
  procent(NrMet1,NrMet2,Procent),
  Procent> MinProcent,
  procent(NrOfMatchesHigh,NrOfMatches,ProcentMatchHigh),
  not(
    assertClassMedium(ClassId1,Name1,ClassId2,Name2,ProcentMatchHigh)
  ).
```


assertClassMedium este predicatul în care decidem dacă clasele ar trebui să fie medium sau dacă le vom memora ca perechi high.

```
assertClassMedium(ClassId1,Name1,ClassId2,Name2,ProcentMatch):-
    profile(Profile),
    classDelta(Profile,[_,_,_,ProcentMatchHigh,_,_]),
    ProcentMatch < ProcentMatchHigh
    retractall(match(ClassId1,_,_,_)),
    retractall(match(_,ClassId2,_,_)),
    assert(match(ClassId1,Name1,ClassId2,Name2,medium)).
```

În această secțiune am prezentat doar deciziile luate la nivel de clasă, modul în care hotărâm dacă clasele sunt copiate cu probabilități ridicate, medii sau scăzute. Deciziile acestea se bazează pe modul în care metodele sunt grupate și probabilitățile lor. Deci, evident, alte puncte de decizie vor fi și la nivel de metodă. Acestea sunt relativ similare, și nu vor mai fi prezentate în continuare. Se pot inspecta cu atenție în program dacă se dorește la predicatele *generateAllMatchingMethods* și *assertMetMatch*.

7. Compararea a doua proiecte

În această parte vom descrie modul în care se lansează în execuție unealtă, și operațiunile pe care le executăm la început. Compararea a două proiecte începe din clauza *run/3* care primește ca termeni numele cu calea către proiecte, și profilul.

La început este golită baza de cunoștințe de faptele memorate de o execuție anterioară și încărcăm în ea profilul ales, interfețele către cele două module și în cele două module încărcăm proiectele comparate.

Înainte de a lansa compararea tuturor claselor verificăm ca numărul claselor din cele două proiecte să nu difere foarte mult, altfel fiind inutilă compararea. *generateAllMatchingClasses* va genera toate perechile de clase copiate în baza de cunoștințe.

Ultimele funcții îndeplinite de predicatul *run* fiind interpretarea rezultatului la nivel de proiect și afișarea lui. Verificarea se face prin compararea diferențelor absolute între numărul de clase împerecheate și numărul de clase din cele două proiecte, ținând cont de toleranțe.

Codul de mai jos arată modul de implementare a predicatului *run*, de unde pornim execuția.

```
run(Proj1,Proj2,Profile):-
    retractall(profile(_)),
    assert(profile(Profile)),
    profile(Profile),
    projectDelta(Profile,[MaxNumberDeltaClasses|_]),
    use_module('./IProject1.pl'),
    use_module('./IProject2.pl'),
    load1(Proj1),
    load2(Proj2),

    findall(Id1,myClass1(Id1,_),ClsPrj1),
    count(ClsPrj1,NrCls1),
    findall(Id2,myClass2(Id2,_),ClsPrj2),
    count(ClsPrj2,NrCls2),
    delta(NrCls1,NrCls2,DeltaClasses),
    DeltaClasses =< MaxNumberDeltaClasses,
    findall(_,generateAllMatchingClasses,_),
    write(Proj1 - Proj2),
```



```

writef("\n"),
listing(match),
findall(Id,match(Id,_,_,high),HighMatchList),
count(HighMatchList,NrHighMatches),
findall(Id,match(Id,_,_,medium),MedMatchList),
count(MedMatchList,NrMedMatches),
findall(Id,match(Id,_,_,low),LowMatchList),
count(LowMatchList,NrLowMatches),
writef("\nHigh matching classes "),
write(NrHighMatches),
writef(" of "),
write(NrCls1),
writef("\n"),
retractall(match(_,_,_,_)),
TotalMatches is NrHighMatches+ NrMedMatches +NrLowMatches,
delta(TotalMatches,NrCls1,NrUnmatched),
TotalMatches> MinTotalMatches,
NrHighMatches>MinHighMatch,
NrUnmatched <MaxUnmatches,
assert(projectMatch(Proj1,Proj2,high-NrHighMatches,medium-NrMedMatches,low-
NrLowMatches,unmatched-NrUnmatched)).

```

Pentru celelalte două moduri de execuție, când rulăm toate proiectele dintr-un director, ne folosim de predicatul *combine2* , care primește lista tuturor fișierelor obținută de la directivă *directory_files* predefinită în prolog. Restul predicatelor ajutătoare le-am descris prin comentariul pus înaintea lor.

```

%generam calea catre doua proiecte din director.
runFromDir(Directory,Profile):-directory_files(Directory,ListFiles),
    combine2(ListFiles,Proj1,Proj2),
    atom_concat(Directory,Proj1,Project1),
    atom_concat(Directory,Proj2,Project2),
    run(Project1,Project2,Profile).

%din lista data ca parametru generez toate combinarile de cate doua proiecte posibile.
combine2(List,Proj1,Proj2):-
    member(Proj1,List),
    qlfExtension(Proj1),
    member(Proj2,List),
    qlfExtension(Proj2),
    nth1(Index1,List,Proj1),
    nth1(Index2,List,Proj2),
    Index1 <Index2

%verificam ca extensia fisierului sa fie qlf
qlfExtension(Name):-file_name_extension(_, ".qlf", Name).

```

VI. Manualul Utilizatorului

În acest capitol vom prezenta pașii ce trebuie executați pentru a rula aplicația, începând cu modul de exportare a bazei de cunoștințe Prolog și încheind cu apelurile posibile pentru lansarea analizei asupra proiectelor suspecte.

Având instalat plugin-ul JTransformer conform ghidului de instalare oferit de [9] se importă în Eclipse proiectele Java ce se doresc exportate. Pentru fiecare proiect va trebui parcurs următorii pași pentru obținerea bazei de cunoștințe Prolog.

1. Crearea bazei de cunoștințe

Se asignează proiectului o bază de date Prolog, ce va fi calculate de JTransformer. Operația se realizează făcând click dreapta pe proiect și selectând din configure opțiunea *Assign JTransformer Factbase* conform figurii 6.1. Va apărea o fereastră de dialog în care selectăm opțiunea OK. Inacest moment baza de cunoștințe e generată în memoria de lucru.

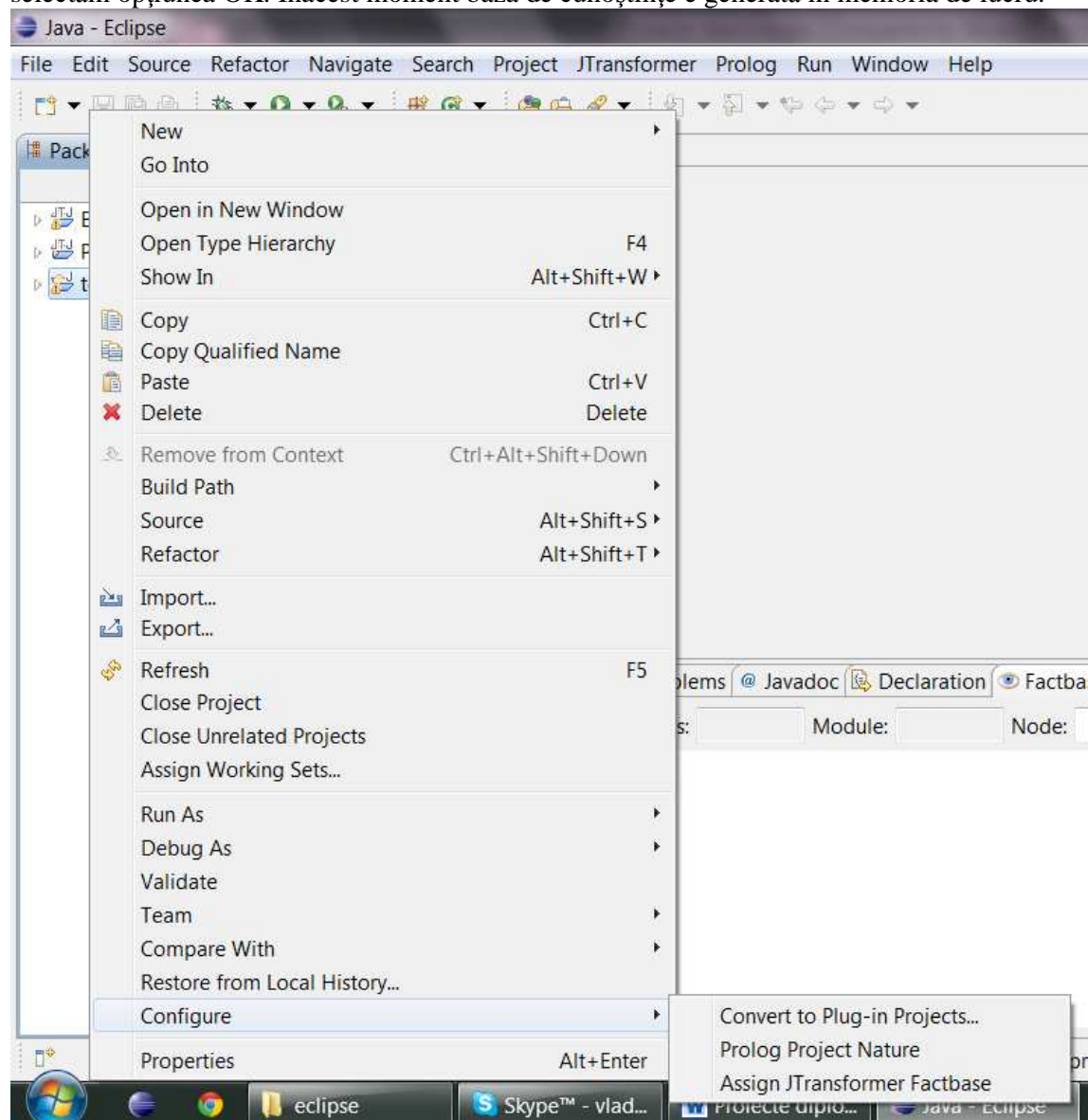


Fig. 7.1 Crearea bazei de cunoștințe

2. Salvarea bazei de cunostinte

În următoarea etapă trebuie să exportăm baza de date creată în memoria de lucru sub formă unor fișiere. Asta se realizează selectând opțiunea Export factbase din meniul contextual al proiectului, conform figurii 6.2. Va apărea fereastra de dialog din figura 6.3. Aici se dă numele bazei de cunoștințe și locul unde va fi salvată. De asemenea, important este bifăm opțiune care ne crează fișierele qlf, cea încercuită cu roșu în imaginea 6.3.

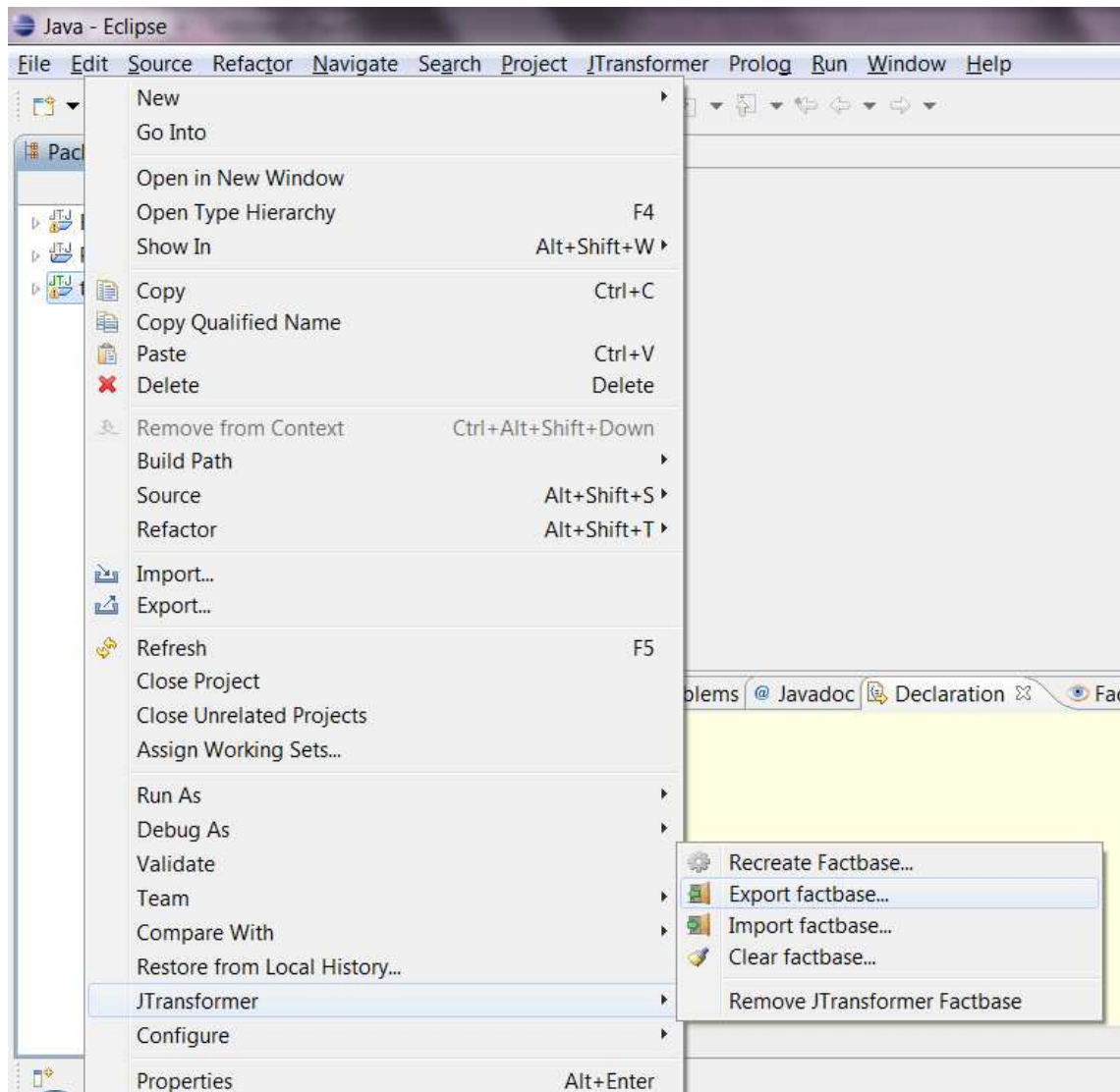


Fig. 6.2 Salvarea bazei de cunostinte în fisier

Baza de cunoștințe va rămâne memorată și în mediul de lucru eclipse unde poate fi interogată prinț facilitățile oferite de JTransformer. Acestea ne-au ajutat mult în dezvoltarea tool-ului dar ca utilizator CloneDetector nu este nevoie să le stăpâniți și nu vor fi descrise în continuare. Dacă doriți să toate facilitățile oferite de JTransformer puteți urmări linkul oferit la [9].

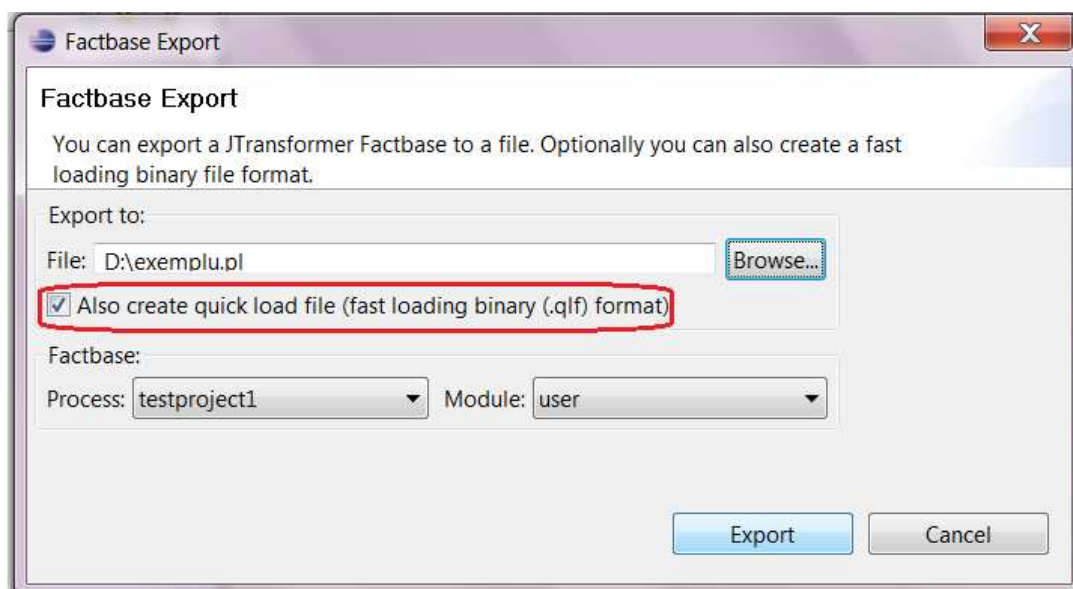


Fig. 6.3 Meniu de salvare

3. Lansarea programului.

Odată create reprezentările proiectelor Java sub forma unor fișiere Prolog, putem lansa în execuție cloneDetector. Având Prolog-ul instalat aceasta se realizează simplu, cu dublu click pe fișierul cloneDetector.pl.

Pentru ca aplicația dezvoltată să fie cât mai utilă profesorilor s-au realizat mai multe moduri de a compara proiectele, apelând diferite predicate.

- Compararea a două proiecte
- Compararea proiectelor dintr-un director, interacționând cu utilizatorul
- Compararea proiectelor dintr-un director, automată

Avem astfel varianta de compara doar două proiecte dând numele lor și profilul. În urma acestei comparări afișându-se clasele ce sunt copiate cu probabilitățile lor. Se specifică cate clase au fost copiate din numărul total. Comandă de lansare în execuție a acestei variante este de forma:

```
run('cale1/numeFisier1.qlf', 'cale2/numeFisier2.qlf', profile).
```

Căile 1 respectiv 2 pot fi absolute sau relative la directorul din care se lansează în execuție CloneDetector. Profilele create la început sunt loose și tight. Se pot crea și altele după cum vom explica la punctul 4. După execuție se va afișa pe ecran ce clase seamănă între ele ca în figura 6.4 . Aceste informații vor fi afișate și în celelalte două variante de execuție, imaginea prezentată fiind pentru cazul doi

Varianta a doua de lansare în execuție a unelei, se face dând ca date de intrare directorul ce conține un set de proiecte sub formă de AST-uri în Prolog și profilul de analiza.

```
runFromDir('caleDirector', profile).
```

Se vor genera grupuri de proiecte combinate două câte două și se va începe analiza. Pentru fiecare două proiecte analizate se vor afișa clasele grupate cu probabilitățile lor și ponderea claselor grupate din numărul total. Nu se trece mai departe la analiza următoarelor două proiecte până ce utilizatorul nu apasă tasta “;”. După analiza tuturor proiectelor se va afișa lista proiectelor găsite suspecte cu o pondere mai mare. Practic această variantă este una

automatizată dar și interactivă, programul așteptând inputul utilizatorului pentru a trece mai departe. Afișarea rezultatelor va fi combinată ca în figura 6.4 pentru fiecare proiect iar la final se va afișa statistica ca în figura 6.5.

```
1 ?- runFromDir('./factbase/', loose).
% ./IPProject1.pl compiled into mod1 0.00 sec, 24,128 bytes
% ./IPProject2.pl compiled into mod2 0.00 sec, 17,424 bytes
% ./factbase/georgescu.qlf loaded into mod1 0.02 sec, 4,770,152 bytes
% ./factbase/ionescu.qlf loaded into mod2 0.02 sec, 4,579,304 bytes
./factbase/georgescu.qlf- ./factbase/ionescu.qlf
:- dynamic match/5.

match(25691, 'Medic', 25693, 'Medic', low).
match(26114, 'Main', 26081, 'Main', high).
match(25701, 'Boala', 25704, 'Boala', high).
match(25725, 'BlackBoard', 25737, 'BlackBoard', high).

High matching classes 3 of 4
true
```

Fig. 6.4 Afișarea rezultatelor pentru compararea a doua proiecte sau a unui întreg director, cu interacțiune din partea utilizatorului

Ultima variantă de a rula programul este una total automatizată. Parametrii de intrare vor fi directorul cu proiectele și profilul selectat.

```
runAll('caleDirector', profile).
```

Programul va genera, asemănător variantei de mai sus, toate combinațiile posibile între proiecte și va începe să le compare două câte două. Diferența între abordarea de sus este că de această dată nu se va mai interacționa cu utilizatorul ci se va continua execuția până se combină toate proiectele. La sfârșit se va afișa lista de proiecte suspecte cu numărul de clase grupate după probabilitățile de a fi copiate. Nu se va mai afișa numele claselor grupate, doar numărul total, de aceea pentru proiectele găsite suspecte se mai poate executa încă odată programul pentru a vedea exact cum au fost asociate clase din cele două proiecte.

```
:- dynamic projectMatch/6.

projectMatch('./factbase/georgescu.qlf', './factbase/ionescu.qlf', high-3, medium-0, low-1, unmatched-0).
projectMatch('./factbase/webserver1.qlf', './factbase/webserver2.qlf', high-4, medium-0, low-0, unmatched-0).
projectMatch('./factbase/passc2.qlf', './factbase/passc2Copy.qlf', high-27, medium-0, low-0, unmatched-0).

true.

2 ?- runAll('./factbase/', loose). █
```

Fig.6.5 Afișarea proiectelor suspecte

În Fig. 6.5 se prezintă modul de afișarea a proiectelor suspecte în cazul evaluării unui director întreg. Se observă că primii doi termeni sunt numele proiectelor urmate de numărul de clase asemănătoare pentru fiecare din cele 3 nivele de similitudine : high, medium, low. Ultimul termen este numărul de clase neimperecheate.

- 4. Editarea sau crearea profilelor.** După cum am menționat de numeroase ori în această lucrare pentru a putea face modificări în modul de analiză ușor în program, în funcție de tipul proiectelor analizate am creat două filtre și opțiunea de a adăuga ușor altele noi. În figura 6.6 se prezintă o parte a fișierului `profiles.pl`, locul unde sunt definite filtrele.

```
%loose profile
projectDelta(loose,
[
    2 % MaxDeltaClasses- diferenta maxima intre nr de clase din primul si cel de-al doilea proiect
]),
classDelta(loose, % profil larg, care genereaza mai multe match-uri [
    5 % delta attributes
    ,0.6% procent attributes
    ,3 % delta methods
    ,1 % delta Interfaces
    ,0.5%ProcentMatchHigh-min procent of method matches for clase to be high
    ,1 %MaxNumberOfUnmathcingMethods - dif. max intre match si nr clase din fiecare proiect
    % pentru a declara o clasa medium
    ,0.5%ProcentMinOfMatchingMethods - min procent pt a declara o clasa med.
]),
methodDelta(loose,
[
    2 %delta if
]),

%-----%

% tight profile. Projects have to be closer mathces
projectDelta(tight,
[
    0 % MaxDeltaClasses- diferenta maxima intre nr de clase din primul si cel de-al doilea proiect
```

Fig. 6.6 Fișierul `profiles.pl`

După cum observăm fiecare profil este compus din trei părți în care sunt definite toleranțele `:projectDelta`, `classDelta` și `methodDelta`. Denumirile sunt sugestive, în zonele acesta descriindu-se toleranțele la nivel de proiect, clase și metode. Primul termen din fiecare clauză este numele profilului, urmat de o listă de numere. Indreptul fiecărui nume există un comentariu care explică ce semnificație are. Unele sunt procente minime altele sunt diferențele maxime pentru anumite toleranțe.

Momentan sunt definite două profile, dar dacă se dorește crearea altora se poate ușor copiind un profil, modificând numele lui într-un nume unic și toleranțele în funcție de necesitate.

- 5. Excluderea claselor din analiza.**

Facilitatea de a exclude clasele din analiza a fost amintită și anterior. În acest paragraf vom explica cum se poate folosi. Pentru setarea claselor ce se doresc excluse din analiză trebuie să modificăm fișierul `commonClasses.pl`. Pentru a exclude o clasă din analiză trebuie să creăm o nouă axiomă în acest fișier care să conțină numele acesteia ca termen. Un exemplu este prezentat în figura 6.6.

```
:-dynamic commonClass/1.
```

```
commonClass('className1.java').
commonClass('className2.java').
commonClass('className3.java').
```

Fig.6.7 Fișierul `commonClasses.pl`

VII. Rezultate

De-a lungul dezvoltării am urmărit să generam perechile de clase suspecte cât mai corect, să nu pierdem perechi și să nu grupăm clase care nu seamănă și să grupăm. Pentru a verifica am rulat unealta pe o serie de proiecte. Incontinuare vom prezenta constatările, cazurile când CloneDetector merge perfect și limitările.

Problema comparării claselor am observat că se poate împărți în 3 cazuri :

1. Compararea claselor cu mai multe de 1-2 metod2. Sunt clase care conțin bucăți din logică proiectului.
2. Compararea claselor foarte mici cu 1 maxim 2 metode. Metodele nu conțin multă logică, nu apelează alte metode, fiind de obicei gettere sau settere.
3. Compararea interfețelor.

Vom discuta primul caz pe exemplu din figura 8.1. Inacest test se compară două proiect cu un număr mic de clase, dar fiecare clasă implementează funcționalitate, având 4-5 metode în medie. Practic clasele acestor proiecte se încadrează în cazul 1. Clasele au suferit modificări de tipul :

- redenumire – majoritatea atributelor și metodelor au fost redenumite. Am redenumit chiar și clasele.
- Schimbarea ordinii. Ordinea metodelor a fost schimbată.
- Ștergerea unor porțiuni. Au fost șterse porțiuni de cod: attribute, mesaje, conditii.

```
match(25696, 'WebServer', 25950, 'AServer', high).
match(26015, 'ConnectionHandler', 25741, 'ManagerConex', high).
match(25700, 'PropertyHolder', 25691, 'ProprietateConex', high).
match(26488, 'Main', 26083, 'Main', high).
```

```
High matching classes 4 of 4
true.
```

```
2 ?-
```

Fig. 7.1 Rezultate dupa compararea unor clase mari.

În ciuda schimbărilor produse asupra codului, CloneDetector a grupat corect toate clasele. Acest lucru s-a întâmplat și în alte cazuri. Deci putem spune că pentru clase de tipul 1 platforma se comportă exemplar.

Pentru a testa celelalte două cazuri am folosit un proiect mai mare, de 27 de clase. Am creat o copie simplă a lui, de tipul 2, în care am modificat doar ordinea claselor în pachete, obținând astfel o baza de date diferită pentru cele două proiecte. Am rulat CloneDetector pe cele două proiecte, în prima fază folosind profilul loose. Toate cele 27 de clase au fost împerecheate cu probabilitate High ,deci verdictul dat de tool este corect, proiectele fiind copiate.

Uitându-ne totuși mai atent am observat că 8 clase din cele 27 nu au fost grupate cu copiile lor din al doilea proiect, ci cu alte clase. Această lucră a avut loc pentru cazurile 2,3. Clasele erau foarte mici, având 1 metodă sau două, cu semănătura la fel, primeau array de bytes și returnau mesajul. Acest lucru era de așteptat, deoarece pentru porțiuni foarte mici de cod, setul de metrice nu le poate diferenția foarte mult iar profilul loose înlătură diferențele între metrice prin toleranțele sale. Aceste clase nu au apeluri către alte clase deci nici compararea apelurilor nu se aplică. De asemenea pentru interfețe folosite doar ca marcă (fără metode) sau care au puține metode cu aceeași semănătură, vom avea aceeași problema și este imposibil chiar în acest caz să diferențiem intefetele pe baza metricilor. Totuși aceste

probleme nu afectează analiza pe ansamblu așa de mult. Pericolul ar putea fi să raportăm prea mulți suspecți pentru un set de proiecte, dar într-un proiect școlar numărul de clase foarte mici și de interfețe este neglijabil în comparație cu numărul total de clase, neinfluentând verdictul. Acest lucru nu se va întâmpla, după cum vom observa și în analiza practică.

```
match(25694, testMars, 26081, testMars, high).
match(25762, 'UnsubscribeMessage', 25768, 'ReplyPriceMessage', high).
match(25793, 'SubscribeMessage', 26173, 'SubscribeMessage', high).
match(25866, 'ReplyAddressMessage', 25904, 'ReplyAddressMessage', high).
match(25723, 'GetAddressMessage', 26144, 'UnsubscribeMessage', high).
match(25741, 'DNSUnmarshaller', 25903, 'ClientUnmarshaller', high).
match(26067, 'StockMarketInterface', 25761, 'StockMarketInterface', high).
match(26074, 'ReplyPriceMessage', 26106, 'GetAddressMessage', high).
match(26111, 'LocateServer', 25805, 'LocateServer', high).
match(26181, 'GetPriceMessage', 25877, 'GetPriceMessage', high).
match(26207, 'ClientUnmarshaller', 25848, 'ClientTransformer', high).
match(26152, 'ClientTransformer', 26123, 'DNSUnmarshaller', high).
match(26353, 'ClientStockMarket', 26051, 'ClientStockMarket', high).
match(26385, 'ServerTransformer', 26448, 'ServerTransformer', high).
match(26386, 'MessageServer', 25695, 'Marshaller', high).
match(26387, 'ServerWithRR', 26450, 'ServerWithRR', high).
match(26113, 'Requestor', 25808, 'Requestor', high).
match(26459, 'Replyer', 26522, 'Replyer', high).
match(26112, 'ByteStreamTransformer', 25807, 'ByteStreamTransformer', high).
match(26453, 'Registry', 26516, 'Registry', high).
match(25697, 'Entry', 25806, 'Entry', high).
match(26862, 'NamingServiceTransformer', 26975, 'NamingServiceTransformer', high).
match(26906, 'NamingServiceServer', 27019, 'NamingServiceServer', high).
match(26442, 'Configuration', 26505, 'Configuration', high).
match(25696, 'Message', 25691, 'Message', high).
match(25695, 'Marshaller', 26449, 'MessageServer', high).
match(25763, 'Address', 25749, 'Address', high).
```

High matching classes 27 of 27
true.

Fig. 7.2. Rezultate cu profilul loose pe un proiect mediu

O soluție pentru problema descrisă mai sus a fost schimbarea profilului. Prin schimbare pe profilul tight am redus numărul de clase grupate greșit de la 8 la 4, după cum putem vedea și în figura 7.3. Totuși schimbarea profilului trebuie să fie făcută în cunoștință de cauză. Prin schimbarea lui, toleranțele vor fi mai mici sau chiar zero, depistându-se copiile exacte. Se poate pierde astfel suspecți care au schimbări în cod. O altă posibilă soluție ar fi adăugarea unor metrice caracteristice pentru metodele mici.

```
match(25694, testMars, 26081, testMars, high).
match(25762, 'UnsubscribeMessage', 26144, 'UnsubscribeMessage', high).
match(25793, 'SubscribeMessage', 26173, 'SubscribeMessage', high).
match(25866, 'ReplyAddressMessage', 25904, 'ReplyAddressMessage', high).
match(25723, 'GetAddressMessage', 25768, 'ReplyPriceMessage', high).
match(25741, 'DNSUnmarshaller', 25903, 'ClientUnmarshaller', high).
match(26067, 'StockMarketInterface', 25761, 'StockMarketInterface', high).
match(26074, 'ReplyPriceMessage', 26106, 'GetAddressMessage', high).
match(26111, 'LocateServer', 25805, 'LocateServer', high).
match(26181, 'GetPriceMessage', 25877, 'GetPriceMessage', high).
match(26207, 'ClientUnmarshaller', 26123, 'DNSUnmarshaller', high).
match(26152, 'ClientTransformer', 25848, 'ClientTransformer', high).
match(26353, 'ClientStockMarket', 26051, 'ClientStockMarket', high).
match(26385, 'ServerTransformer', 26448, 'ServerTransformer', high).
match(26386, 'MessageServer', 26449, 'MessageServer', high).
match(26387, 'ServerWithRR', 26450, 'ServerWithRR', high).
match(26113, 'Requestor', 25808, 'Requestor', high).
match(26459, 'Replyer', 26522, 'Replyer', high).
match(26112, 'ByteStreamTransformer', 25807, 'ByteStreamTransformer', high).
match(26453, 'Registry', 26516, 'Registry', high).
match(25697, 'Entry', 25806, 'Entry', high).
match(26862, 'NamingServiceTransformer', 26975, 'NamingServiceTransformer', high).
match(26906, 'NamingServiceServer', 27019, 'NamingServiceServer', high).
match(26442, 'Configuration', 26505, 'Configuration', high).
match(25696, 'Message', 25691, 'Message', high).
match(25695, 'Marshaller', 25695, 'Marshaller', high).
match(25763, 'Address', 25749, 'Address', high).
```

High matching classes 27 of 27
true.

Fig. 7.3. Rezultate cu profilul tight pe un proiect mediu

Proiectul CloneDetector a fost testat în mai multe circumstanțe și în comparație cu DUDE[10]. DUDE este o unealtă realizată folosind algoritmi de comparare la nivel de șiruri de caracter.

În primă instanță am testat cele două unealte pe teste concepute de noi iar rezultate obținute sunt afișate în tabel.

Caz de test	CloneDetector	Dude
Clona tipul 1	Detectată	Detectată
Clona tipul 2 cu porțiuni redenumite	Detectată	Detectată
Clona tipul 2 total redenumită	Detectată	Nedetectată
Clona tipul 3 partial redenumită	Detectată	Detectată
Clona tipul 3 total redenumită	Detectată	Neraportată
Proiecte diferite	Neraportată	Neraportată
Proiecte mari diferite	Neraportată + timp mare de execuție	Neraportată

Analizând aceste proiecte am observat că platforma DUDE depistează repede orice tip de clonă, dacă există bucăți de cod neredenumite. Dacă clonă a fost total redenumită, DUDE nu va mai depista proiectul plagiat.

CloneDetector a reușit să depisteze clonele pentru testele create, indiferent de numărul redenumirilor.

Ultimele două teste au fost create pentru a vedea cum interpretează unelele proiectele total diferite. S-a observat că ambele unealte nu au raportat proiectele că fiind copii, lucru corect. În ultimul caz au fost date sper comparare proiecte mari, nu școlare, pentru a testa timpii de execuție. În acest caz, pentru un proiect de câteva mii de clase, DUDE a reușit să îl analizeze în aproximativ 30 de secunde timp mai bun față de platforma noastră, care a terminat inspecția în 3 minute. Lucrul acesta era de așteptat, DUDE fiind integrat în alte unealte de analiză software de mari dimensiuni, în vreme ce programul nostru a fost gândit pentru proiecte școlare.

În cea de-a doua etapă am testat cele două proiecte pe un set de teme de laborator de Proiectarea Arhitecturală a Software Complexe. Au fost analizate 16 proiecte cu cele două unealte. Analiza proiectelor a durat mai puțin de 30 de secunde cu CloneDetector și aproximativ 10 cu DUDE, timpi foarte buni pentru ambele.

DUDE a raportat 1 pereche de proiecte copiate, depistând trei porțiuni de cod identice de lungime 7,8 și 9 rânduri. În urma analizării proiectelor s-a confirmat că acestea au fost copiate.

Folosind CloneDetector cu profilul mai larg (acceptă diferențe mai mari între proiecte) au fost depistate 2 perechi de proiecte. Prima a fost cea depistată și de DUDE. A doua pereche a fost analizată și deși structural proiectele erau asemănătoare ele nu au fost copiate.

Schimbând profilul pe cel strict, platforma a raportat doar prima pereche de clone, cea de-a doua fiind exclusă.

VIII. Concluzii

CloneDetector a fost conceput pentru a depista clone în rândul programelor școlare. Pentru a satisface aceste cerințe cât mai bine și a ușura munca utilizatorilor, au fost dezvoltate diverse facilități, fără a afecta scalabilitatea algoritmului.

O facilitate menționată și anterior este aceea de a exclude unele clase, pachete sau foldere din analiză. Aceasta se dovedește foarte utilă în cazul în care toate proiectele au porțiuni de cod date de la început.

După cum am menționat, setarea procentelor sau variațiilor minime/maxime pentru a declara două clase copiate sunt esențiale. Variații mici determină raportarea doar a proiectelor foarte asemănătoare, în vreme ce variațiile mari largesc aria de raportare. În funcție de proiectele propuse și direcțiile de dezvoltare date de profesori unele proiecte pot varia mai mult sau mai puțin. Pentru a putea analiza diferit în funcție de tema dată, am creat două profile de analiză unul larg și al doilea mai exact. Aceste profile pot fi selectate la rulare sau noi profile pot fi ușor definite în fișierul `profiles.pl`.

Varianta actuală a platformei poate fi ușor modificată și îmbunătățită în funcție de nevoile apărute sau a constatărilor utilizatorilor, prin adăugarea unor noi filtre. Aceasta se poate face relativ ușor și cu modificări puține.

Un avantaj dobândit prin folosirea platformei SWI-Prolog este facilitatea acesteia de a lucra multi-threading, sporind viteza de analiză a proiectelor.

Plugin-ul JTransformer are un rol semnificativ dar momentan nu este integrat cu platforma creată. Acest impediment se poate rezolva în viitor integrând CloneDetector cu Eclipse, ca plugin. Până atunci soluția comodă pentru un profesor ar fi de a primi de la studenți odată cu tema și fișierul `qlf` generat.

Pe viitor CloneDetector, poate fi integrat cu serviciile web, după modelul marilor universități Stanford, MIT s.a. Astfel profesorii ar avea acces la unele online dacă își fac un cont. Pentru evaluare ar putea să trimită setul de proiecte primind răspunsul pe email.

CloneDetector se dorește a fi un tool ușor de folosit și eficient în depistarea proiectelor software copiate. În urma testelor realizate asupra unor proiecte concrete, suntem încrezători în utilitatea acestei platforme.

S-a observat că programele ce folosesc algoritmi de comparare a șirurilor de caractere, deși foarte rapide, pot depista doar clonele de tipul 1, eventual 2 și 3 dacă redenumirile nu sunt complete. Instrumentele ce se bazează pe grafuri de apeluri pot semnaliza toate tipurile de clone dar datorită algoritmilor complexi, au o durată de procesare mai mare. CloneDetector poate depista toate clonele de tipul 1 și 2 indiferent de numărul de redenumiri efectuate, semnalând și clone de tipul 3 cu probabilitatea lor de a fi copiate.

Bibliografie

- [1] Stefan Bellon, Rainer Koschke, Giuliano Antoniol, Jens Krinke, Ettore Merlo, “Comparison and Evaluation of Clone Detection Tools”, vol. 33, no.9, IEEE Transaction on Software Engineering, 2007, pp.1–3.
- [2] S. Ducasse, M. Rieger, and S. Demeyer, “A Language Independent Approach for Detecting Duplicated Code,” Proc. Int’l Conf. Software Maintenance (ICSM ’99), 1999.
- [3] B.S. Baker, “On Finding Duplication and Near-Duplication in Large Software Systems,” Proc. Second Working Conf. Reverse Eng., L. Wills, P. Newcomb, and E. Chikofsky, eds., pp. 86-95, July 1995.
- [4] K. Kontogiannis, R.D. Mori, E. Merlo, M. Galler, and M. Bernstein, “Pattern Matching for Clone and Concept Detection,” Automated Software Eng., vol. 3, nos. 1-2, pp. 79-108, June 1996.
- [5] Nghi Troung, Paul Roe, Peter Bancroft, “Static Analysis of Student’s Java Programs”, Queensland University of Technology, Australia, pp.4-5
- [6] J. Krinke, “Identifying Similar Code with Program Dependence Graphs,” Proc. Eighth Working Conf. Reverse Eng. (WCRE’ 01), 2001.
- [7] R. Komondoor and S. Horwitz, “Using Slicing to Identify Duplication in Source Code,” Proc. Int’l Symp. Static Analysis, pp. 40-56, July 2001.
- [8] Peter Bulychev, Marius Minea, “Duplicate code detection using anti-unification”, pp. 2-3
- [9] Proiectul JTransformer. [Online]. Tutorial 2007. Se gaseste la adresa de internet: <http://roots.iai.uni-bonn.de/research/jtransformer/>
- [10] Radu Marinescu, [Online] <http://intooitus.com>
- [11] Paul Clough “Plagiarism in natural and programming languages: an overview of current tools and technologies”, University of Sheffield, Iulie 2000
- [12] Balazinska, M., Merlo, E., Dagenais, M., Lague, B., Kontogiannis, K. : “Measuring clone based reengineering opportunities”, IEEE Symposium on Software Metrics, 1999, pp. 292 – 303.
- [13] Kapser, C., Godfrey, M.W.: “Toward a taxonomy of clones in source code: A case study. In: Evolution of Large Scale Industrial Software Architectures.” (2003)
- [14] Kapser, C., Godfrey, M.: “A taxonomy of clones in source code: The reengineers most wanted list. In: Working Conference on Reverse Engineering “, IEEE Computer Society Press (2003)
- [15] Tudor Girba, Uni. Bern, www.tudorgirba.com
- [16] Xinran Wang, YoonChan Jhi, Sencun Zhu, Peng Liu: “Detecting Software Theft via System Call Based Birthmarks”, Computer Security Applications , 2009
- [17] M. H. Halstead, “Elements of software science”, North Holland, New York, 1977.
- [18] T. J. McCabe, A complexity measure, IEEE Transactions on Software Engineering, SE-2 (4), pp(308-320), 1976.
- [19] Xin Chen, Brent Francia, Ming Li, Brian Mckinnon, Amit Seker,” Shared Information and Program Plagiarism Detection”, Universitatea din California 2003.
- [20] L. M. Ottenstein, “Quantitative estimates of debugging requirements”, IEEE Transactions of Software Engineering, Vol. SE-5, pp(504-514), 1979.
- [21] G. H. Gonnet and R. Baeza-Yates, An Analysis of Karp-Rabin String Matching Algorithm, Information Processing Letters 34, pp(271-274), 1990.
- [22] Stefan Bellon, Rainer Koschke, Giuliano Antoniol, Jens Krinke, Ettore Merlo “ Comparison and Evaluation of Clone Detection Tools”, IEEE TRANSACTIONS ON SOFTWARE ENGINEERING, VOL. 33, Nr. 9, Septembrie 2007
- [23] Paul Clough, Plagiarism in natural and programming languages: an overview of current tools and technologies, University of Sheffield, Iulie 2000